

University of
**Southern
Queensland**

**REVOLUTIONIZING HEALTHCARE WITH
FEDERATED REINFORCEMENT LEARNING:
FROM MACHINE LEARNING TO MACHINE
UNLEARNING**

A Thesis submitted by

Thanveer Basha Shaik
BTech , MSc

For the award of

Doctor of Philosophy

2023

ABSTRACT

The landscape of healthcare is undergoing a transformative shift with the emergence of artificial intelligence (AI) and machine learning (ML) technologies, particularly in remote patient monitoring systems. These systems offer real-time data on patients' health conditions, enabling healthcare professionals to make informed decisions and improve patient outcomes. This doctoral thesis presents a comprehensive investigation into the role of AI in enhancing patient monitoring systems, focusing on innovations in federated learning, reinforcement learning, and machine unlearning across various healthcare settings including remote patient monitoring, personalized activity tracking, mental health facilities, and predictive monitoring. The research outcomes reveal significant advancements in remote patient monitoring through AI-powered systems, enabling early anomaly detection and personalized care. FedStack, a novel federated learning architecture designed for personalized activity monitoring in remote patient monitoring systems, is introduced. Experimental results demonstrate its effectiveness in surpassing traditional approaches and optimizing sensor placement for activity recognition. Furthermore, multi-agent deep reinforcement learning models empower healthcare professionals to predict patient behaviors and take proactive interventions. The exploration of multimodality fusion and graph-enabled techniques demonstrates the potential of comprehensive smart healthcare systems that integrate diverse data sources and enable informed decision-making. Additionally, the thesis introduces a Graph-enabled Reinforcement Learning framework for time series forecasting, leveraging graphical neural networks to outperform traditional models in dynamic environments. The thesis also explores the emerging field of machine unlearning, investigating techniques to address privacy and security concerns. Explainable AI frameworks, such as QXAI, contribute to the reliability and interpretability of patient monitoring systems, fostering trust and collaboration between AI and human experts. FRAMU, a federated reinforcement learning framework with attention-based machine unlearning, is introduced, demonstrating its effectiveness in improving model performance and preserving privacy in dynamic data environments. QXAI, an explainable AI framework for quantitative analysis in patient monitoring, achieves state-of-the-art results in heart rate prediction and activity classification, enhancing model interpretability. While the research demonstrates promising outcomes, it acknowledges certain limitations, including data scale, explainability, and data privacy concerns. Future directions such as dynamic clustering, predictive vital sign monitoring, ensemble methods, and continued progress in machine unlearning are proposed to address these limitations and propel AI-driven patient monitoring systems further. This thesis makes significant contributions to the domain of AI-driven patient monitoring systems, paving the way for personalized, proactive, and effective healthcare delivery globally. It posits that the transformative potential of AI in healthcare is within reach, with continued research and innovation shaping the future of patient care, where AI-driven monitoring becomes an indispensable tool in enhancing patient well-being and transforming healthcare practices.

CERTIFICATION OF THESIS

I, Thanveer Basha Shaik, declare that the Ph.D. thesis entitled *Revolutionizing Healthcare with Federated Reinforcement Learning: From Machine Learning to Machine Unlearning* is not more than 100,000 words in length including quotes and exclusive of tables, figures, appendices, bibliography, references, and footnotes.

This Thesis is the work of Thanveer Basha Shaik except where otherwise acknowledged, with the majority of the contribution to the papers presented as a Thesis by Publication undertaken by the student. The work is original and has not previously been submitted for any other award, except where acknowledged.

Date: 18/10/2023

Endorsed by:

Prof. Xiaohui Tao

Principal Supervisor

Prof. Raj Gururajan

Associate Supervisor

Prof. Xujuan Zhou

Associate Supervisor

A/Prof. Niall Higgins

External Supervisor

STATEMENT OF CONTRIBUTION

In this section, we provide an overview of the publications authored by the student during his candidacy that have been included in this thesis.

1. Shaik, T., Tao, X., Higgins, N., Li, L., Gururajan, R., Zhou, X., Acharya, U. R. (2023). Remote patient monitoring using artificial intelligence: Current state, applications, and challenges. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 13(2), e1485, doi: 10.1002/widm.1485, URL:<https://doi.org/10.1002/widm.1485>. (Q1, Impact Factor 7.8, 2022)

Author	Paper Contribution	Tasks Performed
T. Shaik	75%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	25%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
N. Higgins		
L. Li		
R. Gururajan		
X. Zhou		
U.R. Acharya		

2. Shaik, T., Tao, X., Higgins, N., Gururajan, R., Li, Y., Zhou, X., Acharya, U. R. (2022). FedStack: Personalized activity monitoring using stacked federated learning. *Knowledge-Based Systems*, 257, 109929, doi: 10.1016/j.knosys.2022.109929, URL: <https://doi.org/10.1016/j.knosys.2022.109929>. (Q1, Impact Factor 8.8, 2022)

Author	Paper Contribution	Tasks Performed
T. Shaik	70%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	30%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
N. Higgins		
R. Gururajan		
Y. Li		
X. Zhou		
U.R. Acharya		

3. Shaik, T., Tao, X., Li, L., Higgins, N., Gururajan, R., Zhou, X., & Yong, J. (2024). Clustered FedStack: Intermediate Global Models with Bayesian Information Criterion. *Pattern Recognition Letters*, 177, 121-127, doi: 10.1016/j.patrec.2023.12.004, URL: <http://dx.doi.org/10.1016/j.patrec.2023.12.004>. (Q1, Impact Factor 5.1, 2022)

Author	Paper Contribution	Tasks Performed
T. Shaik	70%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	30%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
N. Higgins		
R. Gururajan		
Y. Li		
X. Zhou		
U.R. Acharya		

4. Shaik, T., Tao, X., Li, L., Xie, H., Dai, H., Yong, J., (2023). Adaptive Multi-Agent Deep Reinforcement Learning for Timely Healthcare Interventions. Submitted to Knowledge-Based Systems (February, 2024). [Under Review]. Preprint available at: <https://arxiv.org/abs/2309.10980>

Author	Paper Contribution	Tasks Performed
T. Shaik	75%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	25%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
H. Xie		
L. Li		
J. Yong		
H. Dai		

5. Shaik, T., Tao, X., Li, L., Xie, H., Acharya, R., Gururajan, R., Zhou, X., (2023) PDRL: Multi-Agent based Reinforcement Learning for Predictive Monitoring. Submitted to Neural Computing and Applications (May, 2023). [Under Review]. Preprint available at: <https://arxiv.org/abs/2309.10576>

Author	Paper Contribution	Tasks Performed
T. Shaik	70%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	30%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
L. Li		
H. Xie		
U. R. Acharya		
R. Gururajan		
X. Zhou		

6. Shaik, T., Tao, X., Li, L., Xie, H., & Velásquez, J. D. (2023). A survey of multimodal information fusion for smart healthcare: Mapping the journey from data to wisdom. *Information Fusion*, 102040, doi:10.1016/j.inffus.2023.102040, URL:<https://doi.org/10.1016/j.inffus.2023.102040> (Q1, Impact Factor 18.6, 2022)

Author	Paper Contribution	Tasks Performed
T. Shaik	70%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	30%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
L. Li		
H. Xie		
J. D. Velásquez		

7. Shaik, T., Tao, X., Xie, H., Li, L., Yong, J., Li, Y. (2023). Graph-enabled Reinforcement Learning for Time Series Forecasting with Adaptive Intelligence. Submitted to *IEEE Transactions on Emerging Topics in Computational Intelligence* (March, 2023). [Under Review]. Preprint available at: <https://arxiv.org/abs/2309.10186>

Author	Paper Contribution	Tasks Performed
T. Shaik	70%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	30%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
H. Xie		
L. Li		
J. Yong		
Y. Li		

8. Shaik, T., Tao, X., Xie, H., Li, L., Velásquez, J., Higgins, N. (2023). QXAI: Explainable AI Framework for Quantitative Analysis in Patient Monitoring Systems. Submitted to *Information Sciences* (January, 2024). [Under Review]. Preprint available at: <https://arxiv.org/abs/2309.10293>

Author	Paper Contribution	Tasks Performed
T. Shaik	65%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	35%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
H. Xie		
L. Li		
J. D. Velásquez		
N. Higgins		

9. Shaik, T., Tao, X., Xie, H., Li, L., Zhu, X., Li, Q. (2023). Exploring the Landscape of Machine Unlearning: A Survey and Taxonomy. Submitted to IEEE Transactions on Information Forensics & Security (February, 2024). [Under Review]. Preprint available at: <https://arxiv.org/abs/2305.06360>

Author	Paper Contribution	Tasks Performed
T. Shaik	75%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	25%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
H. Xie		
L. Li		
X. Zhu		
Q. Li		

10. Shaik, T., Tao, X., Xie, H., Li, L., Cai, T., Zhu, X., Li, Q. (2023). FRAMU: Attention-based Machine Unlearning using Federated Reinforcement Learning. Submitted to IEEE Transactions on Knowledge and Data Engineering (September, 2023). [Under Review]. Preprint available at: <https://arxiv.org/abs/2309.10283>

Author	Paper Contribution	Tasks Performed
T. Shaik	65%	Conceptualization; data curation; formal analysis; investigation; methodology; writing – original draft; writing – review and editing.
X. Tao	35%	Conceptualization; formal analysis; investigation; methodology; project administration; supervision; writing – original draft; writing – review and editing.
L. Li		
H. Xie		
T. Cai		
X. Zhu		
Q. Li		

ACKNOWLEDGEMENTS

First and foremost, I would like to express my deepest gratitude to my principal supervisor, Prof. Xiaohui Tao, for his exceptional guidance, unwavering support, and transformative mentorship throughout my academic journey. Prof. Tao's dedication and belief in my abilities have been instrumental in shaping not only this research project but also my entire career trajectory. Initially uncertain about the path ahead, under Prof. Tao's expert guidance and encouragement, I discovered my true potential and passion for the field of Artificial Intelligence. His profound knowledge, visionary insights, and innovative approach have inspired me to strive for excellence in my academic pursuits. Prof. Tao's mentorship extended beyond the academic realm. He has been a source of inspiration and a pillar of support during challenging times, instilling confidence in me and motivating me to overcome obstacles and push the boundaries of my research. I am grateful to Prof. Tao for providing countless opportunities for growth as a researcher and an academic and for fostering a collaborative and intellectually stimulating environment. His approachability and openness to ideas have made him not just a supervisor but also a trusted mentor. With his guidance, I have successfully published four articles in Q1 journals, and have another six articles under review in Q1 journals for my Ph.D. thesis.

I am also thankful to my associate supervisors, Prof. Raj Gururajan and Prof. Xujian Zhou, for their valuable insights, constructive feedback, and continuous encouragement. Their expertise has significantly enriched this research journey. I extend my thanks to A/Prof. Niall Higgins, my external supervisor, for his valuable input and expertise, which has further enriched the scope of this thesis.

I want to acknowledge the technical and administrative staff at the University of Southern Queensland for their support and assistance in providing the necessary resources and facilities for this research. Additionally, I am grateful to the Graduate Research School and the University of Southern Queensland for awarding me a partial fee scholarship, which has alleviated the financial burden and allowed me to focus on my research.

Special thanks to Dr. Tianning Li, Prof. Jianming Yong, Dr. Ganesh Pant, Dr. Christopher Dann, Prof. Lyn Alderman (Dean, Academic Transformation Portfolio, UniSQ), the Academic Quality Unit, the Academic Integrity Unit, Kishore Bulchandani, Dr. Utsav Bhattarai, Dr. Kishore Aryal, Hamid Ali, Rajan Budhathoki, Hasith Perera, Prasad Surisetty, Niteen, and Neha Rani Thakur.

Finally, I must extend my heartfelt thanks to my family and friends for their unwavering support, love, and understanding throughout this challenging yet rewarding journey. Their constant encouragement and belief in my abilities have been a driving force behind my accomplishments.

DEDICATION

To my beloved family and future wife,

This dissertation is dedicated to you all with all my heart and gratitude. Throughout this challenging and rewarding journey of pursuing my Ph.D., you have been my unwavering source of love, encouragement, and support.

To my parents, your belief in my potential and unwavering encouragement have been the driving force behind my pursuit of knowledge. Your sacrifices, love, and guidance have shaped the person I am today, and I am forever grateful for the values you instilled in me.

To my siblings, your words of encouragement, motivation, and belief in me have been a constant source of inspiration. Your presence in my life has enriched my journey and made it all the more meaningful.

To my future wife, who will join me in the chapters yet to be written, your love and unwavering support mean the world to me. Together, we will embrace the future with shared dreams and aspirations.

This achievement would not have been possible without your unwavering support, understanding, and sacrifices. Your belief in me has been the pillar of strength that kept me going during the most challenging moments.

As I walk this path of knowledge, I carry with me the love, values, and life lessons that you have bestowed upon me. This dissertation is a testament to our collective efforts, and I dedicate it to each one of you with heartfelt gratitude and love.

With love and appreciation,
Thanveer Basha Shaik

TABLE OF CONTENTS

ABSTRACT	i
CERTIFICATION OF THESIS	ii
STATEMENT OF CONTRIBUTION	iii
ACKNOWLEDGEMENTS	vii
DEDICATION	viii
CHAPTER 1: INTRODUCTION	1
1.1 Background	1
1.2 Challenges in AI-driven remote patient monitoring	3
1.2.1 Privacy and security in healthcare	3
1.2.2 Personalized activity monitoring and heterogeneous data analysis	4
1.2.3 Predictive monitoring and model validation	4
1.2.4 Ethical considerations in healthcare	4
1.3 Research motivation	5
1.4 Exploring key technologies	6
1.4.1 Federated learning	6
1.4.2 Reinforcement learning	7
1.4.3 Explainable AI	8
1.4.4 Attention mechanism	9
1.5 Research questions	10
1.6 Objectives	10
1.7 Thesis structure	11
CHAPTER 2: PAPER 1 - REMOTE PATIENT MONITORING USING ARTIFICIAL INTELLIGENCE: CURRENT STATE, APPLICATIONS, AND CHALLENGES	15
2.1 Introduction	15
2.2 Summary	47
CHAPTER 3: PAPER 2 - FEDSTACK: PERSONALIZED ACTIVITY MONITOR- ING USING STACKED FEDERATED LEARNING	48
3.1 Introduction	48
3.2 Summary	63
CHAPTER 4: PAPER 3 - CLUSTERED FEDSTACK: INTERMEDIATE GLOBAL MODELS WITH BAYESIAN INFORMATION CRITERION	64
4.1 Introduction	64

4.2 Summary	72
CHAPTER 5: PAPER 4 - ADAPTIVE MULTI-AGENT DEEP REINFORCEMENT LEARNING FOR TIMELY HEALTHCARE INTERVENTIONS	73
5.1 Introduction	73
5.2 Summary	84
CHAPTER 6: PAPER 5 - PDRL: MULTI-AGENT BASED REINFORCEMENT LEARN- ING FOR PREDICTIVE MONITORING	85
6.1 Introduction	85
6.2 Summary	97
CHAPTER 7: PAPER 6 - A SURVEY OF MULTIMODAL INFORMATION FUSION FOR SMART HEALTHCARE: MAPPING THE JOURNEY FROM DATA TO WIS- DOM	98
7.1 Introduction	98
7.2 Summary	117
CHAPTER 8: PAPER 7 - GRAPH-ENABLED REINFORCEMENT LEARNING FOR TIME SERIES FORECASTING WITH ADAPTIVE INTELLIGENCE	118
8.1 Introduction	118
8.2 Summary	129
CHAPTER 9: PAPER 8 - QXAI: EXPLAINABLE AI FRAMEWORK FOR QUAN- TITATIVE ANALYSIS IN PATIENT MONITORING SYSTEMS	130
9.1 Introduction	130
9.2 Summary	145
CHAPTER 10: PAPER 9 - EXPLORING THE LANDSCAPE OF MACHINE UN- LEARNING: A COMPREHENSIVE SURVEY AND TAXONOMY	146
10.1 Introduction	146
10.2 Summary	171
CHAPTER 11: PAPER 10 - FRAMU: ATTENTION-BASED MACHINE UNLEARN- ING USING FEDERATED REINFORCEMENT LEARNING	172
11.1 Introduction	172
11.2 Summary	187
CHAPTER 12: CONCLUSIONS	188
12.1 Contributions	189
12.2 Limitations	190
12.3 Future directions	191
REFERENCES	193

LIST OF FIGURES

- 1.1 Journey from Machine Learning to Machine Unlearning 3
- 1.2 Thesis Structure 11

- 12.1 Contributions in Revolutionizing Healthcare 189
- 12.2 Machine Unlearning Taxonomy 191

CHAPTER 1: INTRODUCTION

1.1 Background

The contemporary healthcare landscape is undergoing a profound and extensive transformation, primarily driven by the rapid emergence of artificial intelligence (AI) and machine learning (ML) technologies [1]. This transformation is marked by a fundamental paradigm shift that transforms various healthcare domains, with patient monitoring systems prominently positioned as a particularly auspicious application [2]. The advent of these cutting-edge technologies has prompted a fundamental reconfiguration of conventional healthcare practices and processes. In this introductory chapter, we provide an overview of this thesis. We delve into the challenges that have spurred our research motivation, explore the key technologies essential for addressing these challenges, discuss our research questions and objectives, and outline the structure of this thesis.

The significance of AI and ML in this transformative journey is of utmost importance and merits comprehensive examination [3]. These state-of-the-art tools exemplify the technological innovation within the healthcare sector. AI, in particular, has demonstrated an exceptional capacity for pattern recognition, data analysis, and predictive modelling, rendering it an invaluable asset in unravelling the intricacies of patient health [4]. As a subset of AI, ML excels in the development of algorithms and statistical models that facilitate automatic learning from data, subsequently empowering the generation of insights and predictions [5]. Collectively, these technologies serve as the catalysts for propelling innovative solutions within the healthcare landscape.

Amidst the manifold applications of AI and ML in healthcare, patient monitoring systems have emerged as pivotal focal points of transformation [2]. These systems, underpinned by AI and ML, have transcended the boundaries of conventional patient care. They now function as indispensable conduits for the delivery of real-time and continuous patient health data, endowing healthcare professionals with an unprecedented ability to make informed decisions [6].

Machine Learning (ML) and its specialized subset, Deep Learning, have played instrumental roles in the revolutionization of patient monitoring. ML, as a subset of AI, entails the development of algorithms and statistical models that empower computer systems to learn from data and make predictions or decisions grounded in that data [7]. Deep Learning, a subfield of ML, leverages neural networks comprising multiple layers to automatically extract intricate features from complex data, rendering it exceptionally well-suited for tasks encompassing image and speech recognition [8].

In the healthcare domain, these technologies have enabled the analysis of vast datasets encompassing medical images and clinical records. This analytical capability aids in the detection of anomalies, prediction of disease outcomes, and tailoring of personalized treatment plans [9]. Deep Learning, in particular, has garnered notable success in image-based diagnostics, including the detection of tumours in medical images such as X-rays and MRIs [10].

The trajectory of patient monitoring traces its origins back to the early 20th century when rudimentary techniques for measuring physiological parameters laid the foundation for today's sophisticated technologies [11]. After this nascent phase, remarkable advancements in medical sensors, data acquisition methods, and communication technologies have fueled the evolution of patient monitoring systems [12]. These innovations have empowered healthcare professionals to continuously monitor vital signs and health parameters, thereby facilitating the early detection of anomalies and timely interventions [13].

Reinforcement Learning, as a subfield of machine learning, plays a pivotal role in the realm of patient monitoring. Diverging from traditional machine learning approaches where models are trained on labelled data, reinforcement learning centres on training intelligent agents to make sequential decisions by interacting with a dynamic environment [14]. In the context of healthcare, reinforcement learning is harnessed to develop intelligent agents capable of adapting treatment strategies based on patient responses and evolving health conditions [15]. These agents undergo training to optimize actions, including medication dosages or treatment schedules, with the overarching objective of maximizing patient well-being while simultaneously minimizing associated risks and costs [16]. Reinforcement learning also finds application in optimizing resource allocation within healthcare settings, such as staff and equipment assignment, to maximize patient outcomes [17].

Patient monitoring systems hold immense promise due to their inherent capacity to facilitate proactive interventions. Traditionally, healthcare has predominantly adhered to a reactive approach, with treatment initiated in response to the manifestation of symptoms or complications [18]. However, AI-enabled patient monitoring systems harbour the potential to recalibrate this paradigm by enabling the early detection of deviations from baseline health parameters [19]. This early detection, underpinned by AI's predictive capabilities, lays the foundation for timely and precisely targeted interventions. Consequently, patient outcomes stand to be significantly enhanced, and the overall quality of healthcare is poised for elevation [20].

The recent digital revolution has further transformed patient monitoring by seamlessly integrating monitoring devices into electronic health records and hospital information systems [21]. This integration serves to streamline data management processes and augments overall patient care efficiency. Furthermore, the ascent of AI technologies, propelled by groundbreaking advancements in machine learning and deep learning algorithms, has propelled patient monitoring systems to unprecedented heights [22]. AI-driven remote patient monitoring systems have evolved beyond their initial role as mere data collectors; they have now transitioned into intelligent platforms capable of generating predictive insights [23]. These systems, through the analysis of extensive datasets, identification of patterns, and precise predictions, harbour the potential to revolutionize patient care through the delivery of personalized and proactive interventions [24].

A nascent yet critical concept within the healthcare landscape is Machine Unlearning, particularly relevant in the context of patient monitoring. While AI models excel in their ability to learn from data, they may inadvertently accumulate outdated, sensitive, or irrelevant information over time, compromising both privacy and model performance. Machine unlearning emerges as a pragmatic solution to address this challenge by developing techniques to systematically remove such data, thereby ensuring compliance with stringent privacy regulations while concurrently maintaining model relevance [25].

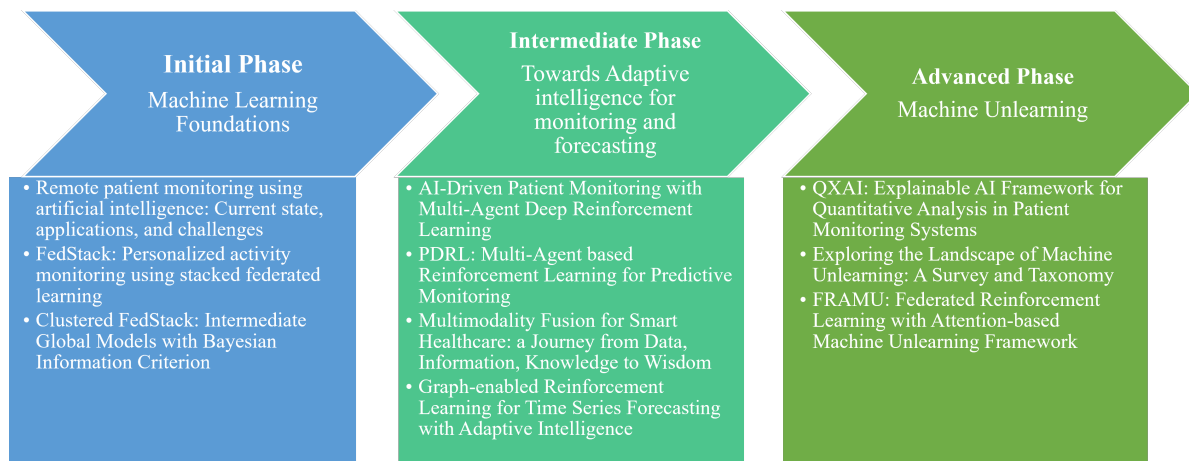


Figure 1.1: Journey from Machine Learning to Machine Unlearning

In the context of healthcare, where data privacy and model accuracy are of paramount importance, machine unlearning assumes a vital role in safeguarding patient information and ensuring that AI-driven patient monitoring systems remain adaptable and responsive to the ever-evolving dynamics of healthcare environments [26]. This concept inherently empowers models to shed obsolete data, thereby perpetuating their effectiveness within dynamic healthcare ecosystems.

The integration of AI and ML technologies into patient monitoring systems signifies an important journey in the evolution of healthcare [27]. These technologies transcend the limitations of traditional healthcare paradigms by facilitating real-time data-driven decision-making and enabling interventions that are not only timely but also profoundly personalized [28]. Subsequent sections of this discourse will delve into the intricate ways in which AI and ML are reshaping patient monitoring, focusing on machine learning and deep learning, historical antecedents of patient monitoring, the role of reinforcement learning, the digital revolution in patient monitoring, and the nascent yet critical concept of machine unlearning within the healthcare landscape. Each of these facets contributes uniquely to the overarching transformation, thereby reinforcing the assertion that AI and ML constitute the fulcrum upon which modern healthcare pivots. Fig. 1.1 presents the journey from machine learning to machine unlearning.

1.2 Challenges in AI-driven remote patient monitoring

While AI-driven remote patient monitoring systems offer immense promise, they also present several critical challenges that must be addressed to fully harness their potential and ensure their effective integration into healthcare practices.

1.2.1 Privacy and security in healthcare

Remote patient monitoring (RPM) holds great potential for revolutionizing healthcare by enabling real-time monitoring of patients in non-clinical environments. However, the decentralized nature of RPM, where data is collected from various remote sources, raises

concerns about patient privacy and data security [29]. Ensuring the confidentiality and integrity of sensitive patient information is paramount, as any breach or unauthorized access could lead to severe consequences for patient trust and overall healthcare operations [30]. Developing robust and secure solutions that protect patient data throughout its transmission and storage in RPM systems is a pressing challenge that demands advanced cryptographic techniques [31], secure communication protocols, and compliance with stringent regulatory standards like HIPAA (Health Insurance Portability and Accountability Act) and GDPR (General Data Protection Regulation) [32].

1.2.2 Personalized activity monitoring and heterogeneous data analysis

AI-driven remote patient monitoring systems have the potential to offer personalized insights and interventions based on individual patient data [33]. However, achieving personalized activity monitoring is challenging due to the diversity of data sources and the heterogeneity of patient data [34]. Healthcare data may come from various devices, electronic health records, wearable sensors, and patient-reported information, leading to variations in data formats, resolutions, and quality [35]. Addressing this challenge requires the development of advanced algorithms capable of analyzing heterogeneous data sources and generating personalized recommendations [36]. Effective solutions call for the implementation of robust data preprocessing techniques, feature extraction methods, and machine learning algorithms adept at handling the intricacies of diverse data modalities [37].

1.2.3 Predictive monitoring and model validation

The realm of predictive monitoring in patient care strives to anticipate critical health events and facilitate proactive interventions [13]. However, the development of accurate and reliable predictive models presents formidable challenges that necessitate extensive research, large-scale data collection, and rigorous model validation. Healthcare data is often susceptible to noise, missing values, and imbalanced distributions, posing additional complexities in training robust predictive models [38]. Ensuring the generalizability and reliability of these predictive models across diverse patient populations and healthcare settings is essential to mitigate potential biases and provide predictions that are trustworthy [39]. The interpretability of predictive models also assumes a pivotal role, enabling healthcare professionals to comprehend and trust the model's predictions, thereby fostering a symbiotic AI-human partnership in patient care [40].

1.2.4 Ethical considerations in healthcare

The application of AI-driven remote patient monitoring in mental health facilities introduces distinctive ethical considerations [41]. While AI technologies offer valuable insights and the potential for early detection of mental health issues, striking a delicate balance between providing personalized care and upholding patient privacy and consent is of paramount importance [42]. Ensuring that AI models are ethically responsible, culturally sensitive, and respectful of patient autonomy and confidentiality holds great significance in mental health monitoring [43].

1.3 Research motivation

The challenges encountered in the realm of AI and machine learning have led to the emergence of a novel paradigm known as machine unlearning (MU). In sectors such as healthcare and finance, the training data for AI models often contains sensitive information. Stricter regulations like GDPR and HIPAA necessitate data handling in compliance with privacy laws. Here, machine unlearning becomes indispensable to remove specific data points, ensuring alignment with data protection regulations.

Another significant challenge arises as new data accumulates over time. Models must dynamically adapt to integrate fresh information. However, reliance on outdated or irrelevant data can significantly hamper performance. This is where machine unlearning steps in as a solution, enabling models to shed obsolete data and thus remain adaptable and responsive. To illustrate, consider a weather forecasting model learning from historical weather patterns. As time passes, some historical data might become irrelevant due to changing climate patterns. Machine unlearning can assist the model in discarding outdated data, ensuring accuracy in predicting current weather trends.

In contexts like healthcare, transparency, and interpretability of AI decisions hold paramount importance. A lack of insight into AI decision-making processes can erode trust. Here, machine unlearning presents an avenue to enhance interpretability by focusing on significant and easily understandable features. For example, in a credit scoring model, machine unlearning could identify and eliminate obscure variables contributing to biased credit decisions, promoting a more transparent and equitable assessment process.

Furthermore, the adaptability of AI models to shifting circumstances is crucial. Unlearning can facilitate the process of "untraining" models from behaviours or patterns that are no longer relevant. Consider a recommendation system for streaming platforms. If user preferences change due to evolving trends, unlearning can help the model adjust from outdated user behaviour patterns, ensuring it offers more accurate and current content suggestions.

Machine unlearning directly addresses the pressing need to enhance privacy and accuracy within AI methodologies. As the research community continues to apply AI methods across diverse domains, challenges related to explainability, interpretability, and privacy are becoming increasingly apparent. These challenges require solutions at a granular level of implementation. Consequently, a pivotal question arises: how can machine unlearning effectively tackle these challenges? The answer lies in leveraging established methods like Federated Learning (FL), Reinforcement Learning, Explainable AI, and Attention mechanisms, as machine unlearning offers a promising approach.

In the context of this thesis, the journey begins by exploring the efficacy of machine learning methodologies and their performance in patient monitoring systems. This exploration then extends to the innovative concept of machine unlearning, as depicted in Fig. 1.1. This concept adeptly addresses challenges linked with outdated, private, and irrelevant data. Simultaneously, it strives to achieve a heightened degree of explainability in model outcomes. By navigating this trajectory, the thesis aims to shed light on the transformative potential of machine unlearning and its pivotal role in propelling forward AI-driven remote patient monitoring systems.

1.4 Exploring key technologies

In the landscape of AI-powered patient monitoring systems, a set of crucial technologies profoundly shape its contours. Together, these technologies chart the course of this thesis journey from machine learning to machine unlearning. This section sheds light on the essence of each technology, laying the groundwork for their integration into AI-driven remote patient monitoring systems.

1.4.1 Federated learning

Federated learning is an innovative approach within the realm of machine learning that offers unique collaborative capabilities. It allows multiple entities, such as hospitals or healthcare institutions, to work together to train a machine learning model without the need to centralize or share their data. This approach has become increasingly important due to growing concerns about data privacy, especially when dealing with sensitive healthcare information [44].

Traditional machine learning models often require centralized data [45]. This means that data from various sources needs to be collected, stored, and processed in a single location. While this centralized approach can be effective for model training, it raises significant privacy concerns [46]. Healthcare data, in particular, is highly sensitive and subject to strict privacy regulations. Centralizing this data increases the risk of data breaches, unauthorized access, and privacy violations [47].

Federated learning offers a solution to this privacy challenge. Instead of pooling all the data into a central server, federated learning conducts model training using localized data. Here's how it works:

- Initialization: A global machine learning model is initialized on a central server.
- Local Training: The local data, which remains on the individual entities' servers (e.g., hospitals), is used to train the global model locally. This means that each entity trains the model using its own data without sharing it with others.
- Model Update: After local training, only the model updates (not the data) are shared with the central server.
- Aggregation: The central server aggregates these model updates to improve the global model. This aggregation process typically involves mathematical operations like averaging or weighted averaging.
- Iterative Process: Steps 2 to 4 are repeated iteratively, with the model getting better after each round of updates.

In healthcare, federated learning has gained significant traction for applications such as remote patient monitoring and personalized healthcare recommendations. Here's how it benefits the healthcare domain:

- Preserving Data Privacy: Healthcare institutions can keep patient data on their premises, ensuring data privacy and compliance with regulations like HIPAA (Health Insurance Portability and Accountability Act) in the United States. Data never leaves the entity's server, reducing the risk of data breaches [48].

- **Collaborative Insights:** Federated learning allows multiple healthcare providers to collaborate on model training without sharing sensitive patient data [49]. This collaboration results in a collective pool of knowledge that caters to the unique health profiles and conditions of different patient populations.
- **Personalized Healthcare:** By training models on a diverse range of patient data from various sources, federated learning enables the development of highly personalized healthcare recommendations and interventions. These recommendations can be tailored to an individual's specific needs and medical history [50].

1.4.2 Reinforcement learning

Reinforcement learning is a dynamic and powerful approach to intelligent decision-making in environments that are constantly changing [51]. At its core, it involves an agent, which can be thought of as an AI entity, interacting with an environment to maximize cumulative rewards [52]. What sets reinforcement learning apart from other machine learning paradigms, like supervised learning, is that it doesn't rely on pre-labelled data [53]. Instead, it learns and improves by trial and error through continuous interactions with its environment.

Reinforcement learning consists of several key components [54]:

1. **Agent:** This is the AI entity or system that is making decisions within the environment.
2. **Environment:** The environment is everything the agent interacts with. It could be a physical space, a virtual world, or even a software application.
3. **Actions:** These are the choices or decisions the agent can make within the environment. Actions can range from simple movements to complex strategies.
4. **Rewards:** Rewards are numerical values that the agent receives from the environment based on the actions it takes. They serve as feedback to the agent, indicating how good or bad its decisions were.
5. **Policy:** The policy is the strategy or set of rules that the agent uses to determine its actions. The goal of the agent is to learn an optimal policy that maximizes its cumulative rewards over time.

In the field of healthcare, reinforcement learning holds great promise. Here are some ways in which it can be applied [55]:

- **Optimizing Patient Treatment:** Reinforcement learning can be used to find the most effective treatment strategies for individual patients. By considering patient data, medical history, and responses to previous treatments, the agent can adapt and optimize treatment plans over time.
- **Fine-Tuning Drug Dosages:** Determining the correct dosage of medication for a patient can be a complex task, as it depends on various factors. Reinforcement learning agents can learn to adjust drug dosages based on patient feedback, ensuring that treatments are both effective and safe.

- **Resource Allocation:** In healthcare settings, resources such as staff, equipment, and hospital beds need to be allocated efficiently. Reinforcement learning can help in making real-time decisions about resource allocation to maximize patient outcomes and minimize costs.
- **Patient Monitoring and Predictive Interventions:** Within patient monitoring systems, reinforcement learning agents analyze patterns of patient behaviour and vital signs. By learning from this data, they can predict when interventions are needed and initiate them proactively. This personalized and preemptive care approach can lead to improved health outcomes and better patient experiences.

The core equation in reinforcement learning is the Bellman equation, which describes how an agent should update its value estimates based on the rewards it receives [56]. It's often written in recursive form as:

$$Q(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \cdot \max_{a'} Q(s', a') \quad (1.1)$$

Where:

- $Q(s, a)$ is the expected cumulative reward of taking action a in state s .
- $R(s, a)$ is the immediate reward of taking action a in state s .
- γ is the discount factor, which represents how much the agent values future rewards.
- $P(s'|s, a)$ is the probability of transitioning to state s' when taking action a in state s .
- $\max_{a'} Q(s', a')$ represents the maximum expected cumulative reward in the next state s' over all possible actions a' .

Reinforcement learning algorithms, like Q-learning and deep reinforcement learning using neural networks, use variations of the Bellman equation to learn optimal policies and make decisions that maximize cumulative rewards in complex and dynamic environments.

1.4.3 Explainable AI

Explainable AI (XAI) stands as a beacon of transparency in the midst of the opacity often associated with the decision-making processes of AI models [57]. The "black box" nature of deep learning models, a term used to describe their inscrutable functioning, can impede the widespread adoption of AI systems, especially in critical domains like healthcare [58].

At the heart of Explainable AI is the aspiration to enable AI models to articulate the rationale behind their decisions [59]. This transparency is achieved through various techniques, including feature attribution, attention mechanisms, and rule-based explanations [60]. These techniques unveil the factors that influence the model's choices, making the decision-making process more comprehensible.

One common method within XAI is feature attribution, which quantifies the importance of each feature in the model's decision [61]. This can be represented mathematically as:

$$\text{Feature Attribution}(x_i) = \frac{\partial \text{Prediction}}{\partial x_i} \quad (1.2)$$

Where:

- Feature Attribution(x_i) represents the contribution of feature x_i to the model's prediction.
- $\frac{\partial \text{Prediction}}{\partial x_i}$ is the partial derivative of the model's prediction with respect to feature x_i .

In the context of neural networks, attention mechanisms are used to identify which parts of the input data are most relevant for a particular task [62]. Mathematically, attention mechanisms can be expressed as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) \cdot V \quad (1.3)$$

Where:

- Attention(Q, K, V) computes the weighted sum of values (V) based on the compatibility of queries (Q) and keys (K).
- The softmax function normalizes the dot product between queries and keys, making it a set of attention scores.
- $\sqrt{d_k}$ is a scaling factor to stabilize the gradients.

Rule-based explanations involve creating human-understandable rules or conditions that the AI model follows when making decisions. For example, a rule might state that if a patient's temperature exceeds a certain threshold, a specific medical action should be taken.

In healthcare, the transparency facilitated by XAI establishes a harmonious partnership between AI algorithms and human caregivers [63]. This partnership instils trust and promotes more enlightened decisions, as healthcare professionals can understand and validate the reasoning behind AI-generated recommendations.

1.4.4 Attention mechanism

The attention mechanism, a cornerstone in the architecture of deep learning models, possesses the remarkable ability to direct the model's focus towards relevant aspects of input data [62]. This selective focus enhances both model performance and interpretability [64]. Widely applied across domains such as natural language processing, computer vision, and sequential data analysis [65, 66], attention mechanisms have gained significant recognition.

At its core, the attention mechanism aims to weigh different parts of input data to varying degrees, emphasising the most relevant elements. This is typically achieved through a process involving queries (Q), keys (K), and values (V), often referred to as the "query-key-value" mechanism.

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) \cdot V \quad (1.4)$$

Where:

- Attention(Q, K, V) computes the weighted sum of values (V) based on the compatibility of queries (Q) and keys (K).
- The softmax function normalizes the dot product between queries and keys, creating a set of attention scores.

- $\sqrt{d_k}$ is a scaling factor used to stabilize gradients.

In the realm of patient monitoring systems, attention mechanisms assume a pivotal role in elucidating the reasoning behind AI-driven decisions. By identifying significant features or temporal aspects within patient data, attention mechanisms offer valuable insights to healthcare professionals. This understanding fosters trust in AI-generated recommendations, creating a collaborative environment where human expertise seamlessly integrates with AI capabilities.

Complementing the technologies discussed above, the innovative concept of machine unlearning emerges as a game-changer. This novel approach tackles the hurdles posed by outdated, sensitive, or irrelevant data. Its essence lies in recalibrating models, discarding obsolete information, and fortifying their adaptability to evolving circumstances. By aligning these technologies in focus, the subsequent chapters of this thesis embark on an expedition to showcase their transformative impact within AI-driven remote patient monitoring systems. This journey delves into pioneering methodologies and applications that leverage these technologies in machine learning and machine unlearning, heralding a revolution in patient care and elevating healthcare practices to unprecedented heights.

1.5 Research questions

In light of the challenges and opportunities presented by AI-driven remote patient monitoring systems, this doctoral thesis embarks on a comprehensive investigation of the role of AI in enhancing patient monitoring systems across various healthcare settings. The research endeavours to answer the following key questions:

1. How can artificial intelligence be leveraged to enhance patient monitoring systems, addressing challenges in remote patient monitoring, personalized activity tracking, and predictive monitoring in healthcare?
2. What are the advancements and implications of federated learning and reinforcement learning techniques in patient monitoring systems, particularly in personalized activity monitoring, human behaviour monitoring, and time series forecasting?
3. What is the potential of multimodality fusion and graph-enabled techniques in creating comprehensive smart healthcare systems that integrate data, information, and knowledge to support informed decision-making?
4. How do explainable AI frameworks contribute to the reliability and interpretability of patient monitoring systems, and how does attention-based machine unlearning further enhance federated reinforcement learning?

1.6 Objectives

With the aim of advancing knowledge and practice in AI-driven remote patient monitoring, this thesis is designed to achieve the following objectives:

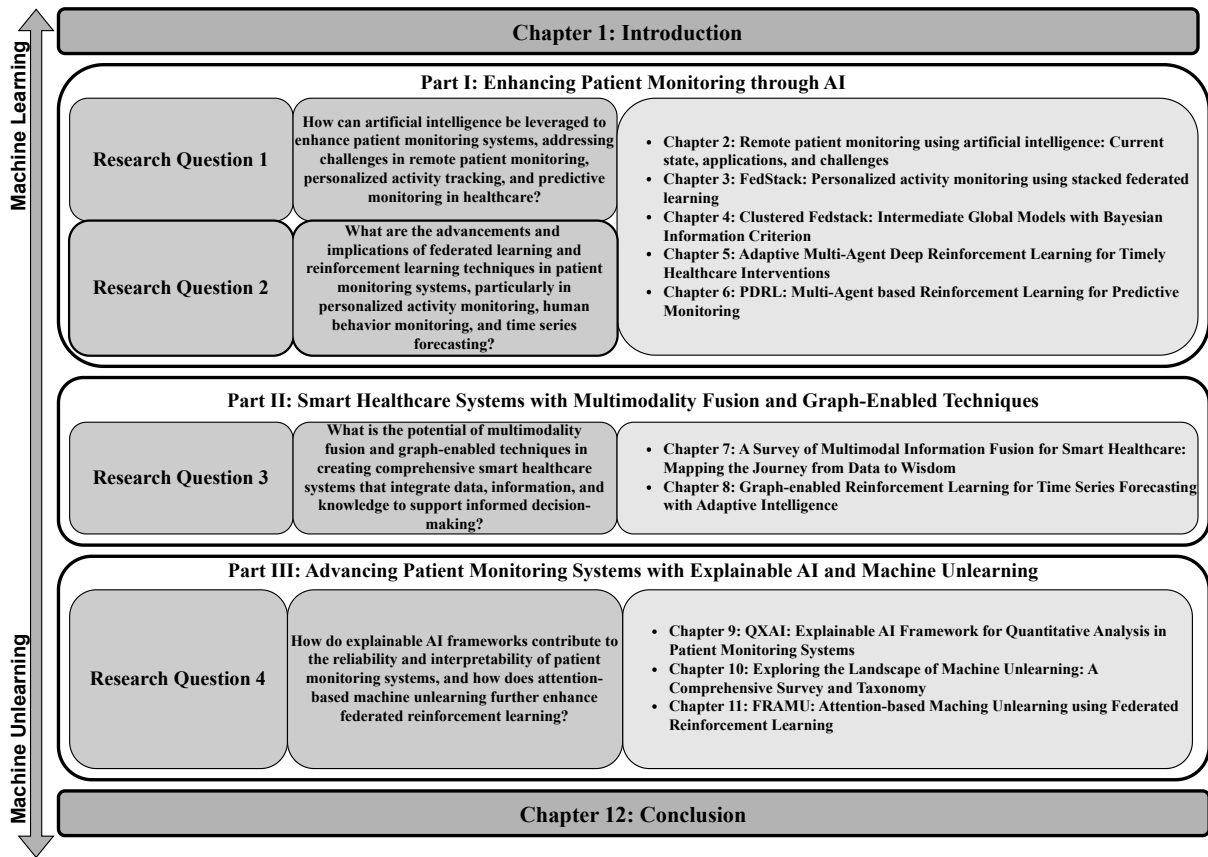


Figure 1.2: Thesis Structure

- To critically evaluate the potential of artificial intelligence in advancing patient monitoring systems, with a focus on devising strategies for safeguarding patient privacy and ensuring effective model evolution.
- To explore and analyze the advancements and implications of reinforcement learning techniques within patient monitoring systems, specifically concentrating on their application in the dynamic context of domain adaptation.
- To investigate the role of explainable AI frameworks in enhancing the trustworthiness and comprehensibility of patient monitoring systems, thereby mitigating the challenge of model opacity inherent in machine learning.
- To assess the viability and effectiveness of machine unlearning as a solution to address privacy-related challenges and to facilitate model refinement, leveraging the advancements offered by federated learning, reinforcement learning, and attention mechanisms.

1.7 Thesis structure

This thesis is structured into three parts each addressing specific aspects of AI-driven remote patient monitoring.

Part I: Enhancing patient monitoring through AI

Part I focuses on exploring the potential of AI to enhance patient monitoring systems, addressing challenges in remote patient monitoring and personalized activity tracking in healthcare via privacy preservation techniques and model update techniques. The chapters in this part delve into the innovative approaches and advancements that AI brings to patient monitoring.

Chapter 2: Remote patient monitoring using artificial intelligence: Current state, applications, and challenges

This chapter provides an in-depth survey of the current state of remote patient monitoring (RPM) in modern healthcare. It highlights the diverse applications of RPM, its impact on patient care, and the challenges faced in implementing effective RPM systems. The chapter also explores the potential of AI-driven RPM solutions in enabling real-time monitoring and early detection of health anomalies, with a particular focus on patient privacy and data security concerns.

Chapter 3: FedStack: Personalized activity monitoring using stacked federated learning

In this chapter, the research introduces an innovative approach to personalized activity tracking using stacked federated learning. The chapter explains the concept of federated learning and its application to personalized activity monitoring. It presents the FedStack architecture, which ensembles local models into a robust global model, surpassing traditional federated learning approaches. The chapter demonstrates how FedStack can provide tailored insights to individual clients, enabling personalized care and interventions based on specific patient needs.

Chapter 4: Clustered Fedstack: Intermediate Global Models with Bayesian Information Criterion

This chapter explores the Clustered FedStack framework, which introduces intermediate global models with Bayesian Information Criterion (BIC). By incorporating BIC, the chapter demonstrates how Clustered FedStack can optimize the selection of intermediate global models, improving the overall efficiency and effectiveness of personalized activity monitoring. This advancement further empowers healthcare professionals with more accurate and reliable insights for proactive patient care.

Chapter 5: Adaptive Multi-Agent Deep Reinforcement Learning for Timely Healthcare Interventions

Building on the concepts of reinforcement learning and AI-driven remote patient monitoring, this chapter introduces the paradigm of multi-agent deep reinforcement learning (DRL) for patient monitoring. The chapter presents AI-driven agents capable of learning behaviour patterns and predicting appropriate actions in dynamic healthcare environments. The novel AI-driven patient monitoring framework has the potential to revolutionize patient care by enabling proactive interventions and personalized treatment strategies.

Chapter 6: PDRL: Multi-Agent based Reinforcement Learning for Predictive Monitoring

This chapter introduces the novel PDRL framework, which extends predictive monitoring capabilities using multi-agent reinforcement learning. By empowering AI agents

to forecast future vital signs and health events, PDRL enhances the potential for early detection and preventive measures in patient care. The chapter highlights how PDRL contributes to reshaping patient monitoring systems by providing real-time, context-aware predictions for improved healthcare outcomes.

Part II: Smart Healthcare Systems with Multimodality Fusion and Graph-Enabled Techniques

Part II delves into the realms of multimodality fusion and graph-enabled techniques, envisioning comprehensive smart healthcare systems that integrate data, information, and knowledge. This encompasses domain adaptation through Reinforcement Learning. The chapters in this part explore how AI can transform patient monitoring through multimodal data integration and dynamic reinforcement learning.

Chapter 7: A Survey of Multimodal Information Fusion for Smart Healthcare: Mapping the Journey from Data to Wisdom

This chapter explores the transformative potential of multimodality fusion in creating comprehensive smart healthcare systems. It discusses the integration of various AI techniques, including feature selection, machine learning, and natural language processing, to extract valuable insights from diverse data sources. The chapter emphasizes how this advancement facilitates informed decision-making, fosters knowledge accumulation, and lays the foundation for predictive, preventive, personalized, and participatory healthcare approaches.

Chapter 8: Graph-enabled Reinforcement Learning for Time Series Forecasting with Adaptive Intelligence

In this chapter, the research investigates the integration of graph-enabled reinforcement learning techniques to revolutionize time series forecasting in patient monitoring. The chapter introduces the GraphRL framework, utilizing Temporal Graphical Convolutional Networks (T-GCN) to forecast dynamic reinforcement learning scenarios. The chapter showcases the potential of GraphRL beyond healthcare, presenting versatile solutions for diverse prediction tasks.

Part III: Advancing Patient Monitoring Systems with Explainable AI and Machine Unlearning

Part III embarks on an exploration of explainable AI and the profound impact of machine unlearning on patient monitoring systems. One of the chapters in this part focuses on enhancing the reliability and interpretability of AI-driven decisions, fostering a collaborative AI-human partnership. It delves deeply into the taxonomy, challenges, and emerging trends of machine unlearning. Building upon the knowledge acquired in preceding chapters, this section endeavours to synthesize the insights gained into a comprehensive framework. The primary goal of this framework is to facilitate the systematic removal of obsolete, confidential, and irrelevant data. This inclusive framework is designed to be applicable in various contexts, encompassing both single-modality and multi-modality scenarios, with a particular emphasis on the healthcare domain.

Chapter 9: QXAI: Explainable AI Framework for Quantitative Analysis in Patient Monitoring Systems

This chapter in Part III introduces the pioneering QXAI framework, which significantly contributes to enhancing the explainability and interpretability of AI-driven remote patient monitoring systems. By providing transparent and interpretable predictions, QXAI fosters trust in AI-driven recommendations, enabling a collaborative AI-human partnership. The chapter emphasizes the paramount importance of explainable AI in critical domains like healthcare, where decisions have significant implications for patient well-being.

Chapter 10: Exploring the Landscape of Machine Unlearning: A Comprehensive Survey and Taxonomy

This chapter embarks on a profound journey to explore the uncharted territory of machine unlearning. It establishes the significance of dynamic model refinement, enabling AI models to shed outdated information, adapt to changing patient conditions, and maintain relevance in the face of evolving healthcare dynamics. The chapter presents a comprehensive survey and taxonomy of machine unlearning techniques, paving the way for future research and advancements in this field.

Chapter 11: FRAMU: Attention-based Maching Unlearning using Federated Reinforcement Learning

In this chapter, the research introduces the groundbreaking FRAMU framework, a novel intersection of federated reinforcement learning and attention-based machine unlearning. The framework elucidates the potential to create AI models that continuously learn and unlearn, ensuring robustness and reliability across disparate healthcare environments. The chapter showcases how FRAMU can contribute to enhancing patient monitoring systems and fostering trust in AI-driven recommendations.

As we embark on this academic journey, this thesis seeks to contribute to the growing body of knowledge in the domain of AI-driven remote patient monitoring systems. By addressing the research questions and objectives outlined, this research endeavours to shed light on the challenges and opportunities in machine learning and machine unlearning, guiding the seamless integration of AI technologies into healthcare practices. The subsequent chapters present a detailed exploration of each research question, drawing upon rigorous methodologies, experimental investigations, and meticulous analysis of findings.

CHAPTER 2: PAPER 1 - REMOTE PATIENT MONITORING USING ARTIFICIAL INTELLIGENCE: CURRENT STATE, APPLICATIONS, AND CHALLENGES

2.1 Introduction

This chapter marks the beginning of an insightful journey into the realm of RPM, situated within the expansive landscape of digital health transformation. It underscores RPM's crucial role in transcending conventional clinical boundaries to integrate healthcare into the fabric of everyday life. Through the adoption of state-of-the-art technologies like the Internet of Things (IoT), wearable sensors, and AI, RPM stands out as a key element in the interconnected narrative of this thesis. This narrative weaves through the diverse impacts of digital innovations on patient care, data analytics, and healthcare delivery models. The introduction outlines RPM's scope and capabilities while subtly drawing connections to forthcoming chapters that delve into AI's role in data analysis, patient engagement, and the ethical considerations in digital health. By exploring the current state, applications, and challenges of AI-enhanced patient monitoring, the chapter sets the stage for addressing critical gaps and fostering innovative solutions in the field.

ADVANCED REVIEW



WILEY

Remote patient monitoring using artificial intelligence: Current state, applications, and challenges

Thanveer Shaik¹ | Xiaohui Tao¹ | Niall Higgins^{2,3} | Lin Li⁴ |
Raj Gururajan⁵ | Xujuan Zhou⁵ | U. Rajendra Acharya⁶

¹School of Mathematics, Physics and Computing, University of Southern Queensland, Toowoomba, Australia

²Metro North Hospital and Health Service, Royal Brisbane and Women's Hospital, Brisbane, Australia

³School of Nursing, Queensland University of Technology, Brisbane, Australia

⁴School of Computer and Artificial Intelligence, Wuhan University of Technology, Wuhan, China

⁵School of Business, University of Southern Queensland, Springfield, Australia

⁶School of Science and Technology, Singapore University of Social Sciences, Singapore

Correspondence

Thanveer Shaik, School of Mathematics, Physics and Computing, University of Southern Queensland, Toowoomba, Queensland, Australia.

Email: thanveer.shaik@usq.edu.au

Edited by: Tianrui Li, Associate Editor and Witold Pedrycz, Editor in Chief

Abstract

The adoption of artificial intelligence (AI) in healthcare is growing rapidly. Remote patient monitoring (RPM) is one of the common healthcare applications that assist doctors to monitor patients with chronic or acute illness at remote locations, elderly people in-home care, and even hospitalized patients. The reliability of manual patient monitoring systems depends on staff time management which is dependent on their workload. Conventional patient monitoring involves invasive approaches which require skin contact to monitor health status. This study aims to do a comprehensive review of RPM systems including adopted advanced technologies, AI impact on RPM, challenges and trends in AI-enabled RPM. This review explores the benefits and challenges of patient-centric RPM architectures enabled with Internet of Things wearable devices and sensors using the cloud, fog, edge, and blockchain technologies. The role of AI in RPM ranges from physical activity classification to chronic disease monitoring and vital signs monitoring in emergency settings. This review results show that AI-enabled RPM architectures have transformed healthcare monitoring applications because of their ability to detect early deterioration in patients' health, personalize individual patient health parameter monitoring using federated learning, and learn human behavior patterns using techniques such as reinforcement learning. This review discusses the challenges and trends to adopt AI to RPM systems and implementation issues. The future directions of AI in RPM applications are analyzed based on the challenges and trends.

This article is categorized under:

Application Areas > Health Care
Technologies > Artificial Intelligence
Technologies > Internet of Things

KEYWORDS

artificial intelligence, IoT, noninvasive technology, remote patient monitoring

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2023 The Authors. *WIREs Data Mining and Knowledge Discovery* published by Wiley Periodicals LLC.

1 | INTRODUCTION

Remote patient monitoring (RPM) is a rapidly growing field in healthcare that is designed to assist clinicians with additional support to provide care in a range of general hospital medical and surgical wards and using flexible materials for wearable sensors (Joshi et al., 2021; Liu, Wang, et al., 2022; Weenk et al., 2020). This is achieved by incorporating new Internet of Things (IoT) methodologies in healthcare such as telehealth applications (Heijmans et al., 2019), wearable devices (Dias & Cunha, 2018), and contact-based sensors (Malasinghe et al., 2017). RPM is commonly used to measure vital signs or other physiological parameters such as motion recognition that can assist with clinical judgments or treatment plans for conditions such as movement disorders or psychological conditions (Shaik, Tao, Higgins, Gururajan, et al., 2022; Shaik, Tao, Higgins, Xie, et al., 2022).

Artificial intelligence (AI) algorithms have been employed to perform analysis of medical images and correlate symptoms and biomarkers from clinical data to characterize an illness and its prognosis (Miller & Brown, 2018; Schnyer et al., 2017). There is immense potential for AI to benefit healthcare service delivery and clinicians are exploring a variety of practical issues for assessing the risk of disease, ongoing patient care, and how AI can help clinicians to alleviate or reduce complications in illness progression (Torous et al., 2018). Medical research is also benefitting from AI by helping to expedite genome sequencing and the development of new drugs and treatments from the knowledge that previously was not possible to obtain or observe from such complex data. Machine learning, a subset of AI, can potentially assist clinicians in interpreting complex data in a relatively short period using specialized algorithms (Helm et al., 2020; Krittanawong et al., 2022). They can assist with a patient assessment to help predict early deterioration of their health status and even classify their types of motion or activities (Z. Liu, Zhu, et al., 2022; Huang et al., 2022). These AI algorithms can process large datasets to recognize and learn complex patterns for decision-making (Dean et al., 2022). Recent increases in computational speed have led to the development of even more powerful artificial neural networks and deep learning algorithms that can handle and optimize very complex datasets (Bini, 2018; Kalfa et al., 2020). Many routine tasks can be automated by incorporating an IoT model with a centralized control unit and interface. This could potentially avoid human errors, and increase patient safety (Tandel et al., 2022).

RPM has traditionally been applied to monitoring patients in rural areas remotely using telehealth technology, monitoring chronically ill people, and the elderly at home using wearable devices or sensors, but the non-intrusive aspects are also attractive for use in hospitals for post-surgery patients, and those in intensive care units using wireless body sensors. It is possible to enhance these monitoring systems to the next level by introducing noninvasive digital technologies which permit patients' daily activities. To support healthcare professionals in visualizing the health status of patients based on vital signs and activity recognition, machine learning (ML) and AI can be implemented as shown in Figure 1. These types of applications can present data related to diagnosing and predicting patient health status and assist with clinical decision-making. This review is motivated by potential advancements in healthcare using AI and machine learning to transform existing traditional medical practices.

This review aims to investigate technologies adopted in current RPM systems for noninvasive techniques. Current trends in RPM and applications of AI to monitor vital signs, physical activities, emergency events, and chronic diseases of patients and assist clinicians to diagnose and provide efficient care. The impact of AI on RPM applications for early detection of health deterioration, personalized monitoring, and adaptive learning are discussed. Finally, the current challenges to the widespread adoption of remote monitoring with AI or machine learning in healthcare are presented and identify what is being done to address these. The contributions of this study are:

- AI impact of RPM applications is investigated and stressed the need for early detection of health deterioration.
- Traditional machine learning and deep learning applications in RPM are investigated.
- Comprehensive review of advanced technologies such as video-based monitoring, IoT-enabled devices, cloud, edge, fog, and blockchain and AI methodologies such as reinforcement learning, and federated learning adopted by RPM systems.
- Challenges in adopting AI-enabled RPM are investigated, and their trends are explored.

The review is organized as follows: Section 2 presents the research explored in this study, search strategies, and inclusion criteria. Section 3 presents advanced RPM architectures including telehealth, IoT, cloud, fog, edge, and blockchain technologies. The scope of AI in RPM applications like monitoring vital signs, physical activities, emergencies, and chronic diseases are discussed in Section 3. In Section 4, the impact of AI on RPM has been

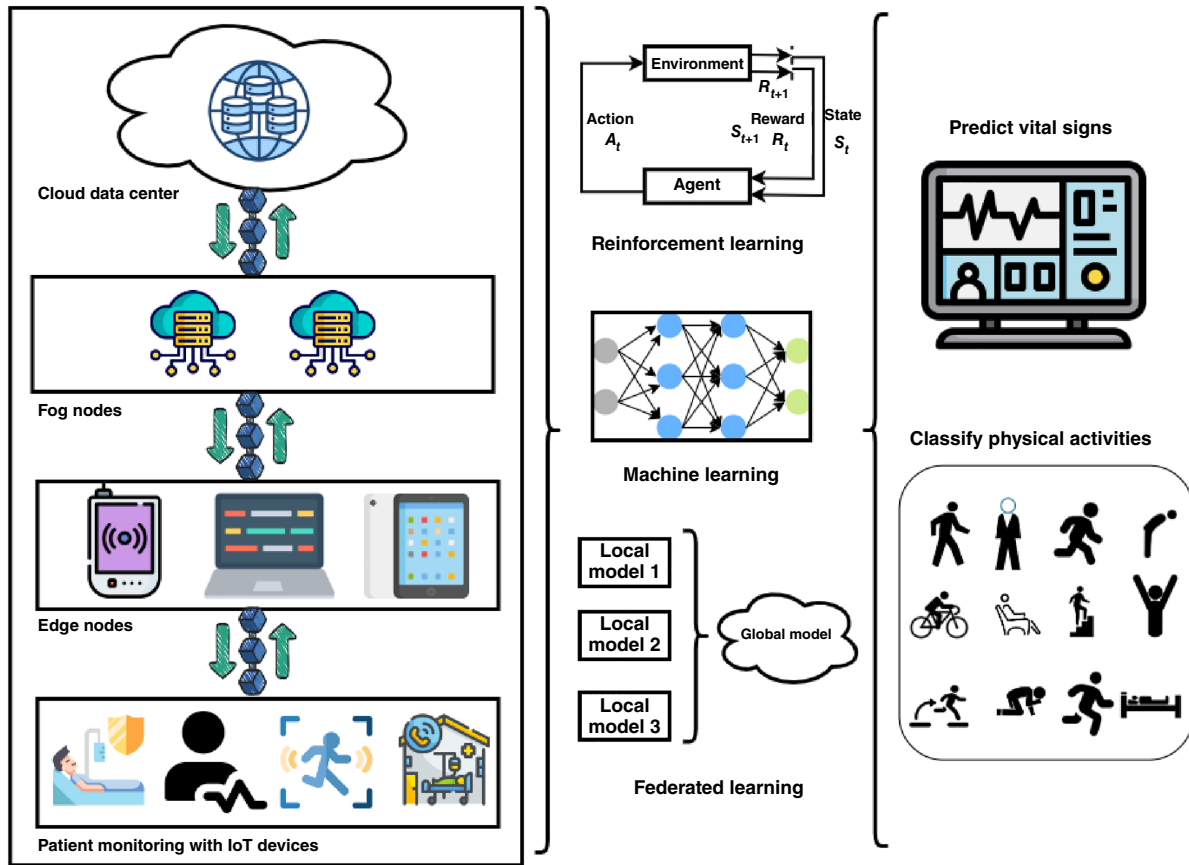


FIGURE 1 Artificial intelligence-enabled remote patient monitoring architectures

discussed. Section 5 describes the challenges involved in adopting AI or machine learning to monitor patients. Section 6 concludes the article with a finding summary and future work with recommendations.

2 | SEARCH STRATEGY AND SELECTION CRITERIA

The objective of this review is to identify journal articles, review articles, and conference papers related to the role of AI in monitoring a patient's health status. This can be done using IoT devices in geographically remote settings or more locally through nontouch techniques. In doing so, the review will seek to address the following research questions:

- RQ1.** What technologies have transformed conventional manual patient monitoring in hospitals?
- RQ2.** How has AI transformed the RPM with its advancements and their impacts?
- RQ3.** What are the challenges in adopting AI for RPM systems and learning healthcare data?
- RQ4.** What are the existing trends in RPM systems for using AI?

TABLE 1 Selected limits for database searches

Inclusion criteria	Exclusion criteria
Journal Article	Books
Review Article	Book chapters
Conference Paper	Abstracts
Conference Paper Review	Short survey Editorial
Published between 2016 and 2022	Letter
Literature in English	Research works related to infants, neonates
Outpatients and inpatients	Experiments on animals
Employs AI & ML	Research work without AI & ML
Experiments on adult and elderly patients	Image processing techniques

Abbreviations: AI, artificial intelligence; ML, machine learning.

2.1 | Information sources

Literature was retrieved from the following bibliographic databases: Web of Science, Scopus, Springer, ACM Digital Library, IEEE Xplore, Pub-Med, Science Direct, and Multidisciplinary Digital Publishing Institute (MDPI). Search strategies were defined using keywords, Boolean operators, truncation, and wildcards. Each database was filtered to search the keywords and their combinations in the title, abstract, and keywords. Results were sorted by relevance, and the first 10 results were checked to ensure that a combination of search terms retrieved articles relevant to the research questions. Finally, the results were exported to EndNote and grouped for each database. Furthermore, the EndNote citations were exported to software called Rayyan (Ouzzani, Hammady, Fedorowicz, & Elmagarmid, 2016) to facilitate the screening and selection process. As databases could host articles published elsewhere, duplicate articles were excluded.

2.2 | Search strategy

Before defining the keywords, a random search was conducted to identify keywords that have been previously used to retrieve relevant articles on these topics. Once this was completed, the main concepts in each article title were categorized into five areas, and synonyms were created using a thesaurus. Table I presents the keywords used for each concept's search, with truncation and wildcards.

Boolean operators AND, and OR were used to form different combinations of keywords within the five key concepts. The final search string used in all the databases was:

(patient? OR victim? OR case? OR subject? OR human?) AND (observ* OR monitor* OR audit* OR detect* OR estimat* OR forecast* OR check*) AND (remote OR distan* OR isolated OR inaccessib* OR outlying) AND (RFID OR sensor* OR wire* OR accelerometer OR doppler OR ECG OR radio* OR polysomno*) AND (“artificial intelligence” OR AI OR “machine learning” OR “neural networks”)

2.3 | Selection criteria

As the review focused on the implementation of AI or machine learning in collaboration with information systems infrastructures, this study excluded research articles with continuous monitoring without AI or machine learning. Table 1 presents the chosen limits used to retrieve articles published between 2016 and 2021 and Figure 2 presents a PRISMA flowchart for the review process.

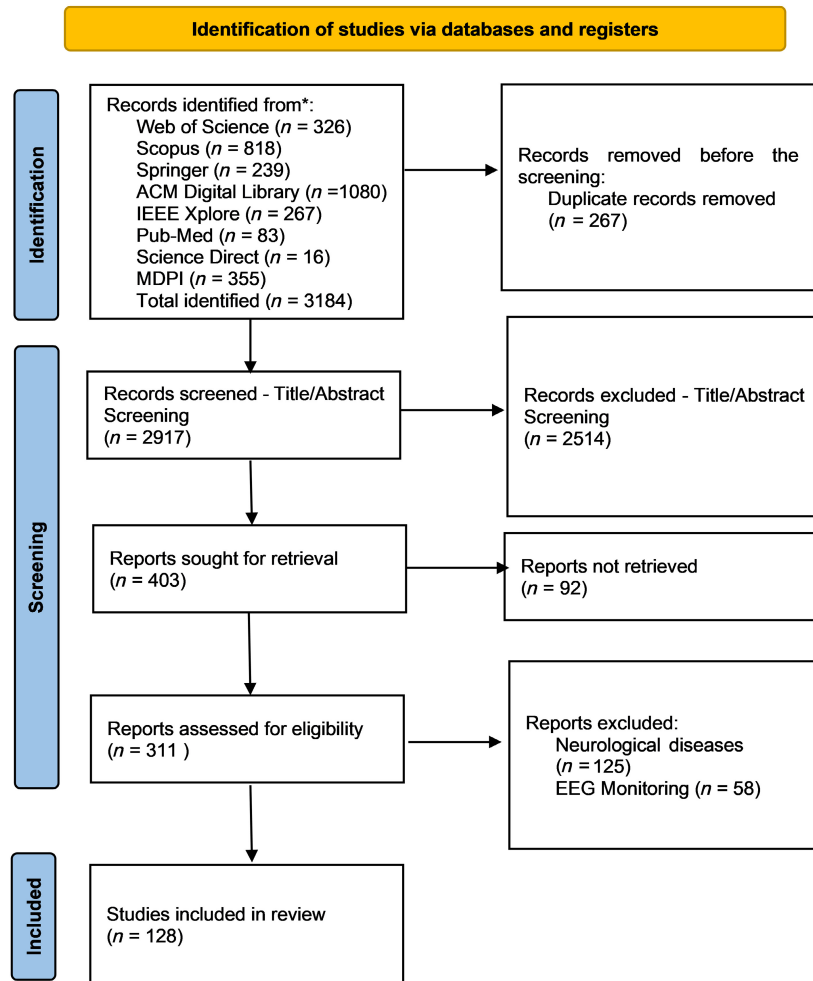


FIGURE 2 Study flow diagram PRISMA-ScR (Page et al., 2021)

3 | REMOTE PATIENT MONITORING ARCHITECTURES

In hospitals, medical staff monitored a patient's health status regularly and manually maintained the records. Collecting patients' vital signs manually in hospitals depends on factors like clinical workload, staff working hours, patients' diagnosis, clinical leadership, and national guidance (Smith et al., 2017) and was limited due to the lack of resources. The patient monitoring was with invasive devices requiring patients' skin contact to estimate their vital signs. Technological advancements in data transmission have disrupted the healthcare industry with noninvasive devices without touching patients' bodies and provided opportunities to monitor patients continuously. The innovations have transformed the traditional patients' health status monitoring patterns and enabled to monitor of patients remotely in hospitals, patient care hospitals, age care facilities, and even in their homes. In this section, technology-enabled RPM architectures are discussed.

3.1 | Video-based monitoring

Telehealth monitoring allows patients to contact their doctors or medical staff via audio call or video call using smart devices. Snoswell et al. conducted a systematic review to measure the clinical effectiveness of telehealth applications. In

their study, Snoswell et al. (2021) reviewed 38 meta-analysis articles published between 2010 and 2019 and covered 10 medical disciplines including multidisciplinary care and other specialized disciplines like cardiovascular disease and pulmonary diseases. The authors reported that the usage of telehealth has exponentially increased over the last decade and demonstrated examples from mental health support, pain management, blood pressure and glucose control, stroke management, and diagnostic services like dermatological and ophthalmic conditions.

As a result of the COVID-19 pandemic, telehealth became a common strategy for maintaining patients' and clinicians' safety. Machine learning and image processing techniques played a vital role in telehealth monitoring. The AI methods are capable of monitoring patients' vital signs such as heart rate, respiratory rate, oxygen saturation (SpO₂), cough analysis, and blood pressure. Rohmetra surveyed AI-enabled telehealth monitoring of vital signs and compared these with traditional methods of monitoring vital signs (Rohmetra et al., 2021). The image and video processing techniques in ML helped identify a region of interest (ROI) on the patient such as facial landmarks and then focused on the selected ROI to estimate vital signs that included heart rate, and respiratory rate. Bousefsaf et al. (2019) monitored the patterns of patients' pulse rates in an ROI of a video frame based on the fluctuations of movement during breathing. Based on the breathing patterns detected in video monitoring, Cho et al. (2017) developed a deep learning model, convolutional neural networks (CNN), which recognizes people's psychological stress levels. Cough analysis was performed based on auscultation sounds by employing a pretrained 3D ResNet18 neural network model to classify the sounds into disease categories. The model achieved 94.57% accuracy, 100% sensitivity, and 94.11% specificity. Heart rate, blood volume pulse, and SpO₂ were measured based on remote Photoplethysmography (rPPG) detected in a video frame captured by a standard smartphone camera (Khalid et al., 2022). The change in blood volume pulse causes blood absorption during a heartbeat was measured by focusing on forehead ROI using the Viola-Jones algorithm. The PPG signal extracted from the video and the ground truth blood pressure from the algorithm was fed into a feed-forward neural network model. The model achieved 85% accuracy in extracting the blood pressure. Laurie et al. further demonstrated how an algorithm specifically designed to control exposure time during video capture improves the accuracy of rPPG (Larue et al., 2021).

Studies that explored the advantages and disadvantages of telehealth are also presented in Table 2. Telehealth cut down travel time, clinic visits, and extended time off work (Nord, Rising, Band, Carr, & Hollander, 2019). However, there are challenges associated with the benefits of telehealth monitoring. Overutilization or misuse of telehealth services has increased healthcare costs to providers (Busso et al., 2022). Telehealth monitoring has widened the disparities between rural and urban populations due to the accessibility of the internet and technology (Drake et al., 2019). Patient data security is another challenge in telehealth monitoring, which jeopardizes patients' health information without an end-to-end encrypted communication service (Fang et al., 2020).

Telehealth patient monitoring techniques have the potential to diagnose patients' health status. AI-enabled telehealth monitoring would be the more enhanced approach to classifying or predicting patients' vital signs.

3.2 | IoT-enabled devices

An IoT based real-time remote patient monitoring system would help achieve continuous patient monitoring (Yew et al., 2020). The majority of IoT technology systems have been developed for use in a hospital setting or a private dwelling. However, there are examples where a single system could be readily applied to both. Figure 3 presents an example of typical architecture that could be used for patient monitoring. The architecture is breakdown into three sections (Pan et al., 2020). Section A illustrates the wearable devices connected to patients to collect vital signs such as heart rate, pulse rate, respiratory rate, breathing rate, body temperature, and so on. In Section B, the collection will be stored in cloud services (Neto et al., 2017; Shao et al., 2020; Shi et al., 2020) for further analysis using machine learning methodologies that could predict or classify the patient data. The process could then estimate any abnormal events in the near future based on known threshold values of the vital signs and update medical staff or healthcare professionals (Ankita et al., 2021; Bekiri et al., 2020; C. Liu et al., 2019; Lin et al., 2018; Devi & Kalaivani, 2019; Efat et al., 2020; Shao et al., 2020) in Section C of the architecture. IoT has the potential to interconnect wearable sensors and their reader-antennas with a patient body to the monitoring network. The types of wearable vital signs sensing technologies, their architectures, and specifications range from physiological measurements, including electrocardiogram, blood oxygen saturation, blood glucose, skin perspiration, and capnography, to motion evaluation and cardiac implantable devices (Dias & Cunha, 2018). The devices can also take the form of wearable t-shirts, chest straps, or adhesive patches. Medical staff or healthcare professionals would take appropriate actions to treat the patient and avoid abnormal events.

TABLE 2 Telehealth monitoring

References	Algorithm	Technology	Advantages	Disadvantages
Bousefsaf et al. (2019)	3D-CNN	RGB camera	• Improved access and timeliness of care	• Overutilization or misuse of telehealth services increase healthcare costs to providers.
Cho et al. (2017)	CNN	Thermal camera	• Emergency preparedness	• Disparities between rural and urban populations.
Khalid et al. (2022)	DFT, CWT	Sensor, RGB camera	• Cost-effectiveness	
Laurie et al. (2021)	Exposure control	Sensor, RGB camera	• Reduced doctor-patient supply-demand mismatch	

Abbreviation: CNN, convolutional neural network.

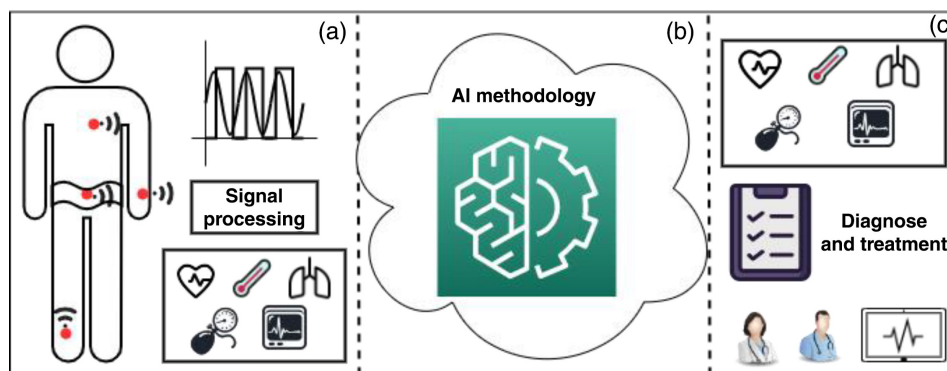


FIGURE 3 Patient monitoring architecture (inspired by Shao et al., 2020)

El-Rashidy et al. (2021) discussed trends and challenges of adopting a wireless body area network (WBAN), a sub-domain of IoT which connects wireless sensors on a patient's body to the network. The WBAN challenges are transmission protocols, data privacy and security, interoperability, and integration. To transmit data from wireless sensors on a human body to local and global networks would need standard data transmission networks like ZigBee, Lora, Wi-Fi, and Bluetooth (Tripathi et al., 2020), and these networks have limitations in terms of energy and range of transmission. Data privacy and security challenges are inevitable in technology-enabled applications, and vast research on secured data transmission or data transaction processes and decentralized or distributed technology is being conducted to overcome the security and privacy challenges. Integration of sensors and remote devices led to sensor interoperability and combining heterogeneous data resources led to data interoperability in wireless sensor networks.

Sensors in RPMs require patients' skin contact to retrieve vital signs, limiting the daily activities of chronically diseased patients. Adopting radio frequency identification (RFID) technology overcomes the challenges of incision in RPM. RFID technology has battery-powered active tags and battery-less passive tags that work on Near-field Coherent Sensing (NCS) principle. To enhance patient comfort and be less restrictive to their daily activities, noninvasive digital technology with NCS has been proposed and developed by Cornell University researchers (Hui & Kan, 2017; Sharma & Kan, 2018). NCS is a method developed by Hui and Kan (2017) that directly modulates the mechanical motion on the surface and inside a body onto multiplexed radio signals. This is integrated with a unique digital identification. In this mechanism, electromagnetic energy is directed into body tissue which reflects back-scattered signals from internal organs and is implicitly amplified. Small mechanical motions inside the body that have a shorter wavelength can be rendered into a large phase variation to improve sensitivity. NCS mechanisms were deployed to monitor vital signs, score sleep (Sharma & Kan, 2018), and accurately extract heartbeat intervals (Hui & Kan, 2018).

Passive RFID tags can be deployed into garments at the chest and wrist areas. This is where the two multiplexed far-field back-scattering waveforms are collected at the reader to retrieve the blood pressure, heart rate, and respiration rate. This could minimize deployment and maintenance costs. Hui and Kan (2017) found that to maximize reading range and immunity to multipath interference caused by indoor occupant motion, active tags could be placed in the

TABLE 3 Internet of Things (IoT) devices monitoring

References	Algorithm	Technology	Advantages	Disadvantages
El-Rashidy et al., 2021		Wireless sensors, ZigBee, Lora, Wi-Fi, Bluetooth	<ul style="list-style-type: none"> • Enable personalized health monitoring 	<ul style="list-style-type: none"> • Patient privacy concerns • High dependence on the internet
Sharma & Kan, 2018	SVM model	NCS, RFID Passive tags	<ul style="list-style-type: none"> • Noninvasive monitoring 	
Hui & Kan, 2018	SVM model	NCS, RFID Passive tags	<ul style="list-style-type: none"> • Continuous health monitoring 	

Abbreviation: SVM, support vector machines.

front pocket. Also, placing the tags in the wrist cuff to measure the antenna reflection due to NCS. With this, vital signals can be sampled and transmitted digitally. Table 3 presents the IoT devices used by the research community for patients' monitoring, and their advantages and disadvantages are explained.

3.3 | Cloud computing

Cloud computing is an essential component of continuous patient monitoring systems. The huge amount of data generated for each patient from the IoT devices in RPMs needs storage to share data between different parties and analyze trends (Zamanifar, 2021). Cloud computing technology is a powerful platform that holds servers, databases, networking, software, and intelligence online (over the internet) for faster innovation and flexible resources. Iranpak et al. (2021) used the features of cloud computing features and developed a patient monitoring system based on IoT devices. The transmission of the data in the IoT platform to the cloud using the Fifth Generation Internet (5G) network. A deep learning neural network model, long short-term memory (LSTM), was used to monitor patients and classify their health conditions. The proposed deep learning model outperformed baseline models with an accuracy of 97.13%. Cloud computing uses a centralized data server to manage large amounts of data from all IoT devices. Integrating the IoT platform with cloud computing raises concerns about latency, real-time response delays, bandwidth overuse, and data security. This led to decentralized or distributed computing approaches like fog computing and edge computing, in which cloud services are brought close to IoT networks and overcome the cloud computing challenges (Pareek et al., 2021). The major concern of data security in cloud computing can be addressed by encrypting the data at the IoT device level and then sending data to cloud storage for data analysis. Siam et al. (2021) proposed the advanced encryption standard (AES) algorithm to encrypt the vital signs retrieved from a patient and send the encrypted data to the cloud. This allows only trusted medical organization servers to access the data with the appropriate decryption key, which will be kept secret between the system and the healthcare center. The proposed approach outperforms the commercial devices available on the market with minimal root-mean-square error, mean absolute error, and mean relative error of 0.012, 0.009, and 0.003, respectively.

3.4 | Fog and edge computing

Fog computing is an extension of cloud computing, which takes cloud computing services closer to IoT devices. Advancements in applications of IoT, and integrated cloud computing such as real-time monitoring of patient vital signs, physical activities have increased threats like security, performance, latency, and network breakdown to cloud computing (Sabireen & Venkataraman, 2021). Fog computing is a distributed or decentralized virtual network to act as a medium between IoT devices and the cloud (Alwakeel, 2021). Pareek et al. (2021) discussed IoT-Fog-based system architectures in healthcare. The delay in real-time responses and latency issues in cloud computing can be addressed by deploying fog nodes that analyze the data from the IoT platform with minimum delay time (Q. Qi & Tao, 2019). The fog computing architecture provides real-time analysis and security of data by preserving sensitive data and performing calculations closer to the IoT platform.

Similarly, cloud computing services are further pushed closer to the edge of the networks or IoT devices by introducing another decentralized or distributed concept called edge computing. The edge computing operations are executed in intelligent devices like programmable controllers, which read IoT devices (Alwakeel, 2021). Edge computing nodes

TABLE 4 Cloud/Fog/Edge monitoring

References	Algorithm	Technology	Advantages	Disadvantages
Uddin (2019)	LSTM model	IoT devices, edge computing	<ul style="list-style-type: none"> Real-time automated analytics 	<ul style="list-style-type: none"> Privacy and security of patient data
Vimal et al. (2021)	CNN	IoT devices, Edge/Fog/Cloud computing	<ul style="list-style-type: none"> Cloud computing services closer to patients 	<ul style="list-style-type: none"> Domain knowledge training for medical staff
Siam et al. (2021)	AES algorithm	IoT devices, cloud computing	<ul style="list-style-type: none"> Decentralized network support personalized monitoring 	

Abbreviations: AES, advanced encryption standard; CNN, convolutional neural network; LSTM, long short-term memory.

deploy intermediate nodes closer to the network with storage and computation capabilities. The cloud/fog/edge monitoring architectures enable real-time monitoring with a decentralized approach for personalized care, but the technologies have their disadvantages, as shown in Table 4.

Uddin (2019) proposed a wearable sensor-based system with an AI-enabled edge device for patients' physical activity prediction. The graphics processing unit in the edge device was used for faster computation results. A deep learning LSTM model was used in the edge for the physical activity classification. The model achieved an accuracy of 99.69% mean prediction performance compared to 92.01% mean recognition performance of traditional approaches like the hidden Markov model and deep belief network. An AI-enabled fog/edge computing approach for fall detection by Vimal et al. (2021) to process binary images of elderly patients in a remote health monitoring setup. The proposed approach has five layers a sensor elder patient body, edge gateway, fog node layer with LoRa connectivity, cloud layer, and application layer for user accessibility. A deep learning convolutional neural networks (CNN) model was used for image processing and compared its performance with support vector machines (SVM) and artificial neural networks (ANN) models. The proposed deep learning model achieved an accuracy of 98% with a minimal processing time of <200 s but with a higher power consumption of >65 decibels.

All the research works discussed in this section have been summarized with their application, algorithms, and technologies used in the RPM system, as shown in Table 4.

3.5 | Blockchain monitoring

Virtual technologies like fog computing and edge computing are prone to security and privacy challenges (Aliyu et al., 2021). Hathaliya et al. (2019) proposed a Permissioned blockchain-based healthcare architecture to overcome these challenges. The study focused on integrating decentralized AI with blockchain networks and discussed blockchain applications in healthcare. Blockchain is a shared, decentralized, immutable ledger that connects multiple parties and records transactions. In RPM applications, blockchain technology can secure data transactions between patients and monitoring technologies like cloud, fog, and edge computing. Faruk et al. (2021) proposed an Ethereum-based data repository for RPM electronic health records data management. The data repository enabled secure upload, storage, analysis, retrieval, and transmit patient data according to the patient's instructions. The proposed decentralized blockchain system supports hospitalized patients and outpatients. The cloud computing challenge of interoperability can be addressed using MedHypChain proposed by Kumar and Chand (2021). MedHypChain is a privacy-preserving medical data-sharing system based on Hyperledger Fabric, in which each data transaction is secured via an Identity-based broadcast group encryption scheme. Another interesting patient-centric secured data recording and remote patient monitoring application SynCare was proposed by Pighini et al. (2022). The study focused on interconnecting patients, healthcare professionals, and caregivers, building secure data-sharing channels, and allowing patients to manage their health data. Blockchain architectures are known for their robust security features that record each transaction throughout the system and cannot be altered. The architecture has the disadvantage of high implementations with complex integration and high energy dependence. The research works adapted to blockchain technology in RPM have been outlined in Table 5.

The technology-enabled RPMs are more concentrated on data acquisition and securing the data transmission to different parties involved in RPM. Adopting AI to the RPM architectures empowers the monitoring process with

TABLE 5 Blockchain architectures

References	Algorithm	Technology	Advantages	Disadvantages
Hathaliya et al. (2019)	AI methods	Blockchain	<ul style="list-style-type: none">• Network security at all levels of data collection• Verification and identification of patients• Authorize patients' EHR data	<ul style="list-style-type: none">• High energy dependence• Integration complexity• High implementation costs
Faruk et al., 2021		Blockchain, Ethereum		
Kumar and Chand (2021)		Blockchain, Hyperledger Fabric		

Abbreviation: AI, artificial intelligence.

capabilities of prediction and classification of the patient data acquired. Each RPM architecture can be enhanced by adding AI modeling to the data analytics step.

4 | AI IN RPM APPLICATIONS

In RPM applications, traditional machine learning and deep learning are common AI methods adopted to detect and predict vital signs and classify patients' physical activities. Malasinghe et al. (2017) present contact and noncontact-based methodologies in RPM. Irrespective of contact or noncontact monitoring systems, all methodologies focus on human vital signs extraction, such as heart rate, pulse rate, respiration rate, blood pressure, and oxygen volume in blood, as the deterioration of these vital signs affects the human health system. Along with vital human signs, the authors reviewed studies on the activity detection of patients like fall detection and mobility-related diseases. The challenges involved mainly include discerning the difference between deliberate quick movements and accidental drops. Apart from wearable devices, the authors reviewed ambient device-based and vision-based fall detection systems but identified significant problems that remain for contactless monitoring. This section discusses applications of machine learning and deep learning methodologies in RPM. The year-wise distribution of the AI-enabled RPM works discussed in this section are presented in Figure 4.

4.1 | Vital signs monitoring

Wearable devices like smartwatches are new technological innovations that continuously track people's vital signs. A system was developed by Bekiri et al. (2020) to monitor the health status of individuals at all times using connected smartwatches. The smartwatches collect patient vital signs and send them to the administrator to analyze for decision-making. The administrator used the SVM model to build a decision model. The results of the patient's status will be informed to the doctors. The machine learning model achieved an accuracy of 90% and a recall is 99%. The proposed system can identify 99% of patients affected by cardiovascular diseases. Shao et al. (2020) also designed an RPM system to detect ECG signals. In that study, a decision tree ensemble classifier was trained using the CatBoost learning kit. The classifier was trained with 20-fold cross-validation and 31 features. Feature importance was extracted from the trained CatBoost model. The top-importance features were used to evaluate the performance based on the feature importance ranking. The CatBoost model processed the 30 s ECG data in 0.5 s and achieved a sensitivity of 99.61%, a specificity of 99.64%, and an accuracy of 99.62% in detecting AF. A novel IoT-based wearable 12-lead ECG SmartVest system based on the SVM model to assess signal quality has achieved an average accuracy of 97.9% and 96.4% for acceptable and unacceptable ECG segments, respectively. Verified the model efficiency to choose good or exclude poor quality ECG segments in the wearable.

ECG monitoring. An SVM model-based ECG telemetry system to monitor cardiac arrhythmia, which processes the ECG signal, was designed by Devi and Kalaivani (2019) to send alerts to a physician in an emergency. Statistical features of ECG signals were combined with dynamic features like heart rate variability (HRV) features from RR intervals to classify cardiac arrhythmia. The SVM classifier model was trained and validated using 10-fold cross-validation. The proposed classification model achieved the effectiveness of 88.9%, 90.8%, and 92.2% for statistical features, HRV features, and statistical and HRV features, respectively.

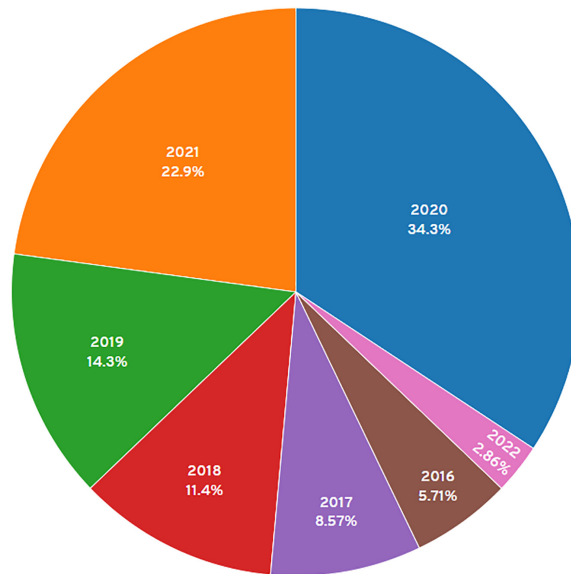


FIGURE 4 Distribution of artificial intelligence-enabled remote patient monitoring applications in this section

Neto et al. (2017) designed an RPM system with a portable ECG device to assist remote electrocardiographic diagnosis and send the data to cloud service, where an intelligent arrhythmia detector (IDAH-ECG) detected abnormal heartbeats and informed physicians. Here, discrete wavelet transforms feature extraction, and principal component analysis (PCA) dimensionality reduction was performed as part of data preprocessing. Multilayer perceptron (MLP) neural network (MLP) classifier was trained using backpropagation and Gradient Descent techniques. This study implemented 10-fold cross-validation with the Monte Carlo testing scheme. The classifier achieved an average accuracy of 96.48%, a sensitivity of 98.70%, and a specificity of 94.45%. Elola et al. (2019) used deep learning methods for sensor-based pulse recognition from short electrocardiogram (ECG) segments that included a deep convolutional neural network (DCNN), auto-encoder, restricted Boltzmann machine (RBM), and recurrent neural network (RNN). The system was designed to detect a pulse in someone who had a heart attack during resuscitation efforts. Although the deep neural network (DNN) architectures outperformed current methods, pulse detection during this scenario remains an unsolved problem. J. Yang et al. (2020) developed a wireless nonline-of-sight (NLOS) bio-radar device that was used to collect physiological parameters such as heart rate and respiratory rate. The device is portable, contactless, and interference-free. A deep learning LSTM was employed at edge nodes to predict the physiological parameters. The authors used LSTM to predict future short-term respiratory rates and patients' heartbeats based on current data within just a few minutes. W. Qi and Aliverti (2020) proposed a wearable respiratory and activity monitoring system to predict breathing patterns during daily activities based on a novel multimodal fusion architecture, respiratory and exercise parameters and human activity. A hybrid hierarchical classification algorithm combining an LSTM model with a threshold-based approach to classify nine breathing patterns while performing 15 physical activities. The hybrid model achieved an accuracy of 97.22% and outperformed the other models' K-nearest neighbor (KNN), multiclass SVM, and artificial neural network (ANN) in terms of classification. The proposed model outperformed LSTM, bidirectional LSTM (Bi-LSTM), and DCNN with a minimal computational time of 0.0094 s. The research works related to vital signs monitoring with AI discussed are consolidated in Table 6.

4.2 | Physical activities monitoring

Pan et al. (2020) designed a fall detection system for older people based on multisensor fusion with multiple three-axis acceleration sensors placed on the waist. In this study, SVM and random forest (RF) algorithms were implemented on the dataset with 100 healthy young volunteers simulating falls and daily activities to compare their recognition time and recognition rate. The authors state that the model's accuracy is based on a large amount

TABLE 6 Vital signs monitoring with artificial intelligence

Applications	Algorithm	Technology	References
Vital signs monitoring	SVM Model	Smartwatches, smart vest, ECG telemetry	Bekiri et al. (2020), Shao et al. (2020), Devi and Kalaivani (2019)
	CNN, LSTM, DCNN, DNN, RNN, ANN auto-encoder	Potable ECG device, sensors, radar device	Neto et al. (2017), Elola et al. (2019), J. Yang et al. (2020), W. Qi and Aliverti (2020)

of valid data, but because SVM has fewer training and recognition times, it may be better suited to this task. Hsieh et al. (2021) proposed a novel multiphase falls identification algorithm combining fragment modification algorithm and machine learning techniques to identify prefall, free-fall, impact, resting, and recovery phases. The fragment modification algorithm adopts rule-based fall identification and five machine learning techniques, SVM, KNN, naive Bayes, decision tree, and adaptive boosting to identify the five phases. Out of the five models, the KNN algorithm achieved the best performance with an accuracy of 90.28%, sensitivity of 82.17%, precision of 85.74%, and Jaccard coefficient of 73.51%. The authors intend to further develop their model with real-world data and a greater range and type of falls. Y. Wang and Zheng (2018) designed a framework for an RPM system to monitor human activities and movement based on a signal reflection model. This framework detected the presence of human activities by analyzing the RSSI patterns from an RFID tag array and segmented phase values using the variance of phase readings, which were used as an indicator for activity segmentation. Six machine learning classifiers RF, multilayer Perceptron-based Neural Network, Decision Tree, SVM, Naive Bayes, and Quadratic Discriminant Analysis, were trained to classify activities raise a hand, drop hand, walk, sit, stand, fall, rotation, get-up, and non-activity. The experiment results show that TACT is robust under different experimental settings and can achieve an average recognition precision of up to 93.5%.

Salah et al. (2022) designed a resource-constrained microcontroller at the edge of the network using a wearable accelerometer to overcome issues such as latency, high power consumption and poor performance in areas with unstable internet. The authors designed three layers edge layer, fog layer, and cloud layer to collect, analyze, and transmit to an IoT gateway via long-range communication technology. Five AI models, KNN, SVM, LSTM, and CNN, were trained to detect falls. The LSTM identified falls from daily activities with high accuracy of 96.78%, while sensitivity and specificity were 97.87%, and 95.21%, respectively. S. Yu et al. (2021) proposed a computational method with a Hierarchical Attention-based Convolutional Neural Network (HACNN) model to detect falls based on wearable sensor data. The novel deep learning model integrated a hierarchical attention mechanism into a CNN model and added two attention layers beyond CNN to interpret which part of the sensor data contributed to the decision of fall or nonfall made by the system. The CNN model outperformed deep learning models like CNN, LSTM, CNN-LSTM, MLP, and HALSTM. Accuracy depended on the two data sets used and their static nature.

To overcome the limitations of existing elderly fall detection methods requiring specialized hardware or invading people's daily lives, Zhu et al. (2017) presented the design and implementation of a motion detection system based on passive radio frequency identification tags. The received signal strength indicator (RSSI) value and Doppler frequency value impacted by static, regular action, sudden falls, elderly movements, and fall actions were estimated. Wavelet transform was implemented for the signal preprocessing, and the machine learning algorithm SVM was adopted to classify the actions into fall detection or other actions. RFID technology could track their motion and fall detection without any hindrance to the daily activities of elderly people. Gesture recognition or motion detection has gained attention to enhance the user experience for human-computer interaction. An application can be used in healthcare to recognize patient gestures or motions in-home or in the hospital using a device-free system. Z. Wang et al. (2019) proposed RF-finger, a device-free system based on Commercial-Off-The-Shelf (COTS) RFID, which leverages a tag array on a letter-size paper to sense the fine-grained finger movements performed in front of the paper presented. Machine learning algorithms were implemented, such as the KNN model to pinpoint the finger position and the CNN model to identify the multitouch gestures based on reflective images. Both the machine learning algorithms yielded 88% and 92% accuracy for finger tracking and multitouch gesture recognition, respectively. Estimate the correlation between RF phase values and human activities by modeling intrinsic characteristics of signal reflection in contact-free scenarios. The research works related to human activity recognition and fall detection are presented in Table 7.

TABLE 7 Human activity recognition with artificial intelligence

Applications	Algorithm	Technology	References
Physical Activities Monitoring	SVM, RF, KNN, naive Bayes, decision tree, adaptive boosting	Sensors, RFID Tags	Pan et al. (2020), Hsieh et al. (2021), Y. Wang and Zheng (2018), Zhu et al. (2017); Z. Wang et al. (2019)
	KNN, SVM, LSTM, CNN, HACNN	Wearable accelerometer, cloud, fog, edge	Salah et al. (2022), S. Yu et al. (2021)

4.3 | Chronic disease monitoring

4.3.1 | Diabetes monitoring

Mujumdar and Vaidehi (2019) proposed a diabetes prediction model to classify diabetes which includes external factors responsible other than regular factors like glucose, body mass index (BMI), age, insulin, and so on. In this study, machine learning algorithms, including Support Vector Classifier, Decision Tree classifier, Extra Tree Classifier, Ada Boost algorithm, Neural Networks, RF Classifier, Linear Discriminant Analysis algorithm, Logistic Regression, K-Nearest Neighbor, Gaussian Naïve Bayes, Bagging algorithm, and Gradient Boost Classifier was implemented for diabetes prediction. All the models were trained and evaluated using a confusion matrix and classification report. Out of all, the logistic regression model was able to classify diabetic and nondiabetic with an accuracy of 96%. To continuously evaluate diabetic patients' data such as sugar level, sleep time, heart pulse, food intake, and exercise collected through sensors. Analyze the data using neural networks and classify the health risk status into four modes: low, medium, high, and extreme (Efat et al., 2020). Based on national physical examination, the risk factors for type II diabetes mellitus (T2DM) were computed using machine learning algorithms. A logistic regression model was implemented on physical measurement and a questionnaire. The 14 risk factors selected in logistics regression were combined and implemented with tree-based machine learning algorithms like decision tree, RF, AdaBoost, and XGBoost. Out of 4 algorithms, XGBoost had achieved an accuracy of 90.6%, precision 91.0%, recall 90.2%, F1 score 90.6% and AUC 96.8%. XGBoost model was used to output feature importance scores. BMI was the most important feature, followed by age, waist circumference, systolic pressure, ethnicity, smoking amount, fatty liver, hypertension, physical activity, drinking status, dietary ratio (meat to vegetables), drink amount, smoking status, and diet habit (oil-loving) (Xue et al., 2020). Several ML techniques could classify different states of diabetic patients with high accuracy.

4.3.2 | Mental health monitoring

Mental health illness is one of the most underestimated human states, which shortens the life span by 10–20 years (McGinty et al., 2021). It would be difficult to manually monitor people with mental health illnesses such as schizophrenia, bipolar disorder, major depressive disorder, and suicidal tendency. Thieme et al. (2020) conducted a systematic review on how implementing machine learning can assist in detecting, diagnosing, and treating mental health problems. Machine learning techniques can offer new routes for learning patterns of human behavior. It helps in identifying mental health symptoms and risk factors. Also, it assists in predicting disease progression and personalizing and optimizing therapies. Advances in machine learning will attempt to predict suicide based on the analysis of relevant data and inform clinical practice. Adamou et al. (2019) proposed a text-mining approach to support risk assessment. Latent Dirichlet allocation (LDA) is a special case of topic modeling for natural language processing. The technology was used to process different types of information like demographics, appointments, progress notes, comprehensive assessments, referrals, and Inpatient stays. Statistically equivalent signatures (SES) mechanism for feature selection. In this study, support vector regression (SVR), RF, and linear ridge regression (RLR) models were implemented along with K-fold cross-validation. The approach achieved the highest area under curve (AUC) value of 0.705. Continuous monitoring would enable the record of abnormal vital sign measurements, and ML techniques have the potential to analyze the data to detect underlying data patterns to take appropriate treatment steps. Diabetes and Mental Health monitoring discussed in this section are outlined in Table 8.

TABLE 8 Chronic disease monitoring with artificial intelligence

Applications	Algorithm	Technology	References
Diabetes monitoring	Decision Tree, RF, AdaBoost, XGBoost, RF, KNN, SVM	EHRs, Sensors	Mujumdar and Vaidehi (2019), Efat et al. (2020), Xue et al. (2020)
Mental health monitoring	LDA, SVM, RF, RLR, SES	Social media text, demographics	Thieme et al. (2020), Adamou et al. (2019), McGinty et al. (2021)

4.4 | Emergency monitoring

4.4.1 | Emergency department

RPMs Decision-making for emergency department patients using machine learning techniques would have helped to improve existing methods. Taylor et al. (2016) compared clinical decision rule to a machine learning approach for predicting in-hospital mortality of patients with sepsis. In this study, machine learning techniques were used to extract a large number of variables through existing emergency department clinical records to predict patient outcomes and facilitate automation and deployment within clinical decision support systems. Patients visiting the emergency department visits were randomly partitioned into an 80%/20% split for training and validation. 500 clinical variables were extracted from the real-time clinical records of four hospitals using an RF model to predict in-hospital mortality. Later, the RF model was compared to the classification and regression tree (CART) and logistic regression models. The random forest model AUC was statistically different from all other models ($p \leq 0.003$ for all comparisons). Kong et al. (2016) designed a decision tool based on rule-based inference methodology using the evidential reasoning approach (RIMER) that was developed and validated to predict trauma outcomes. It helps physicians to predict in-hospital death and intensive care unit (ICU) admission among trauma patients in emergency departments. The prediction performance of the RIMER was compared to logistic regression analysis, support vector machine, artificial neural network, SVM models, and ANN models. Five-fold cross-validation was implemented, and the AUCs of RIMER, logistic regression, SVM, and ANN are 0.952, 0.885, 0.821, and 0.790, respectively. The results show that the RIMER model performs the best. The machine learning techniques could classify near-term mortality based on vital signs analysis in emergency department patients.

The machine learning approach is able to incorporate heart rate variability (HRV) for intensive monitoring, resuscitation facilities, and early intervention for critically ill patients in the emergency department by comparing the area under the curve, sensitivity, and specificity with the modified early warning score (MEWS). In a study (Oh et al., 2018), HRV parameters were generated from a 5-min electrocardiogram (ECG) recording incorporated with age and vital signs to generate the ML score for each patient. The area under the receiver operating characteristic curve (AUROC) for ML scores in predicting cardiac arrest within 72 h is 0.781, compared with 0.680 for MEWS. For in-hospital deaths, the area under the curve for ML score is 0.741, compared with 0.693 for MEWS. A cut-off machine learning score ≥ 60 predicted cardiac arrests with a sensitivity of 84.1%, specificity of 72.3%, and negative predictive value of 98.8%. A cut-off MEWS ≥ 3 predicted cardiac arrest with a sensitivity of 74.4%, a specificity of 54.2%, and a negative predictive value of 97.8% (Blasiak et al., 2020; Ong et al., 2012). Based on the results, machine learning scores were more accurate than the traditional MEWS in predicting cardiac arrest within 72 h.

4.4.2 | RPMs in the ICU

To predict near-term mortality in patients hospitalized with cirrhosis (Antunes et al., 2017), two machine learning approaches (i) logistic regression and (ii) LSTM neural network on medical record entries of 500 patients staying ICU and compared them (Xia et al., 2019). In total, 20 features, such as pulse, respiratory rate, systolic, and diastolic blood pressure, were used for training the algorithm. The machine learning models outperformed the clinical decision tool, a mathematical Chronic Liver Failure (CLIF) Score. A logistic regression model achieved an AUC of 0.80, the RNN-LSTM model achieved an AUC of 0.77, and CLIF achieved an AUC of 0.72 (Harrison et al., 2018). A patient-specific model could analyze vital signs based on historical data. Colopy et al. (2018) proposed Gaussian process regression (GPR) to provide flexible, personalized models of time series of patient vital signs. This study uses a method to build GP models

TABLE 9 Emergency and intensive care unit patients monitoring with artificial intelligence

Applications	Algorithm	Technology	References
Emergency monitoring	RF, RIMER, LSTM, GPR, Logistics Regression, CLIF	EHRs, ECG, RIMER	Taylor et al. (2016), Kong et al. (2016), Oh et al. (2018), Ong et al. (2012), Blasiak et al. (2020), Antunes et al. (2017), Xia et al. (2019), (Harrison et al. (2018), Colopy et al. (2018)

TABLE 10 Facial and emotions recognition with artificial intelligence

Applications	Algorithm	Technology	References
Facial and emotions monitoring	LSTM, SVM	Image Processing Thermal sensors IoT devices RFID signals	Maresh et al. (2021), Chowdary et al. (2021), Zainuddin et al. (2020) Q. Xu et al. (2020)

with varying complexity and regularization using different hyperparameters on a patient-specific level to forecast robust, vital signs. The authors used a random search algorithm to search for patient-specific parameters. Bayesian optimization methods were implemented to accommodate any plausible parameterization in the patient population. Patient-specific parameter optimization using machine learning techniques is the most advanced level of RPM. This helps to build patient-specific models and break down a patient's health status to the lowest level. The Table 9 consolidates the research works related to emergency and ICU patients monitoring.

4.5 | Facial and emotions recognition

AI can classify the patient's emotions based on patient face recognition. A smart integrated patient monitoring system was proposed by Maresh et al. (2021) to detect patients' emotional states and heartbeat levels through face recognition algorithms, heartbeat, and temperature sensors. Their RPM system presented the emotional data of the patients using face recognition algorithms such as image preprocessing, feature extraction, and classification. The facial emotional recognition model is able to identify seven emotions: Anger, Happy, Sad, Neutral, Surprise, Disgust, and Fear. Based on the heart rate sensor and thermal sensor, patients' vital signs were measured with an interval of 5 s. Zainuddin et al. (2020) used IoT technology to communicate facial emotions and vital signs to hospitals. Similarly, Chowdary et al. (2021) designed an RPM system based on deep learning-based facial emotion recognition to overcome problems associated with mutual optimization of feature extraction and classification.

An experimental study to recognize user emotions of users with commercial RFID devices. Q. Xu et al. (2020) designed an RPM system using an emotion recognition framework that first extracts respiration-based features and heartbeat-based features from RFID signals. The extracted features were used in training a classifier that a user's different emotions. In this study, the respiration rate was separated by using filters, whereas the heartbeat signal was retrieved by suppressing the respiratory signal and improving the signal-to-noise ratio. Using their framework, a 2D emotional model divided emotions into four states: joy, pleasure, anger, and sadness. An SVM model was able to classify the four emotion states with an accuracy of 80.65%, 61.29%, 83.87%, and 74.19%, respectively (Q. Xu et al., 2020). The research community developed RPM systems using AI for facial and emotions recognition with technologies, as shown in Table 10.

5 | AI IMPACT ON RPM

5.1 | Early detection of patient deterioration

Early detection of vital signs deterioration is key to the timely invention and avoiding clinical deterioration in acutely ill patients in hospitals. Traditional patient monitoring is to report individual vital signs of patients, which state their

current clinical status. For example, vital signs such as temperature, pulse, respiratory rate, and mean arterial pressure (MAP) are considered continuous predictors for emergency department patients (Asiimwe et al., 2020). New patient monitoring algorithms analyze multiple features from physiological signals. This produces a predictive or prognostic index that measures a specific critical health event or physiological instability (Helman et al., 2022). Posthuma et al. (2020) presented a case series where wireless remote vital signs monitoring systems on surgical wards could reduce the time to detect deteriorating patients. As part of this study, nursing staff found the systems somewhat useful, but still required clinical judgment to assess the patient. They noted that there are still no set standards or guidelines for implementing these types of systems, and the task remains for clinicians to judge which system best meets their needs.

Kellett and Sebat (2017) further elaborated on the need for clinicians to place more importance on regular and accurate recording of vital signs. The authors noted that there is currently no agreement on how often vital signs need to be recorded and that most hospital wards use periodic, manual observation of vital signs. Kellett also highlighted the need for continuous patient monitoring and emphasized how vital this is to predict the onset of abnormal events.

Current approaches by clinicians for early prediction of patient deterioration can be estimated using manually calculated screening metrics called early warning scores (EWS) (Garca-del Valle et al., 2021; Vinegar & Kwong, 2021). Downey et al. (2018) demonstrated that although EWS systems have excellent predictive values, they are limited by their intermittent nature. Until recently, continuous vital signs monitoring was limited to intensive care units that require high staff-to-patient ratios. For example, Alshwaheen et al. (2021) proposed a novel framework of patient deterioration prediction in ICUs based on LSTM-RNNs. The model acquired a significantly better classification performance than the traditional method and could predict deterioration 1 h before onset. Muralitharan et al. (2020) further showed that machine learning based EWS could be applied to a range of acute general medical and surgical wards, including ambulatory and home care settings, and still perform better and with greater accuracy than the traditional manual methods.

Although many studies have focused on the prediction of health outcomes, da Silva et al. (2021) predicted future deteriorating vital signs based on applying RNNs and LSTM to historical data from electronic medical records (EMR). These predicted vital signs were then applied to a clinical prognostic tool that used a combination of laboratory results with vital signs for early diagnosis of worsening health status, with an accuracy of 80%.

Transparency and explainability are essential elements for AI models if they are going to be acceptable to clinicians. Lauritsen et al. (2020) proposed an explainable AI EWS (xAI-EWS) system for the early detection of acute critical illness. The xAI-EWS was composed of a temporal convolutional network (TCN) prediction module and a deep Taylor decomposition (DTD) explanation module tailored to temporal explanations. Clinical experts evaluated the system based on three emergency medicine cases: sepsis, acute kidney injury (AKI), and acute lung injury (ALI). The system facilitated trust in the predictive capability by giving clinicians insights into the internal mechanics of the model without any deep technical knowledge of the mechanisms behind it.

5.2 | Personalized monitoring

Conventional diagnoses of diseases and treatments from doctors are based on population averages and do not consider the individual variability of patients to treatments (G. Chen, Xiao, et al., 2021). The IoT-enabled RPM architecture with cloud computing discussed in previous sections combines patients' data for AI modeling. In contemporary settings, patient-centric or personalized monitoring is critical, particularly for chronic diseases like mental health disorders, diabetes, and so on. Personalized monitoring can be carried out with distributed networks like fog and edge computing, where an edge network is set up for a set of IoT devices on a patient. Mukherjee et al. (2020) proposed an edge-fog-cloud framework for personalized health monitoring to predict patient mobility and advice nearby healthcare centers in case of emergency. However, the data acquired from an IoT platform has to leave the devices and be merged into a centralized cloud server for data analytics. This raises privacy and security concerns about patients' health data. Moreover, it demands huge technological resources and power consumption.

A Federated learning framework in AI methods developed by Google could overcome the challenges by training an AI model across multiple decentralized edge devices with local data available at each patient without exchanging or merging them. The local model weights are aggregated and passed to a cloud server. The aggregated model weights are used to train a robust global AI model. In this decentralized framework, the patient data will not leave their device and ensure data privacy. The robust global model can be passed to local models for better prediction or classification results on local data. The research community has widely adopted the approach for IoT applications. Zheng et al. (2021)

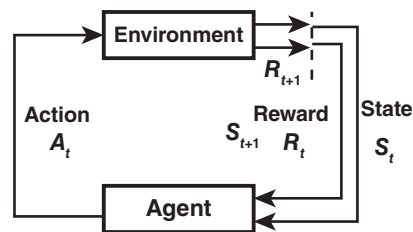


FIGURE 5 Generic reinforcement learning mechanism

proposed a federated transfer learning mechanism for the internet of medical things (IoMT) healthcare. Nguyen et al. (2023) discussed types of federated learning frameworks for smart healthcare, benefits, requirements, federated learning applications in applications, trends, and challenges. Wu et al. (2022) proposed the FedHome framework, a novel cloud-edge-based federated learning framework for in-home health monitoring to training local models. In that study, a generative convolutional autoencoder (GCAE) was designed to process imbalanced and nonidentical distribution data and to achieve accurate results in personalized health monitoring. The proposed approach outperforms baseline models with an accuracy of 95.87% and 95.41% for balanced data and imbalanced data, respectively. In personalized monitoring, physical activity classification task was performed with FedHealth by Y. Chen et al. (2020). The federated learning framework proposed in the study is based on data aggregation and builds personalized models with transfer learning. FedHealth framework was evaluated with two classification problems. One is to classify physical activities and outperform the baseline models with an accuracy of 99.4%. The second one is to classify Parkinson's disease patient's arms droop and postural tremor test and achieved an average accuracy of 84.3% and 74.9%, respectively. Similarly, Shaik, Tao, Higgins, Gururajan, et al. (2022); Shaik, Tao, Higgins, Xie, et al. (2022) proposed a heterogeneous FedStack framework to support diverse architectural local models at patients-end to classify their physical activities and build a robust global model based on predictions of the local models.

5.3 | Adaptive learning

Reinforcement learning, a subset of AI, possess the ability to make a sequence of decision with its reward-driven behavior. The machine learning approach learns to achieve a goal in a potentially uncertain, complex environment. It can employ trial and error to solve a problem and get either rewards or penalties for the steps it executes (C. Yu et al., 2023). In the reinforcement learning approach, a learning agent is deployed in an environment without any prior information or knowledge. The agent has to learn the patterns based on their experience. To transit from the current state (S_t) at time t to the next state (S_{t+1}) at time step $t + 1$, an action (A_t) is taken as shown in Figure 5. For these actions, a predefined reward policy is designed. If the actions chosen is following the policy, the agent gets rewarded (R_t), otherwise penalized. With the sequential decision-making capability, different reinforcement learning schemes are applied to diverse dynamic treatment regimes (Laber et al., 2014) like chronic diseases, mental health diseases, and infectious diseases which need a sequence of decision rules to determine a course of action to suggest treatment type, drug dosage, or re-examination timing. C. Yu et al. (2023) surveyed the applications of reinforcement learning in healthcare. The study has covered treatment strategies built on reinforcement to treat chronic diseases, cancer, diabetes, anemia, HIV, and several common mental illnesses. I. Y. Chen, Joshi, et al. (2021) considered a clinician (learning agent) who monitors the patient (environment) via actions like ventilation and observing the changes in the patient's state (environment) to achieve a goal to discharge the patient successfully. This study provides a practical understanding of the reinforcement learning approach in healthcare. Watts et al. (2020) developed a model to prescribe the timing and dosage of medications using wearable sensors in real-time and deep reinforcement learning. Similarly, Naeem et al. (2021) proposed an intelligent system that relies on algorithms of both Reinforcement Learning and Deep Learning to maximize the successful completion of the patient taking the right pill. Just-in-Time Adaptive Interventions (JITAI) are another healthcare applications which needs timely intervention to provide the right amount of support to patients at right time. This can be achieved by adaptive learning of dynamic health changes in a patient (Nahum-Shani et al., 2017). Wang, Zhang, et al. (2021) adopted reinforcement learning in a data-driven approach for mobile healthcare user and optimize intervention strategies in their context. Similarly, Gönül et al. (2021) proposed a reinforcement

TABLE 11 Artificial intelligence (AI) impact on remote patient monitoring systems

AI impact	Algorithms/technology	Applications	References
Early detection of patient deterioration	explainable AI EWS, LSTM, TCN, DTD	<ul style="list-style-type: none">• Continuous monitoring of emergency patients, sepsis, acute kidney injury, acute lung injury• Early diagnosis with Early Warning Scores and Predictive Prognostic Index	Asiimwe et al. (2020), Helman et al. (2022), Posthuma et al. (2020), Kellett and Sebat (2017), Garca-del Valle et al. (2021), Vinegar and Kwong (2021), Downey et al. (2018), Alshwaheen et al. (2021), da Silva et al. (2021), Lauritsen et al. (2020)
Personalized monitoring	Fog, Edge, IoMT, Cloud GCAE, FedHealth, FedStack	<ul style="list-style-type: none">• Enable personalized monitoring with decentralized learning• Overcome data privacy issues	G. Chen, Xiao, et al. (2021), Mukherjee et al. (2020), Zheng et al. (2021), Nguyen et al. (2023), Wu et al. (2022), Y. Chen et al. (2020), Shaik, Tao, Higgins, Gururajan, et al. (2022); Shaik, Tao, Higgins, Xie, et al. (2022)
Adaptive learning	DRL, A2C, DQN	<ul style="list-style-type: none">• Learn patient behavior patterns• Dynamic treatment regimes• Just-in-time-adaptive-interventions• Sequential decision making tasks	C. Yu et al. (2023), Laber et al. (2014), I. Y. Chen, Joshi, et al. (2021), Watts et al. (2020), Naeem et al. (2021), Nahum-Shani et al. (2017), Wang, Zhang, et al. (2021), Gönül et al. (2021)

learning mechanism to personalize digital adaptive interventions as mobile notifications to the user in coping their health problems. The authors deployed two models, intervention selection and opportune moment identification. With respect to type and frequency, the intervention selection model adopts the intervention delivery. The opportune moment identification is to detect the most opportune moments to intervene. Table 11 presents research works of AI which can transform healthcare applications with advanced mechanism such as reinforcement learning and federated learning.

6 | CHALLENGES AND TRENDS OF AI IN RPM

Implementing a technology-enabled patient monitoring system would require hospital staff support and their views. Ede et al. (2021) did a qualitative study to explore staff expectations of wireless noncontact patient vital signs monitoring, their perception of the utilization of the technology in the ICU, patients, and relative response to introducing the technology. Nine nurses with a median duration experience of 2 years in ICU were interviewed on five different themes such as ICU staff perceptions of the patient and relative monitoring experiences relating to current wired monitoring and expectations of noncontact monitoring, staff expectation of continuous monitoring in ICU, troubleshooting, the hierarchy of monitoring and consensus of trust. Although AI can transform healthcare with its potential to analyze, predict and classify data efficiently, there remains a hesitancy to adopt the technology (Meskó et al., 2017). This section discusses challenges in adopting AI to remote monitoring systems for vital signs precision and activity recognition. Initiatives to overcome the challenges are also presented.

6.1 | AI or ML explainability

The first and foremost challenge is the difficulty associated with interpreting the results generated by an AI or ML model. Current models are better than humans at interpreting complex data and predicting outcomes but lack the capacity to demonstrate how these conclusions were reached or if there were any weaknesses in the algorithm applied by the ML model. This is one of the most challenging barriers for healthcare professionals to adopt AI or machine learning methodologies (Mohanty & Mishra, 2022). Most of the machine learning models, such as neural networks, SVM, and so on are black-box models. These models cannot elaborate their results and provide cause-and-effect relationships between predictor variables and target variables (G. Yang et al., 2022). AI or ML can be adaptable only when interpretable structures and results for healthcare professionals (Sagi & Rokach, 2020).

Sensitivity is one of the promising methods which can explain the cause-and-effect relationship between the input and output variables of a trained neural network. Tree-based methods will not expect a parametric relationship between input and output variables. Both classification and regression trees have been shown to provide interpretable structures that would support clinical practice in decision-making (Jovanovic et al., 2016; Sagi & Rokach, 2020). Other than tree-based methods, pattern-based classification, Naive Bayes, knowledge-based algorithms, logistic regression, and fuzzy models can produce interpretable structures (Shouval et al., 2017).

SHapley Additive exPlanations (SHAP) (Lundberg & Lee, 2017), a method based on cooperative game theory (Shapley, 1953), can increase the transparency and explainability of AI Models. In this method, the impact or contribution of input to the prediction output is represented with Shapley values and the values are calculated for each input feature. The Shapley values of the prediction model were extracted in two forms: a global perspective of factors that required special attention in overall prediction, and a local perspective of each feature in a single prediction. Linardatos et al. (2020) reviewed machine learning interpretability methods. The study reported different scopes of interpretability for deep learning models. It includes gradients explanation technique, integrated gradients, gradient-weighted class activation mapping, DeepLIFT algorithm, deconvolution, and guided back-propagation. Raza et al. (2022) proposed a framework for accurate and efficient personal healthcare using federated transfer learning and an explainable AI (ExAI) model in EEG signal classification. Khodabandehloo et al. (2021) proposed a flexible AI system HealthXAI to predict early symptoms of early decline in smart homes. The anomaly level of behavior was computed based on the anomaly feature vector. The authors built a dashboard to allow clinicians to inspect anomalies, scores, and their automatically generated natural language explanations. Trends in Explainable AI are sensitivity and Shapley values and the research works exploring these trends in Table 12.

6.2 | Privacy

Considering the black-box nature of deep neural networks, it is impossible to predict what neural networks learn from data. The problem with this is that they might unintentionally learn features that discriminate against user information. This increases the risk of information disclosure. Iwasawa et al. (2017) analyzed the features learned by conventional deep neural networks when applied to data of wearable to confirm this phenomenon. A simple logistic regressor could achieve a high user classification accuracy of 84.7% when using the CNN features extracted from basic activity signals. The same classifier could only obtain 35.2% user classification accuracy on raw sensor data. This reveals the privacy leakage potentials of a deep learning model originally used for human activity recognition (K. Chen, Zhang, et al., 2021).

In this study (K. Chen, Zhang, et al., 2021), data transformation and data perturbation techniques, were suggested that could be used to overcome privacy issues with machine learning algorithms. User adversarial neural networks were proposed to integrate an adversarial loss with the standard activity classification loss to minimize the user identification accuracy. However, the adversarial loss technique has a limitation to protecting only private information, such as user identity and gender. To protect all sensitive user identity information, the raw sensor signals were viewed from two perspectives, style and content. D. Zhang et al. (2019) proposed to transform raw sensor data to have the “content” unchanged, but the “style” is similar to random noises. For data perturbation, a deep private auto-encoder (dPA) was proposed by Gati et al. (2021) to perturb the objective functions of the traditional deep auto-encoder to enforce ϵ -differential privacy. In addition to the privacy preservation in feature extraction layers, a ϵ -differential privacy preserving softmax layer was also developed for either classification or prediction. The blockchain technology discussed in Section 3 of this study is one of the trends the research community is adopting to overcome the privacy issue. Hossein et al. (2019) proposed a blockchain-based architecture for e-health applications in which users' data privacy is maintained using features like immutability and anonymity. Ul Hassan et al. (2020) adopted the differential privacy strategy in data perturbation and protect the data in the blockchain. The authors integrated the differential privacy issues in each layer of the blockchain. Another trending AI mechanism, Federated Learning is being adopted for its capacity to collaborate learning and maintain data privacy. The federated learning approach can maintain data privacy by allowing local clients to share only their local AI model parameters, not private data. Singh et al. (2021) combined blockchain and federated learning to propose a secure architecture for privacy-preserving in smart healthcare. The authors take advantage of federated learning features and send only model parameters to the cloud. Data perturbation, blockchain technology and federated learning techniques are being widely adopted to overcome patient privacy and data leakage in healthcare applications as shown in Table 13.

TABLE 12 Trends in explainable artificial intelligence (AI)

Challenge	Trends	References
Explainable AI or ML	Sensitivity	Sagi and Rokach (2020), Jovanovic et al. (2016), (Shouval et al. (2017)
	Shapley values	Lundberg and Lee (2017), Shapley (1953), Linardatos et al. (2020), Raza et al. (2022), Khodabandehloo et al. (2021)

TABLE 13 Trends in protecting privacy

Challenge	Trends	References
Privacy	Differential privacy—Data perturbation	K. Chen, Zhang, et al. (2021), Iwasawa et al. (2017), Gati et al. (2021), D. Zhang et al. (2019),
	Blockchain technology	Hossein et al. (2019), UI Hassan et al. (2020)
	Federated Learning	Singh et al. (2021), Shaik, Tao, Higgins, Gururajan, et al. (2022); Shaik, Tao, Higgins, Xie, et al. (2022), Y. Chen et al. (2020)

6.3 | Uncertainty

There are different uncertainties, such as the data acquisition process, deep neural networks (DNN) building process, and modeling results in adopting AI methodologies to healthcare applications (Gawlikowski et al., 2021). Data acquisition plays a vital role in RPM systems. Still, error, noise in measurement systems, and variability in real-world situations cause uncertainty. While building and training the DNN model with the acquired data would lead to uncertainty in the model structure and training procedure due to a large number of hyperparameters in DNN. The former two uncertainties would lead to uncertainty in the modeling results can be split into data uncertainty (aleatoric uncertainty) and model uncertainty (epistemic uncertainty) (Hüllermeier & Waegeman, 2021). Uncertainty quantification (UQ) can reduce the impact of uncertainties during both the optimization and decision-making process. Abdar et al. (2021) surveyed the research community's work on quantifying uncertainty in machine learning and deep learning models. The review article discussed ensemble techniques and Bayesian techniques like Bayesian deep learning (BDL) (H. Wang & Yeung, 2016) and Bayesian NNs (BNNs) (K. C. Wang et al., 2018) to address the reliability issue of the deep learning models and can interpret their hyperparameters. Begoli et al. (2019) also discussed the need for UQ in machine learning-assisted medical decision-making. The authors discussed four overlapping groups of challenges in UQ especially deep learning models being used in medical applications. The absence of theory in healthcare research is one of the challenges, which means without a fundamental mathematical model, the research is bound to assumptions. The second challenge is the absence of casual models due to limited conclusions from DL models. Sensitivity due to imperfect real-world data while quantifying the uncertainty. The last challenge discussed was computation expense due to deep learning training and re-computation or re-evaluation, causing additional burdens.

6.4 | Signal processing

Most signal processing issues remained with noninvasive RPMs that did not touch the patient. Information system infrastructure like RFID reader-antennas was able to retrieve data from RFID tags placed on different areas of the patients. However, transforming the tags' data into vital signs was a challenging task comprising noise (He et al., 2017; Q. Xu et al., 2020). Environmental noise obscured respiration and heartbeat signals in these device-free scenarios. RFID devices utilize a frequency hopping spread spectrum in many countries and regions, causing a discontinuous phase stream. The signal fluctuation caused by intense emotions can overwhelm the respiration and heartbeat signals, resulting in errors in signal extraction (Hou et al., 2017; Zhao et al., 2018). Signal processing challenges could be handled by taking advantage of frequency differences in vital signs and noisy data. It is evident that the double parameter of the least mean square (LMS) (He et al., 2017) can extract a respiration signal with a fundamental frequency (H. Wang et al., 2016; X. Wang et al., 2017). Contact-less respiration and heartbeat monitoring (CRH) systems (Q. Xu

et al., 2020; Zhao et al., 2018) that were designed to extract vital signs used smoothing, filtering on raw measurements, and used an intense motion detector system to extract the coarse-grained signal. This was further processed to extract respiratory and heartbeat signals. Noninvasive RPM systems have also used smoothing, unwrapping, interpolation, and Fourier transform techniques to extract breathing and heartbeat signals (Hou et al., 2017). He et al. (2017) applied a frequency-modulated continuous wave (FMCW) radar to monitor vital signs for multihuman targets. The data was collected through the chest wall, with periodic vibration to record respiratory and heart rates. The study proposed a vital signal separation method that could obtain accurate respiration and heartbeat signals using a novel double parameter, the least mean square (LMS) filter. The respiration signal was extracted at the fundamental frequency, and the heartbeat signal from the mixed physiological signal was based on the double-parameter LMS filter. Frequency differences can help in signal processing by adopting techniques such as smoothing, interpolation, Fourier transforms, and frequency filters as shown in Table 14.

6.5 | Imbalanced dataset

An imbalanced dataset is a common challenge in AI or ML for data scientists, as it can lead to bias in decision-making. In the supervised machine learning technique, class-imbalanced datasets could affect the predictive ability of the model (Gao et al., 2021). An imbalance in classification categories of a dataset where more samples are from one class is called a majority class, with the other type called a minority class. Conventional machine learning algorithms tend to predict the majority class while ignoring the minority class (Chen, Zhang, et al., 2021; Kaieski et al., 2020). The process could be either using under-sampling techniques like EasyEnsemble and BalanceCascade (Choudhary & Shukla, 2021) to reduce majority class samples or using an over-sampling technique like Synthetic Minority Oversampling Technique (SMOTE) (Hambali & Gbolagade, 2016) to reproduce minority class samples (Alotaibi & Sasi, 2016). Either of these two approaches would adequately deal with class imbalance.

Wang, Yao, and Chen (2021) proposed a long-tail data processing, undersampling-clustering-oversampling algorithm, for heart rate prediction in stroke patients. The authors use a randomly undersampling technique on majority labels and K-Means clustering on minority before applying SMOTE technique on the combined dataset. Kumar et al. (2022) performed a review of class-imbalanced learning situations with six machine-learning classifiers on five imbalanced clinical datasets. The authors explored seven different label balancing techniques such as SMOTE, SVM-SMOTE, ADASYN, Undersampling, Random Oversampling, SMOTETOMEK, and SMOTEEN. Out of all techniques, SMOTEEN with the KNN model achieved the highest accuracy, recall, precision, and F1 score.

Evaluation metrics also play a critical role in addressing bias in class imbalance problems. AI model evaluation results can be misleading. For example, in a multilabel classification problem, considering the overall accuracy of an AI model could show the model's performance in classifying each label. This can be addressed by adopting balanced accuracy, precision, recall, and F1-score (Iwendi et al., 2020) metrics which provide model performance at each label. Evaluation metrics may help to check the bias of model but oversampling and undersampling techniques are adopted by the research community to overcome data imbalance or long tail data as shown in Table 15.

6.6 | Dataset volume

Another challenge in designing an RPM system that uses AI models is the size of the dataset used for its training and predicting purposes. Most machine learning algorithms require large datasets to build a robust model. The size of the dataset matters, as this would hinder the ability of a machine learning model to perform accurately. To analyze hospitalized data or outpatient data, a good model needs to be trained with informative features with a high number of subjects (Ramos et al., 2021). A neural network model could enhance the performance as more data is available (Coppock et al., 2021). Random forests need relatively few training cases to achieve near-peak performance, are computationally cheap to train, and are able to handle large numbers of descriptors well (Teixeira et al., 2016). Data-driven models like logistic regression, SVM, or neural networks have an advantage in model derivation as these models do not require prior knowledge about the relationship between input predictor variables and output target variables. Models like decision trees, random forests, SVM, and Bayesian networks can handle large datasets and integrate background knowledge into the analysis (Awad et al., 2017).

TABLE 14 Trends in signal processing

Challenge	Trends	References
Signal processing	Fourier transforms	Hou et al. (2017), H. Wang et al. (2016), X. Wang et al. (2017)
	Least mean square (LMS) filter	He et al. (2017), Zhao et al. (2018), Q. Xu et al. (2020)

TABLE 15 Trends in data imbalance

Challenge	Trends	References
Imbalanced datasets	OverSampling, SMOTE	Hambali and Gbolagade (2016), Alotaibi and Sasi (2016), Wang, Yao, and Chen (2021), Kumar et al. (2022)
	Undersampling	Choudhary and Shukla (2021), Wang, Yao, and Chen (2021), Kumar et al. (2022)

TABLE 16 Trends in data imbalance

Challenge	Trends	References
Feature extraction	Feature engineering feature learning and representation	K. Chen, Zhang, et al. (2021), Kaieski et al. (2020), X. Zhang et al. (2017), Y. Xu et al. (2018)
	Deep learning	Zhong et al. (2016)

6.7 | Feature extraction

Feature extraction is one of the key steps RPM systems perform in analyzing human vital signs and activity recognition (K. Chen, Zhang, et al., 2021). To generate a model to predict, detect or score the patient's health state, the definition of the features must be included (Kaieski et al., 2020). Lack of efficient feature engineering process, feature selection methods, and the heterogeneity of measured patient data are some challenges limiting the effectiveness of machine learning-based predictive models (X. Zhang et al., 2017). Within ICU, patients are monitored continuously by numerous specialized devices at the bedside, which generates high-density multiple data modalities. As a result, the timestamps, order, and frequency of the measurements may be profoundly different from one patient to another. This type of irregularity and heterogeneity in patient data make feature selection even more challenging (Y. Xu et al., 2018). To overcome the lack of efficient feature engineering, feature selection techniques, and the heterogeneity of measured patient data, feature learning or representation learning techniques can be used. Deep learning algorithms such as RNNs, LSTM, CNN, and other algorithms based on neural networks have the capability to learn this type of data structure (Zhong et al., 2016). The feature extraction challenge can be addressed by adopting feature representation techniques. Deep learning can overcome this issue without any additional framework. The research works presenting the trends with corresponding research works are presented in Table 16.

7 | FUTURE DIRECTION OF AI ON RPM

The future direction of the research is to extend the scope of AI in RPM applications to enhance healthcare services for both providers and patients. To achieve this, the challenges in adopting AI to RPM as well as AI implementation have to be addressed. As discussed in the previous section, the major challenges in adopting AI are AI or ML explainability, privacy, and uncertainty. Explainability of AI results needs to be improved as it assists healthcare professionals in understanding patients' health status better and helps in decision-making. Explainable AI approaches are working in this direction. Explainability techniques such as SHAP, LIME, DeepLIFT, and so on are being adopted to RPM systems widely. This has to be further improved to breakdown the state-of-the-art deep learning and machine learning results to healthcare practitioners to make informed decisions.

Data privacy and the security of patients' health is another major issue that could be addressed with federated learning. However, there is no strong research evidence in the federated learning concept to confirm that the reverse engineering of the local client model parameters would not lead to a patient's private data. Future works should be focused on strengthening the federated learning framework for data privacy and security. Blockchain technology has proved its capacity in maintaining data privacy with transparency and immutability. However, the implementation of blockchain technology demands high implementation costs and high energy dependence. Future works need to concentrate on blockchain engineering issues.

Uncertainty due to model structure and hyperparameters causes uncertainty in results. UQ technique can play a vital role to reduce uncertainties in the model during optimization and decision-making. The aleatoric and epistemic uncertainty can be addressed through different probabilistic and nonprobabilistic, and inverse uncertainty techniques. Focusing on this challenge would help to improve healthcare professionals' trust in machine learning or deep learning model results.

Data-related challenges are inevitable in AI-based applications. It could be due to IoT devices' signal processing, noise, data imbalance, limited labeled data, and feature extraction. Efficient and clean data is the first and most time-consuming part of AI methodology. The challenges such as label imbalance or long-tail data have to be processed with class balancing discussed in this study. The research community should concentrate on input data related at the most to achieve efficient and effective results.

Reinforcement learning has the potential to mimic human behavior and build social assistive robots for patients in a hospital or at home. Although this AI technique has been used for a while, there has not been much research into applying it to healthcare applications. Taking advantage of the sequential decision-making ability of reinforcement learning, healthcare applications such as dynamic treatment regimens and Just-in-Time-Adaptive-Interventions (JITAI) can be further enhanced. However, there have been recent incidents of deploying physical robots causing threats to humans (Stein et al., 2022). It is recommended to build virtual robots using reinforcement learning agents to monitor patients and predict unprecedented events.

8 | CONCLUSION

Healthcare applications have been widely transformed by technological innovations in information systems and AI. In particular, the last decade has revolutionized monitoring patients' health status by tracking their vital signs and physical activities. Advancements in data transmission and data modeling enabled RPM systems to detect patients' health deterioration in advance, customize patient-centric applications, and learn their behavior patterns adaptively. The transformation of RPM systems using noninvasive information system technologies like telehealth, IoT, cloud, fog, edge, and blockchain are explored in this study. The primary focus of this survey article is to present the role of AI in enhancing RPMs with its ability to learn, predict, and classify patients' behavior and vital signs. Applications of AI in monitoring vital signs, physical activities, chronic diseases, and patient emergencies are investigated. Federated learning facilitates a patient-centric monitoring system to focus on their needs while protecting data privacy. Reinforcement learning enhances RPMs to learn patient behavior patterns in a dynamic environment adaptively. The impact of such advanced AI methodologies on RPM systems is detailed with evidence. Even though AI has the potential to transform RPM services, it has certain challenges like explainability, privacy, and uncertainty. Other than this, data learning challenges include feature extraction, imbalanced labels, data volume, and data processing. In this study, the trends and challenges of AI in RPM are discussed in detail.

This study's limitations are that study is focused on RPM systems with vital signs monitoring and physical activities monitoring but not the electroencephalogram (EEG) monitoring and neurological system-related diseases. Also, this study has not explored all chronic disease monitoring research works. While addressing the limitations and challenges discussed in the study, healthcare applications should adopt advanced technological infrastructures like Cloud/Edge/Fog/Blockchain and AI methods such as Federated Learning and Reinforcement Learning. So far, traditional AI methods such as supervised and unsupervised have demonstrated state-of-the-results. However, this is the right time to transform healthcare for preventive, predictive, and personalized monitoring of patients and provide enhanced assistance to healthcare practitioners.

AUTHOR CONTRIBUTIONS

Thanveer Shaik: Conceptualization (equal); data curation (equal); formal analysis (equal); investigation (equal); methodology (equal); writing – original draft (equal); writing – review and editing (equal). **Xiaohui Tao:** Conceptualization (equal); formal analysis (equal); investigation (equal); methodology (equal); project administration (equal); supervision (equal); writing – original draft (equal); writing – review and editing (equal). **Niall Higgins:** Conceptualization (equal); methodology (equal); supervision (equal); writing – original draft (equal); writing – review and editing (equal). **Lin Li:** Formal analysis (equal); investigation (equal); methodology (equal); writing – review and editing (equal). **Raj Gururajan:** Project administration (equal); supervision (equal); writing – review and editing (equal). **Xujuan Zhou:** Supervision (equal); writing – review and editing (equal). **U. Rajendra Acharya:** Writing – review and editing (equal).

ACKNOWLEDGMENT

Open access publishing facilitated by University of Southern Queensland, as part of the Wiley - University of Southern Queensland agreement via the Council of Australian University Librarians.

CONFLICT OF INTEREST

All authors declare there is no conflict of interest in this work.

DATA AVAILABILITY STATEMENT

Data sharing not applicable to this article as no datasets were generated or analyzed during the current study.

ORCID

Thanveer Shaik  <https://orcid.org/0000-0002-9730-665X>

Xiaohui Tao  <https://orcid.org/0000-0002-0020-077X>

Niall Higgins  <https://orcid.org/0000-0002-3260-1711>

Raj Gururajan  <https://orcid.org/0000-0002-5919-0174>

Xujuan Zhou  <https://orcid.org/0000-0002-1736-739X>

U. Rajendra Acharya  <https://orcid.org/0000-0003-2689-8552>

RELATED WIREs ARTICLES

[Internet of Things and data mining: From applications to techniques and systems](#)

[Internet of Things and data analytics: A current review](#)

[Healthcare 4.0: A review of frontiers in digital health](#)

FURTHER READING

Dev, A., & Malik, S. K. (2021). Artificial bee colony optimized deep neural network model for handling imbalanced stroke data. *International Journal of E-Health and Medical Communications*, 12(5), 67–83. <https://doi.org/10.4018/ijehmc.20210901.oa5>

Guillame-Bert, M., Dubrawski, A., Wang, D., Hravnak, M., Clermont, G., & Pinsky, M. R. (2016). Learning temporal rules to forecast instability in continuously monitored patients. *Journal of the American Medical Informatics Association*, 24(1), 47–53. <https://doi.org/10.1093/jamia/ocw048>

Kam, H. J., & Kim, H. Y. (2017). Learning representations for the early detection of sepsis with deep neural networks. *Computers in Biology and Medicine*, 89, 248–255. <https://doi.org/10.1016/j.compbimed.2017.08.015>

Mahtta, D., Daher, M., Lee, M. T., Sayani, S., Shishehbor, M., & Virani, S. S. (2021). Promise and perils of telehealth in the current era. *Current Cardiology Reports*, 23(9), 115. <https://doi.org/10.1007/s11886-021>

Wilcock, A. D., Schwamm, L. H., Zubizarreta, J. R., Zachrisson, K. S., Uscher-Pines, L., Richard, J. V., & Mehrotra, A. (2021). Reperfusion treatment and stroke outcomes in hospitals with telestroke capacity. *JAMA Neurology*, 78(5), 527. <https://doi.org/10.1001/jamaneurol.2021.0023>

Williams, K., Markwardt, S., Kearney, S. M., Karp, J. F., Kraemer, K. L., Park, M. J., Freund, P., Watson, A., Schuster, J., & Beckjord, E. (2021). Addressing implementation challenges to digital care delivery for adults with multiple chronic conditions: Stakeholder feedback in a randomized controlled trial. *JMIR mHealth and uHealth*, 9(2), e23498. <https://doi.org/10.2196/23498>

REFERENCES

Abdar, M., Pourpanah, F., Hussain, S., Rezazadegan, D., Liu, L., Ghavamzadeh, M., Fieguth, P., Cao, X., Khosravi, A., Acharya, U. R., Makarek, V., & Nahavandi, S. (2021). A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76, 243–297. <https://doi.org/10.1016/j.inffus.2021.05.008>

- Adamou, M., Antoniou, G., Greasidou, E., Lagani, V., Charonyktakis, P., Tsamardinos, I., & Doyle, M. (2019). Toward automatic risk assessment to support suicide prevention. *Crisis, 40*(4), 249–256. <https://doi.org/10.1027/0227-5910/a000561>
- Aliyu, F., Sheltami, T., Mahmoud, A., Al-Awami, L., & Yasar, A. (2021). Detecting man-in-the-middle attack in fog computing for social media. *Computers, Materials & Continua, 69*(1), 1159–1181. <https://doi.org/10.32604/cmc.2021.016938>
- Alotaibi, N. N., & Sasi, S. (2016). Stroke in-patients' transfer to the ICU using ensemble based model. In *In 2016 international conference on electrical, electronics, and optimization techniques (ICEEOT)*. IEEE. <https://doi.org/10.1109/iceeot.2016.7755040>
- Alshwaheen, T. I., Hau, Y. W., Ass'Ad, N., & Abualsamen, M. M. (2021). A novel and reliable framework of patient deterioration prediction in intensive care unit based on long short-term memory-recurrent neural network. *IEEE Access, 9*, 3894–3918. <https://doi.org/10.1109/ACCESS.2020.3047186>
- Alwakeel, A. M. (2021). An overview of fog computing and edge computing security and privacy issues. *Sensors, 21*(24), 8226. <https://doi.org/10.3390/s21248226>
- Ankita, R. S., Babbar, H., Coleman, S., Singh, A., & Aljahdali, H. M. (2021). An efficient and lightweight deep learning model for human activity recognition using smartphones. *Sensors, 21*(11), 3845. <https://doi.org/10.3390/s21113845>
- Antunes, A. G., Teixeira, C., Vaz, A. M., Martins, C., Queiro's, P., Alves, A., Velasco, F., Peixe, B., & Guerreiro, H. (2017). Comparison of the prognostic value of chronic liver failure consortium scores and traditional models for predicting mortality in patients with cirrhosis. *Gastroenterologia y Hepatologia (English Edition), 40*(4), 276–285. <https://doi.org/10.1016/j.gastre.2017.03.012>
- Asiimwe, S. B., Vittinghoff, E., & Whooley, M. (2020). Vital signs data and probability of hospitalization, transfer to another facility, or emergency department death among adults presenting for medical illnesses to the emergency department at a large urban hospital in the United States. *The Journal of Emergency Medicine, 58*(4), 570–580. <https://doi.org/10.1016/j.jemermed.2019.11.020>
- Awad, A., Bader-El-Den, M., McNicholas, J., & Briggs, J. (2017). Early hospital mortality prediction of intensive care unit patients using an ensemble learning approach. *International Journal of Medical Informatics, 108*, 185–195. <https://doi.org/10.1016/j.ijmedinf.2017.10.002>
- Begoli, E., Bhattacharya, T., & Kusnezov, D. (2019). The need for uncertainty quantification in machine-assisted medical decision making. *Nature Machine Intelligence, 1*(1), 20–23. <https://doi.org/10.1038/s42256-018-0004-1>
- Bekiri, R., Djeflal, A., & Hettiri, M. (2020). A remote medical monitoring system based on data mining. In *020 1st international conference on communications, control systems and signal processing (CCSSP)*. IEEE. <https://doi.org/10.1109/ccssp49278.2020.9151713>
- Bini, S. A. (2018). Artificial intelligence, machine learning, deep learning, and cognitive computing: What do these terms mean and how will they impact health care? *The Journal of Arthroplasty, 33*(8), 2358–2361. <https://doi.org/10.1016/j.arth.2018.02.067>
- Blasiak, A., Khong, J., & Kee, T. (2020). CURATE.AI: Optimizing personalized medicine with artificial intelligence. *SLAS Technology, 25*(2), 95–105. <https://doi.org/10.1177/2472630319890316>
- Bousefsaf, F., Pruski, A., & Maaoui, C. (2019). 3d convolutional neural networks for remote pulse rate measurement and mapping from facial video. *Applied Sciences, 9*(20), 4364. <https://doi.org/10.3390/app9204364>
- Busso, M., Gonzalez, M. P., & Scartascini, C. (2022). On the demand for telemedicine: Evidence from the COVID-19 pandemic. *Health Economics, 31*(7), 1491–1505. <https://doi.org/10.1002/hec.4523>
- Chen, G., Xiao, X., Zhao, X., Tat, T., Bick, M., & Chen, J. (2021). Electronic textiles for wearable point-of-care systems. *Chemical Reviews, 122*(3), 3259–3291. <https://doi.org/10.1021/acs.chemrev.1c00502>
- Chen, I. Y., Joshi, S., Ghassemi, M., & Ranganath, R. (2021). Probabilistic machine learning for healthcare. *Annual Review of Biomedical Data Science, 4*(1), 393–415. <https://doi.org/10.1146/annurev-biodatasci-092820-033938>
- Chen, K., Zhang, D., Yao, L., Guo, B., Yu, Z., & Liu, Y. (2021). Deep learning for sensor-based human activity recognition. *ACM Computing Surveys, 54*(4), 1–40.
- Chen, Y., Qin, X., Wang, J., Yu, C., & Gao, W. (2020). FedHealth: A federated transfer learning framework for wearable healthcare. *IEEE Intelligent Systems, 35*(4), 83–93. <https://doi.org/10.1109/mis.2020.2988604>
- Cho, Y., Bianchi-Berthouze, N., & Julier, S. J. (2017). DeepBreath: Deep learning of breathing patterns for automatic stress recognition using low-cost thermal imaging in unconstrained settings. In *In 2017 seventh international conference on affective computing and intelligent interaction (ACII)*. IEEE. <https://doi.org/10.1109/acii.2017.8273639>
- Choudhary, R., & Shukla, S. (2021). A clustering based ensemble of weighted kernelized extreme learning machine for class imbalance learning. *Expert Systems with Applications, 164*, 114041. <https://doi.org/10.1016/j.eswa.2020.114041>
- Chowdary, M. K., Nguyen, T. N., & Hemanth, D. J. (2021). Deep learning-based facial emotion recognition for human-computer interaction applications. *Neural Computing and Applications, 1*–18. <https://doi.org/10.1007/s00521-021-06012-8>
- Colopy, G. W., Roberts, S. J., & Clifton, D. A. (2018). Bayesian optimization of personalized models for patient vital-sign monitoring. *IEEE Journal of Biomedical and Health Informatics, 22*(2), 301–310. <https://doi.org/10.1109/jbhi.2017.2751509>
- Coppock, H., Gaskell, A., Tzirakis, P., Baird, A., Jones, L., & Schuller, B. (2021). End-to-end convolutional neural network enables COVID-19 detection from breath and cough audio: A pilot study. *BMJ Innovations, 7*(2), 356–362. <https://doi.org/10.1136/bmjinnov-2021-000668>
- da Silva, D. B., Schmidt, D., da Costa, C. A., da Rosa Righi, R., & Eskofier, B. (2021). DeepSigns: A predictive model based on deep learning for the early detection of patient health deterioration. *Expert Systems with Applications, 165*, 113905. <https://doi.org/10.1016/j.eswa.2020.113905>
- Dean, N. C., Vines, C. G., Carr, J. R., Rubin, J. G., Webb, B. J., Jacobs, J. R., Butler, A. M., Lee, J., Jephson, A. R., Jenson, N., Walker, M., Brown, S. M., Irvin, J. A., Lungren, M. P., & Allen, T. L. (2022). A pragmatic, stepped wedge, cluster-controlled clinical trial of real-time pneumonia clinical decision support. *American Journal of Respiratory and Critical Care Medicine, 205*(11), 1330–1336. <https://doi.org/10.1164/rccm.202109-2092OC>

- Devi, R. L., & Kalaivani, V. (2019). Machine learning and IoT-based cardiac arrhythmia diagnosis using statistical and dynamic features of ECG. *The Journal of Supercomputing*, 76(9), 6533–6544. <https://doi.org/10.1007/s11227-019-02873-y>
- Dias, D., & Cunha, J. P. S. (2018). Wearable health devices—Vital sign monitoring, systems and technologies. *Sensors*, 18(8), 2414. <https://doi.org/10.3390/s18082414>
- Downey, C., Chapman, S., Randell, R., Brown, J., & Jayne, D. (2018). The impact of continuous versus intermittent vital signs monitoring in hospitals: A systematic review and narrative synthesis. *International Journal of Nursing Studies*, 84, 19–27. <https://doi.org/10.1016/j.ijnurstu.2018.04.013>
- Drake, C., Zhang, Y., Chaiyachati, K. H., & Polsky, D. (2019). The limitations of poor broadband internet access for telemedicine use in rural america: An observational study. *Annals of Internal Medicine*, 171(5), 382–384. <https://doi.org/10.7326/m19-0283>
- Ede, J., Vollam, S., Darbyshire, J. L., Gibson, O., Tarassenko, L., & Watkinson, P. (2021). Non-contact vital sign monitoring of patients in an intensive care unit: A human factors analysis of staff expectations. *Applied Ergonomics*, 90, 103149. <https://doi.org/10.1016/j.apergo.2020.103149>
- Efat, M. I. A., Rahman, S., & Rahman, T. (2020). IoT based smart health monitoring system for diabetes patients using neural network. In *Cyber security and computer science* (pp. 593–606). Springer International Publishing. <https://doi.org/10.1007/978-3-030-52856-047>
- Elola, A., Aramendi, E., Irusta, U., Picon, A., Alonso, E., Owens, P., & Idris, A. (2019). Deep neural networks for ECG-based pulse detection during out-of-hospital cardiac arrest. *Entropy*, 21(3), 305. <https://doi.org/10.3390/e21030305>
- El-Rashidy, N., El-Sappagh, S., Islam, S. M. R., El-Bakry, H. M., & Abdelrazek, S. (2021). Mobile health in remote patient monitoring for chronic diseases: Principles, trends, and challenges. *Diagnostics*, 11(4), 607. <https://doi.org/10.3390/diagnostics11040607>
- Fang, J., Liu, Y., Lee, E., & Yadav, K. (2020). Telehealth solutions for in-hospital communication with patients under isolation during COVID-19. *Western Journal of Emergency Medicine*, 21(4), 801. <https://doi.org/10.5811/westjem.2020.5.48165>
- Faruk, M. J. H., Shahriar, H., Valero, M., Sneha, S., Ahamed, S. I., & Rahman, M. (2021). Towards blockchain-based secure data management for remote patient monitoring. In *2021 IEEE international conference on digital health (ICDH)*. IEEE. <https://doi.org/10.1109/icdh52753.2021.00054>
- Gao, L., Lu, P., & Ren, Y. (2021). A deep learning approach for imbalanced crash data in predicting highway rail grade crossings accidents. *Reliability Engineering & System Safety*, 216, 108019. <https://doi.org/10.1016/j.res.2021.108019>
- Garca-del Valle, S., Arnal-Velasco, D., Molina-Mendoza, R., & Gomez-Arnau, J. I. (2021). Update on early warning scores. *Best Practice & Research Clinical Anaesthesiology*, 35(1), 105–113. <https://doi.org/10.1016/j.bpa.2020.12.013>
- Gati, N. J., Yang, L. T., Feng, J., Nie, X., Ren, Z., & Tarus, S. K. (2021). Differentially private data fusion and deep learning framework for cyber-physical-social systems: State-of-the-art and perspectives. *Information Fusion*, 76, 298–314. <https://doi.org/10.1016/j.inffus.2021.04.017>
- Gawlikowski, J., Tassi, C. R. N., Ali, M., Lee, J., Humt, M., Feng, J., Kruspe, A., Triebel, R., Jung, P., Roscher, R., & Shahzad, M. (2021). A survey of uncertainty in deep neural networks. <https://doi.org/10.48550/ARXIV.2107.03342>
- Gönül, S., Namlı, T., Coşar, A., & Toroslu, İ. H. (2021). A reinforcement learning based algorithm for personalization of digital, just-in-time, adaptive interventions. *Artificial Intelligence in Medicine*, 115(0933-3657), 102062. <https://doi.org/10.1016/j.artmed.2021.102062>
- Hambali, A. M., & Gbolagade, D. M. (2016). Ovarian cancer classification using hybrid synthetic minority over-sampling technique and neural network. *Journal of Advances in Computer Research*, 7(4), 109–124 <https://www.jacr.sari.iau.ir/article651012.html>
- Harrison, E., Chang, M., Hao, Y., & Flower, A. (2018). Using machine learning to predict near-term mortality in cirrhosis patients hospitalized at the university of Virginia health system. In *In 2018 systems and information engineering design symposium (SIEDS)*. IEEE. <https://doi.org/10.1109/sieds.2018.8374719>
- Hathaliya, J., Sharma, P., Tanwar, S., & Gupta, R. (2019). Blockchain-based remote patient monitoring in healthcare4.0. In *2019 IEEE 9th international conference on advanced computing (IACC)*. IEEE. <https://doi.org/10.1109/iacc48062.2019.8971593>
- He, M., Nian, Y., & Gong, Y. (2017). Novel signal processing method for vital sign monitoring using FMCW radar. *Biomedical Signal Processing and Control*, 33, 335–345. <https://doi.org/10.1016/j.bspc.2016.12.008>
- Heijmans, M., Habets, J., Kuijff, M., Kubben, P., & Herff, C. (2019). Evaluation of Parkinson's disease at home: Predicting tremor from wearable sensors. In *2019 41st annual international conference of the IEEE engineering in medicine and biology society (EMBC)*. IEEE. <https://doi.org/10.1109/embc.2019.8857717>
- Helm, J. M., Swiergosz, A. M., Haeberle, H. S., Karnuta, J. M., Schaffer, J. L., Krebs, V. E., Spitzer, A. I., & Ramkumar, P. N. (2020). Machine learning and artificial intelligence: Definitions, applications, and future directions. *Current Reviews in Musculoskeletal Medicine*, 13(1), 69–76. <https://doi.org/10.1007/s12178-020-09600-8>
- Helman, S., Terry, M. A., Pellathy, T., Williams, A., Dubrawski, A., Clermont, G., Pinsky, M. R., Al-Zaiti, S., & Hrvanek, M. (2022). Engaging clinicians early during the development of a graphical user display of an intelligent alerting system at the bedside. *International Journal of Medical Informatics*, 159, 104643. <https://doi.org/10.1016/j.ijmedinf.2021.104643>
- Hosseini, K. M., Esmaili, M., Dargahi, T., & Khonsari, A. (2019). Blockchain-based privacy-preserving healthcare architecture. In *2019 IEEE Canadian Conference of Electrical and Computer Engineering (CCECE)*, (2576-7046). IEEE. <https://doi.org/10.1109/CCECE.2019.8861857>
- Hou, Y., Wang, Y., & Zheng, Y. (2017). TagBreathe: Monitor breathing with commodity RFID systems. In *2017 IEEE 37th international conference on distributed computing systems (ICDCS)*. IEEE. <https://doi.org/10.1109/icdcs.2017.76>
- Hsieh, C.-Y., Huang, H.-Y., Liu, K.-C., Liu, C.-P., Chan, C.-T., & Hsu, S. J.-P. (2021). Multiphase identification algorithm for fall recording systems using a single wearable inertial sensor. *Sensors*, 21(9), 3302. <https://doi.org/10.3390/s21093302>

- Hüllermeier, E., & Waegeman, W. (2021). Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. *Machine Learning*, *110*(3), 457–506. <https://doi.org/10.1007/s10994-021-05946-3>
- Huang, C., Fukushi, K., Wang, Z., Nihey, F., Kajitani, H., & Nakahara, K. (2022). Method for estimating temporal gait parameters concerning bilateral lower limbs of healthy subjects using a single in-shoe motion sensor through a gait event detection approach. *Sensors*, *22*(1), 251. <https://doi.org/10.3390/s22010351>
- Hui, X., & Kan, E. C. (2017). Monitoring vital signs over multiplexed radio by near-field coherent sensing. *Nature Electronics*, *1*(1), 74–78. <https://doi.org/10.1038/s41928-017-0001-0>
- Hui, X., & Kan, E. C. (2018). Accurate extraction of heartbeat intervals with near-field coherent sensing. In *2018 IEEE international conference on communications (ICC)*. IEEE. <https://doi.org/10.1109/icc.2018.8423000>
- Iranpak, S., Shahbahrami, A., & Shakeri, H. (2021). Remote patient monitoring and classifying using the internet of things platform combined with cloud computing. *Journal of Big Data*, *8*(1), 120. <https://doi.org/10.1186/s40537-021-00507-w>
- Iwasawa, Y., Nakayama, K., Yairi, I., & Matsuo, Y. (2017). Privacy issues regarding the application of DNNs to activity-recognition using wearables and its countermeasures by use of adversarial training. In *Proceedings of the twenty sixth international joint conference on artificial intelligence*. International Joint Conferences on Artificial Intelligence Organization <https://www.ijcai.org/Proceedings/2017/0268.pdf>
- Iwendi, C., Bashir, A. K., Peshkar, A., Sujatha, R., Chatterjee, J. M., Pasupuleti, S., Mishra, R., Pillai, S., & Jo, O. (2020). COVID-19 patient health prediction using boosted random forest algorithm. *Frontiers in Public Health*, *8*. <https://doi.org/10.3389/fpubh.2020.00357>
- Joshi, M., Archer, S., Morbi, A., Arora, S., Kwasnicki, R., Ashrafian, H., Khan, S., Cooke, G., & Darzi, A. (2021). Short-term wearable sensors for in-hospital medical and surgical patients: Mixed methods analysis of patient perspectives. *JMIR Perioperative Medicine*, *4*(1), e18836. <https://doi.org/10.2196/18836>
- Jovanovic, M., Radovanovic, S., Vukicevic, M., Poucke, S. V., & Delibasic, B. (2016). Building interpretable predictive models for pediatric hospital readmission using tree-lasso logistic regression. *Artificial Intelligence in Medicine*, *72*, 12–21. <https://doi.org/10.1016/j.artmed.2016.07.003>
- Kaieski, N., da Costa, C. A., da Rosa Righi, R., Lora, P. S., & Eskofier, B. (2020). Application of artificial intelligence methods in vital signs analysis of hospitalized patients: A systematic literature review. *Applied Soft Computing*, *96*, 106612. <https://doi.org/10.1016/j.asoc.2020.106612>
- Kalfa, D., Agrawal, S., Goldshtrom, N., LaPar, D., & Bacha, E. (2020). Wireless monitoring and artificial intelligence: A bright future in cardiothoracic surgery. *The Journal of Thoracic and Cardiovascular Surgery*, *160*(3), 809–812. <https://doi.org/10.1016/j.jtcvs.2019.08.141>
- Kellett, J., & Sebat, F. (2017). Make vital signs great again—A call for action. *European Journal of Internal Medicine*, *45*, 13–19. <https://doi.org/10.1016/j.ejim.2017.09.018>
- Khalid, W. B., Anwar, A., & Waheed, O. T. (2022). Contactless vitals measurement robot. In *In 2022 8th international conference on automation, robotics and applications (ICARA)*. IEEE. <https://doi.org/10.1109/icara55094.2022.9738523>
- Khodabandehloo, E., Riboni, D., & Alimohammadi, A. (2021). HealthXAI: Collaborative and explainable AI for supporting early diagnosis of cognitive decline. *Future Generation Computer Systems*, *116*, 168–189. <https://doi.org/10.1016/j.future.2020.10.030>
- Kong, G., Xu, D.-L., Yang, J.-B., Yin, X., Wang, T., Jiang, B., & Hu, Y. (2016). Belief rule-based inference for predicting trauma outcome. *Knowledge-Based Systems*, *95*, 35–44. <https://doi.org/10.1016/j.knosys.2015.12.002>
- Krittanawong, C., Johnson, K. W., Choi, E., Kaplin, S., Venner, E., Murugan, M., Wang, Z., Glicksberg, B. S., Amos, C. I., Schatz, M. C., & Tang, W. W. (2022). Artificial intelligence and cardiovascular genetics. *Life*, *12*(2), 279. <https://doi.org/10.3390/life12020279>
- Kumar, M., & Chand, S. (2021). MedHypChain: A patient-centered interoperability hyperledger-based medical healthcare system: Regulation in COVID-19 pandemic. *Journal of Network and Computer Applications*, *179*, 102975. <https://doi.org/10.1016/j.jnca.2021.102975>
- Kumar, V., Lalotra, G. S., Sasikala, P., Rajput, D. S., Kaluri, R., Lakshmana, K., Shorfuzzaman, M., Alsufyani, A., & Uddin, M. (2022). Addressing binary classification over class imbalanced clinical datasets using computationally intelligent techniques. *Healthcare*, *10*(7), 1293. <https://doi.org/10.3390/healthcare10071293>
- Laber, E. B., Lizotte, D. J., Qian, M., Pelham, W. E., & Murphy, S. A. (2014). Dynamic treatment regimes: Technical challenges and applications. *Electronic Journal of Statistics*, *8*(1), 1225–1272. <https://doi.org/10.1214/14-ejs920>
- Laurie, J., Higgins, N., Peynot, T., Fawcett, L., & Robert, J. (2021). An evaluation of a video magnification-based system for respiratory rate monitoring in an acute mental health setting. *International Journal of Medical Informatics*, *148*, 104378. <https://doi.org/10.1016/j.ijmedinf.2021.104378>
- Lauritsen, S. M., Kristensen, M., Olsen, M. V., Larsen, M. S., Lauritsen, K. M., Jørgensen, M. J., Lange, J., & Thiesson, B. (2020). Explainable artificial intelligence model to predict acute critical illness from electronic health records. *Nature Communications*, *11*(1), 3852. <https://doi.org/10.1038/s41467-020-17431-x>
- Lin, C., Zhang, Y., Ivy, J., Capan, M., Arnold, R., Huddleston, J. M., & Chi, M. (2018). Early diagnosis and prediction of sepsis shock by combining static and dynamic information using convolutional-LSTM. In *2018 IEEE international conference on healthcare informatics (ICHI)*. IEEE. <https://doi.org/10.1109/ichi.2018.00032>
- Linardatos, P., Papastefanopoulos, V., & Kotsiantis, S. (2020). Explainable AI: A review of machine learning interpretability methods. *Entropy*, *23*(1), 18. <https://doi.org/10.3390/e23010018>
- Liu, C., Zhang, X., Zhao, L., Liu, F., Chen, X., Yao, Y., & Li, J. (2019). Signal quality assessment and lightweight QRS detection for wearable ECG SmartVest system. *IEEE Internet of Things Journal*, *6*(2), 1363–1374. <https://doi.org/10.1109/jiot.2018.2844090>

- Liu, H., Wang, L., Lin, G., & Feng, Y. (2022). Recent progress in the fabrication of flexible materials for wearable sensors. *Biomaterials Science*, 10(3), 614–632. <https://doi.org/10.1039/d1bm01136g>
- Liu, Z., Zhu, T., Wang, J., Zheng, Z., Li, Y., Li, J., & Lai, Y. (2022). Functionalized fiber-based strain sensors: Pathway to next generation wearable electronics. *Nano-Micro Letters*, 14(1), 1–39. <https://doi.org/10.1007/s40820-022-00806-8>
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17)* (pp. 4768–4777). Curran Associates Inc., Red Hook, NY, USA. <https://dl.acm.org/doi/epdf/10.5555/3295222.3295230>
- Mahesh, V. G. V., Chen, C., Rajangam, V., Raj, A. N. J., & Krishnan, P. T. (2021). Shape and texture aware facial expression recognition using spatial pyramid zernike moments and law's textures feature set. *IEEE Access*, 9, 52509–52522. <https://doi.org/10.1109/access.2021.3069881>
- Malasinghe, L. P., Ramzan, N., & Dahal, K. (2017). Remote patient monitoring: A comprehensive study. *Journal of Ambient Intelligence and Humanized Computing*, 10(1), 57–76. <https://doi.org/10.1007/s12652-017-0598-x>
- McGinty, E. E., Presskreischer, R., Breslau, J., Brown, J. D., Domino, M. E., Druss, B. G., Horvitz-Lennon, M., Murphy, K. A., Pincus, H. A., & Daumit, G. L. (2021). Improving physical health among people with serious mental illness: The role of the specialty mental health sector. *Psychiatric Services*, 72(11), 1301–1310. <https://doi.org/10.1176/appi.ps.202000768>
- Meskó, B., Drobni, Z., Bényei, É., Gergely, B., & Györfy, Z. (2017). Digital health is a cultural transformation of traditional healthcare. *mHealth*, 3, 38. <https://doi.org/10.21037/mhealth.2017.08.07>
- Miller, D. D., & Brown, E. W. (2018). Artificial intelligence in medical practice: The question to the answer? *The American Journal of Medicine*, 131(2), 129–133. <https://doi.org/10.1016/j.amjmed.2017.10.035>
- Mohanty, A., & Mishra, S. (2022). A comprehensive study of explainable artificial intelligence in healthcare. In *Augmented intelligence in healthcare: A pragmatic and integrated analysis* (pp. 475–502). Springer Nature Singapore. <https://doi.org/10.1007/978-981-19-1076-025>
- Mujumdar, A., & Vaidehi, V. (2019). Diabetes prediction using machine learning algorithms. *Procedia Computer Science*, 165, 292–299. <https://doi.org/10.1016/j.procs.2020.01.047>
- Mukherjee, A., Ghosh, S., Behere, A., Ghosh, S. K., & Buyya, R. (2020). Internet of health things (IoHT) for personalized health care using integrated edge-fog-cloud network. *Journal of Ambient Intelligence and Humanized Computing*, 12(1), 943–959. <https://doi.org/10.1007/s12652-020-02113-9>
- Muralitharan, S., Nelson, W., Di, S., McGillion, M., Devereaux, P., Barr, N. G., & Petch, J. (2020). Machine learning-based early warning systems for clinical deterioration: Systematic scoping review (preprint). <https://doi.org/10.2196/25187>
- Naeem, M., Paragliola, G., & Coronato, A. (2021). A reinforcement learning and deep learning based intelligent system for the support of impaired patients in home treatment. *Expert Systems with Applications*, 168, 114285. <https://doi.org/10.1016/j.eswa.2020.114285>
- Nahum-Shani, I., Smith, S. N., Spring, B. J., Collins, L. M., Witkiewitz, K., Tewari, A., & Murphy, S. A. (2017). Just-in-time adaptive interventions (JITAs) in mobile health: Key components and design principles for ongoing health behavior support. *Annals of Behavioral Medicine*, 52(6), 446–462. <https://doi.org/10.1007/s12160-016-9830-8>
- Neto, L. A. S. M., Pequeno, R., Almeida, C., Galdino, K., Martins, F., & de Moura, A. V. (2017). A method for intelligent support to medical diagnosis in emergency cardiac care. In *2017 international joint conference on neural networks (IJCNN)*. IEEE. <https://doi.org/10.1109/ijcnn.2017.7966438>
- Nguyen, D. C., Pham, Q.-V., Pathirana, P. N., Ding, M., Seneviratne, A., Lin, Z., Dobre, O., & Hwang, W.-J. (2023). Federated learning for smart healthcare: A survey. *ACM Computing Surveys*, 55(3), 1–37. <https://doi.org/10.1145/3501296>
- Nord, G., Rising, K. L., Band, R. A., Carr, B. G., & Hollander, J. E. (2019). On-demand synchronous audio video telemedicine visits are cost effective. *The American Journal of Emergency Medicine*, 37(5), 890–894. <https://doi.org/10.1016/j.ajem.2018.08.017>
- Oh, J., Cho, D., Park, J., Na, S. H., Kim, J., Heo, J., Shin, C. S., Kim, J. J., Park, J. Y., & Lee, B. (2018). Prediction and early detection of delirium in the intensive care unit by using heart rate variability and machine learning. *Physiological Measurement*, 39(3), 035004. <https://doi.org/10.1088/1361-6579/aaab07>
- Ong, M. E. H., Ng, C. H. L., Goh, K., Liu, N., Koh, Z., Shahidah, N., Zhang, T. T., Fook-Chong, S., & Lin, Z. (2012). Prediction of cardiac arrest in critically ill patients presenting to the emergency department using a machine learning score incorporating heart rate variability compared with the modified early warning score. *Critical Care*, 16(3), R108. <https://doi.org/10.1186/cc11396>
- Ouzzani, M., Hammady, H., Fedorowicz, Z., & Elmagarmid, A. (2016). Rayyan—A web and mobile app for systematic reviews. *Systematic Reviews*, 5(1). <https://doi.org/10.1186/s13643-016-0384-4>
- Page, M. J., McKenzie, J. E., Bossuyt, P. M., Boutron, I., Hoffmann, T. C., Mulrow, C. D., Shamseer, L., Tetzlaff, J. M., Akl, E. A., Brennan, S. E., Chou, R., Glanville, J., Grimshaw, J. M., Hróbjartsson, A., Lalu, M. M., Li, T., Loder, E. W., Mayo-Wilson, E., McDonald, S., ... Moher, D. (2021). The PRISMA 2020 statement: An updated guideline for reporting systematic reviews. *BMJ*, 372, n71. <https://doi.org/10.1136/bmj.n71>
- Pan, D., Liu, H., Qu, D., & Zhang, Z. (2020). Human falling detection algorithm based on multisensor data fusion with SVM. *Mobile Information Systems*, 2020, 1–9. <https://doi.org/10.1155/2020/8826088>
- Pareek, K., Tiwari, P. K., & Bhatnagar, V. (2021). Fog computing in healthcare: A review. *IOP Conference Series: Materials Science and Engineering*, 1099(1), 012025. <https://doi.org/10.1088/1757-899x/1099/1/012025>
- Pighini, C., Vezzoni, A., Mainini, S., Migliavacca, A. G., Montanari, A., Guarneri, M. R., Caiani, E. G., & Cesareo, A. (2022). SynCare: An innovative remote patient monitoring system secured by cryptography and blockchain. In *Privacy and identity management. Between data protection and security* (pp. 73–89). Springer International Publishing. <https://doi.org/10.1007/978-3-030-99100-57>

- Posthuma, L., Downey, C., Visscher, M., Ghazali, D., Joshi, M., Ashrafiyan, H., Khan, S., Darzi, A., Goldstone, J., & Preckel, B. (2020). Remote wireless vital signs monitoring on the ward for early detection of deteriorating patients: A case series. *International Journal of Nursing Studies*, *104*, 103515. <https://doi.org/10.1016/j.ijnurstu.2019.103515>
- Qi, Q., & Tao, F. (2019). A smart manufacturing service system based on edge computing, fog computing, and cloud computing. *IEEE Access*, *7*, 86769–86777. <https://doi.org/10.1109/access.2019.2923610>
- Qi, W., & Aliverti, A. (2020). A multimodal wearable system for continuous and real-time breathing pattern monitoring during daily activity. *IEEE Journal of Biomedical and Health Informatics*, *24*(8), 2199–2207. <https://doi.org/10.1109/jbhi.2019.2963048>
- Ramos, G., Gjini, E., Coelho, L., & Silveira, M. (2021). Unsupervised learning approach for predicting sepsis onset in ICU patients. In *2021 43rd annual international conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE. <https://doi.org/10.1109/embc46164.2021.9629559>
- Raza, A., Tran, K. P., Koehl, L., & Li, S. (2022). Designing ECG monitoring healthcare system with federated transfer learning and explainable AI. *Knowledge-Based Systems*, *236*, 107763. <https://doi.org/10.1016/j.knosys.2021.107763>
- Rohmetra, H., Raghunath, N., Narang, P., Chamola, V., Guizani, M., & Lakkani, N. R. (2021). AI-enabled remote monitoring of vital signs for COVID-19: Methods, prospects and challenges. *Computing*, *29*(3), 1–27. PMID: PMC8006120. <https://doi.org/10.1007/s00607-021-00937-7>
- Sabireen, H., & Venkataraman, N. (2021). A review on fog computing: Architecture, fog with IoT, algorithms and research challenges. *ICT Express*, *7*(2), 162–176. <https://doi.org/10.1016/j.icte.2021.05.004>
- Sagi, O., & Rokach, L. (2020). Explainable decision forest: Transforming a decision forest into an interpretable tree. *Information Fusion*, *61*, 124–138. <https://doi.org/10.1016/j.inffus.2020.03.013>
- Salah, O. Z., Selvaperumal, S. K., & Abdulla, R. (2022). Accelerometer-based elderly fall detection system using edge artificial intelligence architecture. *International Journal of Electrical and Computer Engineering (IJECE)*, *12*(4), 4430. <https://doi.org/10.11591/ijece.v12i4.pp4430-4438>
- Schnyer, D. M., Clasen, P. C., Gonzalez, C., & Beevers, C. G. (2017). Evaluating the diagnostic utility of applying a machine learning algorithm to diffusion tensor MRI measures in individuals with major depressive disorder. *Psychiatry Research: Neuroimaging*, *264*, 1–9. <https://doi.org/10.1016/j.pscychresns.2017.03.003>
- Shaik, T., Tao, X., Higgins, N., Gururajan, R., Li, Y., Zhou, X., & Acharya, U. R. (2022). Fedstack: Personalized activity monitoring using stacked federated learning. *Knowledge-Based Systems*, *257*(12), 109929. <https://doi.org/10.1016/j.knosys.2022.109929>, <https://doi.org/10.1016/j.knosys.2022.109929>
- Shaik, T., Tao, X., Higgins, N., Xie, H., Gururajan, R., & Zhou, X. (2022). AI enabled RPM for mental health facility. In *Proceedings of the 1st ACM Workshop on Mobile and Wireless Sensing for Smart Healthcare* (pp. 26–32). Association for Computing Machinery. <https://doi.org/10.1145/3556551.3561191>
- Shao, M., Zhou, Z., Bin, G., Bai, Y., & Wu, S. (2020). A wearable electrocardiogram telemonitoring system for atrial fibrillation detection. *Sensors*, *20*(3), 606. <https://doi.org/10.3390/s20030606>
- Shapley, L. S. (1953). 17. A value for n-person games. *Contributions to the Theory of Games (Am-28)*, *II*, 307–318. <https://doi.org/10.1515/9781400881970-018>
- Sharma, P., & Kan, E. C. (2018). Sleep scoring with a UHF RFID tag by near field coherent sensing. In *2018 IEEE/MTT-s international microwave symposium IMS*. IEEE. <https://doi.org/10.1109/MWSYM.2018.8439216>
- Shi, H., Wang, H., Qin, C., Zhao, L., & Liu, C. (2020). An incremental learning system for atrial fibrillation detection based on transfer learning and active learning. *Computer Methods and Programs in Biomedicine*, *187*, 105219. <https://doi.org/10.1016/j.cmpb.2019.105219>
- Shouval, R., Hadanny, A., Shlomo, N., Iakobishvili, Z., Unger, R., Zahger, D., Alcalai, R., Atar, S., Gottlieb, S., Matetzky, S., Goldenberg, I., & Beigel, R. (2017). Machine learning for prediction of 30-day mortality after ST elevation myocardial infarction: An acute coronary syndrome israeli survey data mining study. *International Journal of Cardiology*, *246*, 7–13. <https://doi.org/10.1016/j.ijcard.2017.05.067>
- Siam, A. I., Almaiah, M. A., Al-Zahrani, A., Elazm, A. A., Banby, G. M. E., El-Shafai, W., El-Samie, F. E. A., & El-Bahnasawy, N. A. (2021). Secure health monitoring communication systems based on IoT and cloud computing for medical emergency applications. *Computational Intelligence and Neuroscience*, *2021*, 1–23. <https://doi.org/10.1155/2021/8016525>
- Singh, S., Rathore, S., Alfarraj, O., Tolba, A., & Yoon, B. (2021). A framework for privacy-preservation of IoT healthcare data using federated learning and blockchain technology. *Future Generation Computer Systems*, *129*, 380–388. <https://doi.org/10.1016/j.future.2021.11.028>
- Smith, G. B., Recio-Saucedo, A., & Griffiths, P. (2017). The measurement frequency and completeness of vital signs in general hospital wards: An evidence free zone? *International Journal of Nursing Studies*, *74*, A1–A4. <https://doi.org/10.1016/j.ijnurstu.2017.07.001>
- Snoswell, C. L., Chelberg, G., Guzman, K. R. D., Haydon, H. H., Thomas, E. E., Caffery, L. J., & Smith, A. C. (2021). The clinical effectiveness of telehealth: A systematic review of meta-analyses from 2010 to 2019. *Journal of Telemedicine and Telecare*. <https://doi.org/10.1177/1357633x211022907>
- Stein, J.-P., Cimander, P., & Appel, M. (2022). Power-posing robots: The influence of a humanoid robot's posture and size on its perceived dominance, competence, eeriness, and threat. *International Journal of Social Robotics*, *14*(6), 1413–1422. <https://doi.org/10.1007/s12369-022-00878-x>
- Tandel, S., Godbole, P., Malgaonkar, M., Gaikwad, R., & Padaya, R. (2022). *An improved health monitoring system using iot*. SSRN 4109039. <https://doi.org/10.2139/ssrn.4109039>
- Taylor, R. A., Pare, J. R., Venkatesh, A. K., Mowafi, H., Melnick, E. R., Fleischman, W., & Hall, M. K. (2016). Prediction of in-hospital mortality in emergency department patients with sepsis: A local big data-driven, machine learning approach. *Academic Emergency Medicine*, *23*(3), 269–278. <https://doi.org/10.1111/acem.12876>

- Teixeira, P. L., Wei, W.-Q., Cronin, R. M., Mo, H., VanHouten, J. P., Carroll, R. J., LaRose, E., Bastarache, L. A., Rosenbloom, S. T., Edwards, T. L., Roden, D. M., & Denny, J. C. (2016). Evaluating electronic health record data sources and algorithmic approaches to identify hypertensive individuals. *Journal of the American Medical Informatics Association*, 24(1), 162–171. <https://doi.org/10.1093/jamia/ocw071>
- Thieme, A., Belgrave, D., & Doherty, G. (2020). Machine learning in mental health. *ACM Transactions on Computer Human Interaction*, 27(5), 1–53. <https://doi.org/10.1145/3398069>
- Torous, J., Nicholas, J., Larsen, M., Firth, J., & Christensen, H. (2018). Clinical review of user engagement with mental health smartphone apps: Evidence, theory and improvements. *Evidence-Based Mental Health*, 21(3), 116–119. <https://doi.org/10.1136/eb-2018-102891>
- Tripathi, G., Ahad, M. A., & Paiva, S. (2020). SMS: A secure healthcare model for smart cities. *Electronics*, 9(7), 1135. <https://doi.org/10.3390/electronics9071135>
- Uddin, M. Z. (2019). A wearable sensor-based activity prediction system to facilitate edge computing in smart healthcare system. *Journal of Parallel and Distributed Computing*, 123, 46–53. <https://doi.org/10.1016/j.jpdc.2018.08.010>
- Ul Hassan, M., Rehmani, M. H., & Chen, J. (2020). Differential privacy in blockchain technology: A futuristic approach. *Journal of Parallel and Distributed Computing*, 145, 50–74. <https://doi.org/10.1016/j.jpdc.2020.06.003>
- Vinegar, M., & Kwong, M. (2021). Taking score of early warning scores. *University of Western Ontario Medical Journal*, 89(2). <https://doi.org/10.5206/uwomj.v89i2.10518>
- Vimal, S., Robinson, Y. H., Kadry, S., Long, H. V., & Nam, Y. (2021). Iot based smart health monitoring with can using edge computing. *Journal of Internet Technology*, 22(1), 173–185. <https://doi.org/10.3966/160792642021012201017>
- Wang, H., & Yeung, D.-Y. (2016). *A survey on Bayesian deep learning*. <https://arxiv.org/abs/1604.01662>
- Wang, M., Yao, X., & Chen, Y. (2021). An imbalanced-data processing algorithm for the prediction of heart attack in stroke patients. *IEEE Access*, 9, 25394–25404. <https://doi.org/10.1109/access.2021.3057693>
- Wang, H., Zhang, D., Ma, J., Wang, Y., Wang, Y., Wu, D., Gu, T., & Xie, B. (2016). Human respiration detection with commodity wifi devices. In *Proceedings of the 2016 ACM international joint conference on pervasive and ubiquitous computing*. ACM. <https://doi.org/10.1145/2971648.2971744>
- Wang, S., Zhang, C., Kröse, B., & van Hoof, H. (2021). Optimizing adaptive notifications in Mobile health interventions systems: Reinforcement learning from a data-driven behavioral simulator. *Journal of Medical Systems*, 45(12), 1–8. <https://doi.org/10.1007/s10916-021-01773-0>
- Wang, X., Yang, C., & Mao, S. (2017). PhaseBeat: Exploiting CSI phase data for vital sign monitoring with commodity WiFi devices. In *2017 IEEE 37th international conference on distributed computing systems (ICDCS)*. IEEE. <https://doi.org/10.1109/icdcs.2017.206>
- Wang, Y., & Zheng, Y. (2018). Modeling RFID signal reflection for contact-free activity recognition. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 2(4), 1–22. <https://doi.org/10.1145/3287071>
- Wang, Z., Xu, M., Ye, N., Wang, R., & Huang, H. (2019). RF-focus: Computer vision-assisted region-of-interest rfid tag recognition and localization in multipath-prevalent environments. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 3(1), 1–30. <https://doi.org/10.1145/3314416>
- Watts, J., Khojandi, A., Vasudevan, R., & Ramdhani, R. (2020). Optimizing individualized treatment planning for parkinson's disease using deep reinforcement learning. In *2020 42nd annual international conference of the IEEE engineering in medicine & biology society (EMBC)*. IEEE. <https://doi.org/10.1109/embc44109.2020.9175311>
- Weenk, M., Bredie, S., Koeneman, M., Hesselink, G., van Goor, H., & van de Belt, T. H. (2020). Continuous monitoring of vital signs in the general ward using wearable devices: Randomized controlled trial. *Journal of Medical Internet Research*, 22(6), e15471. <https://doi.org/10.2196/15471>
- Wu, Q., Chen, X., Zhou, Z., & Zhang, J. (2022). FedHome: Cloud-edge based personalized federated learning for in-home health monitoring. *IEEE Transactions on Mobile Computing*, 21(8), 2818–2832. <https://doi.org/10.1109/tmc.2020.3045266>
- Xia, J., Pan, S., Zhu, M., Cai, G., Yan, M., Su, Q., Yan, J., & Ning, G. (2019). A long short-term memory ensemble approach for improving the outcome prediction in intensive care unit. *Computational and Mathematical Methods in Medicine*, 2019, 1–10. <https://doi.org/10.1155/2019/8152713>
- Xu, Q., Liu, X., Luo, J., & Tang, Z. (2020). Emotion monitoring with RFID: An experimental study. *CCF Transactions on Pervasive Computing and Interaction*, 2(4), 299–313. <https://doi.org/10.1007/s42486-020-00043-1>
- Xu, Y., Biswal, S., Deshpande, S. R., Maher, K. O., & Sun, J. (2018). RAIM. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. ACM. <https://doi.org/10.1145/3219819.3220051>
- Xue, M., Su, Y., Li, C., Wang, S., & Yao, H. (2020). Identification of potential type II diabetes in a large-scale chinese population using a systematic machine learning framework. *Journal of Diabetes Research*, 2020, 1–12. <https://doi.org/10.1155/2020/6873891>
- Yang, G., Ye, Q., & Xia, J. (2022). Unbox the black-box for the medical explainable AI via multi-modal and multi-centre data fusion: A mini-review, two showcases and beyond. *Information Fusion*, 77, 29–52. <https://doi.org/10.1016/j.inffus.2021.07.016>
- Yang, J., Xiao, W., Lu, H., & Barnawi, A. (2020). Wireless high-frequency NLOS monitoring system for heart disease combined with hospital and home. *Future Generation Computer Systems*, 110, 772–780. <https://doi.org/10.1016/j.future.2019.11.001>
- Yew, H. T., Ng, M. F., Ping, S. Z., Chung, S. K., Chekima, A., & Dargham, J. A. (2020). IoT based real-time remote patient monitoring system. In *2020 16th IEEE international colloquium on signal processing & its applications (CSPA)*. IEEE. <https://doi.org/10.1109/CSPA48992.2020.9068699>
- Yu, C., Liu, J., Nemati, S., & Yin, G. (2023). Reinforcement learning in healthcare: A survey. *ACM Computing Surveys*, 55(1), 1–36. <https://doi.org/10.1145/3477600>

- Yu, S., Chai, Y., Chen, H., Brown, R. A., Sherman, S. J., & Nunamaker, J. F. (2021). Fall detection with wearable sensors: A hierarchical attention-based convolutional neural network approach. *Journal of Management Information Systems*, 38(4), 1095–1121. <https://doi.org/10.1080/07421222.2021.1990617>
- Zainuddin, A. A., Superamiam, S., Andrew, A. C., Muraleedharan, R., Rakshys, J., Miriam, J., Bostomi, M. A., Rais, A. M., Khalidin, Z., Mansor, A. F., & Taufik, M. S. M. (2020). Patient monitoring system using computer vision for emotional recognition and vital signs detection. In *2020 IEEE student conference on research and development (SCOREd)*. IEEE. <https://doi.org/10.1109/scored50371.2020.9250950>
- Zamanifar, A. (2021). Remote patient monitoring: Health status detection and prediction in IoT-based health care. In *IoT in healthcare and ambient assisted living* (pp. 89–102). Springer Singapore. <https://doi.org/10.1007/978-981-15-9897-55>
- Zhang, D., Yao, L., Chen, K., Long, G., & Wang, S. (2019). Collective protection: Preventing sensitive inferences via integrative transformation. In *2019 IEEE international conference on data mining (ICDM)*. IEEE. <https://doi.org/10.1109/icdm.2019.00197>
- Zhang, X., Kim, J., Patzer, R. E., Pitts, S. R., Patzer, A., & Schrager, J. D. (2017). Prediction of emergency department hospital admission based on natural language processing and neural networks. *Methods of Information in Medicine*, 56(5), 377–389. <https://doi.org/10.3414/me17-01-0024>
- Zhao, R., Wang, D., Zhang, Q., Chen, H., & Huang, A. (2018). CRH: A contactless respiration and heartbeat monitoring system with COTS RFID tags. In *2018 15th annual IEEE international conference on sensing, communication, and networking (SECON)*. IEEE. <https://doi.org/10.1109/sahcn.2018.8397132>
- Zheng, X., Shah, S. B. H., Ren, X., Li, F., Nawaf, L., Chakraborty, C., & Fayaz, M. (2021). Mobile edge computing enabled efficient communication based on federated learning in internet of medical things. *Wireless Communications and Mobile Computing*, 2021, 1–10. <https://doi.org/10.1155/2021/4410894>
- Zhong, G., Wang, L.-N., Ling, X., & Dong, J. (2016). An overview on data representation learning: From traditional feature learning to recent deep learning. *The Journal of Finance and Data Science*, 2(4), 265–278. <https://doi.org/10.1016/j.jfds.2017.05.001>
- Zhu, L., Wang, R., Wang, Z., & Yang, H. (2017). TagCare: Using RFIDs to monitor the status of the elderly living alone. *IEEE Access*, 5, 11364–11373. <https://doi.org/10.1109/access.2017.2716359>

How to cite this article: Shaik, T., Tao, X., Higgins, N., Li, L., Gururajan, R., Zhou, X., & Acharya, U. R. (2023). Remote patient monitoring using artificial intelligence: Current state, applications, and challenges. *WIREs Data Mining and Knowledge Discovery*, 13(2), e1485. <https://doi.org/10.1002/widm.1485>

2.2 Summary

Chapter 2 provides an in-depth analysis of remote patient monitoring, highlighting its significance and presenting the key challenges that drive the subsequent research endeavours. It sets the stage for the exploration of cutting-edge AI-driven approaches to patient care. The summary bridges the discussion on AI's integration into RPM with a forward-looking perspective on its symbiotic relationship with emerging healthcare technologies, as explored in later chapters. It underscores the challenges and opportunities that AI introduces to RPM, such as ethical data usage and the democratization of healthcare, hinting at the broader thematic explorations of AI ethics, data governance, and the evolution of patient-centric care models in the continuum of the thesis. This cohesive wrap-up not only reflects on the chapter's contributions but also primes the reader for the interconnected discussions that follow, emphasizing the thesis's collective advancement of knowledge in digital health.

CHAPTER 3: PAPER 2 - FEDSTACK: PERSONALIZED ACTIVITY MONITORING USING STACKED FEDERATED LEARNING

3.1 Introduction

This chapter delves into the innovative FedStack framework, a novel approach to personalized activity monitoring within the realm of remote patient monitoring (RPM). It explores how FedStack leverages stacked federated learning to address the challenges of data privacy and model heterogeneity in healthcare applications. By enabling decentralized AI model training across diverse devices, FedStack enhances patient privacy and supports the personalized monitoring essential for individualized healthcare. This introduction outlines the integration of various AI models, including Artificial Neural Networks (ANNs), Convolutional Neural Networks (CNNs), and Bidirectional Long Short-Term Memory (Bi-LSTM) networks, to achieve state-of-the-art performance in activity recognition. The chapter sets the context for understanding how FedStack contributes to the broader thesis by enhancing RPM's efficacy and privacy.

This article cannot be displayed due to copyright restrictions. See the article link in the Related Outputs field on the item record for possible access.

3.2 Summary

In summary, this chapter presents the FedStack framework as a significant advancement in personalized activity monitoring through federated learning. By demonstrating superior performance in recognizing diverse human activities with enhanced privacy, FedStack embodies a significant leap forward in RPM technology. The chapter underscores the successful application of FedStack to a mobile health sensor benchmark dataset, showcasing its effectiveness in leveraging heterogeneous AI models for improved activity recognition. Furthermore, it highlights the optimal sensor placement for accurate data collection, emphasizing the practical implications of FedStack in clinical settings. This chapter's findings not only contribute to the field of RPM but also bridge to subsequent chapters, where the focus shifts to scaling these technologies for broader healthcare applications, maintaining patient privacy, and addressing the challenges of deploying AI-driven health monitoring systems in real-world scenarios.

CHAPTER 4: PAPER 3 - CLUSTERED FEDSTACK: INTERMEDIATE GLOBAL MODELS WITH BAYESIAN INFORMATION CRITERION

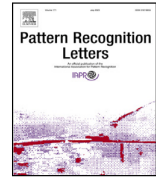
4.1 Introduction

This chapter introduces "Clustered FedStack," an advanced extension of the FedStack framework, designed to address the challenges posed by non-identically and independently distributed (non-IID) data in federated learning environments. By integrating sophisticated clustering techniques and the Bayesian Information Criterion (BIC), this framework innovatively groups local client models based on their output layer weights to form intermediate global models. This approach not only enhances the personalization of federated learning models but also significantly improves their performance in diverse applications, from human activity recognition (HAR) to natural language processing (NLP) tasks. This introduction sets the stage for a detailed exploration of the clustered FedStack framework, its algorithmic underpinnings, and its empirical validation through rigorous experiments.



Contents lists available at ScienceDirect

Pattern Recognition Letters

journal homepage: www.elsevier.com/locate/patrec

Clustered FedStack: Intermediate Global Models with Bayesian Information Criterion

Thanveer Shaik ^{a,*}, Xiaohui Tao ^a, Lin Li ^b, Niall Higgins ^{c,d}, Raj Gururajan ^e, Xujuan Zhou ^e, Jianming Yong ^e

^a School of Mathematics, Physics and Computing, University of Southern Queensland, Toowoomba, Australia

^b Wuhan University of Technology, Wuhan, China

^c Royal Brisbane and Women's Hospital, Brisbane, Australia

^d Queensland University of Technology, Brisbane, Australia

^e School of Business, University of Southern Queensland, Springfield, Australia

ARTICLE INFO

Editor: Li Liu

Dataset link: <https://archive.ics.uci.edu/ml/datasets/PPG-DaLiA>

Keywords:

Federated learning
FedStack
Clustering
Bayesian
Cyclical learning rates

ABSTRACT

Federated Learning (FL) is currently one of the most popular technologies in the field of Artificial Intelligence (AI) due to its collaborative learning and ability to preserve client privacy. However, it faces challenges such as non-identically and non-independently distributed (non-IID) data with imbalanced labels among local clients. To address these limitations, the research community has explored various approaches such as using local model parameters, federated generative adversarial learning, and federated representation learning. In our study, we propose a novel Clustered FedStack framework based on the previously published Stacked Federated Learning (FedStack) framework. Here, the local clients send their model predictions and output layer weights to a server, which then builds a robust global model. This global model clusters the local clients based on their output layer weights using a clustering mechanism. We adopt three clustering mechanisms, namely K-Means, Agglomerative, and Gaussian Mixture Models, into the framework and evaluate their performance. Bayesian Information Criterion (BIC) is used with the maximum likelihood function to determine the number of clusters. Our results show that Clustered FedStack models outperform baseline models with clustering mechanisms. To estimate the convergence of our proposed framework, we use Cyclical learning rates.

1. Introduction

As AI techniques have matured, a vast amount of human data is being generated every second around the world. To manage this huge data, technology giant Google introduced a mechanism that trains a machine learning (ML) algorithm across multiple decentralized devices or servers without exchanging their local data samples. This is called Federated Learning (FL), which is also known as collaborative learning [1]. FL overcomes the issues of data privacy that exist in traditional centralized learning techniques where all device or server data is merged for analysis [2]. FL has garnered significant attention since its introduction by Google as a ML technique for predicting users' input from Gboard (a keypad) on Android devices. This technique has been widely adopted in communication, engineering, and healthcare. However, medical institutes in particular possess a vast amount of patient data that may not be sufficient to train ML or deep learning models, and may even be biased due to a lack of data diversity. FL

addresses this issue through its collaborative learning approach, where local models trained in each medical institute share their model weights with a global model stored in a shared server [3]. This maintains data privacy, as the institute's data remains within its premises. The process can be used at the patient level to monitor their health status by predicting vital signs, such as heart rate and breathing, and classifying their physical activities. It enables personalized patient monitoring with enhanced data privacy.

A heterogeneous stacked FL, FedStack, was proposed by Shaik et al. [4] to overcome the problems of the traditional FL approach, while enabling personalized monitoring of patients' physical actions. The authors achieved state-of-the-art performances using different deep learning models as part of local and global clients. The FedStack approach is confined to building the global model by stacking local clients' predictions heterogeneously and allowing local clients to have

* Corresponding author.

E-mail addresses: Thanveer.Shaik@usq.edu.au (T. Shaik), Xiaohui.Tao@usq.edu.au (X. Tao), cathylilin@whut.edu.cn (L. Li), Niall.Higgins@health.qld.gov.au (N. Higgins), Raj.Gururajan@usq.edu.au (R. Gururajan), Xujuan.Zhou@usq.edu.au (X. Zhou), Jianming.Yong@usq.edu.au (J. Yong).

<https://doi.org/10.1016/j.patrec.2023.12.004>

Received 16 March 2023; Received in revised form 21 October 2023; Accepted 11 December 2023

Available online 14 December 2023

0167-8655/© 2023 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

different architectural models. However, it has a limitation of non-identically and independently distributed (non-IID) data, where the local clients' data distributions may be different. This can be addressed by allowing the global model to group the local clients based on their deep learning model output weights. To avoid any bias in grouping the local clients, unsupervised clustering methods can be adopted.

This study proposes a novel Clustered-FedStack framework to overcome FL's non-IID data challenge [5]. All models trained on local clients pass their predictions and output layer weights to the server, which builds a global server model based on the predictions received from the local models. Later, the global server model clusters local client models with output layer weights received and creates intermediate clustered models between local clients and the server. In this unsupervised process, the server model computes the cosine distance matrix among the local model output layer weights. To determine the number of clusters in this process, the BIC technique is adopted and maximum likelihood estimation is applied to the local model weights in the server. Three types of clustering techniques: centroid-based (k-Means), hierarchical (Agglomerative), and distribution-based (Gaussian Mixture Model) techniques are deployed. Cyclical learning rates are applied to estimate the convergence of the clustered models.

The proposed framework is evaluated with a human activity recognition (HAR) task using the publicly available sensor-based PPG-DALiA dataset [6]. The results show that clustered models have state-of-the-art performance in classifying human activities with the sensor data of 15 subjects. The performance of the clustered FedStack model is compared with four clustered FL baseline models, and the proposed model has outperformed the baseline models in all classification metrics. Moreover, the proposed framework can be scalable to Natural Language Processing (NLP) tasks. This has been evaluated on the drug review dataset [7], where the intermediate clustered models performed better and could handle a huge number of local clients with non-IID data to achieve superconvergence. Thus, the proposed clustered FedStack framework can group local clients and overcome the non-IID challenge in FL. The contributions of the present study include the following:

- A novel Clustered-FedStack framework is proposed to group local clients in an unsupervised approach and overcome the non-IID challenge in FL.
- Improved personalized modeling in FL by building intermediate clustered models between the global server model and local clients.
- Achievement of superconvergence of all clustered-FedStack models using Cyclical learning rates.
- A Clustered-FedStack approach that proves scalable for Natural Language Processing (NLP) tasks, effectively handling a high number of local clients with non-IID data.

Section 2 presents related works on FL and different aggregating techniques developed. Section 3 presents the formulation of the research problem and the proposed Clustered-FedStack framework. In Section 4, the proposed framework is evaluated in HAR and the results are discussed. The framework optimization with Cyclical learning rates is also presented in Section 4. In Section 5, we evaluate the scalability of the proposed framework using a NLP dataset. Section 6 concludes the paper.

2. Related works

Numerous studies have explored the aggregation of local model parameters in FL and passed them to the global model on the server. One of the first proposed aggregation techniques in FL is the Federated Averaging (FedAvg) algorithm, which uses the average function to aggregate local model weights and generate new weights to feed to the global model [8]. However, the FedAvg technique cannot optimize models if a client has a heterogeneous data distribution. To combat this,

Arivazhagan et al. [9] proposed FedPer, which has two layers: a base layer and a personalization layer. FedAvg trains the base layers, while the personalization layers are trained with stochastic gradient descent, helping to mitigate the ill effects of statistical heterogeneity. Wang et al. [10] proposed Federated Matched Averaging (FedMA), which is a layer-wise approach that matches and merges nodes with the same weights, trains them independently, and communicates the layers to the global model.

Osmani et al. [11] proposed a multi-level FL system for HAR, which includes a reconciliation step based on FL aggregation techniques such as FedAvg or Federated Normalized Averaging. Xiao et al. [12] proposed another FL system for HAR with enhanced feature extractions. They designed a Perceptive Extraction Network (PEN) with two networks: a featured network based on the convolutional block to extract local features, and a relation network based on Long Short-Term Memory (LSTM) and an attention mechanism to mine global relationships. Pang et al. [13] proposed a rule-based collaborative framework (CloREF) that allows local clients to use heterogeneous local models. Tian et al. [14] discussed the limitations of traditional FL methods in heterogeneous IoT systems and proposed a novel Weight Similarity-based Client Clustering (WSCC) approach to address the non-IID challenge in FL. The WSCC approach involves splitting clients into different groups based on their data set distributions using an affinity-propagation-based method. Their proposed approach outperformed existing FL schemes under different non-IID settings, achieving up to 20% improvements in accuracy without requiring extra data transmissions or additional models.

Federated Learning in HAR The increasing use of electronic assistive health applications such as smartwatches and activity trackers has led to the emergence of pervasive or ubiquitous computing, where devices can seamlessly exchange data with each other [15]. Although this has the advantage of real-time tracking of human health changes, it is vulnerable to security breaches that compromise data privacy [16]. The advancement of AI techniques as a whole is contributing to the massive amount of human data being generated worldwide every second. To handle such enormous data, Google introduced FL, which trains a ML algorithm across decentralized devices or servers without exchanging their local data samples. FL overcomes the data privacy issues associated with traditional centralized learning techniques, where all device and/or server data is merged for analysis [1]. Sannara et al. [17] evaluated the performance of FL aggregation techniques like Federated Averaging (FedAvg), Federated Learning with Matched Averaging (FedMa), and Federated Personalization Layer (FedPer) against centralized training techniques. They used the CNN model to classify eight physical activities. Zhao et al. [18] designed an activity recognition system based on semi-supervised FL. Ouyang et al. [19] proposed the ClusterFL approach, which exploits the similarity of users' data to minimize the empirical loss of trained models. This improved Federated model accuracy and communication efficiency between local models and global models.

Local clients may have different data distributions, demographics, and model architectures. Passing all the local clients' parameters to build a robust global server model poses challenges such as label imbalance and non-IID. To identify hidden patterns or relationships among the local clients and overcome these challenges, unsupervised clustering techniques can be adopted to improve personalized learning in FL. This study proposes a clustered FL framework to overcome these identified challenges.

3. Methodology

To accommodate heterogeneous architectural models for local clients, we adopt the previously published FedStack framework by Shaik et al. [4]. This study extends the FedStack framework to the clustered-FedStack framework, facilitating the creation of heterogeneous multi-global FL models by clustering individual subjects with local models.

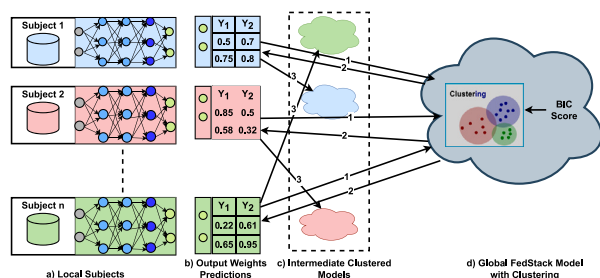


Fig. 1. Clustered FedStack model.

3.1. Research problem

In this study, the research problem is to overcome the non-IID data challenge in a FL environment. Let $S = \{s_1, s_2, \dots, s_N\}$ be the set of subjects, where the data is non-IID. The objective is to divide subjects S into M clusters $C = \{c_1, c_2, \dots, c_M\}$, where each cluster c_m is a subset of subjects S , $c_m \subseteq S$. For each cluster c_m , there exists a local model l_m that can be heterogeneous according to the subject's convenience. The predictions p_m from local models and their corresponding output layer weights are passed to a global model server g . The training process for the global model g is shown in Eq. (1).

$$\text{train}(g) \leftarrow \sum_{m=1}^M c_m \leftarrow \sum_{m=1}^M l_m(p_m) \quad (1)$$

where: $\text{train}(g)$ refers to the training process for the Global Model g using local model predictions of cluster c_m , and $l_m(p_m)$ represents the local model l_m and its predictions p_m for each subject in the cluster c_m .

3.2. Clustered-FedStack framework

In the Clustered-FedStack framework, local clients train their models on private data and then forward their model predictions p and output layer weights Q to the global server model g for training, as shown in Fig. 1. The figure's arrow numbers indicate the framework execution order. After receiving the local model predictions and output layer weights, the global server model determines the number of clusters using the BIC score. It then clusters the local clients based on their output layer weights. For each label i in local model training, an output neuron without a successor is configured to gather the computed and accumulated values from the local model's input and hidden layers. The output neuron value q_i is calculated using Eq. (2), with inputs x_i , weights w_i , and bias b for a local model l_n . By computing all output neuron values, the local model l_n predictions p can be estimated using Eq. (3).

$$q_i = l_n(b, x_i, w_i) \quad (2)$$

$$p = l_n(b + \sum_{i=1}^n x_i \cdot w_i) \quad (3)$$

Output neuron values for each local model l_n are consolidated into a single set Q using Eq. (4). This procedure is repeated for all local models based on their output layer values, forming a large set Q as defined in Eq. (5).

$$Q = \{q_1, q_2, q_3, \dots, q_n\} \quad (4)$$

$$Q = \{Q_{l_1}, Q_{l_2}, \dots, Q_{l_n}\} \quad (5)$$

3.2.1. Clustering technique

Given the set Q from Eq. (5), where each element of the set represents the values of a local model's l_n output layer, the goal is to divide Q into k clusters, where $k \leq n$, represented by $C = \{c_1, c_2, \dots, c_k\}$. There are various techniques that can be applied to clustering, including centroid-based, hierarchical, and distribution-based methods. The general objective of these methods is to minimize the within-cluster sum of squared differences or a related measure of dissimilarity, as described in Eq. (6). The notation $\arg \min_C$ refers to finding the set of clusters C that minimizes the following expression, where the "arg min" stands for the argument of the minimum, i.e., the specific value of the variable that results in the lowest possible value of the given function.

$$\arg \min_C \left(\sum_{i=1}^k \sum_{x \in C_i} \|x - c_i\|^2 \right) \quad (6)$$

Here, C represents the set of clusters, x is a data point, and c_i is the representative point, such as a centroid. The term $\|\cdot\|$ represents a distance measure.

Cosine similarity is utilized to assign each local model's output neuron set to a specific cluster, considering the angle between output neuron sets of two local model l_n as Q_{l_1} and Q_{l_2} , the cosine similarity can be estimated using Eq. (7).

$$S_C(Q_{l_1}, Q_{l_2}) = \frac{Q_{l_1} \cdot Q_{l_2}}{\|Q_{l_1}\| \|Q_{l_2}\|} \quad (7)$$

3.2.2. Bayesian information criterion

The proposed Clustered-FedStack technique enables the global server model to access local models' predictions and layer weights. However, using an unsupervised method to determine the number of clusters in local models is challenging. The BIC technique is utilized to overcome this. BIC calculates its value based on a clustering model \mathcal{M} 's maximum likelihood function M_L , representing the probability that the layer weights data fits the clustering model [20]. This is shown in Eq. (8). BIC values balance the maximum likelihood estimation against the number of model parameters m_p , seeking a model with the fewest parameters that can accurately explain the data clusters, as in Eq. (9).

$$M_L(\mathcal{M}) = -2 \ln(\mathcal{L}) + m_p \ln(n) \quad (8)$$

$$BIC = -2 \ln(\mathcal{L}) + m_p \ln(n) = M_L(\mathcal{M}) \quad (9)$$

The BIC values for each clustering model are compared with the minimum BIC value indicating the optimal clustering model. This process ensures that the global model converges by configuring a suitable number of clusters for local models, resulting in a consolidated global model that represents heterogeneous subject models.

3.3. Clustered-FedStack algorithm

Algorithm 1 presents the proposed Clustered-FedStack process in detail. Line 1 initializes empty sets to collect output layer weights and clustered models, and datasets D and D' for evaluating the global server model. Lines 2–7 detail the FedStack process, where local client model predictions and weights are passed to the global server model g for training and testing. Lines 8–10 present the iteration through all local model weights in g to collect their output layer weights. Lines 11–12 detail the determination of the number of clusters to be formed from the weights W set. Line 13 computes the cosine distance among all the local model weights collected. Lines 14–19 explain the clustering process for all the local models, based on Lines 11–13.

Algorithm 1 Proposed Clustered-FedStack Algorithm**Require:**

Subjects set $S = \{s_1, s_2, \dots, s_n\}$
 Local AI models $M = \{m_1, m_2, \dots, m_m\}$
 Labels set $K = \{1, 2, \dots, k\}$
 Global Server Model g

Ensure: Classification probabilities of labels K for each intermediate cluster model C

1: Initialization:

D : Dataset for training
 D' : Unseen Dataset for testing
 $W = \emptyset$: Set to collect weights
 $CM = \emptyset$: Set for clustered models

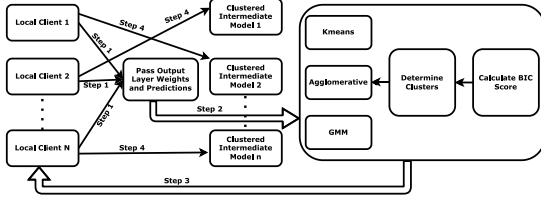
2: $stack = \{(m_1^K, m_1^W), (m_2^K, m_2^W), \dots, (m_m^K, m_m^W)\}$; ▷ Predictions and weights of local AI models3: **for** $m \in M$ **do**4: $g^{train} \leftarrow stack$;5: $g^{test} \leftarrow D'$;6: **end for**7: **for** $m \in G(M)$ **do**8: Collect weights of m : $W \leftarrow \{m, w\}$;9: **end for**10: **Determine Clusters:**11: Compute BIC scores of $CM \geq M$;12: $CM \leftarrow \min(\text{BIC})$;13: Compute cosine distance among $\{(m_1, w_1), (m_2, w_2), \dots, (m_m, w_m)\}$;14: **Assignment:**15: **for** c in C **do**16: $c \leftarrow \arg \min (\sum_{i=1}^k \sum_{x \in c_i} \|x - c_i\|^2)$;17: $CM \leftarrow c$;18: **end for**19: **Return** CM ;

Fig. 2. Experimental design of the proposed framework.

4. Experiments on clustered FedStack in HAR

Conventional FL methods assume that the data distribution is consistent among all clients [21]. However, this assumption may not be valid in FL, as data heterogeneity can be present [22]. This limitation forces clients to have identical data distribution and architectural models to build global models. FedStack [4] addressed the issue of identical architectural models in FL. The goal of this study is to extend the FedStack framework by introducing intermediate clustered models to address the non-IID challenge in FL.

In this study, the objective is to overcome the non-IID challenge in FL. To achieve this, the proposed Clustered-FedStack algorithm is applied to the domain of human activity recognition, where patients' physical activity is classified. The non-IID data distribution of the dataset used in the experiment is presented. The proposed methodology involves passing the output layer weights and predictions of local clients to the global model, which then calculates unsupervised clustering of the local model layer weights to group the local clients and establish clustered intermediate models. The experimental design is presented in Fig. 2. The evaluation results compare the performance of the proposed framework to the baseline models and show clustering results leading to clustered FedStack models. Furthermore, the convergence of the clustered FedStack models is analyzed using Cyclical learning rates.

4.1. Dataset

The proposed Clustered-FedStack algorithm was evaluated on the HAR problem, which involves classifying patients' physical activity.

Table 1

Non-IID data.

Local clients	Distribution	1	2	3	4	5	6	7	8
Subject 1	27 724	2800	1148	1380	1648	3556	9420	3016	4756
Subject 2	22 712	2400	1068	1216	1548	3680	4880	2756	5164
Subject 3	26 900	2400	1740	1172	1516	3640	8640	2952	4840
Subject 4	26 528	2280	2092	1312	1900	4028	7580	2376	4960
Subject 5	26 924	2400	1860	1160	1728	3320	9020	2356	5080
Subject 6	11 812	2532	1720	1236	2132	4192	9020	0	0
Subject 7	28 580	2472	1624	1096	2012	4140	9700	2836	4700
Subject 8	23 992	2400	1648	1292	1680	3080	7200	1924	4768
Subject 9	26 212	2400	1932	1140	2216	3820	7368	2356	4980
Subject 10	28 424	2392	1868	1220	1952	3748	8336	4328	4580
Subject 11	28 052	2400	1828	1296	1960	3440	9632	2616	4880
Subject 12	23 680	2408	1936	1120	1920	3560	5840	2116	4780
Subject 13	26 996	2420	1988	1160	1992	3588	8112	2836	4900
Subject 14	25 584	2432	1824	1300	2008	3816	6924	2460	4820
Subject 15	23 504	2444	1676	1416	1620	3140	5760	2636	4812

The PPG-DALiA [6] dataset, which is publicly accessible and cited in [6], was utilized for this study. This dataset includes physiological and motion data gathered from 15 participants as they engaged in a diverse array of activities, closely mirroring real-life conditions. The data was collected from both a wrist-worn (Empatica E4) and a chest-worn (RespiBAN) device, and includes 11 attributes such as 3-Dimensional (3D) acceleration data, Electrocardiogram (ECG), respiration, Blood Volume Pulse (BVP), Electrothermal Activity (EDA), and body temperature. The 3D acceleration data was labeled with eight different physical activities.

4.2. Non-IID data

Table 1 shows the distribution and activity of local clients in a FL scenario with non-IID data. Each row represents a client, and each column represents a feature. The "Distribution" column shows the number of data points available at each client, which varies across clients, indicating non-IID in the dataset. The remaining columns represent different activities that are each related to the type of data collected or the task being performed. For instance, "Activity 1" to "Activity 8" could be different types of sensor readings or behavioral data collected from different sources. The non-IID nature of this data could potentially impact the performance of the FL algorithm since the data distribution across clients is not uniform, and the model may not generalize well to all clients. Therefore, special attention must be given to handling the non-IID data in FL, by using the technique of personalized FL to improve model performance for each client's unique data distribution.

4.3. Data modeling

In data modeling, three AI models were chosen: Artificial Neural Networks (ANN), Convolutional Neural Networks (CNN), and Bidirectional Long Short-Term Memory (BiLSTM) models, due to their state-of-the-art performances in FL works [12] and activity classification [23]. Each subject trained with one of the chosen models locally and passed their predictions and local model output layer weights to the global server model. The proposed framework clustered the global model without any private information about local clients, based on the output layer weights.

4.4. Baseline models

- **ClusterFL** [19]: A clustering-based FL system for the HAR application. The ClusterFL approach captures the intrinsic clustering relation among local clients and minimizes the training loss.
- **FL+HC** [24]: A hierarchical clustered FL system to separate clusters of clients based on the similarity of their local updates to the global server model.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12	S13	S14	S15
S1	0	0.18	0.21	0.29	0.28	0.35	0.21	0.39	0.2	0.2	0.38	0.52	0.15	0.27	0.14
S2	0.18	0	0.02	0.06	0.1	0.11	0.26	0.11	0.06	0.19	0.21	0.25	0.12	0.06	0.05
S3	0.21	0.02	0	0.06	0.09	0.19	0.28	0.13	0.04	0.24	0.11	0.29	0.08	0.03	0.04
S4	0.29	0.06	0.06	0	0.06	0.15	0.23	0.12	0.06	0.18	0.2	0.26	0.11	0.02	0.13
S5	0.28	0.1	0.09	0.06	0	0.16	0.11	0.15	0.06	0.28	0.16	0.23	0.08	0.05	0.15
S6	0.35	0.11	0.19	0.15	0.16	0	0.33	0.21	0.23	0.37	0.43	0.26	0.31	0.17	0.2
S7	0.21	0.26	0.28	0.23	0.11	0.33	0	0.28	0.18	0.25	0.34	0.33	0.15	0.24	0.3
S8	0.39	0.11	0.13	0.12	0.15	0.21	0.28	0	0.09	0.2	0.26	0.07	0.25	0.15	0.22
S9	0.2	0.06	0.04	0.06	0.06	0.23	0.18	0.09	0	0.16	0.14	0.2	0.06	0.05	0.07
S10	0.2	0.19	0.24	0.18	0.28	0.37	0.25	0.2	0.16	0	0.45	0.32	0.24	0.25	0.25
S11	0.38	0.21	0.11	0.2	0.16	0.43	0.34	0.26	0.14	0.45	0	0.47	0.15	0.11	0.21
S12	0.52	0.25	0.29	0.26	0.23	0.26	0.33	0.07	0.2	0.32	0.47	0	0.39	0.32	0.34
S13	0.15	0.12	0.08	0.11	0.08	0.31	0.15	0.25	0.06	0.24	0.15	0.39	0	0.08	0.08
S14	0.27	0.06	0.03	0.02	0.05	0.17	0.24	0.15	0.05	0.25	0.11	0.32	0.08	0	0.09
S15	0.14	0.05	0.04	0.13	0.15	0.2	0.3	0.22	0.07	0.25	0.21	0.34	0.08	0.09	0

Fig. 3. Cosine distance among local clients.

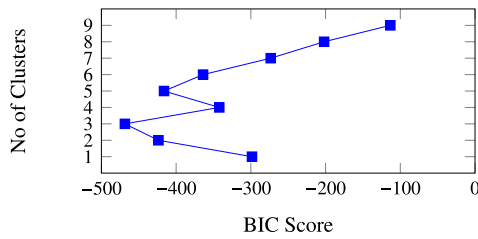


Fig. 4. BIC score to determine the number of clusters.

- **HypCluster** [25]: A hypothesis-based clustering with a stochastic Expectation-Maximization (EM) algorithm adopted for the FL approach, where local clients partition into a certain number of clusters and then the model finds the best hypothesis for each cluster.
- **Dynamic Clustering** [26]: A three-phased data clustering algorithm, namely, generative adversarial network-based clustering, cluster calibration, and cluster division, designed to overcome the fixed shape of clusters, data privacy breaches, and non-adaptive numbers of clusters.

4.5. Results analysis

4.5.1. Clustering results

Before clustering, the cosine distance among all 15 local models trained on clients is calculated to check their similarity in terms of the models' output layer weights, as shown in Fig. 3. The matrix heatmap ranges on a scale from 0 to 0.6 where 0 shows no cosine distance between the client output layer values, and 0.6 shows the maximum cosine distance.

The proposed Clustered-FedStack algorithm employed the BIC approach to calculate the maximum likelihood function on the output layer weights received from the local client models by the global server model, as shown in Fig. 1. This process determines the number of clusters among the 15 local clients. Fig. 4 shows that the lowest BIC score corresponds to three clusters in the global server model. After determining the clusters, three clustering techniques were applied: centroid-based clustering (K-Means) [27], hierarchical clustering (Agglomerative) [28], and distribution-based clustering (Gaussian Mixture Model (GMM)) [29]. Fig. 5 shows that K-Means and Agglomerative clustering produced similar groups of local client models, while GMM clustering grouped all CNN models into the second cluster and distributed other ANN and BiLSTM models in the first and third clusters.

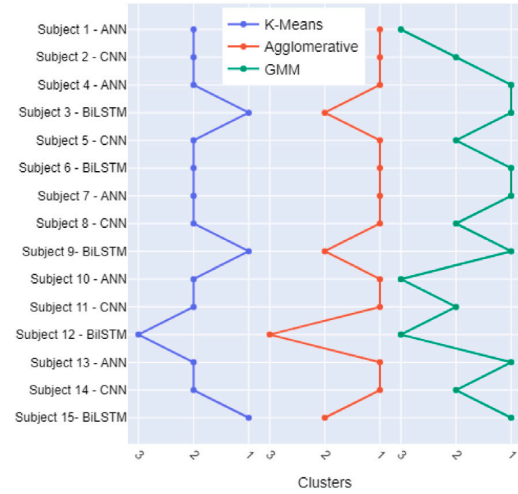


Fig. 5. Clustering results.

4.5.2. Clustered FedStack model performances

After determining the clusters, each cluster of local clients passes their output layer weights to an intermediate Clustered-FedStack model, situated between the local clients and the global server model, as shown in Fig. 1. This approach reduces the load on the global server model and groups similar local models for more efficient AI results. The three clustering techniques generate three Clustered-FedStack models each, and their performance in HAR is shown in Table 2. All nine Clustered-FedStack intermediate global models generated from the clustering techniques have performed well in the HAR task. K-Means and agglomerative clustering, having similar clustering results, showed similar classification accuracy in HAR. While comparing the results, the GMM Clustered FedStack models, which are distribution-based, exhibited slightly better accuracy than the other two clustered models.

4.5.3. Baseline models comparison

The proposed framework was compared against four other baseline models in FL approaches with clustering. All models were trained using 3D acceleration data for HAR tasks, and their evaluation results are presented in Table 3. As K-Means and hierarchical clustering techniques generate similar clusters from the 15 local client models, the table shows three clustered models (Clustered FedStack 1, Clustered FedStack 2, Clustered FedStack 3) built based on K-Means and hierarchical clustering, and three clustered models (Clustered FedStack 7, Clustered FedStack 8, Clustered FedStack 9) built based on the GMM model. The Table presents the mean of four metrics: balanced accuracy, precision, recall, and F1-score in classifying eight activities for six intermediate clustered models. The proposed approach outperformed all other baseline models in terms of all the metrics.

4.6. Convergence analysis

The optimization of the proposed Clustered FedStack framework is estimated using Cyclical learning rates [30] for convergence. The performance of the intermediate Clustered FedStack models shown in Fig. 1 is optimized using the Learning Rate (α) of the deep learning models. In the Cyclical learning rates process, the α values are cycled with an initial learning rate of 0.00001 and a maximum learning rate of 0.001, and stochastic gradient descent is performed. A scale function is defined to control the change from the initial learning rate to the maximal learning rate and back to the initial learning rate. The

Table 2
Clustered FedStack model accuracy in HAR.

Activity	K-Means clusters			Agglomerative clusters			GMM clusters		
	Cluster 1	Cluster 2	Cluster 3	Cluster 1	Cluster 2	Cluster 3	Cluster 1	Cluster 2	Cluster 3
Sitting	0.99	0.96	0.95	0.97	0.99	0.95	0.99	0.99	0.99
Ascending and descending stairs	0.92	0.96	0.94	0.96	0.92	0.94	0.92	0.92	0.92
Table soccer	0.96	0.95	0.95	0.97	0.95	0.95	0.95	0.96	0.97
Cycling	0.94	0.97	0.98	0.95	0.96	0.98	0.96	0.93	0.96
Driving a car	0.89	0.93	0.99	0.89	0.95	0.99	0.95	0.93	0.97
Lunch break	0.87	0.86	0.92	0.87	0.89	0.92	0.9	0.9	0.91
Walking	0.91	0.90	0.89	0.90	0.92	0.89	0.91	0.91	0.92
Working	0.92	0.96	0.95	0.97	0.97	0.95	0.96	0.92	0.97

Table 3
Baseline models comparison.

Model	Balanced accuracy	Precision	Recall	F1-Score
ClusterFL [19]	0.93	0.78	0.86	0.82
FL+HC [24]	0.94	0.85	0.89	0.83
HypCluster [25]	0.9	0.65	0.56	0.65
Dynamic clustering [26]	0.92	0.86	0.75	0.76
Clustered FedStack 1	0.98	0.95	0.91	0.93
Clustered FedStack 2	0.96	0.89	0.9	0.89
Clustered FedStack 3	0.94	0.91	0.92	0.91
Clustered FedStack 7	0.95	0.92	0.91	0.91
Clustered FedStack 8	0.98	0.94	0.93	0.93
Clustered FedStack 9	0.97	0.96	0.95	0.95

Table 4
Clustered FedStack performance in classification of drug ratings.

Model	Accuracy	Precision	Recall	F1-Score
ClusterFL	0.92	0.8	0.88	0.82
FL+HC	0.93	0.87	0.91	0.83
HypCluster	0.89	0.66	0.57	0.65
Dynamic clustering	0.91	0.88	0.77	0.76
Clustered FedStack 1	0.99	0.92	0.93	0.91
Clustered FedStack 2	0.98	0.91	0.92	0.91
Clustered FedStack 3	1	0.96	0.95	0.95
Clustered FedStack 4	0.97	0.94	0.93	0.93
Clustered FedStack 5	0.98	0.97	0.94	0.96
Clustered FedStack 6	1	0.97	0.93	0.95
Clustered FedStack 7	0.99	0.98	0.97	0.97
Clustered FedStack 8	0.98	0.97	0.94	0.96
Clustered FedStack 9	0.96	0.93	0.94	0.93
Clustered FedStack 10	0.94	0.93	0.92	0.91

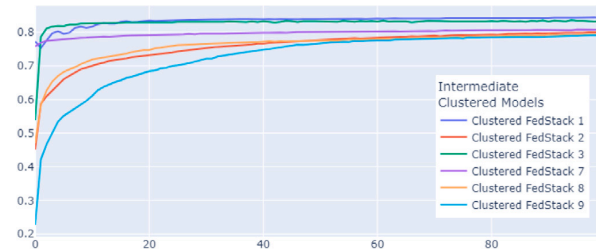


Fig. 6. Convergence of intermediate Clustered FedStack models on PPG-DALIA under the Cyclical learning rates.

scale function, a lambda function shown in Eq. (10), scales the initial amplitude by half with each cycle.

$$\text{lambda } x : \frac{1}{2^{(x-1)}} \quad (10)$$

Fig. 6 presents the convergence curves of six intermediate Clustered FedStack models from the three clustering techniques proposed in this study. The intermediate clustered models built based on K-Means and Agglomerative clustering converge faster than the clustered models built based on GMM clustering. There is not much difference in the number of epochs required for each clustered model to converge. All six models converge in less than 50 epochs. The results show that the proposed Clustered FedStack framework can be implemented with centroid-based, hierarchical or distribution-based clustering. The Clustered FedStack models built based on any of these clustering techniques converge quickly in 50 epochs.

5. Experiments on clustered FedStack scalability in NLP tasks

The scalability of the proposed Clustered FedStack model was rigorously assessed through a targeted evaluation. For this purpose, the drug review dataset [7], containing reviews and ratings, was utilized. This comprehensive dataset encompasses 3677 distinct drugs and 916 different medical conditions. The aim of this experiment is to classify drug ratings (1–10) based on input data such as medical conditions. In alignment with the clustering methodologies proposed in the Clustered FedStack framework, the GMM clustering was employed to perform

the clustering of 2191 drugs, resulting in 78 unique clusters as shown in Supplementary Material. The Supplementary Material also includes information on the cosine distance for 200 local clients (drugs).

The performance comparisons of different clustering models, including the top 10 variations of the Clustered FedStack model, are presented in **Table 4**. The metrics evaluated include accuracy, precision, recall, and F1-score for classifying drug ratings. Four baseline models are included: ClusterFL, FL+HC, HypCluster, and Dynamic Clustering. Their performances are relatively consistent, with accuracy ranging from 0.89 to 0.93. The Clustered FedStack models demonstrated superior performance, with notable improvements in all evaluated metrics. The accuracy for these variations ranged from 0.94 to a perfect 1, highlighting the efficiency and robustness of the model. The first five Clustered FedStack models exhibited particularly impressive results, achieving almost perfect or perfect accuracy. The precision, recall, and F1-score also showcased strong consistency and harmony, reflecting the model's ability to balance both false positives and false negatives.

These results underscore the scalability and effectiveness of the Clustered FedStack model across local clients with non-IID data. The model's scalability and adaptability are evident, maintaining high levels of accuracy and F1-scores regardless of the local clients' variation. This highlights the Clustered FedStack model's potential in managing large and intricate datasets like drug reviews and ratings, validating both its resilience and relevance to real-world applications.

The convergence of the proposed Clustered FedStack on the drug review dataset has been assessed, as shown in **Fig. 7**. The line chart presents a convergence pattern that denotes accuracy in the y-axis across 100 epochs in the x-axis. The values for Clustered FedStack 1 exhibited consistent growth, starting at 0.7882 and reaching 0.8556 by epoch 40. Similarly, other clustered FedStacks demonstrated a progressive increase in values across epochs, such as Clustered FedStack 2, which advanced from 0.5188 to 0.8811, signifying a gradual strengthening of the model. These convergence trends shed light on the efficiency and efficacy of the iterative learning process. Variations in convergence rates among different stacks were observed, reflecting the distinct characteristics of each clustered FedStack. These findings suggest a general trend of convergence towards higher values,

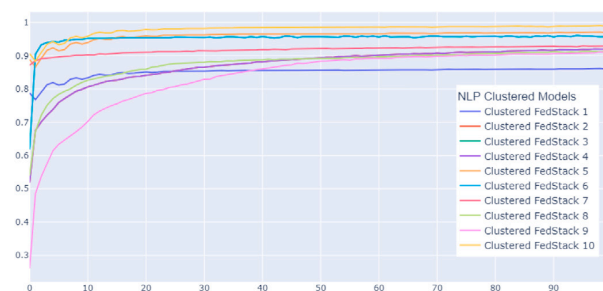


Fig. 7. Convergence of intermediate Clustered FedStack models on drug review dataset under Cyclical learning rates.

though occasional oscillations and fluctuations were detected in specific iterations. This in-depth analysis offers valuable insights into the behavior of clustered FL systems, potentially opening new avenues for enhanced optimization strategies and a more profound understanding of convergence mechanisms within distributed ML frameworks.

6. Conclusion

In the present study, a novel framework named Clustered-FedStack was introduced, designed to cluster local clients within the FL paradigm based on the weights of their output layers. This methodology was devised to address the non-IID challenge inherent to FL. It is important to acknowledge certain limitations of the proposed framework, notably its incompatibility with the application on local clients utilizing conventional Machine Learning models for the training of private data. Moreover, the global server model's process of clustering local clients operates on an unsupervised basis, without access to specific information about local clients, depending solely on the local model rather than client demographics. In light of these considerations, future investigations should aim to develop strategies for the dynamic clustering of local clients, taking into account meta-information that pertains to client similarities.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

Here is the link to publicly available dataset: <https://archive.ics.uci.edu/ml/datasets/PPG-DaLiA>.

Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.patrec.2023.12.004>.

References

- [1] K. Bonawitz, P. Kairouz, B. McMahan, D. Ramage, Federated learning and privacy: Building privacy-preserving systems for machine learning and data science on decentralized data, *Queue* 19 (5) (2021) 87–114.
- [2] Y. Shi, Y. Zhang, P. Zhang, Y. Xiao, L. Niu, Federated learning with l1 regularization, *Pattern Recognit. Lett.* (2023).
- [3] L. Peng, G. Luo, A. Walker, Z. Zaiman, E.K. Jones, H. Gupta, K. Kersten, J.L. Burns, C.A. Harle, T. Magoc, et al., Evaluation of federated learning variations for COVID-19 diagnosis using chest radiographs from 42 US and European hospitals, *J. Amer. Med. Inform. Assoc.* (2022).

- [4] T. Shaik, X. Tao, N. Higgins, R. Gururajan, Y. Li, X. Zhou, U.R. Acharya, FedStack: Personalized activity monitoring using stacked federated learning, *Knowl.-Based Syst.* 257 (2022) 109929.
- [5] M. Arafah, H. Ould-Slimane, H. Otrok, A. Mourad, C. Talhi, E. Damiani, Data independent warmup scheme for non-IID federated learning, *Inform. Sci.* 623 (2023) 342–360.
- [6] A. Reiss, I. Indlekofer, P. Schmidt, K. Van Laerhoven, Deep PPG: large-scale heart rate estimation with convolutional neural networks, *Sensors* 19 (14) (2019) 3079.
- [7] S. Kallumadi, F. Grer, Drug Review Dataset (Drugs.com), *UCI Machine Learning Repository*, 2018, <http://dx.doi.org/10.24432/C5SK5S>.
- [8] B. McMahan, E. Moore, D. Ramage, S. Hampson, B.A. y Arcas, Communication-efficient learning of deep networks from decentralized data, in: *Artificial Intelligence and Statistics*, PMLR, 2017, pp. 1273–1282.
- [9] M.G. Arivazhagan, V. Aggarwal, A.K. Singh, S. Choudhary, Federated learning with personalization layers, 2019, arXiv preprint arXiv:1912.00818.
- [10] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, Y. Khazaeni, Federated learning with matched averaging, in: *International Conference on Learning Representations*, 2020.
- [11] A. Osmani, M. Hamidi, Reduction of the position bias via multi-level learning for activity recognition, in: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, Springer, 2022, pp. 289–302.
- [12] Z. Xiao, X. Xu, H. Xing, F. Song, X. Wang, B. Zhao, A federated learning system with enhanced feature extraction for human activity recognition, *Knowl.-Based Syst.* 229 (2021) 107338.
- [13] Y. Pang, H. Zhang, J.D. Deng, L. Peng, F. Teng, Rule-based collaborative learning with heterogeneous local learning models, in: *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, Springer, 2022, pp. 639–651.
- [14] P. Tian, W. Liao, W. Yu, E. Blasch, WSCC: A weight-similarity-based client clustering approach for Non-IID federated learning, *IEEE Internet Things J.* 9 (20) (2022) 20243–20256.
- [15] A. Alam, S. Qazi, N. Iqbal, K. Raza, Fog, edge and pervasive computing in intelligent internet of things driven applications in healthcare: Challenges, limitations and future use, in: *Fog, Edge, and Pervasive Computing in Intelligent IoT Driven Applications*, Wiley Online Library, 2020, pp. 1–26.
- [16] L.M. Dang, M. Piran, D. Han, K. Min, H. Moon, et al., A survey on internet of things and cloud computing for healthcare, *Electronics* 8 (7) (2019) 768.
- [17] S. Ek, F. Portet, P. Lalanda, G. Vega, Evaluation of federated learning aggregation algorithms: application to human activity recognition, in: *Adjunct Proceedings of the 2020 ACM International Joint Conference on Pervasive and Ubiquitous Computing and Proceedings of the 2020 ACM International Symposium on Wearable Computers*, 2020, pp. 638–643.
- [18] Y. Zhao, H. Liu, H. Li, P. Barnaghi, Semi-supervised federated learning for activity recognition, *ACM Trans. Intell. Syst. Technol.* 1 (1) (2021).
- [19] X. Ouyang, Z. Xie, J. Zhou, J. Huang, G. Xing, Clusterfl: a similarity-aware federated learning system for human activity recognition, in: *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services*, 2021, pp. 54–66.
- [20] C. Xiang, P.C. Yong, L.S. Meng, Design of multiple-level hybrid classifier for intrusion detection system using Bayesian clustering and decision trees, *Pattern Recognit. Lett.* 29 (7) (2008) 918–924.
- [21] L. Yang, J. Huang, W. Lin, J. Cao, Personalized federated learning on non-IID data via group-based meta-learning, *ACM Trans. Knowl. Discov. Data (TKDD)* (2022).
- [22] X. Shang, Y. Lu, Y.-m. Cheung, H. Wang, FEDIC: Federated learning on non-IID and long-tailed data via calibrated distillation, in: *2022 IEEE International Conference on Multimedia and Expo (ICME)*, IEEE, 2022, pp. 1–6.
- [23] J. Wang, Y. Chen, S. Hao, X. Peng, L. Hu, Deep learning for sensor-based activity recognition: A survey, *Pattern Recognit. Lett.* 119 (2019) 3–11.
- [24] C. Briggs, Z. Fan, P. Andras, Federated learning with hierarchical clustering of local updates to improve training on non-IID data, in: *2020 International Joint Conference on Neural Networks (IJCNN)*, IEEE, 2020, pp. 1–9.
- [25] Y. Mansour, M. Mohri, J. Ro, A.T. Suresh, Three approaches for personalization with applications to federated learning, 2020, arXiv preprint arXiv:2002.10619.
- [26] Y. Kim, E. Al Hakim, J. Haraldson, H. Eriksson, J.M.B. da Silva, C. Fischione, Dynamic clustering in federated learning, in: *ICC 2021-IEEE International Conference on Communications*, IEEE, 2021, pp. 1–6.
- [27] A.K. Jain, Data clustering: 50 years beyond K-means, *Pattern Recognit. Lett.* 31 (8) (2010) 651–666.
- [28] A. Sellami, A.B. Abbes, V. Barra, I.R. Farah, Fused 3-D spectral-spatial deep neural networks and spectral clustering for hyperspectral image classification, *Pattern Recognit. Lett.* 138 (2020) 594–600.
- [29] J. Lücke, D. Forster, k-means as a variational EM approximation of Gaussian mixture models, *Pattern Recognit. Lett.* 125 (2019) 349–356.
- [30] L.N. Smith, Cyclical learning rates for training neural networks, in: *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, 2017, pp. 464–472.

4.2 Summary

The chapter concludes by highlighting the notable achievements of the Clustered FedStack framework, demonstrating its superior performance in HAR and NLP tasks compared to traditional federated learning models. The innovative use of clustering techniques and BIC for model optimization has proven effective in overcoming the inherent limitations of non-IID data in federated learning. The successful application of the Clustered FedStack framework in classifying complex human activities and drug review ratings underscores its potential for scalable, efficient, and privacy-preserving AI models in healthcare and beyond. This chapter not only contributes to the advancement of federated learning but also sets a precedent for future research in optimizing federated learning models for diverse and challenging datasets.

CHAPTER 5: PAPER 4 - ADAPTIVE MULTI-AGENT DEEP REINFORCEMENT LEARNING FOR TIMELY HEALTHCARE INTERVENTIONS

5.1 Introduction

This chapter delves into an advanced AI-driven framework for patient monitoring, utilizing the principles of multi-agent deep reinforcement learning (DRL) to offer unprecedented personalisation and responsiveness in healthcare settings. It meticulously examines the deployment of sophisticated DRL agents, each engineered to monitor specific physiological parameters, such as heart rate, respiration, and body temperature, within a unified healthcare ecosystem. These agents are designed to learn and adapt from continuous patient data streams, enabling them to make predictive and prescriptive interventions. This section articulates how this approach surmounts the constraints of conventional monitoring systems by dynamically adjusting to patient-specific conditions, thereby elevating the standard of patient care, safety, and overall healthcare efficiency.

Adaptive Multi-Agent Deep Reinforcement Learning for Timely Healthcare Interventions

Thanveer Shaik, Xiaohui Tao, Lin Li, Haoran Xie, Hong-Ning Dai, and Jianming Yong

Abstract—Effective patient monitoring is vital for timely interventions and improved healthcare outcomes. Traditional monitoring systems often struggle to handle complex, dynamic environments with fluctuating vital signs, leading to delays in identifying critical conditions. To address this challenge, we propose a novel AI-driven patient monitoring framework using multi-agent deep reinforcement learning (DRL). Our approach deploys multiple learning agents, each dedicated to monitoring a specific physiological feature, such as heart rate, respiration, and temperature. These agents interact with a generic healthcare monitoring environment, learn the patients’ behavior patterns, and make informed decisions to alert the corresponding Medical Emergency Teams (METs) based on the level of emergency estimated. In this study, we evaluate the performance of the proposed multi-agent DRL framework using real-world physiological and motion data from two datasets: PPG-DaLiA and WESAD. We compare the results with several baseline models, including Q-Learning, PPO, Actor-Critic, Double DQN, and DDPG, as well as monitoring frameworks like WISEML and CA-MAQL. Our experiments demonstrate that the proposed DRL approach outperforms all other baseline models, achieving more accurate monitoring of patient’s vital signs. Furthermore, we conduct hyperparameter optimization to fine-tune the learning process of each agent. By optimizing hyperparameters, we enhance the learning rate and discount factor, thereby improving the agents’ overall performance in monitoring patient health status. Our AI-driven patient monitoring system offers several advantages over traditional methods, including the ability to handle complex and uncertain environments, adapt to varying patient conditions, and make real-time decisions without external supervision. However, we identify limitations related to data scale and prediction of future vital signs, paving the way for future research directions.

Impact Statement—The proposed approach, which combines artificial intelligence and multi-agent deep reinforcement learning, revolutionizes patient monitoring by offering personalized and real-time solutions. Unlike conventional systems, our approach empowers virtual agents to make autonomous decisions based on raw patient data, adapting to each individual’s physiological patterns and unique health conditions. Its significance lies in its potential to improve patient outcomes, reduce medical errors, and optimize medical resource allocation. The continuous and dynamic monitoring enabled by our approach ensures timely responses to critical situations, enhancing patient safety. Moreover, the framework’s ability to learn from raw data eliminates the need for complex feature engineering and

large labeled datasets, allowing for scalability and application in diverse medical scenarios. Beyond enhancing patient care, our research drives advancements in AI-driven healthcare technology, paving the way for more efficient, accurate, and patient-centric healthcare systems.

Index Terms—Behavior Patterns, Decision Making, Patient Monitoring, Reinforcement Learning, Vital Signs.

I. INTRODUCTION

In the dynamic domain of healthcare, the significance of informed decision-making cannot be overstated. With the advent of continuous patient monitoring systems, it has become possible to remotely track vital signs and physical movements, thereby enhancing the decision-making capabilities of clinicians [1]. The application of machine learning models to analyze transmitted vital sign data has seen a significant uptick in various healthcare applications, ranging from pre-clinical data processing and diagnosis assistance to early warning detection of health deterioration, treatment decision-making, and drug prescription [2], [3]. In this context, the monitoring of human behavior patterns plays a crucial role, especially for remote patient monitoring in hospitals or through Internet of Things (IoT)-enabled home monitoring systems.

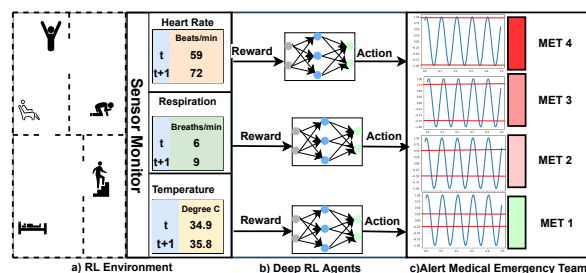


Fig. 1: Human monitoring framework to monitor vital signs of the client during regular activities and alert medical emergency teams accordingly.

Traditionally, methodologies in this field have predominantly relied on unsupervised and supervised learning techniques to identify patterns and classify patients’ activities and vital signs [4], [5]. However, these techniques are limited in their capacity to only observe data and suggest potential decisions without the ability to act upon these observations. In contrast, Reinforcement Learning (RL) introduces a novel paradigm by deploying learning agents within complex and uncertain environments. These agents are empowered to explore and exploit the environment through actions, learning

Thanveer Shaik and Xiaohui Tao are with the School of Mathematics, Physics & Computing, University of Southern Queensland, Queensland, Australia (e-mail: Thanveer.Shaik@usq.edu.au, Xiaohui.Tao@usq.edu.au).

Lin Li is with the School of Computer and Artificial Intelligence, Wuhan University of Technology, China (e-mail: cathylilin@whut.edu.cn)

Haoran Xie is with the Department of Computing and Decision Sciences, Lingnan University, Tuen Mun, Hong Kong (e-mail: hrxie@ln.edu.hk)

Hong-Ning Dai is with the Department of Computer Science, Hong Kong Baptist University, Hong Kong (e-mail: henrydai@hkbu.edu.hk).

Jianming Yong is with the School of Business, University of Southern Queensland, (e-mail: Jianming.Yong@usq.edu.au)

from the outcomes of their actions [6]. A cornerstone of RL is its reward mechanism, which provides the agent with feedback in the form of rewards for its actions. These rewards serve as crucial signals that guide the learning process of the agent, encouraging actions that lead to favorable outcomes and discouraging detrimental ones. This reward system is instrumental in enabling the agent to iteratively refine its strategy based on the consequences of its actions, thereby enhancing its performance over time.

The versatility of RL has been demonstrated in various dynamic domains, such as stock market trading [7], and is increasingly being adapted for healthcare applications, including diagnostic decisions and dynamic treatment regimes that require the consideration of delayed feedback [8]. Specifically, RL-based patient monitoring applications have focused on optimizing the timing and dosage of medications to ensure their correct administration [9], [10]. The analogy of probabilistic machine learning models, such as RL, to an ICU clinician monitoring a patient's state and making subsequent decisions based on observed changes, underscores the potential of RL in healthcare [11].

This study addresses the complex challenge of monitoring multiple vital signs of the human body, tracking health status, and initiating timely interventions during emergencies. We propose an innovative approach that employs multiple deep learning agents within a generic healthcare monitoring environment. Each agent is tasked with monitoring specific vital signs, taking into account different threshold levels for each sign. These agents progressively learn the threshold levels of vital signs based on Modified Early Warning Scores (MEWS) and the rewards accumulated from previous iterations. The well-trained DRL agents are capable of monitoring a patient's heart rate, respiration rate, and temperature, and alerting the clinical team in case of deviations from predefined thresholds [12].

The primary aim of this study is to learn human behavior patterns in the context of clinical health by deploying a DRL agent for each physiological feature. These agents are designed to monitor, learn, and alert the respective clinical teams if any vital signs deviate from the norms established by MEWS, as shown in Fig. 1. We introduce a novel approach for rewarding the actions of RL agents to facilitate the learning of behavior patterns. The generic monitoring environment developed in this study supports multi-agent functionality to monitor various vital signs of a patient, thereby introducing a new paradigm for remotely monitoring patients' health status using a multi-agent DRL environment.

The contributions of this study are as follows:

- Introduction of a novel approach for rewarding RL agents' actions to foster the learning of behavior patterns.
- Development of a generic monitoring environment that accommodates multi-agents for monitoring various vital signs of a patient.
- Establishment of a new paradigm for remotely monitoring patients' health status utilizing the multi-agent DRL environment.

The paper is organized as follows: Section II presents related works in the RL community, specifically focusing on learning

human behavior patterns and applications in the healthcare domain. The research problem formulation and the proposed multi-agent DRL methodology are detailed in Section III. Section IV evaluates the proposed methodology on 10 different subject vital signs, and baseline models are discussed. In Section V, the results of the proposed approach are compared with baseline models, and hyper-parameter optimization of the learning rate and discount factor is discussed. Based on the results, applications of the proposed framework are discussed in Section VI. Section VII concludes the paper, including limitations and future work.

II. RELATED WORKS

A. Machine Learning in Healthcare

Machine learning has transformed healthcare with its ability to predict, detect, and monitor, as noted in [13]. Supervised learning algorithms can learn from labeled data and make predictions or classify based on the input features [14]. For example, machine learning or deep learning techniques can predict human vital signs like heart rate or classify physical activities [15]. In a study by Oyeleye et al. [16], machine learning and deep learning models were used to estimate heart rate using data from wearable devices. The authors tested different regression algorithms including linear regression, k-nearest neighbor, decision tree, random forest, autoregressive integrated moving average, support vector regressor, and long short-term memory recurrent neural networks. Similarly, Luo et al. [17] utilized the LSTM model to predict heart rate based on five factors: heart rate signal, gender, age, physical activities [18], and mental state. Unsupervised learning algorithms learn from unlabeled data and find patterns using association and clustering techniques [19]. Sheng and Huber [19] developed an unsupervised method with an encoder and decoder network to identify similar physical activities using clustering, which achieved a clustering accuracy of 85% based on learning embeddings and behavior clusters. RL, on the other hand, does not require prior knowledge or information and works on an environment-driven approach [20]. The agents learn through receiving rewards or penalties based on their actions which is called as experience. Unlike supervised learning, RL can learn a sequence of actions through exploration and exploitation and does not require extensive labeled data for data-driven models [21].

B. Mimic Human Behavior Patterns

Tirumala et al. [22] studied how to understand human behavior patterns and identify common movements and interactions in a set of related tasks and situations. They used probabilistic trajectory models to develop a framework for hierarchical reinforcement learning (HRL). Janssen et al. [23] suggested breaking down a complex task such as biological behavior into smaller parts, with HRL able to organize sequential actions into a temporary option. They compared biological behavior to options in HRL. Tsiakas et al. [24] proposed a human-centered cyber-physical systems framework for personalized human-robot collaboration and training, focusing on monitoring and evaluating human behavior. The authors aimed

to effectively predict human attention with the minimum and least intrusive sensors. Kubota et al. [25] investigated how robots can adapt to the behavior of people with cognitive impairments for cognitive neuro-rehabilitation. They explored different types of robots for therapeutic, companion, and assistive applications. For health applications, robots must be able to perceive and understand human behavior, which includes high-level behaviors like cognitive abilities and engagement, as well as low-level behaviors like speech, gesture, and physiological signs.

C. RL in Healthcare

Lisowska et al. [26] developed a digital intervention for cancer patients to promote positive health habits and lifestyle changes. They used RL to determine the best time to send intervention prompts to the patients and employed three RL approaches (Deep Q-Learning, Advance Actor Critic, and proximal policy optimization) to create a virtual coach for sending prompts. Other studies have also shown that personalized messages can increase physical activity in type 2 diabetes patients [27]. Li et al. [28] proposed a RL approach based on electronic health records for sequential decision-making tasks. They used a model-free Deep Q Networks (DQN) algorithm to make clinical decisions based on patient data and achieved better results with cooperative multi-agent RL. R decision-making can also be used for human activity recognition, as shown in a study that proposed a dynamic weight assignment network architecture and used a twin delayed deep deterministic algorithm inspired by various other RL algorithms [29].

RL has been widely researched for various applications, including gaming, learning and mimicking human behavior, and deploying socially assistive robots. However, deploying physical robots with human interaction capabilities in sensitive locations such as hospitals, elderly care facilities, and mental health facilities may pose safety risks to patients, carers, and medical staff [?]. Existing health monitoring applications using supervised or unsupervised learning cannot effectively handle uncertain events in the dynamic hospital environment. To address these challenges, virtual robots with adaptive learning abilities can be deployed to monitor and learn human behavior based on their vital signs. In this study, we developed a custom human monitoring environment that can learn behavior patterns from human vital signs and alert the appropriate clinical team in case of an emergency. This work provides a novel approach to learn and monitor human behavior patterns in a safe and clinically effective manner.

III. DRL MONITORING FRAMEWORK

In this section, the design of a human behavior monitoring system, DRL monitoring framework, that uses R is presented in detail. The aim of the system is to monitor vital signs to learn human behavior patterns and ensure clinical safety in an uncertain environment. The proposed framework involves a multi-agent system where each vital sign state is observed by an individual agent, as shown in Fig.2. A DRL algorithm, DQN, is used to learn effective strategies in the sequential decision-making process without prior knowledge through trial-and-error interactions with the environment [30].

A. Problem Formulation

The challenge addressed in this research is the development of a multi-agent framework for real-time health status monitoring by learning and interpreting patterns in vital signs through wearable sensors. The agents must detect deviations from normal vital sign patterns that exceed Modified Early Warning Scores (MEWS) thresholds and alert the emergency team accordingly.

To formulate this problem, we leverage the framework of Markov Decision Processes (MDP), expressed as a 5-tuple $M = (S, A, P, R, \gamma)$. Here, S represents the finite state space, where each state $s_t \in S$ corresponds to a distinct combination of vital sign readings at time t . The action set A comprises potential alerts the agents can issue based on the observed vital signs. The transition function $P(s, a, s')$ models the probability of moving from state s to state s' upon taking action a , reflecting the dynamic nature of human vital signs.

Central to our approach is the reward function $R(s, a)$, which is defined to prioritize actions that lead to the early detection of potential health risks, thereby enabling timely intervention. This is mathematically represented as:

$$R(s_t, a_t) = \sum_{t=0}^{\infty} \gamma^t r_t, \quad (1)$$

where γ is the discount factor that balances the importance of immediate versus future rewards, ensuring the agents' actions are aligned with long-term health monitoring objectives.

The goal is to discover an optimal policy $\pi(s_t)$ that maximizes the expected reward by selecting the most appropriate action a_t in any given state s_t . This optimization is achieved through the iterative update of the Q-function, as outlined in the Bellman equation:

$$Q^{new}(s_t, a_t) \leftarrow (1-\alpha)Q(s_t, a_t) + \alpha \left(r_t + \gamma \max_a Q(s_{t+1}, a) \right), \quad (2)$$

where α represents the learning rate, influencing the integration of new information into the Q-function. Through this process, the agents continually refine their decision-making strategy, enhancing the system's capability to monitor and respond to emerging health risks effectively.

TABLE I: Modified Early Warning Scores [31]

MEWS	4/MET	3	2	1	0	1	2	3	4/MET
Respiratory Rate	≤4	5-8			9-20	21-24	25-30	31-35	≥36
Oxygen Saturation	≤84	85-89	90-92	93-94	≥95				
Temperature	≤34.0	34.1-35.0	35.1-36.0	36.1-37.9	38.0-38.5	≥38.6			
Heart Rate	≤39		40-49	50-99	100-109	110-129	130-139	≥140	
Sedation Score					Awake		Mild	Moderate	Severe

B. Monitoring Environment

A custom RL monitoring system based on MDP has been created to have human vital signs data serve as the observation space S , action space A for learning agents to make decisions, and rewards R for the agents' actions as depicted in Fig. 2.

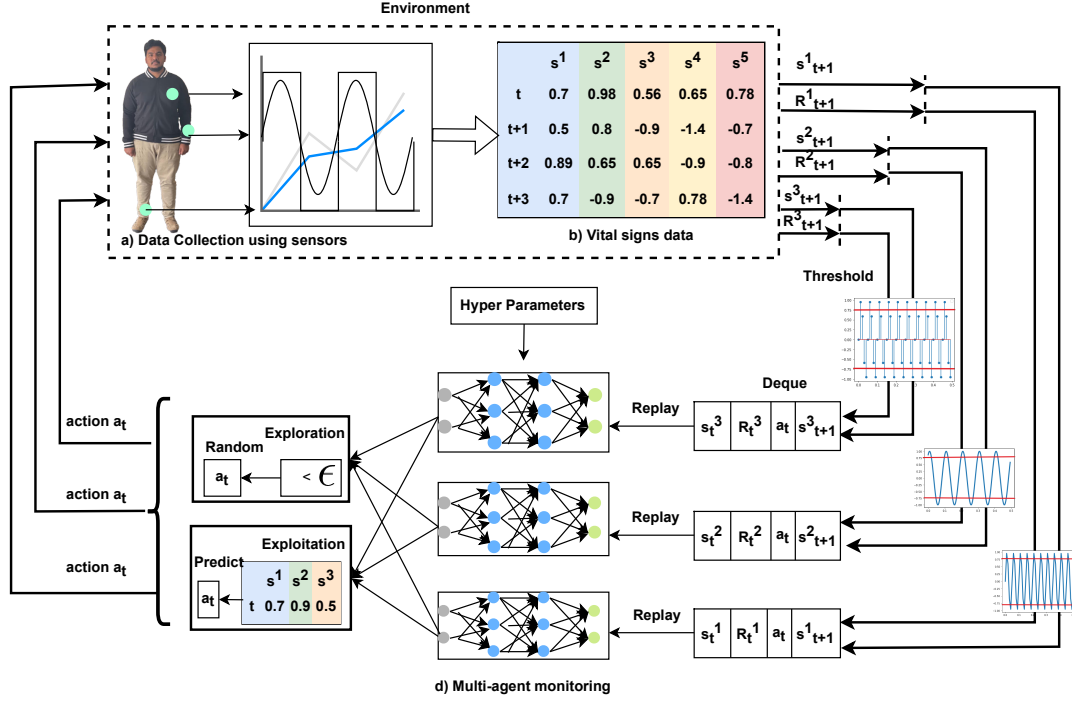


Fig. 2: Multi-agent monitoring framework

This study introduces a novel isolated multi-agent MDP framework that allows multi-agents to share the same environment and make decisions based on the health parameters they are monitoring, receiving rewards without being influenced by the decisions of other agents. The goal of all agents in this environment is to monitor the health of patients using the predefined MEWS, as shown in Tab. I. In healthcare, each vital sign plays a critical role in determining a person's clinical safety.

1) *Observation Space:* The environment in Fig.2 has a state, represented by $s_t^i \in S$, where $i = 0, 1, 2, \dots, n$, refers to observations at time t . The aim is to divide the state into observations and allocate them to multi-agents. Suppose S represents the state of the human body, and there are three observations, $s_t^0, s_t^1, s_t^2 \in S$, that represent different internal vital signs of the human body at time t . The human health status is controlled by multiple internal vital signs, each with a different threshold as shown in MEWS Tab.I. Using a single agent to monitor all the vital signs can result in a sparse rewards challenge [32], where the environment might produce few useful rewards and hinders the learning of an agent. Therefore, multi-agents need to be deployed for each human to monitor the critical vital signs. The expected return E_π of a policy π in a state s can be defined by state-value Eq. 3 in the multi-agent setting, where $i = 0, 1, 2, 3, \dots, n$ is a finite number of observations n in the state.

$$V^\pi(s^i) = E_\pi \left\{ \sum_{t=0, i=0}^{\infty, n} \gamma^t R(s_t, \pi(s_t)) | s_0^i = s \right\} \quad (3)$$

2) *Action Space:* The action space of the monitoring environment is defined based on the MEWS [31] as shown in Tab. I. The table presents early warning scores of adults' vital signs with the appropriate Medical Emergency Team (MET) to contact if any escalations in the health parameters. Based on the MEWS as threshold values, the action space has been segmented to have five discrete actions to communicate the vital signs to MET-0, MET-1, MET-2, MET-3, and MET-4. Each of these actions will be taken by agents based on the current state of the vital signs they are monitoring. The expected return E_π for taking an action a in a state s under a policy π can be measured using the action-value function $Q_\pi(s, a)$ defined in Eq. 4.

$$Q^\pi(s, a) = E_\pi \left\{ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, \pi(s_t)) | s_0 = s, a_0 = a \right\} \quad (4)$$

3) *Rewards:* The goal of RL is to maximize cumulative rewards obtained through the actions of learning agents in an environment. In traditional RL, an agent is rewarded based on its action that leads to a transition from state s_t to s_{t+1} . In this study, the objective of the learning agent is to learn patterns in human vital signs. This is achieved through the design of an effective reward policy. The reward policy, as

TABLE II: Rewards Policy

MEWS	4	3	2	1	0
Action 0	-4	-3	-2	-1	10
Action 1	-4	-3	-2	10	-1
Action 2	-4	-3	10	-1	-2
Action 3	-4	10	-1	-2	-3
Action 4	10	-3	-2	-1	-4

defined in this study, is calculated using Eq. 5. The agents are positively rewarded if they monitor vital signs in a state and take the correct action from the action space to communicate with the correct MET as defined in MEWS Tab.I. On the other hand, if the agent takes the wrong action, it is negatively rewarded. The rewards are split into five categories for the five actions in the action space based on the MET from MEWS Tab.I. The full rewards for each action selected by the agents are presented in Tab.II. The reward policy utilizes the DRL agents' desire to maximize rewards in each learning iteration, making them learn the behavior patterns. Under each category, different levels of rewards were configured. For example, an observation $s_t^1 \in S$ at the time t is related to heart rate falling under MET-4, the rewards are shown in Eq. 6.

$$R(s_t, a_t) = \begin{cases} +reward & \text{if } action = MET \\ -reward & \text{if } action \neq MET \end{cases} \quad (5)$$

$$R(s_t^1, a_t) = \begin{cases} 10 & \text{if } MET = 4 \& action = 4 \\ -1 & \text{if } MET = 4 \& action = 3 \\ -2 & \text{if } MET = 4 \& action = 2 \\ -3 & \text{if } MET = 4 \& action = 1 \\ -4 & \text{if } MET = 4 \& action = 0 \end{cases} \quad (6)$$

C. Learning Agent

In this study, a game learning agent DQN algorithm is employed. The DQN algorithm was first introduced by DeepMind, a subsidiary of Google, for playing Atari games. It allows the agent to play games by simply observing the screen, without any prior training or knowledge about the games. The DQN algorithm approximates the Q-Learning function using neural networks, and the learning agent is rewarded based on the neural network's prediction of the best action for the current state. For this research, the reward policy is described in more detail in Section III-B3.

1) *Function Approximation*: The neural network used in this study to estimate the Q-values for each action has three layers: an input layer, a hidden layer, and an output layer. The input layer has a node for each vital sign in a state and the output layer has a node for each action in the action space. The model is configured with a relu activation function, mean square error as the loss function, and an Adam optimizer. The model is trained on the states and their corresponding rewards and, once trained, it can predict the accumulated reward.

The learning agent takes an action $a_t \in A$ in a transition from state s_t to s_t' and receives a reward R . In this transition, the maximum Q-function value is calculated according to

Eq. 4, and the calculated value is discounted by a discount factor γ to prioritize immediate rewards over future rewards. The discounted future reward is combined with the current reward to obtain the target value. The difference between the prediction from the neural network and the target value forms the loss function, which is a measure of the deviation of the predicted value from the target value and can be estimated using Eq. 7. The square of the loss function penalizes the agent for large loss values.

$$loss = \left(\underbrace{R + \gamma \cdot \max(Q^{\pi^*}(s, a))}_{target_value} - \underbrace{Q^{\pi}(s, a)}_{predicted_value} \right)^2 \quad (7)$$

2) *Memorize and Replay*: The basic neural network model has a limitation in its memory capacity and can forget previous observations as they are overwritten by new observations. To mitigate this issue, a memory array that stores the previous observations including the current state s_t , action a_t , reward R , and next state s_t' is used. This memory array enables the neural network to be retrained using the replay method, where a random sample of previous observations from the memory is selected for training. In this study, the neural network model was retained by using a batch size of 32 previous observations.

3) *Exploration and Exploitation*: The exploration-exploitation trade-off in RL refers to the balancing act between trying out new actions to gather information and exploiting the actions that lead to the highest rewards. This balance can be modeled mathematically using the ϵ -greedy algorithm, which defines a probability ϵ of choosing a random action and a probability $1 - \epsilon$ of choosing the action believed to lead to the highest reward based on the current knowledge of the action-value function $Q(s_t, a)$. The equation to determine the action taken at time t is as follows:

$$a_t = \begin{cases} random(a_t) & \text{with probability } \epsilon \\ greedy(a_t) & \text{with probability } 1 - \epsilon \end{cases} \quad (8)$$

where the greedy action is defined as:

$$a_t = \arg \max_a Q(s_t, a) \quad (9)$$

The value of ϵ determines the level of exploration versus exploitation, with smaller values leading to more exploitation and larger values leading to more exploration. Over time, as the action-value function becomes more accurate, ϵ can be decreased to allow for more exploitation and convergence to the optimal policy.

4) *Hyper Parameters*: Other than the parameters defined for the neural networks, a set of hyperparameters has to supply for the RL process. They are as follows:

- **episodes (\mathcal{M})**: This is a gaming term that means the number of times an agent has to execute the learning process.
- **learning_rate(α)**: Learning rate is to determine much information neural networks learn in an iteration.
- **discount_factor(γ)**: Discount factor ranges from 0 to 1 to limit future rewards and focus on immediate rewards.

Algorithm 1 multi-agents Monitoring

Require: Input: a set of subjects $\mathcal{C} = \{1, 2, \dots, C\}$; a set of vital signs $\mathcal{V} = \{1, 2, \dots, V\}$; Episodes $\mathcal{M} = \{1, 2, \dots, M\}$;

Ensure: Output: Rewards achieved by agents in each episode.

- 1: **Initialization** : $observation_space = s_t \in \mathcal{S}, action_space = a_t \in \mathcal{A}, reward R, \gamma, \epsilon, \epsilon_{decay}, \epsilon_{min}, memory = \emptyset, batch_size$
- 2: Set $monitor_length = N$
- 3: **if** action is appropriate **then**
- 4: $R \leftarrow +reward$
- 5: **else**
- 6: $R \leftarrow -reward$
- 7: **end if**
- 8: **Define** $model \leftarrow NeuralNetworkModel$
- 9: $memory \leftarrow memory \cup (s_t, a_t, R, s_{t+1})$
- 10: **if** $np.random.rand < \epsilon$ **then** ▷ Exploration
- 11: $action_value \leftarrow random(a_t)$
- 12: **else** ▷ Exploitation
- 13: $action_value \leftarrow greedy(a_t)$
- 14: **end if**
- 15: **for** episode $m \in \mathcal{M}$ **do**
- 16: $score = 0$
- 17: **for** time in range($monitor_length$) **do**
- 18: $a_t \leftarrow action(s_t)$
- 19: $s_{t+1}, R, done \leftarrow step(a_t)$
- 20: $memory \leftarrow memory \cup (s_t, a_t, R, s_{t+1})$
- 21: $s_t \leftarrow s_{t+1}$
- 22: **if** done **then**
- 23: $display m, score$
- 24: $break$
- 25: **end if**
- 26: **end for**
- 27: $replay \leftarrow batch_size$
- 28: **end for**

Algorithm 1 implements the proposed multi-agent human monitoring framework. It takes as input a set of subjects $\mathcal{C} = 1, 2, \dots, C$ and a set of vital signs $\mathcal{V} = 1, 2, \dots, V$, along with the number of episodes $\mathcal{M} = 1, 2, \dots, M$. The algorithm outputs the rewards achieved by agents in each episode. Lines 1-2 initializes all the parameters needed for monitoring the environment and learning agent. Lines 3-7 present the reward policy. Lines 8-14 present the function approximation using the neural networks model, memorize & replay, exploration & exploitation of the DRL agent. Lines 15-28 are nested for loops with conditional statements to check if the episode is completed or not. The outer loop is to iterate each episode while resetting the environment to initial values and score to zero. The inner loop is to iterate timesteps which denote the time of the current state and calls the methods.

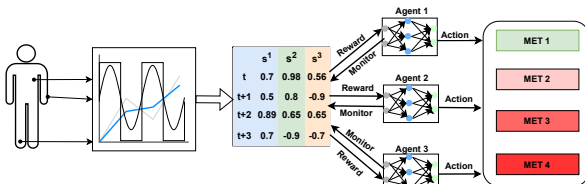


Fig. 3: Experimental Design

IV. EXPERIMENT

In this study, the proposed multi-agent framework was evaluated by deploying an agent for each physiological feature of a different set of subjects. The aim of the learning agents was to monitor their respective vital signs, communicate with the corresponding MET based on the estimated level of emergency, and learn the subjects' behavior patterns. All the experiments were conducted using Python programming language version 3.7.6 and related libraries such as TensorFlow, Keras, Open Gym AI, and stable_baselines3.

A. Dataset

- **PPG-DaLiA [33]:** The dataset contains physiological and motion data of 15 subjects, recorded from both a wrist-worn device and a chest-worn device while the subjects were performing a wide range of activities under conditions close to real life.
- **WESAD [34]:** The WESAD (Wearable Stress and Affect Detection) dataset is a collection of physiological signals recorded from participants while they perform various activities. It includes multi-modal signals such as ECG, PPG, GSR, respiration, and body temperature.

B. Baseline Models

- **WISEML [35]:** Mallozzi et al. proposed an RL framework for runtime monitoring to prevent dangerous and safety-critical actions in safety-critical applications. In this framework, runtime monitoring is used to enforce properties to the agent and shape its reward during learning.
- **CA-MQL [36]:** Chen et al. proposed constrained action-based MQL (CA-MQL) for UAVs to autonomously make flight decisions that consider the uncertainty of the reference point location.
- Existing RL baseline models by Li et al. [28] were deployed to optimize sequential treatment strategies based on Electronic Health Records (EHRs) for chronic diseases using DQN. The multi-agent framework results were compared with Q-Learning and Double DQN.
- Similarly, RL was deployed to recognize human activity using a dynamic weight assignment network architecture with TD3 (a combination of Deep Deterministic Policy Gradient (DDPG), Actor-Critic, and DQN) by Guo et al. [29].
- Yom et al. [27] used Advantage Actor-Critic (A2C) and Proximal Policy Optimization (PPO) algorithms to act as virtual coaches in decision-making and send personalized messages.

C. Performance Measure

Cumulative Rewards is a performance metric used in RL to measure the total rewards obtained by an agent over a specified period of time or number of actions. It is calculated as the sum of all rewards received by the agent over the given period of time or number of actions. The cumulative reward can be used to evaluate the effectiveness of a RL algorithm or to compare different algorithms.

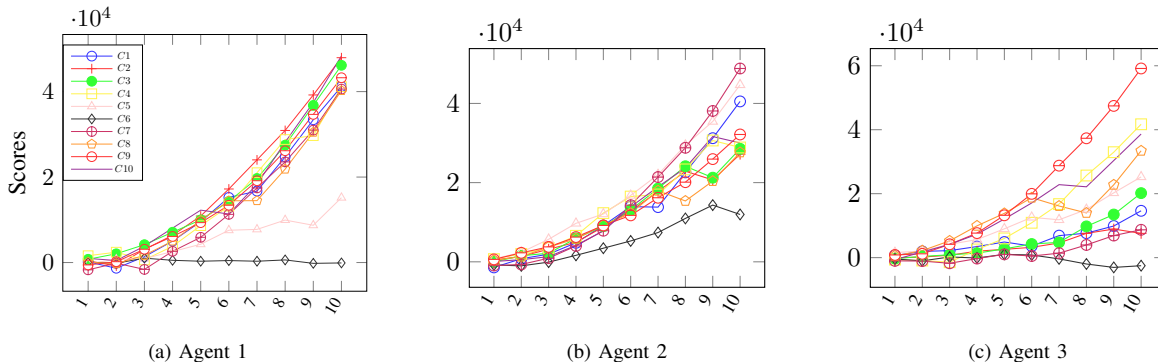


Fig. 4: DQN Agents Performance

TABLE III: DRL Agents Performance

RL Method	PPG-DaLiA Dataset			WESAD Dataset		
	agent 1	agent 2	agent 3	agent 1	agent 2	agent 3
Q Learning	25878	17304	23688	25318	16341	22823
PPO	23688	20367	17688	23128	19404	16823
A2C	24717	13707	24369	24157	12744	23504
Double DQN	25569	15360	20367	25009	14397	19502
DDPG	26760	20754	23967	26200	19791	23102
WISEML	28654	25789	33669	28094	24826	32804
CA-MQL	32985	27856	34685	32425	26893	33820
Proposed DRL Framework	48354	30019	38651	47794	29056	37786

V. EXPERIMENT RESULTS AND ANALYSIS

The advantage of RL for monitoring systems is that it can learn to handle complex, dynamic environments. Many monitoring tasks involve making decisions based on incomplete, uncertain information, and the optimal decision may depend on the context of the situation [37]. RL can learn to make decisions in these types of problems by considering the current state of the system and past experience. In this study, the aim is to leverage the RL capability to optimize the decision-making process in patient monitoring.

A. DRL Agents Performance

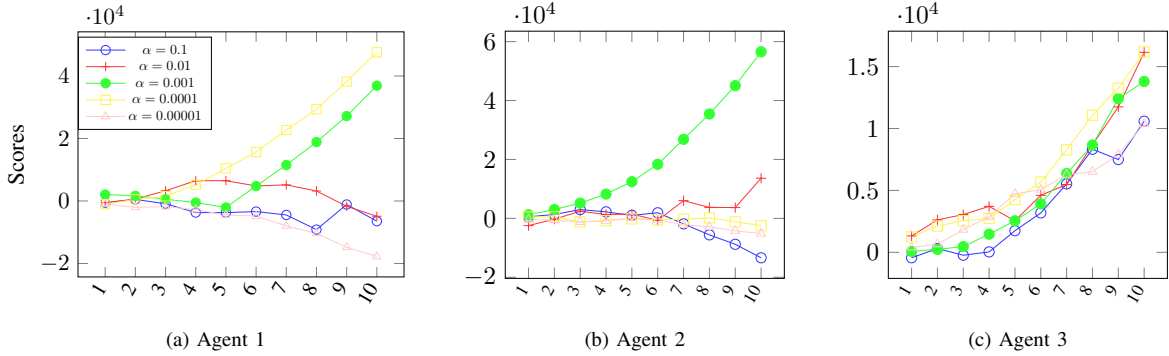
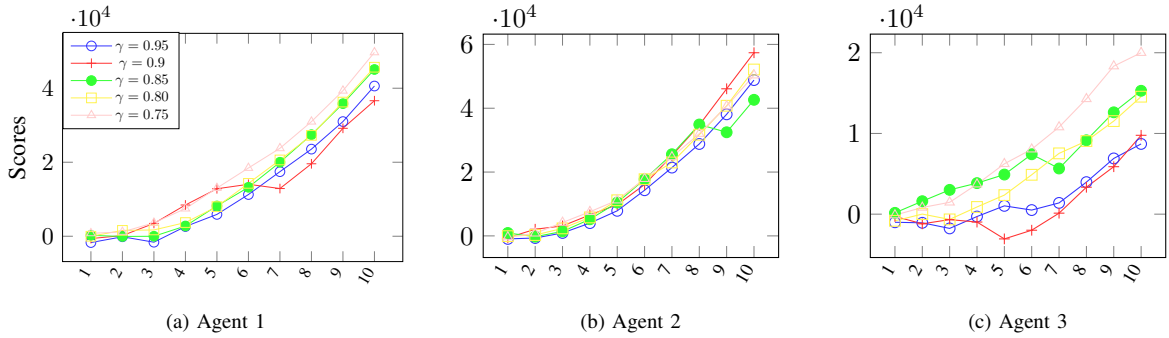
The performance of all baseline models alongside the proposed algorithm was evaluated over 10 episodes using two datasets, with the cumulative rewards at the 10th episode detailed in Table III. The Q-Learning algorithm, despite its foundational role in reinforcement learning, showed limited capability in adapting to the complexity of vital sign pattern recognition, lagging significantly behind more advanced models like the proposed DQN-based framework. This gap underscores the limitations of Q-Learning in handling high-dimensional state spaces typically encountered in health monitoring applications.

Among the baseline models, PPO, A2C, Double DQN, and DDPG demonstrated varied performance across different

agents. While PPO and A2C showed promise, particularly in agent 2's performance on the WESAD dataset, their overall efficacy was inconsistent across different vital signs and datasets. Double DQN and DDPG, on the other hand, presented a more balanced performance, suggesting their robustness but still falling short of the proposed model's results. This indicates that while these algorithms are capable of learning complex patterns, they may require further tuning or modifications to achieve optimal performance in this specific domain.

The WISEML and CA-MQL models, representing more specialized approaches to multi-agent learning in healthcare, came close to matching the proposed framework's performance. Their success points to the potential benefits of tailoring learning algorithms to the specific challenges of health monitoring. However, the proposed DRL framework outshined all baseline models across both datasets and all agents, indicating its superior ability to navigate the complexities of real-time health status monitoring using wearable sensor data. This superior performance could be attributed to the framework's efficient exploration-exploitation balance, effective reward structuring, and its capacity to handle the multidimensional data inherent in monitoring multiple vital signs.

All three learning agents were fed with physiological features such as heart rate, respiration, and temperature, respectively, from the PPG-DaLiA dataset. Based on the observation space, action space, and reward policy defined for a customized gym environment for human behavior monitoring, the learning agents were run for 10 episodes, as shown in Fig. 4. In the results, agent 1 refers to the heart rate monitoring agent, which showed a constant increase in scores for each episode for most of the subjects except subjects 5 and 6. The intermittent low scores in agent 1 performance are due to the exploration rate in DQN learning, where the algorithm tries exploring all the actions randomly instead of relying on neural networks' predictions. Similarly, agent 2 and agent 3 monitor two other physiological features, respiration and temperature, respectively. agent 2 performed better than the other two agents and achieved consistent scores for all subjects. Out of all agents, agent 3, temperature monitoring performance, was poor. This issue was traced back to the data level, where the units of the temperature thresholds in the MEWS table and the

Fig. 5: Hyper Parameters - α optimizationFig. 6: Hyper Parameters - γ optimization

input body temperature data from the dataset were different. Still, agent 3 achieved high scores in monitoring subjects 9, 8, 4, and 10.

The reward policy designed in the proposed multi-agent framework enables agents to learn the human physiological feature patterns. For example, if a subject's heart rate is 139 beats per minute, agent 1 takes Action 3 to communicate the message to MET-3. The agent will get rewarded with +10 points only if Action 3 is taken; otherwise, the agent gets negatively rewarded according to the reward policy (Table II). With this example, the results in Fig. 4 can be interpreted better. An increase in scores episode by episode, with the exception of the exploration rate, actually infers an increase in the learning curve of the agents in terms of human physiological patterns.

B. Hyper-Parameters Optimization

The DRL agents were further evaluated by hyperparameters optimization. Out of all the hyperparameters discussed in this study, two hyperparameters, learning rate (α) and discount factor (γ), were optimized for all three agents, and the results are shown in Figs. 5 and 6. The learning rate determines how much information neural networks learn in an iteration to predict action and approximate the rewards. The discount factor measures how much RL agents focus on future rewards relative to those in the immediate rewards. In Fig. 5, Figs. 5a, 5b, and 5c show the agents' performance while optimizing α of neural networks. The x-axis of the plots represents scores

(cumulative rewards) achieved by an agent in each episode shown on the y-axis. The bar plots show that the learning rate $\alpha = 0.01$ is a more optimized value in all the monitoring agents. Similarly, Figs. 6a, 6b, and 6c present the γ optimization of agent 1, agent 2, and agent 3, respectively. The discount factors $\gamma = 0.9$ and $\gamma = 0.75$ are the more optimized values for agents 2 and 3, respectively, after 10 episodes of training.

VI. DISCUSSION

This study introduces an innovative approach to patient monitoring within the unpredictable environment of healthcare settings, employing adaptive multi-agent deep reinforcement learning (DRL) to ensure timely healthcare interventions. The fluctuating nature of vital signs, crucial indicators of patient health, necessitates a robust system capable of real-time analysis and decision-making. By leveraging the sequential decision-making prowess of RL algorithms, we have established a framework where each vital sign is monitored by a dedicated DRL agent. These agents operate within a cohesive monitoring environment, guided by meticulously defined reward policies to identify and respond to potential health emergencies based on MEWS and MET standards.

A notable aspect of our research is the emphasis on the design of the observation space for each DRL agent. This design is pivotal in ensuring the accuracy and effectiveness of the learning process, as it directly impacts the agent's ability to interpret vital sign data and make informed decisions.

The challenge encountered with DRL agent 3, responsible for monitoring body temperature, underscores the importance of data consistency and the need for a harmonized observation space. The discrepancy between the temperature units in the MEWS table and the dataset highlighted a critical area for improvement, emphasizing the need for standardized data inputs to enhance agent performance.

The autonomous decision-making capability inherent in RL represents a significant advancement in supporting clinicians. By providing real-time updates on patient health, the DRL framework facilitates a proactive approach to patient care, extending its applicability beyond hospital settings to include home and specialized care environments. This adaptability is further enhanced by the strategic optimization of hyperparameters, which fine-tunes the learning process of DRL agents to achieve optimal performance. Our investigation into hyperparameters such as the learning rate and discount rate reveals the critical balance between immediate and future rewards, a balance that is essential for the effective monitoring of patient health.

Comparatively, traditional supervised learning algorithms, while accurate in predicting vital signs, fall short in dynamic healthcare environments due to their reliance on extensive labeled datasets and external supervision. The DRL approach, free from the constraints of labeled data, offers a more flexible and efficient solution for patient monitoring. However, it is essential to acknowledge the considerable effort required in data preparation and model tuning within supervised learning frameworks, which, despite their limitations, contribute significantly to the development of informed clinical decisions.

The adaptive multi-agent DRL framework proposed in this study represents a paradigm shift in patient monitoring, offering a dynamic, efficient, and scalable solution for timely healthcare interventions. The challenges and insights gleaned from this research pave the way for future advancements in the field, promising to enhance the quality of patient care through innovative technological solutions.

VII. CONCLUSION

This study has pioneered an adaptive framework for healthcare interventions using multi-agent DRL to dynamically monitor vital signs, establishing a novel approach in patient care. Through the development of a generic monitoring environment coupled with a strategic reward policy, the DRL agents were empowered to learn from and adapt to vital sign fluctuations, enabling timely interventions by healthcare professionals. Despite its innovative contributions, the research faced challenges, such as discrepancies in body temperature data scales and the absence of predictive capabilities for future vital sign trends, which limited the effectiveness of one DRL agent and the overall predictive potential of the system. Addressing these limitations, future research will focus on enhancing the framework with predictive analytics, allowing DRL agents to forecast vital sign trends and thus revolutionize patient monitoring. This advancement aims to facilitate proactive healthcare measures, significantly reducing the risk of critical health episodes and heralding a new era

in adaptive patient monitoring and healthcare management. Having said that, the future direction of our research will be focused on extending the scope of the research to predict future vital signs using multi-agent DRL.

REFERENCES

- [1] N. El-Rashidy, S. El-Sappagh, S. R. Islam, H. M. El-Bakry, and S. Abdelrazek, "Mobile health in remote patient monitoring for chronic diseases: principles, trends, and challenges," *Diagnostics*, vol. 11, no. 4, p. 607, 2021.
- [2] A. Rana, A. Dumka, R. Singh, M. K. Panda, and N. Priyadarshi, "A computerized analysis with machine learning techniques for the diagnosis of parkinson's disease: Past studies and future perspectives," *Diagnostics*, vol. 12, p. 2708, Nov. 2022.
- [3] A. Bohr and K. Memarzadeh, "The rise of artificial intelligence in healthcare applications," in *Artificial Intelligence in healthcare*, pp. 25–60, Elsevier, 2020.
- [4] R. K. Pattanayak, V. S. Kumar, K. Raman, M. M. Surya, and M. R. Pooja, "E-commerce application with analytics for pharmaceutical industry," in *Advances in Intelligent Systems and Computing*, pp. 291–298, Springer Nature Singapore, Sept. 2022.
- [5] R. Thirunavukarasu, G. P. D. C, G. R, M. Gopikrishnan, and V. Palanisamy, "Towards computational solutions for precision medicine based big data healthcare system using deep learning models: A review," *Computers in Biology and Medicine*, vol. 149, p. 106620, Oct. 2022.
- [6] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. A. Sallab, S. Yogamani, and P. Perez, "Deep reinforcement learning for autonomous driving: A survey," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, pp. 4909–4926, June 2022.
- [7] W. Zhang, N. Zhang, J. Yan, G. Li, and X. Yang, "Auto uning of price prediction models for high-frequency trading via reinforcement learning," *Pattern Recognition*, vol. 125, p. 108543, 2022.
- [8] N. Hong, C. Liu, J. Gao, L. Han, F. Chang, M. Gong, and L. Su, "State of the art of machine learning-enabled clinical decision support in intensive care units: Literature review," *JMIR Medical Informatics*, vol. 10, p. e28781, Mar. 2022.
- [9] J. Watts, A. Khojandi, R. Vasudevan, and R. Ramdhani, "Optimizing individualized treatment planning for parkinson's disease using deep reinforcement learning," in *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, IEEE, July 2020.
- [10] M. Naeem, G. Paragliola, and A. Coronato, "A reinforcement learning and deep learning based intelligent system for the support of impaired patients in home treatment," *Expert Systems with Applications*, vol. 168, p. 114285, Apr. 2021.
- [11] I. Y. Chen, S. Joshi, M. Ghassemi, and R. Ranganath, "Probabilistic machine learning for healthcare," *Annual Review of Biomedical Data Science*, vol. 4, pp. 393–415, July 2021.
- [12] S. Chen, X. Qiu, X. Tan, Z. Fang, and Y. Jin, "A model-based hybrid soft actor-critic deep reinforcement learning algorithm for optimal ventilator settings," *Information Sciences*, vol. 611, pp. 47–64, Sept. 2022.
- [13] M. Rastogi, D. M. Vijarana, and D. N. Goel, "Role of machine learning in healthcare sector," *SSRN Electronic Journal*, 2022.
- [14] B. Mahesh, "Machine learning algorithms-a review," *International Journal of Science and Research (IJSR)*, [Internet], vol. 9, pp. 381–386, 2020.
- [15] S. A. Alsareii, M. Awais, A. M. Alamri, M. Y. AlAsmari, M. Irfan, N. Aslam, and M. Raza, "Physical activity monitoring and classification using machine learning techniques," *Life*, vol. 12, no. 8, p. 1103, 2022.
- [16] M. Oyeleye, T. Chen, S. Titarenko, and G. Antoniou, "A predictive analysis of heart rates using machine learning techniques," *International Journal of Environmental Research and Public Health*, vol. 19, p. 2417, Feb. 2022.
- [17] M. Luo and K. Wu, "Heart rate prediction model based on neural network," *IOP Conference Series: Materials Science and Engineering*, vol. 715, p. 012060, Jan. 2020.
- [18] L. M. Dang, K. Min, H. Wang, M. J. Piran, C. H. Lee, and H. Moon, "Sensor-based and vision-based human activity recognition: A comprehensive survey," *Pattern Recognition*, vol. 108, p. 107561, Dec. 2020.
- [19] T. Sheng and M. Huber, "Unsupervised embedding learning for human activity recognition using wearable sensor data," in *The Thirty-Third International Flairs Conference*, 2020.
- [20] I. H. Sarker, "Machine learning: Algorithms, real-world applications and research directions," *SN computer science*, vol. 2, no. 3, p. 160, 2021.
- [21] R. Lou, Z. Lv, S. Dang, T. Su, and X. Li, "Application of machine learning in ocean data," *Multimedia Systems*, Feb. 2021.

- [22] D. Tirumala, A. Galashov, H. Noh, L. Hasenclever, R. Pascanu, J. Schwarz, G. Desjardins, W. M. Czarnecki, A. Ahuja, Y. W. Teh, *et al.*, "Behavior priors for efficient reinforcement learning," 2020.
- [23] M. Janssen, C. LeWarne, D. Burk, and B. B. Averbeck, "Hierarchical reinforcement learning, sequential behavior, and the dorsal frontostriatal system," *Journal of Cognitive Neuroscience*, vol. 34, no. 8, pp. 1307–1325, 2022.
- [24] K. Tsiakas, M. Papakostas, M. Theofanidis, M. Bell, R. Mihalcea, S. Wang, M. Burzo, and F. Makedon, "An interactive multisensing framework for personalized human robot collaboration and assistive training using reinforcement learning," in *Proceedings of the 10th International Conference on Pervasive Technologies Related to Assistive Environments*, ACM, June 2017.
- [25] A. Kubota and L. D. Riek, "Methods for robot behavior adaptation for cognitive neurorehabilitation," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 5, pp. 109–135, May 2022.
- [26] A. Lisowska, S. Wilk, and M. Peleg, "From personalized timely notification to healthy habit formation: a feasibility study of reinforcement learning approaches on synthetic data.," in *SMARTERCARE@ AI* IA*, pp. 7–18, 2021.
- [27] E. Yom-Tov, G. Feraru, M. Kozdoba, S. Mannor, M. Tennenholtz, and I. Hochberg, "Encouraging physical activity in patients with diabetes: Intervention using a reinforcement learning system," *Journal of Medical Internet Research*, vol. 19, p. e338, Oct. 2017.
- [28] T. Li, Z. Wang, W. Lu, Q. Zhang, and D. Li, "Electronic health records based reinforcement learning for treatment optimizing," *Information Systems*, vol. 104, p. 101878, Feb. 2022.
- [29] J. Guo, Q. Liu, and E. Chen, "A deep reinforcement learning method for multimodal data fusion in action recognition," *IEEE Signal Processing Letters*, vol. 29, pp. 120–124, 2022.
- [30] C. Yu, J. Liu, S. Nemati, and G. Yin, "Reinforcement learning in healthcare: A survey," *ACM Comput. Surv.*, vol. 55, nov 2021.
- [31] V. Signs, "Canberra hospital and health services clinical procedure," 2021.
- [32] C. Wang, J. Wang, J. Wang, and X. Zhang, "Deep-reinforcement-learning-based autonomous uav navigation with sparse rewards," *IEEE Internet of Things Journal*, vol. 7, no. 7, pp. 6180–6190, 2020.
- [33] A. Reiss, I. Indlekofer, P. Schmidt, and K. V. Laerhoven, "Deep PPG: Large-scale heart rate estimation with convolutional neural networks," *Sensors*, vol. 19, p. 3079, July 2019.
- [34] P. Schmidt, A. Reiss, R. Duerichen, C. Marberger, and K. Van Laerhoven, "Introducing wesad, a multimodal dataset for wearable stress and affect detection," in *Proceedings of the 20th ACM international conference on multimodal interaction*, pp. 400–408, 2018.
- [35] P. Mallozzi, E. Castellano, P. Pelliccione, G. Schneider, and K. Tei, "A runtime monitoring framework to enforce invariants on reinforcement learning agents exploring complex environments," in *2019 IEEE/ACM 2nd International Workshop on Robotics Software Engineering (RoSE)*, IEEE, May 2019.
- [36] Y.-J. Chen, D.-K. Chang, and C. Zhang, "Autonomous tracking using a swarm of UAVs: A constrained multi-agent reinforcement learning approach," *IEEE Transactions on Vehicular Technology*, vol. 69, pp. 13702–13717, Nov. 2020.
- [37] M. C. Schippers and D. C. Rus, "Optimizing decision-making processes in times of covid-19: using reflexivity to counteract information-processing failures," *Frontiers in psychology*, vol. 12, p. 650525, 2021.

5.2 Summary

This chapter highlights the transformative potential of the multi-agent DRL framework in redefining the landscape of patient monitoring. It elaborates on the framework's adeptness in harnessing real-time physiological and behavioural data to enhance the precision and timeliness of health assessments and interventions. The narrative underscores this approach's superiority over traditional models in terms of adaptability, scalability, and predictive accuracy. It further reflects on the broader implications of such AI-driven methodologies in improving clinical outcomes, minimizing healthcare errors, and streamlining resource utilization, thereby setting a visionary path for the next generation of healthcare innovations.

CHAPTER 6: PAPER 5 - PDRL: MULTI-AGENT BASED REINFORCEMENT LEARNING FOR PREDICTIVE MONITORING

6.1 Introduction

This chapter introduces a cutting-edge Predictive Deep Reinforcement Learning (PDRL) framework, employing multiple Reinforcement Learning (RL) agents in a time series forecasting environment for predictive monitoring. By integrating Deep Q Network (DQN) agents, the framework innovatively predicts future states of complex environments, leveraging a well-defined reward policy to maximize learning efficiency. The chapter explores the application of this framework in monitoring vital signs such as heart rate, respiration, and temperature, utilizing a Bi-LSTM model for prediction. This novel approach aims to surpass traditional monitoring systems by enabling adaptive decisions in dynamic, uncertain environments, showcasing the PDRL framework's potential in healthcare and beyond.

PDRL: Multi-Agent based Reinforcement Learning for Predictive Monitoring

Thanveer Shaik, Xiaohui Tao, Lin Li, Haoran Xie, U R Acharya, Raj Gururajan, Xujuan Zhou

Abstract—Reinforcement learning has been increasingly applied in monitoring applications because of its ability to learn from previous experiences and can make adaptive decisions. However, existing machine learning-based health monitoring applications are mostly supervised learning algorithms, trained on labels and they cannot make adaptive decisions in an uncertain complex environment. This study proposes a novel and generic system, predictive deep reinforcement learning (PDRL) with multiple RL agents in a time series forecasting environment. The proposed generic framework accommodates virtual Deep Q Network (DQN) agents to monitor predicted future states of a complex environment with a well-defined reward policy so that the agent learns existing knowledge while maximizing their rewards. In the evaluation process of the proposed framework, three DRL agents were deployed to monitor a subject's future heart rate, respiration, and temperature predicted using a BiLSTM model. With each iteration, the three agents were able to learn the associated patterns and their cumulative rewards gradually increased. It outperformed the baseline models for all three monitoring agents. The proposed PDRL framework is able to achieve state-of-the-art performance in the time series forecasting process. The proposed DRL agents and deep learning model in the PDRL framework are customized to implement the transfer learning in other forecasting applications like traffic and weather and monitor their states. The PDRL framework is able to learn the future states of the traffic and weather forecasting and the cumulative rewards are gradually increasing over each episode.

Index Terms—Reinforcement Learning, Timeseries Forecasting, Monitoring, Decision Making, Behavior Patterns

I. INTRODUCTION

Data mining has been widely adopted for analysis and knowledge discovery in databases. This process involves data management, data preprocessing, modeling, and results in inferences and extracting latent data patterns [1]. Early warning systems based on data mining have enabled applications to perform a risk analysis, monitoring and warning, and a response capability [2]. Using existing knowledge or a set of indicators can assist domains such as healthcare [3] and

Thanveer Shaik, Xiaohui Tao are with the School of Mathematics, Physics & Computing, University of Southern Queensland, Toowoomba, Queensland, Australia (e-mail: Thanveer.Shaik@usq.edu.au, Xiaohui.Tao@usq.edu.au).

Lin Li is with the School of Computer and Artificial Intelligence, Wuhan University of Technology, China (e-mail: cathylilin@whut.edu.cn)

Haoran Xie is with the Department of Computing and Decision Sciences, Lingnan University, Tuen Mun, Hong Kong (e-mail: hrxie@ln.edu.hk)

U R Acharya is with School of Mathematics, Physics & Computing, University of Southern Queensland, Toowoomba, Queensland, Australia (e-mail: Rajendra.Acharya@usq.edu.au).

Raj Gururajan is with School of Business, University of Southern Queensland, Springfield, Queensland, Australia (e-mail: Raj.Gururajan@usq.edu.au).

Xujuan Zhou is with School of Business, University of Southern Queensland, Springfield, Queensland, Australia (e-mail: Xujuan.Zhou@usq.edu.au).

education [4] to design decision support systems. Radanliev et al. [5] study used data mining to investigate scientific research response to the COVID-19 pandemic and to review key findings on how early warning systems developed in previous epidemics responded to contain the virus.

Traditional unsupervised learning techniques discover underlying patterns for knowledge discovery in unlabelled data using association rule mining and clustering techniques [6]. Supervised learning strategies learn from labeled data to classify or predict patients' physical activities and vital signs [7]. However, these methodologies are highly dependent on data and can only observe the data and present possible decisions in response, they cannot take actions based on observations. Reinforcement learning (RL) deploys a learning agent in an uncertain, complex environment that explores or exploits the environment with its actions and learns the data based on its experience [8], [9]. This allows the learning agent to gain rewards based on learning and its actions.

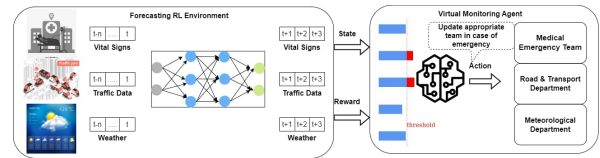


Fig. 1: PDRL framework to monitor different forecasting applications and alert appropriate emergency teams.

RL is used in dynamic domains such as stockmarket trading [10], [11], traffic prediction [12], [13], and weather forecast [14] to tackle decision-making problems using agent-environment interaction samples and potentially delayed feedback [15] that could also be applied to healthcare applications. In the healthcare domain, chronically diseased such as Parkinson's disease [16] and critical care patients often require long-term dynamic treatment regimes with the timely intervention of clinicians to avoid unwanted outcomes [17]. Zeng et al. [18] proposed an RL algorithm to optimize post-operation warfarin anticoagulation dosage. The RL results outperformed conventional clinical practice using rule-based algorithms. Existing patient monitoring applications based on RL primarily focus on prescribing the timing and dosage of medications [19] so that patients are administered take the right medication at the right time [20]. Chen et al. [21] described Probabilistic machine learning models such as RL using the analogy of an ICU clinician (learning agent) to monitor a patient (environment) via actions like ventilation and observing the changes in the environment (patient's state) to

make subsequent decisions that achieve the goal of discharging the patient successfully.

The research problem addressed here is that of being able to monitor the predicted state of an environment and take appropriate actions to avoid an emergency. Traditional supervised learning strategies can classify or predict based on their training but cannot monitor and alert the appropriate team for timely interventions. To assist with tracking the environment state and monitor certain parameters we have designed a virtual generic forecasting environment with observation space, actions, and rewards policy with multiple deep learning agents. Deploying a single learning agent to monitor all the parameters would complicate the environment as there are different thresholds set up for each of the parameters in an environment. For example, the learning agents learn different threshold levels of each vital sign in modified early warning scores (MEWS) [22] based on previous iterations and rewards being accumulated for its actions. Well-trained RL agents are capable of monitoring a patient's vital signs such as heart rate, respiration rate, and temperature, and alerting the corresponding clinical team if the vital signs fall outside any of the predefined thresholds [23].

Modeling forecasting applications, such as vital signs prediction, traffic prediction, and weather prediction, as an RL environment can enable RL agents to monitor tasks. RL agents can learn from historical data and interact with the environment to make real-time decisions based on the predicted states or actions. This approach can be used to develop intelligent monitoring systems that can adapt to changing conditions, optimize actions, and make informed decisions in complex and dynamic environments. By using RL for monitoring, it is possible to automate and optimize monitoring processes in various domains, leading to more efficient and effective monitoring outcomes. In this study, the RL environment is configured with a deep learning model to predict future states, which are then monitored by an RL agent.

The aim of this research is to create a multi-agent framework that utilizes deep reinforcement learning (DRL) agents to monitor and learn data patterns for various parameters. Each parameter will have its own DRL agent, responsible for monitoring, learning, and alerting respective teams if the parameters deviate from predefined thresholds as shown in Fig. 1. Conventional RL methodology is an agent performing a task for a transition from one state to another state, where this action might reward the agent either positively or negatively. In this study, a novel approach is taken to assign rewards so that the RL agents learn data patterns. An agent gets rewarded for predicting an action and performing the action in its current state. The rewards are designed in such a way that the learning agents are penalized for predicting the wrong actions. To learn behaviors we follow the Reward-is-enough hypothesis [24] being that the learning agent always tries to maximize the rewards based on their previous actions. The contributions of this study are as follows:

- A generic monitoring environment accommodates multiple agents to monitor the states of a forecasting environment.

- Proposed a model-free gaming agent to learn the existing knowledge and monitor underlying data patterns adaptively.
- Transfer learning approach for time series forecasting applications such as patients' health status, traffic, and weather using the multi-agents in the PDRL environment.

The paper is organized as follows: Section II presents the related works in the RL community to learning human behavior patterns and application in the healthcare domain. Research problem formulation and the proposed multi-agent PDRL framework have been detailed in Section III. In Section IV, the proposed methodology is evaluated on 10 different subject vital signs, and baseline models are discussed. In Section V, the results of the proposed approach are compared with baseline models, and hyper-parameter optimization of the learning rate and discount factor are discussed. The comparison between the supervised approach and the RL approach is discussed in Section VII. Section VIII concludes the paper with limitations and future work.

II. RELATED WORKS

A. Data Mining in Early Warning Systems

Akçapınar et al. [25] proposed a study that uses interaction data from online learning to predict the academic performance of students at end of term by using the kNN algorithm which predicted unsuccessful students at an 89% rate. It also suggests that performance can be predicted in 3 weeks with 74% accuracy, useful for early warning systems and selecting algorithms for analysis of educational data. Cano et al. [26] presented a multiview early warning system for higher education that uses comprehensible Genetic Programming classification rules to predict student performance, specifically targeting underrepresented and underperforming student populations. The authors integrated various student information sources and have interfaces to provide personalized feedback to students, instructors, and staff. In healthcare, Hussain-Alkhateeb et al. [27] conducted a scoping review to summarize existing evidence of early warning systems for outbreak-prone diseases such as chikungunya, dengue, malaria, yellow fever, and Zika. It found that while many studies showed the quality performance of their prediction models, only a few presented statistical prediction validity of early warning systems. It also found that no assessment of the integration of the early warning systems into a routine surveillance system could be found and that almost all early warning systems tools require highly skilled users with advanced statistics. Spatial prediction remains a limitation with no tool currently able to map high transmission areas at small spatial levels. Liu et al. [28] conducted a study on the use of data mining technology to analyze college students' psychological problems and mental health. The authors use intuitionistic fuzzy reasoning judgment, analytic hierarchy process, and expert scoring method to construct a model for studying college students' online public opinion and use data mining techniques such as the decision tree algorithm and Apriori algorithm to analyze students' psychological problems and provide decision-making support information for the school psychological counseling center.

In [29], the authors talk about the challenges of managing large amounts of data from connected devices and how data mining can be used to extract valuable information. It also mentions the use of fog computing technology to improve the quality of service in healthcare applications. The article suggests wearable clinical devices for continuous monitoring of individual health conditions as a solution for chronic patients. An EWS (Early Warning System) for heavy precipitation using meteorological data from Automatic Weather Stations (AWSs) is proposed by Moon et al. [30] and its performance are measured by various criteria.

B. RL Monitoring

In the healthcare domain, Lisowska et al. [31] developed a digital behavior-change intervention to help cancer patients build positive health habits and enhance their lifestyles. The authors used reinforcement learning to learn the appropriate time to send the intervention prompts to the patients. Furthermore, effective prompt policies to perform an activity have been used in a custom patient environment. Three RL approaches Deep Q-Learning (DQL), Advance Actor-Critic (A2C) and proximal policy optimization (PPO) have been used to suggest a virtual coach for sending a prompt. Similarly, personalized messages enhance physical activity in type 2 diabetes patients [32]. Li et al. [33] proposed an electronic health records (EHRs)-based reinforcement learning approach for sequential decision-making tasks. The authors used a model-free DQN algorithm to learn the patients' data and provide clinical assistance in decision-making. Co-operative multi-agent RL has been deployed using value compositions and achieved better results. RL decision-making ability can be used for human activity recognition as per [34]. The authors proposed a dynamic weight assignment network architecture in which twin delayed deep deterministic (TD3) [35] algorithm was inspired by Deep Deterministic Policy Gradient algorithm (DDPG), Actor-Critic, and DQN algorithms. RL agents tend to learn effective strategies while the sequential decision-making process using trial-and-error interactions with their environments [15].

C. Mimic Human Behavior Patterns

Tirumala et al. [36] discussed learning human behavior patterns and capturing common movement and interaction patterns based on a set of related tasks and contexts. The authors discussed probabilistic trajectory models to learn behavior priors and proposed a generic framework for hierarchical reinforcement learning (HRL) concepts. Janssen et al. [37] suggested breaking a complex task such as biological behavior into more manageable subtasks. HRL is able to combine sequential actions into a temporary option. The authors discussed how biological behavior is conceptually analogous to options in HRL. Tsiakas et al. [38] proposed a human-centric cyber-physical systems (CPS) framework for personalized human-robot collaboration and training. This framework focuses on monitoring and assessment of human behavior. Based on the multi-modal sensing framework, the authors aimed for effective human attention prediction with a minimal and

least intrusive set of sensors. Kubota et al. [39] examined how robots can adapt the behavior of people with cognitive impairments as part of cognitive neuro-rehabilitation. A variety of robots in therapeutic, companion, and assistive applications has been explored in the study.

In addition to RL applications in the gaming industry, a great deal of research is being conducted using RL to learn and mimic human behavior and also to deploy socially assistive robots. However, physical robots with human interaction capability cannot be deployed at sensitive locations like hospitals, educational institutions, elderly home care, and mental health facilities as they might cause safety issues for students, patients, carers, and medical staff [40]. Existing monitoring applications with supervised or unsupervised learning cannot cope with uncertain events in a dynamic and complex environment but the supervised approach is well known for its achievement in regression problems. Virtual robots with adaptive learning abilities can be deployed to overcome these issues to monitor and learn the predicted states from a supervised learning model. In this study, we developed a custom monitoring environment to learn behavior patterns from predicted states by designing rewards according to the applications in certain domains. The framework is capable of alerting the appropriate team based on the level of severity and assisting in timely intervention.

III. PDRL MONITORING FRAMEWORK

In this section, custom behavior forecasting RL environment and monitoring agents of the proposed predictive deep reinforcement learning (PDRL) framework has been discussed in detail along with problem formulation. This study is to learn data patterns in an uncertain environment by monitoring its current state, taking appropriate actions, and getting rewarded for the actions as shown in Fig. 2.

A. Problem Formulation

Knowledge discovery in data mining can be achieved by utilizing previously known and potentially useful information extracted from observations and various data mining techniques. The research problem involves designing a multi-agent framework to monitor data in a forecasting environment and uncover underlying patterns in relation to thresholds established by known knowledge. For example, consider a scenario where a client, $c_n \in C$ where $n = 1, 2, 3, \dots, N$, $C = |N|$ is the number of clients with wearable sensors on their bodies to track and forecast vital signs. Each client c has a set of DRL learning agents to monitor predicted vital signs and alert the appropriate emergency team if health parameters exceed modified early warning scores (MEWS) [22].

To formulate this problem, a customized RL forecasting environment needs to be configured with an innovative reward policy that links the current state and agent actions to learn data patterns while maximizing their rewards. This can be achieved based on a Markov Decision Process (MDP), which can be defined as a 5-tuple $M = (S, A, P, R, \gamma)$, where: S is a finite state space, with $s_t \in S$ denoting the state of an agent at time t , A is a set of actions defined for the agent, with

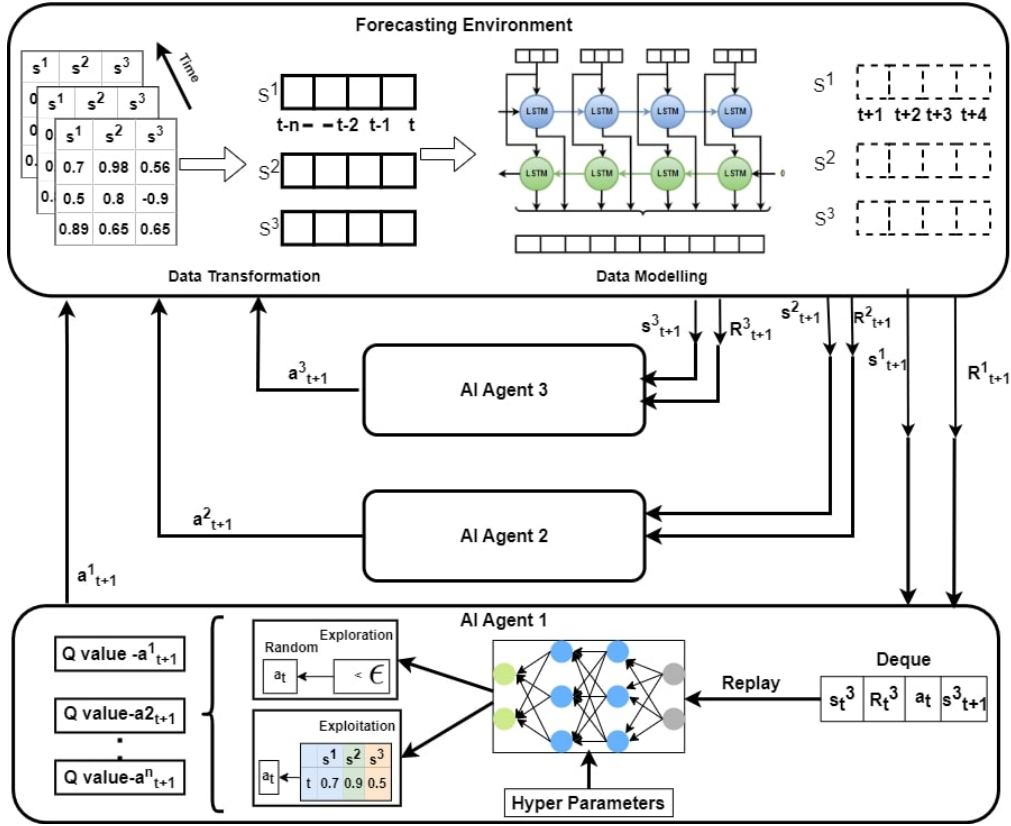


Fig. 2: PDRL monitoring Framework

$a_t \in A$ denoting the action taken by the agent at time t , P is a Markovian transition function as shown in Equation 1, which denotes how the agent transits from state s to state s' while performing an action a , R is a reward function, which returns an immediate reward $R(s, a)$ for the action a taken in a state s defined in Equation 2, γ is a discount factor that focuses on immediate rewards instead of future rewards. It remains between 0 and 1.

$$P_a(s, s') = \Pr(s_{t+1} = s' \mid s_t = s, a_t = a) \quad (1)$$

$$R(s_t, a_t) = \sum_{t=0}^{\infty} \gamma^t r_t, \quad (2)$$

The next step of the research problem is to compute the optimal reinforcement learning policy $\pi : S \times A \rightarrow [0, 1]$, which helps in predicting the probability that an agent selects an appropriate action $a_t \in A$ in a specific state $s_t \in S$ at time t . To do this, the action value (Q-function) needs to be updated in each iteration and can be defined in Equation 3. $Q^{new}(s_t, a_t)$ is the new output of the action a_t and state value s_t . α is the learning rate, which determines how much information from the previously computed Q-value is used for the given state-action pair. γ is a discount factor that focuses on immediate rewards instead of future rewards, and it remains between 0 and 1.

$$Q^{new}(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \cdot (r_t + \gamma \max_a Q(s_{t+1}, a)) \quad (3)$$

B. Forecasting Environment

In this section, forecasting applications are modeled as a customized RL forecasting environment based on MDP has been designed with observation space S , and action space A for learning agents to take appropriate actions, and it rewards R for the agents' actions. The forecasting environment enables a deep learning model to forecast the future states at time $t+1, t+2, t+3, t+4$ based on the training data at previous timesteps $t-n, \dots, t-2, t-1, t$ in the proposed PDRL framework as shown in Fig. 2.

1) **Forecasting States:** In this study, forecasting the future states of an environment is a supervised time series learning approach. For this task, the recurrent neural network (RNN) model variant, the bidirectional LSTM(Bi-LSTM) model is deployed. Mathematically, the Bi-LSTM model is defined in Equation 4. A regularization method, dropout [41] was used to exclude activation and weight updates of recurrent connections from LSTM units probabilistically.

$$y(x) = \sum_{i=1}^n \text{Activation1}(b + w_i x_i) \quad (4)$$

$$\text{Bi-LSTM}(y) = \text{Activation2}\left(\frac{e^{y_i}}{\sum_{j=1}^k e^{y_j}}\right)$$

where b : Bias added on each hidden layer, x : Input value, w : Weights added on each hidden layer, y : Output value from each neuron, *Activation1*: Activation functions on input and hidden layers, *Activation2*: Activation function on the output layer.

Based on the forecasted states, the following components are configured in the forecasting RL environment.

2) **Observation Space**: The environment shown in Fig. 2 has state $s_t^i \in S$ where $i = 0, 1, 2, \dots, n$, observations in a state at time t . The idea is to split the state into observations and forecast the states based on the time series data. The predicted states are getting assigned to multi-agents. Furthermore, considering a single agent to monitor the multiple states of a complex environment might lead to a sparse rewards challenge where the environment rarely produces a useful reward and limits agent learning. Hence, multiple agents need to be deployed to monitor multiple states. To determine the expected return E_π of a policy π in a state s can be defined in state-value Equation 5 adopting multi-agent where $i = 0, 1, 2, 3, \dots, n$ is a finite number of observations n in a state.

$$V^\pi(s^i) = E_\pi \left\{ \sum_{t=0, i=0}^{\infty, n} \gamma^t R(s_t, \pi(s_t)) | s_0^i = s \right\} \quad (5)$$

3) **Action Space**: Defining actions for the RL agent in the environment is the most critical part of the RL process as it directly reflects the capacity of RL agents in adaptive learning. In this study, a discrete set of actions are proposed for a continuous observation space. Each of these actions will be chosen by agents based on the current state of the forecasting environment. The expected return E_π for taking an action a in a state s under a policy π can be measured using action-value function $Q_\pi(s, a)$ Equation 6.

$$Q^\pi(s, a) = E_\pi \left\{ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, \pi(s_t)) | s_0 = s, a_0 = a \right\} \quad (6)$$

Actions within an RL environment vary depending on the application. For instance, in a health monitoring application, the patient's health status may change based on vital signs such as heart rate, blood pressure, respiratory rate, and more. Actions can be configured based on these vital signs, utilizing modified early warning scores [22]. Threshold levels defined for vital signs, such as heart rate, can be used to measure the level of emergency, and appropriate alerts can be triggered accordingly. For example, if the heart rate exceeds a predefined threshold, a high-level emergency alert may be generated, while a lower-level alert may be issued for a heart rate that is moderately deviating from the normal range. This way, RL agents can take actions based on the defined thresholds to

ensure timely and appropriate responses in health monitoring applications.

4) **Rewards**: RL goals can be represented by cumulative rewards achieved by learning agents with their actions in an environment. Conventional reinforcement learning rewards an agent based on their action for a transition from a state s_t to s_{t+1} . In this study, the goal of the learning agent is to learn underlying data patterns in terms of states of a forecasting environment and this can be achieved by the efficient design of a reward policy. The reward policy defined for this study is calculated using Equation 7 where agents get positively rewarded only if an agent monitors the state and predicts the right action from the action space. Otherwise, the agent will be negatively rewarded.

$$R(s_t, a_t) = \begin{cases} +reward & \text{if action is appropriate} \\ -reward & \text{if action is not appropriate} \end{cases} \quad (7)$$

Algorithm 1 Forecasting Environment

Require: time series data $\mathcal{D} = \{s_{t-n}, \dots, s_{t-2}, s_{t-1}, s_t\}$; a set of labels $\mathcal{K} = \{1, 2, \dots, K\}$

Ensure: Predicted time series data of \mathcal{K} , a set of labels, in the form of states $\{s_{t+1}, s_{t+2}, s_{t+3}, s_{t+4}\}$.

```

1: Define forecast_model  $\leftarrow$  Bi-LSTMModel
2: Train(forecast_model)  $\leftarrow$  forecast_model(D)
3:  $\{s_{t+1}, s_{t+2}, s_{t+3}, s_{t+4}\} \leftarrow$  forecast_model(predict)
4: Initialization : observation_space =  $s_t^i \in S$ , action_space =  $a_t \in A$ , reward  $R$ 
5: Set monitor_length =  $N$ 
6: if action is appropriate then
7:    $R \leftarrow +reward$ 
8: else
9:    $R \leftarrow -reward$ 
10: end if
11: monitor_length  $\leftarrow N - 1$ 
12:  $s_{t+1} \leftarrow s_t(\text{monitor\_length})$ 
13: if  $N = 0$  then
14:   done = True
15: else
16:   done = False
17: end if
18: visualize( $a_t, R, \text{vital signs}$ )
19: initial_state  $\leftarrow s_t[0]$  ▷ reset environment

```

The algorithm 1 presents the forecasting environment where observation space, action space, and reward policy have been configured based on the predicted states. Lines 1-3 in the algorithm define the deep learning model to train and predict the time series forecasting data. Lines 4-5 initialize the class and set boundaries for the observation space, action space, rewards, and monitoring length. Lines 6-10 explain the reward policy for the actions of the learning agent and how the agents get rewarded either positively or negatively. Lines 11-19 present monitoring length and visualization of agent performance and reset the environment if needed.

C. Learning Agent

In the proposed PDRL framework, the game learning agent DQN algorithm was used. The algorithm was introduced by Google's DeepMind for playing Atari game to play games by just observing the screen without any training or prior knowledge about those games. In this algorithm, the Q-Learning functions' approximation will be computed using neural networks, and the learning agent gets rewarded based on the neural networks' prediction of the right action for the current state. The reward policy has been discussed in detail in Section III-B.

1) *Q-Function Approximation*: The neural networks model used in this study to approximate the rewards has three layers input layer, a hidden layer, and an output layer. The input layer has a node for each vital sign of a state, the output layer has a node for each action in the action space. The model is configured with parameters such as the relu activation function, mean square error as loss function, and Adam optimizer. The model gets trained with the state and its corresponding reward. Upon training, the model is able to predict reward for the current state.

The learning agent performs an action $a_t \in A$ for a transition from state s_t to s'_t and achieves a reward R for the action. In this transition process, the maximum of the Q-function in Equation 6 is calculated, and the discount of the calculated value uses a discount factor γ to suppress future rewards and focus on immediate rewards. The discounted future reward is added to the current reward to get the target value. The difference between the current prediction from the neural networks and the calculated target value provides a loss function. The loss function is a deviation of the predicted value from the target value and it can be estimated from Equation 8. The square of the loss function allows for the punishment of the agent for a large loss value.

$$loss = \left(\underbrace{R + \gamma \cdot \max(Q^{\pi^*}(s, a))}_{target_value} - \underbrace{Q^{\pi}(s, a)}_{predicted_value} \right)^2 \quad (8)$$

2) *Memorize and Replay*: A simple neural network has the challenge of limited memory and forgetting previous observations once new observations overwrite them. To retrain the model, previous observations can be stored in an array as an experience e that acts as a memory and appends the current state, action, reward, and next state to the memory at time t as $e_t = (s_t, a_t, r_t, s_{t+1})$. A sample of previous observations from the memory is randomly selected to train the neural networks using the replay method. In this study, a batch size of 32 previous observations was to retrain the neural network model.

3) *Exploration and Exploitation*: Exploration and exploitation are two contradictory concepts in RL where exploration is the selection of actions randomly that have never been performed and exploring more possibilities. Exploitation is to select known actions from existing knowledge and previous experiences to maximize the rewards. To balance exploration and exploitation, there are different strategies such as greedy algorithm, epsilon-greedy algorithm, optimistic initialization, and decaying epsilon-greedy algorithm. This study controls

the exploration rate by multiplying decay by the exploration rate. This reduces the number of explorations in the execution as the agents learn the patterns and maximize their rewards to get high scores. While the neural networks model is getting retrained with previous experiences in the replay, the decay gets multiplied by the exploration rate based on how well the agent can predict the right actions. All these parameters are defined as hyper-parameters to DQN learning agents.

Algorithm 2 Learning Agent

```

1: Initialize  $\gamma, \epsilon, \epsilon_{decay}, \epsilon_{min}, memory = \emptyset, batch\_size$ 
2: Define  $model \leftarrow NeuralNetworkModel$ 
3:  $memory \leftarrow append(s_t, a_t, R, s_{t+1})$ 
4: if  $np.random.rand < \epsilon$  then ▷ Exploration
5:    $action\_value \leftarrow random(a_t)$ 
6: else ▷ Exploitation
7:    $action\_value \leftarrow model.predict(s_t)$ 
8: end if
9:  $minibatch \leftarrow random(memory, batch\_size)$ 
10: for  $s_t, a_t, R, s_{t+1}, done$  in  $minibatch$  do
11:    $target \leftarrow R$ 
12:   if not done then
13:      $target \leftarrow R + \gamma(max(model.predict(s_{t+1})))$ 
14:   end if
15:    $target\_f \leftarrow model.predict(s_t)$ 
16:    $target\_f[a_t] \leftarrow target$ 
17:    $model.fit(s_t, target\_f)$ 
18: end for
19: if  $\epsilon \geq \epsilon_{min}$  then
20:    $\epsilon *= \epsilon_{decay}$ 
21: end if

```

The methods to perform function approximation, memorize, replay, exploration, and exploitation are enclosed in a Learning Agent algorithm 2. Line 1-2 initializes all the hyper-parameters required for the agent and a deep learning model for Q-function approximation. Line 3 explains the memorize and replay part to store neural-network experience where state, action, reward, and next state will be stored to retrain the model using the replay method. Lines 4-8 in the algorithm are responsible to predict an action either exploration or exploitation methods. Lines 9-21 explain a batch of previous experiences from memory will be retrieved to process and retrain the neural networks model based on the hyper-parameters defined earlier.

The Algorithm 3 is an extension to the previous two Algorithms 1 2 and implements the proposed generic PDRL monitoring framework. The inputs of the algorithm are a set of subjects and their vital signs along with the number of episodes the agents have to be trained. The algorithm 3 outputs the learning agents score which is the cumulative sum of rewards achieved in each episode. Lines 1-2 create objects of ForecastingEnvironment and LearningAgent. Lines 3-17 are nested for loops with conditional statements to check if the episode is completed or not. The outer loop is to iterate each episode while resetting the environment to initial values and score to zero. The inner loop is to iterate timesteps which denote the time of the current state and call the methods

Algorithm 3 Proposed Multiple Agents Monitoring Framework Implementation

Require: Input:

$\mathcal{C} = 1, 2, \dots, C$: set of subjects
 $\mathcal{V} = 1, 2, \dots, V$: set of vital signs
 $\mathcal{M} = 1, 2, \dots, M$: number of episodes

Ensure: Output: Rewards achieved by Agents in each episode.

```

1:  $env \leftarrow ForecastingEnvironment()$     ▷ Algorithm 1
2:  $agent \leftarrow LearningAgent()$         ▷ Algorithm 2
3: for episode  $m \in \mathcal{M}$  do
4:    $state \leftarrow env.reset()$ 
5:    $score = 0$ 
6:   for time in range(timesteps) do
7:      $a_t \leftarrow agent.action(s_t)$ 
8:      $s_{t+1}, R, done \leftarrow env.step(a_t)$ 
9:      $agent.memorize(s_t, a_t, R, s_{t+1})$ 
10:     $s_t \leftarrow s_{t+1}$ 
11:    if done then
12:       $print(m, score)$ 
13:       $break$ ;
14:    end if
15:  end for
16:   $agent.replay(batch\_size)$ 
17: end for

```

defined in Algorithm 1 2 to predict action for the current state, to reward the agent for predicted action, to retrieve next_state, and to memorize the previous experiences. Finally, the replay method will be called to retain the neural network model with the stored previous experiences.

IV. EXPERIMENT

The multi-agent PDRL monitoring framework proposed in this study has been evaluated with experiments on different datasets related to healthcare, traffic, and weather forecasting. In healthcare, vital signs such as heart rate, respiration, and temperature retrieved from a patient are processed into time series data. The forecasting environment is responsible to learn the time series data and forecast future vital signs in the next 15 minutes, 30 minutes, 45 minutes, and 60 minutes. The predicted data is passed to multiple DRL agents with one vital sign for each agent. The agents are responsible to monitor the vital signs in each iteration and take appropriate actions. For this task, each agent gets rewarded as per the reward policy discussed in the previous section III. As they aim to increase their accumulated rewards, all the agents learn each vital sign pattern of patients and collectively monitor the patient's health status. The isolated DRL agents monitor their vital signs independently and update the corresponding medical emergency team(MET) at the right time.

A. Dataset

- PPG-DaLiA [45]: The dataset contains physiological and motion data of 15 subjects, recorded from both a wrist-worn device and a chest-worn device, while the subjects

TABLE I: Proposed Multi-Agent PDRL framework performance is compared with other baseline models

		MAE	MAPE	RMSE
ELMA [42]	15 Min	6.2	13.91	8.75
	30 Min	6.2	13.91	8.75
	45 Min	6.2	13.91	8.75
	60 Min	6.13	13.91	8.67
GRU [43]	15 Min	0.95	5.47	1.25
	30 Min	0.95	5.48	1.25
	45 Min	0.97	5.51	1.27
	60 Min	0.98	5.5	1.28
GNN-Based Multi-Agent [44]	15 Min	3.64	8	2.46
	30 Min	3.99	3.47	2.58
	45 Min	4.33	4.53	2.69
	60 Min	5.73	5.27	3.09
Multi-Agent PDRL Fraemwork (Ours)	15 Min	0.44	2.6	0.85
	30 Min	0.65	2.67	1.05
	45 Min	0.52	3.17	0.93
	60 Min	0.53	5.39	0.95

were performing a wide range of activities under close to real-life conditions.

- Traffic Dataset [46]: The dataset includes 47 features such as a historical sequence of traffic volume for the last 10 sample points, day of the week, hour of the day, road direction, number of lanes, and name of the road.
- Meteorological Data [47]: This dataset aims to predict the area affected by forest fires using meteorological data such as temperature, relative humidity, wind, and rain. The area affected by the fires is transformed using an $\ln(x+1)$ function before applying various Data Mining methods.

B. Baseline Models

1) RL Algorithms:

- Existing RL baseline models by Li et al. [33] are deployed to optimize sequential treatment strategies based on EHRs for chronic diseases with DQN. The multi-agent framework results are compared with Q Learning and Double DQN.
- Additionally, Guo et al. [34] used RL and a dynamic weight assignment network architecture with TD3 (combination of DDPG, Actor-Critic, and DQN) to recognize human activity.
- Yom et al. [32] used A2C and PPO algorithms to act as virtual coaches in decision-making and send personalized messages.

2) Predictive RL Frameworks:

- Li et al. [42] proposed simultaneous energy-based learning for multi-agent activity forecasting using graph neural networks for activity forecasting based on spatio-temporal data.

- Ma et al. [43] proposed a multi-agent driving behavior prediction across different scenarios based on the agent’s self-supervised domain knowledge.
- Jiang et al. [44] proposed a study on internet traffic prediction with distributed multi-agent learning using LSTM and gated recurrent unit (GRU). GRU-based distributed multi-agent learning scheme achieved the best performance compared to LSTM.

C. Evaluation Metrics

Mean Absolute Error (MAE) is a commonly used regression metric that measures the average magnitude of errors between the predicted and actual values for a set of data. It is calculated as the average of absolute differences between the predicted and actual values and is expressed as a single value. **Root Mean Squared Error (RMSE)** is another commonly used regression metric that measures the average magnitude of the differences between the predicted and actual values. RMSE is calculated as the square root of the mean of the squared differences between the predicted and actual values. **Mean Absolute Percentage Error (MAPE)** is a regression metric that measures the average absolute percentage error between the predicted and actual values. It is calculated as the average of the absolute differences between the predicted and actual values, expressed as a percentage of the actual values. **Cumulative Rewards** is a performance metric used in reinforcement learning to measure the total rewards obtained by an agent over a specified period of time or number of actions. It is calculated as the sum of all rewards received by the agent over the given period of time or number of actions. All the experiments were conducted using Python version 3.7.6 and the TensorFlow, Keras, Open AI Gym, and stable_baselines3 packages.

V. EXPERIMENT RESULTS AND ANALYSIS

In this section, we conduct an analysis and comparison of the performance of the deep learning model within the forecasting environment of the proposed framework against baseline models. Additionally, we evaluate and compare the performance of the monitoring RL agents with baseline models. The proposed framework utilizes deep learning for forecasting the states of the RL environment, while the RL agent monitors the forecasted states. This makes RL a suitable approach for monitoring applications, such as health monitoring, weather monitoring, traffic monitoring, and more. Moreover, it enables the automation of monitoring tasks, reducing the reliance on manual intervention and increasing the efficiency of monitoring processes. RL agents can continuously monitor the environment, make real-time decisions, and take appropriate actions, allowing human operators to focus on other critical tasks.

A. Forecasting Environment Results

Traditional machine learning and deep learning algorithms are capable of predicting heart rate in a supervised learning approach. The baseline models with predicting capability are

TABLE II: Baseline Models Comparison

RL Method	Agent 1	Agent 2	Agent 3
Q Learning	25878	17304	23688
PPO	23688	20367	17688
A2C	24717	13707	24369
Double DQN	25569	15360	20367
DDPG	26760	20754	23967
Proposed DQN (Ours)	48354	30019	38651

adopted in the PDRL framework to replace the proposed DQN algorithm. They are trained with a subject from PPG-DaLiA to forecast heart rate based on physiological features. Tab. I presents the results of various frameworks for time series forecasting in the RL environment. The frameworks being compared are ELMA, GRU, GNN-Based Multi-Agent, and the proposed generic Multi-Agent PDRL Framework. The performance of each framework is evaluated using the three metrics: MAE, MAPE, and RMSE. The results show that the proposed multi-agent PDRL framework performs the best among all the models across all time intervals (15 min, 30 min, 45 min, and 60 min). This can be seen by the lowest values of MAE, MAPE, and RMSE for this model. The GRU model also performs well across all time intervals, with MAE, MAPE, and RMSE values significantly lower than those of the ELMA and GNN-Based Multi-Agent models. It is also worth noting that the performance of the GNN-Based Multi-Agent model is inconsistent across all time intervals, showing varying results for different time steps. The results suggest that the GRU and the proposed Multi-Agent PDRL Framework models, which are specifically designed for time series data, performed much better than ELMA and GNN-Based Multi-Agent which are not well suited for time series forecasting.

B. DRL Agents Performance

The multi-agent in the proposed PDRL monitoring framework has been evaluated with vital signs such as heart rate, respiration, and temperature predicted in the forecasting environment. The agents are responsible to monitor the vital signs in each iteration and take appropriate actions. For this task, each agent gets rewarded as per the reward policy discussed in the previous section III. As they aim to increase their accumulated rewards, all the agents learn each vital sign pattern of patients and collectively monitor the patient’s health status. The isolated PDRL agents monitor their vital signs independently and update the corresponding MET at the right time.

All the baseline models and the proposed algorithm were trained with the same client data for 10 episodes and their cumulative rewards achieved in the 10th episode are shown in Tab. II. Q-Learning algorithm which updates its action-value Bellman equation stopped far behind in learning the patterns of the vital signs compared to the proposed model-free DQN algorithm. The other baseline models PPO, Actor-Critic, Double DQN, and DDPG have considerable performance in all three monitoring agents but couldn’t overshadow the results of

the proposed DQN approach. The multi-model algorithm TD3 was able to score closer to the proposed approach. Overall, the proposed PDRL algorithm has outperformed all other baseline models in all three monitoring agents.

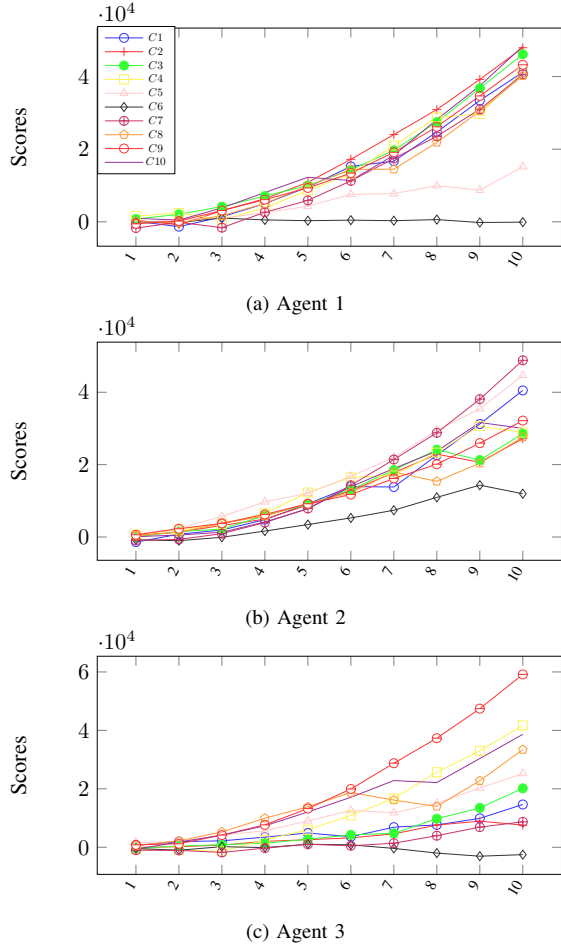


Fig. 3: DQN Agents Performance

All three learning agents have been fed with physiological features such as heart rate, respiration, and temperature respectively. Based on the observation space, action space, and reward policy defined for a customized gym environment for human behavior monitoring, the learning agents were run for 10 episodes shown on the x-axis, and the cumulative rewards have been awarded as scores for each episode shown on the y-axis. The performance of each of the learning agents with respect to each input client (Client 1 to Client 10) data can be seen in Fig. 3. In the results, Agent 1 refers to the heart rate monitoring agent which has a constant increase of scores for each episode for most of the subjects except subjects 5 and 6. The intermittent low scores in Agent 1 performance is due to the exploration rate in DQN learning where the algorithm tries exploring all the actions randomly instead of using neural networks prediction. Similarly, Agent 2 and Agent 3 monitor two other physiological features respiration and temperature

respectively. Agent 2 has performed better than the other two agents and achieved consistent scores for all subjects. Out of all agents, Agent 3, temperature monitoring performance is unsatisfactory. This actually drives us back to the data level and found out the data is with a different scale compared to the MEWS [22]. Still, Agent 3 achieved high scores in monitoring subjects 9, 8, 4, and 10.

VI. OTHER TIME SERIES FORECASTING SYSTEMS

In the healthcare forecasting and monitoring experiment, vital signs such as heart rate, respiration, and temperature are predicted based on the time series data in the forecasting RL environment and the DRL agents monitored the predicted vital signs to communicate with the appropriate medical emergency team in adverse situations. There are different domains such as traffic and weather where time series forecasting is critical and making sequential decisions are essential. In this study, two such time series forecasting systems are evaluated using the proposed PDRL framework. For these experiments, the monitoring agents and forecasting environment trained in the healthcare experiment are adopted for other time series forecasting applications such as traffic and weather by storing the knowledge gained from the health monitoring application. In the transfer learning process, the traffic dataset [46] and meteorological data [47] is used for the evaluation. In the traffic dataset, a DRL agent is deployed for monitoring the traffic forecasting process by customizing the observation space, action space, and rewards in the RL environment. Similarly, a DRL agent is deployed for monitoring the weather forecasting process.

Tab. III appears to show the results of the transfer learning experiment in which different models (ELMA, GRU, GNN-Based Multi-Agent, the proposed Multi-Agent PDRL Framework) are used to predict traffic and meteorological data. The results are shown in terms of several evaluation metrics: MAE, MAPE, and RMSE. The time intervals (15 min, 30 min, 45 min, 60 min) indicate the time intervals for which the predictions are made. GRU and the proposed Multi-Agent PDRL Framework models are performing the best for both traffic and meteorological data forecasting across all time intervals. The ELMA and GNN-Based Multi-Agent models, on the other hand, do not perform as well as the GRU and the proposed Multi-Agent PDRL Framework models.

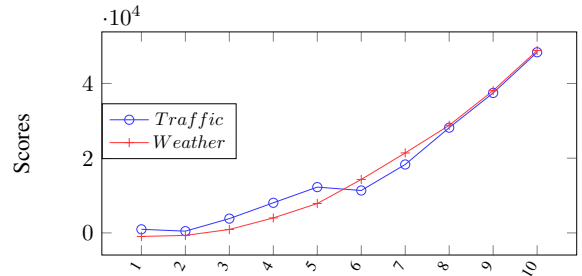


Fig. 4: PDRL Agent results on traffic and weather forecasting

The performance of the proposed PDRL monitoring agents in traffic and weather monitoring applications is presented in

TABLE III: Proposed Multi-Agent PDRL framework performance for traffic and weather prediction

		Traffic Data Forecasting			Meteorological Data Forecasting		
		MAE	MAPE	RMSE	MAE	MAPE	RMSE
ELMA [42]	15 Min	6.73	15.14	9.4	6.69	15.02	9.39
	30 Min	6.73	15.14	9.4	6.69	15.02	9.39
	45 Min	6.73	15.14	9.4	6.69	15.02	9.39
	60 Min	6.72	15.07	9.39	6.65	14.99	9.34
GRU [43]	15 Min	1.04	6.04	1.36	1.03	5.96	1.36
	30 Min	1.04	6.01	1.36	1.03	5.94	1.35
	45 Min	1.04	5.96	1.36	1.04	5.93	1.36
	60 Min	1.04	6.1	1.36	1.04	6.02	1.36
GNN-Based Multi-Agent [44]	15 Min	4.64	6.07	2.88	4.27	7.32	2.76
	30 Min	5.79	5.7	3.21	5.03	4.71	2.99
	45 Min	6.57	4.93	3.43	5.61	4.89	3.16
	60 Min	7.01	6.07	3.54	6.57	5.86	3.43
Multi-Agent PDRL Fraemwork (Ours)	15 Min	0.47	2.94	0.91	0.44	2.83	0.89
	30 Min	0.56	5.61	1.01	0.6	4.21	1.04
	45 Min	0.46	4.04	0.91	0.48	3.69	0.92
	60 Min	0.58	5.71	1.01	0.54	5.71	0.99

Fig. 4. It appears that the PDRL agent is able to perform well on both tasks, as the total rewards for both tasks increase with episode number. However, it is also apparent that the agent performs better on the traffic monitoring task than on the weather monitoring task, as the total rewards for the traffic monitoring task are consistently higher than those for the weather monitoring task. Additionally, the gap between the rewards for the two tasks is increasing as the episode number increases, indicating that the agent is becoming increasingly better at the traffic monitoring task.

VII. DISCUSSION

The primary objective of the proposed study is to design a multi-agent framework to monitor the predicted future states of a dynamic and complex environment. The proposed framework has been adopted to monitor patients in an uncertain hospital environment where the patient's vital signs fluctuate intermittently and might cause health deterioration with delay in treatment. To overcome this challenge, the sequential decision-making capability of RL algorithms was adopted in this study. Each vital sign in the human body has different threshold levels to determine the health emergency as per MEWS [22] and medical emergency teams are predefined for each emergency based on the threshold of the vital sign. In this study, a PDRL agent was deployed for each vital sign and three PDRL agents interacted with the same generic healthcare monitoring environment. The PDRL agents have no prior training or knowledge about patients' vital signs. Based on the reward policy defined in the forecasting environment, the DRL agents learning agents predicted the right action or right MET to communicate the emergency of each vital sign. While designing the environment, setting up the observation space for each DRL agent was critical as it would directly affect the agent learning process and might lead to ambiguity in communicating to the right MET. In this study, PDRL agent 3 was deployed to monitor patients' body temperature,

and its performance was unsatisfactory compared to the other two agents. This raises the question of the sanity of input data and observation space configured in the environment. The issue was the units of the temperature thresholds in the MEWS table and the input body temperature of data from the dataset are different. The proposed multi-agent PDRL framework is a generic framework that can be adapted to different time series forecasting applications and monitored to make sequential decisions. In this study, the experiments on healthcare, traffic, and weather data showed promising results. An added advantage of the proposed framework is multi-agents for multiple monitoring parameters in a dynamic environment. This avoids the sparse rewards challenge and can be easily customized and adapted to different applications.

VIII. CONCLUSION

This study proposes a new paradigm of monitoring forecasted states using multiple DRL agents. A generic PDRL monitoring environment was designed with a reward policy to reward the DRL agents based on their actions in each iteration of monitoring status. The learning agents were compelled to learn the behavior of the data patterns based on the reward policy for all possible actions in the action space for each state in the continuous observation space. Based on the evaluation results, all three DRL agents in the PDRL framework were able to learn the patterns of the vital signs and predict appropriate action to alert corresponding medical emergency teams. Furthermore, the knowledge from health monitoring is stored and performed transfer learning process on traffic and weather monitoring. However, the limitation of this study is the input data scale, or units of states, that the agent is monitoring. This led to the under-performance of DRL agent 3 compared to the other two DRL agents in the health monitoring application. Ensemble methods, such as combining the predictions from multiple DRL agents or combining the PDRL framework with other machine learning approaches, could be explored to further improve the accuracy and robustness of the forecasting and decision-making capabilities of the system.

REFERENCES

- [1] X. Shu and Y. Ye, "Knowledge discovery: Methods from data mining and machine learning," *Social Science Research*, p. 102817, Oct. 2022.
- [2] G. Talari, E. Cummins, C. McNamara, and J. O'Brien, "State of the art review of big data and web-based decision support systems (DSS) for food safety risk assessment with respect to climate change," *Trends in Food Science & Technology*, vol. 126, pp. 192–204, Aug. 2022.
- [3] T. Shaik, X. Tao, N. Higgins, H. Xie, R. Gururajan, and X. Zhou, "AI enabled RPM for mental health facility," in *Proceedings of the 1st ACM Workshop on Mobile and Wireless Sensing for Smart Healthcare*, ACM, Oct. 2022.
- [4] Y.-H. Hu, "Using few-shot learning materials of multiple SPOCs to develop early warning systems to detect students at risk," *The International Review of Research in Open and Distributed Learning*, vol. 23, pp. 1–20, Feb. 2022.
- [5] P. Radanliev, D. D. Roue, and R. Walton, "Data mining and analysis of scientific research data records on covid-19 mortality, immunity, and vaccine development - in the first wave of the covid-19 pandemic," *Diabetes & Metabolic Syndrome: Clinical Research & Reviews*, vol. 14, pp. 1121–1132, Sept. 2020.
- [6] R. K. Pattanayak, V. S. Kumar, K. Raman, M. M. Surya, and M. R. Pooja, "E-commerce application with analytics for pharmaceutical industry," in *Advances in Intelligent Systems and Computing*, pp. 291–298, Springer Nature Singapore, Sept. 2022.

- [7] R. Thirunavukarasu, G. P. D. C. G. R. M. Gopikrishnan, and V. Palanisamy, "Towards computational solutions for precision medicine based big data healthcare system using deep learning models: A review," *Computers in Biology and Medicine*, vol. 149, p. 106020, Oct. 2022.
- [8] D. Qiao, S. Guo, D. Liu, S. Long, P. Zhou, and Z. Li, "Adaptive federated deep reinforcement learning for proactive content caching in edge computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 12, pp. 4767–4782, 2022.
- [9] T. M. Ho, K.-K. Nguyen, and M. Cheriet, "Federated deep reinforcement learning for task scheduling in heterogeneous autonomous robotic system," *IEEE Transactions on Automation Science and Engineering*, pp. 1–13, 2022.
- [10] W. Zhang, N. Zhang, J. Yan, G. Li, and X. Yang, "Auto uning of price prediction models for high-frequency trading via reinforcement learning," *Pattern Recognition*, vol. 125, p. 108543, 2022.
- [11] W. Zhang, L. Wang, L. Xie, K. Feng, and X. Liu, "Tradebot: Bandit learning for hyper-parameters optimization of high frequency trading strategy," *Pattern Recognition*, vol. 124, p. 108490, 2022.
- [12] L. Han, B. Du, L. Sun, Y. Fu, Y. Lv, and H. Xiong, "Dynamic and multi-faceted spatio-temporal deep learning for traffic speed forecasting," in *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pp. 547–555, 2021.
- [13] Z. Fang, Q. Long, G. Song, and K. Xie, "Spatial-temporal graph ode networks for traffic flow forecasting," in *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, pp. 364–373, 2021.
- [14] P. Kumar, R. Chandra, C. Bansal, S. Kalyanaraman, T. Ganu, and M. Grant, "Micro-climate prediction-multi scale encoder-decoder based deep learning framework," in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 3128–3138, 2021.
- [15] C. Yu, J. Liu, S. Nemati, and G. Yin, "Reinforcement learning in healthcare: A survey," *ACM Comput. Surv.*, vol. 55, nov 2021.
- [16] H. Yoon and J. Li, "A novel positive transfer learning approach for tele-monitoring of parkinson's disease," *IEEE Transactions on Automation Science and Engineering*, vol. 16, no. 1, pp. 180–191, 2019.
- [17] N. Hong, C. Liu, J. Gao, L. Han, F. Chang, M. Gong, and L. Su, "State of the art of machine learning-enabled clinical decision support in intensive care units: Literature review," *JMIR Medical Informatics*, vol. 10, p. e28781, Mar. 2022.
- [18] J. Zeng, J. Shao, S. Lin, H. Zhang, X. Su, X. Lian, Y. Zhao, X. Ji, and Z. Zheng, "Optimizing the dynamic treatment regime of in-hospital warfarin anticoagulation in patients after surgical valve replacement using reinforcement learning," *Journal of the American Medical Informatics Association*, vol. 29, pp. 1722–1732, July 2022.
- [19] J. Watts, A. Khojandi, R. Vasudevan, and R. Ramdhani, "Optimizing individualized treatment planning for parkinson's disease using deep reinforcement learning," in *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, IEEE, July 2020.
- [20] M. Naeem, G. Paragliola, and A. Coronato, "A reinforcement learning and deep learning based intelligent system for the support of impaired patients in home treatment," *Expert Systems with Applications*, vol. 168, p. 114285, Apr. 2021.
- [21] I. Y. Chen, S. Joshi, M. Ghassemi, and R. Ranganath, "Probabilistic machine learning for healthcare," *Annual Review of Biomedical Data Science*, vol. 4, pp. 393–415, July 2021.
- [22] V. Signs, "Canberra hospital and health services clinical procedure," 2021.
- [23] S. Chen, X. Qiu, X. Tan, Z. Fang, and Y. Jin, "A model-based hybrid soft actor-critic deep reinforcement learning algorithm for optimal ventilator settings," *Information Sciences*, vol. 611, pp. 47–64, Sept. 2022.
- [24] D. Silver, S. Singh, D. Precup, and R. S. Sutton, "Reward is enough," *Artificial Intelligence*, vol. 299, p. 103535, Oct. 2021.
- [25] G. Akçapınar, A. Altun, and P. Aşkar, "Using learning analytics to develop early-warning system for at-risk students," *International Journal of Educational Technology in Higher Education*, vol. 16, Oct. 2019.
- [26] A. Cano and J. D. Leonard, "Interpretable multiview early warning system adapted to underrepresented student populations," *IEEE Transactions on Learning Technologies*, vol. 12, pp. 198–211, Apr. 2019.
- [27] L. Hussain-Alkhateeb, T. R. Ramirez, A. Kroeger, E. Gozzer, and S. Runge-Ranzinger, "Early warning systems (EWSs) for chikungunya, dengue, malaria, yellow fever, and zika outbreaks: What is the evidence? a scoping review," *PLOS Neglected Tropical Diseases*, vol. 15, p. e0009686, Sept. 2021.
- [28] J. Liu, G. Shi, J. Zhou, and Q. Yao, "Prediction of college students' psychological crisis based on data mining," *Mobile Information Systems*, vol. 2021, pp. 1–7, May 2021.
- [29] E. Moghadas, J. Rezazadeh, and R. Farahbakhsh, "An IoT patient monitoring based on fog computing and data mining: Cardiac arrhythmia usecase," *Internet of Things*, vol. 11, p. 100251, Sept. 2020.
- [30] S.-H. Moon, Y.-H. Kim, Y. H. Lee, and B.-R. Moon, "Application of machine learning to an early warning system for very short-term heavy rainfall," *Journal of Hydrology*, vol. 568, pp. 1042–1054, Jan. 2019.
- [31] A. Lisowska, S. Wilk, and M. Peleg, "From personalized timely notification to healthy habit formation: a feasibility study of reinforcement learning approaches on synthetic data," in *SMARTERCARE@ AI* IA*, pp. 7–18, 2021.
- [32] E. Yom-Tov, G. Feraru, M. Kozdoba, S. Mannor, M. Tennenholtz, and I. Hochberg, "Encouraging physical activity in patients with diabetes: Intervention using a reinforcement learning system," *Journal of Medical Internet Research*, vol. 19, p. e338, Oct. 2017.
- [33] T. Li, Z. Wang, W. Lu, Q. Zhang, and D. Li, "Electronic health records based reinforcement learning for treatment optimizing," *Information Systems*, vol. 104, p. 101878, Feb. 2022.
- [34] J. Guo, Q. Liu, and E. Chen, "A deep reinforcement learning method for multimodal data fusion in action recognition," *IEEE Signal Processing Letters*, vol. 29, pp. 120–124, 2022.
- [35] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *Proceedings of the 35th International Conference on Machine Learning (J. Dy and A. Krause, eds.)*, vol. 80 of *Proceedings of Machine Learning Research*, pp. 1587–1596, PMLR, 10–15 Jul 2018.
- [36] D. Tirumala, A. Galashov, H. Noh, L. Hasenclever, R. Pascanu, J. Schwarz, G. Desjardins, W. M. Czarnecki, A. Ahuja, Y. W. Teh, and N. Heess, "Behavior priors for efficient reinforcement learning," 2020.
- [37] M. Janssen, C. LeWarne, D. Burk, and B. B. Averbeck, "Hierarchical reinforcement learning, sequential behavior, and the dorsal frontostriatal system," *Journal of Cognitive Neuroscience*, vol. 34, no. 8, pp. 1307–1325, 2022.
- [38] K. Tsiakias, M. Papakostas, M. Theofanidis, M. Bell, R. Mihalcea, S. Wang, M. Burzo, and F. Makedon, "An interactive multisensing framework for personalized human robot collaboration and assistive training using reinforcement learning," in *Proceedings of the 10th International Conference on Pervasive Technologies Related to Assistive Environments*, ACM, June 2017.
- [39] A. Kubota and L. D. Riek, "Methods for robot behavior adaptation for cognitive neurorehabilitation," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 5, pp. 109–135, May 2022.
- [40] H. Almohri, L. Cheng, D. Yao, and H. Alemzadeh, "On threat modeling and mitigation of medical cyber-physical systems," in *2017 IEEE/ACM International Conference on Connected Health: Applications, Systems and Engineering Technologies (CHASE)*, IEEE, July 2017.
- [41] M. M. Hassan, S. Ullah, M. S. Hossain, and A. Alelwi, "An end-to-end deep learning model for human activity recognition from highly sparse body sensor data in internet of medical things environment," *The Journal of Supercomputing*, vol. 77, pp. 2237–2250, June 2020.
- [42] Y. Li, P. Wang, L. Chen, Z. Wang, and C.-Y. Chan, "Elma: Energy-based learning for multi-agent activity forecasting," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 36, pp. 1482–1490, 2022.
- [43] H. Ma, Y. Sun, J. Li, and M. Tomizuka, "Multi-agent driving behavior prediction across different scenarios with self-supervised domain knowledge," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pp. 3122–3129, IEEE, 2021.
- [44] W. Jiang, M. He, and W. Gu, "Internet traffic prediction with distributed multi-agent learning," *Applied System Innovation*, vol. 5, no. 6, p. 121, 2022.
- [45] A. Reiss, I. Indlekofer, P. Schmidt, and K. V. Laerhoven, "Deep PPG: Large-scale heart rate estimation with convolutional neural networks," *Sensors*, vol. 19, p. 3079, July 2019.
- [46] L. Zhao, O. Gkoutouna, and D. Pfoser, "Spatial auto-regressive dependency interpretable learning based on spatial topological constraints," *ACM Transactions on Spatial Algorithms and Systems (TSAS)*, vol. 5, no. 3, pp. 1–28, 2019.
- [47] P. Cortez and A. d. J. R. Morais, "A data mining approach to predict forest fires using meteorological data," 2007.

6.2 Summary

The PDRL framework's capability to significantly enhance predictive monitoring through multi-agent deep reinforcement learning is underscored. The framework's application to vital sign monitoring demonstrates not only its ability to learn and adapt to complex patterns but also its superiority over conventional models. The chapter highlights the framework's scalability and adaptability to various domains, including traffic and weather forecasting, through transfer learning. The comprehensive evaluation and successful deployment of this framework suggest a promising avenue for future research and development in intelligent monitoring systems, potentially transforming the landscape of real-time data analysis and decision-making in healthcare and other critical fields.

CHAPTER 7: PAPER 6 - A SURVEY OF MULTIMODAL INFORMATION FUSION FOR SMART HEALTHCARE: MAPPING THE JOURNEY FROM DATA TO WISDOM

7.1 Introduction

This chapter embarks on an in-depth exploration of multimodal information fusion's pivotal role in the evolution of smart healthcare systems. It meticulously examines how the convergence of diverse data modalities—spanning electronic health records, wearable technologies, genomics, and environmental factors—facilitates a holistic view of patient health within the Data-Information-Knowledge-Wisdom (DIKW) hierarchy. This detailed scrutiny not only sheds light on the theoretical foundations and methodologies of multimodal fusion but also highlights its instrumental role in refining clinical decision-making processes, tailoring patient care, and enhancing the efficacy of healthcare interventions.



Contents lists available at ScienceDirect

Information Fusion

journal homepage: www.elsevier.com/locate/inffus

A survey of multimodal information fusion for smart healthcare: Mapping the journey from data to wisdom

Thanveer Shaik ^{a,*}, Xiaohui Tao ^a, Lin Li ^b, Haoran Xie ^c, Juan D. Velásquez ^{d,e}

^a School of Mathematics, Physics and Computing, University of Southern Queensland, Australia

^b School of Computer Science and Artificial Intelligence, Wuhan University of Technology, China

^c Department of Computing and Decision Sciences, Lingnan University, Hong Kong

^d Industrial Engineering Department, University of Chile, Chile

^e Instituto Sistemas Complejos de Ingeniería, Santiago, Chile

ARTICLE INFO

Keywords:

DIKW
Multimodality
Data fusion
p4 medicine
Smart healthcare

ABSTRACT

Multimodal medical data fusion has emerged as a transformative approach in smart healthcare, enabling a comprehensive understanding of patient health and personalized treatment plans. In this paper, a journey from data to information to knowledge to wisdom (DIKW) is explored through multimodal fusion for smart healthcare. We present a comprehensive review of multimodal medical data fusion focused on the integration of various data modalities. The review explores different approaches such as feature selection, rule-based systems, machine learning, deep learning, and natural language processing, for fusing and analyzing multimodal data. This paper also highlights the challenges associated with multimodal fusion in healthcare. By synthesizing the reviewed frameworks and theories, it proposes a generic framework for multimodal medical data fusion that aligns with the DIKW model. Moreover, it discusses future directions related to the four pillars of healthcare: Predictive, Preventive, Personalized, and Participatory approaches. The components of the comprehensive survey presented in this paper form the foundation for more successful implementation of multimodal fusion in smart healthcare. Our findings can guide researchers and practitioners in leveraging the power of multimodal fusion with the state-of-the-art approaches to revolutionize healthcare and improve patient outcomes.

1. Introduction

In the realm of smart healthcare, where cutting-edge technologies and data-driven approaches are revolutionizing the field, the integration of multimodal data has emerged as a transformative tool to enhance decision-making and improve patient outcomes. This paper presents a comprehensive exploration of multimodal medical data fusion for smart healthcare, illustrating the journey from raw data to actionable insights through the four-level pyramid shown in Fig. 1.

The Data Information Knowledge Wisdom (DIKW) model is a conceptual framework that illustrates the hierarchical progression of data into wisdom [1]. Through its process, raw data is transformed into meaningful information, knowledge, and ultimately wisdom, which can be used for informed decision-making and problem-solving [2]. The DIKW model recognizes that data alone is not sufficient to drive insights and actions. Instead, data needs to be processed, organized, and contextualized to extract valuable information. This information is then synthesized and combined with existing knowledge to gain understanding, leading to the development of knowledge. Knowledge,

in turn, can then be applied in practical situations to make informed decisions and solve complex problems, resulting in wisdom.

At the base of the pyramid in Fig. 1, we have the Data level, which encompasses diverse sources of data such as Electronic Health Records (EHRs), medical imaging, wearable devices, genomic data, sensor data, environmental data, and behavioral data. This raw data serves as the foundation for subsequent analysis and interpretation. Moving up the pyramid, we reach the Information level, where the raw data undergoes processing, organization, and structuring to derive meaningful and contextualized information. For instance, heart rate data from a wearable device can be processed to determine average resting heart rates, activity levels, and potential anomalies.

The Knowledge level, situated above Information, represents the interconnected structure of the organized data from various sources. By establishing relationships and connections between entities like patients, diseases, or medical treatments, the Knowledge level enables the identification of patterns, trends, and correlations. It facilitates a holistic understanding of the data and serves as a powerful tool for generating insights.

* Corresponding author.

E-mail address: Thanveer.Shaik@usq.edu.au (T. Shaik).

<https://doi.org/10.1016/j.inffus.2023.102040>

Received 19 June 2023; Received in revised form 3 September 2023; Accepted 23 September 2023

Available online 28 September 2023

1566-2535/© 2023 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

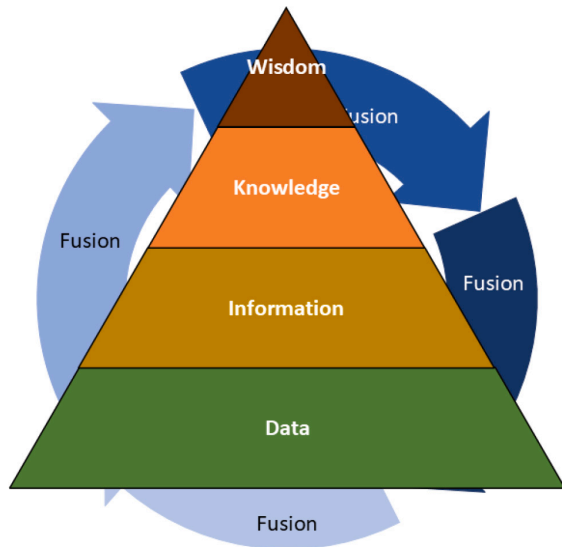


Fig. 1. DIKW Fusion conceptual model.

Finally, at the pinnacle of the pyramid, we have the Wisdom level. This is where actionable insights are derived from the Knowledge level, allowing informed decision-making, prediction of future outcomes [3], and a deeper understanding of complex phenomena. These insights enable personalized treatment plans, predictions about disease progression, and the identification of risk factors.

It is important to point out a key feature of the model, its circular structure, shown by the arrows in Fig. 1. These arrows indicate that combining different types of data assists in the progression of Data to Information, then to Knowledge and Wisdom. This cyclical nature adds flexibility to the model, allowing for constant updates and improvements in how data is processed. In this sense, reaching the level of Wisdom helps to fine-tune the steps and methods used at earlier stages, making future data collection, information gathering, and knowledge creation more effective.

This paper further explores different approaches such as feature selection, rule-based systems, machine learning, deep learning, and natural language processing for multimodal fusion. It also addresses challenges related to data quality, privacy, security, processing, analysis, clinical integration, ethics, and interpretation of results. With its emphasis on the transformative potential of multimodal medical data fusion, this paper sets the stage for future research and advancements in the field of smart healthcare, and paves the way for improved patient care outcomes and personalized healthcare solutions.

The following are the key contributions of this paper:

- The application and adaptation of the existing DIKW conceptual model to describe the journey of data to information to knowledge to wisdom in the context of multimodality fusion for smart healthcare.
- A taxonomy that organizes state-of-the-art techniques in multimodality fusion with the DIKW conceptual model.
- A proposed generic DIKW techniques framework for smart healthcare, that not only highlights the current efforts, but also provides a vision for its future evolution.
- A review of the challenges and recommended solutions associated with multimodal fusion, informed by the existing DIKW conceptual model and the proposed framework, to guide future research directions.

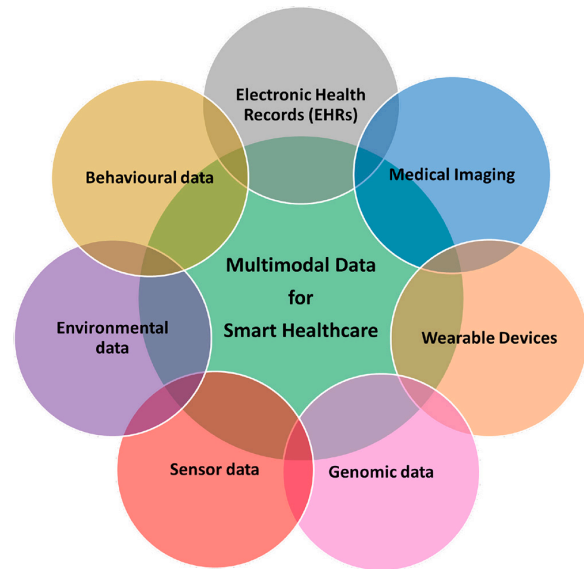


Fig. 2. Overview of multimodality fusion for smart healthcare.

The rest of this paper is organized as follows: Section 2 delves into the representation of data in multimodal applications. Section 3 explores various approaches for integrating information from multiple modalities, then outlines the proposed taxonomy. Section 4 examines the challenges and trends associated with multimodal fusion. Section 5 presents a generic framework for multimodal fusion that aligns with the DIKW model. In Section 6, we outline the future directions for multimodal fusion in smart healthcare, with a particular focus on the 4Ps of healthcare. Finally, the paper concludes in Section 7.

2. Modalities in smart healthcare

There are various data modalities in healthcare, such as EHRs, Medical Imaging, Wearable Devices, Genomic data, Sensor data, Environmental data, and Behavioral data as shown in Fig. 2. These modalities contain unstructured raw data specific to their respective formats. As the data is processed, it is transformed into meaningful information through the involvement of techniques such as structuring EHRs, feature extraction from Medical Imaging, and analysis of wearable device data.

2.1. Electronic Health Records (EHRs)

EHRs serve as a central repository of medical data for healthcare providers, the adoption of which has led to a surge in the amount of intricate patient data [4]. These datasets, although extensive and tailored to each patient, are often fragmented and may lack organization. They encompass diverse variables like medications, laboratory values, imaging results, physiological measurements, and historical notes [5], which lead to increased complexity in analysis. Machine learning (ML) provides a potential solution to this complexity by enabling the exploration of intricate relationships among the diverse variables present in EHR datasets [6].

EHRs play a crucial role in multimodal fusion systems within smart healthcare, especially for multidisciplinary and life-threatening diseases like diabetes [7]. However, managing and analyzing unstructured data collected from sensors and EHRs is challenging. Data fusion is crucial for accurate predictions, and deep learning approaches are effective for larger healthcare datasets. Healthcare datasets can be enhanced

by collecting patients' data through wearable sensors and EHRs. An ensemble ML approach is employed to develop a recommendation system for accurate prediction and timely recommendations for patients with multidisciplinary diabetes.

To optimize multimodal fusion strategies in EHR data, Xu et al. [8] proposed MUFASA, a novel approach that extends Neural Architecture Search (NAS). The authors based their model on the Transformer architecture, which has shown promise in leveraging EHR's internal structure. Experimental results demonstrated that MUFASA outperformed Transformer, Evolved Transformer, RNN variants, and traditional NAS models on public EHR data. MUFASA architectures achieved higher top-5 recall rates compared to Transformer in predicting CCS diagnosis codes, and they outperformed unimodal NAS by customizing modeling for each modality. MUFASA also exhibited effective transfer learning to ICD-9, another EHR task. The representation of EHR data is challenging due to different modalities, such as medical codes and clinical notes, all of which have distinct characteristics.

Another challenge is the extraction of inter-modal correlations, which are often overlooked or not effectively captured by existing models. An et al. [9] proposed the Multimodal Attention-based fusion Networks (MAIN) model, which aims to address two key challenges in healthcare prediction using EHR data. The MAIN model incorporates multiple independent feature extraction modules tailored to each modality, including self-attention and time-aware Transformer for medical codes, and a CNN model for clinical notes. It also introduces an inter-modal correlation extraction module composed of a low-rank multimodal fusion method and a cross-modal attention layer. The model combines the representations of each modality and their correlations to generate visit and patient representations for diagnosis prediction [10], by leveraging attention mechanisms and neural networks. Overall, MAIN offers a comprehensive framework for multimodal fusion and correlation extraction in EHR-based prediction tasks.

2.2. Wearable devices

Wearable devices have become increasingly prevalent in the field of smart healthcare, offering the potential to monitor and track various aspects of an individual's health and well-being. These devices, typically worn on the body or incorporated into clothing or accessories, can collect real-time data about vital signs, physical activity, sleep patterns, and other health-related metrics [11,12]. The data gathered by wearable devices can provide valuable insights into an individual's overall health status, enabling personalized health monitoring and preventive care [13]. Moreover, wearable devices can facilitate remote patient monitoring, allowing healthcare professionals to track patients' health remotely and intervene when necessary. The integration of wearable devices with smart healthcare systems enables continuous monitoring, early detection of health issues, and the ability to deliver personalized interventions and recommendations [14].

2.3. Sensor data

Sensor data plays a crucial role in enabling smart healthcare by providing real-time monitoring and tracking of various physiological parameters and activities. Wearable devices, implantable sensors, and remote monitoring systems collect data such as heart rate, blood pressure, temperature, glucose levels, physical activity, sleep patterns, and so on. Sensor data in smart healthcare enables continuous monitoring of an individual's health status, facilitating early detection and intervention for potential health issues [15]. It allows healthcare providers to gather objective and accurate data, leading to more informed decision-making and personalized treatment plans. For example, sensor data can help in managing chronic conditions like diabetes or cardiovascular diseases by monitoring glucose levels or heart rate variability. Real-time sensor data can also enable remote patient monitoring, telemedicine, and telehealth services, allowing healthcare professionals to monitor

patients' conditions from a distance [16]. This is particularly beneficial for individuals with limited mobility or those residing in remote areas, by providing access to healthcare services without the need for frequent hospital visits [17].

The life cycle for data from EHRs, wearable devices, and sensors in the context of smart healthcare each follow a similar trajectory, encompassing stages such as raw data acquisition, data structuring, data fusion, and ultimately, predictive modeling. This cyclical process is graphically illustrated in Fig. 3.

2.4. Medical imaging

Medical imaging plays a crucial role in smart healthcare by providing valuable diagnostic information and aiding in the management of various medical conditions [18–23]. It involves the use of advanced imaging technologies to capture detailed images of the human body, allowing healthcare professionals to visualize and analyze anatomical structures, detect abnormalities, and monitor the progress of treatments. In smart healthcare, medical imaging is integrated with digital technologies and data analytics to enhance the efficiency, accuracy, and accessibility of healthcare services [18].

2.5. Genomic data

Genomic data plays a significant role in the realm of smart healthcare, offering valuable insights into an individual's genetic makeup and its impact on their health. This type of data includes information about an individual's DNA sequence, genetic variations, and gene expression patterns [24]. With advancements in genomic sequencing technologies, it has become more accessible and affordable to obtain a person's genetic information. In smart healthcare, genomic data can be utilized for various purposes, such as in the diagnosis and prediction of genetic disorders, as well as in the identification of genetic markers associated with increased disease risk or treatment response [25]. Genomic data can also enable personalized medicine by guiding treatment decisions based on an individual's unique genetic profile [26]. For example, it can help determine optimal drug choices and dosages, minimizing adverse reactions and improving treatment outcomes. Integrating genomic data with other health data sources, such as EHRs and wearable devices, can provide a comprehensive overview of an individual's health [27]. This multimodal approach allows for practices such as more accurate assessment of disease risks, personalized prevention strategies, and targeted interventions.

2.6. Environmental data

Environmental data can significantly contribute to smart healthcare by providing insights into the impact of environmental factors on individual health and well-being. Environmental data includes information about air quality, temperature, humidity, pollution levels, noise levels, and other relevant parameters in a person's surroundings. By incorporating environmental data into smart healthcare systems, healthcare providers can better understand the environmental conditions that may influence a person's health outcomes [28]. For example, monitoring air quality can help identify areas with high pollution levels, which is particularly valuable for individuals with respiratory conditions like asthma.

Analyzing this environmental data in real-time allows healthcare professionals to provide personalized recommendations and interventions to mitigate the impact of poor air quality on patients' health [29]. Environmental data can also aid in preventive healthcare by identifying patterns and correlations between environmental factors and specific health conditions [30,31]. For instance, studying the relationship between temperature and heat-related illnesses facilitates the implementation of early warning systems and interventions during heat waves or extreme weather events.

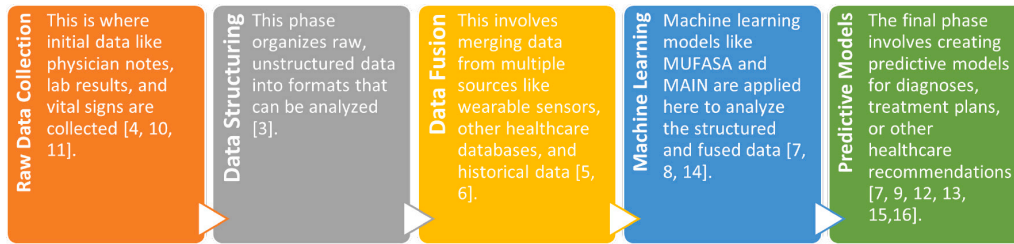


Fig. 3. The lifecycle of electronic health records, wearable devices, and sensors in smart healthcare.

Table 1
Multimodal Datasets for smart healthcare.

Modality	Datatype	Dataset	No. of Instances	No. of Attributes	Task	Popularity ^a
Single Modality	EHR	eICU Collaborative Research Database [32]	200,000 admissions	Varies	Various tasks, mainly diagnosis and prognosis	Medium
		MIMIC-III [33]	40,000 patients	Varies	Various tasks, mainly diagnosis and prognosis	High
	Imaging	MRNet [34]	1370 exams	MRI data	Disease detection	Low
		RSNA Pneumonia Detection Challenge [35]	30,000 images	Pneumonia labels	Disease detection	Low
		MURA [36]	40,895 images	Abnormal/normal	Disease detection	Medium
		Pediatric Bone Age Challenge Dataset [37]	Thousands of images	Bone age	Bone age estimation	Medium
		Indiana University Chest X-ray Collection [38]	8000 images	Chest radiograph DICOM images	Various tasks	Medium
		FastMRI [39]	Thousands of scans	MRI data	Image reconstruction	Medium
		CheXpert [40]	224,316 images	14 labels per image	Disease detection	High
		OASIS Brains Project [41]	Varies with dataset	MRI and clinical data	Brain studies	High
Multimodality	Genomics, Imaging	TCGA [48]	Thousands of patients	Genomic and clinical data	Cancer research	High
		UK Biobank [49]	500,000 individuals	Various data types	Various tasks	Medium
	Imaging, Genomics, EHR	ADNI [50]	Thousands of patients	MRI and clinical data	Alzheimer's research	High
		Imaging, Text	ImageCLEFmed [51]	Varies annually	Various data types	Various tasks
	Openi [52]		4.5 million images	Various data types	Various tasks	Low
	Various modalities	PhysioNet [53]	Various datasets	Various data types	Various tasks	High

^a Popularity is determined by the citation count in Google Scholar as of 05/06/2023. It is categorized as Low (≤ 200 citations), Medium (> 200 and < 1000 citations), and High (≥ 1000 citations).

2.7. Behavioral data

Behavioral data plays a crucial role in smart healthcare by providing valuable insights into individuals' habits, lifestyles, and behaviors, which have a significant impact on their overall health and well-being. Behavioral data encompasses various aspects of human behavior, including physical activity, sleep patterns, dietary habits, stress levels, social interactions, and adherence to medical treatments [54]. By leveraging behavioral data, smart healthcare systems can monitor and analyze individuals' behaviors in real time, allowing healthcare providers to comprehensively understand their patients' daily routines

and habits [55]. This data can help identify patterns, trends, and deviations from normal behavior, enabling early detection of potential health issues and the implementation of timely interventions.

Stress levels and emotional well-being can also be monitored through behavioral data, enabling healthcare providers to identify triggers and patterns that may impact individuals' mental health [56]. This data has been used to guide the development of stress management techniques, relaxation strategies, and personalized interventions to support individuals in maintaining good mental health [57]. Furthermore, behavioral data can facilitate patient engagement and self-management.

Through interactive platforms and feedback mechanisms, individuals can actively participate in monitoring their own behaviors, goal-setting, and receive personalized recommendations based on their provided data. This empowerment can lead to increased motivation and accountability in someone managing their health and well-being. However, the collection and analysis of behavioral data raises important ethical considerations, including privacy, data security, and informed consent. It is crucial to ensure that individuals' privacy is protected, their data is securely stored, and proper consent is obtained for data collection and usage.

2.8. Multimodality data

Multimodality data fusion in smart healthcare involves integrating information from various sources, such as EHRs, medical imaging, wearable devices, genomic data, sensor data, environmental data, and behavioral data. By combining data from different modalities, healthcare professionals can gain a comprehensive understanding of a patient's health, leading to personalized care and informed decision-making. Each modality provides unique insights, and their fusion enhances the accuracy and completeness of the analysis. For example, in a fusion approach, EHRs provide historical medical records, medical imaging offers anatomical details, wearables capture real-time physiological data, genomics reveal genetic predispositions, sensors provide contextual information, and behavioral data reflect lifestyle choices.

By integrating these modalities, healthcare professionals can uncover hidden patterns, correlations, and relationships that contribute towards techniques for optimizing treatment strategies, predicting disease progression, identifying risk factors, and implementing preventive measures. Multimodality data fusion is a crucial step towards a holistic approach for advancing smart healthcare and improving patient outcomes.

2.9. Datasets for multimodal fusion for smart healthcare

An overview of multimodal datasets used in smart healthcare is presented in Table 1. It includes information on the modality, dataset name, number of instances, number of attributes, task, popularity, and reference count. The datasets cover various modalities such as EHRs, Genomics, Imaging, and Text. Examples of datasets include the eICU Collaborative Research Database, TCGA, UK Biobank, MRNet, RSNA Pneumonia Detection Challenge, MURA, and ChestX-ray8. These datasets are used for diagnosis, prognosis, cancer research, disease detection, and image segmentation [58]. The popularity of the datasets depends on the reference counts, and the table provides reference counts for further exploration.

3. SOTA techniques in multimodal fusion for smart healthcare

Multimodal medical data fusion involves combining information from multiple modalities such as medical imaging, genomic data, EHRs, wearable devices, and more. In this section, we explore state-of-the-art (SOTA) techniques for multimodal fusion across these multiple modalities within the context of smart healthcare.

3.1. Feature selection

Feature selection focuses on identifying and selecting relevant features from raw data to transform them into meaningful information. In the context of multimodal medical data fusion for smart healthcare, recent studies have highlighted the importance of feature selection. Albahri et al. [59] emphasized its role in effective decision-making and improved patient care by identifying the most relevant features, reducing dimensionality, and enhancing the accuracy and interpretability of the fusion process. Similarly, Alghowinem et al. [60] emphasized its impact on improving result interpretability. By systematically applying

feature selection techniques, healthcare professionals can extract the most relevant information from the diverse data sources available, to facilitate a more comprehensive understanding of a patient's health conditions and support informed decision-making in personalized healthcare settings. Feature selection ensures efficient and effective data analysis and integration, ultimately enhancing the quality of data-driven insights and improving patient outcomes in smart healthcare environments.

Before fusion, it is important to perform feature selection within each modality separately. This can be achieved using various techniques, such as statistical tests, information gain, correlation analysis, or ML algorithms [61]. By selecting relevant features within each modality, noise and irrelevant information can be reduced, leading to improved fusion outcomes [62]. Certain modalities may contain more inherent noise or provide less relevant information compared to others. In such cases, modality-specific feature selection methods can be employed to identify the most informative features within each modality [63]. This can be done by leveraging domain knowledge, statistical analysis, or ML techniques tailored to the specific modality. After performing feature selection within each modality, the next step is to select features that are relevant across different modalities. Cross-modal feature selection methods aim to identify features that carry complementary information from multiple modalities [64]. These methods can involve techniques such as correlation analysis, mutual information, or joint optimization algorithms [65].

Feature selection and fusion should be performed in a coordinated manner to optimize the overall process. The selected features can be used as input for the fusion algorithm, which combines the information from different modalities [66]. This integration can be achieved through techniques such as early fusion, late fusion, or hybrid fusion approaches, depending on the nature of the data and the problem at hand [67].

It is important to evaluate the performance of the feature selection and fusion methods using appropriate evaluation metrics. Evaluation can involve assessing classification accuracy, regression performance, clustering quality, or other domain-specific evaluation criteria [68]. Cross-validation or independent validation on separate datasets can help validate the effectiveness of the feature selection techniques [69]. In healthcare applications, interpretability and explainability of the selected features and fusion results are crucial for building trust and understanding the decision-making process [70]. Various methods can be employed to enhance interpretability, such as feature importance ranking, visualization techniques, or rule extraction algorithms [71].

It is worth noting that the choice of feature selection methods may vary depending on the specific data characteristics, the fusion task, and the available computational resources. Additionally, the field of multimodal medical data fusion is an active area of research, and new techniques and algorithms are continuously being developed to address its challenges. The feature selection techniques of multimodal fusion are summarized in Fig. 4.

3.2. Rule-based systems

Rule-based systems operate to process and interpret information using predefined rules or logical statements. By employing these rules, these systems can make inferences and derive knowledge from the available data. The defined rules capture relationships and patterns, enabling decision-making based on the processed information. In the context of multimodal medical data fusion for smart healthcare, rule-based systems play a crucial role. They provide a structured approach to decision-making and knowledge representation [72], offering a systematic framework for integrating information from various modalities. By employing a set of predefined rules, these systems can effectively process and integrate data from multiple sources, enabling informed decisions and providing valuable recommendations [73]. The utilization of rule-based systems in multimodal medical data fusion enhances

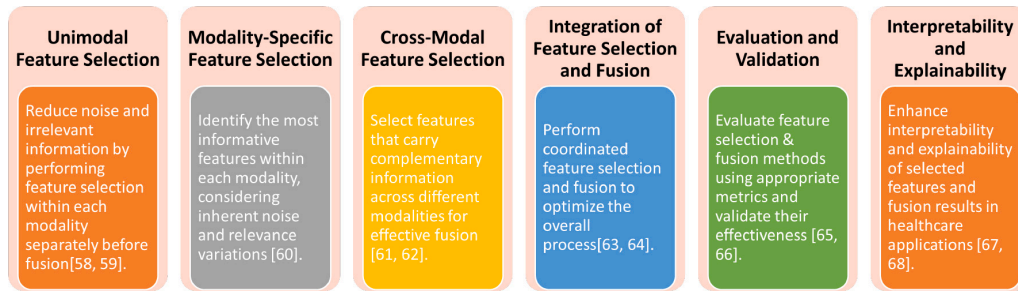


Fig. 4. Multimodal fusion — feature selection.

the overall knowledge generation process, aiding in accurate diagnoses, personalized treatment plans, and improved healthcare outcomes.

Rules in a rule-based system are typically defined using an “if-then” format, meaning that each rule consists of an antecedent (if) and a consequent (then). The antecedent represents the conditions or criteria that need to be satisfied, while the consequent specifies the action or conclusion to be taken if the conditions are met [74]. Rule-based systems can also contribute to feature selection in multimodal fusion. For example, rules can be designed to identify and select relevant features from different modalities based on their contribution to the decision-making process [70]. These rules can incorporate domain knowledge or statistical analysis to determine the importance of the features.

In multimodal medical data fusion, rules can be designed to accommodate multiple modalities. The antecedent of a rule can include conditions from different modalities, allowing the system to consider information from various sources simultaneously [75]. This integration can leverage the complementary nature of different modalities to enhance decision-making [76].

Medical data often contains uncertainty and imprecision. Rule-based systems can leverage fuzzy logic, a mathematical framework that handles uncertainty, to model and reason with uncertain or imprecise data [77]. Fuzzy rules allow for more flexible decision-making by assigning degrees of membership to antecedents and consequents, capturing the inherent uncertainty in medical data fusion [78]. In scenarios where multiple rules are applicable, conflicts may arise. Here, rule-based systems can employ strategies for rule prioritization and conflict resolution. These strategies determine the order in which rules are applied and resolve conflicts when multiple rules have conflicting conclusions [79]. This ensures a systematic and consistent decision-making process.

Rule-based systems also offer transparency and interpretability by providing explicit rules that can be examined and understood by healthcare professionals [80]. The rules provide explanations for the system’s decisions, allowing users to understand the underlying reasoning process. This transparency is crucial in building trust and facilitating collaboration between clinicians and the decision-support system [81]. Furthermore, rule-based systems enable the incorporation of expert knowledge into the decision-making process. Domain experts contribute their expertise by defining the rules that encapsulate their knowledge and clinical guidelines [73]. This allows the system to leverage the collective intelligence of healthcare professionals and enhance the accuracy and reliability of decision-making [82].

By monitoring a system’s performance and collecting feedback from its users, rule-based systems can be adapted or refined over time to improve decision-making [83]. This adaptive capability enables the system to evolve with new insights, changes in medical guidelines, or updates in the underlying data.

Although rule-based systems provide a structured and interpretable framework for multimodal medical data fusion, it is important to consider the limitations of these approaches, such as the challenge of capturing complex relationships or interactions between modalities,

and the potential for a large number of rules to manage. Hybrid approaches that combine rule-based systems with ML techniques can offer more flexibility and scalability in handling multimodal fusion tasks. The rule-based systems discussed in this subsection are outlined in Fig. 5.

3.3. Machine learning

Machine Learning (ML) encompasses the creation of algorithms capable of learning from data, enabling them to make predictions and informed decisions. By analyzing and processing vast amounts of data, ML algorithms can identify patterns, extract knowledge, and generate valuable insights to support decision-making processes. In the context of smart healthcare, ML techniques play a critical role in multimodal medical data fusion. These methods harness the capabilities of algorithms and statistical models to autonomously detect and understand patterns, relationships, and representations within diverse medical data sources. Through this automated learning process, ML contributes to overall knowledge generation within the healthcare sector, facilitating accurate diagnoses, personalized treatment plans, and improved patient outcomes.

Ensemble methods, such as Random Forests, gradient boosting, or AdaBoost, can be employed to combine the predictions or decisions of multiple ML models that have been trained on different modalities. Each modality can be processed independently using suitable algorithms, and their outputs can be fused using ensemble techniques to make a final decision [84]. Ensemble learning helps leverage the diversity and complementary information present in different modalities [85]. Its adaptive weights combination approach works by assigning weights to different modalities based on their relevance or importance for the fusion task. The weights can be learned using various techniques, such as feature selection algorithms, statistical analysis, or ML models [86]. The modalities are then combined using weighted fusion strategies, such as weighted averaging or weighted voting, to generate a fused representation [87].

Bayesian networks provide a probabilistic graphical model framework for representing and reasoning about uncertainty in medical data fusion. Each modality can be treated as a node in the network, and the dependencies between modalities can be modeled using conditional probability distributions [88]. Bayesian networks allow for principled fusion of multimodal information, enabling probabilistic inference and decision-making [89].

Multiple Kernel Learning (MKL) is a technique that combines multiple kernels, which capture different types of information or relationships, into a unified representation. Each modality can be represented using a separate kernel, and MKL methods can learn the optimal combination of these kernels to maximize the performance of the fusion task [90]. MKL allows for flexible and effective integration of information from different modalities. Feature-level fusion techniques combine features extracted from different modalities to create a unified feature representation. This can involve techniques like concatenation,

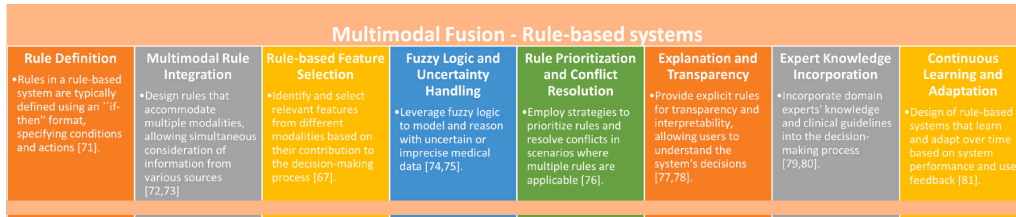


Fig. 5. Multimodal fusion — rule-based systems.

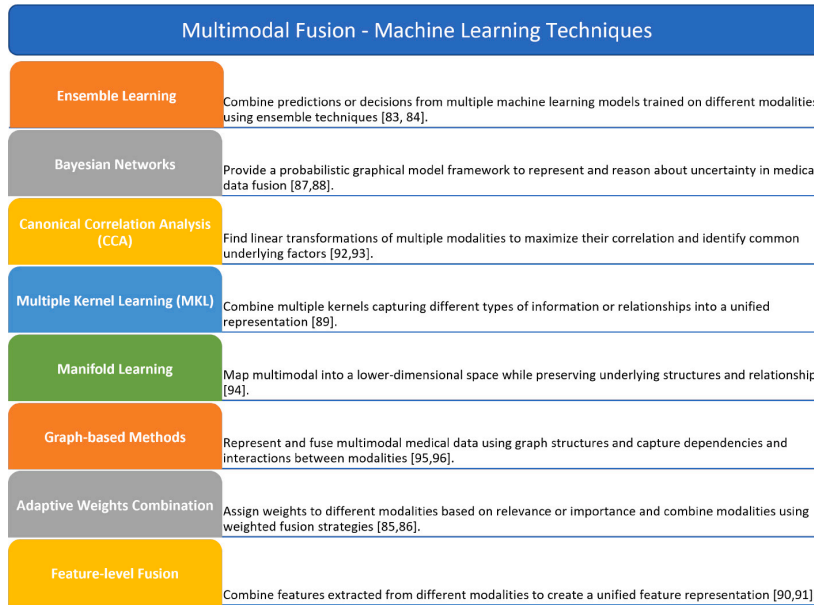


Fig. 6. Multimodal fusion — machine learning techniques.

feature stacking, or feature selection, based on relevance or mutual information [91]. The fused features can then be used as input for traditional ML algorithms, such as support vector machines (SVM), logistic regression, or k-nearest neighbors (k-NN), to perform classification, regression, or clustering tasks [92].

Canonical Correlation Analysis(CCA) is a statistical technique that aims to find linear transformations of multiple modalities to maximize their correlation. It identifies common underlying factors that explain the correlations between modalities [93]. The fused representation can then be used as input for subsequent ML algorithms [94]. Manifold learning techniques, such as t-SNE (t-Distributed Stochastic Neighbor Embedding) or Isomap, can be used to map multimodal data into a lower-dimensional space, while preserving the underlying structures and relationships [95]. By projecting the multimodal data onto a common latent space, these techniques facilitate the fusion of modalities and enable visualization and analysis of the fused data.

Graph-based methods offer a framework for representing and fusing multimodal medical data using graph structures. Each modality can be represented as nodes, and edges can be defined based on the relationships or correlations between modalities [96]. Graph-based algorithms, such as graph convolutional networks (GCNs) or graph regularized non-negative matrix factorization (GNMF), can then be applied to capture the dependencies and interactions between modalities [97].

ML techniques, even without deep learning, can still be effective in multimodal medical data fusion for smart healthcare. It is important to note that the selection of specific ML techniques for multimodal medical data fusion depends on the nature of the data, the fusion

task, and the available computational resources. Careful consideration should be given to feature selection, normalization, and data pre-processing steps to ensure optimal fusion performance. Additionally, model evaluation and validation using appropriate metrics and cross-validation techniques are crucial to assess the effectiveness of the fusion approach in smart healthcare applications. In Fig. 6, the ML techniques in multimodal fusion are presented.

3.4. Deep learning

Deep learning, as a subset of ML, plays a crucial role at the knowledge level of the DIKW framework. It specializes in training neural networks with multiple layers to extract intricate features and representations from data. By processing vast amounts of data, deep learning models can learn complex patterns and generate valuable knowledge for decision-making purposes. In the context of smart healthcare, deep learning has emerged as a potent approach for multimodal medical data fusion. Its unique capability to automatically learn hierarchical representations from diverse and complex medical modalities makes it particularly well-suited for integrating and extracting meaningful information. With its ability to handle diverse data types and capture intricate relationships, deep learning contributes to the overall knowledge-generation process in healthcare, enabling more accurate diagnoses, personalized treatment plans, and improved patient outcomes. Here are some key aspects of deep learning in multimodal medical data fusion.

Deep learning architectures, such as convolutional neural networks (CNNs), recurrent neural networks (RNNs), or transformers, can learn representations directly from multimodal medical data [98]. By jointly processing multiple modalities, these models capture both local and global dependencies, enabling the extraction of rich and informative features. This facilitates the fusion of modalities at various levels, ranging from low-level pixel or waveform data to high-level semantic representations [99]. Deep learning models with recurrent or temporal components, such as RNNs or long short-term memory (LSTM) networks, can handle sequential or temporal aspects of multimodal medical data [100]. These models can capture temporal dependencies, changes over time, or dynamic patterns across modalities. This is particularly relevant for applications such as physiological signal analysis, time-series data fusion, or modeling disease progression [101].

Deep learning models that have been pre-trained on large-scale datasets, such as ImageNet or natural language corpora, can be leveraged for multimodal medical data fusion [102]. Transfer learning techniques allow the transfer of knowledge from pretraining to the medical domain, enabling the models to learn relevant representations from limited medical data. This approach can boost performance, especially when multimodal medical datasets are small or resource-intensive to collect [103]. Attention mechanisms in deep learning models provide a mechanism for focusing on salient regions or modalities within the input data. They learn to allocate attention to the most relevant features or modalities, enhancing the fusion process [104]. Attention mechanisms can be employed within CNNs, RNNs, or transformer architectures to selectively combine or weigh the contributions of different modalities based on their importance for the task at hand [105].

Generative models, such as generative adversarial networks (GANs) or variational autoencoders (VAEs), can be used for multimodal fusion. These models learn to generate new samples from the joint distribution of multiple modalities, capturing their underlying correlations [106]. Generative models can aid in data augmentation, missing data imputation, or synthesis of multimodal data, facilitating improved training and fusion outcomes [107]. Deep fusion architectures combine multiple modalities at different stages of the network, allowing for the explicit integration of multimodal information. For example, early fusion involves combining modalities at the input level, while late fusion integrates modalities at higher layers or during decision-making [108]. Hybrid fusion approaches leverage both early and late fusion strategies to capture complementary information effectively. Deep fusion architectures can enhance the performance and robustness of multimodal fusion tasks [109].

Deep learning models will benefit from the integration of clinical knowledge, domain expertise, or prior medical information. Architectures that incorporate domain-specific constraints, expert rules, or Bayesian priors can enhance the fusion process and align the models with established medical knowledge [110]. Integrating clinical knowledge helps improve the interpretability, reliability, and acceptance of deep learning models in smart healthcare settings [111].

Deep learning models, although powerful, can be challenging to interpret. However, techniques such as attention visualization, saliency mapping, or gradient-based methods can provide insights into the model's decision-making process [112]. Interpretable deep learning architectures, such as CNNs with structured receptive fields or interpretable RNN variants, are also being explored to enhance transparency and explainability in multimodal medical data fusion [113].

Deep learning techniques offer promising avenues for multimodal medical data fusion, but challenges such as the need for large labeled datasets, interpretability, and generalization to new patient populations must be addressed. Collaboration between deep learning researchers, healthcare professionals, and data scientists is crucial to developing effective and reliable deep learning approaches for multimodal medical data fusion in smart healthcare. The deep learning techniques for multimodal fusion are outlined in Fig. 7

3.5. Natural language processing

NLP plays a vital role in transforming textual data into structured information, uncovering insights, and facilitating informed decision-making. In the context of multimodal medical data fusion for smart healthcare, NLP is particularly valuable for processing textual information from clinical notes, reports, and records. It extracts relevant details, identifies relationships, and discovers hidden patterns within the text. By incorporating NLP into the data fusion process, healthcare professionals can gain a comprehensive understanding of patients' health, improve diagnosis accuracy, and enhance personalized treatment planning. NLP is a powerful tool for integrating text-based information with other data modalities, enabling a holistic approach to healthcare decision-making.

NLP techniques are employed to process and extract meaningful information from unstructured textual data. Tasks such as tokenization, sentence segmentation, part-of-speech tagging, named entity recognition, and syntactic parsing help structure and analyze clinical text [114]. NLP enables the extraction of relevant concepts, medical terms, and relationships from textual data, facilitating their integration with other modalities [115]. NLP techniques can extract structured information from clinical narratives, such as diagnoses, medications, procedures, and patient demographics. Named entity recognition and relationship extraction algorithms identify and classify relevant entities and their associations, contributing to the fusion of textual information with other modalities. This extracted information can be used for decision support, clinical coding, or cohort identification [116]. NLP techniques, including semantic parsing, semantic role labeling, and medical concept normalization, each of which enable the understanding of clinical text in a structured manner [117]. This facilitates the extraction of clinical concepts, relations, and contextual knowledge, which can be fused with other modalities for comprehensive analysis, decision support, or knowledge discovery [118].

NLP models can be trained to classify clinical text into various categories, such as disease categories, severity levels, or treatment options [119,120]. Sentiment analysis techniques can also assess the sentiment or opinion expressed in patient feedback, social media data, or clinical notes. Text classification and sentiment analysis provide valuable insights, and can be integrated with other modalities for a comprehensive understanding of the patient's condition [3,121]. NLP can also bridge the gap between textual and visual modalities in medical data fusion. By analyzing textual descriptions or radiology reports, NLP techniques can extract relevant information about anatomical locations, findings, or abnormalities [122]. This information can be linked to corresponding images or visual data, enabling the fusion of text and image modalities for improved diagnosis, treatment planning, or image interpretation [123].

NLP methods aid in identifying and monitoring adverse events by analyzing textual data sources such as EHRs, patient complaints, or pharmacovigilance reports [124]. Sentiment analysis, information extraction, and text mining techniques can automatically detect and categorize adverse events [125], enabling timely interventions and enhancing patient safety [126]. NLP techniques contribute to patient risk assessment by analyzing clinical narratives and extracting relevant information related to patient history, comorbidities, and lifestyle factors [127]. By integrating this textual information with other patient data, such as vital signs, imaging results, or genetic information, multimodal medical data fusion can provide a comprehensive risk assessment strategy for personalized healthcare interventions and preventive measures [128].

NLP techniques support the development of clinical decision support systems by extracting relevant clinical knowledge from medical literature, clinical guidelines, or research articles [129,130].

NLP techniques provide valuable capabilities in extracting, processing, and integrating textual information within multimodal medical data fusion. They improve the comprehension of clinical text, enable the incorporation of clinical knowledge, and facilitate comprehensive and effective analysis in smart healthcare applications. Fig. 8 presents the NLP techniques that can be adopted for multimodal fusion.

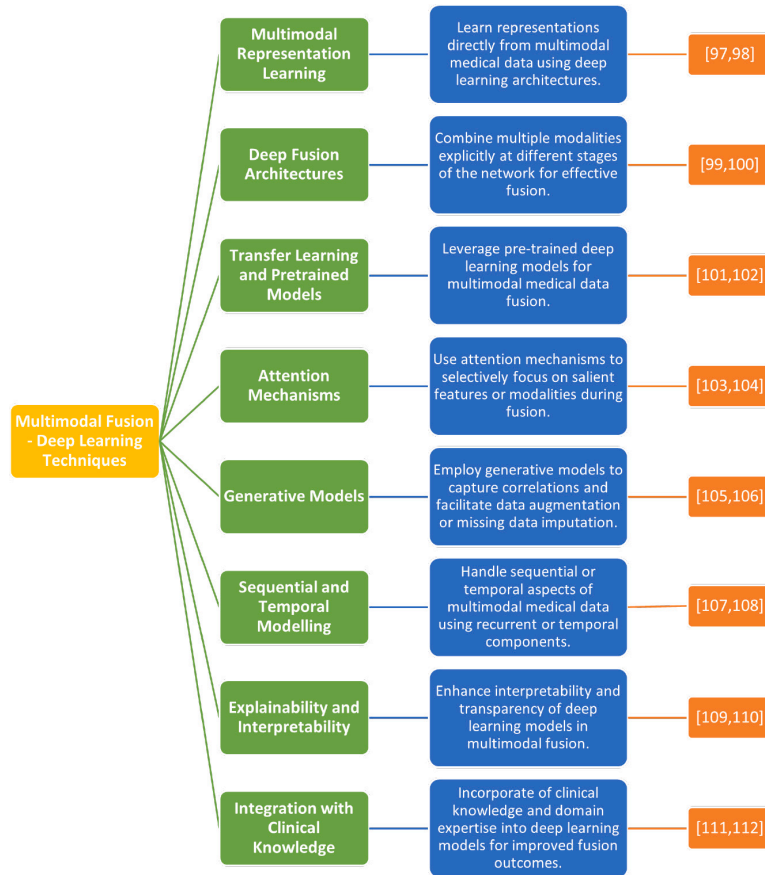


Fig. 7. Multimodal fusion — deep learning techniques.

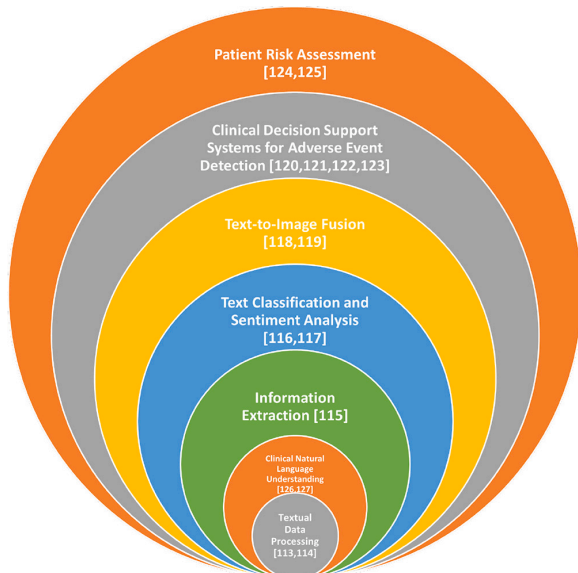


Fig. 8. Multimodal fusion — NLP techniques.

3.6. Taxonomy of approaches in multimodal fusion

The taxonomy of approaches in multimodal fusion for smart healthcare encompasses feature selection, rule-based systems, ML, deep learning, and NLP, as shown in Fig. 9. These techniques play a crucial role in integrating and analyzing diverse data modalities to extract valuable insights and support informed decision-making in healthcare applications.

- Feature selection focuses on identifying relevant features to create a concise representation of the data.
- Rule-based systems utilize predefined rules to process and combine data from multiple modalities.
- ML leverages patterns and relationships in the data to predict, classify, or cluster information from different modalities.
- Deep learning employs neural networks to automatically learn hierarchical representations and capture complex relationships.
- NLP techniques process and analyze textual information, enhancing the understanding of clinical data and facilitating its integration with other modalities.

By leveraging these approaches, researchers and healthcare professionals can gain a deeper understanding of patient health, enable personalized care, and make more informed decisions in intelligent healthcare systems, ultimately advancing patient care and well-being. In Fig. 9, we provide a comprehensive taxonomy of multimodal fusion methods relevant to smart healthcare. The Fig. 9 categorizes SOTA techniques, including feature selection, rule-based systems, machine

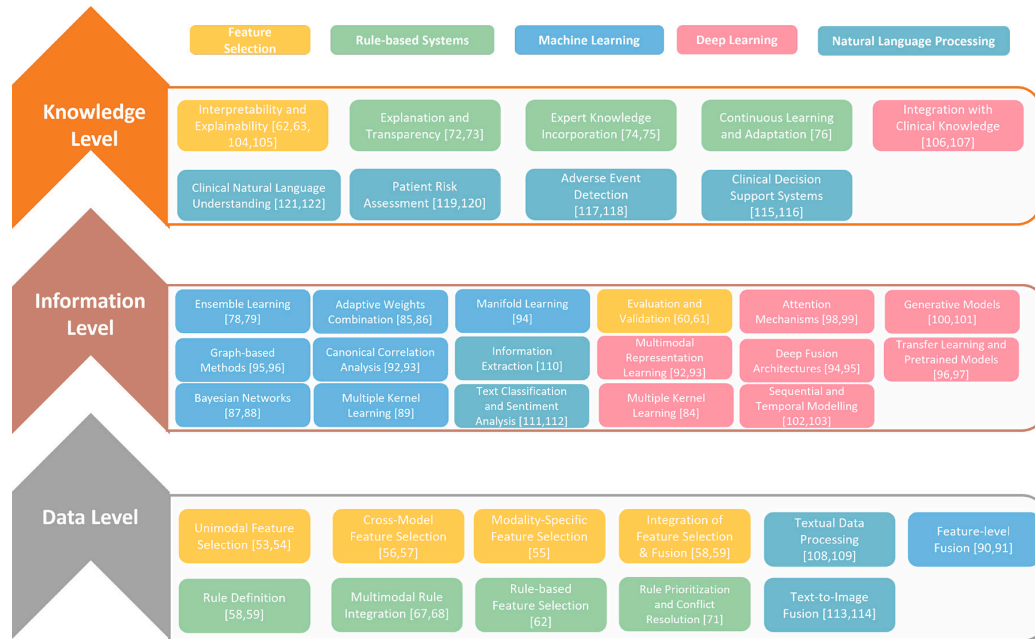


Fig. 9. Taxonomy of SOTA techniques in multimodal fusion for smart healthcare.

learning, deep learning, and NLP, and aligns them with the different levels of the DIKW conceptual model.

4. Challenges in adopting multimodal fusion

There are numerous challenges in adopting multimodal data fusion approaches that influence different stages and aspects of the DIKW framework. These challenges include issues such as data quality and interoperability, privacy and security, data processing and analysis, clinical integration and adoption, ethical considerations, and interpretation of results, all of which impact the transformation of data into meaningful information, knowledge, and wisdom in healthcare.

4.1. Data quality and interoperability

Data quality and interoperability present significant challenges in the context of multimodal fusion for smart healthcare within the DIKW framework. The integration of data from diverse sources and ensuring its quality and compatibility across different healthcare systems and modalities can be complex and time-consuming [131]. Insufficient data quality and a lack of interoperability can result in inaccurate analysis and hinder the effectiveness of the data fusion process [132].

Addressing these challenges necessitates the development and adoption of data standards and protocols. Standardizing data formats, such as HL7 for EHRs or DICOM for medical imaging, facilitates seamless integration and data exchange across healthcare systems and modalities [133]. These standards establish a shared language for data representation, simplifying the processing and integration of data from various sources [134]. Establishing interoperability frameworks plays a vital role in promoting smooth data sharing and exchange among different healthcare systems and modalities [135]. These frameworks provide guidelines and best practices for data integration, harmonization, and transmission, since they define the protocols, data models, and communication standards that enable efficient interoperability across diverse the data sources. Adhering to interoperability frameworks enhances data compatibility and coherence, thereby facilitating effective multimodal fusion [136].

Advancements in ML techniques offer promising avenues for addressing data integration and interoperability challenges [137]. ML algorithms can learn patterns and relationships in data obtained from different sources, facilitating automated data mapping, harmonization, and integration. Leveraging ML enables organizations to streamline the data fusion process, improving efficiency and accuracy in multimodal integration [138].

4.2. Privacy and security

Privacy and security pose significant challenges in the integration of sensitive patient data from multiple sources within the DIKW framework. Protecting patient privacy and ensuring data security are paramount concerns in healthcare, particularly when dealing with sensitive health information [139]. Robust privacy and security measures are necessary to safeguard patient confidentiality and prevent unauthorized access or data breaches in the integration of multimodal medical data.

To address these challenges, implementing data encryption techniques is essential to protect patient data during transmission and storage [140]. Encryption converts data into an unreadable form, ensuring that only authorized individuals with decryption keys can access and interpret the data [141]. Secure storage methods, such as secure servers or cloud platforms with robust access controls, play a crucial role in safeguarding patient data from unauthorized access, loss, or theft [142]. In addition, adopting privacy-preserving techniques is crucial to protect patient privacy during data fusion. Differential privacy, for example, adds noise to aggregated data to prevent individual identification while preserving the utility of the fused data [64,143]. Secure multiparty computation (SMC) techniques allow collaboration on data fusion without revealing individual-level data to any party involved, ensuring privacy during the fusion process [144]. The continuous monitoring and auditing of data access and usage are vital in detecting and preventing unauthorized activities and potential data breaches [145]. Robust auditing mechanisms and log analysis techniques enable organizations to track and investigate any suspicious or anomalous access patterns or data breaches [146]. Real-time monitoring systems can

provide alerts and notifications in case of any unauthorized access attempts or potential security incidents [147]. By implementing these privacy and security measures, healthcare organizations can ensure the protection of patient data, maintain privacy during data fusion, and mitigate the risks associated with unauthorized access or data breaches.

4.3. Data processing and analysis

Data processing and analysis play a critical role in the DIKW framework, particularly in multimodal medical data fusion. Challenges arise in handling large volumes of data, scalability of data processing, and extracting actionable insights from the fused data [148]. To address these challenges, ML algorithms and AI techniques are leveraged to enable efficient processing and analysis of multimodal medical data [70, 149]. Supervised and unsupervised learning algorithms are utilized for classification, prediction, and pattern discovery [150]. Deep learning models, such as CNNs and RNNs, are applied to tasks involving medical imaging and sequential data analysis [151]. Reinforcement learning techniques optimize treatment plans and interventions based on patient outcomes [152].

Collaboration between clinicians and data scientists are crucial in developing effective data processing and analysis solutions that align with clinical needs [153,154]. Integration of multimodal medical data fusion into Clinical Decision Support Systems (CDSS) enhances clinical decision-making and improves patient outcomes [155,156]. Scalable data processing techniques, including distributed computing frameworks and cloud computing platforms, handle large-scale datasets [157, 158]. Real-time data analytics enable immediate insights and proactive interventions [159]. Effective visualization techniques can aid in interpreting and communicating analysis results [160]. By addressing these challenges and utilizing advanced data processing and analysis techniques, healthcare organizations can unlock the full potential of multimodal medical data fusion within the DIKW framework.

4.4. Clinical integration and adoption

Clinical integration and adoption present significant challenges in the successful implementation of multimodal fusion within the DIKW framework in healthcare [70]. To address these challenges, involving clinicians and other stakeholders in the development and implementation process is crucial for both ensuring successful adoption and maximizing the impact of multimodal fusion in clinical practice [161,162]. Their input and feedback contributes towards designing technologies that align with clinical workflows and meet end-users needs [163].

The design of user-friendly interfaces and intuitive workflows is essential to facilitate the integration of multimodal fusion into clinical practice. Applying user-centered design principles and conducting usability testing can identify and address usability issues, enhancing user satisfaction and adoption rates [83,164]. Integrating multimodal fusion technologies with CDSS can enhance clinical decision-making processes by providing real-time recommendations and alerts based on the fused data [155,165]. Embedding these technologies into familiar clinical tools streamlines the integration process and facilitates adoption. To ensure successful adoption, it is crucial to provide adequate training and education to healthcare professionals [166,167]. Training programs should focus not only on the technical aspects, but also on the clinical relevance and potential impacts on patient care. Continual education and support help healthcare professionals remain proficient in using the technologies and stay updated on advancements.

4.5. Ethical considerations

Ethical considerations are of utmost importance in the context of multimodal medical data fusion within the DIKW framework. Ensuring

patient privacy, autonomy, and fairness is essential in utilizing patient data ethically [168]. Obtaining informed consent from patients is a fundamental ethical requirement in multimodal medical data fusion. Patients should be fully informed about the purpose, risks, and benefits of data fusion, and consent processes should be transparent and understandable [169]. Clear mechanisms for patients to withdraw their consent should also be provided.

Defining data ownership and governance policies is crucial. Healthcare organizations should establish guidelines to determine data ownership, usage, and access [170]. Transparent governance mechanisms, such as data access committees, should oversee the ethical use of patient data and compliance with privacy regulations. Respecting patient privacy and ensuring data confidentiality is paramount. Robust security measures, including encryption and access controls, should be implemented [171]. Compliance with privacy regulations like HIPAA or GDPR should be ensured.

Addressing potential biases is essential in multimodal fusion. Efforts should be made to mitigate biases through rigorous data collection processes, algorithmic fairness assessments, and diverse representation in data and development teams [172]. Regular monitoring and auditing can help identify and address biases. Ethical frameworks and guidelines should be developed and followed. These frameworks should outline principles and best practices for ethical data collection, fusion, analysis, and decision-making [173]. Guidelines should cover areas such as data privacy, informed consent, fairness, transparency, and accountability. Engaging the public and stakeholders in discussions about ethical considerations is crucial. Open communication channels should be established to foster trust and incorporate patient perspectives into decision-making processes [174]. By addressing ethical considerations, healthcare organizations can ensure the responsible and ethical use of patient data in multimodal medical data fusion, promoting patient privacy, fairness, and trust within the DIKW framework.

4.6. Interpretation of results

Interpreting the results of multimodal medical data fusion within the DIKW framework can be challenging due to the complexity of integrating multiple modalities and the large volumes of generated data. Effective interpretation is crucial for extracting meaningful insights and making actionable decisions in clinical settings [175].

To facilitate interpretation, the development of visual analytics tools and techniques is essential. Interactive visualizations, such as heatmaps, scatter plots, or network diagrams, can assist clinicians in identifying patterns, correlations, and outliers in the fused data [176]. These visualizations should provide intuitive representations and enable interactive exploration at different levels of detail [177]. Enhancing the interpretability of fusion models and algorithms is another important aspect. Techniques like explainable AI, interpretable machine learning, or rule-based systems can provide transparent explanations for the fusion process and decision-making [178].

Understanding how the fusion models arrive at certain conclusions helps clinicians trust the results and make informed decisions based on the interpretation of the fused data. Clinical validation studies are crucial for assessing the clinical utility and effectiveness of the interpretation results. Real-world evaluation with healthcare professionals provides insights into the practical applicability of the interpretation and helps refine the techniques to ensure meaningful and actionable results are aligned with clinical practice. Involving domain experts, such as clinicians or medical researchers, in the interpretation process is vital. They bring valuable insights into the clinical relevance of the fused data and help interpret the results within the context of patient care [179].

Collaboration between data scientists and domain experts fosters mutual understanding and enables the development of interpretation techniques that meet the specific needs of healthcare professionals [180]. Incorporating clinical guidelines and contextual information into the interpretation of fused data is crucial. Considering

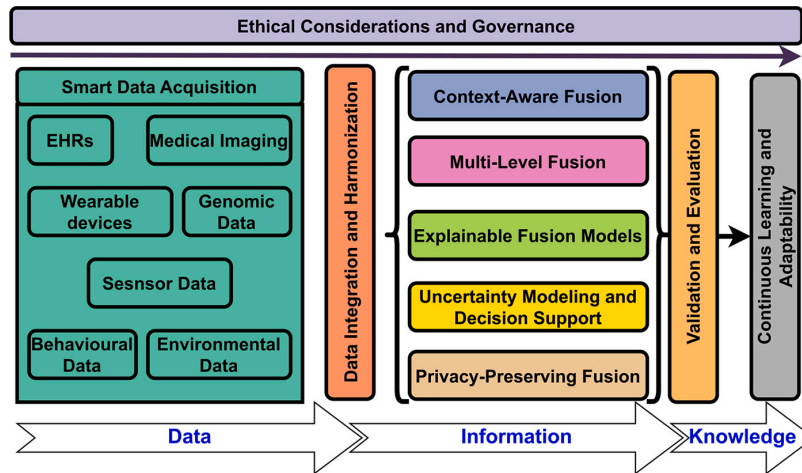


Fig. 10. Generic framework of multimodal fusion for smart healthcare.

patient-specific factors, such as demographics, medical history, or treatment guidelines, helps provide personalized interpretations and recommendations [181]. Aligning the interpretation with existing clinical knowledge and guidelines ensures clinically meaningful and actionable results for healthcare providers.

5. DIKW fusion framework with multimodality

Based on the data representation modalities and multimodal fusion approaches, we present a universal framework that can be applied to diverse applications in Fig. 10. In the framework for multimodal fusion in smart healthcare, several components are identified that facilitate the progression towards wisdom. At the data level, strategies for efficient data acquisition are employed, followed by data integration and harmonization processes to create a unified dataset.

Moving to the information level, context-aware fusion incorporates contextual information to enhance the fusion process, while multi-level fusion techniques capture complex relationships and patterns. Explainable fusion models provide transparency and trust, and uncertainty modeling supports decision-making based on fused data. Privacy-preserving fusion techniques ensure responsible data handling, and validation and evaluation methods assess the performance of the fusion framework. At the knowledge level, continuous learning and adaptability mechanisms enable the framework to stay up-to-date, while ethical considerations and governance frameworks address ethical issues in healthcare fusion.

In the landscape of multimodal data fusion for healthcare applications, the journey towards wisdom can be conceptualized as a hierarchical framework consisting of four integral stages: data fusion, information fusion, knowledge fusion, and ultimately, wisdom, as shown in Fig. 11. In the first stage, data fusion, we employ techniques such as feature selection, ensemble learning, and graph-based methods to combine and select the most relevant features from diverse data sources [182]. This crucial step forms a solid groundwork for effectively integrating and processing multimodal data.

Moving on to the second stage, information fusion, we delve deeper into the data by utilizing advanced techniques such as deep fusion architectures, transfer learning, attention mechanisms, and sequential modeling. These sophisticated approaches enable us to uncover intricate relationships and patterns across modalities, providing a more profound understanding of the data at hand. Additionally, we employ explainability and interpretability techniques to gain valuable insights into the decision-making process of the fusion models.

In the final stage, knowledge fusion, we take integration to the next level by incorporating clinical knowledge and domain expertise. Here, techniques like CDSS, adverse event detection, patient risk assessment, and clinical natural language understanding come into play. By leveraging this wealth of clinical knowledge, we can extract actionable insights and make informed decisions in the healthcare domain.

By following this progressive journey from data fusion to information fusion and knowledge fusion, we empower ourselves to enhance our understanding and analysis of multimodal data. This, in turn, contributes to the development of wisdom in the field of multimodal data fusion, enabling us to make more impactful advancements in smart healthcare and its applications.

6. Future directions of DIKW fusion in smart healthcare

The field of multimodal medical data fusion for smart healthcare is expected to evolve in line with the 4Ps of healthcare — Predictive, Preventive, Personalized, and Participatory [183,184]. Predictive healthcare data fusion aims to anticipate health events and outcomes by combining data from various sources, while Preventive data fusion focuses on identifying risk factors and promoting healthy behaviors. Personalized fusion caters to individual-specific data for customized care, and Participatory fusion involves patients and stakeholders in the data fusion process, enhancing transparency and trust. The progress towards these goals forms the crux of our future research.

6.1. Predictive healthcare

In the context of the DIKW framework and the generic framework discussed, the “Predictive” component of multimodal fusion in smart healthcare focuses on utilizing diverse data sources to anticipate health events and outcomes, thereby enabling proactive interventions and personalized healthcare strategies [183]. Within the DIKW framework, the Predictive component involves leveraging multimodal fusion to generate predictive models for assessing disease likelihoods [113]. By integrating and analyzing data from various modalities, such as genomics, medical imaging, and clinical records, healthcare professionals can identify early indicators or risk factors [185]. ML algorithms play a key role in analyzing the combined data and uncovering patterns indicative of disease risk [186].

In the generic framework, the Predictive component of multimodal fusion aims to identify potential risk factors or biomarkers for disease prediction [128]. Healthcare professionals gain insights into genetic predispositions, imaging abnormalities, and the clinical context [187].

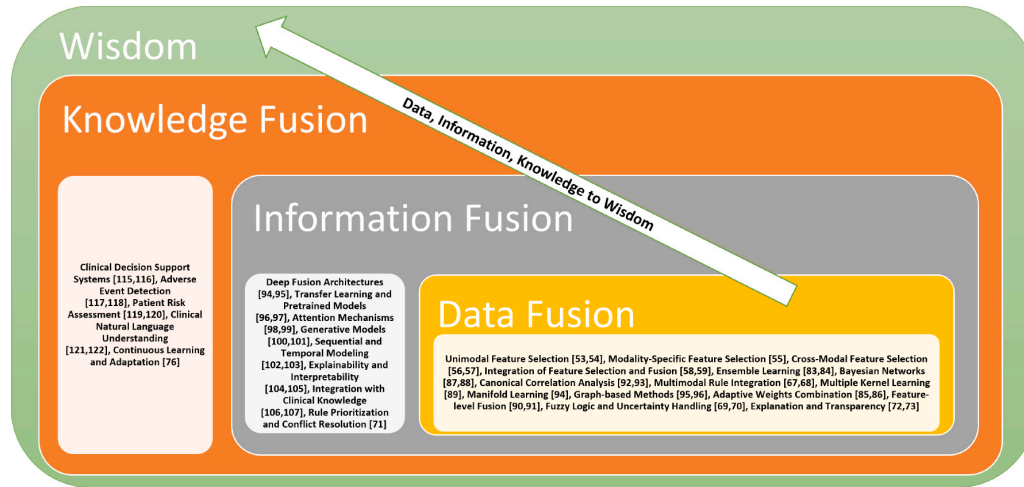


Fig. 11. DIKW conceptual model journey.

ML algorithms enable the development of predictive models by leveraging the combined data and considering variables such as genetic markers, imaging features, clinical parameters, lifestyle factors, and environmental exposures [188,189].

By applying predictive multimodal fusion within the framework, healthcare professionals can proactively identify individuals at higher risk of developing specific diseases or conditions, facilitating preventive interventions and personalized healthcare strategies [27,190]. For instance, early identification of individuals at risk of cardiovascular disease enables targeted lifestyle modifications, medication interventions, and regular monitoring to prevent or manage the condition [15]. Through the integration of multimodal data and the application of predictive analytics, the Predictive component enhances healthcare decision-making and supports proactive interventions, ultimately improving patient outcomes and healthcare delivery [191].

6.2. Preventive healthcare

In the DIKW framework, the “Preventive” component of multimodal fusion focuses on utilizing diverse data sources to develop personalized preventive strategies for patients. By integrating and analyzing data from various sources, such as mHealth devices and EHRs, healthcare providers can gain insights into patients’ lifestyle factors and health issues, enabling them to develop targeted interventions and preventive measures [113,190]. By combining data from sources such as mHealth devices and EHRs, healthcare providers can gain a comprehensive understanding of patients’ health status and develop personalized preventive strategies [72,192].

Multimodal data fusion enables the integration and analysis of data from diverse sources, such as wearable fitness trackers or smart watches for monitoring physical activity, heart rate, and sleep patterns, and EHRs containing medical history, diagnoses, and laboratory results [193]. By combining these data sources, healthcare providers obtain a holistic view of patients’ health and lifestyle factors [194]. Through multimodal fusion, healthcare providers can identify lifestyle factors that may contribute to patients’ health issues. Data from mHealth devices and EHRs may indicate emergent patterns such as sedentary behavior, inadequate sleep, or poor nutrition that may impact the development or progression of certain health conditions [187–189].

Based on the analysis of multimodal data, healthcare providers can develop personalized preventive strategies tailored to each patient’s unique needs. For instance, if data fusion analysis reveals that a patient’s sedentary behavior contributes to their health issues, healthcare

providers may prescribe regular physical activity, provide educational materials on exercise, or suggest behavioral change techniques to promote a more active lifestyle [195]. Furthermore, multimodal fusion facilitates ongoing monitoring and feedback, empowering patients to maintain their preventive strategies and make informed decisions about their health. By leveraging technology and data integration, patients can receive real-time feedback on their health behaviors, track their progress, and receive personalized recommendations to support their preventive efforts [196,197]. By incorporating the Preventive component within the DIKW and generic frameworks, multimodal fusion plays a vital role in developing personalized preventive strategies for patients, aligning with the broader goals of precision medicine and personalized healthcare [184,198].

6.3. Personalized healthcare

In the DIKW framework, the “Personalized” component of multimodal fusion focuses on utilizing diverse data sources to develop personalized treatment plans for patients. By integrating and analyzing data from different modalities, such as imaging and genomics, healthcare providers can gain a deeper understanding of the patient’s molecular profile and tailor treatment strategies based on their individual characteristics [113,190].

Within the generic framework, the Personalized component of multimodal fusion aims to identify specific genetic mutations or variations that underpin the patient’s disease. By combining imaging data with genomics data, healthcare providers can obtain a comprehensive view of the patient’s health condition and uncover genetic markers associated with conditions such as tumor growth [72,192]. Multimodal data fusion allows for the integration and analysis of data from diverse sources, such as magnetic resonance imaging (MRI), computed tomography (CT), positron emission tomography (PET), and genomics data. These modalities provide detailed anatomical, functional, and genetic information about the patient’s body and disease state [194].

By fusing imaging and genomics data, healthcare providers can identify specific genetic mutations or variations that inform the underlying molecular mechanisms of the disease. This knowledge guides the development of personalized treatment plans tailored to the patient’s individual molecular profile [195]. The analysis of multimodal fusion enables healthcare providers to make informed treatment decisions, such as recommending targeted therapies for specific genetic mutations. This personalized approach ensures that patients receive the most effective treatments based on their unique genetic characteristics [196].

Moreover, multimodal fusion facilitates ongoing monitoring of treatment response and enables treatment adjustments over time. By integrating data from imaging, genomics, and other modalities, healthcare providers can assess treatment effectiveness and make informed decisions regarding treatment modifications to optimize patient outcomes [197]. By incorporating the Personalized component within the DIKW and generic frameworks, multimodal fusion can play a significant role in tailoring treatment plans based on individual patient characteristics and molecular profiles. This personalized approach to patient care and treatment outcomes aligns with the broader goals of precision medicine and personalized healthcare [198].

6.4. Participatory healthcare

In the DIKW framework, the “Participatory” component of multimodal fusion in smart healthcare focuses on empowering patients to actively participate in their own healthcare journey, fostering collaboration and shared decision-making with healthcare providers [199, 200]. Within the DIKW framework, the Participatory component involves leveraging multimodal fusion to provide patients with real-time feedback and insights into their health status [201]. By integrating data from mHealth devices and patient-reported information, patients can actively monitor and track their health, enabling them to make informed decisions about their well-being [202].

In the generic framework, the Participatory component of multimodal fusion emphasizes the active engagement of patients in their healthcare by combining data from mHealth devices and patient-reported data [203]. Through real-time access to their health information, patients can receive personalized feedback and recommendations, and participate in discussions about their treatment plans [199]. This collaborative approach enables patients to actively contribute to decision-making based on their preferences, values, and personal health goals. Multimodal data fusion also enables patients to participate in larger-scale initiatives, such as contributing their data to aggregated and anonymized datasets [200]. By participating in research studies, clinical trials, or public health monitoring programs, individuals can contribute to advancements in medical research, personalized interventions, and population health initiatives. By embracing the Participatory component of multimodal fusion, patients become active partners in their own healthcare, and can be further empowered to make proactive choices and actively contribute to their own well-being. This collaborative approach enhances patient-centered care and fosters a stronger partnership between patients and healthcare providers.

7. Conclusion

Multimodal medical data fusion, integrating various modalities like EHRs, medical imaging, wearable devices, genomic data, sensor data, environmental data, and behavioral data, has the potential to revolutionize smart healthcare. By leveraging approaches such as feature selection, rule-based systems, ML, deep learning, and NL, practitioners can extract valuable insights from a wealth of diverse sources, which will advance gains in knowledge and wisdom in healthcare.

However, the challenges related to data quality, interoperability, privacy, security, data processing, clinical integration, and ethical considerations must be addressed. Future research should focus on Predictive, Preventive, Personalized, and Participatory approaches, the implementation or combination of which can enable better anticipation of health events, identify risk factors, deliver tailored interventions, or further empower patients in their healthcare journeys. Embracing these opportunities will transform healthcare by improving patient well-being, treatment outcomes, and the overall function of the healthcare industry.

CRediT authorship contribution statement

Thanveer Shaik: Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Writing – original draft, Visualization, Writing – review & editing. **Xiaohui Tao:** Conceptualization, Formal analysis, Investigation, Methodology, Project administration, Supervision, Writing – original draft, Visualization, Writing – review & editing. **Lin Li:** Conceptualization, Investigation, Methodology, Supervision, Writing – review & editing. **Haoran Xie:** Investigation, Methodology, Writing – review & editing. **Juan D. Velásquez:** Investigation, Methodology, Writing – review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

This is a survey article. The datasets mentioned in the paper has been cited.

Acknowledgments

The authors gratefully acknowledge financial support ANID Fondo-cyt 1231122, PIA/PUENTE AFB220003, Chile.

References

- [1] R.L. Ackoff, From data to wisdom, *J. Appl. Syst. Anal.* 16 (1) (1989) 3–9.
- [2] S.M. Fiore, J. Elias, E. Salas, N.W. Warner, M.P. Letsky, From data, to information, to knowledge: Measuring knowledge building in the context of collaborative cognition, in: *Macro-cognition Metrics and Scenarios*, CRC Press, 2018, pp. 179–200.
- [3] X. Tao, T. Pham, J. Zhang, J. Yong, W.P. Goh, W. Zhang, Y. Cai, Mining health knowledge graph for health risk prediction, *World Wide Web* 23 (2020) 2341–2362.
- [4] J. Liang, Y. Li, Z. Zhang, D. Shen, J. Xu, X. Zheng, T. Wang, B. Tang, J. Lei, J. Zhang, et al., Adoption of electronic health records (EHRs) in China during the past 10 years: consecutive survey data analysis and comparison of sino-american challenges and experiences, *J. Med. Internet Res.* 23 (2) (2021) e24813.
- [5] Z. Zhang, E.P. Navarese, B. Zheng, Q. Meng, N. Liu, H. Ge, Q. Pan, Y. Yu, X. Ma, Analytics with artificial intelligence to advance the treatment of acute respiratory distress syndrome, *J. Evid.-Based Med.* 13 (4) (2020) 301–312.
- [6] E. Hossain, R. Rana, N. Higgins, J. Soar, P.D. Barua, A.R. Pisani, K. Turner, Use of AI/ML-enabled state-of-the-art method in electronic medical records: A systematic review, *Comput. Biol. Med.* (2023) 106649.
- [7] B. Ihnaini, M. Khan, T.A. Khan, S. Abbas, M.S. Daoud, M. Ahmad, M.A. Khan, A smart healthcare recommendation system for multidisciplinary diabetes patients with data fusion based on deep ensemble learning, *Comput. Intell. Neurosci.* 2021 (2021).
- [8] Z. Xu, D.R. So, A.M. Dai, Mufasa: Multimodal fusion architecture search for electronic health records, in: *Proceedings of the AAAI Conf. on Artificial Intelligence*, Vol. 35, (12) 2021, pp. 10532–10540.
- [9] Y. An, H. Zhang, Y. Sheng, J. Wang, X. Chen, MAIN: Multimodal attention-based fusion networks for diagnosis prediction, in: *2021 IEEE Int'l Conf. on Bioinformatics and Biomedicine, (BIBM)*, IEEE, 2021, pp. 809–816.
- [10] S. Malakar, S.D. Roy, S. Das, S. Sen, J.D. Velásquez, R. Sarkar, Computer based diagnosis of some chronic diseases: A medical journey of the last two decades, *Arch. Comput. Methods Eng.* (2022) 1–43.
- [11] A. Papa, M. Mital, P. Pisano, M. Del Giudice, E-health and wellbeing monitoring using smart healthcare devices: An empirical investigation, *Technol. Forecast. Soc. Change* 153 (2020) 119226.
- [12] E. Teixeira, H. Fonseca, F. Diniz-Sousa, L. Veras, G. Boppre, J. Oliveira, D. Pinto, A.J. Alves, A. Barbosa, R. Mendes, et al., Wearable devices for physical activity and healthcare monitoring in elderly people: A critical review, *Geriatrics* 6 (2) (2021) 38.
- [13] A. Sheth, U. Jaimini, H.Y. Yip, How will the internet of things enable augmented personalized health? *IEEE Intell. Syst.* 33 (1) (2018) 89–97.
- [14] T. Shaik, X. Tao, N. Higgins, L. Li, R. Gururajan, X. Zhou, U.R. Acharya, Remote patient monitoring using artificial intelligence: Current state, applications, and challenges, *Wiley Interdiscip. Rev.: Data Min. Knowl. Discov.* (2023) e1485.

- [15] X. Tao, T.B. Shaik, N. Higgins, R. Gururajan, X. Zhou, Remote patient monitoring using radio frequency identification (RFID) technology and machine learning for early detection of suicidal behaviour in mental health facilities, *Sensors* 21 (3) (2021) 776.
- [16] K. Mohammed, A. Zaidan, B. Zaidan, O.S. Albahri, M. Alsaalem, A.S. Albahri, A. Hadi, M. Hashim, Real-time remote-health monitoring systems: A review on patients prioritisation for multiple-chronic diseases, taxonomy analysis, concerns and solution procedure, *J. Med. Syst.* 43 (2019) 1–21.
- [17] L.A. Durán-Vega, P.C. Santana-Mancilla, R. Buenrosto-Mariscal, J. Contreras-Castillo, L.E. Anido-Rifón, M.A. García-Ruiz, O.A. Montesinos-López, F. Estrada-González, An IoT system for remote health monitoring in elderly adults through a wearable device and mobile application, *Geriatrics* 4 (2) (2019) 34.
- [18] S. Tian, W. Yang, J.M. Le Grange, P. Wang, W. Huang, Z. Ye, Smart healthcare: making medical care more intelligent, *Global Health J.* 3 (3) (2019) 62–65.
- [19] M. Senbekov, T. Saliev, Z. Bukeyeva, A. Almbayeva, M. Zhanaliyeva, N. Aitenova, Y. Toishibekov, I. Fakhradiyev, et al., The recent progress and applications of digital technologies in healthcare: A review, *Int. J. Telemed. Appl.* 2020 (2020).
- [20] M.S. Linet, T.L. Slovis, D.L. Miller, R. Kleinerman, C. Lee, P. Rajaraman, A. Berrington de Gonzalez, Cancer risks associated with external radiation from diagnostic imaging procedures, *CA: Cancer J. Clin.* 62 (2) (2012) 75–100.
- [21] X. Liu, L. Faes, A.U. Kale, S.K. Wagner, D.J. Fu, A. Bruynseels, T. Mahendiran, G. Moraes, M. Shandas, C. Kern, et al., A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging: A systematic review and meta-analysis, *Lancet Digital Health* 1 (6) (2019) e271–e297.
- [22] A. Garain, A. Basu, F. Giampaolo, J.D. Velasquez, R. Sarkar, Detection of COVID-19 from CT scan images: A spiking neural network-based approach, *Neural Comput. Appl.* 33 (19) (2021) 12591–12604.
- [23] S. Das, S.D. Roy, S. Malakar, J.D. Velásquez, R. Sarkar, Bi-level prediction model for screening COVID-19 patients using chest X-ray images, *Big Data Res.* 25 (2021) 100233.
- [24] J.B. Awotunde, F.E. Ayo, R.G. Jimoh, R.O. Ogundokun, O.E. Matiluko, I.D. Oladipo, M. Abdulraheem, Prediction and classification of diabetes mellitus using genomic data, in: *Intelligent IoT Systems in Personalized Health Care*, Elsevier, 2021, pp. 235–292.
- [25] H. Yu, H. Yan, L. Wang, J. Li, L. Tan, W. Deng, Q. Chen, G. Yang, F. Zhang, T. Lu, et al., Five novel loci associated with antipsychotic treatment response in patients with schizophrenia: A genome-wide association study, *Lancet Psychiatry* 5 (4) (2018) 327–338.
- [26] S. Pai, G.D. Bader, Patient similarity networks for precision medicine, *J. Mol. Biol.* 430 (18) (2018) 2924–2938.
- [27] J.N. Acosta, G.J. Falcone, P. Rajpurkar, E.J. Topol, Multimodal biomedical AI, *Nature Med.* 28 (9) (2022) 1773–1784.
- [28] O. Taiwo, A.E. Ezugwu, Smart healthcare support for remote patient monitoring during covid-19 quarantine, *Inform. Med. Unlocked* 20 (2020) 100428.
- [29] C. Carlsten, S. Salvi, G.W. Wong, K.F. Chung, Personal strategies to minimise effects of air pollution on respiratory health: advice for providers, patients and the public, *Eur. Respir. J.* 55 (6) (2020).
- [30] M. Hu, J.D. Roberts, G.P. Azevedo, D. Milner, The role of built and social environmental factors in Covid-19 transmission: A look at America's capital city, *Sustainable Cities Soc.* 65 (2021) 102580.
- [31] E.A. Alvarez, M. Garrido, D.P. Ponce, G. Pizarro, A.A. Córdova, F. Vera, R. Ruiz, R. Fernández, J.D. Velásquez, E. Tobar, et al., A software to prevent delirium in hospitalised older adults: development and feasibility assessment, *Age Ageing* 49 (2) (2020) 239–245.
- [32] T.J. Pollard, A.E. Johnson, J.D. Raffa, L.A. Celi, R.G. Mark, O. Badawi, The eICU collaborative research database, a freely available multi-center database for critical care research, *Sci. Data* 5 (1) (2018) 1–13.
- [33] A.E. Johnson, T.J. Pollard, L. Shen, L.-w.H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. Anthony Celi, R.G. Mark, MIMIC-III, a freely accessible critical care database, *Sci. Data* 3 (1) (2016) 1–9.
- [34] D. Azcona, K. McGuinness, A.F. Smeaton, A comparative study of existing and new deep learning methods for detecting knee injuries using the mrnet dataset, 2020, *arXiv preprint arXiv:2010.01947*.
- [35] G. Shih, C.C. Wu, S.S. Halabi, M.D. Kohli, L.M. Prevedello, T.S. Cook, A. Sharma, J.K. Amorosa, V. Arteaga, M. Galperin-Aizenberg, et al., Augmenting the national institutes of health chest radiograph dataset with expert annotations of possible pneumonia, *Radiology: Artif. Intell.* 1 (1) (2019) e180041.
- [36] P. Rajpurkar, J. Irvin, A. Bagul, D. Ding, T. Duan, H. Mehta, B. Yang, K. Zhu, D. Laird, R.L. Ball, et al., Mura: Large dataset for abnormality detection in musculoskeletal radiographs, 2017, *arXiv preprint arXiv:1712.06957*.
- [37] S.S. Halabi, L.M. Prevedello, J. Kalpathy-Cramer, A.B. Mamonov, A. Bilbily, M. Cicero, I. Pan, L.A. Pereira, R.T. Sousa, N. Abdala, et al., The RSNA pediatric bone age machine learning challenge, *Radiology* 290 (2) (2019) 498–503.
- [38] D. Demner-Fushman, M.D. Kohli, M.B. Rosenman, S.E. Shooshan, L. Rodriguez, S. Antani, G.R. Thoma, C.J. McDonald, Preparing a collection of radiology examinations for distribution and retrieval, *J. Am. Med. Inform. Assoc.* 23 (2) (2016) 304–310.
- [39] J. Zbontar, F. Knoll, A. Sriram, T. Murrell, Z. Huang, M.J. Muckley, A. Defazio, R. Stern, P. Johnson, M. Bruno, et al., fastMRI: An open dataset and benchmarks for accelerated MRI, 2018, *arXiv preprint arXiv:1811.08839*.
- [40] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Ilcus, C. Chute, H. Marklund, B. Haghgoo, R. Ball, K. Shpanskaya, et al., Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison, in: *Proceedings of the AAAI Conf. on Artificial Intelligence*, Vol. 33, (01) 2019, pp. 590–597.
- [41] D.S. Marcus, T.H. Wang, J. Parker, J.G. Csernansky, J.C. Morris, R.L. Buckner, Open access series of imaging studies (OASIS): cross-sectional MRI data in young, middle aged, nondemented, and demented older adults, *J. Cogn. Neurosci.* 19 (9) (2007) 1498–1507.
- [42] S.G. Armato III, G. McLennan, L. Bidaut, M.F. McNitt-Gray, C.R. Meyer, A.P. Reeves, B. Zhao, D.R. Aberle, C.I. Henschke, E.A. Hoffman, et al., The lung image database consortium (LIDC) and image database resource initiative (IDRI): A completed reference database of lung nodules on ct scans, *Med. Phys.* 38 (2) (2011) 915–931.
- [43] K. Clark, B. Vendt, K. Smith, J. Freymann, J. Kirby, P. Koppel, S. Moore, S. Phillips, D. Maffitt, M. Pringle, et al., The cancer imaging archive (TCIA): maintaining and operating a public information repository, *J. Digit. Imaging* 26 (2013) 1045–1057.
- [44] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, R.M. Summers, Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases, in: *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2017, pp. 2097–2106.
- [45] B.H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest, et al., The multimodal brain tumor image segmentation benchmark (BRATS), *IEEE Trans. Med. Imaging* 34 (10) (2014) 1993–2024.
- [46] S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, J.S. Kirby, J.B. Freymann, K. Farahani, C. Davatzikos, Advancing the cancer genome atlas glioma MRI collections with expert segmentation labels and radiomic features, *Sci. Data* 4 (1) (2017) 1–13.
- [47] S. Bakas, M. Reyes, A. Jakab, S. Bauer, M. Rempfler, A. Crimi, R.T. Shinohara, C. Berger, S.M. Ha, M. Rozycki, et al., Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BRATS challenge, 2018, *arXiv preprint arXiv:1811.02629*.
- [48] K. Tomczak, P. Czerwińska, M. Wiznerowicz, Review the cancer genome atlas (TCGA): An immeasurable source of knowledge, *Contemp. Oncol./Współczesna Onkologia* 2015 (1) (2015) 68–77.
- [49] N.E. Allen, C. Sudlow, T. Peakman, R. Collins, U. biobank, UK biobank data: come and get it, *Sci. Transl. Med.* 6 (224) (2014) 224ed4.
- [50] C.R. Jack, Jr., M.A. Bernstein, N.C. Fox, P. Thompson, G. Alexander, D. Harvey, B. Borowski, P.J. Britson, J. L. Whitwell, C. Ward, et al., The alzheimer's disease neuroimaging initiative (ADNI): MRI methods, *J. Magn. Reson. Imaging: Official J. Int. Soc. Magn. Reson. Med.* 27 (4) (2008) 685–691.
- [51] A. García Seco de Herrera, R. Schaer, S. Bromuri, H. Müller, Overview of the ImageCLEF 2016 medical task, in: *Working Notes of CLEF 2016 (Cross Language Evaluation Forum)*, 2016.
- [52] D. Demner-Fushman, S. Antani, M. Simpson, G.R. Thoma, Design and development of a multimodal biomedical information retrieval system, *J. Comput. Sci. Eng.* 6 (2) (2012) 168–177.
- [53] A.L. Goldberger, L.A.N. Amaral, L. Glass, J.M. Hausdorff, P.C. Ivanov, R.G. Mark, J.E. Mietus, G.B. Moody, C.-K. Peng, H.E. Stanley, PhysioBank, PhysioToolkit, and PhysioNet: Components of a new research resource for complex physiologic signals, *Circulation* 101 (23) (2000 (June 13)) e215–e220, <http://dx.doi.org/10.1161/01.CIR.101.23.e215>, *Circulation Electronic Pages: http://circ.ahajournals.org/content/101/23/e215.full* PMID:1085218.
- [54] T. Feng, B.M. Booth, B. Baldwin-Rodríguez, F. Osorno, S. Narayanan, A multimodal analysis of physical activity, sleep, and work shift in nurses with wearable sensor data, *Sci. Rep.* 11 (1) (2021) 8693.
- [55] S. Zeadally, O. Bello, Harnessing the power of internet of things based connectivity to improve healthcare, *Internet Things* 14 (2021) 100074.
- [56] K. Woodward, E. Kanjo, D.J. Brown, T.M. McGinnity, B. Inkster, D.J. Macintyre, A. Tsanas, Beyond mobile apps: A survey of technologies for mental well-being, *IEEE Trans. Affect. Comput.* 13 (3) (2020) 1216–1235.
- [57] S. Soklaridis, E. Lin, Y. Lalani, T. Rodak, S. Sockalingam, Mental health interventions and supports during COVID-19 and other medical pandemics: A rapid systematic review of the evidence, *Gen. Hosp. Psychiatry* 66 (2020) 133–146.
- [58] P. Bhowal, S. Sen, J.D. Velasquez, R. Sarkar, Fuzzy ensemble of deep learning models using choquet fuzzy integral, coalition game and information theory for breast cancer histology classification, *Expert Syst. Appl.* 190 (2022) 116167.
- [59] A. Albahri, A.M. Duhaime, M.A. Fadhel, A. Alnoor, N.S. Baqer, L. Alzubaidi, O. Albahri, A. Alamoodi, J. Bai, A. Salhi, et al., A systematic review of trustworthy and explainable artificial intelligence in healthcare: assessment of quality, bias risk, and data fusion, *Inf. Fusion* (2023).
- [60] S.M. Alghowinem, T. Gedeon, R. Goecke, J. Cohn, G. Parker, Interpretation of depression detection models via feature selection methods, *IEEE Trans. Affect. Comput.* (2020).

- [61] J. Zhang, Z. Yin, P. Chen, S. Nichele, Emotion recognition using multi-modal data and machine learning techniques: A tutorial and review, *Inf. Fusion* 59 (2020) 103–126.
- [62] T. Zhang, M. Shi, Multi-modal neuroimaging feature fusion for diagnosis of alzheimer disease, *J. Neurosci. Methods* 341 (2020) 108795.
- [63] T. Zhou, M. Liu, K.-H. Thung, D. Shen, Latent representation learning for alzheimer's disease diagnosis with incomplete multi-modality neuroimaging and genetic data, *IEEE Trans. Med. Imaging* 38 (10) (2019) 2411–2422.
- [64] D. Kim, Y.-H. Tsai, B. Zhuang, X. Yu, S. Sclaroff, K. Saenko, M. Chandraker, Learning cross-modal contrastive features for video domain adaptation, in: *Proceedings of the IEEE/CVF Int'l Conf. on Computer Vision*, 2021, pp. 13618–13627.
- [65] T. Hoang, T.-T. Do, T.V. Nguyen, N.-M. Cheung, Multimodal mutual information maximization: A novel approach for unsupervised deep cross-modal hashing, *IEEE Trans. Neural Netw. Learn. Syst.* (2022).
- [66] S. Qiu, H. Zhao, N. Jiang, Z. Wang, L. Liu, Y. An, H. Zhao, X. Miao, R. Liu, G. Fortino, Multi-sensor information fusion based on machine learning for real applications in human activity recognition: State-of-the-art and research challenges, *Inf. Fusion* 80 (2022) 241–265.
- [67] M. Abdel-Basset, D. El-Shahhat, I. El-Henawy, V.H.C. De Albuquerque, S. Mirjalili, A new fusion of grey wolf optimizer algorithm with a two-phase mutation for feature selection, *Expert Syst. Appl.* 139 (2020) 112824.
- [68] J. Zhou, A.H. Gandomi, F. Chen, A. Holzinger, Evaluating the quality of machine learning explanations: A survey on methods and metrics, *Electronics* 10 (5) (2021) 593.
- [69] X. Hao, Y. Bao, Y. Guo, M. Yu, D. Zhang, S.L. Risacher, A.J. Saykin, X. Yao, L. Shen, A.D.N. Initiative, et al., Multi-modal neuroimaging feature selection with consistent metric constraint for diagnosis of alzheimers disease, *Med. Image Anal.* 60 (2020) 101625.
- [70] G. Yang, Q. Ye, J. Xia, Unbox the black-box for the medical explainable AI via multi-modal and multi-centre data fusion: A mini-review, two showcases and beyond, *Inf. Fusion* 77 (2022) 29–52.
- [71] Y. Zhang, P. Tiño, A. Leonadis, K. Tang, A survey on neural network interpretability, *IEEE Trans. Emerg. Top. Comput. Intell.* 5 (5) (2021) 726–742.
- [72] G. Muhammad, F. Alshehri, F. Karray, A. El Saddik, M. Alsulaiman, T.H. Falk, A comprehensive survey on multimodal medical signals fusion for smart healthcare systems, *Inf. Fusion* 76 (2021) 355–375.
- [73] M. Hussain, F.A. Satti, S.I. Ali, J. Hussain, T. Ali, H.-S. Kim, K.-H. Yoon, T. Chung, S. Lee, Intelligent knowledge consolidation: from data to wisdom, *Knowl.-Based Syst.* 234 (2021) 107578.
- [74] T. Chen, C. Shang, P. Su, E. Keravnou-Papailiou, Y. Zhao, G. Antoniou, Q. Shen, A decision tree-initialised neuro-fuzzy approach for clinical decision support, *Artif. Intell. Med.* 111 (2021) 101986.
- [75] T.K. Mohd, N. Nguyen, A.Y. Javaid, Multi-modal data fusion in enhancing human-machine interaction for robotic applications: A survey, 2022, arXiv preprint arXiv:2202.07732.
- [76] R. Yan, F. Zhang, X. Rao, Z. Lv, J. Li, L. Zhang, S. Liang, Y. Li, F. Ren, C. Zheng, et al., Richer fusion network for breast cancer classification based on multimodal data, *BMC Med. Inform. Decis. Mak.* 21 (1) (2021) 1–15.
- [77] A. Amirkhani, E.I. Papageorgiou, M.R. Mosavi, K. Mohammadi, A novel medical decision support system based on fuzzy cognitive maps enhanced by intuitive and learning capabilities for modeling uncertainty, *Appl. Math. Comput.* 337 (2018) 562–582.
- [78] A. Geramian, A. Abraham, M. Ahmadi Nozari, Fuzzy logic-based FMEA robust design: A quantitative approach for robustness against groupthink in group/team decision-making, *Int. J. Prod. Res.* 57 (5) (2019) 1331–1344.
- [79] A. Alharbi, A. Poujade, K. Malandrakis, I. Petrunin, D. Panagiotakopoulos, A. Tsourdos, Rule-based conflict management for unmanned traffic management scenarios, in: *2020 AIAA/IEEE 39th Digital Avionics Systems Conf., (DASC), IEEE*, 2020, pp. 1–10.
- [80] K. Bahani, M. Moujabbir, M. Ramdani, An accurate fuzzy rule-based classification systems for heart disease diagnosis, *Sci. Afr.* 14 (2021) e01019.
- [81] A.M. Antoniadis, Y. Du, Y. Guendouz, L. Wei, C. Mazo, B.A. Becker, C. Mooney, Current challenges and future opportunities for XAI in machine learning-based clinical decision support systems: A systematic review, *Appl. Sci.* 11 (11) (2021) 5088.
- [82] W.-T. Wang, S.-Y. Wu, Knowledge management based on information technology in response to COVID-19 crisis, *Knowl. Manag. Res. Pract.* 19 (4) (2021) 468–474.
- [83] L. Rundo, R. Pirrone, S. Vitabile, E. Sala, O. Gambino, Recent advances of HCI in decision-making tasks for optimized clinical workflows and precision medicine, *J. Biomed. Inform.* 108 (2020) 103479.
- [84] S. El-Sappagh, F. Ali, T. Abuhmed, J. Singh, J.M. Alonso, Automatic detection of alzheimer's disease progression: An efficient information fusion approach with heterogeneous ensemble classifiers, *Neurocomputing* 512 (2022) 203–224.
- [85] K. Srinivasan, K. Raman, J. Chen, M. Bendersky, M. Najork, Wit: Wikipedia-based image text dataset for multimodal multilingual machine learning, in: *Proceedings of the 44th Int'l ACM SIGIR Conf. on Research and Development in Information Retrieval*, 2021, pp. 2443–2449.
- [86] M. Yan, Z. Deng, B. He, C. Zou, J. Wu, Z. Zhu, Emotion classification with multichannel physiological signals using hybrid feature and adaptive decision fusion, *Biomed. Signal Process. Control* 71 (2022) 103235.
- [87] A. de Souza Brito, M.B. Vieira, S.M. Villela, H. Tacon, H. de Lima Chaves, H. de Almeida Maia, D.T. Concha, H. Pedrini, Weighted voting of multi-stream convolutional neural networks for video-based action recognition using optical flow rhythms, *J. Vis. Commun. Image Represent.* 77 (2021) 103112.
- [88] J. Gaebel, H.-G. Wu, A. Oeser, M.A. Cypko, M. Stoehr, A. Dietz, T. Neumuth, S. Franke, S. Oeltze-Jafra, Modeling and processing up-to-dateness of patient information in probabilistic therapy decision support, *Artif. Intell. Med.* 104 (2020) 101842.
- [89] J. Chen, Y. Liu, Multimodality data fusion for probabilistic strength estimation of aging materials using Bayesian networks, in: *AIAA Scitech 2020 Forum*, 2020, p. 1653.
- [90] P. Cao, X. Liu, J. Yang, D. Zhao, M. Huang, O. Zaiane, l2, 1-11 regularized non-linear multi-task representation learning based cognitive performance prediction of alzheimers disease, *Pattern Recognit.* 79 (2018) 195–215.
- [91] S. Sharma, P.K. Mandal, A comprehensive report on machine learning-based early detection of alzheimer's disease using multi-modal neuroimaging data, *ACM Comput. Surv.* 55 (2) (2022) 1–44.
- [92] K.M.M. Lopez, M.S.A. Magboo, A. Tallón-Ballesteros, C. Chen, A clinical decision support tool to detect invasive ductal carcinoma in histopathological images using support vector machines, Naïve-Bayes, and K-nearest neighbor classifiers, in: *MLIS*, 2020, pp. 46–53.
- [93] Y. Liu, Z. Gu, T.H. Ko, J. Liu, Identifying key opinion leaders in social media via modality-consistent harmonized discriminant embedding, *IEEE Trans. Cybern.* 50 (2) (2018) 717–728.
- [94] C. Zhang, Z. Yang, X. He, L. Deng, Multimodal intelligence: Representation learning, information fusion, and applications, *IEEE J. Sel. Top. Sign. Proces.* 14 (3) (2020) 478–493.
- [95] F. Anowar, S. Sadaoui, B. Selim, Conceptual and empirical comparison of dimensionality reduction algorithms (pca, kpc, lda, mds, svd, lle, isomap, le, ica, t-sne), *Comp. Sci. Rev.* 40 (2021) 100378.
- [96] S. Zheng, Z. Zhu, Z. Liu, Z. Guo, Y. Liu, Y. Yang, Y. Zhao, Multi-modal graph learning for disease prediction, *IEEE Trans. Med. Imaging* 41 (9) (2022) 2207–2216.
- [97] S. Yan, Y. Xiong, D. Lin, Spatial temporal graph convolutional networks for skeleton-based action recognition, in: *Proceedings of the AAAI Conf. on Artificial Intelligence*, Vol. 32, (1) 2018.
- [98] M. Hügler, G. Kalweit, T. Hügler, J. Boedecker, A dynamic deep neural network for multimodal clinical data analysis, *Explain. AI Healthc. Med.: Build. Cult. Transpar. Accountability* (2021) 79–92.
- [99] A. Elboushaki, R. Hannane, K. Afdel, L. Koutti, Multid-CNN: A multi-dimensional feature learning approach based on deep convolutional networks for gesture recognition in RGB-D image sequences, *Expert Syst. Appl.* 139 (2020) 112829.
- [100] K.M. Rashid, J. Louis, Times-series data augmentation and deep learning for construction equipment activity recognition, *Adv. Eng. Inform.* 42 (2019) 100944.
- [101] N. Bahador, J. Jokelainen, S. Mustola, J. Kortelainen, Multimodal spatio-temporal-spectral fusion for deep learning applications in physiological time series processing: A case study in monitoring the depth of anesthesia, *Inf. Fusion* 73 (2021) 125–143.
- [102] X. Wang, G. Chen, G. Qian, P. Gao, X.-Y. Wei, Y. Wang, Y. Tian, W. Gao, Large-scale multi-modal pre-trained models: A comprehensive survey, 2023, arXiv preprint arXiv:2302.10035.
- [103] G. Ayana, K. Dese, S.-w. Choe, Transfer learning in breast cancer diagnoses via ultrasound imaging, *Cancers* 13 (4) (2021) 738.
- [104] A. de Santana Correia, E.L. Colombini, Attention, please! a survey of neural attention models in deep learning, *Artif. Intell. Rev.* 55 (8) (2022) 6037–6124.
- [105] Z. Niu, G. Zhong, H. Yu, A review on the attention mechanism of deep learning, *Neurocomputing* 452 (2021) 48–62.
- [106] Y. Shi, B. Paige, P. Torr, et al., Variational mixture-of-experts autoencoders for multi-modal deep generative models, *Adv. Neural Inf. Process. Syst.* 32 (2019).
- [107] C. Du, C. Du, H. He, Multimodal deep generative adversarial models for scalable doubly semi-supervised learning, *Inf. Fusion* 68 (2021) 118–130.
- [108] H.R.V. Joze, A. Shaban, M.L. Iuzzolino, K. Koishida, MMTM: Multimodal transfer module for CNN fusion, in: *Proceedings of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2020, pp. 13289–13299.
- [109] Y. Zhang, D. Sidibé, O. Morel, F. Mériaudeau, Deep multimodal fusion for semantic image segmentation: A survey, *Image Vis. Comput.* 105 (2021) 104042.
- [110] R. Carvalho, A.C. Morgado, C. Andrade, T. Nedelcu, A. Carreiro, M.J.M. Vasconcelos, Integrating domain knowledge into deep learning for skin lesion risk prioritization to assist teledermatology referral, *Diagnostics* 12 (1) (2021) 36.
- [111] D. Jin, E. Sergeeva, W.-H. Weng, G. Chauhan, P. Szolovits, Explainable deep learning in healthcare: A methodological survey from an attribution view, *WIREs Mech. Dis.* 14 (3) (2022) e1548.

- [112] R. Sevastjanova, F. Beck, B. Ell, C. Turkay, R. Henkin, M. Butt, D.A. Keim, M. El-Assady, Going beyond visualization: Verbalization as complementary medium to explain machine learning models, in: Workshop on Visualization for AI Explainability At IEEE VIS, 2018.
- [113] K.M. Boehm, P. Khosravi, R. Vanguri, J. Gao, S.P. Shah, Harnessing multimodal data integration to advance precision oncology, *Nat. Rev. Cancer* 22 (2) (2022) 114–126.
- [114] T. Shaik, X. Tao, Y. Li, C. Dann, J. McDonald, P. Redmond, L. Galligan, A review of the trends and challenges in adopting natural language processing methods for education feedback analysis, *IEEE Access* (2022).
- [115] Z. Zeng, Y. Deng, X. Li, T. Naumann, Y. Luo, Natural language processing for EHR-based computational phenotyping, *IEEE/ACM Trans. Comput. Biol. Bioinform.* 16 (1) (2018) 139–153.
- [116] P. Bhatia, B. Celikkaya, M. Khalilia, S. Senthivel, Comprehend medical: A named entity recognition and relationship extraction web service, in: 2019 18th IEEE Int'l. Conf. on Machine Learning and Applications, (ICMLA), IEEE, 2019, pp. 1844–1851.
- [117] D. Demner-Fushman, N. Elhadad, C. Friedman, Natural language processing for health-related texts, in: *Biomedical Informatics: Computer Applications in Health Care and Biomedicine*, Springer, 2021, pp. 241–272.
- [118] E. Petrova, P. Pauwels, K. Svidt, R.L. Jensen, Towards data-driven sustainable design: decision support based on knowledge discovery in disparate building data, *Archit. Eng. Des. Manag.* 15 (5) (2019) 334–356.
- [119] T. Pham, X. Tao, J. Zhang, J. Yong, Constructing a knowledge-based heterogeneous information graph for medical health status classification, *Health Inf. Syst. Syst.* 8 (2020) 1–14.
- [120] M. Tang, P. Gandhi, M.A. Kabir, C. Zou, J. Blakey, X. Luo, Progress notes classification and keyword extraction using attention-based deep learning models with BERT, 2019, arXiv preprint arXiv:1910.05786.
- [121] N. Chintalapudi, G. Battineni, M. Di Canio, G.G. Sagaró, F. Amenta, Text mining with sentiment analysis on seafarers' medical documents, *Int. J. Inf. Manag. Data Insights* 1 (1) (2021) 100005.
- [122] S. Bozkurt, E. Alkim, I. Banerjee, D.L. Rubin, Automated detection of measurements and their descriptors in radiology reports using a hybrid natural language processing algorithm, *J. Digit. Imaging* 32 (2019) 544–553.
- [123] X. Pei, K. Zuo, Y. Li, Z. Pang, A review of the application of multi-modal deep learning in medicine: Bibliometrics and future directions, *Int. J. Comput. Intell. Syst.* 16 (1) (2023) 1–20.
- [124] L. Wang, M. Rastegar-Mojarad, Z. Ji, S. Liu, K. Liu, S. Moon, F. Shen, Y. Wang, L. Yao, J.M. Davis III, et al., Detecting pharmacovigilance signals combining electronic medical records with spontaneous reports: A case study of conventional disease-modifying antirheumatic drugs for rheumatoid arthritis, *Front. Pharmacol.* 9 (2018) 875.
- [125] M.F. Guinazú, V. Cortés, C.F. Ibáñez, J.D. Velásquez, Employing online social networks in precision-medicine approach using information fusion predictive model to improve substance use surveillance: A lesson from Twitter and marijuana consumption, *Inf. Fusion* 55 (2020) 150–163.
- [126] A. Choudhury, O. Asan, et al., Role of artificial intelligence in patient safety outcomes: systematic literature review, *JMIR Med. Inform.* 8 (7) (2020) e18599.
- [127] A. Le Glaz, Y. Haralambous, D.-H. Kim-Dufor, P. Lenca, R. Billot, T.C. Ryan, J. Marsh, J. Devylder, M. Walter, S. Berrouiguet, et al., Machine learning and natural language processing in mental health: systematic review, *J. Med. Internet Res.* 23 (5) (2021) e15708.
- [128] J. Lipkova, R.J. Chen, B. Chen, M.Y. Lu, M. Barbieri, D. Shao, A.J. Vaidya, C. Chen, L. Zhuang, D.F. Williamson, et al., Artificial intelligence for multimodal data integration in oncology, *Cancer Cell* 40 (10) (2022) 1095–1110.
- [129] B.N. Hiremath, M.M. Patil, Enhancing optimized personalized therapy in clinical decision support system using natural language processing, *J. King Saud Univ.-Comput. Inf. Sci.* 34 (6) (2022) 2840–2848.
- [130] I. Spasic, G. Nenadic, et al., Clinical text data in machine learning: systematic review, *JMIR Med. Inform.* 8 (3) (2020) e17984.
- [131] I. Perez-Pozuelo, B. Zhai, J. Palotti, R. Mall, M. Aupetit, J.M. Garcia-Gomez, S. Taheri, Y. Guan, L. Fernandez-Luque, The future of sleep health: A data-driven revolution in sleep science and medicine, *NPJ Digit. Med.* 3 (1) (2020) 42.
- [132] J. Sun, H. Shi, J. Zhu, B. Song, Y. Tao, S. Tan, Self-attention-based multi-block regression fusion neural network for quality-related process monitoring, *J. Taiwan Inst. Chem. Eng.* 133 (2022) 104140.
- [133] F.A. Reegu, H. Abas, Y. Gulzar, Q. Xin, A.A. Alwan, A. Jabbari, R.G. Sonkamble, R.A. Dziyuddin, Blockchain-based framework for interoperable electronic health records for an improved healthcare system, *Sustainability* 15 (8) (2023) 6337.
- [134] C. Lyketos, S. Roberts, E.K. Swift, A. Quina, G. Moon, I. Kremer, P. Tariot, H. Fillit, D. Bovenkamp, P. Zandi, et al., Standardizing electronic health record data on AD/ADRD to accelerate health equity in prevention, detection, and treatment, *J. Prev. Alzheimers Dis.* 9 (3) (2022) 556–560.
- [135] G. Diraco, G. Rescio, P. Siciliano, A. Leone, Review on human action recognition in smart living: Multimodality, real-time processing, interoperability, resource-constrained processing, and sensing technology, 2023.
- [136] C. Mwangi, C. Mukanya, C. Maghanga, Assessing the interoperability of mlab and ushauri mhealth systems to enhance care for HIV/AIDS patients in Kenya, *J. Intellect. Prop. Inf. Technol. Law (JIPIIT)* 2 (1) (2022) 83–116.
- [137] M. Kor, I. Yitmen, S. Alizadehsalehi, An investigation for integration of deep learning and digital twins towards construction 4.0, *Smart Sustain. Built Environ.* 12 (3) (2023) 461–487.
- [138] X. Tao, J.D. Velásquez, Multi-source information fusion for smart health with artificial intelligence, *Inf. Fusion* 83–84 (2022) 93–95.
- [139] M. Paul, L. Maglaras, M.A. Ferrag, I. AlMamani, Digitization of healthcare sector: A study on privacy and security concerns, *ICT Express* (2023).
- [140] I. Yasser, A.T. Khalil, M.A. Mohamed, A.S. Samra, F. Khalifa, A robust chaos-based technique for medical image encryption, *IEEE Access* 10 (2021) 244–257.
- [141] P.B. Regade, A.A. Patil, S.S. Koli, R.B. Gokavi, M. Bhandigare, Survey on secure file storage on cloud using hybrid cryptography, *Int. Res. J. Modern. Eng. Technol. Sci.* 4 (06) (2022).
- [142] Y. Al-Issa, M.A. Ottom, A. Tamrawi, Ehealth cloud security challenges: A survey, *J. Healthc. Eng.* 2019 (2019).
- [143] M. Mohammed, S. Desyansah, S. Al-Zubaidi, E. Yusuf, An internet of things-based smart homes and healthcare monitoring and management system, *J. Phys.: Conf. Ser.* 1450 (1) (2020) 012079.
- [144] J.J. Hathaliya, S. Tanwar, P. Sharma, Adversarial learning techniques for security and privacy preservation: A comprehensive review, *Secur. Privacy* 5 (3) (2022) e209.
- [145] N.N. Neto, S. Madnick, A.M.G. de Paula, N. Malara Borges, A case study of the capital one data breach: why didn't compliance requirements help prevent it? *J. Inf. Syst. Secur.* 17 (1) (2021).
- [146] R. Kumar, R. Goyal, On cloud security requirements, threats, vulnerabilities and countermeasures: A survey, *Comp. Sci. Rev.* 33 (2019) 1–48.
- [147] V.R. Kebande, N.M. Karie, R.A. Ikuesan, Real-time monitoring as a supplementary security component of vigilantism in modern network environments, *Int. J. Inf. Technol.* 13 (2021) 5–17.
- [148] R. Bokade, A. Navato, R. Ouyang, X. Jin, C.-A. Chou, S. Ostadabbas, A.V. Mueller, A cross-disciplinary comparison of multimodal data fusion approaches and applications: Accelerating learning through trans-disciplinary information sharing, *Expert Syst. Appl.* 165 (2021) 113885.
- [149] A.M. Flores, F. Demsas, N.J. Leeper, E.G. Ross, Leveraging machine learning and artificial intelligence to improve peripheral artery disease detection, treatment, and outcomes, *Circ. Res.* 128 (12) (2021) 1833–1850.
- [150] M. Swathy, K. Saruladha, A comparative study of classification and prediction of cardio-vascular diseases (CVD) using machine learning and deep learning techniques, *ICT Express* 8 (1) (2022) 109–116.
- [151] I. Banerjee, Y. Ling, M.C. Chen, S.A. Hasan, C.P. Langlotz, N. Moradzadeh, B. Chapman, T. Amrhein, D. Mong, D.L. Rubin, et al., Comparative effectiveness of convolutional neural network (CNN) and recurrent neural network (RNN) architectures for radiology text report classification, *Artif. Intell. Med.* 97 (2019) 79–88.
- [152] A. Coronato, M. Naeem, G. De Pietro, G. Paragliola, Reinforcement learning for intelligent healthcare applications: A survey, *Artif. Intell. Med.* 109 (2020) 101964.
- [153] D. Wang, J.D. Weisz, M. Muller, P. Ram, W. Geyer, C. Dugan, Y. Tausczik, H. Samulowitz, A. Gray, Human-AI collaboration in data science: Exploring data scientists' perceptions of automated AI, *Proc. ACM Hum.-Comput. Interact.* 3 (CSCW) (2019) 1–24.
- [154] I.H. Sarker, Data science and analytics: An overview from data-driven smart computing, decision-making and applications perspective, *SN Comput. Sci.* 2 (5) (2021) 377.
- [155] S. Steyaert, M. Pizurica, D. Nagaraj, P. Khandelwal, T. Hernandez-Boussard, A.J. Gentles, O. Gevaert, Multimodal data fusion for cancer biomarker discovery with deep learning, *Nat. Mach. Intell.* (2023) 1–12.
- [156] D. Wang, L. Wang, Z. Zhang, D. Wang, H. Zhu, Y. Gao, X. Fan, F. Tian, "Brilliant AI doctor" in rural clinics: Challenges in AI-powered clinical decision support system deployment, in: Proceedings of the 2021 CHI Conf. on Human Factors in Computing Systems, 2021, pp. 1–18.
- [157] E. Nazari, M.H. Shahriari, H. Tabesh, BigData analysis in healthcare: apache hadoop, apache spark and apache flink, *Front. Health Inform.* 8 (1) (2019) 14.
- [158] A. Kaur, P. Singh, A. Nayyar, Fog computing: Building a road to IoT with fog analytics, *Fog Data Anal. IoT Appl.: Next Generation Process Model State Art Technol.* (2020) 59–78.
- [159] R. Dwivedi, D. Mehrotra, S. Chandra, Potential of internet of medical things (IoMT) applications in building a smart healthcare system: A systematic review, *J. Oral Biol. Craniofac. Res.* 12 (2) (2022) 302–318.
- [160] Q. Qi, F. Tao, T. Hu, N. Anwer, A. Liu, Y. Wei, L. Wang, A. Nee, Enabling technologies and tools for digital twin, *J. Manuf. Syst.* 58 (2021) 3–21.
- [161] P.K.R. Maddikunta, Q.-V. Pham, B. Prabadevi, N. Deepa, K. Dev, T.R. Gadekallu, R. Ruby, M. Liyanage, Industry 5.0: A survey on enabling technologies and potential applications, *J. Ind. Inf. Integr.* 26 (2022) 100257.
- [162] A. Vakil, J. Liu, P. Zulch, E. Blasch, R. Ewing, J. Li, A survey of multimodal sensor fusion for passive RF and EO information integration, *IEEE Aerosp. Electron. Syst. Mag.* 36 (7) (2021) 44–61.
- [163] L. You, M. Danaf, F. Zhao, J. Guan, C.L. Azevedo, B. Atasoy, M. Ben-Akiva, A federated platform enabling a systematic collaboration among devices, data and functions for smart mobility, *IEEE Trans. Intell. Transp. Syst.* 24 (4) (2023) 4060–4074.

- [164] R. Dabiz, S.K. Poon, A. Ritchie, R. Burke, J. Penm, Usability evaluation of an integrated electronic medication management system implemented in an oncology setting using the unified theory of acceptance and use of technology, *BMC Med. Inform. Decis. Mak.* 21 (1) (2021) 1–11.
- [165] B.N. Limketkai, K. Mauldin, N. Manitus, L. Jalilian, B.R. Salonen, The age of artificial intelligence: use of digital technology in clinical nutrition, *Curr. Surg. Rep.* 9 (7) (2021) 20.
- [166] X. Chen, H. Xie, Z. Li, G. Cheng, M. Leng, F.L. Wang, Information fusion and artificial intelligence for smart healthcare: a bibliometric study, *Inf. Process. Manage.* 60 (1) (2023) 103113.
- [167] V.M. O'Hara, S.V. Johnston, N.T. Browne, The paediatric weight management office visit via telemedicine: pre-to post-COVID-19 pandemic, *Pediatr. Obes.* 15 (8) (2020) e12694.
- [168] A. Holzinger, M. Dehmer, F. Emmert-Streib, R. Cucchiara, I. Augenstein, J. Del Ser, W. Samek, I. Jurisica, N. Díaz-Rodríguez, Information fusion as an integrative cross-cutting enabler to achieve robust, explainable, and trustworthy medical artificial intelligence, *Inf. Fusion* 79 (2022) 263–278.
- [169] L. Baum, M. Johns, M. Poikela, R. Möller, B. Ananthasubramaniam, F. Prasser, Data integration and analysis for circadian medicine, *Acta Physiol.* 237 (4) (2023) e13951.
- [170] S.M. van Rooden, O. Aspevall, E. Carrara, S. Gubbels, A. Johansson, J.-C. Lucet, S. Mookerjee, Z.R. Palacios-Baena, E. Presterl, E. Tacconelli, et al., Governance aspects of large-scale implementation of automated surveillance of healthcare-associated infections, *Clin. Microbiol. Infect.* 27 (2021) S20–S28.
- [171] C. Thapa, S. Camtepe, Precision health data: Requirements, challenges and existing techniques for data security and privacy, *Comput. Biol. Med.* 129 (2021) 104130.
- [172] N. Gaw, S. Yousefi, M.R. Gahrooei, Multimodal data fusion for systems improvement: A review, *IJSE Trans.* 54 (11) (2022) 1098–1116.
- [173] J. Mökander, J. Morley, M. Taddeo, L. Floridi, Ethics-based auditing of automated decision-making systems: Nature, scope, and limitations, *Sci. Eng. Ethics* 27 (4) (2021) 44.
- [174] L. Belgodère, D.P. Bertrand, M.C. Jaulent, V. Rabeharisoa, W. Janssens, V. Rollason, J. Barbot, J.P. Vernant, W.O. Gonin, P. Maison, et al., Patient and public involvement in the benefit–risk assessment and decision concerning health products: position of the scientific advisory board of the french national agency for medicines and health products safety (ANSM), *BMJ Glob. Health* 8 (5) (2023) e011966.
- [175] S. Ali, T. Abuhmed, S. El-Sappagh, K. Muhammad, J.M. Alonso-Moral, R. Confalonieri, R. Guidotti, J. Del Ser, N. Díaz-Rodríguez, F. Herrera, Explainable artificial intelligence (XAI): What we know and what is left to attain trustworthy artificial intelligence, *Inf. Fusion* (2023) 101805.
- [176] N. Rostamzadeh, S.S. Abdullah, K. Sedig, Visual analytics for electronic health records: A review, in: *Informatics*, Vol. 8, (1) MDPI, 2021, p. 12.
- [177] T. Höllt, A. Vilanova, N. Pezzotti, B.P. Lelieveldt, H. Hauser, Focus+ context exploration of hierarchical embeddings, in: *Computer Graphics Forum*, Vol. 38, (3) Wiley Online Library, 2019, pp. 569–579.
- [178] A.B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bannetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, et al., Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI, *Inf. Fusion* 58 (2020) 82–115.
- [179] A. Holzinger, M. Dehmer, F. Emmert-Streib, R. Cucchiara, I. Augenstein, J. Del Ser, W. Samek, I. Jurisica, N. Díaz-Rodríguez, Information fusion as an integrative cross-cutting enabler to achieve robust, explainable, and trustworthy medical artificial intelligence, *Inf. Fusion* 79 (2022) 263–278.
- [180] Y. Mao, D. Wang, M. Muller, K.R. Varshney, I. Baldini, C. Dugan, A. Mojsilović, How data ScientistsWork together with domain experts in scientific collaborations, *Proc. ACM Hum.-Comput. Interact.* 3 (GROUP) (2019) 1–23.
- [181] L. Müller, A. Srinivasan, S.R. Abeles, A. Rajagopal, F.J. Torriani, E. Aronoff-Spencer, A risk-based clinical decision support system for patient-specific antimicrobial therapy (iBiogram): Design and retrospective analysis, *J. Med. Internet Res.* 23 (12) (2021) e23571.
- [182] T. Pham, X. Tao, J. Zhang, J. Yong, Y. Li, H. Xie, Graph-based multi-label disease prediction model learning from medical data and domain knowledge, *Knowl.-Based Syst.* 235 (2022) 107662.
- [183] G. Collatuzzo, P. Boffetta, Application of P4 (predictive, preventive, personalized, participatory) approach to occupational medicine, *Med. Lavoro* 113 (1) (2022).
- [184] R.B. Ruiz, J.D. Velásquez, Artificial intelligence for the future of medicine, in: *Artificial Intelligence and Machine Learning for Healthcare: Vol. 2: Emerging Methodologies and Trends*, Springer, 2022, pp. 1–28.
- [185] A. Esteva, A. Robicquet, B. Ramsundar, V. Kuleshov, M. DePristo, K. Chou, C. Cui, G. Corrado, S. Thrun, J. Dean, A guide to deep learning in healthcare, *Nature Med.* 25 (1) (2019) 24–29.
- [186] M.A. Myszczyńska, P.N. Ojames, A.M. Lacoste, D. Neil, A. Saffari, R. Mead, G.M. Hautbergue, J.D. Holbrook, L. Ferraiuolo, Applications of machine learning to diagnosis and treatment of neurodegenerative diseases, *Nature Rev. Neurol.* 16 (8) (2020) 440–456.
- [187] Y.-D. Zhang, Z. Dong, S.-H. Wang, X. Yu, X. Yao, Q. Zhou, H. Hu, M. Li, C. Jiménez-Mesa, J. Ramirez, et al., Advances in multimodal data fusion in neuroimaging: overview, challenges, and novel orientation, *Inf. Fusion* 64 (2020) 149–187.
- [188] X.-a. Bi, X. Hu, Y. Xie, H. Wu, A novel CERNNE approach for predicting parkinson's disease-associated genes and brain regions based on multimodal imaging genetics data, *Med. Image Anal.* 67 (2021) 101830.
- [189] L.A. Vale-Silva, K. Rohr, Long-term cancer survival prediction using multimodal deep learning, *Sci. Rep.* 11 (1) (2021) 13505.
- [190] R. Nabbut, M. Kuchenbuch, Impact of predictive, preventive and precision medicine strategies in epilepsy, *Nature Rev. Neurol.* 16 (12) (2020) 674–688.
- [191] T. Shaik, X. Tao, N. Higgins, H. Xie, R. Gururajan, X. Zhou, AI enabled RPM for mental health facility, in: *Proceedings of the 1st ACM Workshop on Mobile and Wireless Sensing for Smart Healthcare*, 2022, pp. 26–32.
- [192] M.C. Liefwaard, E.H. Lips, J. Wesseling, N.M. Hylton, B. Lou, T. Mansi, L. Pusztai, The way of the future: personalizing treatment plans through technology, *Am. Soc. Clin. Oncol. Educ. Book* 41 (2021) 12–23.
- [193] T. Shaik, X. Tao, N. Higgins, R. Gururajan, Y. Li, X. Zhou, U.R. Acharya, Fedstack: Personalized activity monitoring using stacked federated learning, *Knowl.-Based Syst.* 257 (2022) 109929.
- [194] A.A.T. Naqvi, K. Fatima, T. Mohammad, U. Fatima, I.K. Singh, A. Singh, S.M. Atif, G. Hariprasad, G.M. Hasan, M.I. Hassan, Insights into SARS-CoV-2 genome, structure, evolution, pathogenesis and therapies: Structural genomics approach, *Biochim. Biophys. Acta (BBA)-Mol. Basis Dis.* 1866 (10) (2020) 165878.
- [195] D. Horgan, T. Čufer, F. Gatto, I. Lugowska, D. Verbanac, Â. Carvalho, J.A. Lal, M. Kozaric, S. Toomey, H.Y. Ivanov, et al., Accelerating the development and validation of liquid biopsy for early cancer screening and treatment tailoring, in: *Healthcare*, Vol. 10, (9) MDPI, 2022, p. 1714.
- [196] H. Cai, Z. Qu, Z. Li, Y. Zhang, X. Hu, B. Hu, Feature-level fusion approaches based on multimodal EEG data for depression recognition, *Inf. Fusion* 59 (2020) 127–138.
- [197] J. Mateo, L. Steuten, P. Aftimos, F. André, M. Davies, E. Garralda, J. Geissler, D. Husereau, I. Martinez-Lopez, N. Normanno, et al., Delivering precision oncology to patients with cancer, *Nature Med.* 28 (4) (2022) 658–665.
- [198] G. Aceto, V. Persico, A. Pescapé, Industry 4.0 and health: Internet of things, big data, and cloud computing for healthcare 4.0, *J. Ind. Inform. Integr.* 18 (2020) 100129.
- [199] K.M. Boehm, E.A. Aherne, L. Ellenson, I. Nikolovski, M. Alghamdi, I. Vázquez-García, D. Zamarin, K. Long Roche, Y. Liu, D. Patel, et al., Multimodal data integration using machine learning improves risk stratification of high-grade serous ovarian cancer, *Nature Cancer* 3 (6) (2022) 723–733.
- [200] P. Carayon, A. Wooldridge, P. Hoonakker, A.S. Hundt, M.M. Kelly, SEIPS 3.0: Human-centered design of the patient journey for patient safety, *Appl. Ergon.* 84 (2020) 103033.
- [201] H. Dhayne, R. Haque, R. Kilany, Y. Taher, In search of big medical data integration solutions—a comprehensive survey, *IEEE Access* 7 (2019) 91265–91290.
- [202] A. El Saddik, F. Laamarti, M. Alja'afreh, The potential of digital twins, *IEEE Instrum. Meas. Mag.* 24 (3) (2021) 36–41.
- [203] K. Walker, J. Yates, T. Denning, B. Völm, J. Tomlin, C. Griffiths, Quality of life, wellbeing, recovery, and progress for older forensic mental health patients: A qualitative investigation based on the perspectives of patients and staff, *Int. J. Qual. Stud. Health Well-being* 18 (1) (2023) 2202978.

7.2 Summary

This chapter consolidates the comprehensive insights derived from an extensive review of multimodal information fusion in the context of smart healthcare. It emphasizes the critical contribution of this technology in realizing the aspirations of predictive, preventive, personalized, and participatory healthcare. The nuanced integration of heterogeneous data streams through advanced fusion techniques marks a paradigm shift towards intelligent healthcare ecosystems. The chapter advocates for sustained innovation and interdisciplinary research to navigate the complexities of multimodal data, envisioning a future where integrated, data-driven insights form the cornerstone of proactive and patient-centric healthcare.

CHAPTER 8: PAPER 7 - GRAPH-ENABLED REINFORCEMENT LEARNING FOR TIME SERIES FORECASTING WITH ADAPTIVE INTELLIGENCE

8.1 Introduction

This chapter delves into the innovative GraphRL framework, which marries Graph Neural Networks (GNNs) with Reinforcement Learning (RL) to pioneer time series forecasting in dynamic environments. By leveraging the structural and temporal insights provided by GNNs, and the adaptive decision-making prowess of RL, the framework sets a new benchmark in predictive analytics. The introduction outlines the motivation behind GraphRL, its architectural nuances, and its application across diverse domains such as healthcare, traffic, and weather forecasting, setting the stage for a detailed exploration of its capabilities and contributions to smart monitoring and predictive analytics.

Graph-enabled Reinforcement Learning for Time Series Forecasting with Adaptive Intelligence

Thanveer Shaik, Xiaohui Tao, Haoran Xie, Lin Li, Jianming Yong, and Yuefeng Li

Abstract—Reinforcement learning is well known for its ability to model sequential tasks and learn latent data patterns adaptively. Deep learning models have been widely explored and adopted in regression and classification tasks. However, deep learning has its limitations such as the assumption of equally spaced and ordered data, and the lack of ability to incorporate graph structure in terms of time-series prediction. Graphical neural network (GNN) has the ability to overcome these challenges and capture the temporal dependencies in time-series data. In this study, we propose a novel approach for predicting time-series data using GNN and monitoring with Reinforcement Learning (RL). GNNs are able to explicitly incorporate the graph structure of the data into the model, allowing them to capture temporal dependencies in a more natural way. This approach allows for more accurate predictions in complex temporal structures, such as those found in healthcare, traffic and weather forecasting. We also fine-tune our GraphRL model using a Bayesian optimisation technique to further improve performance. The proposed framework outperforms the baseline models in time-series forecasting and monitoring. The contributions of this study include the introduction of a novel GraphRL framework for time-series prediction and the demonstration of the effectiveness of GNNs in comparison to traditional deep learning models such as RNNs and LSTMs. Overall, this study demonstrates the potential of GraphRL in providing accurate and efficient predictions in dynamic RL environments.

Index Terms—Graph Neural Networks, Reinforcement Learning, Intelligent Monitoring, Bayesian Optimization

I. INTRODUCTION

The emergence of Machine Learning (ML) in healthcare signifies a paradigm shift towards automating clinician tasks and augmenting patient care capabilities [1]. Amidst the evolving ML landscape, Federated Learning has gained traction for preserving data privacy while constructing sophisticated server models [2]. Reinforcement Learning (RL), another ML strategy, has demonstrated substantial improvements in prediction performance and decision-making tasks [3], [4]. RL's application is particularly noteworthy in controlling autonomous systems, such as robots and drones, training them to make optimal decisions in real-time based on environmental sensor data.

Thanveer Shaik and Xiaohui Tao are with the School of Mathematics, Physics & Computing, University of Southern Queensland, Toowoomba, Queensland, Australia (e-mail: Thanveer.Shaik@usq.edu.au, Xiaohui.Tao@usq.edu.au).

Haoran Xie is with the Department of Computing and Decision Sciences, Lingnan University, Tuen Mun, Hong Kong (e-mail: hrxie@ln.edu.hk)

Lin Li is with the School of Computer and Artificial Intelligence, Wuhan University of Technology, China (e-mail: cathylilin@whut.edu.cn)

Jianming Yong is with the School of Business at the University of Southern Queensland, Queensland, Australia (e-mail: Jianming.Yong@usq.edu.au)

Yuefeng Li is with the School of Computer Science, Queensland University of Technology, Brisbane, Australia (e-mail: y2.li@qut.edu.au).

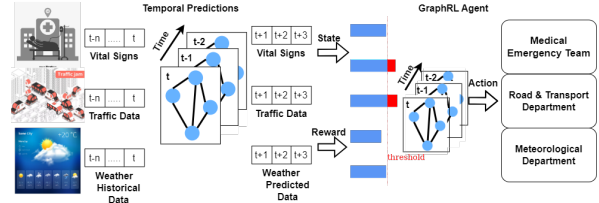


Fig. 1: Graphical Abstract

In various sectors, including healthcare, traffic, and weather forecasting, Early Warning Systems (EWS) play a pivotal role. They analyze real-time monitoring data and issue alerts for potential issues, facilitating proactive responses. RL-based EWS can adapt over time, refining their predictions and supporting clinical decision-making. This adaptability has proven effective in applications like predicting hospital readmissions and sepsis detection.

Time-series data modeling, vital in monitoring and predicting future states, has seen advancements with deep learning models like Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks [5]. These models are adept at capturing temporal dependencies, yet face limitations in handling irregularly structured data and integrating complex graph structures. This study introduces the GraphRL framework, an innovative amalgamation of RL and Temporal Graphical Convolutional Networks (T-GCN), aiming to surpass the constraints of traditional deep learning models in time-series prediction. GraphRL's design facilitates handling complex temporal structures and incorporates additional information such as node and edge attributes, as depicted in Fig. 1. The GraphRL framework's core contributions include:

- A versatile framework capable of providing early warnings and monitoring in complex settings.
- A customizable RL environment designed for effective forecasting in dynamic domains like healthcare and traffic systems.
- A novel approach to virtual monitoring of predicted states in RL, enhancing decision-making and intervention capabilities.

Our comparative analysis with state-of-the-art models across various datasets showcases GraphRL's superior performance, underscoring its potential as a versatile solution for time-series prediction challenges.

The paper is organized as follows: Section II reviews existing literature on self-learning systems and prediction tasks. Section III outlines the research problem. The proposed

GraphRL framework and its algorithm are detailed in Section IV. Section V describes the datasets and baseline models used for evaluation. The performance of the predictive RL environment and agent is compared with baseline models in Section VI. Section VII discusses the fine-tuning of the framework’s hyperparameters using Bayesian optimization. The paper concludes in Section VIII.

II. RELATED WORKS

A. Self Learning Systems

Self-learning systems, particularly those utilizing Reinforcement Learning (RL), have seen significant advancements in various applications. For instance, Shin et al. [6] introduced a dual-agent framework in mobile health, effectively demonstrating user modeling and behavior intervention strategies. This work underscores the potential of RL in personalizing user experiences, a concept that aligns with our GraphRL framework’s goal of adaptive learning in dynamic environments. Similarly, Taylor et al. [7] applied RL in modeling maladaptive eating behaviors, further showcasing RL’s versatility in behavior prediction and modification. Chen et al. [8] developed the MIRROR framework, emphasizing the rapid learning capabilities of RL in human behavior modeling. These advancements set a precedent for our work in complex sequential decision-making tasks. Zhou et al.’s [9] CalFit app and Li et al.’s [10] method in autonomous driving highlight RL’s efficacy in personalized goal setting and complex urban scenario navigation, respectively, which are foundational to our GraphRL framework’s approach in handling dynamic and intricate patterns in data.

B. Early Detection of Patient Deterioration

In the healthcare domain, early detection of patient deterioration is vital. Traditional vital signs monitoring, as discussed by Asiimwe et al. [11] and evaluated by Scully et al. [12], Baig et al. [13], and others, has laid the groundwork for our study. These works highlight the importance of continuous monitoring and early warning systems (EWS), which are integral to GraphRL’s objective. The limitations in existing methods, such as the need for manual calculations and the inability to handle large, unstructured data effectively, are addressed in our framework through the integration of GNNs, which can process complex temporal data more efficiently.

C. Vital Signs Prediction

Vital signs prediction has been explored through various machine learning models. Alghatani et al. [14] and Youssef et al. [15] demonstrated the use of traditional machine learning in mortality prediction and vital signs forecasting, respectively. Harerimana et al.’s [16] work with multi-head transformers and Xie et al.’s [17] DeepVS model highlight the potential of deep learning in this domain. However, these methods often assume equally spaced and ordered data and lack the ability to incorporate complex graph structures, limitations our GraphRL framework aims to overcome.

D. Temporal Graphical Convolutional Networks (T-GCN)

The integration of T-GCN within our GraphRL framework is pivotal. T-GCNs, known for their ability to capture temporal dependencies and complex relationships in graph-structured data, offer significant enhancements in processing time-series data [18]. This technology addresses limitations in traditional deep learning models by effectively managing irregular time intervals and integrating additional contextual information (such as node and edge attributes) for richer data representation [19]. The inclusion of T-GCN in GraphRL allows for a more nuanced understanding and prediction of dynamic systems [20], making it highly suitable for applications in healthcare monitoring, traffic forecasting, and weather prediction. The capability of T-GCNs to handle non-linear and complex temporal patterns [21] aligns with the core objectives of GraphRL, pushing the boundaries of current self-learning systems in real-world scenarios.

In summary, while existing works in self-learning systems, patient deterioration detection, and vital signs prediction have laid a strong foundation, our GraphRL framework aims to address their limitations by introducing a novel approach that combines the strengths of GNNs and RL. This approach allows for a more sophisticated handling of temporal dependencies and real-time monitoring, which is crucial in dynamic environments such as healthcare, traffic management, and weather forecasting.

The motivation behind the use of RL in our framework primarily arises from the need to tackle the challenges of multi-step time series prediction, where traditional supervised learning approaches may encounter limitations. Although supervised learning methods like GNN+Bert and GNN+TCN are indeed common and effective for time series forecasting, RL offers a unique advantage in dealing with situations where errors can accumulate over time, especially in dynamic environments. RL enables our predictive GraphRL Environment not only to forecast future states but also to actively influence decision-making, a capability particularly valuable in applications such as healthcare monitoring and the gaming industry.

III. RESEARCH PROBLEM

The research problem addresses deep learning challenges in predicting future states of a complex and dynamic Reinforcement Learning (RL) environment and adaptively learning latent behavior patterns of data.

Definition 1 (Vital Parameters and Time-Series Forecasting): In the context of our framework, we consider a set V of n vital parameters, denoted as $V_t = \{v^1, v^2, \dots, v^n\}$, which represent continuous time-series data reflecting the health status of a subject S . These vital parameters are dynamic and change over time, providing valuable insights into the subject’s well-being. To facilitate time-series forecasting, we segment these continuous vital parameters into time windows, denoted as T , which encompasses data points from the past (V_{t-2}, V_{t-1}, V_t) and extends into the future ($V_{t+1}, V_{t+2}, \dots, V_{t+n}$). Non-linear models are trained on historical data within the time windows $\{V_{t-2}, V_{t-1}, V_t\}$ to predict future values $\{V_{t+1}, V_{t+2}, \dots, V_{t+n}\}$.

Definition 2 (Learning Agents and Markov Decision Process): Following the training phase, subject S is associated with a group of learning agents that operate based on the principles of the Markov Decision Process (MDP). This MDP is a 5-tuple denoted as $M = (S, A, P, R, \gamma)$, and it forms the foundation for continuous monitoring and pattern learning of the vital parameters $V_t = \{v^1, v^2, \dots, v^n\}$. Here's a breakdown of each component:

- S represents a finite state space, where $s_t \in S$ signifies the state of an agent at a specific time t ,
- A is the set of actions available to each agent, and $a_t \in A$ represents the action taken by the agent at time t ,
- P is a Markovian transition function $P(s, a, s')$ that quantifies the probability of the agent transitioning from state s to state s' while executing action a ,
- R is a reward function $R : S \times A \rightarrow \mathbb{R}$ that provides an immediate reward $R(s, a)$ for the action a performed in state s ,
- γ is a discount factor, ranging between 0 and 1, which emphasizes immediate rewards over future rewards.

$$R(s_t, a_t) = \sum_{t=0}^{\infty} \gamma^t r_t, \quad (1)$$

This equation returns the immediate reward $R(s, a)$ for the action taken in state s , as defined in Eq. 1.

IV. GRAPHRL FRAMEWORK

To address the research problem, a novel graphical neural network (GNN) enabled reinforcement learning (RL) framework is proposed. In the graph-enabled RL framework (GraphRL), two GNNs are deployed: one for forecasting time-series data and another for Q-function approximation as shown in Fig. 2. The proposed framework is demonstrated in Fig. 2 in which the interaction between an environment and an AI agent is illustrated. As discussed in the research problem, finite MDP is adopted to formulate the process of modelling current and past states.

A. Predictive GraphRL Environment

The primary objective of the proposed study is to learn from the past and current states of a dynamic environment and predict the future states of the complex environment. To achieve this objective, we propose a predictive monitoring environment which is responsible for defining the observation space with state $s_t^i \in S$ where $i = 0, 1, 2, \dots, n$, action space with actions $a_t^j \in A$ where $j = 1, 2, 3, \dots, m$, and rewards R for each action taken by the agent as it transitions from a state s_t to s_{t+1} in a real-world scenario. For example, consider a subject in a dynamic environment whose current state is denoted as $V_t = v^1, v^2, \dots, v^n$ at time t . Similarly, the subject holds historical data of their state at times $t-1, t-2, t-3, \dots, t-n$. In traditional reinforcement learning formulations, the monitoring environment is a static entity that cannot forecast future states, which might affect the subjects in the environment.

1) *T-GCN Forecast:* Forecasting the future states of a subject before a few time steps can revolutionise the most dynamic industries such as gaming, healthcare, and so on by identifying the deteriorating state of the subject in the environment. To predict the future state in a reinforcement learning environment, a temporal graph convolutional network (T-GCN) is adopted. The graphical network is trained with past and current states at their timestamps in a supervised approach as shown in Eq. 3. The training process also includes the features leading to those states.

$$y = f(b + \sum_{i=1}^n v_i \cdot w_i) \quad (2)$$

$$y(v) = \sum_{i=1}^n \text{Activation1}(b + w_i v_i) \quad (3)$$

$$y(v) = \text{Activation1}\left(\frac{e^{y_i}}{\sum_{j=1}^k e^{y_j}}\right)$$

Eq. 2 describes a basic neural network neuron, with y as the output, f as the activation function, b as the bias, v_i as the input features, and w_i as the weights; while Eq. 3 involves an 'Activation1' function, computing a weighted sum of inputs and normalizing these outputs, possibly into a probability distribution, akin to a softmax function.

2) *Static Spatial-temporal modelling:* A two-layered graphical network is adopted for Spatial-temporal modelling, a spatial modelling layer is based on a graphical convolutional network (GCN), and a temporal layer based on recurrent neural networks (RNN) is configured. The spatial layer is responsible to capture spatial features among nodes which are input states $s_t^i \in S$. This can be achieved by constructing Fourier transform filter and it acts on the graph nodes and its first-order neighbourhood. In this study, a static graph temporal signal is adopted in which the node positions in the graph remain the same and the label information is dynamic. The spatial layer is to set the static graph with nodes as input states $s_t^i \in S$. The two-layered GCN model is defined in Eq. 4.

$$f(X, A) = \sigma(\hat{A} \text{Relu}(\tilde{A} X W_0) W_1) \quad (4)$$

Where X is the input matrix, A represents the graph adjacency matrix, \hat{A} and \tilde{A} represent the preprocessing step and self-connection structure respectively. W_0, W_1 represents weights of the first and second layers of ST-GCN, and $\sigma(\cdot)$. $\text{Relu}()$ is an activation function.

Temporal modelling is based on the RNN variant gated recurrent unit (GRU) [22] which has a simple structure and faster training ability. In the GRU model, an update gate z_t controls the degree of information retrieved from the previous state and a reset gate r_t controls the degree of ignoring the status information at the previous moment are configured as shown in Eq. 5.

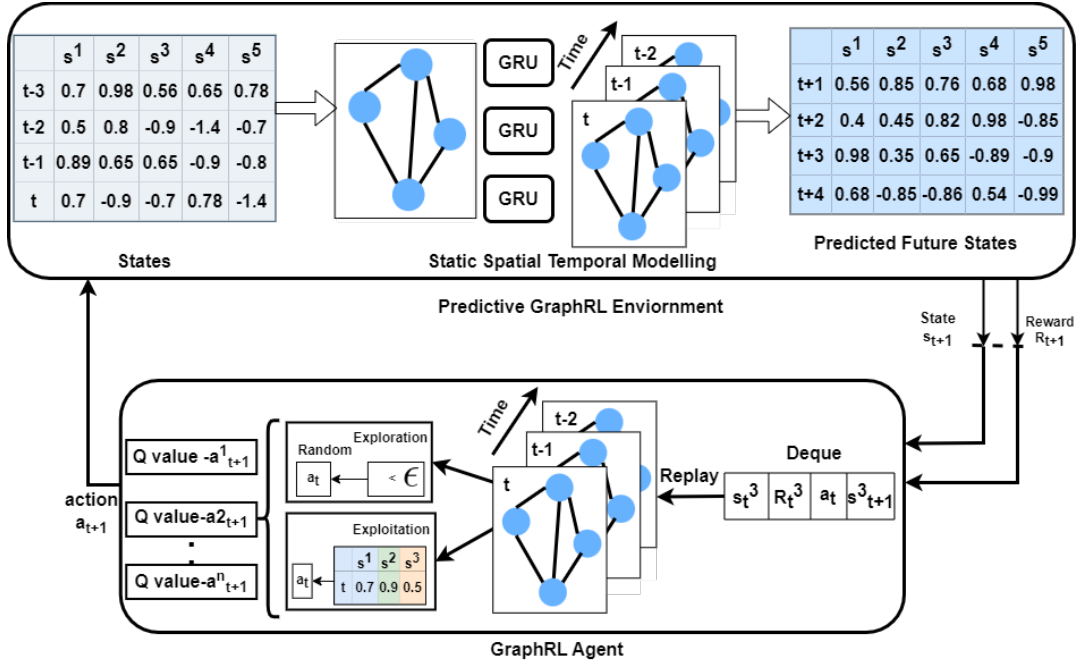


Fig. 2: GraphRL Framework

$$\begin{aligned}
 z_t &= \sigma_g(W_z[f(X_t, A), h_{t-1}] + b_z) \\
 r_t &= \sigma_g(W_r[f(X_t, A), h_{t-1}] + b_r) \\
 \hat{h}_t &= \phi_h(W_h[f(X_t, A), h_{t-1}](r_t \odot h_{t-1}) + b_h) \\
 h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \hat{h}_t
 \end{aligned} \quad (5)$$

x_t : input vector, h_t : output vector, \hat{h}_t : candidate activation vector, z_t : update gate vector, r_t : reset gate vector, W, U, b : parameter matrices and vector, σ_g : The original is a sigmoid function, ϕ_h The original is a hyperbolic tangent.

In the training process, the T-GCN model predicts the future states at time $t+1, t+2, t+3, \dots, t+n$ as \hat{Y}_t and compares with the real-data Y_t . This determines the loss function of the graph network [23] as shown in Eq. 6. To avoid over-fitting problems in the training process, the loss function is optimised with L2 regularisation L_{reg} and a hyperparameter λ .

$$loss = \|\hat{Y}_t - Y_t\| + \lambda L_{reg} \quad (6)$$

B. Predictive GraphRL Environment Algorithm

Algorithm 1 outlines the creation of the Predictive GraphRL Environment, a crucial component of the proposed GraphRL framework. This environment leverages the T-GCN, chosen for its effectiveness in capturing spatial-temporal dynamics essential for complex systems. The T-GCN's ability to forecast future states, in addition to analyzing current and historical data, makes it invaluable for critical applications like health monitoring and traffic management, where early detection and timely response are key. A significant feature of this

Algorithm 1 Predictive GraphRL Environment

Ensure: Input: time series data $\mathcal{D} = \{s_{t-n}, \dots, s_{t-2}, s_{t-1}, s_t\}$; a set of labels $\mathcal{K} = \{1, 2, \dots, K\}$
Ensure: Output: Predicted time series data of \mathcal{K} , a set of labels, in the form of states $\{s_{t+1}, s_{t+2}, s_{t+3}, s_{t+4}\}$

- 1: Define $forecast_model \leftarrow T - GCNModel$
- 2: $Train(forecast_model) \leftarrow forecast_model(\mathcal{D})$
- 3: $\{s_{t+1}, s_{t+2}, s_{t+3}, s_{t+4}\} \leftarrow forecast_model(predict)$
- 4: **Initialization** : $observation_space = \{s_i^t \in S\}, action_space = \{a_t \in A\}, rewardR$
- 5: $Set\ monitor_length = N$
- 6: **if** action is appropriate **then**
- 7: $R \leftarrow +reward$
- 8: **else**
- 9: $R \leftarrow -reward$
- 10: **end if**
- 11: $monitor_length \leftarrow N - 1$
- 12: $s_{t+1} \leftarrow s_t(monitor_length)$
- 13: **if** $N = 0$ **then**
- 14: $done \leftarrow True$
- 15: **else**
- 16: $done \leftarrow False$
- 17: **end if**
- 18: $visualize(a_t, R, vital\ signs)$
- 19: $initial_state \leftarrow s_t[0]$ ▷ reset environment

algorithm is its reward mechanism, which is instrumental in guiding the learning process. By awarding rewards based on the suitability of the agent's actions, the environment ensures that the agent's policy is aligned with the primary objectives of accurate forecasting and effective intervention. This design was motivated by the need for a proactive system, capable of not only forecasting but also informing real-time decision-making. The algorithm is meticulously structured to set up the observation space, action space, and reward policy based on predicted states. The initial lines (1-3) justify the use of the T-GCN model, especially for its applicability in

dynamic and nonlinear data contexts like vital sign monitoring. The subsequent lines (4-5) are dedicated to initializing the environment, forming the basis for RL-driven decision-making. The reward policy, detailed in lines 6-10, aligns with standard RL practices, promoting actions that yield beneficial outcomes. Finally, lines 11-19 focus on continuous monitoring and adaptation, a critical aspect for applications that demand real-time responsiveness, such as in healthcare scenarios. This algorithm represents a significant step in advancing the field of RL, moving from passive observation to an active role in shaping decisions.

C. GraphRL Agent

In this study, the Deep Q-Networks (DQN) algorithm is used. The Deep Q-Networks (DQN) algorithm, developed by Google's DeepMind, was initially designed for playing Atari games. This algorithm enabled the AI to learn game strategies directly from visual input, without requiring pre-programmed rules or prior game-specific training. In this algorithm, the Q-Learning functions are approximated using the proposed T-GCN model, and the learning agent is rewarded based on the graph network prediction of the right action for the current state.

1) *Q-Function Approximation:* T-GCN model used in this study to approximate the Q-Function for each action in the action space as shown in Fig. 2. The model is configured with parameters such as the relu activation function, mean square error as loss function, and Adam optimiser. The model gets trained with the state and its corresponding action. The learning agent performs an action $a_t \in A$ for a transition from state s_t to s'_t and achieves a reward R for the action. In this transition process, the maximum of the Q-function in Eq. 7 is calculated, and the discount of the calculated value uses a discount factor γ to suppress future rewards and focus on immediate rewards. The discounted future reward is added to the current reward to get the target value. The difference between the current prediction from the neural networks and the calculated target value provides a loss function. The loss function is a deviation of the predicted value from the target value and it can be estimated from Eq. 8. The square of the loss function allows for the punishment of the agent for a large loss value.

$$Q^\pi(s, a) = E_\pi \left\{ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, \pi(s_t)) | s_0 = s, a_0 = a \right\} \quad (7)$$

$$\text{loss} = \left(\underbrace{R + \gamma \cdot \max(Q^{\pi^*}(s, a))}_{\text{target_value}} - \underbrace{Q^\pi(s, a)}_{\text{predicted_value}} \right)^2 \quad (8)$$

2) *Exploration and Exploitation:* The concepts of exploration and exploitation are at odds with each other. Exploration involves randomly selecting actions that have not been performed before to uncover more possibilities and enhance the agent's understanding. Exploitation, on the other hand, entails selecting actions based on past experiences and knowledge to

maximize rewards. To balance the trade-off between exploration and exploitation, different strategies such as the greedy algorithm, epsilon-greedy algorithm, optimistic initialization, and decaying epsilon-greedy algorithm are employed. This study proposes controlling the exploration rate by multiplying the decay by the exploration rate. This approach reduces the number of explorations during execution as the agent learns patterns and maximizes rewards to achieve high scores. As the T-GCN model is retrained with previous experiences in the replay, the decay is multiplied by the exploration rate depending on the agent's ability to predict the right actions. All these parameters are defined as hyper-parameters for DQN learning agents.

Algorithm 2 Learning Agent

```

1: Initialize  $\gamma, \epsilon, \epsilon_{decay}, \epsilon_{min}, memory = \emptyset, batch\_size$ 
2: Define  $model \leftarrow T - GCN\_model$ 
3:  $memory \leftarrow \text{append}(s_t, a_t, R, s_{t+1})$ 
4: if  $np.random.rand() < \epsilon$  then ▷ Exploration
5:    $action\_value \leftarrow \text{random}(a_t)$ 
6: else ▷ Exploitation
7:    $action\_value \leftarrow model.predict(s_t)$ 
8: end if
9:  $minibatch \leftarrow \text{random}(memory, batch\_size)$ 
10: for  $s_t, a_t, R, s_{t+1}, done$  in  $minibatch$  do
11:    $target \leftarrow R$ 
12:   if not done then
13:      $target \leftarrow R + \gamma \cdot \max(model.predict(s_{t+1}))$ 
14:   end if
15:    $target\_f \leftarrow model.predict(s_t)$ 
16:    $target\_f[a_t] \leftarrow target$ 
17:    $model.fit(s_t, target\_f)$ 
18: end for
19: if  $\epsilon \geq \epsilon_{min}$  then
20:    $\epsilon \leftarrow \epsilon \cdot \epsilon_{decay}$ 
21: end if

```

3) *GraphRL Agent Algorithm:* Algorithm 2 introduces the GraphRL Agent, presents the functionality of the GraphRL Agent within a complex action-state environment, utilizing T-GCN for Q-function approximation. This integration enables effective handling of spatial-temporal data complexities, enhancing decision-making. The algorithm's design is founded on a strategic balance between exploration and exploitation, achieved via an epsilon-greedy strategy, crucial for adaptive learning and continual improvement in decision-making. It starts with initializing key parameters (Line 1), defining a T-GCN model for handling complex data structures (Line 2), and storing memories for experience replay (Line 3). The agent's learning process involves iterative learning from minibatches of experiences, computing target Q-values, and adjusting its policy (Lines 8-21), with a dynamic adjustment of the exploration rate (Lines 16-18). The predict() function, pivotal in the exploitation phase, utilizes the T-GCN model's predictions to guide actions, showcasing the algorithm's advanced approach in navigating dynamic environments through a blend of exploration and strategic exploitation.

4) *Implementation Algorithm:* Algorithm 3 serves as the comprehensive implementation of the GraphRL framework, intricately combining the Predictive GraphRL Environment (Algorithm 1) with the GraphRL Agent (Algorithm 2). This pivotal algorithm orchestrates the real-time interactions between the agent and the environment, thus forming the operational core of the framework. It outlines the simulation scope, input parameters (subjects \mathcal{C} , vital signs \mathcal{V} , and episodes

Algorithm 3 Proposed GraphRL Framework Implementation

Require: Input:
 $C = \{1, 2, \dots, C\}$: set of subjects
 $\mathcal{V} = \{1, 2, \dots, V\}$: set of vital signs
 $\mathcal{M} = \{1, 2, \dots, M\}$: number of episodes

Ensure: Output: Rewards achieved by Agents in each episode.

```

1:  $env \leftarrow ForecastingEnvironment()$  ▷ Algorithm 1
2:  $agent \leftarrow LearningAgent()$  ▷ Algorithm 2
3: for episode  $m \in \mathcal{M}$  do
4:    $state \leftarrow env.reset()$ 
5:    $score \leftarrow 0$ 
6:   for time in range(timesteps) do
7:      $a_t \leftarrow agent.action(s_t)$ 
8:      $s_{t+1}, R, done \leftarrow env.step(a_t)$ 
9:      $agent.memorize(s_t, a_t, R, s_{t+1})$ 
10:     $s_t \leftarrow s_{t+1}$ 
11:    if done then
12:       $print(m, score)$ 
13:       $break$ ;
14:    end if
15:  end for
16:   $agent.replay(batch\_size)$ 
17: end for

```

\mathcal{M}), and the output in the form of cumulative rewards. The initialization phase prepares the environment env and the agent for interaction. The episodic loop, encompassing the agent’s action-response cycle, is vital for continuous learning and adaptation. Crucial to this process is the `memorize()` function, which stores experiences for later recall during experience replay, allowing the agent to learn from past actions and refine its decision-making strategy. This algorithm thus encapsulates the dynamic and iterative nature of the GraphRL framework, highlighting the importance of memory and experience in the realm of advanced reinforcement learning, and showcasing its functionality in complex, evolving environments.

D. Bayesian Optimisation

Bayesian optimisation is a global optimisation method that uses a probabilistic model to guide the search for the optimal solution. The model is updated as new data points are sampled and evaluated, allowing the algorithm to improve its predictions over time by fine-tuning the L_{reg} , λ and minimising the loss function defined in Eq. 6. The basic idea behind Bayesian optimisation is to model the objective function, $f(x)$, as a Gaussian process (GP). The GP model is used to predict the objective function value at any point x , given the observations of the function at other points. The prediction is given by the posterior distribution of the GP, which is a Gaussian distribution with mean and variance given by Eq. 9.

$$\begin{aligned} \mu(x) &= k(x, X)^T (K + \sigma^2 I)^{-1} y \\ \sigma^2(x) &= k(x, x) - k(x, X)^T (K + \sigma^2 I)^{-1} k(X, x) \end{aligned} \quad (9)$$

Where X is the matrix of previously sampled points, y is the vector of corresponding function values, K is the Gram matrix of the covariance function evaluated at X , and σ^2 is the noise level in the function evaluations. The next point to sample is chosen based on an acquisition function, which balances the trade-off between exploration and exploitation. Common acquisition functions include the probability of improvement and the expected improvement. Given a set of observed points (X, y) and a Gaussian process prior, Bayesian optimisation seeks the point x^* that minimises the loss function value, given

by Eq. 10. The optimisation process continues iteratively, sampling new points and updating the GP model until a stopping criterion is met.

$$EI(x) = \mathbf{E}[max(0, f(x) - f(x^*))] \quad (10)$$

Where x^* is the current best point.

V. EXPERIMENT

The primary objective of this study is to overcome deep learning challenges such as the assumption of equally spaced and ordered data [24] and the lack of ability to incorporate graph structure where the data has a complex temporal structure [25]. These challenges are particularly relevant in domains such as health, weather, and traffic where it is important to analyze temporal patterns and make accurate forecasts for early warning systems. However, traditional deep learning models often fail to capture these complex patterns, limiting their effectiveness in these critical domains.

The proposed GraphRL framework is evaluated on three different forecasting applications: heart rate prediction, traffic forecast, and weather forecast. The framework predicts future events in the form of states and optimizes actions based on those predictions. The observation space is customized for each application and actions for the agent are pre-defined. The agent receives a reward for correctly predicting a state and communicating with the relevant team. The proposed approach is a generic framework that can be applied to monitor and predict time-series data and train an RL agent to learn the latent patterns of the monitoring process. The baseline models for comparison include traditional deep learning models such as GRU, LSTM, and RNNs.

A. Datasets

The GraphRL framework’s testing with datasets from healthcare, traffic, and weather domains was a deliberate strategy to assess its versatility and robustness in handling diverse time-series data. The choice of these varied domains was intended to demonstrate the framework’s adaptability and efficacy in different contexts. Each domain poses unique challenges: healthcare data’s complexity and sensitivity, traffic data’s dynamic patterns requiring real-time analysis, and weather data’s intricate interplay of environmental factors. Successfully navigating these distinct datasets underscores the framework’s capability for widespread real-world application. Additionally, using datasets from different fields facilitates a thorough evaluation of the framework, ensuring its versatility and effectiveness across various problem types and data structures. This comprehensive approach is vital for a tool designed for extensive applications in data analysis and prediction. Three datasets utilized for evaluating the GraphRL framework as shown in Fig. I, each from a different domain: healthcare, traffic, and weather. In healthcare, the WESAD dataset, containing electrocardiogram (ECG) and photoplethysmogram (PPG) data from 17 participants, offers a rich source of biometric time-series data for pattern recognition analysis. The Los Angeles (LA) Traffic dataset, sourced from the Los

Angeles Department of Transportation (LADOT), provides real-time urban traffic data like traffic counts and speeds, while the Large-Scale Traffic and Weather Events (LSTW) dataset, with data across the United States, uniquely combines traffic conditions and weather events, posing a multifaceted challenge for the framework.

TABLE I: Datasets

Dataset	Domain	Key Features	Statistics	Suitability
WESAD [26]	Healthcare	Physiological data (ECG, PPG), motion data (accelerometers)	17 participants	Rich biometric time-series data, ideal for testing pattern recognition in health-related data.
LA Traffic [27]	Traffic	Traffic counts, speeds, travel times	Data from LADOT	Real-time urban traffic data, useful for analyzing and predicting dynamic traffic flows.
LSTW [28]	Weather	Traffic conditions, weather events	Data across the United States	Combines traffic and weather data, challenging the framework to handle complex, multifactorial scenarios.

B. Baseline Models

In our study, we selected three baseline models for comparison, each epitomizing state-of-the-art approaches in multi-agent forecasting, graph neural networks, and traffic prediction. These models were chosen based on their innovative methodologies and proven effectiveness in areas closely aligned with our research objectives.

- **ELMA Method** [29]: Developed by Li et al. [29], the ELMA method utilizes graph neural networks for forecasting multi-agent activities, particularly adept at handling spatiotemporal data. Its novelty lies in the use of energy-based learning, making it an excellent benchmark against our framework, which similarly leverages graph-based techniques in complex environments.
- **Self-Supervised Technique** [30]: This technique is pioneering in self-supervised learning for predicting multi-agent driving behavior. Its relevance to our study comes from its focus on behavior prediction in diverse scenarios, using self-supervised domain knowledge—an advanced trend in multi-agent learning.
- **Internet Traffic Prediction Study** [31]: It involves internet traffic prediction using distributed multi-agent learning, employing LSTM and GRU models. GRU’s superior performance in their study provides a valuable point of comparison for our research, which focuses on sophisticated learning techniques in traffic prediction.

Each of these models represents a significant stride in their respective fields. Their selection for comparison in our study is justified by their alignment with our research goals and their benchmark status in handling complex, dynamic datasets. By comparing our GraphRL framework against these models, we aim to demonstrate our approach’s novelty and effectiveness in diverse real-world applications.

C. Evaluation Metrics

Mean Absolute Error (MAE) is a widely-utilized regression metric that gauges the average magnitude of errors between predicted and actual values in a dataset. It is calculated by averaging the absolute differences between these values, yielding a singular metric. Root Mean Squared Error (RMSE)

TABLE II: Performance of the proposed framework in health forecasting

		15 Min	30 Min	45 Min	60 Min
ELMA [29]	MAE	6.2	6.2	6.2	6.13
	MAPE	13.91	13.91	13.91	13.91
	RMSE	8.75	8.75	8.75	8.67
GRU [30]	MAE	0.95	0.95	0.97	0.98
	MAPE	5.47	5.48	5.51	5.5
	RMSE	1.25	1.25	1.27	1.28
GRU-Based	MAE	1.02	1.02	1.25	1.65
	MAPE	8	3.47	4.53	5.27
Multi-Agent [31]	RMSE	2.46	2.58	2.69	3.09
	MAE	0.56	0.87	0.68	0.7
GraphRL (Ours)	MAPE	2.8	2.9	2.65	3.98
	RMSE	1.18	1.47	1.3	1.32

is another prominent regression metric, assessing the average magnitude of differences between predicted and actual values. RMSE is computed as the square root of the mean of these squared differences. Mean Absolute Percentage Error (MAPE) represents yet another regression metric, quantifying the average absolute percentage error between predicted and actual values. It is derived by averaging the absolute differences between these values, expressed as a percentage of the actual values. Conversely, Cumulative Rewards is a performance metric specific to reinforcement learning. It measures the total rewards an agent accumulates over a specified timeframe or across a set number of actions, calculated by summing all rewards received during this period.

In the context of the experiments conducted for this study, Python version 3.7.6 served as the programming environment, with the deployment of several packages including TensorFlow, Keras, OpenAI Gym, and stable_baselines3.

VI. RESULTS AND ANALYSIS

In this section, the proposed framework performance in terms of time series forecasting and RL monitoring is compared to the baseline models in each application.

A. Predictive GraphRL Performance

Healthcare Forecasting: The proposed framework is evaluated to monitor health status by predicting future vital signs such as heart rate. Based on the sensor data and other clinical parameters such as ECG, Respiration, the time series prediction of the heart rate is conducted. The predicted values of heart for the next one hour are break-down into different time intervals (15 minutes, 30 minutes, 45 minutes, 60 minutes). Each of these time interval values acts as an observation for the GraphRL agent to monitor and communicate with the appropriate emergency team. The observation space of the vital sign, action space of different emergency teams and rewards for the agent actions in the predictive GraphRL environment are defined based on the modified early warning scores (MEWS) [32]. For the evaluation process, the WESAD dataset

TABLE III: Performance of the proposed framework in traffic forecasting

		15 Min	30 Min	45 Min	60 Min
ELMA [29]	MAE	6.73	6.73	6.73	6.72
	MAPE	6.73	15.14	15.14	15.07
	RMSE	6.72	9.4	9.4	9.39
GRU [30]	MAE	1.04	1.04	1.04	1.04
	MAPE	6.04	6.01	5.96	6.1
	RMSE	1.36	1.36	1.36	1.36
GRU-Based	MAE	1.85	1.85	1.96	1.82
Multi-Agent [31]	MAPE	6.07	5.7	4.93	6.07
	RMSE	2.88	3.21	3.43	3.54
GraphRL (Ours)	MAE	0.65	0.78	0.64	0.8
	MAPE	4.1	7.85	5.65	7.99
	RMSE	1.27	1.22	1.26	1.41

is adopted to conduct time series forecasting of heart rate. The proposed T-GCN in the predictive GraphRL environment performs better than the other baseline frameworks ELMA, GRU, and GRU-Based Multi-Agent as shown in Tab. II. It achieves the lowest MAE, MAPE and RMSE values in all the time intervals.

Traffic Forecasting: The goal of the proposed framework is to predict traffic using the predictive GraphRL environment. The system takes in data with the following features: EventId, Type, Severity, TMC, Description, StartTime, EndTime, TimeZone, LocationLat, LocationLng, Distance, AirportCode, Number, Street, Side, City, County, State, and ZipCode. The observation space includes the current traffic state, which is represented by the traffic events and their severity in a particular region. The actions referred to possible traffic management strategies, such as altering traffic light timings or changing the speed limit. The rewards are defined based on the efficiency of the chosen strategy, such as reduced travel time or decreased congestion. For all the baseline models and the proposed framework, the MAE, MAPE, and RMSE values are reported for forecasting at 15, 30, 45, and 60-minute intervals. As shown in Tab. III, T-GCN outperforms the other models for all the forecasting intervals with the lowest MAE, MAPE, and RMSE values. The second-best performer is the GRU-Based Multi-Agent model, followed by GRU and ELMA.

Weather Forecasting: In weather forecasting, the goal of the proposed framework is to use past weather data to predict future weather events and to optimise actions based on those predictions. In the predictive environment, the observation space is configured based on both the traffic and weather events datasets, including the event type, severity, start time, end time, location (latitude and longitude), and timezone. The actions represent the decisions the RL agent can take based on the observation space. For example, the agent could decide to issue a warning or alert for severe weather, adjust traffic signals or road signs, or change the speed limit on certain roads. The agent could receive a reward for correctly predicting severe weather and issuing a timely warning. Using

TABLE IV: Performance of the proposed framework in weather forecasting

		15 Min	30 Min	45 Min	60 Min
ELMA [29]	MAE	6.69	6.69	6.69	6.65
	MAPE	6.69	15.02	15.02	14.99
	RMSE	6.69	9.39	9.39	9.34
GRU [30]	MAE	1.03	1.03	1.04	1.04
	MAPE	5.96	5.94	5.93	6.02
	RMSE	1.36	1.35	1.36	1.36
GRU-Based	MAE	1.65	1.65	1.85	2.02
Multi-Agent [31]	MAPE	7.32	4.71	4.89	5.86
	RMSE	2.76	2.99	3.16	3.43
GraphRL (Ours)	MAE	0.61	0.83	0.66	0.75
	MAPE	3.95	5.88	5.15	7.99
	RMSE	1.23	1.12	1.28	1.26

TABLE V: Proposed GraphRL Performance

AI Agents	WESAD	LAM Traffic Forecasting	US Weather Forecasting
Q Learning	43130	28840	39480
PPO	39480	33945	29480
A2C	41195	22845	40615
Double DQN	42615	25600	33945
DDPG	44600	34590	39945
DQN	41986	35219	40985
GraphRL	48790	36195	53145

the proposed GraphRL framework allows modelling the relationships between different weather events and their impact on traffic in a more efficient way than traditional machine learning methods. The GraphRL agent learns from these relationships to make better decisions and improve its predictions over time. Comparing the different models, T-GCN had the best performance across all metrics and different time intervals: 15, 30, 45, and 60 minutes, followed by GRU-Based Multi-Agent, GRU, and ELMA. The results show that the T-GCN model had the lowest MAE, MAPE, and RMSE values for all forecasting intervals, indicating its superior forecasting performance compared to the other models as shown in Tab. IV.

B. GraphRL Agent Performance

The proposed RL agent was enabled with T-GCN and its performance is compared with other traditional RL agents as shown in Tab. V. The table provides a comparison of different AI agents and their performance on three different datasets: WESAD, LAM Traffic Forecasting, and US Weather Forecasting. The performance of each agent is measured by a score, which is the total score achieved by the agent on the task over ten episodes. From the table, it can be seen

that the proposed GraphRL agent is the most efficient agent on the WESAD dataset, as it scored the highest score. The DDPG and Q-Learning agents have the second-highest score on the WESAD dataset. On the LAM Traffic Forecasting dataset, the Q-Learning agent scored the lowest, and the proposed GraphRL agent scored the highest. On the US Weather Forecasting dataset, the A2C agent scored the lowest, while the GraphRL agent scored the highest. The GraphRL agent has outperformed other RL agents in all three predictive and monitoring applications.

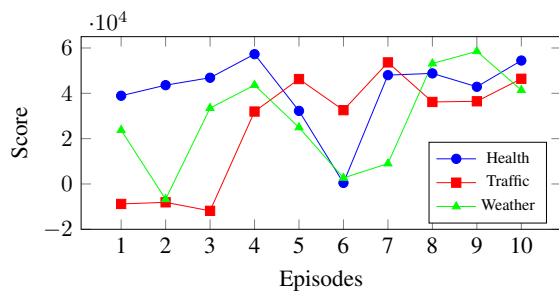


Fig. 3: GraphRL Agent Rewards Distribution

The performance of the GraphRL agent is measured by the episode score, which appears to be the total score achieved by the agent after a certain number of episodes. The breakdown of the proposed agent's score in each episode of the three applications is presented and compared in Fig. 3. The agent's performance on the WESAD dataset is relatively consistent, with the scores fluctuating between 32245 and 57280. On the LAM Traffic Forecasting dataset, the agent's performance is relatively inconsistent, with the scores fluctuating between -11845 and 46295. On the US Weather Forecasting dataset, the agent's performance is also relatively inconsistent, with the scores fluctuating between -6765 and 58530. This inconsistency of the scores is due to the exploration rate where the algorithm tries exploring all the actions randomly instead of using T-GCN model predictions.

VII. BAYESIAN OPTIMISATION RESULTS

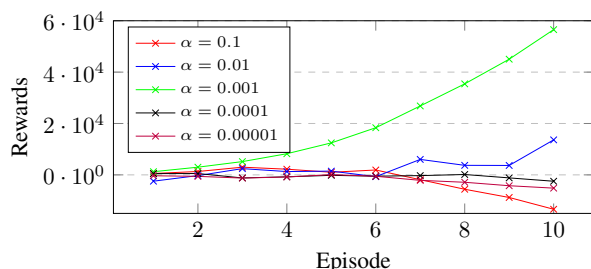


Fig. 4: Bayesian optimisation of α for GraphRL Agent

The results of Bayesian optimisation for the proposed GraphRL agent using different values of the learning rate parameter, α , during different episodes are shown in Fig. 4.

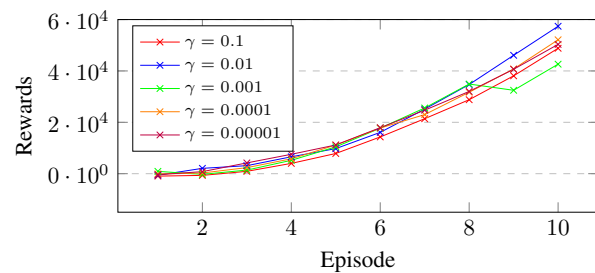


Fig. 5: Bayesian optimisation of γ for GraphRL Agent

The values in the y-axis of the line chart represent the scores or rewards obtained by the agent during each episode. It can be observed that the performance of the agent varies for different values of alpha. For example, in episode 1, the agent performs better with alpha = 0.001 (20630) compared to the other values. Similarly, in episode 10, the agent performs better with alpha = 0.001 (97380) compared to the other values. These results suggest that the optimal value of alpha for the agent is $\alpha = 0.001$, and Bayesian optimisation can be used to find the best value of alpha for a given task.

These results in Fig. 5 show the performance of an RL agent using temporal GCN for Q function approximation, using different values of the discount factor gamma. As we can see, the performance varies greatly depending on the value of the gamma chosen. A high value of gamma (0.95) results in poor performance, while lower values (0.75) result in better performance. This suggests that a lower discount factor is more appropriate, as it gives more weight to immediate rewards and less to future rewards. It also suggests that there is an optimal value of gamma, which would need to be further explored through more extensive experimentation.

VIII. CONCLUSION

The GraphRL framework, introduced in this study, embodies an innovative amalgamation of T-GCN and RL. It is specifically engineered to augment the prediction of future states in dynamic environments. Rigorous evaluations, utilizing an array of datasets such as WESAD, LA Traffic Forecasting, and US Weather Forecasting, have substantiated the framework's enhanced performance compared to conventional RL agents. Nonetheless, it is imperative to acknowledge that the efficacy of GraphRL is significantly contingent upon the caliber of the input data and necessitates substantial computational resources. The framework's reliance on data of high quality and structure constitutes a considerable limitation, with its accuracy and effectiveness being closely bound to the data's integrity. Additionally, the computational requisites, predominantly due to the T-GCN model integration, present challenges in scalability and broader applicability.

Future enhancements of the GraphRL framework will be directed towards surmounting these constraints and broadening its functional scope. Prospective developments entail the incorporation of spatial data processing, aimed at bolstering the framework's analytical prowess, particularly in processing

data with spatial or geographical dimensions. Investigating a spectrum of graph-based models could yield insights for enhancing both the efficiency and efficacy of the framework. Furthermore, the exploration of real-time adaptive learning strategies presents a promising avenue for subsequent research. Such advancements are anticipated to enable the framework to dynamically adapt to evolving data patterns and environmental shifts. In summation, the GraphRL framework signifies a substantial advancement in the domain of time-series prediction and monitoring. Its adeptness in managing complex temporal data surpasses traditional RL methodologies, heralding innovative applications in sectors such as healthcare, traffic management, and environmental forecasting. As the framework undergoes continued refinement and evolution, it is positioned to emerge as an instrumental component in the progression of predictive analytics and intelligent monitoring systems, with extensive applicability across diverse fields.

REFERENCES

- [1] A. Zhang, L. Xing, J. Zou, and J. C. Wu, "Shifting machine learning for healthcare from development to deployment and from models to data," *Nature Biomedical Engineering*, July 2022.
- [2] T. Shaik, X. Tao, N. Higgins, R. Gururajan, Y. Li, X. Zhou, and U. R. Acharya, "Fedstack: Personalized activity monitoring using stacked federated learning," *Knowledge-Based Systems*, vol. 257, p. 109929, 2022.
- [3] G. Gao, Q. Gao, X. Yang, M. Pajic, and M. Chi, "A reinforcement learning-informed pattern mining framework for multivariate time series classification," 2022.
- [4] E. M. Forman, S. G. Kerrigan, M. L. Butryn, A. S. Juarascio, S. M. Manasse, S. Ontañón, D. H. Dallal, R. J. Crochiere, and D. Moskow, "Can the artificial intelligence technique of reinforcement learning use continuously-monitored digital data to optimize treatment for weight loss?," *Journal of behavioral medicine*, vol. 42, no. 2, pp. 276–290, 2019.
- [5] T. Wang, B. Lu, W. Wang, W. Wei, X. Yuan, and J. Li, "Reinforcement learning-based optimization for mobile edge computing scheduling game," *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2022.
- [6] E. Shin, S. Swaroop, W. Pan, S. Murphy, and F. Doshi-Velez, "Modeling mobile health users as reinforcement learning agents," *arXiv preprint arXiv:2212.00863*, 2022.
- [7] V. A. Taylor, I. Moseley, S. Sun, R. Smith, A. Roy, V. U. Ludwig, and J. A. Brewer, "Awareness drives changes in reward value which predict eating behavior change: Probing reinforcement learning using experience sampling from mobile mindfulness training for maladaptive eating," *Journal of Behavioral Addictions*, vol. 10, no. 3, pp. 482–497, 2021.
- [8] K. Chen, J. Fong, and H. Soh, "Mirror: Differentiable deep social projection for assistive human-robot communication," *arXiv e-prints*, pp. arXiv-2203, 2022.
- [9] M. Zhou, Y. Mintz, Y. Fukuoka, K. Goldberg, E. Flowers, P. Kaminsky, A. Castillejo, and A. Aswani, "Personalizing mobile fitness apps using reinforcement learning," in *CEUR workshop proceedings*, vol. 2068, NIH Public Access, 2018.
- [10] Z. Li, "A hierarchical autonomous driving framework combining reinforcement learning and imitation learning," in *2021 International Conference on Computer Engineering and Application (ICCEA)*, pp. 395–400, IEEE, 2021.
- [11] S. B. Asimwe, E. Vittinghoff, and M. Whooley, "Vital signs data and probability of hospitalization, transfer to another facility, or emergency department death among adults presenting for medical illnesses to the emergency department at a large urban hospital in the united states," *The Journal of Emergency Medicine*, vol. 58, pp. 570–580, apr 2020.
- [12] C. G. Scully and C. Daluwatte, "Evaluating performance of early warning indices to predict physiological instabilities," *Journal of Biomedical Informatics*, vol. 75, pp. 14–21, nov 2017.
- [13] M. M. Baig, H. G. Hosseini, and M. Linden, "Machine learning-based clinical decision support system for early diagnosis from real-time physiological data," in *2016 IEEE Region 10 Conference (TENCON)*, IEEE, nov 2016.
- [14] K. Alghatani, N. Ammar, A. Rezgui, A. Shaban-Nejad, *et al.*, "Predicting intensive care unit length of stay and mortality using patient vital signs: machine learning model development and validation," *JMIR medical informatics*, vol. 9, no. 5, p. e21347, 2021.
- [15] A. Youssef Ali Amer, F. Wouters, J. Vranken, P. Dreesen, D. de Korte-de Boer, F. van Rosmalen, B. C. van Bussel, V. Smit-Fun, P. Duflo, J. Guiot, *et al.*, "Vital signs prediction for covid-19 patients in icu," *Sensors*, vol. 21, no. 23, p. 8131, 2021.
- [16] G. Harerimana, J. W. Kim, and B. Jang, "A multi-headed transformer approach for predicting the patient's clinical time-series variables from charted vital signs," *IEEE Access*, vol. 10, pp. 105993–106004, 2022.
- [17] Z. Xie, H. Wang, S. Han, E. Schoenfeld, and F. Ye, "Deepvps: A deep learning approach for rf-based vital signs sensing," in *Proceedings of the 13th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*, pp. 1–5, 2022.
- [18] A. M. T. Elsir, A. Khaled, and Y. Shen, "Hlgt: Hybrid local-global spatio-temporal model for travel time estimation using siamese graph convolutional with triplet networks," *Expert Systems with Applications*, vol. 229, p. 120502, 2023.
- [19] S. Wang, J. Cao, and S. Y. Philip, "Deep learning for spatio-temporal data mining: A survey," *IEEE transactions on knowledge and data engineering*, vol. 34, no. 8, pp. 3681–3700, 2020.
- [20] Z. Chen, Y. Wan, Y. Liu, and A. Valera-Medina, "A knowledge graph-supported information fusion approach for multi-faceted conceptual modelling," *Information Fusion*, vol. 101, p. 101985, 2024.
- [21] F. B. Hüttel, F. Rodrigues, and F. C. Pereira, "Mind the gap: Modelling difference between censored and uncensored electric vehicle charging demand," *Transportation Research Part C: Emerging Technologies*, vol. 153, p. 104189, 2023.
- [22] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," *arXiv preprint arXiv:1409.1259*, 2014.
- [23] B. Yu, H. Yin, and Z. Zhu, "Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting," *arXiv preprint arXiv:1709.04875*, 2017.
- [24] G. Karagiorgi, G. Kasieczka, S. Kravitz, B. Nachman, and D. Shih, "Machine learning in the search for new fundamental physics," *Nature Reviews Physics*, vol. 4, no. 6, pp. 399–412, 2022.
- [25] P. Almasan, J. Suárez-Varela, K. Rusek, P. Barlet-Ros, and A. Cabellos-Aparicio, "Deep reinforcement learning meets graph neural networks: exploring a routing optimization use case," *Computer Communications*, vol. 196, pp. 184–194, 2022.
- [26] P. Schmidt, A. Reiss, R. Duerichen, C. Marberger, and K. Van Laerhoven, "Introducing wesad, a multimodal dataset for wearable stress and affect detection," in *Proceedings of the 20th ACM international conference on multimodal interaction*, pp. 400–408, 2018.
- [27] Y. Li, R. Yu, C. Shahabi, and Y. Liu, "Diffusion convolutional recurrent neural network: Data-driven traffic forecasting," *arXiv preprint arXiv:1707.01926*, 2017.
- [28] S. Moosavi, M. H. Samavatian, A. Nandi, S. Parthasarathy, and R. Ramnath, "Short and long-term pattern discovery over large-scale geospatiotemporal data," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 2905–2913, 2019.
- [29] Y. Li, P. Wang, L. Chen, Z. Wang, and C.-Y. Chan, "Elma: Energy-based learning for multi-agent activity forecasting," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 36, pp. 1482–1490, 2022.
- [30] H. Ma, Y. Sun, J. Li, and M. Tomizuka, "Multi-agent driving behavior prediction across different scenarios with self-supervised domain knowledge," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pp. 3122–3129, IEEE, 2021.
- [31] W. Jiang, M. He, and W. Gu, "Internet traffic prediction with distributed multi-agent learning," *Applied System Innovation*, vol. 5, no. 6, p. 121, 2022.
- [32] V. Signs, "Canberra hospital and health services clinical procedure," 2021.

8.2 Summary

The GraphRL framework is showcased as a significant leap forward in the realm of time series forecasting, especially in applications requiring nuanced interpretation of complex, dynamic data. Through rigorous comparative analysis, GraphRL's superior performance over traditional and contemporary models is highlighted, underscoring its efficacy in capturing intricate temporal patterns and its versatility across varied application domains. The chapter emphasizes the framework's potential in revolutionizing early warning systems and real-time monitoring solutions, paving the way for future advancements in intelligent, data-driven decision-making processes.

CHAPTER 9: PAPER 8 - QXAI: EXPLAINABLE AI FRAMEWORK FOR QUANTITATIVE ANALYSIS IN PATIENT MONITORING SYSTEMS

9.1 Introduction

This chapter presents the QAXI framework, a pioneering approach in explainable AI tailored for quantitative data analysis in healthcare. It intricately combines Shapley values and attention mechanisms within deep learning models to demystify AI predictions, offering both local and global interpretability. This novel framework is designed to bridge the gap between AI's predictive power and the clinical need for understandable and trustworthy decision-making support, particularly in patient monitoring systems where interpreting AI-driven predictions is crucial for clinical interventions.

QXAI: Explainable AI Framework for Quantitative Analysis in Patient Monitoring Systems

Thanveer Shaik, Xiaohui Tao, Haoran Xie, Lin Li, Juan D. Velásquez, and Niall Higgins

Abstract—Artificial Intelligence techniques can be used to classify a patient’s physical activities and predict vital signs for remote patient monitoring. Regression analysis based on non-linear models like deep learning models has limited explainability due to its black-box nature. This can require decision-makers to make blind leaps of faith based on non-linear model results, especially in healthcare applications. In non-invasive monitoring, patient data from tracking sensors and their predisposing clinical attributes act as input features for predicting future vital signs. Explaining the contributions of various features to the overall output of the monitoring application is critical for a clinician’s decision-making. In this study, an Explainable AI for Quantitative analysis (QXAI) framework is proposed with post-hoc model explainability and intrinsic explainability for regression and classification tasks in a supervised learning approach. This was achieved by utilizing the Shapley values concept and incorporating attention mechanisms in deep learning models. We adopted the artificial neural networks (ANN) and attention-based Bidirectional LSTM (BiLSTM) models for the prediction of heart rate and classification of physical activities based on sensor data. The deep learning models achieved state-of-the-art results in both prediction and classification tasks. Global explanation and local explanation were conducted on input data to understand the feature contribution of various patient data. The proposed QXAI framework was evaluated using PPG-DaLiA data to predict heart rate and mobile health (MHEALTH) data to classify physical activities based on sensor data. Monte Carlo approximation was applied to the framework to overcome the time complexity and high computation power requirements required for Shapley value calculations.

Index Terms—Explainability, Shapley, Attention, Monte Carlo, Vital Signs, Physical Activities

I. INTRODUCTION

In the realm of modern healthcare, the integration of cutting-edge technology, notably through remote monitoring systems, represents a pivotal advancement in patient care and the management of diseases. These systems play an essential role in the prompt detection and averting of grave health events, chiefly through their capacity to precisely monitor

Thanveer Shaik and Xiaohui Tao are with the School of Mathematics, Physics & Computing, University of Southern Queensland, Toowoomba, Queensland, Australia (e-mail: Thanveer.Shaik@usq.edu.au, Xiaohui.Tao@usq.edu.au).

Haoran Xie is with the Department of Computing and Decision Sciences, Lingnan University, Tuen Mun, Hong Kong (e-mail: hrxie@ln.edu.hk)

Lin Li is with the School of Computer and Artificial Intelligence, Wuhan University of Technology, China (e-mail: cathylin@whut.edu.cn)

Juan D. Velásquez is with the Industrial Engineering Department at University of Chile, Chile (e-mail: jvelasqu@dii.uchile.cl)

Niall Higgins is with Metro North Hospital and Health Service, Royal Brisbane and Women’s Hospital, and also with School of Nursing, Queensland University of Technology, Brisbane, Australia (e-mail: Niall.Higgins@health.qld.gov.au).

and scrutinize vital signs such as temperature, pulse, respiratory rate, and mean arterial pressure [1, 2]. However, the scope of traditional monitoring systems is often constrained to displaying a patient’s current health status, which limits their effectiveness in preemptively predicting and managing potential health complications.

The advent of Artificial Intelligence (AI) and deep learning heralds a new era in healthcare, transcending the boundaries of traditional methods by offering predictive insights that are indispensable for early and effective medical interventions [3, 4]. Nevertheless, these advanced methodologies come with their own set of complexities, chief among them being the lack of transparency and comprehensibility in deep learning models. These models, often labeled as “black-box” models, pose a significant challenge in elucidating how input factors correlate with the predictive outcomes [5, 6]. This issue is particularly critical in healthcare, where understanding the rationale behind AI-driven decisions is vital for their acceptance in clinical settings and for ensuring ethical applications of such technologies.

In response to these challenges, our research presents an innovative Explainable AI framework tailored for Quantitative data (QXAI), ingeniously amalgamating the Shapley values concept [7] with an attention mechanism in the realm of deep learning models. Our approach is uniquely poised to demystify AI predictions on both granular (local) and aggregate (global) scales. It provides insightful revelations on how each individual feature contributes to specific input records and offers a comprehensive overview of feature contributions throughout the entire model. This dual-level explanation capability is adeptly employed in our framework for the purpose of predicting human vital signs and classifying physical activities, utilizing two advanced deep learning models: Artificial Neural Networks (ANN) and attention-based Bidirectional LSTM (BiLSTM). The empirical evidence from our study highlights the framework’s proficiency in delivering detailed Shapley values and attention weights for each input feature, thereby clarifying their respective impacts on the outcomes of deep learning models.

Recognizing the computational demands in calculating Shapley values for extensive datasets, we have judiciously integrated the Monte Carlo method of approximation with random sampling. This strategic addition not only mitigates the computational complexities but also augments the practical utility of our framework across a spectrum of real-world applications.

Overall, our study represents a significant advancement in the field of explainable AI within healthcare. The key

contributions of our research are:

- Development of an innovative, adaptable Explainable AI framework (QXAI) for quantitative data analysis in healthcare. This framework uniquely combines attention layer mechanisms with Shapley values within deep learning models, setting a new standard in AI explainability.
- Comprehensive evaluation of the framework's explainability capabilities, focusing on the importance of features and providing both local and global explanations. This dual approach significantly enhances the understanding of AI models, offering insights into their cognitive and behavioral aspects.
- Adoption of the Monte Carlo method to address the computational challenges in calculating Shapley values, especially for large datasets. This method significantly reduces the computational overhead, making the framework more practical for real-world applications.
- Establishment of a new paradigm in patient monitoring systems for interpreting and explaining AI predictions related to vital signs and physical activity classification. This advancement is pivotal for clinical decision-making, offering a more nuanced and in-depth understanding of patient health dynamics.

The remainder of the article is organized as follows: Section II presents related works on explainability in healthcare applications. Section III presents a formal definition of the research problem addressed. Section IV details the novel QXAI framework to explain prediction and classification problems proposed in this study. Experimental design, dataset description, data modelling, and traditional models are discussed in Section V. In Section VI, experimental results of the QXAI framework are discussed, along with its explainability and feature identification performance. Section VII discusses the random sampling approximation using the Monte Carlo method. In Section VIII, we discuss implications, strengths, and limitations of the study. Finally, the paper concludes with Section IX.

II. RELATED WORK

In the realm of remote patient monitoring systems, the primary objective is to promptly identify high-risk patients, enabling clinicians to allocate resources effectively and intervene in a timely manner. The integration of machine learning and AI in these systems has led to significant advancements in predictive healthcare.

A. Machine Learning in Healthcare Prediction

Gong et al. [8] developed a machine learning-based framework for predicting acute kidney injury (AKI), showcasing an end-to-end decision support system that encompasses data pre-processing, risk prediction, and model explanation. This framework utilized logistic regression, random forest, and a voting-based ensemble model, along with gradient boosting algorithms, to address the challenges posed by imbalanced datasets. The model's prediction capability within 48 hours was complemented by SHapley Additive exPlanations (SHAP) values for a dual perspective: a global view highlighting

critical factors and a local view detailing individual patient-level feature contributions. In addition, Wu et al. [9] compared eight feature selection methods to enhance AKI prediction, underlining the importance of feature selection stability and similarity.

B. Assessment of Interpretability Techniques

ElShawi et al. [10] proposed quantitative measures to assess the quality of several model-agnostic interpretability techniques, including LIME, SHAP, Anchors, and others. Their study utilized a random forest model to predict mortality and diabetes risk, evaluating the performance of these interpretability techniques in terms of similarity, bias detection, execution time, and trust. In a separate study, Elshawi et al. [11] applied global and local explainability techniques to predict the risk of hypertension, enhancing the transparency of machine learning outcomes. Ilic et al. [12] introduced an explainable boosted linear regression (EBLR) algorithm for time series forecasting, demonstrating that maintaining interpretability does not necessarily compromise model performance.

C. Attention Mechanism in Deep Learning

The attention mechanism, initially a breakthrough in machine translation tasks, has been adapted for healthcare applications. Bari et al. [13] conducted an empirical evaluation of attention-based deep neural networks, assessing prediction performance, explainability correctness, and sensitivity. Their results indicated that multi-variable LSTM models with explainability features performed well with complex data. Kaji et al. [14] implemented an attention-based LSTM model for predicting medical conditions like sepsis and myocardial infarction, using MIMIC-III dataset patient data. They highlighted the importance of the attention layer in extracting influential input features for better explainability. Chen et al. [15] further advanced this field by proposing bilateral asymmetric skewed Gaussian attention (bi-SGA) to improve the performance and interpretability of deep convolutional neural networks.

D. Gap in Literature and Study Contribution

The literature reveals that while deep learning is capable of predicting vital signs with minimal healthcare domain knowledge, its lack of explainability remains a significant drawback. This underscores the need for explainable AI methods to demystify the results produced by these "black-box" models. Our study addresses this gap by introducing a novel framework that not only estimates feature importance for enhancing explainability but also provides both global and local interpretations of deep learning model predictions. This comprehensive framework aims to balance the trade-off between deep learning model performance and its explainability, thereby contributing significantly to the field of predictive healthcare.

III. RESEARCH PROBLEM

The central research problem tackled in this study is the elucidation of deep learning model results, particularly the interpretation of predictions based on independent feature inputs

in healthcare settings. This task involves comprehending the causal relationships between input factors and their effect on model predictions. It's crucial for healthcare professionals to grasp the rationale behind AI-driven predictions, understanding how variations in input feature values can influence these predictions. In a scenario where a deep learning model M uses N features, denoted as x_j (where $j = 1, \dots, N$), to predict an output y , the research aims to elucidate how each input feature x contributes to this prediction. This understanding is vital for models where weights w_j are applied to respective features j at different layers of model M . This process can be mathematically represented as:

$$y \leftarrow f_M(w_j \cdot x_j) \quad (1)$$

To enhance the explainability of predictions from complex, non-linear models such as neural networks and deep learning, it is essential to quantify the contribution of each feature, φ_{x_j} , in a comprehensible manner. To enhance the explainability of non-linear model predictions, the contribution of each feature φ_{x_j} can be estimated into two patterns.

$$\varphi_{x_j} = w_j * x_j - E(w_j * X_j) \quad (2)$$

$$\sum_{j=1}^N \varphi_{x_j} = \sum_{j=1}^N w_j * x_j - E(w_j * x_j) \quad (3)$$

- The first pattern estimates the model output with each feature and subtracts the output with the average effect of all the features, $E(w_j * X_j)$ as shown in Equation 2. The same approach can estimate the contributions of all features. Summing up all the features' contribution in a prediction instance is, where Equation 3 shows the predicted value $f_M(x)$ minus the average predicted value $E(f_M(x))$ for the instance x .
- The second pattern adds an attention layer to the non-linear model and enables the model to focus on certain important features contributing to the output. This pattern creates a representation h_j with $j = 1, \dots, N$ of each input in vector space, and the weighted sum of the representation act as context vectors as shown in Equation 4. Extracting the weights for each input feature can influence output feature contribution φ_{x_j} .

$$c = \sum_{j=1}^N \alpha_j h_{x_j} \quad (4)$$

In this current study, the two patterns estimate feature contribution to explain the prediction process of the deep learning model.

IV. EXPLAINABLE AI FOR QUANTITATIVE DATA (QXAI)

In this section, Explainable AI for Quantitative data (QXAI) is proposed to estimate input feature importance in deep learning model results that could be prediction or classification tasks. The proposed framework can provide explainability at two levels, one is post-hoc explainability using Shapley values and the other is intrinsic explainability using attention mechanism as shown in Fig. 1.

A. Shapley Values Calculation

To explain the contribution of input features, the Shapley value concept based on a coalition game was adopted [7]. The coalition game theory can be defined by designating a value for each coalition game with a limited set of players N , $S \subseteq N$ to be a subset of $|S|$ players and a characteristic function $v : 2^N \rightarrow \mathbb{R}$ from the set of all possible coalitions of players to a set of players that satisfies $v(\emptyset) = 0$ where (\emptyset) is an empty set. This function determines each player's contribution to the outcome, and the game can be called a profit game or value game.

The profit game or value game can be adapted to the proposed QXAI framework to determine players (features) contributing to the prediction capacity of a trained deep learning model. To attribute a value to the contribution of each feature, the Shapley value concept can be adapted to explain the contribution in terms of expected marginal contribution. Shapley values assume that all the features contribute to the outcome, and the amount that each feature x_j contributes in a coalition game (v, N) is shown in Equation 5.

$$\varphi_{x_j}(v) = \sum_{S \subseteq N \setminus \{x_j\}} \frac{|S|! (n - |S| - 1)!}{n!} (v(S \cup \{x_j\}) - v(S)) \quad (5)$$

where the sum extends over all subsets S of N not containing feature i and n is the total number of features.

The above Equation 5 can further break-down to have individual feature contribution as $v(S \cup x_j) - v(S)$. The characteristic function $v(S)$ can be calculated by using Kernel SHAP.

$$\varphi_{x_j}(v) = \frac{1}{n!} \sum_R [v(P_{x_j}^R \cup \{x_j\}) - v(P_{x_j}^R)] \quad (6)$$

where the sum iterates over all $n!$ orders R of the features and $P_{x_j}^R$ is the set of features in N which precedes the order R .

In simple terms, Shapley of a feature x_j can be defined as below, Equation 7:

$$\varphi_{x_j}(v) = \frac{1}{n} \sum_K \frac{\varphi(x_j)}{Z} \quad (7)$$

Where n is a number of features, $\varphi(x_j)$ is marginal contribution of feature x_j to coalition, K is coalitions excluding x_j , Z is a number of coalitions excluding x_j .

Shapley proposed four conditions (or axioms) below that must be satisfied to have fair contribution of features to a prediction. Equations 5,6 obey these conditions while estimating the contribution value of each feature.

- The summation of Shapley values of all agents equals the value of the total coalition.
- All features have a fair chance to participate in a prediction outcome by including in all permutations and combinations of the features.
- If a participated feature x_j contributes nothing to a prediction outcome, then zero value is attributed to the feature's contribution.

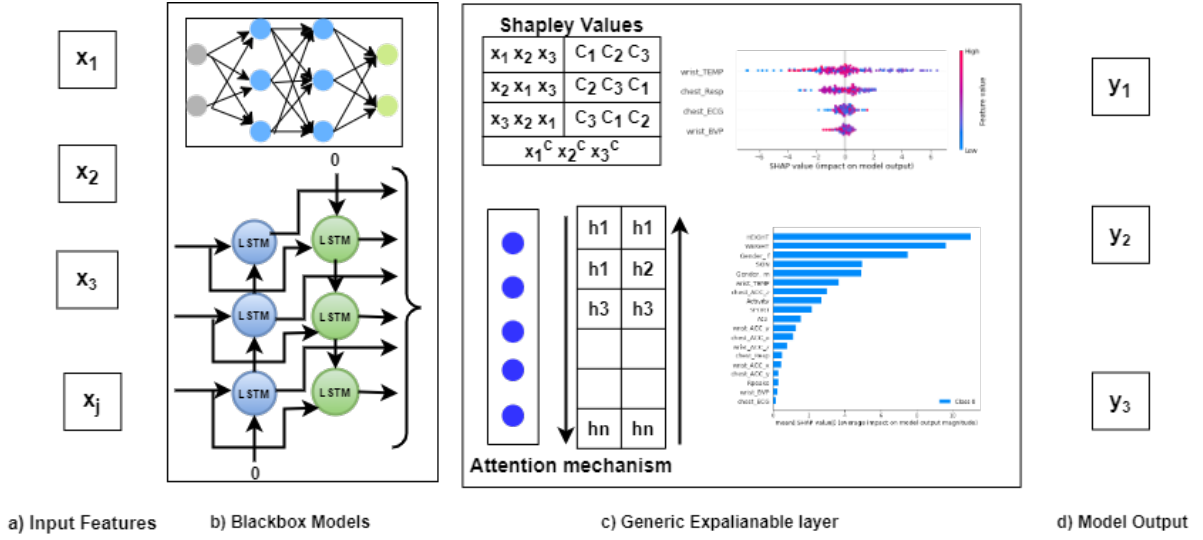


Fig. 1: Explainable AI (QXAI) framework

- For any pair of predictions v, w : $\varphi(v + w) = \varphi(v) + \varphi(w)$ in which the values are based on additive property $(v + w) = v(S) + w(S)$ for all subsets S .

B. Kernel SHAP

Kernel SHAP is a model-agnostic method from the combination of classical Shapley values discussed in Equations 5,6 and local explainable model-agnostic explanations (LIME) to approximate SHAP values. Instead of retraining models with a subset of features $|S|$, the full model f can be used which is already trained, while replacing missing features with marginalized features. Considering an instance with three features x_1, x_2, x_3 and following Equation 8 estimates a partial model with x_3 being missed. However, $p(x_3)$ is still required to approximate the missing x_3 feature. To address this, a custom proximity function π from LIME as shown in Equation 9 and SHAP similarity kernel equation 10 can be used.

$$f_{x_1, x_2}(x_1, x_2) \rightarrow \int f(x_1, x_2, x_3) p(x_3) dx_3 \quad (8)$$

$$\pi_x^{LIME}(z) = \exp(-D(x, z)^2 / \sigma^2) \quad (9)$$

$$\pi_x^{SHAP}(z') = \frac{(p-1)}{\binom{p}{|z'|} |z'| (p-|z'|)} \quad (10)$$

Equation 9 penalizes the distance between sample points and the original features' data, for which explainability is being estimated. In Equation 10, coalitions with a number of features that are far from 0 and p will be penalized. The equation adds more weight to coalitions with a small set of features or almost all the features to highlight the independent behavior of the feature or the impact of the features in interaction with others. The choice of this SHAP similarity kernel is based on three properties of additive feature attribution

methods local accuracy, missingness, and consistency [7]. In this study, Kernel SHAP is used to estimate the contributions of each feature x_j value to the prediction. It consists of five steps: 1) Sample coalitions with features and without features. 2) Get prediction for each sample coalition by first converting to the original feature space and applying the machine learning model. 3) Estimate the weight for each coalition with the SHAP kernel. 4) Fit the weighted linear model. 5) Return Shapley value $\varphi_{x_j}(v)$ and the coefficients of the model.

C. Attention Mechanism

The attention mechanism is a widely adopted concept in Natural Language Processing (NLP) tasks like neural machine translations and extracting the cause-effect of input features to model output [16, 17]. The attention mechanism predicts the outcome with better accuracy because its cognitive capability can enhance certain parts of important input data for deep learning model training. The idea of using the attention mechanism to model explainability is to identify the weights being assigned to each input feature in predicting the outcome. This assists in decoding the importance of each feature and enables human explanation of the cause-effect of the input features.

An attention layer added to a deep learning model can mimic the cognitive capability of the attention mechanism. Given a set of input features N , x_j is a feature value, with $j = 1, \dots, N$ to predict an output value y . A Bidirectional Long Short-Term Memory (BiLSTM) model can generate vector representations of the input features, such as h_j with, $j = 1, \dots, N$ based on the forward and backward hidden states in the deep learning model. A generic encoder-decoder model focuses on the last state of the encoder LSTM model and uses it as a context vector. This would cost the information loss of previous states. Attention acts as an interface between the encoder and decoder states of the BiLSTM model and provides a context vector to the decoder with information

from every encoder's hidden states. For each prediction value y , a context vector c is generated using the weighted sum of the vector representations, as shown in Equation 11. The weights α_j are computed using a softmax function as shown in Equation 12. The output score e_j is calculated in a feedforward neural network described by a function f to capture alignment between input feature x_j and output y . The input features are then multiplied (dot product) with (w_j+B) where w_j is weight and B bias followed by a tan hyperbolic function to estimate the score e_j as shown in Equation 13

$$c = \sum_{j=1}^N \alpha_j h_{x_j} \quad (11)$$

$$\alpha_j = \text{softmax}(e_j) = \frac{\exp(e_j)}{\sum_{j=1}^N \exp(e_j)} \quad (12)$$

$$e_j = f(x_j, h_{x_j}) = \tanh(x_j \cdot (w_j + B)) \quad (13)$$

For input features x_1, x_2, x_3, x_4 , let the weights α_j be, [0.2, 0.4, 0.6, 0.1] then the context vector would be as shown in equation 14. This can assist in estimating the importance of each input feature in the context vector, which will be fed to the decoder network for model predictions.

$$c = 0.2 \times x_1 + 0.4 \times x_2 + 0.6 \times x_3 + 0.1 \times x_4 \quad (14)$$

D. Global and Local explanation

Two different forms of explanation perspectives such as global explanation and local explanation are proposed in this study. The global explanation can provide the contribution of each feature in the prediction of vital sign. This is designed to assist clinicians by providing holistic information about the prediction and to identify which clinical factors or features need special attention. To estimate the global importance of the features in the prediction, the absolute Shapley values calculated from Equation 5 are averaged for each feature across the data, as shown in Equation 15. Based on this calculation, the features can have their importance sorted in descending order.

$$I_i = \frac{1}{n} \sum_{i=1}^n |\varphi_i| \quad (15)$$

Although feature importance can provide an overview of all selected features' importance towards a prediction, it cannot uncover the correlation of the features with a target variable and estimate contributing and non-contributing data points of a feature. This, however, can be achieved by using Shapley values of each feature on a summary plot showing the level of positive and negative contribution to a target variable.

In the case of local explanation, vital signs prediction at each time step can be decrypted. This can summarize features that are aiding the patient's health in terms of vital signs and can enable personalized monitoring, which is critical in healthcare applications. The Shapley values of each feature can be positive or negative, and each value is considered a force that either increases or decreases the prediction value.

This helps to explain individual features that are forcing the prediction value to either increase or decrease. The local explanation concept can be applied to an individual record in a prediction or a group of records related to a specific subject or activity.

Algorithm 1 Feature contribution estimation

Require: a set of features $\mathcal{F} = \{1, 2, \dots, N\}$; a set of deep learning models $\mathcal{M} = \{m_1, m_2\}$ where m_1 is without attention and m_2 is with attention; an input dataset D

Ensure: Contributions of the features $\mathcal{F} = \{1, 2, \dots, N\}$ in the form of Shapley values and attention weights;

1: Split dataset: $D = D^{train} \cup D^{test}$

Global explanation

2: $m_1^{train} \leftarrow D^{train}$

3: $m_1^{test} \leftarrow D^{test}$

4: $Shapley_values \leftarrow \text{kernelshap}(m_1^{train}, D^{test})$

5: $m_2^{train} \leftarrow D^{train}$

6: $m_2^{test} \leftarrow D^{test}$

7: $attention_weights \leftarrow \text{model.attention_weights}()$

Local explanation

8: **for** d in D **do**

9: $Shapley_values \leftarrow \text{kernelshap}(m_1^{train}, d)$

10: $attention_weights \leftarrow m_2.attention_weights()$

11: **end for**

E. QXAI Algorithm

The proposed QXAI framework comprises two deep learning model approaches, one with model attention and the second without. The framework can be implemented with the Algorithm 1 and can be adapted to execute global and local explanations. In Algorithm 1, line 1 splits the input data into test and train sets to train and evaluate the deep learning models. Lines 2-7 present the global explanation using kernel SHAP and attention layer weights. Lines 2-4 train a deep learning model without an attention layer and pass it to the kernel SHAP explainer to extract Shapley values of the input features. Lines 5-7 present the attention-based deep learning model and extracts the attention layer weights, thus defining input feature importance. Lines 8-11 present the local explanation for each input record d from data D .

V. EXPERIMENT

The two key aspects of an explainable AI framework are the understanding phase and the explaining phase [18]. The former is concerned with improving models during training by interpreting critical features and building robust models, while the latter involves deploying and providing human-readable explanations to end users. Striking a balance between model performance and explainability is always a challenge in AI applications. In AI applications, there is always a trade-off between model performance and explainability [19]. According to Zacharias et al. [20], the preprocessing stage, specifically feature selection, has been overlooked in explainable AI applications and requires attention. The importance of each feature to the outcome can be used for semantic labeling

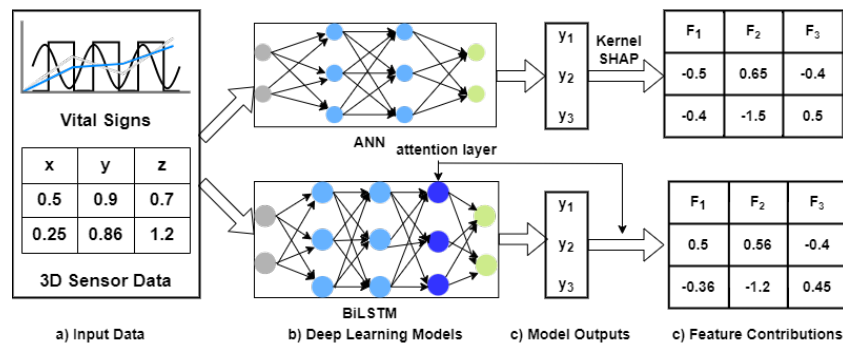


Fig. 2: Experimental Design

and to improve cognitive understanding, as it provides positive framing and direction (positive or negative contribution).

To address the limitations of explainable AI, the proposed QXAI framework in this study focuses on feature selection and provides local and global explanations through post-hoc models and intrinsic weights. The study evaluates feature importance in the QXAI framework, which can reduce dimensionality, improve cognitive understanding, and help with decision-making. Good explanations are crucial for making informed decisions, especially in dynamic domains like healthcare. In addition to the feature importance step in explainability, this study further breaks down the explainability into local and global explanations for each supervised learning task in classification and regression. Local and global explanations help in understanding the positive or negative contribution of features at the model and individual levels. The study used publicly available benchmark datasets for evaluation. Figure 2 illustrates the experimental design of the proposed framework.

A. Datasets

- **PPG-DaLiA** [21]: This dataset from 15 subjects comprised physiological and motion data while performing a wide range of activities under close to real-life conditions. The collected data were from both a wrist-worn (Empatica E4) and a chest-worn (RespiBAN) device. The dataset consists of 11 attributes, including 3-dimensional (3D) acceleration data, electrocardiogram (ECG), respiration, blood volume pulse (BVP), electrothermal activity (EDA), and body temperature.
- **MHEALTH** [22]: This dataset comprises the body motion of ten volunteers while performing 12 physical activities recorded from three sensors at the chest, left ankle, and right lower arm. There were 21 independent attributes including acceleration, gyroscope, and magnetometer of the three sensors. A dependent variable classifying the 12 activities was based on the sensor data.

B. Data Modelling

Datasets consisted of preprocessed raw data from the sensor's signal and features were stored in different CSV files. In

this step of data preparation, the dataset was further preprocessed to have a single structured file with a set of features for each subject. The datasets were prepared for two different tasks: regression and classification. The regression task was to predict the heart rate of the subjects based on their sensor readings. The classification task was to classify the physical activities of the subjects based on their motion data recorded from three axes of sensors. The physical activities label was preprocessed to have a multi-label classification. Each of these datasets was split into an 80:20 ratio for 80% of data for training and 20% of data for testing.

In this study, two deep learning models artificial neural networks (ANN), and Bidirectional LSTM (BiLSTM) models were adopted. The ANN model was configured with an input layer, hidden layers, and an output layer. The traditional activation function rectified linear unit (ReLU) has a limitation of defining negative inputs to zero which deactivates the nodes or neurons. Considering the negative values in 3D sensor data, the ANN model used the activation function LeakyReLU in input and hidden layers to avoid the zero input values of the negative attributes. The output layer was configured with the traditional activation function ReLU to predict the target variable heart rate greater than zero based on the activation function property. The loss function used for the regression study was mean absolute error, which also acted as a performance metric for the model. For the classification task, binary cross entropy acted as a loss function along with metrics like accuracy. The Adam adaptive optimizer [23] was chosen for the model for its quick computational time, it requires fewer parameters for tuning compared to other optimizers. The attention mechanism discussed in the proposed framework was added to the BiLSTM model, which has encoder and decoder states to generate vector representations. The preprocessed data was fed to the attention-based BiLSTM model and extracted the attention layer weights. This determined the input feature importance in the deep learning model prediction.

The datasets in this study were created by preprocessing raw data from sensor signals and storing the features in separate CSV files. These datasets were then combined into a single structured file for each subject, with separate datasets prepared for regression and classification tasks. The regression task involved predicting the subject's heart rate based on sensor

TABLE I: Implementation details

	Regression		Classification	
	Shapley Values	Attention Mechanism	Shapley Values	Attention Mechanism
Models	ANN	BiLSTM	ANN	BiLSTM
No of Layers	5	4+attention layer	5	4+attention layer
Activation Functions	relu, sigmoid	relu, Softmax	LeakyReLU, Sigmoid	relu, Softmax
Optimizers	Adam		Adam	
loss Functions	mean_absolute_error		binary_crossentropy	
Epochs	100		100	
Batch Size	64		64	

readings, while the classification task involved categorizing the subject’s physical activities using motion data from three axes of sensors. The datasets were split into 80% for training and 20% for testing. Two deep learning models, ANN and BiLSTM, were used in this study as shown in Tab. I. The table presents implementation details of ANN and BiLSTM models in regression and classification tasks.

For regression tasks, the models use the Shapley Values and attention mechanism. The ANN model has 5 layers with the activation functions of relu and sigmoid. The BiLSTM model has 4 layers with an additional attention layer with the activation functions of ReLU and softmax. The optimizer used is Adam, and the loss function is mean_absolute_error. For classification tasks, the models also use ANN and BiLSTM architectures, Shapley Values, and attention mechanisms. The ANN model has 5 layers with the activation functions of LeakyReLU and sigmoid. The optimizer used is Adam, and the loss function is binary_crossentropy. The BiLSTM model has 4 layers with an additional attention layer. The activation functions used are ReLU and softmax. For both prediction and classification tasks, the models are trained for 100 epochs with a batch size of 64.

C. Traditional Models

By comparing the feature importance estimated using Shapley values and intrinsic weights of the attention mechanism with the traditional machine learning models, the explainability of the proposed framework was evaluated. The two deep learning models in the framework, ANN and BiLSTM, were also evaluated to ensure high performance and robustness with explainability. This allowed the study to evaluate the effectiveness of the framework in explainability without compromising model performance.

The proposed approach was evaluated with models with state-of-art performances. The deep learning models adopted in the proposed approach were compared with heart rate prediction and human activity recognition performances. The feature importance was compared with traditional machine learning models, which had the capability to produce feature

importance for prediction and classification results.

Prediction

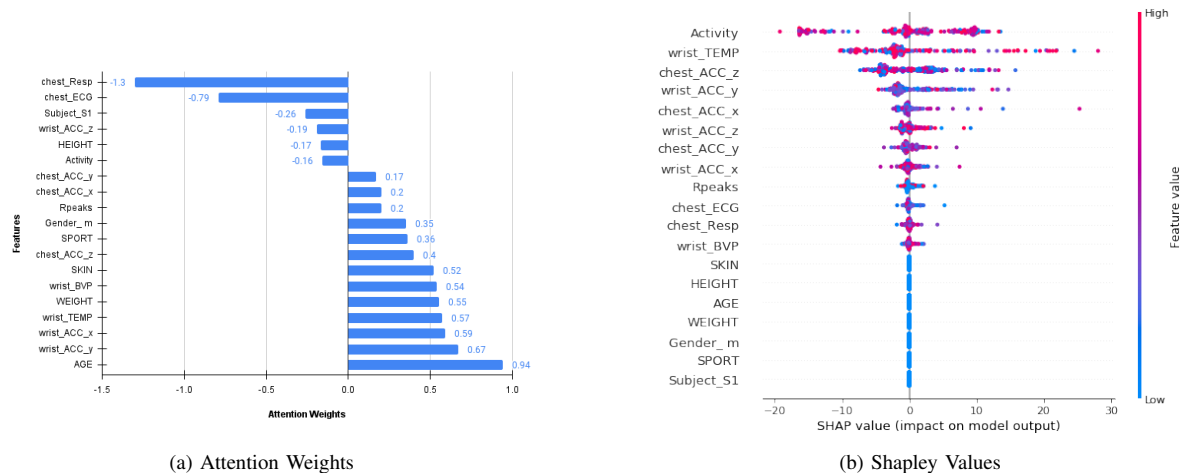
- Ni et al. [24] proposed context-aware sequential models to capture personalized fitness data and forecast heart rate to recommend suitable activities. The authors used a multi-layer perceptron model to forecast heart rate.
- Zhu et al. [25] proposed four LSTM models for an optimization training system to predict heart rate under three different types of exercises walking, rope jumping, and running. Three of the four LSTM models were used for heart rate prediction and one for human activity recognition.

Classification

- In a previous study, we proposed FedStack [26], a novel federated framework to classify patients’ physical activities. We adopted deep learning models such as CNN, ANN, and BiLSTM for the classification.
- Bozkurt et al. [27] compared deep learning model performance with traditional machine learning models for human activity recognition. Deep Neural Network (DNN) model achieved an accuracy of 96.81% and outperformed other models.

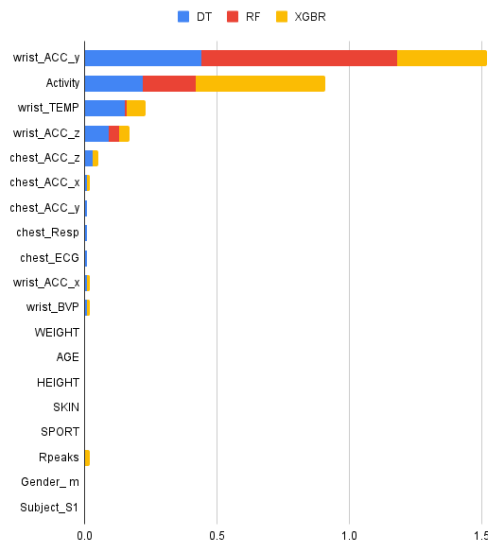
Feature Importance

- Li et al. [28] proposed an explainable machine learning model named cardiac arrest prediction index for early detection of cardiac arrest. The authors used the XGBoost model for the prediction and achieved an area under the receiver operating characteristic curve (AUROC) of 0.94.
- Gong et al. [8] used XGBoost and voting ensemble method combining random forest and logistic regression to predict acute kidney injury. For explanation, the SHAP technique was used to understand important predictors and relationships among the predictors.
- Ali et al. [29] proposed supervised machine learning algorithms such as Random Forest, Decision Tree, and KNN for heart disease prediction. Feature importance scores for each feature were computed with Decision Tree and Random Forest [30].



(a) Attention Weights

(b) Shapley Values



(c) Traditional Models

Fig. 3: Regression Model—Feature Importance Plots

D. Performance Metrics

Explainability is a multifaceted concept, and there is no single metric to measure it. The evaluation of explainability involves comparing the feature importance provided by different models, such as comparing the explanations of ANN and BiLSTM with those of traditional models. In this study, another two sets of performance metrics were adapted to evaluate deep learning models' prediction and classification results. For prediction, mean absolute error (MAE) and mean squared error (MSE) was used to evaluate the performance of the prediction model. Both metrics measure the deviation or difference of a predicted value from its actual value. For classification, a traditional confusion matrix was used to calculate precision, F1-Score, recall, and balanced accuracy metrics of deep learning results on multi-label classification.

VI. EXPERIMENTAL RESULTS AND ANALYSIS

In this section, we analyze the evaluation results of the proposed QXAI framework. The results are focused on explainability in terms of feature importance for positive framing, local explanations for semantic labelling to explain the positive or negative contributions of each input feature to the deep learning model's prediction, and global explanations that can explain a model's overall predictions with interactive plots. To address the trade-off between explainability and model performance in AI applications [31], the performance of the deep learning models, ANN and BiLSTM-attn, in the framework for both regression and classification tasks was evaluated and compared with those of traditional machine learning models.

TABLE II: QXAI Prediction Performance

Model	MAE	MSE
ANN	3.33	24.51
BiLSTM-attn	4.40	43.72
MLP [24]	4.71	47.95
LSTM [25]	5.54	69.03

A. QXAI in Regression Problem

The proposed QXAI approach was evaluated on its ability to predict heart rate based on sensor data and clinical indicators. Other vital signs retrieved from human subjects were in the PPG-DaLiA dataset. The two deep learning models ANN and attention-based BiLSTM proposed in the framework were trained on the data to predict the vital signs. The models' performance was compared with other traditional models shown in Tab. II. The ANN model performed better than the attention-based BiLSTM, MLP, and LSTM models with MAE and MSE of 3.33 and 24.51 respectively.

1) *Feature Importance*: Feature importance of input features was estimated using the proposed QXAI approach and compared with traditional machine learning model feature importance. The three feature importance plots shown in Fig. 3a, 3b, and 3c present attention weights retrieved from the BiLSTM model, Shapley values estimated from Kernel SHAP, and traditional model feature importance respectively. The y-axes in each subplot hold the input features, with x-axes showing the importance of each feature to the respective model's prediction. The large value of the x-axis determines the importance or contribution of a feature to model performance in predicting heart rate. Activity, chest, and wrist sensors data had high feature importance for heart rate prediction compared to other input vital signs like wrist_BVP, chest_Resp, and chest_ECG. The Shapley values plot 3b and attention weights plot 3a presented the negative dimensions of each feature's contribution.

2) *Explainability*: As discussed in Section IV, global and local explanations both contribute to presenting a patient's health status at different levels. The local explanation assists the clinician to explain the health status at a particular time step of patient monitoring. Fig. 4a presents feature contribution to the ANN model label prediction for a selected random record. The randomly selected record is of a male subject aged 25 years, height 168 centimeters, weight 57 kilograms with fitness level 5 on a scale 1-6 where 1 refers to them exercising less than once a month and 6 refers to 5-7 times a week. The subject's activity was measured during his lunch break, and his heart rate prediction was 71.24. The red highlighted features in Fig. 4a indicated a negative contribution and pushed the prediction value to the right (higher) side of the scale, whereas the blue features positively contributed and pushed the prediction value to the left (lower) side of the scale. This infers activity, wrist_ACC_y, and wrist_ACC_x features are negatively contributing and trying to decrease the heart rate value. The Rpeaks and wrist_TEMP features are balanced by increasing the heart rate to the expected value of 72.95. The SHAP values of each feature can be positive or negative. Sim-

ilarly, Fig. 4b presents a subject-level explanation of features' contribution to their heart rate prediction based on 200 records. The chart is related to a subject and presents each predicted value on the y-axis with its feature contribution spread on the x-axis in blue and red highlight. This is an interactive plot with dropdowns on the x-axis and y-axis changing and shows the impact of individual features on all 200 predictions. The plot is a screenshot of a prediction value of 107.9 in which the feature activity from wrist_ACC_x and wrist_TEMP are negatively contributing to the heart rate prediction.

TABLE III: QXAI Classification Performance

	Precision	Recall	F1-score	Balanced Accuracy
ANN	1	1	1	1
BiLSTM-Atten	0.92	0.78	0.77	0.88
CNN [26]	0.99	0.98	0.98	0.98
DNN [27]	0.97	0.97	0.97	0.97

B. QXAI in Classification Problem

The proposed QXAI approach was also used to explain the classification of human physical activities. Both the deep learning models ANN and attention-based BiLSTM were trained on the MHEALTH dataset. Model classification performance was compared to DNN and CNN, as shown in Tab. III. The ANN model had the best performance, with all evaluation metric values equalling 100%. CNN and DNN models also performed better than the attention-based BiLSTM model. The proposed framework disclosed the intrinsic weights of each feature in classification and post-hoc model explanations with Shapley values.

1) *Feature Importance*: The Shapley values and attention weights computed from the deep learning models determined the input feature importance in classifying human physical activities. Feature importance from the deep learning model was compared with feature importance from traditional machine learning models as shown in Fig. 5. The y-axes in all three subplots, 5a, 5b, and 5c refer to the 21 input features passed to the deep learning and the x-axes present the importance of a feature to model classification results. The attention-based BiLSTM model assigned more negative weights to all the input features. The sensor attributes at the wrist and ankle area were assigned with more weights in terms of magnitude to classify human physical activities as shown in Fig. 5a. The Shapley values plot 5b shows full body motion activities such as climbing stairs, jogging, walking, running, and jump front & back rely on left ankle sensor gyroscope data. The feature importance metrics from traditional machine learning models could not differentiate labels in their plot, as shown in Fig. 5c, but the results show that gyroscope data features contribute more to physical activity classification.

2) *Explainability*: The patients' physical activity classification can be explained in detail by breaking down the Shapley values with force plots as shown in Fig. 6a, 6b. The local explanation at each input record level can assist clinicians to explain physical activity classification and can explain which

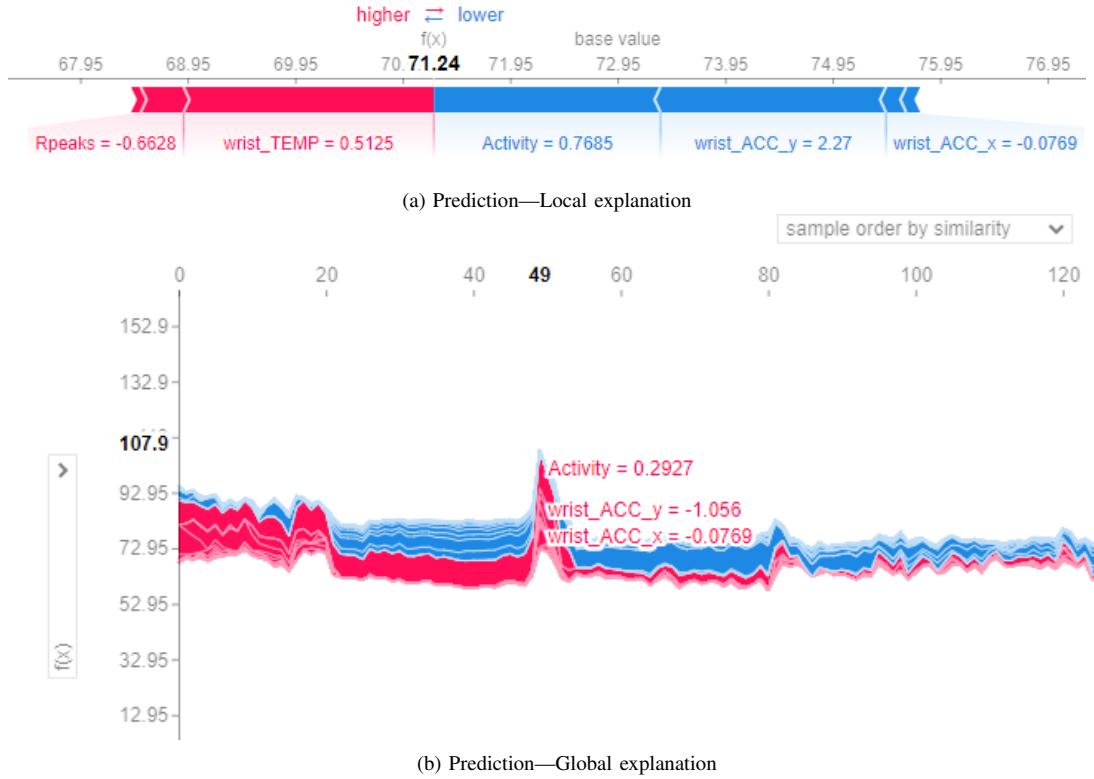


Fig. 4: Explanations for prediction: (a) Local explanation illustrating individual feature contributions. (b) Global explanation showing overall feature contributions.

sensor features are actually contributing to the classification. In plot 6a, the ANN model prediction probability of an arbitrarily selected record shows that the x, and z dimensions of the left ankle gyroscope try to push the model probability higher but the y-axis of the right lower arm and left ankle sensors are pushing the probability negatively according to Shapley values' feature importance. Similarly, Fig. 6b presents a subject-level interpretation of features that contribute to their physical activity classification based on 200 records. The chart is related to a subject and presents each predicted value on the y-axis with its feature contribution spread on the x-axis in blue and red highlights. This is an interactive plot with dropdowns on the x-axis and y-axis to change to see the impact of the individual feature on all 200 predictions. The plot is a screenshot of a predicted value 1 in which chest sensor acceleration positively contributes and left ankle and right lower arm sensor features negatively contribute to the heart rate prediction.

VII. MONTE CARLO APPROXIMATION

Feature contributions in model prediction can be estimated based on Shapley value computed using Equation 5 proposed in Section IV. These computations have an exponential time complexity and increase in number of features makes the Shapley value calculation unfeasible. In this study, Monte

Algorithm 2 Monte Carlo Approximation on Feature contribution estimation

Require: a set of features $x_j = \{1, 2, \dots, N\}$; a set of deep learning models $\mathcal{M} = \{m_1, m_2\}$ where m_1 is without attention and m_2 is with attention; input data \mathcal{D}

Ensure: Contribution of the features $x_j = \{1, 2, \dots, N\}$

- 1: marginal contribution $\phi_{x_j} \leftarrow \emptyset$
- 2: **for all** $x_j = \{1, 2, \dots, N\}$ **do**
- 3: $z \leftarrow$ random sample from \mathcal{D}
- 4: $x \leftarrow$ random sample from N
- 5: choose random permutation o of the feature x_j
- 6: $x : x_o = x_1, \dots, x_j$
- 7: $z : z_o = z_1, \dots, z_j$
- Build two new samples
- 8: with factor F :
- 9: $x_{+j} = (x_1, \dots, x_{j-1}, z_o = z_1, \dots, z_{j-1})$
- 10: without factor F :
- 11: $x_{-j} = (x_1, \dots, x_{j+1}, z_o = z_1, \dots, z_{j+1})$
- Compute marginal contribution of feature F :
- 12: $\phi_{x_j} \leftarrow m_1(x_{+j}) - m_1(x_{-j})$
- 13: **end for**
- 14: $\hat{\phi}_{x_j} \leftarrow \frac{1}{x_j} \sum_{m=1}^{x_j} \phi_{x_j}$

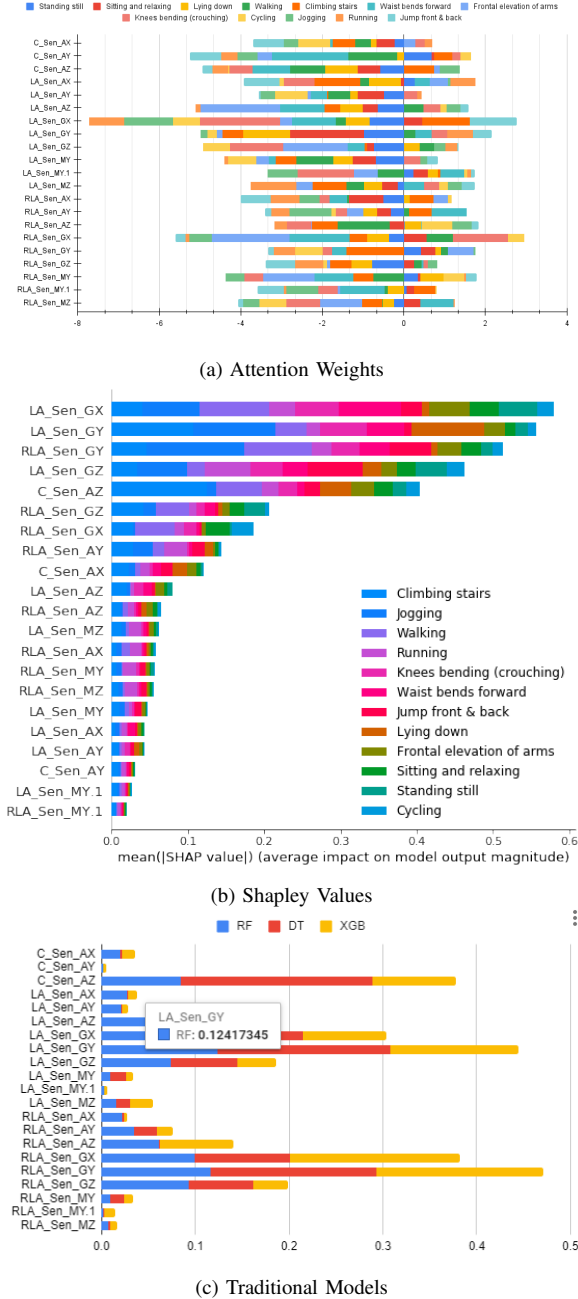


Fig. 5: Classification Model—Feature Importance Plots

Carlo approximation was adopted to calculate each feature contribution as shown in Equation 16. This approximation technique can extract Shapley values for each feature for both deep learning models. The results have been discussed in this section.

$$\varphi_i = \frac{1}{N} \sum_{n=1}^N (f(d_{+i}^m) - f(d_{-i}^m)) \quad (16)$$

where $f(\cdot)$ is the contribution of subset features. The d_{+i}^m and d_{-i}^m is the subset of with and without factor i in subset n features, respectively.

The implementation of the Monte Carlo approximation is presented in algorithm 2. Lines 3-7 obtain sampled data from the input data D . Lines 8-11 build new samples with or without consideration of a feature x_j . Line 12 calculates the marginal contribution ϕ_{x_j} of feature x_j . Lines 2-13 are a loop iterating to calculate the contribution of each feature one by one. Finally, line 14 calculates the Shapley value by averaging the outputs of multiple runs.

Monte Carlo approximation was applied to both deep learning models in the proposed QXAI framework. In ANN model prediction results, the approximation technique estimate wrist_TEMP is the most contributing feature in terms of magnitude, but the negative value shows that the feature is inversely proportional to the heart rate prediction. The other chest_ACC_z, chest_ACC_x, and chest_Resp features contributed positively towards the heart rate prediction. The attention weights from the BiLSTM model show that most of the input features are inversely proportional to the model output with negative values. The heat map shows that wrist_TEMP, wrist_ACC_z, and chest_ACC_z are the most contributing features to the heart rate prediction as shown in Fig. 7a. Similarly, the Monte Carlo approximation was applied to the deep learning classification models. The 3D axes of the sensor inputs were merged to have chest sensor acceleration, left ankle sensors' acceleration, gyroscope, and magnetometer, and right lower arm sensors' acceleration, gyroscope, and magnetometer as shown in Fig. 7b. The figure shows ANN model classification Shapley values for the consolidated input feature in the top heat map. The bottom heat map shows the attention-based BiLSTM model Shapley values. The full body activity like climbing stairs classification was more contributed by gyroscope data of the left ankle and right lower arm sensors and acceleration data of the chest sensor.

VIII. DISCUSSION

The research presented in this paper makes a significant contribution to the emerging field of explainable AI (XAI) in healthcare, particularly by addressing the challenge of interpretability in deep learning models for vital sign prediction and physical activity classification. The proposed Explainable AI for Quantitative data (QXAI) framework is noteworthy for its innovative approach that combines Shapley values and attention mechanisms, offering a comprehensive dual perspective on both post-hoc and intrinsic explainability. This discussion delves into the implications, strengths, limitations, and future directions of this study.

Implications and Contributions: The QXAI framework addresses a critical gap in healthcare AI by providing a solution to the 'black-box' nature of deep learning models. This is crucial as the explainability of AI models is increasingly becoming a requirement, especially in high-stakes fields like

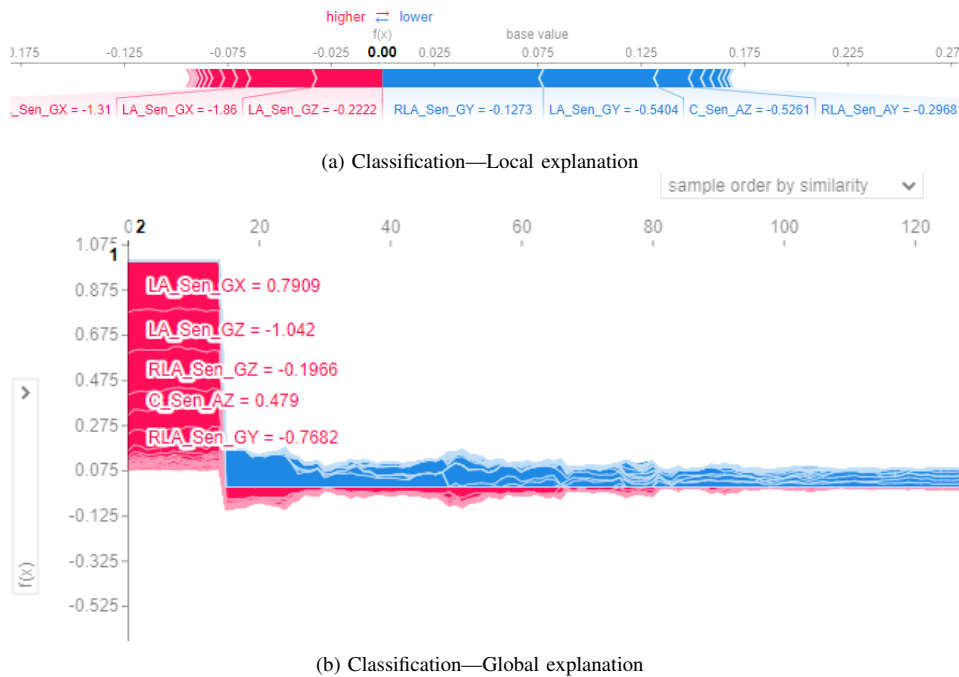


Fig. 6: Explanations for classification: (a) Local explanation illustrating individual feature contributions. (b) Global explanation showing overall feature contributions.

healthcare. The integration of Shapley values for post-hoc explainability and attention mechanisms for intrinsic understanding allows for a nuanced interpretation of AI decisions. This dual perspective of explainability not only enhances the trustworthiness of AI models but also makes them more practical and useful for clinicians. By enabling healthcare professionals to understand the reasoning behind AI-driven predictions, the framework facilitates informed decision-making in patient care.

Strengths of the Study: One of the major strengths of this study is the robust performance of the QXAI framework in both vital sign prediction and physical activity classification tasks. The superior performance, as compared to traditional models, highlights the potential of deep learning in enhancing healthcare diagnostics and monitoring. Furthermore, the comprehensive nature of the explainability approach employed in this study marks a significant advancement over existing methods that typically focus on either post-hoc or intrinsic explainability. The practical application of the framework, demonstrated through its effectiveness on real-world datasets like PPG-DaLiA and MHEALTH, underscores its potential for implementation in real healthcare settings.

Limitations and Future Directions: Despite its strengths, the study is not without limitations. The computational demands, particularly with large datasets due to the use of kernel SHAP, highlight the need for more efficient XAI algorithms. Additionally, while the framework shows promise, its generalizability across a broader range of healthcare scenarios remains to be tested. Future research should aim at scaling

the framework for different types of healthcare data and conditions. Another area for future improvement is the user-centric design of the framework. Tailoring explanations to be intuitive for healthcare practitioners, with varying levels of technical expertise, could enhance its clinical adoption. Moreover, the integration of the framework within existing clinical workflows and ensuring data privacy and ethical AI use are crucial considerations for future development.

IX. CONCLUSION

In healthcare applications, the explainability of machine learning model predictions or results is critical. This can assist clinicians in understanding the results to assist with clinical decisions that take appropriate steps for treatment. Existing deep learning models have a limitation in the explainability or interpretability of their results. The prediction or classification capacity of the proposed QXAI framework is outstanding compared to traditional machine learning models, with minimal knowledge of the healthcare domain knowledge to address the research problem. To utilize the advantage of the prediction capacity, this study proposed to adopt the Shapley values concept to vital signs prediction and decode global explanation at the overall population and local explanation at the subject level. However, the study was limited by the kernel SHAP method, which required significant memory and storage for large datasets. Future directions include incorporating more diverse feature inputs to enhance remote monitoring systems for clinical decision support.

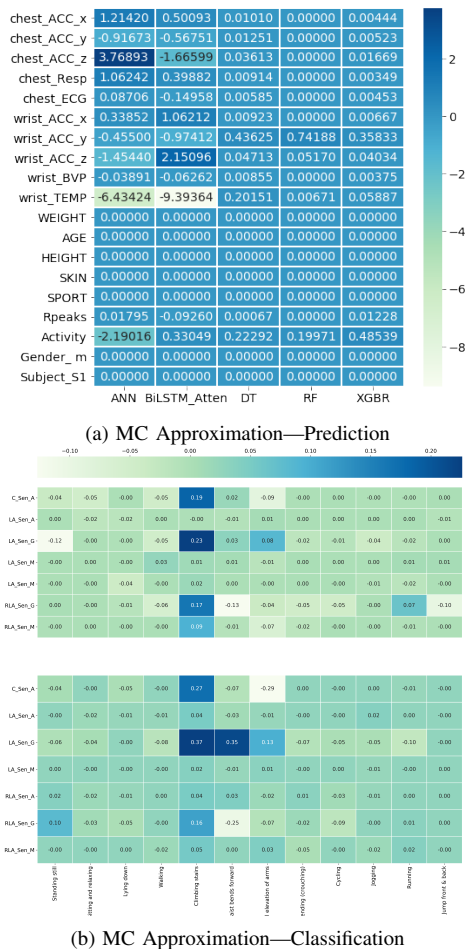


Fig. 7: Monte Carlo Approximation for Feature Importance Analysis in Prediction and Classification.

REFERENCES

- [1] L. P. Malasinghe, N. Ramzan, and K. Dahal, “Remote patient monitoring: a comprehensive study,” *Journal of Ambient Intelligence and Humanized Computing*, vol. 10, no. 1, pp. 57–76, 2019.
- [2] S. B. Asimwe, E. Vittinghoff, and M. Whooley, “Vital signs data and probability of hospitalization, transfer to another facility, or emergency department death among adults presenting for medical illnesses to the emergency department at a large urban hospital in the united states,” *The Journal of Emergency Medicine*, vol. 58, pp. 570–580, apr 2020.
- [3] X. Tao and J. D. Velasquez, “Multi-source information fusion for smart health with artificial intelligence,” 2022.
- [4] N. Prakash, A. Manconi, and S. Loew, “Mapping landslides on eo data: Performance of deep learning models vs. traditional machine learning models,” *Remote Sensing*, vol. 12, no. 3, p. 346, 2020.
- [5] S. M. Muddamsetty, M. N. Jahromi, A. E. Ciontos, L. M. Fenoy, and T. B. Moeslund, “Visual explanation of black-box model: similarity difference and uniqueness (sidu) method,” *Pattern recognition*, vol. 127, p. 108604, 2022.
- [6] A. Adadi and M. Berrada, “Explainable AI for healthcare: From black box to interpretable models,” in *Embedded Systems and Artificial Intelligence*, pp. 327–337, Springer Singapore, 2020.
- [7] S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” *Advances in neural information processing systems*, vol. 30, 2017.
- [8] K. Gong, H. K. Lee, K. Yu, X. Xie, and J. Li, “A prediction and interpretation framework of acute kidney injury in critical care,” *Journal of Biomedical Informatics*, vol. 113, p. 103653, Jan. 2021.
- [9] L. Wu, Y. Hu, X. Liu, X. Zhang, W. Chen, A. S. L. Yu, J. A. Kellum, L. R. Waitman, and M. Liu, “Feature ranking in predictive models for hospital-acquired acute kidney injury,” *Scientific Reports*, vol. 8, Nov. 2018.
- [10] R. ElShawi, Y. Sherif, M. Al-Mallah, and S. Sakr, “Interpretability in healthcare: A comparative study of local machine learning interpretability techniques,” *Computational Intelligence*, vol. 37, pp. 1633–1650, Nov. 2020.
- [11] R. Elshawi, M. H. Al-Mallah, and S. Sakr, “On the interpretability of machine learning-based model for predicting hypertension,” *BMC Medical Informatics and Decision Making*, vol. 19, July 2019.
- [12] I. Ilic, B. Görgülü, M. Cevik, and M. G. Baydoğan, “Explainable boosted linear regression for time series forecasting,” *Pattern Recognition*, vol. 120, p. 108144, 2021.
- [13] D. Barić, P. Fumić, D. Horvatić, and T. Lipic, “Benchmarking attention-based interpretability of deep learning in multivariate time series predictions,” *Entropy*, vol. 23, p. 143, Jan. 2021.
- [14] D. A. Kaji, J. R. Zech, J. S. Kim, S. K. Cho, N. S. Dangayach, A. B. Costa, and E. K. Oermann, “An attention based deep learning model of clinical events in the intensive care unit,” *PLOS ONE*, vol. 14, p. e0211057, Feb. 2019.
- [15] C. Chen and B. Li, “An interpretable channelwise attention mechanism based on asymmetric and skewed gaussian distribution,” *Pattern Recognition*, p. 109467, 2023.
- [16] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [17] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” 2014.
- [18] R. Dwivedi, D. Dave, H. Naik, S. Singhal, R. Omer, P. Patel, B. Qian, Z. Wen, T. Shah, G. Morgan, *et al.*, “Explainable ai (xai): Core ideas, techniques, and solutions,” *ACM Computing Surveys*, vol. 55, no. 9, pp. 1–33, 2023.
- [19] D. Gunning and D. Aha, “Darpa’s explainable artificial intelligence (xai) program,” *AI magazine*, vol. 40, no. 2, pp. 44–58, 2019.

- [20] J. Zacharias, M. von Zahn, J. Chen, and O. Hinz, "Designing a feature selection method based on explainable artificial intelligence," *Electronic Markets*, pp. 1–26, 2022.
- [21] A. Reiss, I. Indlekofer, P. Schmidt, and K. Van Laerhoven, "Deep ppg: large-scale heart rate estimation with convolutional neural networks," *Sensors*, vol. 19, no. 14, p. 3079, 2019.
- [22] O. Banos, R. Garcia, J. A. Holgado-Terriza, M. Damas, H. Pomares, I. Rojas, A. Saez, and C. Villalonga, "mHealthDroid: A novel framework for agile development of mobile health applications," in *Ambient Assisted Living and Daily Activities*, pp. 91–98, Springer International Publishing, 2014.
- [23] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [24] J. Ni, L. Muhlstein, and J. McAuley, "Modeling heart rate and activity data for personalized fitness recommendation," in *The World Wide Web Conference*, pp. 1343–1353, 2019.
- [25] Z. Zhu, H. Li, J. Xiao, W. Xu, and M.-C. Huang, "A fitness training optimization system based on heart rate prediction under different activities," *Methods*, vol. 205, pp. 89–96, 2022.
- [26] T. Shaik, X. Tao, N. Higgins, R. Gururajan, Y. Li, X. Zhou, and U. R. Acharya, "FedStack: Personalized activity monitoring using stacked federated learning," *Knowledge-Based Systems*, vol. 257, p. 109929, Dec. 2022.
- [27] F. Bozkurt, "A comparative study on classifying human activities using classical machine and deep learning methods," *Arabian Journal for Science and Engineering*, vol. 47, pp. 1507–1521, July 2021.
- [28] L. Yijing, Y. Wenyu, Y. Kang, Z. Shengyu, H. Xianliang, J. Xingliang, W. Cheng, S. Zehui, and L. Mengxing, "Prediction of cardiac arrest in critically ill patients based on bedside vital signs monitoring," *Computer Methods and Programs in Biomedicine*, vol. 214, p. 106568, 2022.
- [29] M. M. Ali, B. K. Paul, K. Ahmed, F. M. Bui, J. M. Quinn, and M. A. Moni, "Heart disease prediction using supervised machine learning algorithms: Performance analysis and comparison," *Computers in Biology and Medicine*, vol. 136, p. 104672, Sept. 2021.
- [30] S. Malakar, S. D. Roy, S. Das, S. Sen, J. D. Velásquez, and R. Sarkar, "Computer based diagnosis of some chronic diseases: A medical journey of the last two decades," *Archives of Computational Methods in Engineering*, pp. 1–43, 2022.
- [31] L.-V. Herm, K. Heinrich, J. Wanner, and C. Janiesch, "Stop ordering machine learning algorithms by their explainability! a user-centered investigation of performance and explainability," *International Journal of Information Management*, vol. 69, p. 102538, 2023.

9.2 Summary

This chapter underscores the transformative impact of the QAXI framework in enhancing the explainability of AI models within healthcare. By providing clear insights into the contribution of individual features to model predictions, QAXI facilitates a deeper understanding of AI-driven decisions, fostering trust and transparency in clinical settings. The successful application of QAXI in vital sign prediction and physical activity classification exemplifies its potential to revolutionize patient monitoring by offering a more nuanced and interpretable AI-driven approach.

CHAPTER 10: PAPER 9 - EXPLORING THE LANDSCAPE OF MACHINE UNLEARNING: A COMPREHENSIVE SURVEY AND TAXONOMY

10.1 Introduction

This chapter presents an in-depth survey on Machine Unlearning (MU), an emerging field addressing the need to efficiently remove specific data or knowledge from trained machine learning models. MU is driven by increasing privacy concerns, regulatory requirements, and the dynamic nature of data, making it imperative to develop models capable of "forgetting." The introduction explores various MU techniques, challenges, and their implications in ensuring model compliance with privacy standards and enhancing model adaptability.

Exploring the Landscape of Machine Unlearning: A Comprehensive Survey and Taxonomy

Thanveer Shaik, Xiaohui Tao, Haoran Xie, Lin Li, Xiaofeng Zhu, and Qing Li

Abstract—Machine unlearning (MU) is gaining increasing attention due to the need to remove or modify predictions made by machine learning (ML) models. While training models have become more efficient and accurate, the importance of unlearning previously learned information has become increasingly significant in fields such as privacy, security, and ethics. This paper presents a comprehensive survey of MU, covering current state-of-the-art techniques and approaches, including data deletion, perturbation, and model updates. In addition, commonly used metrics and datasets are presented. The paper also highlights the challenges that need to be addressed, including attack sophistication, standardization, transferability, interpretability, training data, and resource constraints. The contributions of this paper include discussions about the potential benefits of MU and its future directions. Additionally, the paper emphasizes the need for researchers and practitioners to continue exploring and refining unlearning techniques to ensure that ML models can adapt to changing circumstances while maintaining user trust. The importance of unlearning is further highlighted in making Artificial Intelligence (AI) more trustworthy and transparent, especially with the increasing importance of AI in various domains that involve large amounts of personal user data.

Index Terms—Machine Unlearning, Privacy, right to be forgotten, Federated Unlearning, Graph Unlearning

I. INTRODUCTION

Machine learning (ML) refers to the process of training an algorithm to make predictions or decisions based on data [1]. ML has become increasingly important in applications such as health, higher education, and other relevant domains. In healthcare, ML models can be used to predict patient outcomes, identify high-risk patients and personalize treatment plans [2]. For higher education, ML has been used to improve student outcomes and enhance the learning experience, or even used to analyze student data and predict their online class engagement [3].

In ML, an algorithm is trained on a dataset to learn patterns and relationships in the data. Once the algorithm has been trained, it can be used to make predictions on new data. Thus, the goal of ML is to create accurate models that can generalize well onto new data [4]. On the other hand, machine

unlearning (MU) is the process of removing certain data points or features from a trained ML model without affecting its performance [5]. MU is a relatively new and challenging field of research that is concerned with developing techniques for removing sensitive or irrelevant data from trained models. The goal of MU is to ensure that trained models are free from biases and sensitive information that could lead to negative outcomes [6].

MU was first introduced by Cao et al. [7], who recognized the need for a “forgetting system” and developed one of the initial unlearning algorithms called *machine unlearning*. This approach efficiently removes data traces by converting learning algorithms into a summation form, which can help counter data pollution attacks. The increasing need for regulatory compliance with modern privacy regulations led to the creation of MU, which involves deleting data not only from storage archives but also from ML models [8]. Existing studies update the model weights for unlearning using either the whole training data, a subset of training data, or some metadata stored during training [9]. Although strict regulatory compliance requires the timely deletion of data, there are instances where data about the training process may not be available for unlearning purposes. Companies and organizations commonly employ user data to train ML models, but legal frameworks like GDPR, CCPA, and CPPA demand that user data be erased when requested [10]. The question is whether merely deleting the data is sufficient, or if the models trained on this data should also be adjusted [11]. However, straightforward techniques like retraining models from scratch or check-pointing can be computationally costly and require significant storage resources [12]. With MU, we can modify models to exclude specific data points more efficiently [13].

Following the introduction of ML, the discussion on MU emerges as an important counterpart, especially with the current privacy regulations. The surveys by Nguyen et al. [9] and Xu et al. [14] offer initial insights into MU, but they differ in depth and scope compared to this proposed survey. Nguyen et al. [9] provide a broad overview of machine unlearning with a focus on privacy, while Xu et al. [14] mainly categorize existing unlearning solutions. Unlike these surveys, the proposed survey explores detailed challenges such as attack sophistication and lack of standardization, going beyond the general challenges discussed by Nguyen et al. [9]. Moreover, this survey introduces a structured discussion through Data-centric and Model-centric approaches, enriching the technical dialogue beyond the basic methods presented in the referenced surveys. The addition of Machine Unlearning Evaluation Metrics in this survey creates a solid framework for assessing

Thanveer Shaik and Xiaohui Tao are with the School of Mathematics, Physics and Computing, University of Southern Queensland, Queensland, Australia (e-mail: Thanveer.Shaik@usq.edu.au, Xiaohui.Tao@usq.edu.au).

Haoran Xie is with the Department of Computing and Decision Sciences, Lingnan University, Tuen Mun, Hong Kong (e-mail: hrxie@ln.edu.hk)

Lin Li is with the School of Computer and Artificial Intelligence, Wuhan University of Technology, China (e-mail: cathylin@whut.edu.cn)

Xiaofeng Zhu is with the University of Electronic Science and Technology of China (e-mail: seanzhuxf@gmail.com)

Qing Li is with the Department of Computing, Hong Kong Polytechnic University, Hong Kong Special Administrative Region of China (e-mail: qing-prof.li@polyu.edu.hk).

unlearning techniques, an aspect not well-covered in the referenced surveys. This survey also expands the discussion to domain-specific unlearning scenarios like Natural Language Processing (NLP), Computer Vision (CV), and Recommender Systems, broadening the survey’s scope and applicability. Further, it delves into detailed solutions to the identified challenges and provides an expansive outlook on future directions, thus offering a more thorough exploration of the MU domain. Through these aspects, this survey not only addresses the identified challenges in a structured manner but also suggests potential solutions, advancing the MU discourse beyond the foundational inquiries made by Nguyen et al. [9] and Xu et al. [14], and aligning with the modern privacy-centric regulatory frameworks mentioned in the previous discourse on MU.

Various techniques have been suggested for managing user data deletion requests, such as optimization, clustering, and regression methods [15]. By conducting a comprehensive survey of existing literature on managing user data deletion requests we can identify gaps and trends in the field, which in turn will both guide future research and provide insights for the organizations handling such requests. In this study, we address the following research questions:

- 1) What are the most effective techniques for unlearning data from ML models?
- 2) How can the impact of unlearning on model performance be measured and evaluated?
- 3) What are the challenges in MU, and how can these challenges be addressed?

The contributions of this study are:

- A comprehensive and up-to-date taxonomy about the emerging field of MU, including an explanation of its importance and potential applications;
- A detailed taxonomy of the various techniques and approaches that have been developed for unlearning data from ML models, such as data deletion, data perturbation, and model update techniques;
- A discussion of different evaluation methods for assessing the effectiveness of MU techniques, such as measuring the degree of forgetting or their impact on model performance;
- A taxonomy of several key challenges in the field of MU, including attack sophistication, standardization, transferability, interpretability, training data, and resource constraints;
- Finally, a discussion of the potential benefits of MU and its future directions in natural language processing (NLP), computer vision, and recommender systems.

The remainder of the paper is organized as follows. Section II outlines the aims and objectives of MU. In Section III, we delve into data deletion, data perturbation, and model update techniques in greater depth. Section IV details the evaluation metrics of MU, while Section V discusses the challenges associated with the field and proposes potential solutions. In Section VI, we explore the future directions of MU in NLP, computer vision, and recommender systems. Finally, Section VII concludes the paper.

II. OVERVIEW OF MACHINE UNLEARNING

The “right to be forgotten” is an evolving concept, emphasizing the need for individuals to have personal data expunged from online platforms in specific circumstances [16], [17]. While its definition and classification as a human right remain contested, countries like Argentina, the European Union (EU), and the Philippines are inching towards regulatory frameworks around this proposal ¹.

Past information, even if outdated or resolved, can significantly impact an individual’s present reputation. An illustrative case is the 2018 incident involving Disney’s dismissal of James Gunn over previously tweeted controversial content [18]. Similarly, removal requests, such as those lodged against Google, underscore the intricate challenges surrounding data persistence on the Internet and the growing demands for data erasure [19].

MU sits at the nexus of these discussions, especially within the broader artificial intelligence (AI) landscape. It seeks a harmonious equilibrium between retaining model efficacy and adapting to the shifting data paradigms, regulatory mandates, and ethical considerations [20]. This practice, termed as “selective amnesia” [21], has myriad applications and goals, spanning various dimensions. We derive a definition of MU, building upon comprehensive literature reviews:

General Definition: *Machine unlearning is the recalibration of machine learning models by selectively discarding specific data points, patterns, or predictions. While traditional machine learning accentuates pattern recognition and prediction from data, machine unlearning modulates these patterns or predictions in response to data shifts, privacy imperatives, or model performance enhancements.*

MU aims to achieve multiple objectives, including:

Privacy-Preserving Adaptation: In an era of rigid data privacy norms, MU aids models in adhering to directives such as the “right to be forgotten”. This usually implies excluding specific data instances linked to personal or sensitive details rather than erasing entire learned patterns [22]–[29].

Accuracy and Fairness Enhancement: MU can rectify biases in ML models, enhancing fairness and accuracy by unlearning certain patterns or data [30]–[33].

Adaptive Learning: With evolving data landscapes, MU ensures models remain pertinent by shedding obsolete or irrelevant information [34]–[36].

Reducing Computation Costs: MU offers computational efficiency by updating models based on data changes, avoiding the need for complete retraining and thereby saving resources [37]–[42].

While companies invest substantially in AI model training, rising regulatory cautions - especially from bodies in the EU and the U.S. - indicate potential data and model deletions. For instance, the U.S. Federal Trade Commission recently directed Paravision to delete data and associated models derived from inappropriately gathered facial photos ². Although retraining

¹<https://link.library.eui.eu/portal/The-Right-To-Be-Forgotten-A-Comparative-Study/tw0VHCyGcDc/>

²<https://www.wired.com/story/startup-nix-algorithms-ill-gotten-facial-data/>

remains the most straightforward strategy for data point removal, it is also a resource-intensive one, as evidenced by the significant costs incurred for training models like GPT-3 [43]. This emphasizes the need for cost-effective solutions in the realm of MU.

Balancing privacy and the right to expression is paramount to prevent misuse of the “right to be forgotten” [44]. With emerging technologies like blockchain, this balance becomes even more intricate. As data privacy concerns surge, companies (exemplified by Google’s recent policy expansion) are taking proactive steps. However, the actual challenge lies in ensuring that AI models are cleansed appropriately once data points are removed to prevent biases or sensitive information propagation. MU, although intricate, is pivotal in this endeavor. As data privacy regulations intensify, the role of MU in fostering transparent and ethical AI models will become even more critical.

III. TECHNIQUES AND APPROACHES

This section discusses the taxonomy of MU techniques, which are categorized as data-centric unlearning approaches, model-centric unlearning approaches, federated unlearning approaches, and graph unlearning, as shown in Fig. 1. In this section, the first research question “*What are the most effective techniques for unlearning data from ML models?*” will be addressed.

A. Data-Centric Unlearning Approaches

1) Data Deletion

The concept of data deletion within the domain of MU presents a complex and multifaceted challenge, drawing significant attention from researchers. This process transcends the simple removal of data points, aiming to eradicate their influence on machine learning models while preserving the models’ utility and ensuring data privacy.

Chourasia et al. [45] highlight the limitations of straightforward deletion or retraining approaches, noting that subtle influences or traces of the deleted data might persist in retrained models, thereby compromising privacy. They propose leveraging noisy gradient descent to illustrate the interplay between deletion privacy and differential privacy, emphasizing the criticality of existing data’s privacy for the authentic deletion of removed data. This underscores a pivotal aspect of MU: genuine data deletion necessitates deep, algorithmic interventions beyond mere model retraining.

Conversely, Garg et al. [46] provide a structured framework for data deletion, emphasizing the need for a technically sound foundation for MU methodologies. Their work not only showcases the intricacies involved in ensuring genuine data deletion but also delves into the legal and ethical dimensions, proposing a blueprint that could steer future unlearning methodologies. This emphasizes the importance of harmonizing technical approaches with regulatory and ethical considerations in MU.

Exploring the vulnerabilities inherent in the data deletion process, Gao et al. [47] shed light on potential risks associated with MU. They advocate for a Deletion Compliance framework, highlighting the possibility of exploiting the deletion process to reconstruct or infer deleted data, thus breaching

privacy. This underscores the necessity for unlearning methodologies to be resilient against such exploitations, ensuring both the genuineness and security of the deletion process.

Focusing on the operational perspective, Wang et al. [48] and Ginart et al. [49] discuss the computational and practical complexities involved in updating or ‘unlearning’ models efficiently in response to data deletion requests. They navigate the delicate balance between model utility, compliance with deletion requests, and computational resource management, advocating for unlearning methodologies that are efficient and effective.

Collectively, these studies navigate the intricate landscapes of MU and data deletion, highlighting the need to bridge the gaps between privacy, utility, operational efficiency, and regulatory compliance. They convey a unified message: data deletion in MU is a complex, multidimensional challenge that requires a cohesive, robust, and multifaceted approach. This approach must integrate algorithmic, operational, legal, and ethical considerations into a comprehensive methodology, underscoring the need for holistic research and development in MU to develop secure, efficient, and comprehensive data deletion methodologies.

2) Mitigating Data Poisoning

Data poisoning represents a formidable challenge in MU, where adversaries intentionally corrupt the training dataset with malicious data to degrade the model’s performance or induce biased decision-making. This nefarious activity is particularly problematic in privacy-centric systems or automated decision-making processes where integrity and accuracy are paramount.

Consider a dataset $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$. An adversary aims to compromise the model by introducing a poisoned data point (x', y') with the intention of causing the model to misclassify a specific target label y_{target} . The resulting poisoned dataset is denoted as:

$$D' = \{(x_1, y_1), (x_2, y_2), \dots, (x_i, y_{\text{target}}), \dots, (x_n, y_n)\} \quad (1)$$

where (x_i, y_{target}) represents the adulterated data point.

The objective is to minimize the model’s loss function $L(\theta; D')$ while ensuring the model’s accuracy on the original dataset D does not fall below a specified threshold Acc_0 :

$$\min L(\theta; D') \quad \text{subject to} \quad Acc(\theta; D) \geq Acc_0 \quad (2)$$

Data poisoning exploits vulnerabilities in data collection and processing, injecting data that, while appearing legitimate, is designed to bias or corrupt the model’s outputs.

Marchant et al. [50] propose a solution using projected gradient descent (PGD) to address data poisoning, underscoring the challenges of adhering to data protection regulations such as the right to erasure. They unveil a novel vulnerability in ML systems—poisoning attacks—that not only compromise accuracy but also resemble denial-of-service attacks by impeding the unlearning process.

Sun et al. [51] explore the threats posed by attackers leveraging federated learning (FL) to conduct poisoning attacks across various nodes. They introduce the Attack on Federated Learning (AT^2FL) framework, which employs systems-aware optimization techniques to discern and mitigate the effects

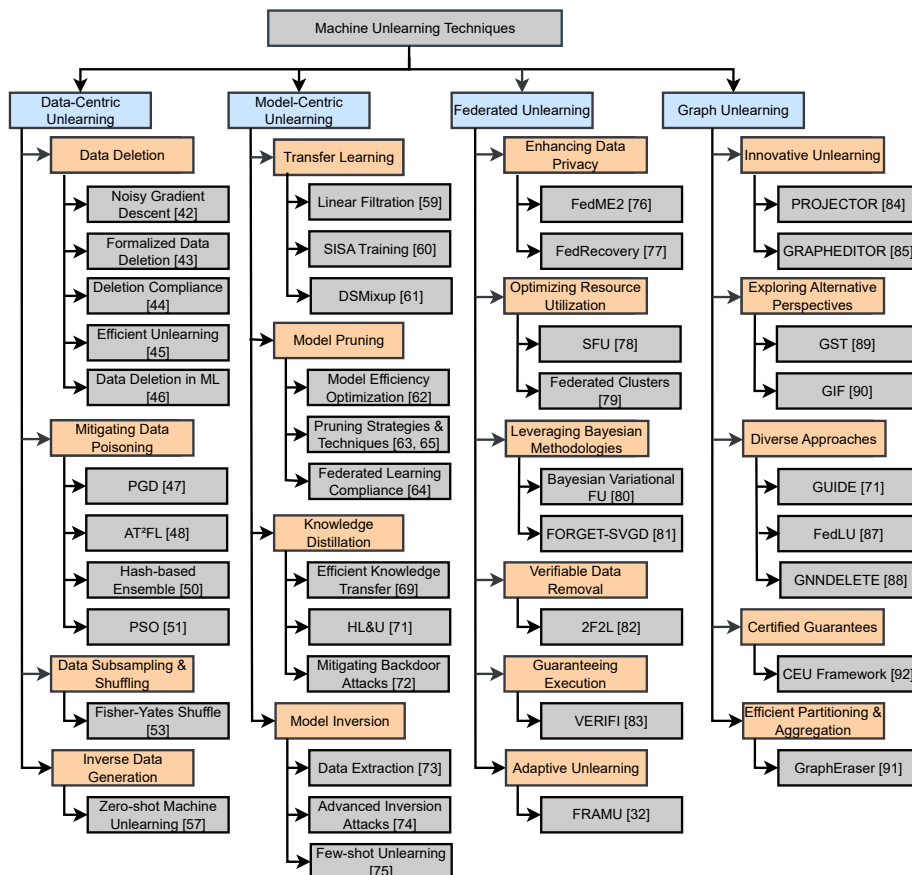


Fig. 1. Machine Unlearning - Taxonomy

of poisoned data. Their framework demonstrates enhanced resilience against both direct and indirect poisoning attacks in a federated multitask learning setting, highlighting its potential in safeguarding FL systems.

Furthermore, data poisoning is not limited to direct attacks but can also manifest through adversarial manipulations such as random label flipping and distance-based label flipping attacks. Yerlikaya et al. [52] assess the susceptibility of six ML algorithms to these adversarial tactics, noting variability in algorithm performance based on dataset characteristics. Anisetti et al. [53] evaluate the robustness of a hash-based ensemble approach against data poisoning in tabular datasets, demonstrating that even modestly sized ensembles can offer significant protection against poisoning attacks. This suggests that ensemble methods possess inherent strengths that can be harnessed to counteract adversarial attacks.

Maabreh et al. [54] propose the development of deep learning (DL) models optimized with the particle swarm optimizer (PSO), designed to perform effectively even when confronted with fake or poisoned data samples. However, they caution that PSO's efficacy may be compromised in scenarios where the dataset is heavily contaminated with malicious data.

These studies collectively emphasize the complexity and

multi-dimensional nature of mitigating data poisoning in MU. They advocate for a holistic approach that encompasses algorithmic innovations, robust optimization techniques, and the leveraging of ensemble methods to enhance the resilience of ML models against sophisticated poisoning attacks, ensuring the integrity and reliability of automated decision-making processes in the face of adversarial threats.

3) Data Subsampling and Shuffling

In the realm of MU, safeguarding the integrity of models and ensuring the privacy of user data necessitate innovative and robust techniques. Data subsampling and shuffling stand out as key strategies in this context, offering distinct but complementary approaches to mitigating risks associated with data management and model training.

Data Subsampling is particularly advantageous when dealing with extensive datasets or when computational resources are constrained. This technique involves selecting a random subset S from the original dataset $X = \{x_1, x_2, \dots, x_n\}$ with corresponding labels $Y = \{y_1, y_2, \dots, y_n\}$, thereby creating reduced training sets X' and Y' :

$$\begin{aligned} X' &= X \setminus S, \\ Y' &= Y \setminus S. \end{aligned} \quad (3)$$

Iterative application of subsampling with varying subsets S can diminish the model’s dependency on the initial training data, thereby bolstering its resilience, particularly against Membership Inference (MI) attacks, which aim to deduce individual data points’ presence in the training dataset [55].

Data Shuffling, on the other hand, aims to obfuscate sensitive information by randomizing the sequence of data points within the dataset, thus preserving its statistical properties while protecting individual data entries. For a labeled dataset X , shuffling reorganizes the data to form a new dataset X' with the same elements but in a different order. This shuffled dataset is then partitioned into training and validation sets, X_{tr} and X_{val} , respectively, upon which the model is trained and validated. Repeating this shuffling process ensures diverse model initializations and data presentations, contributing to the development of models that are more robust and generalizable.

When combined, *data subsampling and shuffling* provide a multi-layered defense mechanism for ML models. Subsampling addresses computational efficiency and reduces the risk of overfitting by minimizing the reliance on extensive or potentially biased datasets. Shuffling, by disrupting the data order, further complicates malicious attempts to reverse-engineer or compromise the dataset, enhancing data security.

The implementation of these strategies is facilitated by algorithms such as the Fisher-Yates shuffle [56] and utilities provided in ML libraries like Scikit-learn, which offer practical tools for executing data shuffling and subsampling efficiently. Integrating these techniques within a broader MU framework, especially alongside other strategies like data deletion and poisoning mitigation, could pave the way for a comprehensive solution that not only ensures the integrity of the learning and unlearning processes but also robustly secures both the data and the models involved. This integrated approach represents a significant stride towards realizing secure, efficient, and transparent MU paradigms, aligning with the overarching goals of user privacy protection and model reliability in the evolving landscape of machine learning.

4) Inverse Data Generation

Inverse Data Generation (IDG) has emerged as a pivotal technique within the domain of MU, particularly in its capacity to synthesize datasets that retain the statistical essence of the original data while meticulously excluding sensitive information. This technique leverages the prowess of generative models to fabricate data points that not only mirror the characteristics of the original dataset but also ensure the omission of privacy-compromising elements.

Generative Models in IDG: The cornerstone of IDG lies in the utilization of advanced generative models, notably Generative Adversarial Networks (GANs) [57], Variational Autoencoders (VAEs) [58], and Deep Belief Networks (DBNs) [59]. These models are adept at learning the intricate distributions of the original data and generating new instances that adhere to these learned distributions, thereby facilitating the creation of a sanitized dataset that closely resembles the original without encroaching upon sensitive information.

Zero-shot Machine Unlearning: A notable advancement in the field is the concept of zero-shot MU, as proposed by Chundawat et al. [60], which addresses scenarios where the

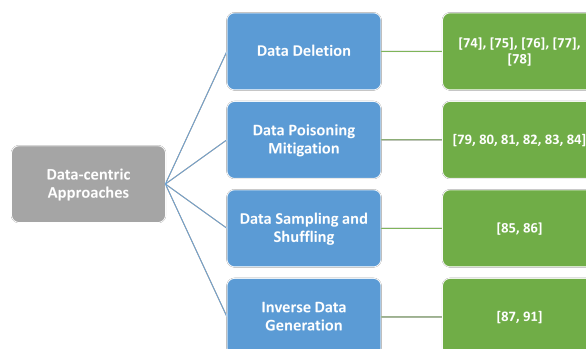


Fig. 2. Data-centric Unlearning Approaches

original data samples are not accessible. The authors introduce innovative solutions centered around error-minimizing-maximizing noise and gated knowledge transfer, aiming to eradicate the model’s reliance on the data slated for unlearning while maintaining the model’s performance on the residual data. This approach is particularly significant as it offers a robust defense mechanism against potential model inversion attacks and membership inference attacks, thereby enhancing the privacy-preserving capabilities of MU methodologies.

As MU continues to evolve, the exploration of IDG and its integration with advanced generative models will play a crucial role in shaping the future of privacy-preserving machine learning. The ability to generate data that is both representative of the original dataset and devoid of sensitive attributes will not only enhance the unlearning process but also pave the way for more secure and adaptable AI systems.

B. Model-Centric Unlearning Approaches

1) Transfer Learning

Transfer learning emerges as a pivotal strategy in MU, offering a pathway to efficiently repurpose the knowledge acquired by one ML model to facilitate or expedite the training of another model. This approach harnesses the capabilities of a pre-trained model, typically a Deep Neural Network (DNN), to serve as a foundational basis for a new model tailored to a different, yet related, task [61].

MU involves the strategic removal of specific training data from a model, simulating a scenario where the model is oblivious to the said data. This process poses significant computational challenges, especially when retraining models from scratch is deemed impractical due to resource constraints or the real-time demands of certain applications. In this context, transfer learning not only offers a solution to these challenges but also enhances the efficiency and applicability of MU across various contexts.

Baumhauer et al. [62] introduce an innovative approach to modify and cleanse classification models by employing a method termed “linear filtration.” This technique involves adjusting the model’s output predictions through a series of mathematical operations, effectively removing unwanted class-specific information. This process, coined as “sanitize classifi-

cation,” requires minimal computational resources and can be applied in a “black-box” manner, maintaining the model’s internal workings concealed while ensuring the effective removal of sensitive information. Their methodology’s practicality is further demonstrated through its resilience against attacks aimed at compromising ML model privacy.

SISA Training and Transfer Learning: Bourtole et al. [63] explore the synergy between MU and Sharded, Isolated, Sliced, and Aggregated (SISA) training, revealing that constraining individual data points’ influence during training can significantly accelerate the unlearning process. This synergy is particularly beneficial for algorithms such as stochastic gradient descent used in DNNs, where optimizing performance is crucial. The integration of SISA training with transfer learning has shown to enhance retraining speed by 1.36 times for complex tasks like ImageNet classification, albeit with a slight compromise in accuracy.

Kochno et al. [42] delve into the effects of SISA unlearning within contexts characterized by imbalanced class distributions and a correlation between class membership and unlearning probability. Their findings suggest that while SISA training can expedite unlearning, it may disproportionately affect the performance of minority classes. They propose that in scenarios with significant class imbalance, simpler strategies such as down-sampling could outperform SISA in maintaining unlearning efficiency without sacrificing model fairness.

Dynamic Selection in MU: Zhou et al. [64] propose the Dynamically Selected Mix-up (DSMixup) strategy, building upon the SISA framework to enhance MU’s efficiency. DSMixup dynamically selects mix-up data augmentation to merge shards requiring retraining, thereby reducing the need for comprehensive retraining. This approach not only boosts unlearning efficiency but also maintains system stability. Through empirical evaluations, DSMixup has demonstrated superior performance over traditional SISA in both unlearning cost and overall model performance.

In essence, transfer learning stands as a cornerstone in the evolution of MU, offering a robust framework for repurposing pre-existing model knowledge to facilitate the unlearning process. Its integration with MU methodologies like linear filtration, SISA training, and dynamic selection underscores its potential to address the computational challenges inherent in MU, paving the way for more efficient, scalable, and adaptable unlearning processes in the ever-evolving landscape of machine learning.

2) Model Pruning

Model pruning stands as a critical technique within the domain of ML, primarily aimed at optimizing model efficiency by judiciously removing non-essential parameters, thus preserving computational resources without significantly compromising the model’s performance [65]. Within the context of MU, model pruning assumes a pivotal role in facilitating the removal of sensitive or private data from a trained model, ensuring that the model’s accuracy remains largely intact.

Pruning Techniques: The essence of model pruning lies in the iterative elimination of weights or neurons considered least crucial for the model’s output. This selection is typically based on specific criteria, such as the magnitude of weights

or their contribution to the model’s output predictions [66]. Among the array of pruning techniques, weight magnitude pruning and sensitivity-based pruning are notably prevalent. Weight magnitude pruning targets weights with the smallest absolute values for elimination, while sensitivity-based pruning focuses on removing weights whose absence minimally impacts the model’s output, thereby preserving the integrity and performance of the model post-pruning.

Pruning in Federated Learning: Wang et al. [67] introduce an innovative ML unlearning method tailored to comply with the General Data Protection Regulation (GDPR), particularly focusing on the removal of specific categories from trained Convolutional Neural Network (CNN) models within a Federated Learning (FL) framework. By employing Term Frequency - Inverse Data Frequency (TF-IDF), the federated server evaluates the relevance scores between channels and categories, facilitating the construction of a pruner that targets the most discriminative channels associated with the category in question. This approach underscores the significance of model pruning as an integral component of FL, aligning with both legal and ethical standards.

Pruning Strategies for Unlearning: Jia et al. [68] propose a novel “prune first, then unlearn” paradigm, positing that initiating the unlearning process on an already sparse model can minimize unlearning errors and enhance the overall efficiency of MU. Their discussion encompasses various pruning methodologies, including one-shot magnitude pruning (OMP) [69], which involves a single iteration of pruning based on weight magnitudes, pruning at random initialization [70], which suggests pruning weights before the commencement of model training, and iterative magnitude pruning (IMP) [71], a method that combines pruning with concurrent training to iteratively refine the model’s sparsity. They advocate for the selection of a pruning technique that minimally relies on the dataset intended for forgetting, thus ensuring that the pruning process does not compromise the model’s generalization capabilities.

Model pruning emerges as an indispensable tool in the arsenal of MU techniques, offering a strategic pathway to optimize model efficiency while safeguarding data privacy and compliance with regulatory standards. The integration of model pruning within MU, particularly in conjunction with advanced methodologies like federated learning, opens new vistas for developing ML models that are not only efficient and accurate but also adaptable and compliant with evolving data privacy norms. As the field of ML continues to evolve, the role of model pruning in facilitating effective and efficient unlearning processes is poised to become increasingly significant, heralding a new era of privacy-conscious and resource-efficient ML models.

3) Knowledge Distillation

Knowledge distillation emerges as a strategic approach in MU, aimed at condensing the knowledge of complex, large-scale models into more compact and efficient counterparts while preserving the essence of their predictive capabilities [72]. This technique is predicated on the paradigm of a teacher-student relationship, wherein a smaller, less complex model (the student) is trained to mimic the predictive behavior

exhibited by a larger, pre-trained model (the teacher), thereby encapsulating the core knowledge in a more resource-efficient framework [73].

Distillation Process: The distillation process commences with the training of the teacher model on an extensive dataset, post which it generates soft targets for the student model. These soft targets, representing probability distributions over the output space of the teacher model, encapsulate richer information than traditional hard labels, thereby facilitating a more nuanced learning experience for the student model. The student model's training objective is to minimize the divergence between its predictive probabilities and the soft targets derived from the teacher model, effectively absorbing the distilled knowledge while maintaining a reduced computational footprint.

Knowledge Distillation in Federated Unlearning: Wu et al. [74] introduce a pioneering federated unlearning methodology that leverages knowledge distillation to eliminate the influence of specific clients from a global model post-federated training. This approach utilizes the global model as a teacher to guide the training of an unlearning model, effectively mitigating the attacker's influence while preserving the integrity and performance of the global model. This method offers significant advantages, including the reduction of client-side computational demands and the enhancement of model generalization, thereby bolstering the robustness and efficacy of the federated learning ecosystem.

Knowledge Distillation for Heterogeneous Learning: Zhu et al. [75] propose the FedLU framework, tailored for heterogeneous knowledge graph (KG) embedding learning and unlearning within a federated setting. By employing mutual knowledge distillation, the framework facilitates the bidirectional transfer of knowledge between local and global levels, enabling the systematic unlearning of specific knowledge components from local embeddings. This methodology underscores the versatility of knowledge distillation in managing data heterogeneity and fostering coherent knowledge integration and unlearning across diverse data partitions.

Mitigating Backdoor Attacks: Addressing the challenge of backdoored Deep Neural Networks (DNNs), Li et al. [76] present the Neuron Attention Distillation (NAD) framework, designed to neutralize backdoor triggers embedded within DNNs. NAD employs a fine-tuning process guided by a teacher network, ensuring the alignment of intermediate-layer attentions between the teacher and the backdoored student network. This alignment, achieved through fine-tuning on a subset of clean training data, effectively diverts the student network's attention away from the malicious triggers, thereby restoring the network's integrity. The empirical validation of NAD demonstrates its superior efficacy in mitigating backdoor influences compared to conventional fine-tuning and pruning techniques.

Knowledge distillation stands as a cornerstone in the advancement of MU, offering a pathway to encapsulate and transfer essential knowledge across models in a resource-efficient manner. Its application spans federated unlearning, heterogeneous learning environments, and the mitigation of backdoor attacks [77], highlighting its versatility and po-

tential in enhancing the privacy, efficiency, and security of machine learning models. As the landscape of MU continues to evolve, knowledge distillation is poised to play an integral role in shaping the next generation of privacy-preserving and resource-efficient machine learning paradigms.

4) Model Inversion

Model inversion has emerged as a sophisticated technique within MU, primarily aimed at elucidating sensitive information encoded within ML models. This technique capitalizes on the capacity to reverse-engineer a model's predictions to infer the characteristics of the underlying training data, thereby posing potential risks to data privacy and security [78].

Inversion Mechanism: The crux of model inversion lies in manipulating a model's inputs in such a manner as to elicit a particular output, which, in turn, sheds light on the original input data's attributes. This process essentially enables the extraction of confidential information from the model, potentially compromising data privacy. For instance, attackers can utilize model inversion to deduce specific characteristics or features that a model heavily relies on for classification tasks, subsequently exploiting this knowledge to circumvent model predictions or to reconstruct sensitive data.

Mitigating Inversion Attacks: To counteract the threats posed by model inversion attacks, various defensive strategies have been proposed, including the adoption of differential privacy measures, regularization techniques, and adversarial training methodologies. These approaches aim to fortify model resilience against inversion attempts, thereby safeguarding sensitive data from unauthorized reconstruction or inference [78].

Innovative Inversion Approaches: Graves et al. [79] introduce a refined inversion attack strategy, extending beyond conventional model inversion paradigms. Their methodology commences with an initial feature vector assigned null values, subsequently iteratively modified to mirror what the model perceives as representative of a given class. This iterative process is augmented with periodic application of image processing techniques, enhancing the fidelity of the reconstructed data. Such advancements underscore the evolving complexity of inversion attacks, necessitating continual enhancement of defensive measures.

Few-shot Unlearning via Model Inversion: Yoon et al. [80] explore the application of model inversion in the context of few-shot unlearning, presenting a novel framework that leverages inversion techniques to approximate and subsequently eliminate specific data distributions from a model. This approach not only facilitates the efficient unlearning of targeted data but also introduces a mechanism for identifying and filtering out noisy samples, thereby enhancing the precision of the unlearning process.

The exploration of model inversion within MU encapsulates a broad spectrum of techniques and methodologies, aimed at both exploiting and defending against inversion attacks. The strategies encompass not only the direct inversion of model predictions but also the utilization of inversion techniques to facilitate targeted unlearning and data extraction processes, as illustrated in Fig. 3.

Model inversion represents a critical aspect of MU, offering both challenges and opportunities in the realm of data privacy

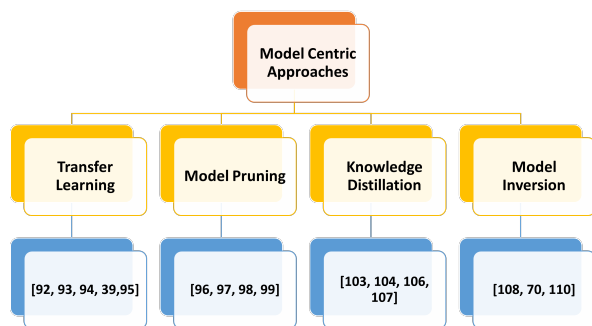


Fig. 3. Model-centric Unlearning Approaches

and model integrity. As ML models continue to pervade diverse domains, the imperative to develop and refine inversion-resistant models becomes increasingly paramount. The ongoing advancements in inversion techniques and countermeasures herald a dynamic and evolving landscape in MU, underscoring the need for vigilant and innovative approaches to safeguard sensitive data within ML models.

C. Federated Unlearning

Federated Unlearning represents a cutting-edge paradigm within the broader MU domain, focusing on the meticulous erasure of specific data contributions from models trained in a Federated Learning (FL) framework, while simultaneously preserving the model's overall functionality and accuracy. This subsection delves into the multifaceted methodologies and frameworks that embody federated unlearning, as summarized in Table I, each tailored to harmonize the triad of privacy, efficiency, and model performance.

Enhancing Data Privacy: A pivotal aspect of federated unlearning is its commitment to bolstering data privacy without compromising the efficacy of the model. Illustrative of this commitment is the FedME2 framework [81], which pioneers a pathway to privacy-centric unlearning by orchestrating data forgetting with model accuracy within Digital Twins for Mobile Networks (DTMN) via a nuanced multi-loss training approach. Complementarily, FedRecovery [82] leverages differential privacy to obfuscate a client's data footprint from the global model, circumventing the need for retraining and ensuring a seamless congruence between unlearned and retrained models. These innovative methodologies underscore the quintessential challenge federated unlearning aims to address: ensuring robust privacy protections while minimizing the computational demands associated with model retraining.

Optimizing Resource Utilization: In scenarios characterized by limited server storage or computational resources, the significance of efficient unlearning mechanisms becomes paramount. For instance, the Subspace-based Federated Unlearning (SFU) approach [83] negates the necessity for additional storage by conducting gradient ascent in the orthogonal complement of input gradient spaces, effectively nullifying a

target client's influence on the model. Similarly, the Federated Clusters method [84] significantly expedites the unlearning process, offering a substantial reduction in time compared to traditional retraining methods. These strategies highlight the critical need for judicious resource management, a cornerstone for achieving scalable federated unlearning solutions.

Leveraging Bayesian Methodologies: Bayesian approaches offer a structured framework for federated unlearning, ensuring a systematic and principled removal of data while maintaining model integrity. The Bayesian Variational Federated Learning and Unlearning technique [85] exemplifies this approach by employing federated variational inference to facilitate efficient unlearning via local free energy minimization within exponential-family models. Additionally, FORGET-SVGD [86] introduces a particle-based Bayesian unlearning method, which enacts local updates on agents desiring to unlearn, punctuated by communication rounds with a parameter server, thereby providing a non-parametric strategy for federated unlearning.

Assuring Verifiable Data Removal: The aspect of ensuring verifiable and credible data removal, without impinging on model performance, is a critical facet of federated unlearning. The Forgettable Federated Linear Learning (2F2L) framework [87] epitomizes this by employing a linear approximation on model parameters and introducing an efficient removal strategy that constrains the variances in model weights, thereby ensuring reliable data removal. This facet of providing validated and certified unlearning is paramount for upholding the reliability and trustworthiness of federated unlearning processes.

Guaranteeing Execution of Unlearning: The VERIFI framework [88] represents a seminal advancement in federated unlearning by not only facilitating data removal but also enabling the verification of the unlearning effect, thus ensuring both the execution and quantifiability of the unlearning process. This framework empowers departing participants with the right to verify (RTV) the efficacy of unlearning, thereby instilling a tangible assurance in the "right to be forgotten."

Adaptive Learning Mechanisms: The FRAMU framework [35] adopts a holistic stance by integrating adaptive learning mechanisms with privacy preservation strategies and optimization techniques, ensuring the model's adeptness in unlearning obsolete or irrelevant data, thereby supporting continual model evolution in accordance with dynamic data landscapes without compromising privacy norms.

The exploration of federated unlearning traverses through the intricate challenges of ensuring privacy, computational efficiency, and model integrity, weaving a comprehensive narrative that shapes the future trajectory of federated unlearning within the domain of privacy-preserving ML. This dynamic landscape beckons for innovative, balanced solutions that adeptly navigate the complexities of data privacy and efficient model unlearning.

D. Graph Unlearning

Graph unlearning, as an emerging field within MU, addresses the complex challenge of retracting specific nodes, edges, or attributes from models trained on graph-structured data.

TABLE I
SUMMARY OF FEDERATED UNLEARNING FRAMEWORKS

Framework	Description	Reference
FedME2	Managing data forgetting and model accuracy in DTMN through a multi-loss training approach.	[81]
FedRecovery	Utilizes differential privacy to erase a client's data influence without retraining, ensuring statistical indistinguishability between models.	[82]
SFU	Performs gradient ascent in orthogonal space of input gradient spaces, negating a target client's contribution without additional storage.	[83]
Federated Clusters Method	Accelerates the unlearning process, providing a significant speed-up compared to retraining.	[84]
Bayesian Variational Federated Learning and Unlearning	Utilizes federated variational inference solutions, offering an efficient unlearning mechanism via local free energy minimization.	[85]
FORGET-SVGD	Employs a particle-based Bayesian unlearning method, providing a non-parametric strategy for federated unlearning.	[86]
2F2L	Ensures certified data removal by employing a linear approximation of the model parameter space and introducing an efficient removal strategy.	[87]
VERIFI	Facilitates federated unlearning and enables verification of the unlearning effect, ensuring both execution and quantifiability of removal of a participant's gradients from the global model.	[88]
FRAMU	Amalgamates adaptive learning mechanisms, privacy preservation, and optimization strategies for model evolution without infringing upon privacy while unlearning outdated data.	[35]

Given the inherent interconnectedness of graph data, graph unlearning requires nuanced approaches that ensure the efficient removal of data without degrading the model's performance. This section explores the forefront methodologies and frameworks devised to address these challenges, as summarized in Table II.

Innovative Unlearning Methods: Leading the charge in graph unlearning are methodologies like PROJECTOR and GRAPHEDITOR, developed by Cong et al. [89], [90]. PROJECTOR employs projection techniques to effectively erase specific nodes from GNNs, ensuring that the unlearned node features are completely purged from the model parameters. Conversely, GRAPHEDITOR offers a dynamic solution to graph unlearning by facilitating various operations such as node/edge deletion, addition, and node feature modification, thereby showcasing the adaptability required in graph unlearning [89], [90].

Guaranteeing Certified Unlearning: The realm of certified graph unlearning has witnessed significant contributions, with frameworks providing theoretical guarantees for unlearning processes. Noteworthy among these is the CEU framework, which introduces a single-step update methodology for the

removal of specific edges, backed by robust theoretical underpinnings [91]. Such certified approaches are crucial for applications requiring verifiable assurances of data removal.

Diverse Approaches to Graph Unlearning: The diversity in graph unlearning methodologies is further exemplified by GUIDE, FedLU, and GNNDELETE. GUIDE focuses on inductive learning and unlearning within dynamic graphs, emphasizing the need for adaptable models in evolving graph environments [92]. FedLU, on the other hand, addresses the unlearning challenges in federated settings for heterogeneous KG embeddings, underscoring the complexities of data heterogeneity in graph unlearning [75]. GNNDELETE introduces optimization strategies for node and edge deletions, ensuring the preservation of learned knowledge post-unlearning [93].

Exploring Alternative Perspectives: Frameworks like GST and GIF provide alternative viewpoints on graph unlearning. GST, or Graph Scattering Transform, focuses on the mathematical robustness of unlearning processes, while GIF, or Graph Influence Function, highlights the role of influence functions in graph unlearning [94], [95]. These perspectives enrich the discourse on graph unlearning, offering novel insights into the unlearning mechanisms.

GraphEraser: A Paradigm Shift: GraphEraser emerges as a paradigm-shifting framework, emphasizing efficient partitioning and aggregation mechanisms for unlearning in graph data. This framework illustrates the potential of MU in addressing the privacy and integrity concerns inherent in graph-structured data [96].

These methodologies and frameworks signify a pivotal evolution in graph unlearning, each contributing unique solutions to the challenges posed by graph-structured data. The amalgamation of these diverse approaches underscores the multidimensional nature of graph unlearning, driving forward the agenda of privacy preservation and the "right to be forgotten" in the context of Graph Neural Networks (GNNs).

IV. ADVANCED METRICS FOR EVALUATING MACHINE UNLEARNING

Evaluating the efficacy and integrity of MU methodologies necessitates a comprehensive set of metrics that cater to various aspects of model performance and data privacy post-unlearning. This section delineates an array of advanced metrics, each offering unique insights into the MU process, thereby addressing the critical research question: "How can the impact of unlearning on model performance and privacy be measured and evaluated comprehensively?" The metrics discussed herein are summarized in Table III, offering a holistic view of their applications and implications within the MU paradigm. Supplementary Material provides a curated list of publicly accessible datasets, detailing their categories, instances, attributes, tasks, and citation frequency, facilitating empirical evaluations.

Accuracy and Precision: Accuracy, a fundamental metric, offers a straightforward assessment of a model's predictive performance post-unlearning [98]. While its simplicity is commendable, it may not fully capture the nuanced effects of unlearning on specific data classes or distributions. Precision,

TABLE II
SUMMARY OF GRAPH UNLEARNING FRAMEWORKS

Topic	Methodologies	Reference
Unique Methods	1. PROJECTOR - Uses projection techniques to remove specific nodes, ensuring no trace in model parameters. 2. GRAPHEEDITOR - Manages dynamic graphs and supports operations like node/edge deletion and feature updates.	[89], [90]
Certified Guarantees	Emphasis on theoretical insights and performance guarantees in unlearning. CEU framework provides a single-step update mechanism with theoretical support for removing specific edges.	[91], [97]
Diverse Approaches	1. GUIDE - Inductive graph learning/unlearning in dynamic graphs. 2. FedLU - Federated learning for KG embedding, addressing data heterogeneity. 3. GNNDELETE - Optimization strategies for node/edge deletions without loss of knowledge.	[75], [92] [93]
Alternative Perspectives	1. GST (Graph Scattering Transform) - Focus on mathematical robustness in unlearning. 2. GIF (Graph Influence Function) - Highlights influence functions in unlearning.	[94], [95]
Innovative Framework	GraphEraser - Stresses efficient partitioning and aggregation mechanisms for unlearning in graph data.	[96]

on the other hand, provides a finer granularity by measuring the model’s ability to correctly predict positive instances among all predicted positives, which becomes crucial in imbalanced datasets where specific classes are more sensitive to unlearning.

Anamnesis Index (AIN): The Anamnesis Index quantifies the extent to which a model retains information about unlearned data, serving as a direct measure of the unlearning effectiveness [99]. This metric is particularly relevant in scenarios where regulatory compliance demands verifiable evidence of data removal. However, it might not fully encapsulate the model’s retained knowledge on the remaining dataset.

Internal Representation Distances: Metrics such as Activation Distance and Layer-wise Distance delve into the model’s internal state changes post-unlearning, shedding light on the structural and functional modifications within the network [98]. These metrics are invaluable for understanding the deeper implications of unlearning on model architecture. Their complexity, however, might obscure their interpretability in relation to direct model output or user-facing performance metrics.

Membership Inference (MI) Attacks Vulnerability: Assessing the model’s susceptibility to MI attacks post-unlearning is crucial for ensuring data privacy [28]. This metric evaluates the model’s propensity to leak information about whether a specific data point was part of the training set, a critical aspect in the context of privacy-preserving MU. While highly relevant for security assessments, this metric might not directly correlate with the model’s predictive accuracy.

Behavioral Divergence Metrics: Metrics like JS-Divergence and Zero Retrain Forgetting Metric offer quantitative measures to evaluate the divergence in model behavior before and after unlearning [100]. These metrics are instrumental in assessing the stability and consistency of

model predictions, ensuring that unlearning does not lead to erratic or significantly altered model behavior.

Epistemic Uncertainty: Quantifying epistemic uncertainty post-unlearning provides insights into the model’s confidence in its predictions, reflecting the impact of unlearning on the model’s knowledge base [101]. This metric is especially pertinent in high-stakes applications where decision-making relies on model certainty. However, the computation of uncertainty metrics can be resource-intensive and may require sophisticated probabilistic models.

Specialized Unlearning Metrics: Metrics such as Fisher Forgetting and Variational Forgetting delve into the model’s ability to effectively “forget” specific tasks or data points, crucial for targeted unlearning scenarios [102], [103]. These metrics are tailored to assess the model’s resilience to adversarial manipulations and its capacity for task-specific unlearning. Their specialized nature, however, may limit their applicability across diverse unlearning contexts.

In essence, the selection and application of evaluation metrics in MU must be context-driven, taking into account the specific goals of unlearning, the nature of the data and model, and the regulatory and ethical standards governing data privacy. The comprehensive assessment facilitated by these metrics enables a nuanced understanding of MU’s impact, guiding the development of more robust, efficient, and privacy-preserving unlearning methodologies.

V. CHALLENGES AND POTENTIAL SOLUTIONS

The landscape of MU is riddled with complex challenges, necessitating innovative solutions to ensure robust and effective unlearning mechanisms. This section delves into the prevalent challenges within MU, offering insights into potential strategies for addressing these hurdles, thereby answering the pivotal research question: “*What are the predominant challenges in MU, and how can these challenges be effectively addressed?*” Figure 4 provides a visual roadmap linking these challenges to their corresponding solutions.

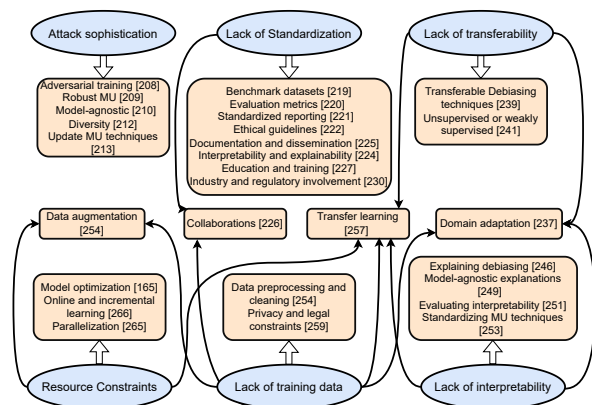


Fig. 4. A Roadmap of Machine Unlearning Challenges and Corresponding Solutions

TABLE III
MACHINE UNLEARNING - EVALUATION METRICS

Metric	Description	Equation	References
Accuracy	The proportion of correctly classified instances in the test set	$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$	[98]
Anamnesis Index	Measures the extent to which unlearned data has been forgotten	$AI(T) = (Acc(M, T) - Acc(M - T, T)) / Acc(Naive, T)$	[99]
Activation Distance	Measures the difference between the activations of two models on a given input	$AD(A(x), B(x)) = \ A(x) - B(x)\ _2$	[98]
Layer-wise Distance	Measures the difference between the weight matrices of two models	$LWD_l = \ W_l - W'_l\ _F$	[104]
Membership Inference (MI) Attacks	Measures the degree to which an attacker can infer whether a particular instance was included in the training data	N/A	[28]
JS-Divergence	Jensen-Shannon divergence between the predictions of the unlearned and retrained model	$JS(M(x), Td(x)) = 0.5 \times KL(M(x) \ m) + 0.5 \times KL(Td(x) \ m)$	[100]
Zero Retrain Forgetting Metric	Measures the change in accuracy on the unlearned data after retraining the model without the data to be unlearned	$ZFR = 1 - \frac{1}{n_f} \sum_{i=0}^{n_f} JS(M(x_i), Td(x_i))$	[100]
Reconstruction Error	Measures the difference between the original input and the reconstructed input	$RE = \ x - f(g(x))\ _2$	[105]
Completeness	Ensures that all traces of the unlearned data have been removed from the model	NA	[106]
Unlearn time	substantial number of epochs for the network to unlearn the forget samples	NA	[7]
Relearn time	substantial number of epochs for the unlearned model to achieve accuracy of original model	NA	[7], [104]
Epistemic Uncertainty	Measures the model's uncertainty regarding its predictions due to a lack in knowledge (model uncertainty)	$efficacy(\theta; D) = \begin{cases} 1/i(\theta; D), & \text{if } i(\theta; D) > 0 \\ \infty, & \text{otherwise} \end{cases}$	[101]
Model Inversion attack	Quantifies the success of reconstructing input data from model outputs by an adversary	NA	[102]
Fisher Forgetting	Measures the forgetting of task A while learning task B using Fisher Information	$S(w) = w + (\lambda\sigma^2 h)^{\frac{1}{4}} F^{-\frac{1}{4}}$	[103]
Variational Forgetting	Quantifies forgetting by measuring the variational distance between posterior distributions of the parameters before and after learning new data	$VF = D_{KL}(q(\theta D_{old}) q(\theta D_{new}))$	[98]

A. Complexities of Attack Sophistication

As MU evolves, so does the sophistication of potential adversarial attacks aimed at compromising the privacy and integrity of ML models. These advanced threats often combine techniques such as data poisoning, model inversion, and adversarial attacks to exploit vulnerabilities within MU mechanisms [107], [108]. Addressing the challenge of attack sophistication in MU encompasses a multifaceted approach, focusing on enhancing the resilience of MU techniques against a diverse array of sophisticated adversarial strategies.

Evolving Attack Strategies: Adversaries continually refine their techniques, crafting more complex and stealthy attacks that can bypass conventional MU defenses, challenging the detection and mitigation capabilities of these systems [34].

Adaptive Adversaries: Attackers adeptly modify their strategies in response to the deployment of new MU defenses, creating a dynamic adversarial landscape that necessitates continuous adaptation from MU mechanisms [109].

Stealth and Subtlety of Attacks: Advanced attacks are designed to be less detectable by MU systems, subtly manipulating data or model behavior in ways that evade traditional detection methods [110].

Robustness Limitations: Current MU techniques may lack the robustness required to counteract sophisticated attacks, especially those exploiting specific vulnerabilities in the unlearning processes [111].

Addressing the nuanced challenge of attack sophistication necessitates the development of MU defenses that are not only robust but also adaptive, capable of evolving in tandem with the changing adversarial tactics.

Adversarial Training: Enhancing the resilience of ML models through adversarial training, which involves incorporating adversarial examples into the training process, thereby fortifying the model against potential attacks [112].

Development of Robust MU Techniques: Crafting MU mechanisms that can withstand sophisticated attacks, possibly by integrating advanced adversarial training methods to bolster their efficacy against evolving threats [113].

Model-Agnostic Approaches: Leveraging techniques that are independent of specific model architectures to detect and mitigate biases or vulnerabilities, thereby offering a broader defense mechanism against complex attacks [114], [115].

Diversifying Training Data: Ensuring a wide-ranging and diverse dataset for training to reduce the model's susceptibil-

ity to targeted attacks, thereby enhancing the robustness of MU processes [116].

Continuous Re-evaluation and Updating: Implementing a dynamic framework for the regular assessment and enhancement of MU techniques, ensuring they remain effective against the continuously evolving strategies of adversaries [117].

B. Lack of standardization

The absence of standardization in the application of MU techniques presents a significant challenge, hindering the ability to effectively compare, evaluate, and validate different approaches within the field [41], [118]. This gap in standardization spans across various dimensions, including methodologies, evaluation metrics, data and model compatibility, ethical considerations, and the interpretability of unlearning outcomes. This lack of standardization can pose challenges across several areas:

Diverse Methodologies: The field currently lacks a universally accepted framework for implementing unlearning, leading to a proliferation of disparate techniques and algorithms. This diversity complicates the comparison and reproduction of results across different studies, potentially resulting in inconsistent outcomes [118].

Evaluation Metrics: The evaluation of MU techniques is encumbered by the lack of standardized benchmarks and metrics. This makes it challenging to objectively assess the effectiveness, fairness, and robustness of various unlearning methods [41].

Data and Model Compatibility: The absence of a standardized approach for ensuring compatibility across different data types, ML models, and applications complicates the adaptation of unlearning methods to diverse contexts [119].

Ethical and Transparency Concerns: MU raises ethical issues related to fairness, accountability, and transparency. The lack of standardized ethical guidelines can lead to varied interpretations and implementations, potentially resulting in biases or unintended consequences in unlearning outcomes [120], [121].

Interpretability and Explainability: There is a notable absence of standardized methods for explaining the modifications made by unlearning techniques to ML models, which is crucial for ensuring the trustworthiness and comprehensibility of the unlearning process [122].

Addressing the challenge of standardization in MU necessitates collaborative efforts aimed at developing common frameworks, benchmarks, and guidelines, thereby fostering a cohesive and transparent approach to unlearning. Some potential approaches to promote standardization in the MU community include:

Community-wide Collaborations: Fostering dialogue and cooperation among researchers, practitioners, and stakeholders to establish consensus on standards, guidelines, and best practices for MU. Such efforts could include dedicated workshops, symposiums, and forums to facilitate knowledge exchange [123].

Benchmark Datasets and Evaluation Metrics: Creating benchmark datasets and defining clear evaluation metrics

tailored for MU. These benchmarks will enable the objective assessment of unlearning methods and facilitate comparisons across different studies [124].

Standardized Reporting: Advocating for transparent and uniform reporting of unlearning studies, encompassing methodologies, algorithms, datasets, and evaluation metrics used. This approach will enhance the reproducibility and comparability of research findings in the MU domain [125].

Ethical Guidelines: Developing comprehensive ethical guidelines that address the unique challenges posed by MU, including considerations of fairness, accountability, and transparency. These guidelines can guide practitioners in addressing ethical concerns associated with unlearning techniques [126], [127].

Interpretability and Explainability Frameworks: Establishing standardized methods for interpreting and explaining the outcomes of unlearning processes. Techniques for model introspection, visualization, and explanation can contribute to a better understanding and communication of unlearning results [128].

Documentation and Dissemination: Encouraging thorough documentation and open sharing of unlearning methodologies, techniques, and findings. Open repositories and code libraries can promote knowledge sharing and the development of common practices in MU [129].

Cross-disciplinary Collaboration: Engaging with adjacent fields such as machine learning, ethics, and fairness research to leverage established standards and practices. This cross-pollination can enrich the MU field with proven methodologies and foster interdisciplinary collaboration [130].

Education and Training Initiatives: Integrating MU concepts, methodologies, and best practices into educational and training programs. This effort will cultivate a shared understanding and widespread adoption of standardized approaches within the MU community [131], [132].

Engagement with Industry and Regulatory Bodies: Involving industry stakeholders and regulatory authorities in the standardization process to ensure the practical relevance and regulatory compliance of MU standards. This collaboration can facilitate the responsible and ethical application of unlearning techniques in real-world scenarios [133], [134].

C. Lack of transferability

The challenge of transferability in MU encompasses the difficulty in applying unlearning techniques developed for specific models or datasets to other contexts effectively. This lack of transferability hampers the scalability and adaptability of MU methods across diverse ML models and real-world scenarios [135], [136]. Several factors contribute to this challenge:

Model-Specific Biases: The efficacy of MU techniques can vary significantly across different ML models due to model-specific biases. Techniques optimized for deep neural networks, for instance, may not be directly applicable to simpler models like decision trees [137], [138].

Dataset-Specific Biases: The success of unlearning approaches can depend heavily on the characteristics of the training data. Techniques designed for a particular dataset may lose effectiveness when applied to datasets with different

distributions or types of biases [31], [139].

Domain-Specific Challenges: MU methods developed within a specific domain (e.g., healthcare or finance) might not translate well to other domains, which could have unique data distributions and domain-specific biases [30], [140].

Lack of Labeled Data: The transferability of MU techniques is also hindered by the scarcity of labeled data in new domains or datasets, which is crucial for training, validating, or fine-tuning unlearning models [141].

Overcoming the transferability challenge in MU involves developing strategies that enhance the adaptability of unlearning techniques across various models, datasets, and domains. Some potential approaches to mitigate this challenge include:

Domain Adaptation: Leveraging domain adaptation and generalization methods can facilitate the adaptation of MU techniques to new domains, even with limited labeled data in the target domain. These methods focus on bridging the gap between source and target domains, enhancing the applicability of unlearning methods across different contexts [142], [143].

Transfer Learning: Employing transfer learning strategies can enable the transfer of unlearning knowledge from one model or dataset to another. This approach can be particularly effective in utilizing pre-trained models or algorithms to adapt unlearning techniques to new models with varying architectures.

Development of Transferable Unlearning Techniques: Crafting debiasing and unlearning methods with built-in transferability, potentially through the incorporation of generalization principles or domain-independent features, can facilitate their application across diverse settings. Model-agnostic approaches can also contribute to this adaptability, ensuring that unlearning techniques remain effective irrespective of model architecture [144], [145].

Exploration of Unsupervised and Weakly Supervised Methods: Investigating unsupervised or weakly supervised MU techniques can mitigate the dependency on labeled data, enhancing the transferability of these methods. Approaches such as unsupervised learning or self-supervised learning can provide viable pathways for applying MU techniques in scenarios where labeled data is scarce [146], [147].

Utilization of Benchmark Datasets: Establishing and employing benchmark datasets that cover a wide range of models, datasets, and domains can aid in the standardized evaluation of MU techniques. This approach can offer valuable insights into the transferability and generalizability of unlearning methods, facilitating their refinement and adaptation to new contexts [148], [149].

D. Lack of interpretability

The challenge of interpretability in MU revolves around the difficulty in understanding and explaining the mechanisms and outcomes of unlearning techniques, especially when involving complex models like deep neural networks (DNNs) or generative models [150]. The complexity of these techniques, coupled with the potential loss of transparency during the unlearning process, complicates efforts to achieve interpretability. This section outlines several strategies aimed at addressing

these challenges and enhancing the interpretability of MU techniques.

Addressing this challenge requires careful consideration of the specific de-biasing or unlearning techniques used and their impact on model interpretability. Some potential approaches to enhance interpretability include:

Explaining the Unlearning Process: Providing detailed explanations about the data modifications, feature alterations, or the mechanisms employed during the unlearning process can help stakeholders comprehend the rationale behind the model's revised outputs [32], [33], [151]. This involves elucidating the specific actions taken to remove biases or undesired information from the model, enhancing transparency.

Simplification of Unlearning Techniques: Developing MU techniques that are inherently simpler and more interpretable can make it easier for both developers and end-users to understand how unlearning is achieved [97]. Simplification might involve adopting less complex models or mechanisms that are easier to explain and validate.

Model-Agnostic Explanations: Employing model-agnostic interpretability tools, such as feature importance measures or partial dependence plots, can offer insights into how different features or data points influence the model's predictions after unlearning [152], [153]. These tools can provide a layer of interpretability that is independent of the underlying model architecture or the specific unlearning method applied.

Incorporating Interpretability in Evaluation: Making interpretability an explicit criterion during the development and evaluation of MU techniques can ensure that it is not compromised in the pursuit of performance [154], [155]. This involves assessing the interpretability of models post-unlearning, alongside traditional performance metrics like accuracy or precision.

Standardization of Unlearning Methods: Establishing standardized approaches and guidelines for implementing and evaluating MU techniques can foster consistency and interpretability within the field [156]. Standardization can also facilitate the replication of unlearning studies and the comparison of different unlearning approaches, contributing to a more transparent and accountable MU ecosystem.

E. Lack of training data

The effectiveness of MU techniques is significantly influenced by the availability and quality of training data. However, obtaining sufficient and relevant training data for MU, particularly in contexts where original training data is scarce or sensitive, presents a notable challenge [157]. This section explores several approaches to mitigate the impact of limited training data on MU.

Data Augmentation: Utilizing data augmentation techniques, such as synthetic data generation or oversampling of underrepresented groups, can help address issues of data scarcity and imbalance [158], [159]. This approach can expand the available training data pool, making it more representative and balanced for effective unlearning.

Data Pre-processing and Cleaning: Implementing rigorous data pre-processing and cleaning procedures can enhance the quality and reliability of training data [160], [161]. This

involves removing inaccuracies, inconsistencies, or noise from the data, thereby improving the foundation upon which MU techniques are applied.

Leveraging Transfer Learning and Domain Adaptation: Employing transfer learning and domain adaptation techniques can facilitate the application of knowledge from data-rich source domains to target domains with limited data [160], [162]. This approach can enhance the performance of MU models in new contexts, despite the scarcity of labeled data.

Collaborative Data Sharing and Pooling: Fostering collaborative efforts among researchers, organizations, and stakeholders to share and pool data resources can collectively overcome data limitations [163]. Shared datasets and benchmarks can provide a valuable resource for the MU community, facilitating research and development efforts.

Adherence to Privacy and Legal Guidelines: Ensuring compliance with privacy regulations and ethical guidelines is paramount when addressing data scarcity [164]. This may involve anonymizing data, obtaining necessary consents, or adhering to data-sharing agreements, ensuring that MU techniques are applied responsibly and ethically.

F. Resource constraints

Resource constraints represent a significant challenge in the domain of MU, referring to the limitations associated with computational resources, time, and data availability. These constraints can hinder the effective and efficient implementation of unlearning techniques, particularly in scenarios requiring intensive computations, real-time processing, or substantial labeled data [10], [165]. Some common types of resource constraints in MU include:

Computational Resources: The requirement for high processing power, memory, or storage can be a barrier, especially for deep learning-based MU methods that necessitate substantial computational resources for training and application [166], [167].

Time Constraints: The significant time required for training or applying MU techniques can be problematic in real-time or online learning contexts where prompt unlearning is essential [79], [168].

Data Availability: The need for extensive labeled data for training, validation, or fine-tuning poses challenges, particularly in situations involving sensitive, proprietary, or hard-to-collect data [169], [170].

Addressing resource constraints requires careful consideration of the available resources and the specific requirements of the unlearning techniques. Some potential approaches to mitigate resource constraints in MU include:

Model Optimization: Techniques such as model compression, quantization, or approximation can mitigate the computational demands of MU methods, facilitating deployment in environments with limited computational capabilities [171], [172].

Parallelization: Employing parallel computing strategies, including distributed computing and the use of multi-core processors or GPUs, can expedite MU processes, addressing computational and time constraints [83], [173].

Online and Incremental Learning: Updating MU models in an online or incremental fashion, as opposed to batch retraining, can enhance time efficiency and resource utilization by incrementally adjusting models with new data [174]–[176].

Transfer Learning: Applying transfer learning, through techniques like fine-tuning or domain adaptation, leverages pre-trained models or knowledge from related tasks, reducing the dependency on extensive labeled data for MU [177].

Data Augmentation: Generating synthetic data points through data synthesis, generation, or simulation techniques can augment existing datasets, thereby alleviating the challenges posed by limited data availability [109], [178], [179].

VI. FUTURE DIRECTIONS

A. Machine Unlearning in Natural Language Processing (NLP)

The dynamic landscape of Natural Language Processing (NLP) presents unique challenges and opportunities for MU. NLP models, designed to understand, interpret, and generate human language, often rely on vast datasets to learn patterns, semantics, and syntax [11], [180]. However, the mutable nature of language, influenced by cultural shifts, emerging terminologies, and evolving social contexts, necessitates continuous adaptation of these models to remain relevant and accurate [181].

1) Challenges in NLP MU

Rapid Evolution of Language: The fast-paced evolution of linguistic patterns, especially in online and informal communication, requires NLP models to constantly update to understand contemporary usage and semantics [182].

Bias and Ethical Considerations: NLP models can inadvertently learn and perpetuate biases present in the training data. MU is essential for identifying and rectifying these biases to ensure that models do not propagate harmful stereotypes or misinformation [14].

Adaptation to New Contexts: The application of NLP models across diverse domains necessitates their adaptation to domain-specific terminologies and contexts, which can significantly differ from the data on which the models were originally trained [183].

2) Opportunities and Future Directions

- **Continual Learning:** Integrating MU with continual learning frameworks can enable NLP models to adapt to new linguistic trends and data streams without the need for complete retraining, thereby maintaining their relevance and accuracy over time.
- **Ethical AI and Fairness:** Developing MU techniques that specifically target and mitigate biases in language models can contribute to the advancement of ethical AI, ensuring that NLP applications promote fairness and inclusivity.
- **Cross-Domain Adaptability:** Enhancing MU methodologies to facilitate seamless adaptation of NLP models across different domains can significantly broaden their applicability and utility, enabling more accurate and context-aware language processing.

- **Interactive and Explainable MU:** Fostering research in interactive MU techniques that involve human-in-the-loop approaches can enhance the interpretability and trustworthiness of NLP models, allowing for more nuanced unlearning and relearning based on expert feedback.

MU in NLP is pivotal for addressing the dynamic challenges posed by the evolving nature of language and societal norms. Future research in this domain should focus on developing robust, adaptable, and ethical MU strategies that ensure the continual relevance and fairness of NLP models.

B. Machine Unlearning in Computer Vision

Computer Vision (CV) stands at the forefront of fields benefiting from MU, addressing challenges inherent to visual data processing and interpretation. The integration of MU in CV is driven by the need for models that are adaptable, interpretable, and ethically aligned [184].

1) Challenges and Future Directions in CV MU

Adaptability to Evolving Visual Data: The ever-expanding diversity in visual content necessitates CV models that can dynamically adapt to new data types, environmental variations, and evolving user needs. Future MU research should focus on developing models capable of incremental learning and real-time adaptation, especially crucial for applications like autonomous driving and real-time surveillance [185].

Enhanced Model Robustness: Robustness in CV models entails resilience against adversarial attacks, noise, and data corruption. MU can contribute to enhancing robustness by allowing models to unlearn vulnerabilities and adapt to counteract adversarial techniques [128].

Bias Mitigation and Ethical AI: CV applications, particularly those involving facial recognition and social analytics, are prone to biases that can perpetuate stereotypes and discrimination. MU offers a pathway to identify and rectify these biases, ensuring models reflect ethical considerations and fairness [104].

Explainability and Transparency: As CV systems become integral to critical decision-making processes, explainability becomes paramount. Future MU initiatives should aim at developing techniques that not only allow models to unlearn but also provide insights into how and why certain decisions are made, facilitating human oversight and trust [186].

Privacy and Security in MU: With the increasing use of personal and sensitive visual data, ensuring the privacy and security of this information during the unlearning process is essential. Future research must address the development of secure MU processes that safeguard data integrity while enabling effective unlearning [187].

Human-in-the-loop Unlearning: In complex or ambiguous scenarios, human expertise may be required to guide the unlearning process. Incorporating human feedback can enhance the relevance and accuracy of unlearning, making models more aligned with real-world needs and ethical standards.

2) Implications for Real-world Applications

MU in CV holds significant implications for a broad spectrum of applications, from enhancing the ethical deployment of facial recognition technologies to ensuring the safety and reliability of autonomous systems. By embedding MU principles,

CV models can become more adaptable, ethically conscious, and aligned with user expectations and societal norms. This not only enhances the performance and reliability of CV systems but also fosters trust and acceptance among users and stakeholders [188].

The future of MU in CV is intertwined with advancements in model adaptability, robustness, ethical AI, and explainability. Addressing these challenges through innovative MU techniques will pave the way for CV models that are not only technically proficient but also ethically responsible and user-centric.

C. Machine Unlearning in Recommender Systems

The evolution of MU within Recommender Systems (RS) is poised at a crucial juncture, shaped by rapid advancements in ML methodologies and the escalating emphasis on core principles such as interpretability, fairness, and privacy [189], [190]. The dynamic and personalized nature of RS, coupled with the complexities of user data, necessitates the integration of MU to ensure the adaptability, accuracy, and ethical integrity of recommendations [191].

1) Adaptability and Personalization

MU's role in enhancing the adaptability and personalization of RS is paramount. With user preferences and needs continually evolving, RS must leverage MU to remain relevant and accurate. Unlearning mechanisms enable the iterative refinement of RS models, allowing them to discard outdated or irrelevant user interactions and adapt to current user behaviors and preferences [192]. This dynamic adjustment ensures that recommendations stay aligned with user interests, thereby enhancing user engagement and satisfaction.

2) Transparency, Fairness, and Explainability

The pursuit of transparency and fairness in RS is significantly bolstered by MU. By facilitating the removal of biased, misleading, or irrelevant features from the recommendation algorithms, MU contributes to the development of more transparent and fair RS. This not only improves the explainability of the recommendations provided but also ensures that the RS adheres to ethical standards, mitigating the risks of perpetuating biases or inaccuracies [193], [194].

3) Privacy Preservation

MU emerges as a critical tool in safeguarding user privacy within RS. Through the strategic removal or obfuscation of sensitive or identifying information from the training datasets, MU ensures that personal user data is handled with the utmost integrity, preventing inappropriate use or unauthorized disclosure [195].

4) Challenges in MU for RS

Implementing MU in RS entails navigating a series of challenges related to the model's complexity, data characteristics, learning algorithms, volume of data removal, and initial model training strategies [196]–[200]. The interplay of these factors determines the model's post-unlearning performance and efficiency, highlighting the need for a comprehensive approach to optimize MU.

MU in RS holds the promise of revolutionizing the way recommendations are curated and presented, ensuring they are not only reflective of current user preferences but also ethically

sound and privacy-conscious. As RS continues to evolve, the integration of MU will be crucial in overcoming the inherent challenges of adaptability, transparency, fairness, and privacy, thereby shaping a future where recommendations are not only accurate but also responsible and user-centric.

VII. CONCLUSION

MU is a relatively new and rapidly evolving field that has gained increasing attention in recent years. While the process of training ML models to recognize patterns and make predictions has become more efficient and accurate, the need to remove or modify these predictions is now equally important. Unlearning, as the name implies, refers to the process of removing previously learned information from a model, and it has important applications in areas such as privacy, security, and fairness. As our literature survey has shown, there are a variety of approaches and techniques being developed for unlearning data, ranging from regularization methods to model inversion techniques. However, there are still challenges that must be addressed in this area, such as scalability to larger datasets, the ability to unlearn specific subsets of data, and the impact of unlearning on model performance. However, despite these challenges, the benefits of MU are significant, and we expect to see continued progress in this field in the coming years as researchers develop even more effective and efficient methods for unlearning data from ML models. Researchers and practitioners must continue to explore and refine unlearning techniques to ensure that ML models can adapt to changing circumstances and maintain the trust of their users. With the increasing importance of AI in various domains, unlearning will play a crucial role in making AI more trustworthy and transparent.

REFERENCES

- [1] I. H. Sarker, "Machine learning: Algorithms, real-world applications and research directions," *SN computer science*, vol. 2, no. 3, p. 160, 2021.
- [2] T. Davenport and R. Kalakota, "The potential for artificial intelligence in healthcare," *Future healthcare journal*, vol. 6, no. 2, p. 94, 2019.
- [3] R. M. Martins and C. Gresse Von Wangenheim, "Findings on teaching machine learning in high school: A ten-year systematic literature review," *Informatics in Education*, vol. 22, no. 3, pp. 421–440, 2023.
- [4] C. Janiesch, P. Zschech, and K. Heinrich, "Machine learning and deep learning," *Electronic Markets*, vol. 31, no. 3, pp. 685–695, 2021.
- [5] Z. Izzo, M. Anne Smart, K. Chaudhuri, and J. Zou, "Approximate data deletion from machine learning models," in *Proceedings of The 24th Int'l conf. on Artificial Intelligence and Statistics* (A. Banerjee and K. Fukumizu, eds.), vol. 130 of *Proceedings of Machine Learning Research*, pp. 2008–2016, PMLR, 13–15 Apr 2021.
- [6] H. Zhang, T. Nakamura, T. Isohara, and K. Sakurai, "A review on machine unlearning," *SN Computer Science*, vol. 4, no. 4, p. 337, 2023.
- [7] Y. Cao and J. Yang, "Towards making systems forget with machine unlearning," in *2015 IEEE Symposium on Security and Privacy*, 2015.
- [8] T. Eisenhofer, D. Riepel, V. Chandrasekaran, E. Ghosh, O. Ohrimenko, and N. Papernot, "Verifiable and provably secure machine unlearning," *arXiv preprint arXiv:2210.09126*, 2022.
- [9] T. T. Nguyen, T. T. Huynh, P. L. Nguyen, A. W.-C. Liew, H. Yin, and Q. V. H. Nguyen, "A survey of machine unlearning," *arXiv preprint arXiv:2209.02299*, 2022.
- [10] A. Sekhari, J. Acharya, G. Kamath, and A. T. Suresh, "Remember what you want to forget: Algorithms for machine unlearning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 18075–18086, 2021.
- [11] R. Mehta, S. Pal, V. Singh, and S. N. Ravi, "Deep unlearning via randomized conditionally independent Hessians," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10422–10431, 2022.
- [12] S. Verkijk and P. Vossen, "Medroberta. nl: a language model for dutch electronic health records," *Computational Linguistics in the Netherlands Journal*, vol. 11, 2021.
- [13] Q. P. Nguyen, R. Oikawa, D. M. Divakaran, M. C. Chan, and B. K. H. Low, "Markov chain monte carlo-based machine unlearning: Unlearning what needs to be forgotten," in *Proceedings of the 2022 ACM on Asia Conference on Computer and Communications Security, ASIA CCS '22*, (New York, NY, USA), p. 351–363, ACM, 2022.
- [14] H. Xu, T. Zhu, L. Zhang, W. Zhou, and P. S. Yu, "Machine unlearning: A survey," *ACM Computing Surveys*, vol. 56, no. 1, pp. 1–36, 2023.
- [15] W. Wang, C. Zhang, Z. Tian, and S. Yu, "Machine unlearning via representation forgetting with parameter self-sharing," *IEEE Transactions on Information Forensics and Security*, 2023.
- [16] J. A. S. A. Guimarães, "Preserving personal dignity: the vital role of the right to be forgotten," *Brazilian Journal of Law, Technology and Innovation*, vol. 1, no. 1, pp. 163–186, 2023.
- [17] D. Zhang, S. Pan, T. Hoang, Z. Xing, M. Staples, X. Xu, L. Yao, Q. Lu, and L. Zhu, "To be forgotten or to be fair: Unveiling fairness implications of machine unlearning methods," *AI and Ethics*, pp. 1–11, 2024.
- [18] C. N. Fortner, "Decision making within a cancel culture environment," tech. rep., US Army Command and General Staff College, 2020.
- [19] J. Globocnik, "The right to be forgotten is taking shape: Cjeu judgments in gc and others (c-136/17) and google v cnil (c-507/17)," *GRUR International*, vol. 69, no. 4, pp. 380–388, 2020.
- [20] J. Weng, S. Yao, Y. Du, J. Huang, J. Weng, and C. Wang, "Proof of unlearning: Definitions and instantiation," *IEEE Transactions on Information Forensics and Security*, 2024.
- [21] S. Greengard, "Can ai learn to forget?," *Communications of the ACM*, vol. 65, no. 4, pp. 9–11, 2022.
- [22] Y. Huang and C. L. Canonne, "Tight bounds for machine unlearning via differential privacy," *arXiv preprint arXiv:2309.00886*, 2023.
- [23] A. Goyal, V. Hassija, and V. H. C. d. Albuquerque, "Revisiting machine learning training process for enhanced data privacy," in *2021 Thirteenth International Conference on Contemporary Computing (IC3-2021)*, pp. 247–251, 2021.
- [24] A. Golatkar, A. Achille, A. Ravichandran, M. Polito, and S. Soatto, "Mixed-privacy forgetting in deep networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 792–801, 2021.
- [25] D. Yu, H. Zhang, W. Chen, J. Yin, and T.-Y. Liu, "How does data augmentation affect privacy in machine learning?," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, pp. 10746–10753, 2021.
- [26] S. Shan, E. Wenger, J. Zhang, H. Li, H. Zheng, and B. Y. Zhao, "Fawkes: Protecting privacy against unauthorized deep learning models," in *29th USENIX security symposium (USENIX Security 20)*, pp. 1589–1604, 2020.
- [27] A. Thudi, G. Deza, V. Chandrasekaran, and N. Papernot, "Unrolling sgd: Understanding factors influencing machine unlearning," in *2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P)*, pp. 303–319, IEEE, 2022.
- [28] A. Golatkar, A. Achille, A. Ravichandran, M. Polito, and S. Soatto, "Mixed-privacy forgetting in deep networks," in *CVPR 2021*, 2021.
- [29] D. M. Sommer, L. Song, S. Wagh, and P. Mittal, "Athena: Probabilistic verification of machine unlearning," *Proceedings on Privacy Enhancing Technologies*, vol. 3, pp. 268–290, 2022.
- [30] N. K. Dinsdale, M. Jenkinson, and A. I. Namburete, "Fedharmony: unlearning scanner bias with distributed data," in *Medical Image Computing and Computer Assisted Intervention—MICCAI 2022: 25th Int'l conf., Singapore, September 18–22, 2022, Proceedings, Part VIII*, pp. 695–704, Springer, 2022.
- [31] A. Ashraf, S. Khan, N. Bhagwat, M. Chakravarty, and B. Taati, "Learning to unlearn: Building immunity to dataset bias in medical imaging studies," *arXiv preprint arXiv:1812.01716*, 2018.
- [32] H. He, S. Zha, and H. Wang, "Unlearn dataset bias in natural language inference by fitting the residual," in *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, pp. 132–142, 2019.
- [33] P. Bevan and A. Atapour-Abarghouei, "Skin deep unlearning: Artefact and instrument debiasing in the context of melanoma classification," in *Int'l conf. on Machine Learning*, pp. 1874–1892, PMLR, 2022.

- [34] V. Gupta, C. Jung, S. Neel, A. Roth, S. Sharifi-Malvajardi, and C. Waites, "Adaptive machine unlearning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 16319–16330, 2021.
- [35] T. Shaik, X. Tao, L. Li, H. Xie, T. Cai, X. Zhu, and Q. Li, "Framu: Attention-based machine unlearning using federated reinforcement learning," *arXiv preprint arXiv:2309.10283*, 2023.
- [36] S. Gallacher, E. Papadopoulou, N. K. Taylor, and M. H. Williams, "Learning user preferences for adaptive pervasive environments: An incremental and temporal approach," *ACM Transactions on Autonomous and Adaptive Systems (TAAAS)*, vol. 8, no. 1, pp. 1–26, 2013.
- [37] J. Foster, S. Schoepf, and A. Brintrup, "Fast machine unlearning without retraining through selective synaptic dampening," *arXiv preprint arXiv:2308.07707*, 2023.
- [38] Y. Wu, E. Dobriban, and S. Davidson, "Deltagrad: Rapid retraining of machine learning models," in *International Conference on Machine Learning*, pp. 10355–10366, PMLR, 2020.
- [39] D. Dai, L. Dong, Y. Hao, Z. Sui, B. Chang, and F. Wei, "Knowledge neurons in pretrained transformers," *arXiv preprint arXiv:2104.08696*, 2021.
- [40] J. Kim and S. S. Woo, "Efficient two-stage model retraining for machine unlearning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4361–4369, 2022.
- [41] Y. Liu, L. Xu, X. Yuan, C. Wang, and B. Li, "The right to be forgotten in federated learning: An efficient realization with rapid retraining," in *IEEE INFOCOM 2022-IEEE Conference on Computer Communications*, pp. 1749–1758, IEEE, 2022.
- [42] K. Koch and M. Soll, "No matter how you slice it: Machine unlearning with sisa comes at the expense of minority classes," in *First IEEE Conference on Secure and Trustworthy Machine Learning*.
- [43] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, et al., "A survey of large language models," *arXiv preprint arXiv:2303.18223*, 2023.
- [44] S. Kulathoor and S. Gurpur, "'right to be forgotten and re-referencing' and the freedom of speech and expression: Indian resonances to google vs. cnil," *International Journal of Early Childhood Special Education*, vol. 14, no. 4, 2022.
- [45] R. Chourasia, N. Shah, and R. Shokri, "Forget unlearning: Towards true data-deletion in machine learning," *arXiv preprint arXiv:2210.08911*, 2022.
- [46] S. Garg, S. Goldwasser, and P. N. Vasudevan, "Formalizing data deletion in the context of the right to be forgotten," in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, pp. 373–402, Springer, 2020.
- [47] J. Gao, S. Garg, M. Mahmoody, and P. N. Vasudevan, "Deletion inference, reconstruction, and compliance in machine (un) learning," *arXiv preprint arXiv:2202.03460*, 2022.
- [48] B. L. Wang and S. Schelter, "Efficiently maintaining next basket recommendations under additions and deletions of baskets and items," *arXiv preprint arXiv:2201.13313*, 2022.
- [49] A. Ginart, M. Guan, G. Valiant, and J. Y. Zou, "Making ai forget you: Data deletion in machine learning," *Advances in neural information processing systems*, vol. 32, 2019.
- [50] N. G. Marchant, B. I. Rubinstein, and S. Alfeld, "Hard to forget: Poisoning attacks on certified machine unlearning," in *Proc. of the AAAI Conf. on Artificial Intelligence*, vol. 36, pp. 7691–7700, 2022.
- [51] G. Sun, Y. Cong, J. Dong, Q. Wang, L. Lyu, and J. Liu, "Data poisoning attacks on federated machine learning," *IEEE Internet of Things Journal*, vol. 9, pp. 11365–11375, July 2022.
- [52] F. A. Yerlikaya and Ş. Bahtiyar, "Data poisoning attacks against machine learning algorithms," *Expert Systems with Applications*, vol. 208, p. 118101, Dec. 2022.
- [53] M. Anisetti, C. A. Ardagna, A. Balestrucci, N. Bena, E. Damiani, and C. Y. Yeun, "On the robustness of ensemble-based machine learning against data poisoning," 2022.
- [54] M. Maabreh, O. Darwish, O. Karajeh, and Y. Tashtoush, "On developing deep learning models with particle swarm optimization in the presence of poisoning attacks," in *2022 Int'l Arab Conference on Information Technology (ACIT)*, IEEE, Nov. 2022.
- [55] A. Thudi, I. Shumailov, F. Boenisch, and N. Papernot, "Bounding membership inference," *arXiv preprint arXiv:2202.12232*, 2022.
- [56] R. A. Fisher, F. Yates, et al., *Statistical tables for biological, agricultural and medical research, edited by ra fisher and f. yates*. Edinburgh: Oliver and Boyd, 1963.
- [57] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [58] D. P. Kingma, M. Welling, et al., "An introduction to variational autoencoders," *Foundations and Trends® in Machine Learning*, vol. 12, no. 4, pp. 307–392, 2019.
- [59] A.-r. Mohamed, G. E. Dahl, and G. Hinton, "Acoustic modeling using deep belief networks," *IEEE transactions on audio, speech, and language processing*, vol. 20, no. 1, pp. 14–22, 2011.
- [60] V. S. Chundawat, A. K. Tarun, M. Mandal, and M. Kankanalli, "Zero-shot machine unlearning," *arXiv preprint arXiv:2201.05629*, 2022.
- [61] M. Casimiro, P. Romano, D. Garlan, G. A. Moreno, E. Kang, and M. Klein, "Self-adaptation for machine learning based systems," in *ECSA (Companion)*, 2021.
- [62] T. Baumhauer, P. Schöttle, and M. Zeppelzauer, "Machine unlearning: Linear filtration for logit-based classifiers," *Machine Learning*, vol. 111, no. 9, pp. 3203–3226, 2022.
- [63] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot, "Machine unlearning," in *2021 IEEE Symposium on Security and Privacy (SP)*, May 2021.
- [64] Z. Zhou, X. Liu, J. Li, J. Ruan, and M. Fan, "Dynamically selected mixup machine unlearning," in *2022 IEEE Int'l conf. on Trust, Security and Privacy in Computing and Communications (TrustCom)*, pp. 514–524, IEEE, 2022.
- [65] M. Zhao, M. Li, S.-L. Peng, and J. Li, "A novel deep learning model compression algorithm," *Electronics*, vol. 11, no. 7, p. 1066, 2022.
- [66] S.-K. Yeom, P. Seegerer, S. Lapuschkin, A. Binder, S. Wiedemann, K.-R. Müller, and W. Samek, "Pruning by explaining: A novel criterion for deep neural network pruning," *Pattern Recognition*, vol. 115, p. 107899, 2021.
- [67] J. Wang, S. Guo, X. Xie, and H. Qi, "Federated unlearning via class-discriminative pruning," in *Proceedings of the ACM Web Conference 2022*, pp. 622–632, 2022.
- [68] J. Jia, J. Liu, P. Ram, Y. Yao, G. Liu, Y. Liu, P. Sharma, and S. Liu, "Model sparsification can simplify machine unlearning," *arXiv preprint arXiv:2304.04934*, 2023.
- [69] X. Ma, G. Yuan, X. Shen, T. Chen, X. Chen, X. Chen, N. Liu, M. Qin, S. Liu, Z. Wang, et al., "Sanity checks for lottery tickets: Does your winning ticket really win the jackpot?," *Advances in Neural Information Processing Systems*, vol. 34, pp. 12749–12760, 2021.
- [70] H. Tanaka, D. Kunin, D. L. Yamins, and S. Ganguli, "Pruning neural networks without any data by iteratively conserving synaptic flow," *Advances in neural information processing systems*, vol. 33, pp. 6377–6389, 2020.
- [71] J. Frankle and M. Carbin, "The lottery ticket hypothesis: Finding sparse, trainable neural networks," *arXiv preprint arXiv:1803.03635*, 2018.
- [72] S. Lin, X. Zhang, C. Chen, X. Chen, and W. Susilo, "Erm-ktp: Knowledge-level machine unlearning via knowledge transfer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20147–20155, 2023.
- [73] L. Wang and K.-J. Yoon, "Knowledge distillation and student-teacher learning for visual intelligence: A review and new outlooks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [74] C. Wu, S. Zhu, and P. Mitra, "Federated unlearning with knowledge distillation," *arXiv preprint arXiv:2201.09441*, 2022.
- [75] X. Zhu, G. Li, and W. Hu, "Heterogeneous federated knowledge graph embedding learning and unlearning," *arXiv preprint arXiv:2302.02069*, 2023.
- [76] Y. Li, X. Lyu, N. Koren, L. Lyu, B. Li, and X. Ma, "Neural attention distillation: Erasing backdoor triggers from deep neural networks," *arXiv preprint arXiv:2101.05930*, 2021.
- [77] Y. Guo, Y. Zhao, S. Hou, C. Wang, and X. Jia, "Verifying in the dark: Verifiable machine unlearning by using invisible backdoor triggers," *IEEE Transactions on Information Forensics and Security*, 2023.
- [78] Z. He, T. Zhang, and R. B. Lee, "Model inversion attacks against collaborative inference," in *Proceedings of the 35th Annual Computer Security Applications Conference*, pp. 148–162, 2019.
- [79] L. Graves, V. Nagisetty, and V. Ganesh, "Amnesiac machine learning," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, pp. 11516–11524, May 2021.
- [80] Y. Yoon, J. Nam, H. Yun, D. Kim, and J. Ok, "Few-shot unlearning by model inversion," *arXiv preprint arXiv:2205.15567*, 2022.
- [81] H. Xia, S. Xu, J. Pei, R. Zhang, Z. Yu, W. Zou, L. Wang, and C. Liu, "Fedme 2: Memory evaluation & erase promoting federated unlearning in dtm," *IEEE Journal on Selected Areas in Communications*, 2023.
- [82] L. Zhang, T. Zhu, H. Zhang, P. Xiong, and W. Zhou, "Fedrecovery: Differentially private machine unlearning for federated learning frameworks," *IEEE Transactions on Information Forensics and Security*, 2023.

- [83] G. Li, L. Shen, Y. Sun, Y. Hu, H. Hu, and D. Tao, "Subspace based federated unlearning," *arXiv preprint arXiv:2302.12448*, 2023.
- [84] C. Pan, J. Sima, S. Prakash, V. Rana, and O. Milenkovic, "Machine unlearning of federated clusters," *arXiv preprint arXiv:2210.16424*, 2022.
- [85] J. Gong, O. Simeone, and J. Kang, "Bayesian variational federated learning and unlearning in decentralized networks," in *2021 IEEE 22nd International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, pp. 216–220, IEEE, 2021.
- [86] J. Gong, J. Kang, O. Simeone, and R. Kassab, "Forget-svkd: Particle-based bayesian federated unlearning," in *2022 IEEE Data Science and Learning Workshop (DSLW)*, pp. 1–6, IEEE, 2022.
- [87] R. Jin, M. Chen, Q. Zhang, and X. Li, "Forgettable federated linear learning with certified data removal," *arXiv preprint arXiv:2306.02216*, 2023.
- [88] X. Gao, X. Ma, J. Wang, Y. Sun, B. Li, S. Ji, P. Cheng, and J. Chen, "Verifi: Towards verifiable federated unlearning," *arXiv preprint arXiv:2205.12709*, 2022.
- [89] W. Cong and M. Mahdavi, "Efficiently forgetting what you have learned in graph representation learning via projection," in *International Conference on Artificial Intelligence and Statistics*, pp. 6674–6703, PMLR, 2023.
- [90] W. Cong and M. Mahdavi, "Grapheditor: An efficient graph representation learning and unlearning approach," 2022.
- [91] K. Wu, J. Shen, Y. Ning, T. Wang, and W. H. Wang, "Certified edge unlearning for graph neural networks," in *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2606–2617, 2023.
- [92] C.-L. Wang, M. Huai, and D. Wang, "Inductive graph unlearning," *arXiv preprint arXiv:2304.03093*, 2023.
- [93] J. Cheng, G. Dasoulas, H. He, C. Agarwal, and M. Zitnik, "Gnndelete: A general strategy for unlearning in graph neural networks," *arXiv preprint arXiv:2302.13406*, 2023.
- [94] C. Pan, E. Chien, and O. Milenkovic, "Unlearning graph classifiers with limited data resources," in *Proceedings of the ACM Web Conference 2023*, pp. 716–726, 2023.
- [95] J. Wu, Y. Yang, Y. Qian, Y. Sui, X. Wang, and X. He, "Gif: A general graph unlearning strategy via influence function," in *Proceedings of the ACM Web Conference 2023*, pp. 651–661, 2023.
- [96] M. Chen, Z. Zhang, T. Wang, M. Backes, M. Humbert, and Y. Zhang, "Graph unlearning," in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pp. 499–513, 2022.
- [97] E. Chien, C. Pan, and O. Milenkovic, "Certified graph unlearning," *arXiv preprint arXiv:2206.09140*, 2022.
- [98] A. Golatkar, A. Achille, and S. Soatto, "Eternal sunshine of the spotless net: Selective forgetting in deep networks," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, (Los Alamitos, CA, USA), pp. 9301–9309, IEEE Computer Society, 2020.
- [99] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," *Neural networks*, vol. 113, pp. 54–71, 2019.
- [100] V. S. Chundawat, A. K. Tarun, M. Mandal, and M. Kankanhalli, "Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher," 2022.
- [101] A. Becker and T. Liebig, "Evaluating machine unlearning via epistemic uncertainty," *arXiv preprint arXiv:2208.10836*, 2022.
- [102] M. Fredrikson, S. Jha, and T. Ristenpart, "Model inversion attacks that exploit confidence information and basic countermeasures," in *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pp. 1322–1333, 2015.
- [103] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, et al., "Overcoming catastrophic forgetting in neural networks," *Proceedings of the national academy of sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [104] A. K. Tarun, V. S. Chundawat, M. Mandal, and M. Kankanhalli, "Fast yet effective machine unlearning," 2023.
- [105] K. Y. Tan, L. Yueming, Y.-S. Ong, and I. Tsang, "Unfolded self-reconstruction lsh: Towards machine unlearning in approximate nearest neighbour search," *arXiv preprint arXiv:2304.02350*, 2023.
- [106] J. Xu, Z. Wu, C. Wang, and X. Jia, "Machine unlearning: Solutions and challenges," *arXiv preprint arXiv:2308.07061*, 2023.
- [107] J. Lin, L. Dang, M. Rahouti, and K. Xiong, "ML attack models: Adversarial attacks and data poisoning attacks," *arXiv preprint arXiv:2112.02797*, 2021.
- [108] N. Carlini, F. Tramèr, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. B. Brown, D. Song, U. Erlingsson, et al., "Extracting training data from large language models," in *USENIX Security Symposium*, vol. 6, 2021.
- [109] D. M. Sommer, L. Song, S. Wagh, and P. Mittal, "Towards probabilistic verification of machine unlearning," *arXiv preprint arXiv:2003.04247*, 2020.
- [110] J. Guo, A. Li, and C. Liu, "Backdoor detection in reinforcement learning," *arXiv preprint arXiv:2202.03609*, 2022.
- [111] S. Schelter, S. Grafberger, and T. Dunning, "Hedgecut: Maintaining randomised trees for low-latency machine unlearning," in *Proceedings of the 2021 Int'l conf. on Management of Data*, pp. 1545–1557, 2021.
- [112] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.
- [113] Y.-H. Wu, C.-H. Yuan, and S.-H. Wu, "Adversarial robustness via runtime masking and cleansing," in *Int'l conf. on Machine Learning*, pp. 10399–10409, PMLR, 2020.
- [114] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork, "Learning fair representations," in *Int'l conf. on machine learning*, pp. 325–333, PMLR, 2013.
- [115] G. Pleiss, M. Raghavan, F. Wu, J. Kleinberg, and K. Q. Weinberger, "On fairness and calibration," *Advances in neural information processing systems*, vol. 30, 2017.
- [116] A. Noack, I. Ahern, D. Dou, and B. Li, "An empirical study on the relation between network interpretability and adversarial robustness," *SN Computer Science*, vol. 2, pp. 1–13, 2021.
- [117] S. Wachter, B. Mittelstadt, and L. Floridi, "Transparent, explainable, and accountable ai for robotics," *Science robotics*, vol. 2, no. 6, p. eaan6080, 2017.
- [118] A. Tahilian, V. Hassija, V. Chamola, and M. Guizani, "Machine unlearning: Its need and implementation strategies," in *the 13th Int'l conf. on Contemporary Computing (IC3-2021)*, pp. 241–246, 2021.
- [119] A. Achille, M. Kearns, C. Klingenberg, and S. Soatto, "Ai model disgorgement: Methods and choices," *arXiv preprint arXiv:2304.03545*, 2023.
- [120] R. Pradhan, J. Zhu, B. Glavic, and B. Salimi, "Interpretable data-based explanations for fairness debugging," in *Proceedings of the 2022 Int'l conf. on Management of Data*, pp. 247–261, 2022.
- [121] A. Warnecke, L. Pirch, C. Wressnegger, and K. Rieck, "Machine unlearning of features and labels," *arXiv preprint arXiv:2108.11577*, 2021.
- [122] A. Thudi, H. Jia, I. Shumailov, and N. Papernot, "On the necessity of auditable algorithmic definitions for machine unlearning," in *31st USENIX Security Symposium (USENIX Security 22)*, pp. 4007–4022, 2022.
- [123] M. C. Vu and L. A. Nguyen, "Mindful unlearning in unprecedented times: Implications for management and organizations," *Management Learning*, vol. 53, no. 5, pp. 797–817, 2022.
- [124] S. Moon, S. Cho, and D. Kim, "Feature unlearning for generative models via implicit feedback," *arXiv preprint arXiv:2303.05699*, 2023.
- [125] S. Park and E.-J. Kim, "Exploring linkages between unlearning and human resource development: Revisiting unlearning cases," *Human Resource Development Quarterly*, vol. 31, Mar. 2020.
- [126] M. Sand, J. M. Durán, and K. R. Jongsma, "Responsibility beyond design: Physicians' requirements for ethical medical AI," *Bioethics*, vol. 36, pp. 162–169, June 2021.
- [127] Z. Tóth, R. Caruana, T. Gruber, and C. Loebbecke, "The dawn of the AI robots: Towards a new framework of AI robot accountability," *Journal of Business Ethics*, vol. 178, pp. 895–916, Mar. 2022.
- [128] S. Sajid, "A methodology to build interpretable machine learning models in organizations," Master's thesis, University of Twente, 2023.
- [129] L. Takeuchi, C. K. Martin, and B. Barron, "Learning together: Adapting methods for family and community research during a pandemic," in *Joan Ganz Cooney Center at Sesame Workshop*, ERIC, 2021.
- [130] A. Hanif, X. Zhang, and S. Wood, "A survey on explainable artificial intelligence techniques and challenges," in *2021 IEEE 25th Int'l enterprise distributed object computing workshop (EDOCW)*, pp. 81–89, IEEE, 2021.
- [131] M. M. Bühler, T. Jelinek, and K. Nübel, "Training and preparing tomorrow's workforce for the fourth industrial revolution," *Education Sciences*, vol. 12, p. 782, Nov. 2022.
- [132] V.-D. Păvăloaia and S.-C. Necula, "Artificial intelligence as a disruptive technology—a systematic literature review," *Electronics*, vol. 12, p. 1102, Feb. 2023.

- [133] L. Manning, W. Morris, and I. Birchmore, "Organizational unlearning: A risky food safety strategy?," *Comprehensive Reviews in Food Science and Food Safety*, Mar. 2023.
- [134] A. Kumar, B. Finley, T. Braud, S. Tarkoma, and P. Hui, "Sketching an AI marketplace: Tech, economic, and regulatory aspects," *IEEE Access*, vol. 9, pp. 13761–13774, 2021.
- [135] R. R. Fletcher, A. Nakeshimana, and O. Olubeko, "Addressing fairness, bias, and appropriate use of artificial intelligence and machine learning in global health," 2021.
- [136] A. Ross and F. Doshi-Velez, "Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, 2018.
- [137] J. Brophy, "Exit through the training data: A look into instance-attribution explanations and efficient data deletion in machine learning," *Technical report*, 2020.
- [138] V. Buhmester, D. Münch, and M. Arens, "Analysis of explainers of black box deep neural networks for computer vision: A survey," *Machine Learning and Knowledge Extraction*, vol. 3, no. 4, pp. 966–989, 2021.
- [139] E. Sarkar and M. Maniatakos, "Trapdoor: Repurposing backdoors to detect dataset bias in machine learning-based genomic analysis," *arXiv preprint arXiv:2108.10132*, 2021.
- [140] H. Hwang, G.-H. Kim, S. Hong, and K.-E. Kim, "Variational interaction information maximization for cross-domain disentanglement," *Advances in Neural Information Processing Systems*, vol. 33, pp. 22479–22491, 2020.
- [141] Z. Ma, Y. Liu, X. Liu, J. Liu, J. Ma, and K. Ren, "Learn to forget: Machine unlearning via neuron masking," *IEEE Transactions on Dependable and Secure Computing*, 2022.
- [142] Y. Yin, B. Huang, Y. Wu, and M. Soleymani, "Speaker-invariant adversarial domain adaptation for emotion recognition," in *Proceedings of the 2020 Int'l conf. on Multimodal Interaction*, pp. 481–490, 2020.
- [143] J. Corral Acero, V. Sundaresan, N. Dinsdale, V. Grau, and M. Jenkinson, "A 2-step deep learning method with domain adaptation for multi-centre, multi-vendor and multi-disease cardiac magnetic resonance segmentation," in *Statistical Atlases and Computational Models of the Heart. M&Ms and EMIDEC Challenges: 11th Int'l Workshop, STACOM 2020, joint with MICCAI 2020, Lima, Peru, October 4, 2020, Revised Selected Papers 11*, pp. 196–207, Springer, 2021.
- [144] X. Wang, G. Chen, G. Qian, P. Gao, X.-Y. Wei, Y. Wang, Y. Tian, and W. Gao, "Large-scale multi-modal pre-trained models: A comprehensive survey," *arXiv preprint arXiv:2302.10035*, 2023.
- [145] Z. Zhao, *Using Pre-trained Language Models for Toxic Comment Classification*. PhD thesis, University of Sheffield, 2022.
- [146] W. Zhang, L. Deng, L. Zhang, and D. Wu, "A survey on negative transfer," *IEEE/CAA Journal of Automatica Sinica*, 2022.
- [147] H. Cheng, X. Liu, L. Pereira, Y. Yu, and J. Gao, "Posterior differential regularization with f-divergence for improving model robustness," in *Proc. of the 2021 Conf. of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1078–1089, ACL, 2021.
- [148] A. Kalinowska, P. M. Pilariski, and T. D. Murphey, "Embodied communication: How robots and people communicate through physical interaction," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 6, 2023.
- [149] M. Goldblum, D. Tsipras, C. Xie, X. Chen, A. Schwarzschild, D. Song, A. Madry, B. Li, and T. Goldstein, "Dataset security for machine learning: Data poisoning, backdoor attacks, and defenses," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 2, pp. 1563–1580, 2022.
- [150] J. Brophy and D. Lowd, "Machine unlearning for random forests," in *Int'l conf. on Machine Learning*, pp. 1092–1104, PMLR, 2021.
- [151] A. K. Tarun, V. S. Chundawat, M. Mandal, and M. Kankanhalli, "Deep regression unlearning," *arXiv preprint arXiv:2210.08196*, 2022.
- [152] J. Lu, T. Issarano, and D. A. Forsyth, "Feature-guided black-box safety testing of deep neural networks," in *Proc. ICCV*, pp. 446–454, 2017.
- [153] H. Jia, H. Chen, J. Guan, A. S. Shamsabadi, and N. Papernot, "A zest of lime: Towards architecture-independent model distances," in *Int'l conf. on Learning Representations*, 2022.
- [154] S. Mohseni, J. E. Block, and E. Ragan, "Quantitative evaluation of machine learning explanations: A human-grounded benchmark," in *26th Int'l conf. on Intelligent User Interfaces*, pp. 22–31, 2021.
- [155] S. Mohseni, J. E. Block, and E. D. Ragan, "A human-grounded evaluation benchmark for local explanations of machine learning," *arXiv preprint arXiv:1801.05075*, 2018.
- [156] D. Kaur, S. Uslu, K. J. Rittichier, and A. Duresi, "Trustworthy artificial intelligence: a review," *ACM Computing Surveys (CSUR)*, vol. 55, no. 2, pp. 1–38, 2022.
- [157] D. Saunders, "Domain adaptation and multi-domain adaptation for neural machine translation: A survey," *Journal of Artificial Intelligence Research*, vol. 75, pp. 351–424, 2022.
- [158] E. Strelcenia and S. Prakoornwit, "A survey on GAN techniques for data augmentation to address the imbalanced data issues in credit card fraud detection," *Machine Learning and Knowledge Extraction*, vol. 5, pp. 304–329, Mar. 2023.
- [159] P. Rana, A. Sowmya, E. Meijering, and Y. Song, "Data augmentation with improved regularisation and sampling for imbalanced blood cell image classification," *Scientific Reports*, vol. 12, Oct. 2022.
- [160] C. Fu, Y. Zheng, Y. Liu, Q. Xuan, and G. Chen, "NES-TL: Network embedding similarity-based transfer learning," *IEEE Transactions on Network Science and Engineering*, vol. 7, pp. 1607–1618, July 2020.
- [161] Z. Zhang, Y. Zhou, X. Zhao, T. Che, and L. Lyu, "Prompt certified machine unlearning with randomized gradient smoothing and quantization," in *Advances in Neural Information Processing Systems* (S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, eds.), vol. 35, pp. 13433–13455, Curran Associates, Inc., 2022.
- [162] K. Bu, Y. He, X. Jing, and J. Han, "Adversarial transfer learning for deep learning based automatic modulation classification," *IEEE Signal Processing Letters*, vol. 27, pp. 880–884, 2020.
- [163] A. K. Kar and Y. K. Dwivedi, "Theory building with big data-driven research – moving away from the "what" towards the "why"?" *International Journal of Information Management*, vol. 54, p. 102205, Oct. 2020.
- [164] A. Aljeraisy, M. Barati, O. Rana, and C. Perera, "Privacy laws and privacy by design schemes for the internet of things," *ACM Computing Surveys*, vol. 54, pp. 1–38, May 2021.
- [165] P. K. R. Maddikunta, Q.-V. Pham, D. C. Nguyen, T. Huynh-The, O. Aouedi, G. Yenduri, S. Bhattacharya, and T. R. Gadekallu, "Incentive techniques for the internet of things: a survey," *Journal of Network and Computer Applications*, p. 103464, 2022.
- [166] A. Tuladhar, D. Rajashekar, and N. D. Forkert, "Distributed learning in healthcare," in *Trends of Artificial Intelligence and Big Data for E-Health*, pp. 183–212, Springer International Publishing, 2022.
- [167] N. Rathi, I. Chakraborty, A. Kosta, A. Sengupta, A. Ankit, P. Panda, and K. Roy, "Exploring neuromorphic computing based on spiking neural networks: Algorithms to hardware," *ACM Computing Surveys*, vol. 55, pp. 1–49, Mar. 2023.
- [168] S. Mondal and Z. Rehena, "Challenges and limitations of social data analysis approaches," in *Smart Computing and Intelligence*, pp. 307–323, Springer Nature Singapore, 2022.
- [169] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?," in *Advances in Neural Information Processing Systems* (Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, eds.), vol. 27, Curran Associates, Inc., 2014.
- [170] A. Hsu, W. Khoo, N. Goyal, and M. Wainstein, "Next-generation digital ecosystem for climate data mining and knowledge discovery: A review of digital data collection technologies," *Frontiers in Big Data*, vol. 3, Sept. 2020.
- [171] N. Aldaghri, H. Mahdaviifar, and A. Beirami, "Coded machine unlearning," *IEEE Access*, vol. 9, pp. 88137–88150, 2021.
- [172] K. Chen, Y. Huang, and Y. Wang, "Machine unlearning via gan," *arXiv preprint arXiv:2111.11869*, 2021.
- [173] W. Yuan, H. Yin, F. Wu, S. Zhang, T. He, and H. Wang, "Federated unlearning for on-device recommendation," in *Proceedings of the 16th ACM Int'l conf. on Web Search and Data Mining*, pp. 393–401, 2023.
- [174] G. Cauwenberghs and T. Poggio, "Incremental and decremental support vector machine learning," *Advances in neural information processing systems*, vol. 13, 2000.
- [175] E. Romero, I. Barrio, and L. Belanche, "Incremental and decremental learning for linear support vector machines," in *Artificial Neural Networks–ICANN 2007: 17th Int'l conf., Porto, Portugal, September 9–13, 2007, Proceedings, Part I 17*, pp. 209–218, Springer, 2007.
- [176] C.-H. Tsai, C.-Y. Lin, and C.-J. Lin, "Incremental and decremental training for linear classification," in *Proc. of 20th ACM SIGKDD Int'l conf. on Knowledge discovery and data mining*, pp. 343–352, 2014.
- [177] W. Xu, J. He, and Y. Shu, "Transfer learning and deep domain adaptation," *Advances and applications in deep learning*, vol. 45, 2020.
- [178] S. Deepanjali, S. Dhivya, and S. Monica Catherine, "Efficient machine unlearning using general adversarial network," in *Artificial Intelligence Techniques for Advanced Computing Applications: Proceedings of ICACT 2020*, pp. 487–494, Springer, 2021.

- [179] J. Z. Di, J. Douglas, J. Acharya, G. Kamath, and A. Sekhari, "Hidden poison: Machine unlearning enables camouflaged poisoning attacks," in *NeurIPS ML Safety Workshop*, 2022. 40cm
- [180] T. Shaik, X. Tao, Y. Li, C. Dann, J. McDonald, P. Redmond, and L. Galligan, "A review of the trends and challenges in adopting natural language processing methods for education feedback analysis," *IEEE Access*, 2022.
- [181] K. Maharana, S. Mondal, and B. Nemade, "A review: Data pre-processing and data augmentation techniques," *Global Transitions Proceedings*, 2022.
- [182] V. B. Kumar, R. Gangadharaiah, and D. Roth, "Privacy adhering machine un-learning in nlp," *arXiv preprint arXiv:2212.09573*, 2022.
- [183] J. Stacey, Y. Belinkov, and M. Rei, "Supervising model attention with human explanations for robust natural language inference," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, pp. 11349–11357, 2022.
- [184] I. H. Sarker, "Multi-aspects ai-based modeling and adversarial learning for cybersecurity intelligence and robustness: A comprehensive overview," *Security and Privacy*, p. e295, 2022.
- [185] J. L. Lobo, J. Del Ser, A. Bifet, and N. Kasabov, "Spiking neural networks and online learning: An overview and perspectives," *Neural Networks*, vol. 121, pp. 88–100, 2020.
- [186] M. Du, Z. Chen, C. Liu, R. Oak, and D. Song, "Lifelong anomaly detection through unlearning," in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1283–1297, 2019.
- [187] J. Fan, *Machine learning and unlearning for IoT anomaly detection*. PhD thesis, 2023.
- [188] G. Li, H. Hsu, R. Marculescu, et al., "Machine unlearning for image-to-image generative models," *arXiv preprint arXiv:2402.00351*, 2024.
- [189] G. A. Kaissis, M. R. Makowski, D. Rückert, and R. F. Braren, "Secure, privacy-preserving and federated machine learning in medical imaging," *Nature Machine Intelligence*, vol. 2, no. 6, pp. 305–311, 2020.
- [190] C. Chen, F. Sun, M. Zhang, and B. Ding, "Recommendation unlearning," in *Proc. of the ACM Web Conference 2022*, pp. 2768–2777, 2022.
- [191] Y. Li, C. Chen, X. Zheng, J. Liu, and J. Wang, "Making recommender systems forget: Learning and unlearning for erasable recommendation," *Knowledge-Based Systems*, vol. 283, p. 111124, 2024.
- [192] M. Ramsar, "A discriminative account of the learning, representation and processing of inflection systems," *Language, Cognition and Neuroscience*, pp. 1–25, 2021.
- [193] H. Yin, Y. Sun, G. Xu, and E. Kanoulas, "Trustworthy recommendation and search: Introduction to the special issue-part 1," 2023.
- [194] Y. Yao, C. Wang, and H. Li, "Counterfactually evaluating explanations in recommender systems," *arXiv preprint arXiv:2203.01310*, 2022.
- [195] T. Schnitzler, *Analyzing privacy and end user information exposure in digital communication environments*. PhD thesis, Dissertation, Bochum, Ruhr-Universität Bochum, 2022, 2023.
- [196] B. Chandrasekaran, *Intelligence as adaptive behavior: An experiment in computational neuroethology*, vol. 6. Academic press, 2013.
- [197] S. Cha, S. Cho, D. Hwang, H. Lee, T. Moon, and M. Lee, "Learning to unlearn: Instance-wise unlearning for pre-trained classifiers," *arXiv preprint arXiv:2301.11578*, 2023.
- [198] S. A. Cambo and D. Gergle, "Model positionality and computational reflexivity: Promoting reflexivity in data science," in *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pp. 1–19, 2022.
- [199] J. Chen, S. Sathe, C. Aggarwal, and D. Turaga, "Outlier detection with autoencoder ensembles," in *Proceedings of the 2017 SIAM international conference on data mining*, pp. 90–98, SIAM, 2017.
- [200] E. Cagli, C. Dumas, and E. Prouff, "Convolutional neural networks with data augmentation against jitter-based countermeasures: Profiling attacks without pre-processing," in *Cryptographic Hardware and Embedded Systems—CHES 2017: 19th International Conference, Taipei, Taiwan, September 25-28, 2017, Proceedings*, pp. 45–68, Springer, 2017.

Exploring the Landscape of Machine Unlearning: A Comprehensive Survey and Taxonomy (Supplementary Material)

Thanveer Shaik, Xiaohui Tao, Haoran Xie, Lin Li, Xiaofeng Zhu, and Qing Li

SUPPLEMENTARY MATERIAL

appeared in four to six sources, and poorly popular if it was mentioned in three or fewer sources.

In this section, we present supplementary details for Table IV “Public Datasets for Machine Unlearning” in the survey, which includes an approximation of their popularity based on the frequency of their references. The popularity of datasets is classified into three categories, namely, “High”, “Moderate”, or “Low”, depending on the frequency they were referenced. Specifically, a dataset is considered highly popular if it was referred to in seven or more sources, moderately popular if it

TABLE I
PUBLIC DATASETS FOR MACHINE UNLEARNING WITH SUPPLEMENTARY INFORMATION FOR POPULARITY

Modality	Dataset	No. of Instances	No. of Attributes	Task	Popularity	References	References Count
Image	SVHN [1]	600,000	3072	Object recognition	High	[2]–[10]	9
	CIFAR-100 [11]	60,000	3072	Object recognition	High	[12]–[18]	7
	Imagenet [19]	1.2 million	1,000	Object recognition	Medium	[20]–[23]	4
	Mini-Imagenet [24]	100,000	784	Object recognition	Low	[25]	1
	LSUN [26]	1.2 million	varies	Scene recognition	Low	[27], [28]	2
	MNIST [29]	70,000	784	Object recognition	High	[20], [30]–[48]	20
Text	IMDB [49]	50,000	varies	Sentiment analysis	Medium	[50]–[54]	5
	Newsgroup [55]	19,188	varies	Text classification	Low	[56]	1
	Reuters [57]	10,788	varies	Text classification	Low	[58], [59]	2
	SQuAD [60]	100,000	Varies	Question answering	Low	[61]–[63]	3
Tabular	Adult [64]	48,842	14	Income prediction	Low	[65]–[67]	3
	Breast Cancer [68]	286	9	Cancer diagnosis	Low	[69], [70]	2
	Diabetes [71]	768	8	Diabetes diagnosis	Low	[72], [73]	2
Time series	Epileptic Seizure [74]	11,500	178	Seizure prediction	Low	[32], [75]	2
	Activity Recognition [76]	10,299	561	Activity Classification	Low	[75], [77], [78]	3
Graph	OGB [79]	1.2 million	varies	Graph classification	Low	[80]	1
	Cora [81]	2,708	1,433	Graph classification	Low	[82], [83]	2
	Yelp Dataset [84]	8,282,442	Varies	Recommendation	Low	[85], [86]	2
Computer Vision	Fashion-MNIST [87]	70,000	784	Image classification	Medium	[88]–[91]	4
	Caltech-101 [92]	9,146	Varies	Object recognition	Low	[93]–[95]	3
	COCO [96]	330,000	Varies	Object detection	Medium	[97]–[101]	5
	YouTube Faces [102]	3,425	2,622	Face recognition	Medium	[103]–[107]	5
	EuroSAT [108]	27,000	13	Land use classification	Low	[10], [109]	2
Transaction	Purchase [110]	39,624	8	Purchase prediction	Medium	[111]–[115]	5
Sequence	Human Activity Recognition [116]	10,299	561	Activity recognition	Low	[117]	1
Recommendation	MovieLens [118]	100,000	varies	Movie recommendation	High	[119]–[125]	7

REFERENCES

- [1] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [2] Z. Zhang, Y. Zhou, X. Zhao, T. Che, and L. Lyu, "Prompt certified machine unlearning with randomized gradient smoothing and quantization," in *Advances in Neural Information Processing Systems* (S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, eds.), vol. 35, pp. 13433–13455, Curran Associates, Inc., 2022.
- [3] J. Jia, J. Liu, P. Ram, Y. Yao, G. Liu, Y. Liu, P. Sharma, and S. Liu, "Model sparsification can simplify machine unlearning," *arXiv preprint arXiv:2304.04934*, 2023.
- [4] B. Casella, A. B. Chisari, S. Battiato, and M. V. Giuffrida, "Transfer learning via test-time neural networks aggregation," *arXiv preprint arXiv:2206.13399*, 2022.
- [5] A. Setlur, B. Eysenbach, V. Smith, and S. Levine, "Adversarial unlearning: Reducing confidence along adversarial directions," *arXiv preprint arXiv:2206.01367*, 2022.
- [6] Y. Li, Y. Jiang, Z. Li, and S.-T. Xia, "Backdoor learning: A survey," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [7] L. Feng, Z. Qian, S. Li, and X. Zhang, "Trojaning semi-supervised learning model via poisoning wild images on the web," *arXiv preprint arXiv:2301.00435*, 2023.
- [8] A. S. Rakin, Z. He, and D. Fan, "Tbt: Targeted neural network attack with bit trojan," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13198–13207, 2020.
- [9] P. Pathak, J. Zhang, and D. Samaras, "Local learning on transformers via feature reconstruction," *arXiv preprint arXiv:2212.14215*, 2022.
- [10] J. Loedeman, M. C. Stol, T. Han, and Y. M. Asano, "Prompt generation networks for efficient adaptation of frozen vision transformers," *arXiv preprint arXiv:2210.06466*, 2022.
- [11] A. Krizhevsky, G. Hinton, *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [12] Y. Dukler, B. Bowman, A. Achille, A. Goltkar, A. Swaminathan, and S. Soatto, "Safe: Machine unlearning with shard graphs," *arXiv preprint arXiv:2304.13169*, 2023.
- [13] A. Thudi, G. Deza, V. Chandrasekaran, and N. Papernot, "Unrolling sgd: Understanding factors influencing machine unlearning," in *2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P)*, pp. 303–319, IEEE, 2022.
- [14] A. Mahadevan and M. Mathioudakis, "Certifiable unlearning pipelines for logistic regression: An experimental study," *Machine Learning and Knowledge Extraction*, vol. 4, no. 3, pp. 591–620, 2022.
- [15] Y. Jiang, S. Liu, T. Zhao, W. Li, and X. Gao, "Machine unlearning survey," in *Fifth Int'l conf. on Mechatronics and Computer Technology Engineering (MCTE 2022)*, vol. 12500, pp. 1596–1603, SPIE, 2022.
- [16] M. Chen, W. Gao, G. Liu, K. Peng, and C. Wang, "Boundary unlearning," *arXiv preprint arXiv:2303.11570*, 2023.
- [17] X. Gao, X. Ma, J. Wang, Y. Sun, B. Li, S. Ji, P. Cheng, and J. Chen, "Verifi: Towards verifiable federated unlearning," *arXiv preprint arXiv:2205.12709*, 2022.
- [18] N. Carlini, F. Tramèr, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. B. Brown, D. Song, U. Erlingsson, *et al.*, "Extracting training data from large language models," in *USENIX Security Symposium*, vol. 6, 2021.
- [19] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255, IEEE, 2009.
- [20] S. Poppi, S. Sarto, M. Cornia, L. Baraldi, and R. Cucchiara, "Multi-class explainable unlearning for image classification via weight filtering," *arXiv preprint arXiv:2304.02049*, 2023.
- [21] S. Cha, S. Cho, D. Hwang, H. Lee, T. Moon, and M. Lee, "Learning to unlearn: Instance-wise unlearning for pre-trained classifiers," *arXiv preprint arXiv:2301.11578*, 2023.
- [22] W. Chen, B. Wu, and H. Wang, "Effective backdoor defense by exploiting sensitivity of poisoned samples," *Advances in Neural Information Processing Systems*, vol. 35, pp. 9727–9737, 2022.
- [23] S. Shan, W. Ding, E. Wenger, H. Zheng, and B. Y. Zhao, "Post-breach recovery: Protection against white-box adversarial examples for leaked dnn models," in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pp. 2611–2625, 2022.
- [24] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, *et al.*, "Matching networks for one shot learning," *Advances in neural information processing systems*, vol. 29, 2016.
- [25] M. Riemer, I. Cases, R. Ajemian, M. Liu, I. Rish, Y. Tu, and G. Tesauro, "Learning to learn without forgetting by maximizing transfer and minimizing interference," *arXiv preprint arXiv:1810.11910*, 2018.
- [26] F. Yu, Y. Zhang, S. Song, A. Seff, and J. Xiao, "Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop," *arXiv preprint arXiv:1506.03365*, 2015.
- [27] Z. Li, A. Hoogs, and C. Xu, "Discover and mitigate unknown biases with debiasing alternate networks," in *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIII*, pp. 270–288, Springer, 2022.
- [28] Y. Zhou and I. Nwogu, "Learning to generate high resolution images with bilateral adversarial networks," in *Proceedings of the Int'l conf. on Advances in Image Processing*, pp. 113–117, 2017.
- [29] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [30] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot, "Machine unlearning," in *2021 IEEE Symposium on Security and Privacy (SP)*, May 2021.
- [31] V. Gupta, C. Jung, S. Neel, A. Roth, S. Shariif-Malvajerdi, and C. Waites, "Adaptive machine unlearning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 16319–16330, 2021.
- [32] T. T. Nguyen, T. T. Huynh, P. L. Nguyen, A. W.-C. Liew, H. Yin, and Q. V. H. Nguyen, "A survey of machine unlearning," *arXiv preprint arXiv:2209.02299*, 2022.
- [33] D. M. Sommer, L. Song, S. Wagh, and P. Mittal, "Towards probabilistic verification of machine unlearning," *arXiv preprint arXiv:2003.04247*, 2020.
- [34] V. S. Chundawat, A. K. Tarun, M. Mandal, and M. Kankanalli, "Zero-shot machine unlearning," *arXiv preprint arXiv:2201.05629*, 2022.
- [35] T. Baumhauer, P. Schöttle, and M. Zeppelzauer, "Machine unlearning: Linear filtration for logit-based classifiers," *Machine Learning*, vol. 111, no. 9, pp. 3203–3226, 2022.
- [36] N. G. Marchant, B. I. Rubinstein, and S. Alfeld, "Hard to forget: Poisoning attacks on certified machine unlearning," in *Proc. of the AAAI Conf. on Artificial Intelligence*, vol. 36, pp. 7691–7700, 2022.
- [37] Z. Lu, H. Liang, M. Zhao, Q. Lv, T. Liang, and Y. Wang, "Label-only membership inference attacks on machine unlearning without dependence of posteriors," *Int'l Journal of Intelligent Systems*, vol. 37, no. 11, pp. 9424–9441, 2022.
- [38] P.-F. Zhang, G. Bai, Z. Huang, and X.-S. Xu, "Machine unlearning for image retrieval: A generative scrubbing approach," in *Proceedings of the 30th ACM Int'l conf. on Multimedia*, pp. 237–245, 2022.
- [39] A. Becker and T. Liebig, "Evaluating machine unlearning via epistemic uncertainty," *arXiv preprint arXiv:2208.10836*, 2022.
- [40] Z. Zhou, X. Liu, J. Li, J. Ruan, and M. Fan, "Dynamically selected mixup machine unlearning," in *2022 IEEE Int'l conf. on Trust, Security and Privacy in Computing and Communications (TrustCom)*, pp. 514–524, IEEE, 2022.
- [41] K. Chen, Y. Wang, and Y. Huang, "Lightweight machine unlearning in neural network," *arXiv preprint arXiv:2111.05528*, 2021.
- [42] Y. Yoon, J. Nam, H. Yun, D. Kim, and J. Ok, "Few-shot unlearning by model inversion," *arXiv preprint arXiv:2205.15567*, 2022.
- [43] A. Halimi, S. Kadhe, A. Rawat, and N. Baracaldo, "Federated unlearning: How to efficiently erase a client in fl?," *arXiv preprint arXiv:2207.05521*, 2022.
- [44] C. Wu, S. Zhu, and P. Mitra, "Federated unlearning with knowledge distillation," *arXiv preprint arXiv:2201.09441*, 2022.
- [45] Y. Fraboni, R. Vidal, L. Kamani, and M. Lorenzi, "Sequential informed federated unlearning: Efficient and provable client unlearning in federated optimization," *arXiv preprint arXiv:2211.11656*, 2022.
- [46] Z. Lu, Y. Wang, Q. Lv, M. Zhao, and T. Liang, "Fp 2-mia: A membership inference attack free of posterior probability in machine unlearning," in *Provable and Practical Security: 16th Int'l conf., ProvSec 2022, Nanjing, China, November 11–12, 2022, Proceedings*, pp. 167–175, Springer, 2022.
- [47] J. Gao, S. Garg, M. Mahmoody, and P. N. Vasudevan, "Deletion inference, reconstruction, and compliance in machine (un) learning," *arXiv preprint arXiv:2202.03460*, 2022.
- [48] J. Stock, J. Wettlaufer, D. Demmler, and H. Federrath, "Property unlearning: A defense strategy against property inference attacks," *arXiv preprint arXiv:2205.08821*, 2022.
- [49] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics*:

- Human Language Technologies*, (Portland, Oregon, USA), pp. 142–150, Association for Computational Linguistics, June 2011.
- [50] E. Chien, C. Pan, and O. Milenkovic, “Certified graph unlearning,” *arXiv preprint arXiv:2206.09140*, 2022.
- [51] R. Zhu, D. Tang, S. Tang, X. Wang, and H. Tang, “Selective amnesia: On efficient, high-fidelity and blind suppression of backdoor effects in trojaned machine learning models,” *arXiv preprint arXiv:2212.04687*, 2022.
- [52] M. Abolfazli, A. Host-Madsen, and J. Zhang, “Differential description length for hyperparameter selection in machine learning,” *arXiv preprint arXiv:1902.04699*, 2019.
- [53] A. R. Nelakurthi, R. Maciejewski, and J. He, “Source free domain adaptation using an off-the-shelf classifier,” in *2018 IEEE Int'l conf. on big data (Big Data)*, pp. 140–145, IEEE, 2018.
- [54] D. Friedman, A. Wettig, and D. Chen, “Finding dataset shortcuts with grammar induction,” *arXiv preprint arXiv:2210.11560*, 2022.
- [55] T. Joachims, “A probabilistic analysis of the rocchio algorithm with tfidf for text categorization,” tech. rep., Carnegie-mellon univ pittsburgh pa dept of computer science, 1996.
- [56] H. Oxhammar, “Document classification with uncertain data using statistical and rule-based methods,” *Uppsala University: Department of Linguistics*, 2003.
- [57] C. Apté, F. Damerou, and S. M. Weiss, “Automated learning of decision rules for text categorization,” *ACM Transactions on Information Systems (TOIS)*, vol. 12, no. 3, pp. 233–251, 1994.
- [58] R. Liere, *Active learning with committees: An approach to efficient learning in text categorization using linear threshold algorithms*. Oregon State University, 2000.
- [59] D. H. Widyantoro, T. R. Ioerger, and J. Yen, “Learning user interest dynamics with a three-descriptor representation,” *Journal of the American Society for Information Science and Technology*, vol. 52, no. 3, pp. 212–225, 2001.
- [60] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, “Squad: 100,000+ questions for machine comprehension of text,” *arXiv preprint arXiv:1606.05250*, 2016.
- [61] G. Shen, Y. Liu, G. Tao, Q. Xu, Z. Zhang, S. An, S. Ma, and X. Zhang, “Constrained optimization with dynamic bound-scaling for effective nlp backdoor defense,” in *Int'l conf. on Machine Learning*, pp. 19879–19892, PMLR, 2022.
- [62] R. Zellers, A. Holtzman, E. Clark, L. Qin, A. Farhadi, and Y. Choi, “Turingadvice: A generative and dynamic evaluation of language use,” *arXiv preprint arXiv:2004.03607*, 2020.
- [63] A. Ravichander, Y. Belinkov, and E. Hovy, “Probing the probing paradigm: Does probing accuracy entail task relevance?,” *arXiv preprint arXiv:2005.00719*, 2020.
- [64] R. Kohavi *et al.*, “Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid,” in *Kdd*, vol. 96, pp. 202–207, 1996.
- [65] H. Zeng, Z. Yue, Z. Kou, Y. Zhang, L. Shang, and D. Wang, “Fairness-aware training of face attribute classifiers via adversarial robustness,” *Knowledge-Based Systems*, vol. 264, p. 110356, 2023.
- [66] X. Han, T. Baldwin, and T. Cohn, “Everybody needs good neighbours: An unsupervised locality-based method for bias mitigation,” in *The Eleventh Int'l conf. on Learning Representations*, 2023.
- [67] M. Du, R. Tang, W. Fu, and X. Hu, “Towards debiasing dnn models from spurious feature influence,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, pp. 9521–9528, 2022.
- [68] R. S. Michalski, I. Mozetic, J. Hong, and N. Lavrac, “The multi-purpose incremental learning system aq15 and its testing application to three medical domains,” in *Proc. AAAI*, vol. 1986, pp. 1–041, 1986.
- [69] K. Sullivan, A. ElMolla, B. Squires, and S. Luke, “Unlearning from demonstration,” in *Twenty-Third Int'l Joint Conference on Artificial Intelligence*, Citeseer, 2013.
- [70] S. Gallacher, E. Papadopoulou, N. K. Taylor, and M. H. Williams, “Learning user preferences for adaptive pervasive environments: An incremental and temporal approach,” *ACM Transactions on Autonomous and Adaptive Systems (TAAS)*, vol. 8, no. 1, pp. 1–26, 2013.
- [71] J. W. Smith, J. E. Everhart, W. Dickson, W. C. Knowler, and R. S. Johannes, “Using the adap learning algorithm to forecast the onset of diabetes mellitus,” in *Proceedings of the annual symposium on computer application in medical care*, p. 261, American Medical Informatics Association, 1988.
- [72] A. Tuladhar, D. Rajashekar, and N. D. Forkert, “Distributed learning in healthcare,” *Trends of Artificial Intelligence and Big Data for E-Health*, pp. 183–212, 2023.
- [73] A. Abolfazli and E. Ntoutsis, “Drift-aware multi-memory model for imbalanced data streams,” in *2020 IEEE Int'l conf. on Big Data (Big Data)*, pp. 878–885, IEEE, 2020.
- [74] R. G. Andrzejak, K. Lehnertz, F. Mormann, C. Rieke, P. David, and C. E. Elger, “Indications of nonlinear deterministic and finite-dimensional structures in time series of brain electrical activity: Dependence on recording region and brain state,” *Physical Review E*, vol. 64, no. 6, p. 061907, 2001.
- [75] V. S. Chundawat, A. K. Tarun, M. Mandal, and M. Kankanhalli, “Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher,” 2022.
- [76] D. Anguita, A. Ghio, L. Oneto, X. Parra, J. L. Reyes-Ortiz, *et al.*, “A public domain dataset for human activity recognition using smartphones,” in *Esann*, vol. 3, p. 3, 2013.
- [77] X. Cao, J. Jia, Z. Zhang, and N. Z. Gong, “Fedrecover: Recovering from poisoning attacks in federated learning using historical information,” *arXiv preprint arXiv:2210.10936*, 2022.
- [78] L. Song, *Understanding and Measuring Privacy Risks in Machine Learning*. PhD thesis, Princeton University, 2021.
- [79] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec, “Open graph benchmark: Datasets for machine learning on graphs,” *Advances in neural information processing systems*, vol. 33, pp. 22118–22133, 2020.
- [80] P. Tarau, “A gaze into the internal logic of graph neural networks, with logic,” *arXiv preprint arXiv:2208.03093*, 2022.
- [81] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, “Collective classification in network data,” *AI magazine*, vol. 29, no. 3, pp. 93–93, 2008.
- [82] H. Hu, Z. Salcic, L. Sun, G. Dobbie, P. S. Yu, and X. Zhang, “Membership inference attacks on machine learning: A survey,” *ACM Computing Surveys (CSUR)*, vol. 54, no. 11s, pp. 1–37, 2022.
- [83] E. Chien, C. Pan, and O. Milenkovic, “Efficient model updates for approximate unlearning of graph-structured data,” in *The Eleventh Int'l conf. on Learning Representations*, 2023.
- [84] X. Zhang, J. Zhao, and Y. LeCun, “Character-level convolutional networks for text classification,” *Advances in neural information processing systems*, vol. 28, 2015.
- [85] D. M. Sommer, L. Song, S. Wagh, and P. Mittal, “Athena: Probabilistic verification of machine unlearning,” *Proceedings on Privacy Enhancing Technologies*, vol. 3, pp. 268–290, 2022.
- [86] M. A. A. Gonzalez, M. J. Abe, B. F. Amadeu, S. L. Sayuri, *et al.*, “Pann component for use in pattern recognition in medical diagnostics decision-making,” *Procedia Computer Science*, vol. 192, pp. 1750–1759, 2021.
- [87] H. Xiao, K. Rasul, and R. Vollgraf, “Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms,” *arXiv preprint arXiv:1708.07747*, 2017.
- [88] K. Chen, Y. Huang, and Y. Wang, “Machine unlearning via gan,” *arXiv preprint arXiv:2111.11869*, 2021.
- [89] Y. Liu, G. Shen, G. Tao, Z. Wang, S. Ma, and X. Zhang, “Complex backdoor detection by symmetric feature differencing,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15003–15013, 2022.
- [90] J. H. Lee, “Dyngmat, a network that can learn after learning,” *Neural Networks*, vol. 116, pp. 88–100, 2019.
- [91] S. Lee, W. Song, S. Jana, M. Cha, and S. Son, “Evaluating the robustness of trigger set-based watermarks embedded in deep neural networks,” *IEEE Transactions on Dependable and Secure Computing*, 2022.
- [92] L. Fei-Fei, R. Fergus, and P. Perona, “Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories,” in *2004 conference on computer vision and pattern recognition workshop*, pp. 178–178, IEEE, 2004.
- [93] M. Du, Z. Chen, C. Liu, R. Oak, and D. Song, “Lifelong anomaly detection through unlearning,” in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1283–1297, 2019.
- [94] S. S. Shivagunde and V. V. Saradhi, “View incremental decremental multi-view discriminant analysis,” *Applied Intelligence*, pp. 1–15, 2022.
- [95] T. Diethe and M. Girolami, “Online learning with (multiple) kernels: A review,” *Neural computation*, vol. 25, no. 3, pp. 567–625, 2013.
- [96] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pp. 740–755, Springer, 2014.
- [97] E. Zhang, K. Wang, X. Xu, Z. Wang, and H. Shi, “Forget-me-not: Learning to forget in text-to-image diffusion models,” *arXiv preprint arXiv:2303.17591*, 2023.

- [98] S. Yu, J. Guo, R. Zhang, Y. Fan, Z. Wang, and X. Cheng, "A rebalancing strategy for class-imbalanced classification based on instance difficulty," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 70–79, 2022.
- [99] A. Hoblitzell, M. Babbar-Sebens, and S. Mukhopadhyay, "Uncertainty-based deep learning networks for limited data wetland user models," in *2018 IEEE Int'l conf. on Artificial Intelligence and Virtual Reality (AIVR)*, pp. 19–26, IEEE, 2018.
- [100] J. Autz, S. K. Mishra, L. Herrmann, and J. Hertzberg, "The pitfalls of transfer learning in computer vision for agriculture," *42. GIL-Jahrestagung, Künstliche Intelligenz in der Agrar-und Ernährungswirtschaft*, 2022.
- [101] H. Bansal, N. Singhi, Y. Yang, F. Yin, A. Grover, and K.-W. Chang, "Cleanclip: Mitigating data poisoning attacks in multimodal contrastive learning," in *ICLR 2023 Workshop on Trustworthy and Reliable Large-Scale Machine Learning Models*.
- [102] L. Wolf, T. Hassner, and I. Maoz, "Face recognition in unconstrained videos with matched background similarity," in *CVPR 2011*, pp. 529–534, IEEE, 2011.
- [103] A. Bendale and T. Boulton, "Reliable posterior probability estimation for streaming face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 56–63, 2014.
- [104] S. Kaviani and I. Sohn, "Defense against neural trojan attacks: A survey," *Neurocomputing*, vol. 423, pp. 651–667, 2021.
- [105] Y. Zhao, H. Zhu, K. Chen, and S. Zhang, "Ai-lancet: Locating error-inducing neurons to optimize neural networks," in *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pp. 141–158, 2021.
- [106] A. K. Veldanda, K. Liu, B. Tan, P. Krishnamurthy, F. Khorrani, R. Karri, B. Dolan-Gavitt, and S. Garg, "Nnoculation: Catching badnets in the wild," in *Proceedings of the 14th ACM Workshop on Artificial Intelligence and Security*, pp. 49–60, 2021.
- [107] I. Serma, A. Morales, J. Fierrez, and N. Obradovich, "Sensitive loss: Improving accuracy and fairness of face representations with discrimination-aware deep learning," *Artificial Intelligence*, vol. 305, p. 103682, 2022.
- [108] P. Helber, B. Bischke, A. Dengel, and D. Borth, "Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 12, no. 7, pp. 2217–2226, 2019.
- [109] L. Risser, A. Picard, L. Hervier, and J.-M. Loubes, "A survey of identification and mitigation of machine learning algorithmic biases in image analysis," *arXiv preprint arXiv:2210.04491*, 2022.
- [110] C. O. Sakar, S. O. Polat, M. Katircioglu, and Y. Kastro, "Real-time prediction of online shoppers' purchasing intention using multilayer perceptron and lstm recurrent neural networks," *Neural Computing and Applications*, vol. 31, pp. 6893–6908, 2019.
- [111] S. Mercuri, R. Khraishi, R. Okhrati, D. Batra, C. Hamill, T. Ghasempour, and A. Nowlan, "An introduction to machine unlearning," *arXiv preprint arXiv:2209.00939*, 2022.
- [112] J. Weng, S. Yao, Y. Du, J. Huang, J. Weng, and C. Wang, "Proof of unlearning: Definitions and instantiation," *arXiv preprint arXiv:2210.11334*, 2022.
- [113] K. Koch and M. Soll, "No matter how you slice it: Machine unlearning with sisa comes at the expense of minority classes," in *First IEEE Conference on Secure and Trustworthy Machine Learning*.
- [114] Y. He, G. Meng, K. Chen, J. He, and X. Hu, "Deepoblivate: a powerful charm for erasing data residual memory in deep neural networks," *arXiv preprint arXiv:2105.06209*, 2021.
- [115] N. Su and B. Li, "Asynchronous federated unlearning."
- [116] D. Anguita, A. Ghio, L. Oneto, X. Parra, and J. L. Reyes-Ortiz, "Human activity recognition on smartphones using a multiclass hardware-friendly support vector machine," in *Ambient Assisted Living and Home Care: 4th Int'l Workshop, IWAAL 2012, Vitoria-Gasteiz, Spain, December 3-5, 2012. Proceedings 4*, pp. 216–223, Springer, 2012.
- [117] M. N. K. Boulos and S. P. Yang, "Mobile physical activity planning and tracking: a brief overview of current options and desiderata for future solutions," *Mhealth*, vol. 7, 2021.
- [118] F. M. Harper and J. A. Konstan, "The movielens datasets: History and context," *ACM Trans. Interact. Intell. Syst.*, vol. 5, Dec. 2015.
- [119] Z. Chen, F. Sun, Y. Tang, H. Chen, J. Gao, and B. Ding, "Studying the impact of data disclosure mechanism in recommender systems via simulation," *ACM Transactions on Information Systems*, vol. 41, no. 3, pp. 1–26, 2023.
- [120] Z. Chen, F. Sun, Y. Tang, H. Chen, J. Gao, and B. Ding, "Proactively control privacy in recommender systems," *arXiv preprint arXiv:2204.00279*, 2022.
- [121] C. Ganhör, D. Penz, N. Rekabsaz, O. Lesota, and M. Schedl, "Unlearning protected user attributes in recommendations with adversarial training," in *Proceedings of the 45th Int'l ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 2142–2147, 2022.
- [122] Y. Yao, C. Wang, and H. Li, "Learning to counterfactually explain recommendations," *arXiv preprint arXiv:2211.09752*, 2022.
- [123] W. Yuan, H. Yin, F. Wu, S. Zhang, T. He, and H. Wang, "Federated unlearning for on-device recommendation," in *Proceedings of the 16th ACM Int'l conf. on Web Search and Data Mining*, pp. 393–401, 2023.
- [124] M. Xu, J. Sun, X. Yang, K. Yao, and C. Wang, "Netflix and forget: Efficient and exact machine unlearning from bi-linear recommendations," *arXiv preprint arXiv:2302.06676*, 2023.
- [125] B. Omidvar-Tehrani and S. Amer-Yahia, "User group analytics survey and research opportunities," *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 10, pp. 2040–2059, 2019.

10.2 Summary

The chapter concludes by synthesizing the state-of-the-art in MU, highlighting its significance in addressing privacy, security, and ethical considerations in AI. It underscores the challenges faced in MU, such as attack sophistication, lack of standardization, and resource constraints, and proposes future directions for research. The summary emphasizes the need for continued innovation in MU techniques to ensure AI systems are not only efficient and effective but also align with evolving privacy norms and ethical standards.

CHAPTER 11: PAPER 10 - FRAMU: ATTENTION-BASED MACHINE UNLEARNING USING FEDERATED REINFORCEMENT LEARNING

11.1 Introduction

This chapter introduces the FRAMU framework, an innovative solution that combines federated learning with reinforcement learning and attention mechanisms to address the challenges of machine unlearning in dynamic data environments. It delves into the necessity for models that can adaptively forget outdated, irrelevant, or private data, ensuring data privacy and model efficiency. The chapter explores FRAMU's unique approach to handling single-modality and multimodality data, emphasizing its potential to revolutionize privacy-preserving machine learning by enabling models to dynamically adapt to changing data distributions and privacy requirements.

FRAMU: Attention-based Machine Unlearning using Federated Reinforcement Learning

Thanveer Shaik*, Xiaohui Tao*, Lin Li, Haoran Xie, Taotao Cai, Xiaofeng Zhu, and Qing Li

Abstract—Machine Unlearning, a pivotal field addressing data privacy in machine learning, necessitates efficient methods for the removal of private or irrelevant data. In this context, significant challenges arise, particularly in maintaining privacy and ensuring model efficiency when managing outdated, private, and irrelevant data. Such data not only compromises model accuracy but also burdens computational efficiency in both learning and unlearning processes. To mitigate these challenges, we introduce a novel framework, Attention-based Machine Unlearning using Federated Reinforcement Learning (FRAMU). This framework incorporates adaptive learning mechanisms, privacy preservation techniques, and optimization strategies, making it a well-rounded solution for handling various data sources, either single-modality or multi-modality, while maintaining accuracy and privacy. FRAMU’s strength lies in its adaptability to fluctuating data landscapes, its ability to unlearn outdated, private, or irrelevant data, and its support for continual model evolution without compromising privacy. Our experiments, conducted on both single-modality and multi-modality datasets, revealed that FRAMU significantly outperformed baseline models. Additional assessments of convergence behaviour and optimization strategies further validate the framework’s utility in federated learning applications. Overall, FRAMU advances Machine Unlearning by offering a robust, privacy-preserving solution that optimizes model performance while also addressing key challenges in dynamic data environments.

Index Terms—Machine Unlearning, Privacy, Reinforcement Learning, Federated Learning, Attention Mechanism.

I. INTRODUCTION

The widespread availability of decentralized and heterogeneous data sources has created a demand for Machine Learning models that can effectively leverage this data while preserving privacy and ensuring accuracy [1]. Traditional approaches struggle to handle the continual influx of new data streams, and the accumulation of outdated or irrelevant information hinders their adaptability in dynamic data environments [2], [3]. Moreover, the presence of sensitive or private data introduces concerns regarding data breaches and unauthorized

*Corresponding authors: Thanveer Shaik (email: Thanveer.Shaik@usq.edu.au) and Xiaohui Tao (email: Xiaohui.Tao@usq.edu.au) are with the School of Mathematics, Physics, and Computing, University of Southern Queensland, Queensland, Australia

Taotao Cai is with the School of Mathematics, Physics, and Computing, University of Southern Queensland, Queensland, Australia (e-mail: Taotao.Cai@usq.edu.au).

Lin Li is with the School of Computer and Artificial Intelligence, Wuhan University of Technology, China (e-mail: cathylin@whut.edu.cn)

Haoran Xie is with the Department of Computing and Decision Sciences, Lingnan University, Tuen Mun, Hong Kong (e-mail: hrxie@ln.edu.hk)

Xiaofeng Zhu is with the University of Electronic Science and Technology of China (e-mail: seanzhuf@gmail.com)

Qing Li is with the Department of Computing, Hong Kong Polytechnic University, Hong Kong Special Administrative Region of China (e-mail: qing-li@polyu.edu.hk).

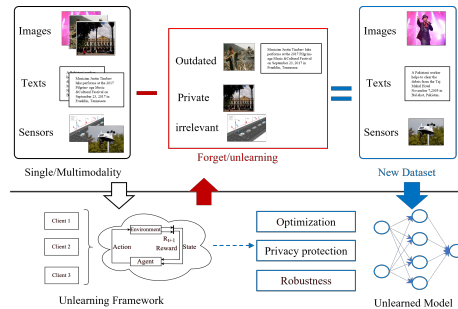


Fig. 1: Graphical abstract depicts the evolution of the FRAMU framework

access, necessitating the development of privacy-preserving techniques [4]. The concept of the “right to be forgotten” allows individuals to have their personal information removed from online platforms, although there’s no universal agreement on its definition or its status as a human right [5]. Despite this, countries like Argentina, the Philippines, and large parts of the EU are working on regulations¹. Therefore, there is a pressing need to advance the field of Machine Unlearning to ensure both adaptability and privacy in Machine Learning applications.

Example 1. In a landmark 2014 decision that underscored the pressing need for Machine Unlearning, a Spanish court ruled in favor of an individual who sought the removal of specific, outdated Google search results related to a long-settled debt [6]. This verdict not only led to Google taking down the search results but also influenced broader European Union policies on the subject, emphasizing the urgent need for mechanisms that can efficiently erase outdated or private information from Machine Learning models without sacrificing accuracy. This critical requirement for Machine Unlearning is further highlighted by high-profile cases such as that of James Gunn, the famed writer and director, who was dismissed by Disney in 2018 when old, inappropriate tweets resurfaced [7]. Although social media platforms like Facebook offer features like “Off-Facebook Activity” to disconnect user data from third-party services, this does not guarantee the complete erasure of that data from the internet². Together, these instances accentuate the growing imperative for the development of robust Machine Unlearning technologies, especially in an era where data privacy regulations are continuously evolving

¹<https://link.library.eui.eu/portal/The-Right-To-Be-Forgotten-A-Comparative-Study/tw0VHCyGcDc/>

²<https://www.facebook.com/help/2207256696182627>

and the "right to be forgotten" is increasingly recognized as essential.

Challenges. In today's digitally connected environment, data is distributed in various forms and from different sources, such as sensors, text documents, images, and time series data. For unlearning outdated or private data, Machine Unlearning presents unique challenges depending on whether it is a single type of data (known as single-modality) or multiple types of data (referred to as multimodality) [8]. With single-modality data, the issue primarily lies in the build-up of outdated or irrelevant information, which can negatively affect the model's effectiveness and precision [9], [10]. On the other hand, multimodality situations are even more complicated. Here, each type of data can have different characteristics and varying contributions to the overall model's performance [11], [12]. As we discussed in example 1, the need to unlearn outdated or private data is most important. This ensures individuals have the "right to be forgotten" about their information in publicly available platforms. However, the unlearning needs to happen in both single-modality and multimodality data to make it a holistic unlearning.

Distributed learning systems, particularly federated learning, have made significant strides forward in enabling Machine Learning models to train on decentralized data, offering the dual advantage of reduced communication costs and enhanced privacy [13], [14]. Notable efforts have been made to incorporate Differential Privacy (DP) into these systems [15], ensuring robust privacy safeguards through techniques like DP-SGD and DP-FedAvg [16], [17]. However, these existing frameworks face limitations when confronted with the dynamic nature of data distribution, an intrinsic challenge in distributed learning [18]. Although some efforts have been made in Machine Unlearning to address data irrelevancy over time, such as Sharded, Isolated, Sliced, and Aggregated(SISA) training methods, these solutions often operate in isolation from privacy-preserving mechanisms [19], [20]. This bifurcation leaves a crucial research gap: the absence of a unified approach that addresses both privacy concerns and the adaptability requirements in the face of ever-changing data landscapes. There is a need to bridge this gap by providing an integrated solution for robust privacy measures and efficient selective unlearning, thereby enabling Machine Learning models to be both secure and adaptable in dynamic, distributed environments.

The primary challenges in Machine Unlearning involve addressing the buildup of outdated or irrelevant information in single-modality data, which affects model precision, and handling the complexity of multimodality data where each type contributes differently to model performance. Additionally, current distributed learning systems, while advancing privacy and reducing communication costs, face limitations in adapting to dynamic data distributions and integrating robust privacy measures with efficient unlearning mechanisms, highlighting a need for a unified approach that ensures both security and adaptability in rapidly evolving data environments.

To address these challenges, we propose an Attention-based Machine Unlearning using Federated Reinforcement Learning (FRAMU) as shown in Fig. 1. By integrating federated learn-

ing, adaptive learning mechanisms, and privacy preservation techniques, FRAMU aims to leverage the diverse and dynamic nature of data in both single-modality and multimodality scenarios, while upholding privacy regulations and optimizing the learning process. An attention mechanism is incorporated into FRAMU to ensure responsible and secure handling of sensitive information across modalities. FRAMU leverages reinforcement learning and adaptive learning mechanisms to enable models to dynamically adapt to changing data distributions and individual participant characteristics in both single-modality and multimodality scenarios. This adaptability facilitates ongoing model evolution and improvement in a privacy-preserving manner, accommodating the dynamic nature of the data present in federated learning scenarios. In addition to addressing the challenges associated with unlearning outdated, private, and irrelevant data in both single-modality and multimodality scenarios, FRAMU offers valuable insights into the convergence behaviour and optimization of the federated learning process. The major contributions of our work are as follows:

- We propose an adaptive unlearning algorithm using an attention mechanism to adapt to changing data distributions and participant characteristics in single-modality and multimodality scenarios.
- We develop a novel design to personalize the unlearning process using the FedAvg mechanism [21] and unlearn the outdated, private, and irrelevant data.
- We propose an efficient unlearning algorithm that demonstrates fast convergence and achieves optimal solutions within a small number of communication rounds.
- We conduct extensive experiments to demonstrate the efficiency and effectiveness of the proposed approach using real-world datasets.

Organization. In Section II, we review related works. Section III outlines the problem addressed in this study. We present the proposed framework FRAMU in Section IV. The applications of FRAMU in single-modality and multimodality are discussed in Section V. In Section VI, we present the experimental setup and the evaluation results of the proposed framework, along with convergence and optimization analysis. Section VII delves into the implications of the proposed framework. Finally, in Section VIII, we conclude the paper.

II. RELATED WORKS

The importance of data privacy in distributed learning systems has garnered significant attention, especially when handling sensitive types of data like medical or behavioral information [22]. Differential Privacy (DP), a mathematically rigorous framework for ensuring individual privacy, has been widely adopted for this purpose [23], [24]. Efforts to integrate DP within distributed learning environments, particularly in federated learning, have been increasing [13], [14]. Abadi et al. [16] developed a seminal approach called Deep Learning with Differential Privacy (DP-SGD), which adapts the Stochastic Gradient Descent (SGD) algorithm to meet DP standards by clipping gradients and injecting noise, thereby offering stringent privacy safeguards during deep neural network (DNN) training. Building on this, McMahan et al. [17] further

tailored DP mechanisms for federated learning through an extension called DP-FedAvg. While these methods effectively address privacy concerns, they often fall short in dealing with dynamic data distributions, a prevalent issue in distributed learning [18]. Specifically, data sets can evolve over time, rendering some information outdated or irrelevant, and the persistence of such data in the learning process can compromise model efficacy. Although Machine Unlearning approaches like Sharded, Isolated, Sliced, and Aggregated (SISA) training [19] have emerged to tackle this issue by enabling efficient selective forgetting of data, these methods are not yet designed to work synergistically with privacy-preserving techniques like DP [20].

Federated learning has substantially revolutionized distributed learning, enabling the training of Machine Learning models on decentralized networks while preserving data privacy and minimizing communication costs [25]. Among the pioneering works in this area is the FedAvg algorithm by McMahan et al. [21], which relies on model parameter averaging across local models and a central server. However, FedAvg is not without its limitations, particularly when handling non-IID data distributions [26]. Solutions like FedProx by Li et al. [27] have sought to address this by introducing a proximal term for improved model convergence. While other researchers like Sahu et al. [28] and Konečný et al. [29] have made strides in adaptive learning rates and communication efficiency, the realm of federated learning still faces significant challenges in dynamic adaptability and efficient Machine Unlearning. While privacy has been partially addressed through Differential Privacy [30] and Secure Multiparty Computation [31], these techniques often compromise on model efficiency. Additionally, the applicability of federated learning in diverse sectors like healthcare and IoT emphasizes the unmet need for a model capable of dynamically adapting to varied data distributions, while preserving privacy and efficiency [32], [33].

Reinforcement Learning has garnered much attention for its ability to train agents to make optimal decisions through trial-and-error interactions with their environments [34], [35]. Several pivotal advancements have shaped the field, including the development of Deep Q-Networks (DQNs) [36]. DQNs combine traditional reinforcement learning techniques with DNNs, significantly enhancing the system's ability to process high-dimensional inputs such as images. Furthermore, experience replay mechanisms have been integrated into them to improve learning stability by storing and reusing past experiences [37]. Mnih et al. [38] significantly accelerated the reinforcement learning field by implementing DQNs that achieved human-level performance on a variety of complex tasks. However, there are evident gaps in addressing challenges posed by non-stationary or dynamic environment situations where the statistical properties of the environment change over time. Under such conditions, a reinforcement learning agent's ability to adapt quickly is paramount. Several approaches like meta-learning [39] and attention mechanisms [40], [41] have sought to remedy these issues to some extent. Meta-learning, for example, helps models quickly adapt to new tasks by training them on a diverse range of tasks. However, the technique does not offer a robust solution for unlearning or forgetting outdated

or irrelevant information, which is crucial for maintaining performance in dynamic environments. In a similar vein, attention mechanisms help agents focus on important regions of the input space, but they also fail to address the need for efficient unlearning of obsolete or irrelevant data. This leaves us with a significant research gap: the lack of mechanisms for efficient unlearning and adaptability in reinforcement learning agents designed for non-stationary, dynamic environments.

A key challenge for federated learning when faced with dynamic data distributions and the accumulation of outdated or irrelevant information is its adaptability in evolving environments. Reinforcement learning has been instrumental in training agents for optimal decision-making in dynamic environments, yet it too grapples with the need to efficiently unlearn outdated or irrelevant data. These challenges underscore the importance of integrating attention mechanisms into the Machine Unlearning process. Unlike selective data deletion, attention mechanisms assign reduced weights to outdated, private, or irrelevant information. The dynamic adjustment of attention scores allows these models to prioritize relevant data while disregarding obsolete or extraneous elements. By bridging the worlds of federated learning and reinforcement learning with attention mechanisms, our study addresses the pressing need for an integrated solution that optimizes decision-making in distributed networks with changing data landscapes [42]. In addition, this approach must preserve data privacy and adaptively forget outdated, private, or irrelevant information.

III. PRELIMINARIES & PROBLEM DEFINITION

This section establishes the foundational concepts and mathematical notations essential for the discussions and analyses presented in this paper. These concepts are summarized in Tab. I and form the basis for understanding the subsequent problem definitions and solution approaches. Our research is centered around the exploration of unlearning mechanisms in Machine Learning models, focusing on maintaining accuracy and computational efficiency while addressing the challenges posed by outdated or irrelevant data.

The problem is defined by two distinct settings: single-modality and multimodality. The single-modality setting is simpler and widely applicable in scenarios with uniform data types, such as sensor networks or content recommendation systems. However, it may lack the context provided by different types of data, potentially leading to less nuanced decisions. On the other hand, the multimodality setting is more complex but highly relevant in fields like healthcare, where a range of data types (e.g., medical imaging, patient history, etc.) can be used for more comprehensive understanding and decision-making. By exploring the problem in both these settings, we offer solutions that are both versatile and contextually rich.

A. Problem Definition - Single Modality

Problem Definition 1. Let $AG = \{ag_1, ag_2, \dots, ag_n\}$ be a set of agents, where each agent $ag \in AG$ represents an entity like an IoT device, traffic point, wearable device, edge computing node, or content recommendation system. Each agent ag observes states $S_i = \{s_1, s_2, \dots, s_m\}$ and takes actions $A = \{a_1, a_2, \dots, a_n\}$ based on a policy $\pi_i(s, a)$. Rewards $R_i(s, a)$ evaluate the quality of actions taken in

TABLE I: Summary of Notations and Descriptions

Symbol	Description
AG	Set of agents in the model
ag	An individual agent in the set AG
S_i	States observed by an agent ag
A	Set of possible actions
$\pi_i(s, a)$	Policy followed by the agent
$R_i(s, a)$	Rewards for actions in different states
θ_{ag}	Parameters of local models for agent ag
θ_g	Parameters of global model
$w_{i, ag}$	Attention score for a data point i in agent ag
M	Set of modalities in multimodality setting
X_m	Data vectors for modality m
θ_m	Parameters for modality m
$w_{i, m}$	Attention scores within a modality m
t	Time step
s_t	State at time step t
a_t	Action at time step t
r_t	Reward at time step t
R_t	Cumulative reward
$\pi(a_t s_t)$	Policy function
$Q(s_t, a_t)$	Q -function
γ	Discount factor
α_i	Attention score for feature i
$\Delta\theta_{ag}$	Update sent by agent ag
f	Function for calculating attention scores
$w_{g, ag}$	Global attention score for update from agent ag
AG	Number of local agents
α_{avg}	Average attention score
δ	Predetermined threshold for attention score
$ag \in AG$	A specific agent within the set of all agents AG
m	Number of modalities
x_1, x_2, \dots, x_m	Data vectors for each modality
v_i	Feature vector for modality i
\bar{w}_j	Averaged attention score across modalities for data point j
λ	Mixing factor
T	The total number of training rounds
α	Learning rate for Q -value function updates
η	Scaling factor for attention score updates
β	Mixing factor for combining global and local model parameters
ε	Convergence threshold for global model parameters
w_{ag}	Local model parameters for agent ag
W	Global model parameters
A_i	Attention score for data point i
$A_{i, ag}$	Attention score for data point i within agent ag
N	Total number of data points across all agents
n_{ag}	Number of data points in agent ag

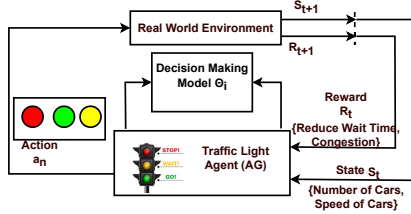


Fig. 2: Single Modality Example

different states. Agents possess local models with parameters θ_i , while a central server maintains a global model with parameters θ_g .

Example 2. In the single-modality setting shown in Fig. 2, let $AG = \{ag_1, ag_2, \dots, ag_n\}$ be a set of agents. An agent ag can represent a real-world entity such as a traffic light in a city. These traffic lights observe various states $S_i = \{s_1, s_2, \dots, s_m\}$, such as the number and speed of passing cars, and the change of colors (actions $A = \{a_1, a_2, \dots, a_n\}$) according to an algorithmic policy $\pi_i(s, a)$. The system evaluates the effectiveness of the traffic light changes in reducing wait time or congestion (rewards $R_i(s, a)$). Each traffic light has its own local decision-making model characterized by parameters θ_i , and there is a global model for optimizing city-wide traffic flow with parameters θ_g .

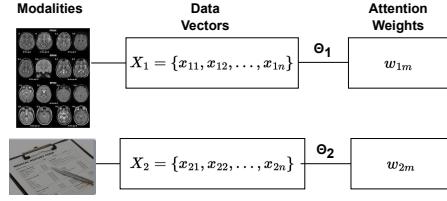


Fig. 3: Multimodality Example

To address the challenge of preserving data privacy and adaptively forgetting private, outdated, or irrelevant information, attention scores w_{ij} are assigned to each data point j in the local dataset of agent $ag \in AG$. These attention scores, computed using a function f that considers the current model state or contextual information, guide the learning and unlearning process within each agent. By assigning higher attention scores to relevant data and potentially forgetting or down-weighting irrelevant data, the agents can effectively focus on the most informative and up-to-date information.

B. Problem Definition - Multimodality

Problem Definition 2. In the multimodality setting, let $M = \{1, 2, \dots, m\}$ represent the set of modalities, where m is the total number of modalities. Each modality $m \in M$ is associated with a set of data vectors $X_m = \{x_{m1}, x_{m2}, \dots, x_{mn}\}$, and has its local model with parameters θ_k . Attention scores w_{im} are assigned to individual data points x_{im} within each modality to guide the learning and unlearning process.

Example 3. In the multimodality setting shown in Fig. 3, consider a healthcare system as a collection of agents in set $M = \{1, 2, \dots, m\}$, where m represents different types of medical data (modalities) such as medical imaging and patient history. For instance, medical imaging (modality M_1) would have a set of MRI scans represented as data vectors $X_1 = \{x_{11}, x_{12}, \dots, x_{1n}\}$. Likewise, patient history (modality M_2) might involve a set of past diagnosis records that are represented as data vectors $X_2 = \{x_{21}, x_{22}, \dots, x_{2n}\}$. Each modality has a specialized model with parameters θ_1 for medical imaging and θ_2 for patient history. These models use attention mechanisms to weigh the importance of each data point, represented by attention scores w_{1m} for MRI scans and w_{2m} for patient history records. These scores guide the decision-making process in diagnosis and treatment.

In the multimodality setting, the complexity is elevated by the integration of heterogeneous data types and the application of specialized machine learning models for each modality. For example, in a healthcare system, combining data from disparate sources like medical imaging and patient history presents a unique challenge. Each data type, or modality, not only varies in format but also in the nature of the information it conveys, necessitating distinct processing and analysis methods. The key challenge here is to develop an integrated approach that effectively synthesizes these diverse data streams into a coherent understanding, enhancing decision-making in critical applications such as patient diagnosis and treatment. Attention mechanisms play a crucial role in this context, determining the relevance of each data point across different

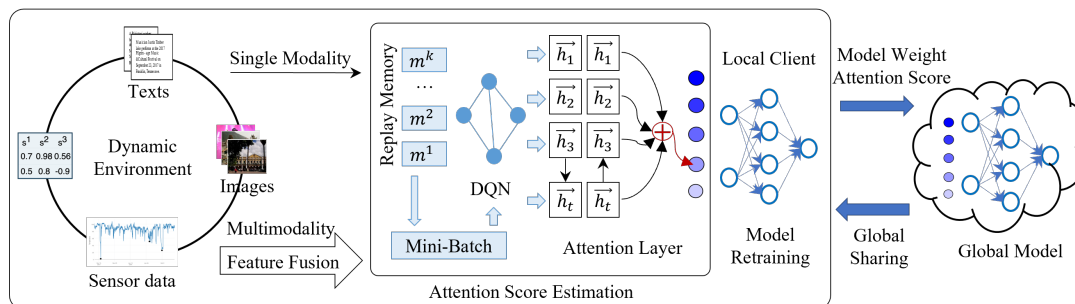


Fig. 4: An overview of the proposed FRAMU framework, illustrating its end-to-end adaptive algorithm that incorporates an attention mechanism. The figure is divided into multiple components, each corresponding to a specific phase in the federated learning process. Starting from the left, the diagram begins with data collection from diverse modalities. The framework applies an adaptive learning algorithm that not only updates the global model, but also incorporates an efficient unlearning mechanism for discarding outdated, private, or irrelevant data.

modalities. However, assigning and calibrating these attention scores is non-trivial and introduces an additional layer of complexity. The successful implementation of multimodal systems has profound implications, particularly in improving the accuracy and efficacy of decision-making processes.

IV. FRAMU FRAMEWORK

In an era marked by an ever-increasing influx of data, the need for adaptive Machine Learning models that can efficiently unlearn outdated, private, or irrelevant information is paramount. The methodology proposed in this paper addresses this necessity by introducing two key technical contributions. First, we propose an adaptive unlearning algorithm that utilizes attention mechanisms to tailor the learning and unlearning processes in a single-modality, and then extend the process to multimodality. This innovative approach allows the model to adapt to dynamic changes in data distributions, as well as variations in participant characteristics such as demographic information, behavioural patterns, and data contribution frequencies among others. Second, we put forth a novel design that employs the FedAvg mechanism [21] to personalize the unlearning process. This design ensures that the model is able to discard data that has become irrelevant, outdated, or potentially invasive from a privacy perspective, thus preserving the integrity of the learning model while adapting to new or changing data. The following sections will elaborate on these contributions, providing a detailed discussion of the proposed framework as depicted in Fig. 4.

The FRAMU framework adopts a federated learning architecture comprising Local Agents and a Central Server, each with distinct roles in model training, unlearning, and adaptation. It employs a reinforcement learning paradigm where each agent iteratively learns from its environment. This integration of federated learning and reinforcement learning is termed federated reinforcement learning. However, what sets FRAMU apart is the integration of attention mechanisms to weigh the relevance of each data point in learning and unlearning. The attention scores are then aggregated and processed at the Central Server to refine the global model.

- **Local Agents:** Responsible for collecting real-time data and performing local model updates. They observe states, take actions, and calculate rewards to update their Q-values and attention scores.
- **Central Server:** Aggregates local models and attention scores, filters out irrelevant data points, and updates the global model.
- **Attention Mechanism:** Dynamically calculates attention scores for each data point to inform the unlearning process.
- **FedAvg Mechanism:** Utilized for global model updates, ensuring that the global model represents a consensus across all agents.

The FRAMU framework, as outlined in Algorithm 1, has been carefully designed to facilitate adaptive decision-making in distributed networks through federated reinforcement learning. Each step within the algorithm is crafted with specific intentions: The initialization stage (Lines 1-3) sets the groundwork by initializing local and global model parameters, as well as attention scores. These initializations are crucial for ensuring that both local and global perspectives are considered right from the start of the learning process. The iterative learning process (Lines 4-24) involves several key components. Local Agent Decision-Making (Lines 5-11) enables each local agent to observe states, take actions, and update its Q-values and attention scores, ensuring that local knowledge is continuously updated to reflect the dynamic nature of the agents' environments. Central Server Aggregation (Lines 12-17) plays a pivotal role in integrating local updates and refining the global model. By assessing the attention scores, the server can identify and diminish the influence of less relevant data points, thereby enhancing the model's focus on significant information. Model Synchronization (Lines 18-24) involves the dissemination of global model parameters back to local agents for fine-tuning, ensuring a bi-directional flow of information that keeps local models informed by their immediate environment and aligned with the broader objectives of the global model.

Algorithm 1: FRAMU Framework

Input: a set of Local Agents, a Central Server, $T, \theta, \alpha, \eta, \gamma, \beta, \varepsilon$

Output: \bar{W} : Trained global model parameters for federated reinforcement learning

- 1 Initialize local model parameters w_{ag} for each agent ag ;
- 2 Initialize global model parameters W at the central server;
- 3 Initialize attention scores $A_{i,ag,m}$ for each data point i in agent ag and modality m ;
- 4 **while** $t \leq T$ **do**
- 5 **foreach** local agent ag **do**
- 6 Observe current states $s_{i,m}$ for each modality m ;
- 7 Take action a_t based on policy derived from $Q(s, a; w_{ag})$;
- 8 Observe reward r_t and next states $s'_{i,m}$ for each modality m ;
- 9 Compute TD error $\delta = r_t + \gamma \max_a Q(s'_{i,m}, a; w_{ag}) - Q(s_{i,m}, a_t; w_{ag})$;
- 10 Update $Q(s_{i,m}, a_t; w_{ag}) \leftarrow Q(s_{i,m}, a_t; w_{ag}) + \alpha \delta$;
- 11 Update attention scores $A_{i,ag,m} \leftarrow A_{i,ag,m} + \eta |\delta|$;
- 12 Send local model parameters w_{ag} and attention scores $A_{i,ag,m}$ to Central Server;
- 13 **foreach** data point i in modality m **do**
- 14 **if** $\sum_{ag} \frac{1}{m} \sum_m A_{i,ag,m} / N_{ag} < \theta$ **then**
- 15 Reduce influence of data point i in the global model;
- 16 Aggregate local model parameters to update global parameters: $W \leftarrow \sum_{ag} \left(\frac{N_{ag}}{N} \right) w_{ag}$;
- 17 Send updated global model parameters W to local agents;
- 18 **foreach** local agent ag **do**
- 19 Fine-tune local model with global model:
- 20 $w'_{ag} \leftarrow \beta W + (1 - \beta) w_{ag}$;
- 21 **if** $|P(W_{t+1}) - P(W_t)| < \varepsilon$ **then**
- 22 Break;
- 23 Increment t ;
- 24 **return** W

V. APPLICATIONS OF FRAMU

This section explores the practical applications of the FRAMU framework across different settings, single-modality and multimodality, and its continuous adaptation and learning.

A. FRAMU with Single Modality

Central to FRAMU is an attention layer that functions as a specialized approximator, augmenting the learning capability of individual agents. This attention layer distinguishes itself by assigning attention scores to individual data points during the function approximation process. These scores serve as indicators of each data point's relevance in the agent's local learning. The agent updates these scores as it interacts with its environment and receives either rewards or penalties, thereby continually refining its model. Specifically, an agent operates in discrete time steps, current state s_t , taking action a_t , and receiving reward r_t , at each time step t . The ultimate goal is to determine an optimal policy $\pi(a_t|s_t)$ that maximizes the accumulated reward R_t . The Q -function, which quantifies expected accumulated rewards with a discount factor γ , is given by Equation 1.

$$Q(s_t, a_t) = \mathbb{E}[R_t | s_t, a_t] = r_t + \gamma \mathbb{E}[Q(s_{t+1}, a_{t+1}) | s_t, a_t] \quad (1)$$

The attention layer further characterizes each state s_t by its features $[x_1, x_2, \dots, x_n]$, and assigns attention scores α_i as per:

$$\alpha_i = \text{Attention}(x_i, \text{context}) \quad (2)$$

Here, the context may include additional data such as previous states or actions. The Q -function is then approximated using a weighted sum of these features:

$$Q(s_t, a_t) \approx \sum (\alpha_i \cdot x_i) \quad (3)$$

After completing their respective learning cycles, agents forward their model updates θ and attention scores α to the Central Server as a tuple (θ, α) .

1) Local and Global Attention Score Estimation

FRAMU estimates attention scores both locally and globally. On the local front, each agent employs its attention mechanism to compute scores for individual data points based on their relevance to the task at hand. For an agent ag with local model parameters θ_{ag} , the attention score w_{ij} for data point j is given by:

$$w_{ij} = f(s_j, \theta_{ag}) \quad (4)$$

At the global level, these scores assist the Central Server in prioritizing updates or pinpointing data points for global unlearning. For global parameters θ_g , the global attention score derived from the updates of agent ag is:

$$w_{g,ag} = f(\Delta\theta_{ag}, \theta_g) \quad (5)$$

In this equation, $\Delta\theta_{ag}$ is the model update from agent ag , and the function f calculates attention scores while taking into account the aggregated local scores and other global contextual cues.

2) Global Model Refinement and Unlearning

Model updates from local agents are aggregated at the Central Server using FedAvg [43]. The attention scores are instrumental in the global unlearning process, with the average attention score calculated as:

$$\alpha_{\text{avg}} = \frac{1}{AG} \sum \alpha_{ag} \quad (6)$$

When α_{avg} falls below a predetermined threshold δ , the server adjusts the contribution of the respective feature in the global model as given by Equation 7:

$$\theta_{\text{global}}' = g(\theta_{\text{global}}, \alpha_{\text{avg}}) \quad (7)$$

Once refined, this global model is sent back to the local agents. The enhanced model shows improved adaptability and robustness to changes in data distributions due to the integration of aggregation and unlearning mechanisms. Consequently, the local agents are better positioned to excel within their particular operational environments. These revised global model parameters, denoted as θ_{global}' , are then dispatched from the Central Server to the local agents, where $\theta_k = \theta_{\text{global}}'$.

B. FRAMU with Multimodality

The multimodal FRAMU Framework extends its capabilities to seamlessly incorporate various data types, including images, text, audio, and sensor readings. This integration not only enriches decision-making but also optimizes the performance of local agents. By fine-tuning their models to multiple data types, agents are better equipped to operate in complex environments.

1) Modality-Specific Attention Mechanisms

To effectively manage data from diverse sources, the framework employs specialized attention mechanisms for each modality. These mechanisms generate unique attention scores for data points within a given modality, aiding in both learning and unlearning processes. By doing so, the framework allows local agents to focus on the most relevant and informative aspects of each modality.

The attention scores for a specific modality j for an agent $ag \in AG$ can be mathematically represented as:

$$w_{ij} = f_j(s_{ij}, \theta_i), \quad (8)$$

Here, s_{ij} signifies a data point from modality j related to agent $ag \in AG$, while θ_i represents that agent's local model parameters. The function f_j considers modality-specific attributes and context to compute these attention scores.

For a feature vector v_i derived from modality j within agent $ag \in AG$, feature-level fusion can be represented as:

$$v_i = [x_{i1}, x_{i2}, \dots, x_{im}] \quad (9)$$

2) Unlearning and Adaptation across Modalities

In a multimodal setup, attention scores from all modalities collectively inform the unlearning process. If a data point consistently receives low attention scores across different modalities, it indicates that the point is either irrelevant or outdated. The Central Server uses this multimodal insight to refine the global model.

The average attention score across all modalities for a specific data point is:

$$\bar{w}_j = \frac{1}{m} \sum_{i=1}^m w_{ij} \quad (10)$$

If \bar{w}_j falls below a predefined threshold, the Central Server de-emphasizes or removes that data point from the global model, ensuring that only current and relevant data contribute to decision-making.

During the adaptation phase, local agents utilize the updated global model to enhance their local models. The interplay between global and local parameters is regulated by a mixing factor, which allows local agents to leverage shared insights while preserving modality-specific skills. This relationship can be denoted by:

$$\theta_i^{\text{new}} = \lambda \theta_{\text{global}} + (1 - \lambda) \theta_i^{\text{old}} \quad (11)$$

Here, θ_i^{new} represents the updated local model parameters, θ_{global} signifies the global model parameters, θ_i^{old} is the previous local parameters, and λ serves as the mixing factor.

Through this, the multimodal FRAMU framework maintains an up-to-date and relevant global model, while enabling local agents to make better decisions across a range of data types.

C. Continuous Adaptation and Learning in the FRAMU Framework

Continuous adaptation and learning are critical in the FRAMU framework, enabling it to thrive in dynamic and changing environments. These processes create an iterative exchange of knowledge between local agents and a Central Server, which leads to consistent model refinement on both local and global scales.

1) Local-Level Adaptation

Local agents need the ability to adapt in real time to changes in their operational environments. Within reinforcement learning paradigms, agents continually update their policies in response to actions taken and rewards observed. Furthermore, attention scores allocated to data points or features can vary dynamically based on new data or shifts in relevance. This adaptability ensures that the models of individual local agents remain current. Let s_t denote the state of the environment at time t , and a_t represent the action taken by the agent. After receiving a reward r_t and transitioning to a new state s_{t+1} , the agent aims to maximize the expected cumulative reward. The Q-value function $Q(s, a)$ serves as a proxy for this cumulative reward, and it is updated using temporal-difference learning algorithms as follows:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right] \quad (12)$$

Here, α is the learning rate, and γ is the discount factor.

Attention scores, denoted by A_i for data point i , are updated based on the temporal-difference error δ :

$$A_i \leftarrow A_i + \eta |\delta|, \quad (13)$$

where η is a scaling factor, and $\delta = r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)$.

2) Global Model Aggregation and Adaptation

As local agents continuously update their models, these adaptations are communicated to the Central Server. It aggregates this information to refine the global model while also tracking the attention scores from local agents. If these scores reveal diminishing importance for certain data points, the server may initiate global unlearning. This ensures the global model remains current and avoids obsolescence. Local agents send their updated model parameters, w_{ag} for agent ag , and attention scores $A_{i,ag}$ to the Central Server. The server aggregates these to update the global model parameters W as follows:

$$W \leftarrow \frac{1}{AG} \sum_{ag} w_{ag}, \quad (14)$$

where AG represents the total number of local agents.

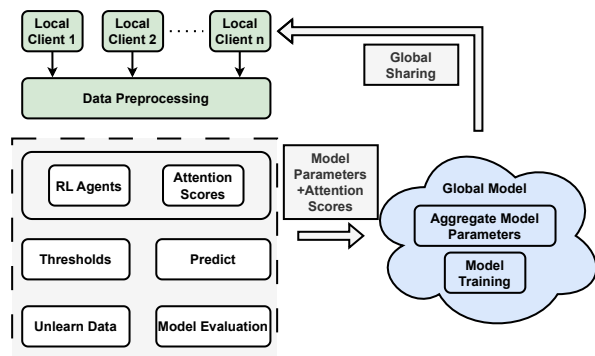


Fig. 5: Experimental Setup: This diagram showcases the architecture of the FRAMU framework, detailing the interaction between local and global models within a federated learning environment.

3) Feedback Mechanisms

After the global model is updated, it is disseminated back to local agents through a feedback loop. This cyclic interaction allows local agents to either initialize or further refine their models based on the global one. This is particularly beneficial when local agents confront new or unfamiliar data points that other agents have encountered. Through this mechanism, the global model acts as a repository of shared knowledge, enhancing the decision-making capabilities of all local agents. The global model parameters W are sent to local agents, who then adjust their local models using a mixing factor β as follows:

$$w'_{ag} \leftarrow \beta W + (1 - \beta)w_{ag}, \quad (15)$$

where β ranges from 0 to 1 and regulates the influence of the global model on local models.

VI. EXPERIMENTAL SETUP AND RESULTS ANALYSIS

To effectively evaluate the performance of the FRAMU framework, we undertook comprehensive experiments using real-world datasets. These experiments were designed to validate not only the efficiency and effectiveness of our approach but also to establish the practical utility of FRAMU in real-world applications. Our experimental setup encompassed several components, including datasets, baseline models, evaluation metrics, and specific FRAMU configurations, as depicted in Fig. 5. A critical aspect of our experimentation involved fine-tuning key thresholds to guide the unlearning process, particularly the `outdated_threshold` and `irrelevant_threshold`. These parameters were adjusted based on domain expertise and sensitivity analysis, with the `outdated_threshold` defining the time frame for data obsolescence and the `irrelevant_threshold` setting criteria for data's statistical insignificance. Additionally, we introduced a `privacy_epsilon` parameter to balance data utility with privacy preservation, aligning with GDPR regulations.

Deep learning methods are known for their ability to learn features autonomously and automate model-building processes. Despite criticisms of neural network family algorithms for their 'black box' nature, deep learning models are

renowned for their robust and efficient performance. These models are widely adopted by the research community [44], [45]. In our work, we utilized a Convolutional Neural Network (CNN) for image and sensor data, and a Long Short-Term Memory (LSTM) network for time series and text data, specifically tailored for federated learning scenarios. This model choice was made to efficiently handle both single and multimodal data types, integrating attention mechanisms and unlearning processes to enhance overall functionality.

We found the tuning of parameters such as the outdated threshold, irrelevant threshold, and the β value for local model fine-tuning to be crucial. It was essential to strike the right balance in the outdated threshold to prevent premature data discarding or retention of outdated information, which could affect model accuracy and relevancy. Similarly, careful calibration of the irrelevant threshold was necessary to maintain a balance between data comprehensiveness and quality, ensuring useful data was not excluded nor excessive noise retained. The β value, crucial in determining the extent of global model influence on local models, required fine-tuning to ensure an optimal balance between local and global learning. This was key for local models to benefit from global insights while preserving their unique learning characteristics. The interplay of these hyperparameters significantly influenced FRAMU's performance, particularly in its ability to adapt to new data and retain relevant historical information. Through sensitivity analyses, we determined their optimal ranges, aiming to maximize FRAMU's efficiency and adaptability in various real-world scenarios.

TABLE II: Datasets for evaluation

Modality	Dataset	OD*	PD*	ID*	Description
Single Modality	AMPds2 [46]	✓	✓	✓	Electricity, water, and natural gas consumption data from a Canadian household.
	METR-LA [47]	✓	✗	✓	Traffic speed data from over 200 sensors in Los Angeles Metropolitan area.
	MIMIC-III [48]	✓	✓	✓	Health-related data from critical care units, including demographics, vital signs, laboratory results, and medications.
Multi Modality	NYPD [49]	✓	✓	✓	Records of complaints filed with the New York City Police Department.
	MIMIC-CXR [50]	✓	✓	✓	Chest radiographs with associated radiology reports for medical image analysis tasks.
	Smart Home EnergyDataset (SHED) [51]	✓	✓	✓	Energy consumption data from smart home devices and appliances.

*OD - Outdated Data, PD - Privacy Data, ID - Irrelevant Data.

A. Datasets

In this study, publicly available datasets that encompass various modalities and address specific challenges related to outdated, private, and irrelevant data are adopted. Tab. II provides detailed information about each dataset, including the data modality, number of instances, attributes, target variables, and specific characteristics pertinent to our study. In order to evaluate FRAMU, we conducted a comprehensive comparison of its performance against several contemporary baseline models.

B. Baseline Models

In the evaluation of the FRAMU framework’s performance and robustness, we have carefully selected several baseline models for comparison. The models in baseline models were adopted from the original work. The rationale behind choosing each model and its relevance to our study is elaborated below:

• Single-modality

- **FedLU [52]**: FedLU represents a significant advance in federated learning, integrating knowledge graph embedding with mutual knowledge distillation. Its selection as a baseline is due to its innovative approach to collaborative learning, which is closely aligned with FRAMU’s objectives in single-modality settings. FedLU’s methodology provides a comparative framework for assessing FRAMU’s efficiency in knowledge synthesis and distribution.
- **Zero-shot MU [53]**: Zero-shot MU specializes in Machine Unlearning, employing error-minimizing-maximizing noise and gated knowledge transfer. This model was chosen for its novel approach to unlearning, providing a benchmark to evaluate FRAMU’s capability in effectively removing learned information without extensive retraining, a crucial aspect in dynamic environments.
- **SISA Training [19]**: The SISA Training framework is a strategic model that limits data points for optimized unlearning. Its inclusion as a baseline allows us to compare FRAMU’s efficiency in data management and unlearning processes, especially in scenarios where data minimization is key to performance and privacy.

• Multimodality

- **MMoE [54]**: The MMoE model, optimized for handling multimodal data via ensemble learning, serves as a benchmark for evaluating FRAMU’s performance in multimodality settings. Its approach, employing expert networks for different data modalities, provides a comparative perspective for FRAMU’s adaptability and efficiency in handling diverse data types.
- **CleanCLIP [55]**: CleanCLIP, a fine-tuning framework that mitigates spurious associations from backdoor attacks, is pivotal for comparing FRAMU’s robustness against data security threats. Its focus on weakening spurious correlations offers insights into FRAMU’s capabilities in maintaining data integrity and security.
- **Privacy-Enhanced Emotion Recognition (PEER) [56]**: The PEER model, utilizing adversarial learning for privacy-preserving emotion recognition, aligns well with FRAMU’s privacy objectives. Its comparison with FRAMU highlights the effectiveness of FRAMU in safeguarding privacy while performing complex analytical tasks.

C. Evaluation Metrics

The FRAMU framework is evaluated using several important metrics: Mean Squared Error (MSE) [57], Mean Absolute Error (MAE) [58], Reconstruction Error (RE) [59], and Activation Distance (AD) [60]. A lower MSE or MAE score

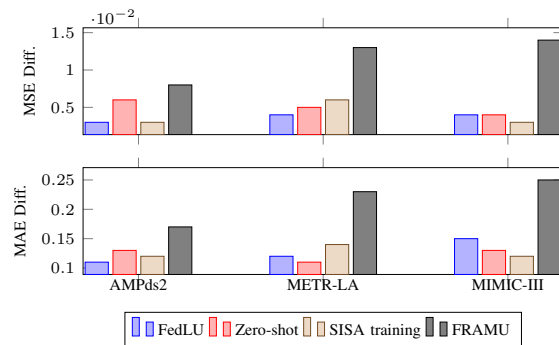


Fig. 6: Comparative Analysis of MSE and MAE Differences between Original and Unlearned Single Modality Data

shows that the unlearning process is closely aligned with what was expected, indicating a high quality of unlearning. The RE measures how well the model can rebuild data that it has unlearned, with a lower score being better. AD measures the average distance between the predictions of the model before and after unlearning, using what is known as L2-distance, on a specific set of forgotten data. These metrics together give a well-rounded evaluation of how well the unlearning process is working.

All the experiments were run using Python programming language (version 3.7.6) and related TensorFlow, Keras, Open Gym AI, and stable_baselines3 packages.

D. FRAMU Unlearning Results in Single Modality Context

To assess the effectiveness of FRAMU in unlearning outdated, private, and irrelevant data, we analyzed the results from various experiments. FRAMU’s performance was benchmarked against that of established baseline models: FedLU, Zero-shot MU, and SISA Training. It’s important to note that the METR-LA dataset [47] was excluded from the private data unlearning evaluation due to its lack of privacy-sensitive data. For a thorough comparison, we present the performance metrics of FRAMU in unlearning outdated, private, and irrelevant data alongside the results from baseline models in Tab. III. The p-values in these comparisons are indicative of the statistical significance of FRAMU’s performance improvements.

1) Outdated Data

The unlearning of outdated data is vital for maintaining model accuracy and relevance. Outdated data might introduce noise, biases, or outdated patterns. By selectively unlearning such data, FRAMU aims to align the model with the latest data distribution. FRAMU consistently achieved lower MSE and MAE than the baseline models in unlearning outdated data across various datasets. This improvement, evident from the low p-values in Tab. III, demonstrates FRAMU’s statistically significant superiority in adapting models to current data distributions.

2) Private Data

The retention of private data in models can pose significant privacy and legal risks. To mitigate this, FRAMU incorporates techniques for unlearning private data while preserving privacy. Excluding the METR-LA dataset from this analysis, FRAMU consistently outperformed the baseline models in

TABLE III: FRAMU - Evaluation Results in Single Modality Context

Unlearning	Dataset		FedLU [52]			Zero-shot [53]			SISA [19]			FRAMU (Ours)	
			MSE	MAE	p-value	MSE	MAE	p-value	MSE	MAE	p-value	MSE	MAE
Outdated Data	Original	AMPds2	0.063	6.740	0.024	0.061	6.890	0.031	0.059	6.760	0.041	0.046	5.570
		METR-LA	0.079	7.140	0.016	0.082	7.210	0.038	0.078	7.090	0.029	0.065	5.930
		MIMIC-III	0.099	12.800	0.031	0.102	12.930	0.045	0.097	12.680	0.032	0.083	10.650
	Unlearned	AMPds2	0.060	6.630	0.015	0.055	6.860	0.029	0.056	6.690	0.036	0.038	4.670
		METR-LA	0.075	7.020	0.029	0.077	7.100	0.025	0.072	6.960	0.032	0.052	4.910
		MIMIC-III	0.095	12.650	0.023	0.098	12.820	0.041	0.094	12.580	0.017	0.069	8.900
Private Data	Original	AMPds2	0.052	6.780	0.014	0.054	6.930	0.037	0.053	6.810	0.041	0.041	5.540
		MIMIC-III	0.078	12.870	0.035	0.080	13.010	0.043	0.079	12.760	0.045	0.064	10.600
		AMPds2	0.049	6.670	0.011	0.052	6.910	0.035	0.051	6.740	0.015	0.033	4.590
	Unlearned	MIMIC-III	0.075	12.720	0.031	0.077	12.900	0.038	0.076	12.650	0.016	0.053	8.860
		AMPds2	0.047	6.700	0.035	0.050	6.850	0.044	0.048	6.730	0.031	0.037	5.440
		METR-LA	0.054	7.100	0.027	0.056	7.170	0.041	0.055	7.050	0.025	0.043	5.830
Irrelevant Data	Original	MIMIC-III	0.070	12.730	0.038	0.072	12.870	0.031	0.071	12.620	0.039	0.057	10.410
		AMPds2	0.045	6.590	0.011	0.047	6.830	0.036	0.046	6.660	0.029	0.030	4.510
		METR-LA	0.052	6.980	0.014	0.054	7.070	0.019	0.053	6.930	0.022	0.035	4.750
	Unlearned	MIMIC-III	0.068	12.580	0.029	0.070	12.760	0.024	0.069	12.510	0.027	0.047	8.690

TABLE IV: Comparative analysis of FRAMU's performance in single modality against baseline models in RE and AD metrics.

Unlearning	Dataset	FedLU [52]		Zero-shot MU [53]		SISA training [19]		FRAMU (Ours)	
		RE	AD	RE	AD	RE	AD	RE	AD
Outdated Data	AMPds2	0.03	0.66	0.029	0.68	0.028	0.67	0.024	0.57
	METR-LA	0.038	0.7	0.039	0.71	0.037	0.69	0.033	0.59
	MIMIC-III	0.048	1.26	0.049	1.28	0.047	1.25	0.043	1.15
Private Data	AMPds2	0.031	0.67	0.032	0.69	0.03	0.67	0.026	0.57
	MIMIC-III	0.049	1.27	0.051	1.29	0.048	1.27	0.044	1.17
Irrelevant Data	AMPds2	0.028	0.66	0.029	0.68	0.027	0.66	0.023	0.56
	METR-LA	0.034	0.7	0.035	0.71	0.033	0.69	0.029	0.59
	MIMIC-III	0.05	1.26	0.052	1.28	0.049	1.25	0.045	1.15

both MSE and MAE metrics in scenarios involving private data. For example, in the AMPds2 dataset, FRAMU's superior performance in MSE (0.038) and MAE (4.670) is a testament to its effective federated reinforcement learning approach that respects privacy concerns. The significance of these performance gains is reinforced by the associated p-values.

3) Irrelevant Data

Unlearning irrelevant data helps reduce noise and interference from non-contributory data points, enhancing model accuracy and prediction. FRAMU showed exceptional performance in unlearning irrelevant data, recording the lowest MSE and MAE values across all datasets in comparison to the baseline models. For instance, in the AMPds2 dataset, FRAMU's MSE of 0.033 and MAE of 5.600 surpassed other models. The low p-values validate FRAMU's significant advantage in discarding irrelevant data.

Fig. 6 visually compares the differences in MSE and MAE between original and unlearned data across various datasets and models. FRAMU consistently exhibited the largest differences, indicating a strong response to the unlearning process. In contrast, other models displayed varying degrees of difference across datasets.

Moreover, in the comparison of RE and AD metrics as illustrated in Tab. IV, FRAMU consistently outperformed its counterparts. Specifically, in the AMPds2 dataset, FRAMU's RE and AD values (0.024 and 0.57, respectively) were superior to those of FedLU (0.03 and 0.66). Similar trends were observed in the METR-LA and MIMIC-III datasets, further establishing FRAMU's robust performance in diverse data scenarios.

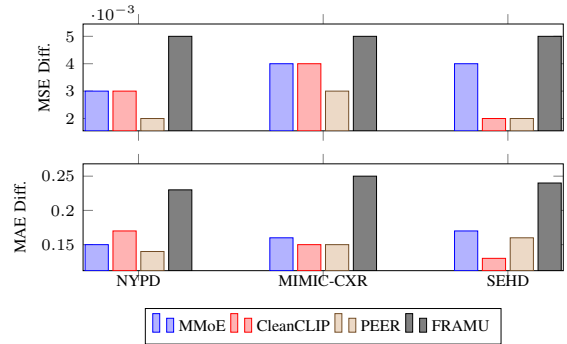


Fig. 7: Comparative Analysis of MSE and MAE Differences between Original and Unlearned multimodality Data

E. FRAMU Unlearning Results in Multimodality Context

In the multimodality experiment, the FRAMU framework demonstrated its capability to handle diverse data types, including images, text, and sensor data. The aim was to assess FRAMU's effectiveness in unlearning outdated, private, and irrelevant data in a multimodal context. For this, we utilized benchmark datasets like MIMIC-CXR [50], NYPD Complaint Data [49], and SHED [51]. The key focus was on evaluating error reduction and performance improvements in comparison to baseline models, with p-values highlighting the statistical significance of FRAMU's advancements.

1) Outdated Data

FRAMU consistently outperformed baseline models across all datasets in handling outdated data. In the NYPD Complaint Data [49], for instance, it achieved a lower MSE (0.047) and MAE (5.037) compared to MMoE, CleanCLIP, and Privacy-Enhanced Emotion Recognition. Similar trends were observed in the MIMIC-CXR [50] and SHED [51] datasets. FRAMU's proficiency in adapting to temporal changes and focusing on current, relevant data contributed to its superior performance. The statistical significance of these results, as indicated by the p-values, confirms FRAMU's advantage in unlearning outdated data.

2) Private Data

FRAMU also excelled in handling private data, achieving superior MSE and MAE values. In the NYPD Complaint Data,

TABLE V: FRAMU - Evaluation Results in Multimodality Context

Unlearning	Dataset		MMoE [54]			CleanCLIP [55]			PEER [56]			FRAMU (ours)	
			MSE	MAE	p-value	MSE	MAE	p-value	MSE	MAE	p-value	MSE	MAE
Outdated Data	Original	NYPD	0.064	7.28	0.024	0.062	6.95	0.031	0.06	6.41	0.041	0.055	5.77
		MIMIC-CXR	0.075	8.71	0.016	0.079	8.31	0.038	0.074	7.67	0.029	0.071	6.9
		SHED	0.095	11.27	0.031	0.098	10.76	0.045	0.093	9.92	0.032	0.089	8.93
	Unlearned	NYPD	0.061	7.13	0.015	0.059	6.78	0.029	0.058	5.71	0.036	0.042	4.54
		MIMIC-CXR	0.071	8.55	0.029	0.075	8.12	0.025	0.07	6.84	0.032	0.052	5.45
		SHED	0.091	11.1	0.023	0.094	10.54	0.041	0.09	9.76	0.017	0.067	7.07
Private Data	Original	NYPD	0.053	7.33	0.014	0.055	7	0.037	0.054	6.45	0.041	0.051	5.81
		MIMIC-CXR	0.063	8.76	0.035	0.065	8.36	0.043	0.064	7.71	0.045	0.062	6.94
		SHED	0.078	11.34	0.035	0.08	10.82	0.044	0.079	9.98	0.031	0.077	8.98
	Unlearned	NYPD	0.051	7.17	0.011	0.053	6.82	0.035	0.052	6.31	0.015	0.039	4.57
		MIMIC-CXR	0.06	8.6	0.031	0.062	8.17	0.038	0.061	7.56	0.016	0.046	5.48
		SHED	0.075	11.17	0.011	0.077	10.61	0.036	0.076	9.81	0.029	0.058	7.11
Irrelevant Data	Original	NYPD	0.047	7.25	0.027	0.05	6.92	0.041	0.048	6.38	0.025	0.046	5.74
		MIMIC-CXR	0.054	8.66	0.038	0.056	8.27	0.031	0.055	7.63	0.039	0.053	6.87
		SHED	0.07	11.21	0.045	0.072	10.7	0.032	0.071	9.87	0.042	0.069	8.88
	Unlearned	NYPD	0.045	7.1	0.014	0.047	6.74	0.019	0.046	6.24	0.022	0.034	4.52
		MIMIC-CXR	0.052	8.5	0.029	0.054	8.08	0.024	0.053	7.48	0.027	0.04	5.42
		SHED	0.068	11.04	0.025	0.07	10.49	0.022	0.069	9.71	0.021	0.052	7.04

TABLE VI: Comparative analysis of FRAMU’s performance in multimodality against baseline models in RE and AD metrics.

Unlearning	Dataset	MMoE [54]		CleanCLIP [55]		PEER [56]		FRAMU (Ours)	
		RE	AD	RE	AD	RE	AD	RE	AD
Outdated Data	NYPD	0.029	0.71	0.028	0.68	0.029	0.57	0.022	0.45
	MIMIC-CXR	0.035	0.85	0.037	0.81	0.034	0.68	0.027	0.54
	SHED	0.045	1.11	0.047	1.05	0.045	0.97	0.035	0.7
Private Data	NYPD	0.031	0.71	0.031	0.68	0.031	0.63	0.023	0.46
	MIMIC-CXR	0.038	0.86	0.04	0.81	0.039	0.75	0.028	0.54
	SHED	0.046	1.11	0.048	1.06	0.047	0.98	0.036	0.71
Irrelevant Data	NYPD	0.028	0.71	0.029	0.67	0.028	0.62	0.021	0.45
	MIMIC-CXR	0.033	0.85	0.034	0.8	0.032	0.74	0.027	0.54
	SHED	0.043	1.1	0.044	1.04	0.043	0.97	0.035	0.7

it showed notable performance with an MSE of 0.043 and an MAE of 5.067. This trend was consistent in the MIMIC-CXR and SHED datasets. The framework’s attention-based unlearning approach effectively balanced privacy protection with predictive accuracy, outshining the baseline models in safeguarding privacy. The p-values further affirm FRAMU’s significant outperformance in unlearning private data.

3) Irrelevant Data

Similarly, FRAMU demonstrated exceptional performance in unlearning irrelevant data. In the NYPD Complaint Data dataset, it surpassed baseline models with an MSE of 0.038 and an MAE of 5.012. This pattern persisted in the MIMIC-CXR and SHED datasets. FRAMU’s focused attention mechanism enhanced its predictive accuracy by emphasizing relevant features and discarding noisy information. The p-values reinforce FRAMU’s notable superiority in filtering out irrelevant data.

Fig. 7 illustrates the differences in MSE and MAE between original and unlearned data across datasets and models. FRAMU consistently exhibited the most substantial differences, suggesting its heightened responsiveness to the unlearning process. Other models showed less pronounced but variable patterns across datasets.

In Tab. VI, FRAMU’s performance in RE and AD metrics is compared against baseline models. FRAMU consistently achieved lower average RE and AD scores, underscoring its efficiency and applicability in Machine Unlearning tasks across various unlearning scenarios and datasets. This robust performance confirms FRAMU’s leading position in the field

of multimodal Machine Unlearning.

F. Convergence Analysis

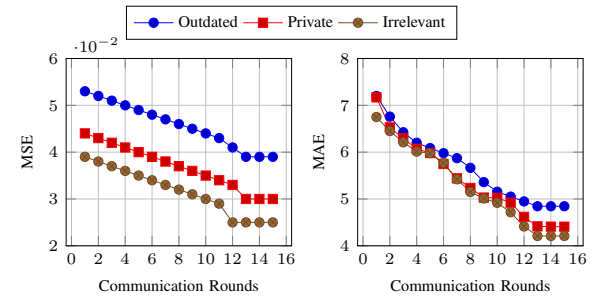


Fig. 8: Convergence Analysis

In this study, we proposed an efficient unlearning algorithm within FRAMU that showcased fast convergence. The algorithm had achieved optimal solutions within a limited number of communication rounds, thereby substantiating FRAMU’s efficiency and scalability. The convergence analysis of FRAMU, as shown in Fig. 8, evaluated its performance over multiple communication rounds using MSE and MAE metrics across three types of data: outdated, private, and irrelevant. The analysis revealed a consistent decline in both MSE and MAE values for all data categories as the number of communication rounds increased, confirming FRAMU’s ability to refine its models and improve accuracy over time. Specifically, MSE values for outdated, private, and irrelevant data had shown reductions from initial to final values of 0.053 to 0.039, 0.044 to 0.030, and 0.039 to 0.025, respectively. Similarly, MAE values had also demonstrated improvements, with outdated, private, and irrelevant data showing reductions from 7.201 to 4.845, 7.17 to 4.409, and 6.75 to 4.210, respectively.

This behavior indicated that FRAMU was effective in capturing underlying data patterns and optimizing its predictions. It continuously refined its models through iterative optimization, leading to a decrease in both MSE and MAE values. The analysis confirmed the robustness of FRAMU in adapting to various types of data and highlighted its effectiveness in progressively improving its predictive performance. Overall,

FRAMU’s strong convergence characteristics across different data categories have demonstrated its versatility and capability in minimizing errors, making it a robust choice for various federated learning applications.

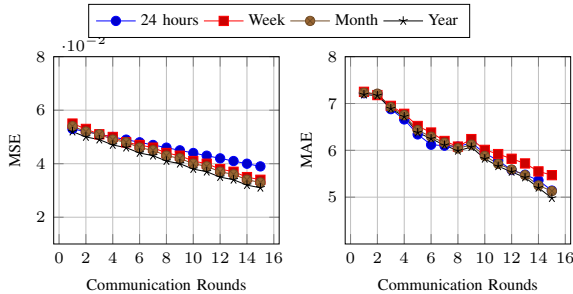


Fig. 9: Optimization Analysis - Outdated Data

G. Optimization

The performance of the FRAMU framework is evaluated through MSE and MAE metrics across various communication rounds and thresholds, as presented in Fig. 9 and Fig. 10. Fig. 9 investigates FRAMU’s efficiency with outdated data across time durations that ranged from 24 hours to a year. Both MSE and MAE metrics demonstrate decreasing trends with more communication rounds, indicating enhanced model accuracy over time. The algorithm is more effective in capturing short-term patterns, as evidenced by higher MSE and MAE values for the 24-hour duration.

Fig. 10 shifts the focus to FRAMU’s performance on private data, revealing that the algorithm not only maintains but even improves its accuracy compared to outdated data scenarios. Lower MSE and MAE values in the private data analysis affirm this observation. Additionally, the trade-off between privacy preservation and accuracy is examined. Although increasing privacy guarantees (lower ϵ values) generally leads to higher MSE and MAE, FRAMU still manages to maintain reasonable accuracy levels. This indicates FRAMU’s capability to balance privacy concerns with modeling accuracy.

VII. RESEARCH IMPLICATIONS

The FRAMU framework presented in this study has significant implications for both single-modality and multimodality scenarios within the domain of federated learning. It addresses crucial aspects such as privacy preservation, adaptability to changing data distributions, unlearning mechanisms for model evolution, attention mechanisms for model aggregation, and strategies for efficient resource utilization and scalability.

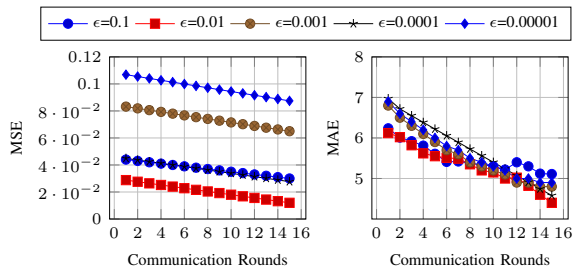


Fig. 10: Optimization Analysis - Private Data

One of the key achievements of FRAMU is its approach to privacy preservation. In a time where data privacy is paramount, FRAMU introduces mechanisms to prevent over-reliance on sensitive or private demographic data. Importantly, this emphasis on privacy does not detract from accuracy. Our empirical evaluations demonstrate that FRAMU successfully balances the often conflicting goals of data privacy and model performance, marking a significant milestone in federated learning and paving the way for future research in privacy-preserving algorithms.

Adaptability is another strength of FRAMU. Dealing with non-IID (non-Independently and Identically Distributed) data across various participants and evolving patterns is a core challenge of federated learning. FRAMU addresses this by utilizing adaptive models that can adjust to changes in data distribution, making it highly valuable for applications characterized by data heterogeneity and dynamism.

The unlearning mechanisms within FRAMU are also noteworthy. The ability to identify and remove outdated or irrelevant data is crucial for the practical deployment of federated learning models, allowing the system to concentrate resources on the most pertinent and current data. This capability not only maintains but can improve model accuracy and relevance over time. Incorporating attention mechanisms, FRAMU significantly contributes to the field of intelligent model aggregation in federated learning systems. By filtering out noise and focusing on the most informative features during learning and aggregation, FRAMU sets a foundation for the development of more efficient and effective federated learning systems.

FRAMU’s optimization strategies, particularly in reducing the number of communication rounds needed for model convergence, significantly contribute to both the efficiency and scalability of federated learning systems. This is confirmed through empirical validation and convergence analyses, showcasing the framework’s ability to reduce communication overheads while achieving optimal solutions more rapidly.

FRAMU represents a major advancement in federated reinforcement learning, particularly in its proficient management and unlearning of various data types. Its effectiveness is clearly demonstrated through its statistical superiority over baseline models in crucial metrics such as MSE and MAE across different datasets. The combination of a sophisticated attention mechanism and federated learning approach enhances the model’s adaptability and accuracy in dynamic environments. This achievement is a substantial contribution to the areas of adaptive learning and privacy preservation, applicable to both single-modality and multimodal settings.

VIII. CONCLUSION

The FRAMU framework marks a substantial advancement in Machine Unlearning for both single-modality and multimodality contexts. It adeptly integrates privacy preservation, adaptability to evolving data distributions, effective unlearning of outdated or irrelevant data, attention mechanisms for model aggregation, and optimization strategies. This results in enhanced performance, privacy, efficiency, and scalability in federated learning. Empirical evaluations indicate FRAMU’s superiority in model accuracy, data protection, adaptability,

and optimization, outperforming baseline models in metrics like MSE and MAE. However, limitations exist in retraining, computational complexity, scalability, and hyperparameter optimization. Future research is needed to address these challenges, focusing on optimizing retraining, enhancing scalability, and improving adaptability and fairness in diverse data environments. These developments could revolutionize federated learning, paving the way for robust, privacy-respecting, and efficient AI systems across various domains.

REFERENCES

- [1] P. Kumar, G. P. Gupta, and R. Tripathi, "Tp2sf: A trustworthy privacy-preserving secured framework for sustainable smart cities by leveraging blockchain and machine learning," *Journal of Systems Architecture*, vol. 115, p. 101954, 2021.
- [2] R. Nian, J. Liu, and B. Huang, "A review on reinforcement learning: Introduction and applications in industrial process control," *Computers & Chemical Engineering*, vol. 139, p. 106886, 2020.
- [3] O. A. Wahab, A. Mourad, H. Otok, and T. Taleb, "Federated machine learning: Survey, multi-level classification, desirable criteria and future directions in communication and networking systems," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 2, pp. 1342–1397, 2021.
- [4] M. Al-Rubaie and J. M. Chang, "Privacy-preserving machine learning: Threats and solutions," *IEEE Security & Privacy*, vol. 17, no. 2, pp. 49–58, 2019.
- [5] E. Politou, E. Alepis, and C. Patsakis, "Forgetting personal data and revoking consent under the gdpr: Challenges and proposed solutions," *Journal of cybersecurity*, vol. 4, no. 1, p. tty001, 2018.
- [6] J. Globocnik, "The right to be forgotten is taking shape: Cjeu judgments in gc and others (c-136/17) and google v cnil (c-507/17)," *GRUR International*, vol. 69, no. 4, pp. 380–388, 2020.
- [7] C. N. Fortner, "Decision making within a cancel culture environment," tech. rep., US Army Command and General Staff College, 2020.
- [8] K. Vasilevski, "Meta-learning for clinical and imaging data fusion for improved deep learning inference," 2023.
- [9] M. Sun, J. Xiao, E. G. Lim, C. Zhao, and Y. Zhao, "Unified multi-modality video object segmentation using reinforcement learning," *IEEE Transactions on Circuits and Systems for Video Technology*, 2023.
- [10] A. Malekloo, E. Ozer, M. AlHamaydeh, and M. Girolami, "Machine learning and structural health monitoring overview with emerging technology and high-dimensional data source highlights," *Structural Health Monitoring*, vol. 21, no. 4, pp. 1906–1955, 2022.
- [11] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, "Multimodal deep learning," in *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 689–696, 2011.
- [12] T. Zhou, S. Ruan, and S. Canu, "A review: Deep learning for medical image segmentation using multi-modality fusion," *Array*, vol. 3, p. 100004, 2019.
- [13] J. D. Fernández, S. P. Menci, C. M. Lee, A. Rieger, and G. Fridgen, "Privacy-preserving federated learning for residential short-term load forecasting," *Applied energy*, vol. 326, p. 119915, 2022.
- [14] D. C. Nguyen, Q.-V. Pham, P. N. Pathirana, M. Ding, A. Seneviratne, Z. Lin, O. Dobre, and W.-J. Hwang, "Federated learning for smart healthcare: A survey," *ACM Computing Surveys (CSUR)*, vol. 55, no. 3, pp. 1–37, 2022.
- [15] Q. Li, Z. Wen, Z. Wu, S. Hu, N. Wang, Y. Li, X. Liu, and B. He, "A survey on federated learning systems: Vision, hype and reality for data privacy and protection," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [16] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pp. 308–318, 2016.
- [17] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang, "Learning differentially private recurrent language models," *arXiv preprint arXiv:1710.06963*, 2017.
- [18] X. Chen and B. Wujek, "A unified framework for automatic distributed active learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 12, pp. 9774–9786, 2021.
- [19] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot, "Machine unlearning," in *2021 IEEE Symposium on Security and Privacy (SP)*, pp. 141–159, IEEE, 2021.
- [20] M. Jegorova, C. Kaul, C. Mayor, A. Q. O'Neil, A. Weir, R. Murray-Smith, and S. A. Tsafaris, "Survey: Leakage and privacy at inference time," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [21] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*, pp. 1273–1282, PMLR, 2017.
- [22] Z. Wu, H. Wang, Z. Wang, H. Jin, and Z. Wang, "Privacy-preserving deep action recognition: An adversarial learning framework and a new dataset," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 4, pp. 2126–2139, 2020.
- [23] J. Liang, Z. Liu, J. Zhou, X. Jiang, C. Zhang, and F. Wang, "Model-protected multi-task learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 2, pp. 1002–1019, 2020.
- [24] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pp. 265–284, Springer, 2006.
- [25] P. Zhou, K. Wang, L. Guo, S. Gong, and B. Zheng, "A privacy-preserving distributed contextual federated online learning framework with big data support in social recommender systems," *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 3, pp. 824–838, 2019.
- [26] Z. Ma, Y. Liu, Y. Miao, G. Xu, X. Liu, J. Ma, and R. H. Deng, "Flgan: Gan-based unbiased federated learning under non-iid settings," *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [27] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE signal processing magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [28] A. K. Sahu, T. Li, M. Sanjabi, M. Zaheer, A. Talwalkar, and V. Smith, "On the convergence of federated optimization in heterogeneous networks," *ArXiv*, vol. abs/1812.06127, 2018.
- [29] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," *arXiv preprint arXiv:1610.05492*, 2016.
- [30] T. Zhu, D. Ye, W. Wang, W. Zhou, and S. Y. Philip, "More than privacy: Applying differential privacy in key areas of artificial intelligence," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 6, pp. 2824–2843, 2020.
- [31] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pp. 1175–1191, 2017.
- [32] T. Shaik, X. Tao, N. Higgins, R. Gururajan, Y. Li, X. Zhou, and U. R. Acharya, "Fedstack: Personalized activity monitoring using stacked federated learning," *Knowledge-Based Systems*, vol. 257, p. 109929, 2022.
- [33] T. S. Brisimi, R. Chen, T. Mela, A. Olshesky, I. C. Paschalidis, and W. Shi, "Federated learning of predictive models from federated electronic health records," *International journal of medical informatics*, vol. 112, pp. 59–67, 2018.
- [34] W. Huang, J. Liu, T. Li, T. Huang, S. Ji, and J. Wan, "Feddsr: Daily schedule recommendation in a federated deep reinforcement learning framework," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [35] T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey, "Meta-learning in neural networks: A survey," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 9, pp. 5149–5169, 2021.
- [36] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," 2013.
- [37] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," *arXiv preprint arXiv:1511.05952*, 2015.
- [38] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fiedjeland, G. Ostrovski, et al., "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [39] Z. Bing, D. Lerch, K. Huang, and A. Knoll, "Meta-reinforcement learning in non-stationary and dynamic environments," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 3476–3491, 2022.
- [40] G. Brauwers and F. Frasincar, "A general survey on attention mechanisms in deep learning," *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [41] I. Sorokin, A. Seleznev, M. Pavlov, A. Fedorov, and A. Ignateva, "Deep attention recurrent q-network," *arXiv preprint arXiv:1512.01693*, 2015.

- [42] Z. Guan, Y. Li, Z. Pan, Y. Liu, and Z. Xue, "Rfdg: Reinforcement federated domain generalization," *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [43] H. Miao, X. Zhong, J. Liu, Y. Zhao, X. Zhao, W. Qian, K. Zheng, and C. S. Jensen, "Task assignment with efficient federated preference learning in spatial crowdsourcing," *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [44] S. Ek, F. Portet, P. Lalanda, and G. Vega, "A federated learning aggregation algorithm for pervasive computing: Evaluation and comparison," in *2021 IEEE International Conference on Pervasive Computing and Communications (PerCom)*, pp. 1–10, 2021.
- [45] X. Ouyang, Z. Xie, J. Zhou, J. Huang, and G. Xing, "ClusterFL," in *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services*, ACM, June 2021.
- [46] S. Makonin, B. Ellert, I. V. Bajić, and F. Popowich, "Electricity, water, and natural gas consumption of a residential house in Canada from 2012 to 2014," *Scientific data*, vol. 3, no. 1, pp. 1–12, 2016.
- [47] Y. Li, R. Yu, C. Shahabi, and Y. Liu, "Diffusion convolutional recurrent neural network: Data-driven traffic forecasting," in *International Conference on Learning Representations (ICLR '18)*, 2018.
- [48] A. E. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. Anthony Celi, and R. G. Mark, "Mimic-iii, a freely accessible critical care database," *Scientific data*, vol. 3, no. 1, pp. 1–9, 2016.
- [49] P. D. (NYPD), "Nypd complaint data current (year to date): Nyc open data." *NYPD Complaint Data Current (Year To Date) — NYC Open Data*, Apr 2023.
- [50] A. E. Johnson, T. J. Pollard, S. J. Berkowitz, N. R. Greenbaum, M. P. Lungren, C.-y. Deng, R. G. Mark, and S. Horng, "Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports," *Scientific data*, vol. 6, no. 1, p. 317, 2019.
- [51] K. Dataset, "Smart home dataset with weather information," 2019.
- [52] X. Zhu, G. Li, and W. Hu, "Heterogeneous federated knowledge graph embedding learning and unlearning," in *Proceedings of the ACM Web Conference 2023*, pp. 2444–2454, 2023.
- [53] V. S. Chundawat, A. K. Tarun, M. Mandal, and M. Kankanhalli, "Zero-shot machine unlearning," *IEEE Transactions on Information Forensics and Security*, 2023.
- [54] J. Ma, Z. Zhao, X. Yi, J. Chen, L. Hong, and E. H. Chi, "Modeling task relationships in multi-task learning with multi-gate mixture-of-experts," in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1930–1939, 2018.
- [55] H. Bansal, N. Singhi, Y. Yang, F. Yin, A. Grover, and K.-W. Chang, "Cleanclip: Mitigating data poisoning attacks in multimodal contrastive learning," *arXiv preprint arXiv:2303.03323*, 2023.
- [56] M. Jaiswal and E. M. Provost, "Privacy enhanced multimodal neural representations for emotion recognition," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 7985–7993, 2020.
- [57] C. Tofallis, "A better measure of relative prediction accuracy for model selection and model estimation," *Journal of the Operational Research Society*, vol. 66, pp. 1352–1362, 2015.
- [58] T. Chai and R. R. Draxler, "Root mean square error (rmse) or mean absolute error (mae)," *Geoscientific model development discussions*, vol. 7, no. 1, pp. 1525–1534, 2014.
- [59] K. Y. Tan, L. Yueming, Y.-S. Ong, and I. Tsang, "Unfolded self-reconstruction lsh: Towards machine unlearning in approximate nearest neighbour search," *arXiv preprint arXiv:2304.02350*, 2023.
- [60] A. Golatkar, A. Achille, and S. Soatto, "Eternal sunshine of the spotless net: Selective forgetting in deep networks," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, (Los Alamitos, CA, USA), pp. 9301–9309, IEEE Computer Society, 2020.

11.2 Summary

The chapter concludes by highlighting FRAMU's significant contributions to the field of machine unlearning, showcasing its effectiveness in addressing the dual challenges of maintaining model accuracy while ensuring data privacy. Through extensive experiments, FRAMU demonstrates superior performance in unlearning outdated, private, and irrelevant data across various datasets. The summary emphasizes the framework's adaptability, privacy preservation, and efficiency in federated learning environments, offering valuable insights into future directions for research and development in adaptive and privacy-preserving machine learning models.

CHAPTER 12: CONCLUSIONS

This doctoral thesis embarked on a transformative journey into the realm of patient monitoring systems, harnessing the potential of AI and cutting-edge techniques such as federated learning, reinforcement learning, and machine unlearning. Through a comprehensive investigation of various AI in healthcare scenarios, the thesis made significant strides in enhancing patient care and revolutionizing healthcare practices.

The first part of this research journey delved into using Artificial Intelligence (AI) to enhance patient monitoring systems. It focused on remote patient monitoring and personalized activity tracking, addressing the significant challenges with innovative solutions like FedStack and Clustered FedStack. These models used stacked federated learning to improve personalized care, allowing healthcare providers to offer more specific and effective interventions based on individual patient needs. Additionally, the development of Multi-Agent Deep Reinforcement Learning and PDRL frameworks was crucial, offering innovative methods for predictive monitoring and enabling the early detection of health anomalies, which is vital for providing proactive and personalized treatments.

The second part of this research extended the exploration into multimodality fusion and graph-enabled techniques, aiming to develop a holistic smart healthcare system. The synthesis of diverse AI techniques, as outlined in this section, emphasized the transformative potential of integrating various data sources for informed decision-making in healthcare. Furthermore, the GraphRL framework showcased the capability of Temporal Graphical Convolutional Networks (T-GCN) in enhancing dynamic reinforcement learning scenarios in patient monitoring, indicating a broad spectrum of potential applications beyond the healthcare sector.

In the third part, the focus shifted to explainable AI and machine unlearning, unveiling the importance of transparency and interpretability in AI-driven systems. The development of the QXAI framework highlighted the necessity for explainability in healthcare applications of AI, aiming to strengthen trust and facilitate more effective collaboration between AI systems and human practitioners. Moreover, the exploration of machine unlearning underscored its role in adapting AI models to the continuously evolving landscape of healthcare, with the FRAMU framework illustrating the convergence of federated reinforcement learning and attention-based machine unlearning to ensure the robustness of AI models.

This thesis represents a comprehensive exploration into the realms of AI-driven patient monitoring systems. It has traversed the landscapes of enhanced patient monitoring, smart healthcare systems, and the subtleties of explainable AI and machine unlearning. By integrating innovative AI technologies into healthcare, this research illuminates the possibilities of creating a seamless and symbiotic relationship between technological advancements and healthcare needs. The insights and frameworks presented herein are hoped to serve as catalysts for future research, outlining a progressive path forward in the pursuit of more personalized and impactful healthcare solutions.

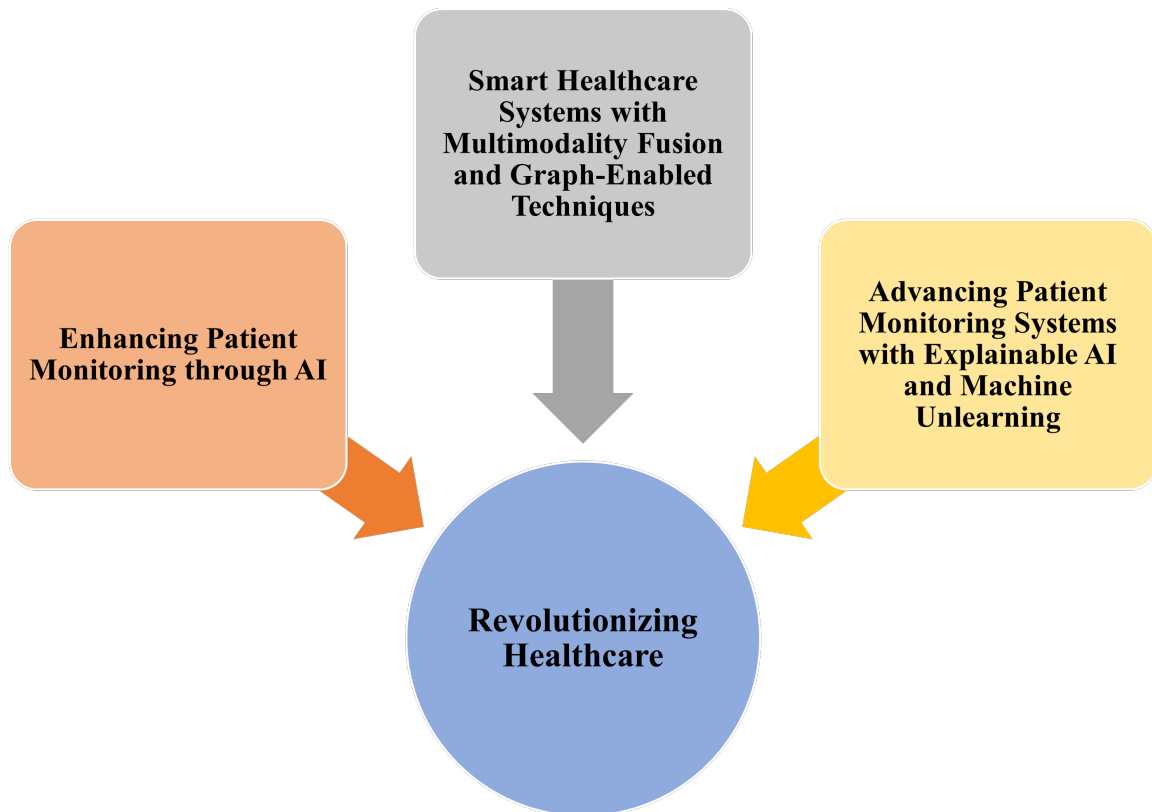


Figure 12.1: Contributions in Revolutionizing Healthcare

12.1 Contributions

This research explores various aspects of AI methodologies, encompassing federated learning, reinforcement learning, attention mechanisms, and machine unlearning, segmented into three parts within this thesis, with an aim to revolutionize healthcare and refine patient monitoring systems As shown in Fig.12.1. The objective is to enable the prediction of vital signs and the classification of physical activities while ensuring the provision of transparent, personalized, and patient-centric AI mechanisms. This study stands as a multi-dimensional contribution to the domain of AI in Healthcare. It primarily focuses on amplifying the efficacy of remote patient monitoring using advanced AI techniques and underscores the importance of multimodal fusion in the context of smart healthcare systems. Furthermore, this research promotes the integration of transparent and explainable AI to advance patient monitoring systems and employs machine unlearning to safeguard patient privacy. This approach also works to enhance model accuracy by unlearning outdated and irrelevant data, thereby aligning the system with contemporary healthcare requirements. The major contributions of our work in each part of this thesis are as follows:

Part I: Enhancing Patient Monitoring through AI

- We present a comprehensive review of AI's impact on remote patient monitoring systems and identify challenges in adopting AI-enabled remote patient monitoring systems
- We proposed novel federated learning approaches such as FedStack and Clustered FedStack enabling personalized monitoring in human activity recognition.

- We designed multi-agent deep reinforcement learning frameworks that have re-shaped patient monitoring by learning behavior patterns and predicting appropriate responses, allowing medical teams to act proactively during emergencies.

Part II: Smart healthcare systems with multimodality fusion and graph-enabled techniques

- We present the journey of data to information to knowledge to wisdom in the context of multimodality fusion for smart healthcare.
- We proposed a generic DIKW techniques framework for smart healthcare, that not only highlights the current efforts but also provides a vision for its future evolution.
- We proposed a novel and generic GraphRL framework with predictive and monitoring capabilities for early warnings in a complex environment.

Part III: Advancing patient monitoring systems with explainable AI and machine unlearning

- We propose a new paradigm to interpret and explain the vital sign prediction and physical activity classification in patient monitoring systems through a generic Explainable AI framework(QXAI).
- We present a detailed taxonomy of techniques in machine unlearning that can be adopted in natural language processing (NLP), computer vision, and recommender systems which act as important roles in patient monitoring systems.
- We present a novel adaptive unlearning framework to unlearn outdated, private, and irrelevant data to protect patient privacy and enhance model accuracy by removing outdated and irrelevant data.

12.2 Limitations

The doctoral thesis acknowledges certain limitations that need attention:

Scope of applications: The study primarily focused on remote patient monitoring, personalized activity tracking, and predictive monitoring in healthcare. Future research should encompass a broader range of healthcare applications, including neurological system-related diseases and other chronic conditions.

Data scale and explainability: Some frameworks faced challenges related to data scale and explainability. Ensuring the robustness and reliability of AI-driven decisions in patient monitoring requires addressing these limitations. Models should not only deliver accurate predictions but also provide clear explanations to healthcare professionals for informed decision-making.

Data privacy and security: AI-driven patient monitoring relies on diverse data sources, making data privacy and security critical concerns. Future research should focus on developing privacy-preserving techniques and protocols to safeguard patient data while enabling efficient and effective AI model training.

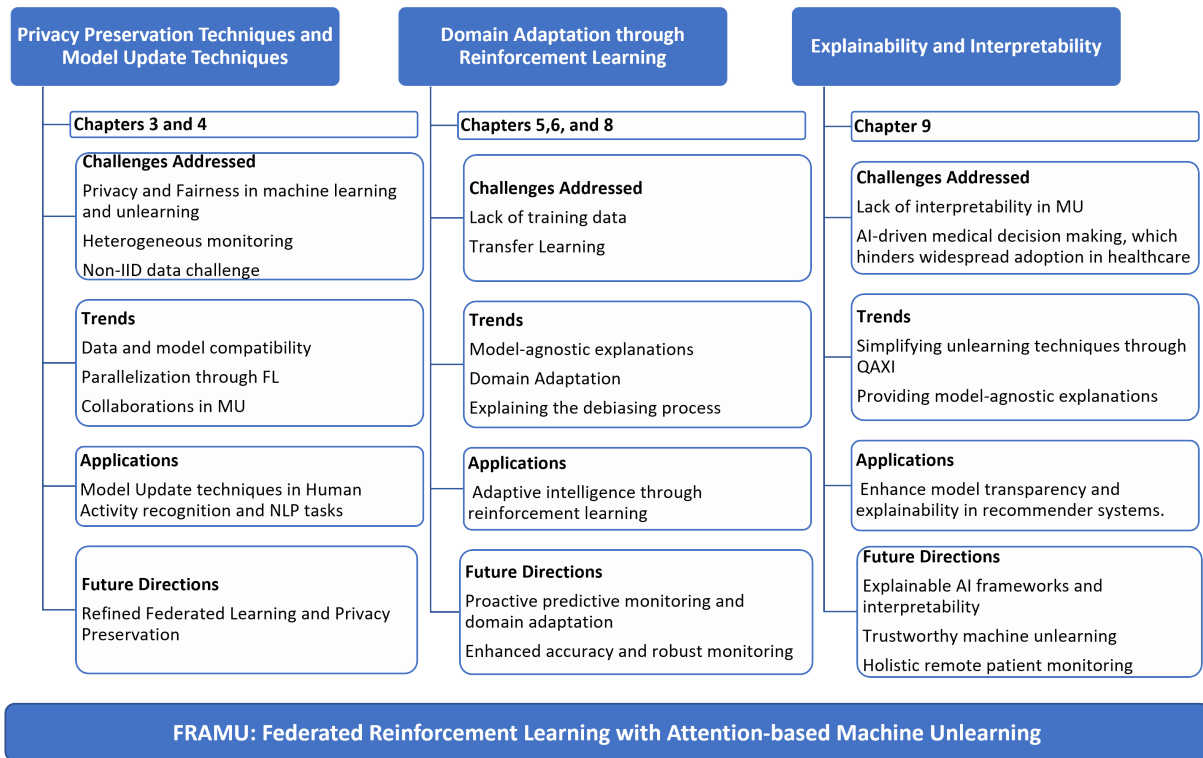


Figure 12.2: Machine Unlearning Taxonomy

12.3 Future directions

The doctoral thesis, by incorporating insights from the multifaceted machine learning and machine unlearning taxonomy, sets the stage for groundbreaking research and advancements in AI-driven patient monitoring systems as shown in Fig. 12.2, illustrating convergent paths between future directions and the taxonomy’s core components:

Refined Federated Learning and Privacy Preservation: Integrating insights from the first axis of the taxonomy, focused on Privacy Preservation and Model Update Techniques, future endeavours will delve into refining federated learning frameworks such as Clustered-FedStack through dynamic clustering and enhanced incorporation of client demographics [50]. These nuanced advancements are anticipated to orchestrate more personalized, improved outcomes and advance privacy preservation, a pivotal aspect of machine unlearning.

Proactive predictive monitoring and domain adaptation: Aligned with the second axis emphasizing Domain Adaptation through Reinforcement Learning, exploring proactive predictive capabilities in patient monitoring using multi-agent DRL exemplifies an innovative frontier in preventive patient care [67]. By resolving challenges related to limited training data and leveraging transfer learning, this direction is poised to offer adaptive intelligence, making monitoring frameworks more versatile and domain-adaptable [68].

Enhanced accuracy and robust monitoring: The synergistic exploration of ensem-

ble methods and transfer learning resonates with the aspirations of the taxonomy's second axis. The union of diverse insights from multiple DRL agents is set to fortify the robustness and accuracy of patient monitoring systems [69], addressing the ever-evolving needs and challenges in healthcare.

Holistic remote patient monitoring: Envisioned advancements in the integration of a broader spectrum of feature inputs and patient data sources align with the overarching goals of the machine unlearning taxonomy [70], aiming to provide comprehensive remote monitoring solutions and augmented clinical decision support, thereby pushing the boundaries of AI in healthcare.

Trustworthy machine unlearning: Continuing advancements in machine unlearning, focusing on enhanced transparency and trustworthiness of AI models, align with the intrinsic objectives of the taxonomy [71]. Addressing aspects such as scalability, selective unlearning, and performance impacts is pivotal for the widespread and ethical deployment of machine unlearning across various sectors.

Explainable AI frameworks and interpretability: Further refinement and exploration of frameworks like QXAI, drawing insights from the third axis concentrating on Explainability and Interpretability, aim to bring clarity and transparency to AI-driven predictions, fostering enriched collaborations between AI and human insights and enhancing the overall healthcare experience [72].

This doctoral thesis has significantly contributed to advancing knowledge and practice in AI-driven patient monitoring systems. By addressing limitations and embracing future directions, the transformative potential of AI in healthcare is within reach. Continued research and innovation in AI, combined with a patient-centric approach, will usher in a new era of personalized, proactive, and effective healthcare delivery, benefiting individuals worldwide. The collective efforts of researchers and practitioners in the field will shape the future of healthcare, where AI-driven patient monitoring becomes an indispensable tool in enhancing patient well-being and transforming healthcare practices.

REFERENCES

- [1] Y. K. Dwivedi, L. Hughes, E. Ismagilova, G. Aarts, C. Coombs, T. Crick, Y. Duan, R. Dwivedi, J. Edwards, A. Eirug, V. Galanos, P. V. Ilavarasan, M. Janssen, P. Jones, A. K. Kar, H. Kizgin, B. Kronemann, B. Lal, B. Lucini, R. Medaglia, K. L. Meunier-FitzHugh, L. C. L. Meunier-FitzHugh, S. Misra, E. Mogaji, S. K. Sharma, J. B. Singh, V. Raghavan, R. Raman, N. P. Rana, S. Samothrakis, J. Spencer, K. Tamilmani, A. Tubadji, P. Walton, and M. D. Williams, "Artificial intelligence (AI): Multidisciplinary perspectives on emerging challenges, opportunities, and agenda for research, practice and policy," *International Journal of Information Management*, vol. 57, p. 101994, Apr. 2021. [Online]. Available: <https://doi.org/10.1016/j.ijinfomgt.2019.08.002>
- [2] F. D. Carlo, A. Sociali, E. Picutti, M. Pettoruso, F. Vellante, V. Verrastro, G. Martinotti, and M. di Giannantonio, "Telepsychiatry and other cutting-edge technologies in COVID-19 pandemic: Bridging the distance in mental health assistance," *International Journal of Clinical Practice*, vol. 75, no. 1, Oct. 2020. [Online]. Available: <https://doi.org/10.1111/ijcp.13716>
- [3] A. Mahajan, T. Vaidya, A. Gupta, S. Rane, and S. Gupta, "Artificial intelligence in healthcare in developing nations: The beginning of a transformative journey," *Cancer Research, Statistics, and Treatment*, vol. 2, no. 2, p. 182, 2019. [Online]. Available: https://doi.org/10.4103/crst.crst_50_19
- [4] R. I. Sifat and U. Bhattacharya, "Transformative potential of artificial intelligence in global health policy," Jun. 2023. [Online]. Available: <https://doi.org/10.1080/20016689.2023.2230660>
- [5] A. Panesar, "What is artificial intelligence?" in *Machine Learning and AI for Healthcare*. Apress, Dec. 2020, pp. 1–18. [Online]. Available: https://doi.org/10.1007/978-1-4842-6537-6_1
- [6] H. Keserwani, S. V. Kakade, S. K. Sharma, M. Manchanda, and G. F. Nama, "Real-time analysis of wearable sensor data using iot and machine learning in healthcare," *International Journal of Intelligent Systems and Applications in Engineering*, vol. 11, no. 7s, p. 85 –, Jul. 2023. [Online]. Available: <https://ijisae.org/index.php/IJISAE/article/view/2934>
- [7] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso, "Machine learning interpretability: A survey on methods and metrics," *Electronics*, vol. 8, no. 8, p. 832, Jul. 2019. [Online]. Available: <https://doi.org/10.3390/electronics8080832>
- [8] Y. Bengio, A. Lodi, and A. Prouvost, "Machine learning for combinatorial optimization: A methodological tour d'horizon," *European Journal of Operational*

- Research*, vol. 290, no. 2, pp. 405–421, Apr. 2021. [Online]. Available: <https://doi.org/10.1016/j.ejor.2020.07.063>
- [9] A. Esteva, A. Robicquet, B. Ramsundar, V. Kuleshov, M. DePristo, K. Chou, C. Cui, G. Corrado, S. Thrun, and J. Dean, “A guide to deep learning in healthcare,” *Nature Medicine*, vol. 25, no. 1, pp. 24–29, Jan. 2019. [Online]. Available: <https://doi.org/10.1038/s41591-018-0316-z>
- [10] D. Shen, G. Wu, and H.-I. Suk, “Deep learning in medical image analysis,” *Annual Review of Biomedical Engineering*, vol. 19, no. 1, pp. 221–248, Jun. 2017. [Online]. Available: <https://doi.org/10.1146/annurev-bioeng-071516-044442>
- [11] T. F. Collura, “History and evolution of electroencephalographic instruments and techniques,” *Journal of Clinical Neurophysiology*, vol. 10, no. 4, pp. 476–504, Oct. 1993. [Online]. Available: <https://doi.org/10.1097/00004691-199310000-00007>
- [12] S. Majumder, T. Mondal, and M. Deen, “Wearable sensors for remote health monitoring,” *Sensors*, vol. 17, no. 12, p. 130, Jan. 2017. [Online]. Available: <https://doi.org/10.3390/s17010130>
- [13] C. Chen, S. Ding, and J. Wang, “Digital health for aging populations,” *Nature Medicine*, vol. 29, no. 7, pp. 1623–1630, Jul. 2023. [Online]. Available: <https://doi.org/10.1038/s41591-023-02391-8>
- [14] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. A. Sallab, S. Yogamani, and P. Perez, “Deep reinforcement learning for autonomous driving: A survey,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 6, pp. 4909–4926, Jun. 2022. [Online]. Available: <https://doi.org/10.1109/tits.2021.3054625>
- [15] C. Yu, J. Liu, S. Nemati, and G. Yin, “Reinforcement learning in healthcare: A survey,” *ACM Computing Surveys*, vol. 55, no. 1, pp. 1–36, Nov. 2021. [Online]. Available: <https://doi.org/10.1145/3477600>
- [16] A. F. Lehman, L. B. Dixon, T. H. McGlashan, A. L. Miller, and D. O. Perkins, “Treatment of patients with schizophrenia,” *American Psychiatric Association*, vol. 1, 2010. [Online]. Available: <http://ajp.psychiatryonline.org/cgi/content/abstract/161/2/1>
- [17] S. Srinivas and A. R. Ravindran, “Optimizing outpatient appointment system using machine learning algorithms and scheduling rules: A prescriptive analytics framework,” *Expert Systems with Applications*, vol. 102, pp. 245–261, Jul. 2018. [Online]. Available: <https://doi.org/10.1016/j.eswa.2018.02.022>
- [18] D. Sow, D. S. Turaga, and M. Schmidt, “Mining of sensor data in healthcare: A survey,” in *Managing and Mining Sensor Data*. Springer US, Dec. 2012, pp. 459–504. [Online]. Available: https://doi.org/10.1007/978-1-4614-6309-2_14
- [19] S. Tinelli and I. Juran, “Artificial intelligence-based monitoring system of water quality parameters for early detection of non-specific bio-contamination in water distribution systems,” *Water Supply*, vol. 19, no. 6, pp. 1785–1792, Apr. 2019. [Online]. Available: <https://doi.org/10.2166/ws.2019.057>

- [20] K. Lawrence, N. Singh, Z. Jonassen, L. L. Groom, V. A. Arias, S. Mandal, A. Schoenthaler, D. Mann, O. Nov, and G. Dove, "Operational implementation of remote patient monitoring within a large ambulatory health system: Multimethod qualitative case study," *JMIR Human Factors*, vol. 10, p. e45166, Jul. 2023. [Online]. Available: <https://doi.org/10.2196/45166>
- [21] M. M. Baig, H. GholamHosseini, A. A. Moqem, F. Mirza, and M. Lindén, "A systematic review of wearable patient monitoring systems – current challenges and opportunities for clinical adoption," *Journal of Medical Systems*, vol. 41, no. 7, Jun. 2017. [Online]. Available: <https://doi.org/10.1007/s10916-017-0760-1>
- [22] M. Azimi, A. Eslamlou, and G. Pekcan, "Data-driven structural health monitoring and damage detection through deep learning: State-of-the-art review," *Sensors*, vol. 20, no. 10, p. 2778, May 2020. [Online]. Available: <https://doi.org/10.3390/s20102778>
- [23] V. Patel, A. Chesmore, C. M. Legner, and S. Pandey, "Trends in workplace wearable technologies and connected-worker solutions for next-generation occupational safety, health, and productivity," *Advanced Intelligent Systems*, vol. 4, no. 1, Sep. 2021. [Online]. Available: <https://doi.org/10.1002/aisy.202100099>
- [24] Z. Ahmed, K. Mohamed, S. Zeeshan, and X. Dong, "Artificial intelligence with multi-functional machine learning platform development for better healthcare and precision medicine," *Database*, vol. 2020, Jan. 2020. [Online]. Available: <https://doi.org/10.1093/database/baaa010>
- [25] J. Xu, Z. Wu, C. Wang, and X. Jia, "Machine unlearning: Solutions and challenges," 2023. [Online]. Available: <https://arxiv.org/abs/2308.07061>
- [26] P. Solanki, J. Grundy, and W. Hussain, "Operationalising ethics in artificial intelligence for healthcare: a framework for AI developers," *AI and Ethics*, vol. 3, no. 1, pp. 223–240, Jul. 2022. [Online]. Available: <https://doi.org/10.1007/s43681-022-00195-z>
- [27] Z. Ahmed, K. Mohamed, S. Zeeshan, and X. Dong, "Artificial intelligence with multi-functional machine learning platform development for better healthcare and precision medicine," *Database*, vol. 2020, Jan. 2020. [Online]. Available: <https://doi.org/10.1093/database/baaa010>
- [28] R. Najjar, "Redefining radiology: A review of artificial intelligence integration in medical imaging," *Diagnostics*, vol. 13, no. 17, p. 2760, Aug. 2023. [Online]. Available: <https://doi.org/10.3390/diagnostics13172760>
- [29] M. J. Coye, A. Haselkorn, and S. DeMello, "Remote patient management: Technology-enabled innovation and evolving business models for chronic disease care," *Health Affairs*, vol. 28, no. 1, pp. 126–135, Jan. 2009. [Online]. Available: <https://doi.org/10.1377/hlthaff.28.1.126>
- [30] B. C. Stahl, N. F. Doherty, and M. Shaw, "Information security policies in the UK healthcare sector: a critical evaluation," *Information Systems Journal*, vol. 22, no. 1, pp. 77–94, Jul. 2011. [Online]. Available: <https://doi.org/10.1111/j.1365-2575.2011.00378.x>

- [31] A. K. Tyagi, S. Aswathy, G. Aghila, and N. Sreenath, "AARIN: Affordable, accurate, reliable and innovative mechanism to protect a medical cyber-physical system using blockchain technology," *International Journal of Intelligent Networks*, vol. 2, pp. 175–183, 2021. [Online]. Available: <https://doi.org/10.1016/j.ijin.2021.09.007>
- [32] M. Shuaib, S. Alam, M. S. Alam, and M. S. Nasir, "WITHDRAWN: Compliance with HIPAA and GDPR in blockchain-based electronic health record," *Materials Today: Proceedings*, Mar. 2021. [Online]. Available: <https://doi.org/10.1016/j.matpr.2021.03.059>
- [33] D. C. Nguyen, Q.-V. Pham, P. N. Pathirana, M. Ding, A. Seneviratne, Z. Lin, O. Dobre, and W.-J. Hwang, "Federated learning for smart healthcare: A survey," *ACM Computing Surveys*, vol. 55, no. 3, pp. 1–37, Feb. 2022. [Online]. Available: <https://doi.org/10.1145/3501296>
- [34] J. Qi, P. Yang, G. Min, O. Amft, F. Dong, and L. Xu, "Advanced internet of things for personalised healthcare systems: A survey," *Pervasive and Mobile Computing*, vol. 41, pp. 132–149, Oct. 2017. [Online]. Available: <https://doi.org/10.1016/j.pmcj.2017.06.018>
- [35] P.-Y. Wu, C.-W. Cheng, C. D. Kaddi, J. Venugopalan, R. Hoffman, and M. D. Wang, "–omic and electronic health record big data analytics for precision medicine," *IEEE Transactions on Biomedical Engineering*, vol. 64, no. 2, pp. 263–273, 2017.
- [36] M. Verma, R. Hontecillas, N. Tubau-Juni, V. Abedi, and J. Bassaganya-Riera, "Challenges in personalized nutrition and health," *Frontiers in Nutrition*, vol. 5, Nov. 2018. [Online]. Available: <https://doi.org/10.3389/fnut.2018.00117>
- [37] H. F. Nweke, Y. W. Teh, M. A. Al-garadi, and U. R. Alo, "Deep learning algorithms for human activity recognition using mobile and wearable sensor networks: State of the art and research challenges," *Expert Systems with Applications*, vol. 105, pp. 233–261, Sep. 2018. [Online]. Available: <https://doi.org/10.1016/j.eswa.2018.03.056>
- [38] P. C. Tiwari, R. Pal, M. J. Chaudhary, and R. Nath, "Artificial intelligence revolutionizing drug development: Exploring opportunities and challenges," *Drug Development Research*, Sep. 2023. [Online]. Available: <https://doi.org/10.1002/ddr.22115>
- [39] K. G. Moons, R. F. Wolff, R. D. Riley, P. F. Whiting, M. Westwood, G. S. Collins, J. B. Reitsma, J. Kleijnen, and S. Mallett, "PROBAST: A tool to assess risk of bias and applicability of prediction model studies: Explanation and elaboration," *Annals of Internal Medicine*, vol. 170, no. 1, p. W1, Jan. 2019. [Online]. Available: <https://doi.org/10.7326/m18-1377>
- [40] C. Díaz and M. Neubert, "Application of artificial intelligence in international decision-making processes in the healthcare industry," *International Journal of Teaching and Case Studies*, vol. 13, no. 4, p. 341, 2022. [Online]. Available: <https://doi.org/10.1504/ijtcs.2022.130320>

- [41] Y. Y. M. Aung, D. C. S. Wong, and D. S. W. Ting, "The promise of artificial intelligence: a review of the opportunities and challenges of artificial intelligence in healthcare," *British Medical Bulletin*, vol. 139, no. 1, pp. 4–15, Aug. 2021. [Online]. Available: <https://doi.org/10.1093/bmb/ldab016>
- [42] J. Sun, Q.-X. Dong, S.-W. Wang, Y.-B. Zheng, X.-X. Liu, T.-S. Lu, K. Yuan, J. Shi, B. Hu, L. Lu, and Y. Han, "Artificial intelligence in psychiatry research, diagnosis, and therapy," *Asian Journal of Psychiatry*, vol. 87, p. 103705, Sep. 2023. [Online]. Available: <https://doi.org/10.1016/j.ajp.2023.103705>
- [43] D. D. Luxton, S. L. Anderson, and M. Anderson, "Ethical issues and artificial intelligence technologies in behavioral and mental health care," in *Artificial Intelligence in Behavioral and Mental Health Care*. Elsevier, 2016, pp. 255–276. [Online]. Available: <https://doi.org/10.1016/b978-0-12-420248-1.00011-8>
- [44] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [45] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning," *ACM Transactions on Intelligent Systems and Technology*, vol. 10, no. 2, pp. 1–19, Jan. 2019. [Online]. Available: <https://doi.org/10.1145/3298981>
- [46] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, "Federated learning in mobile edge networks: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2031–2063, 2020. [Online]. Available: <https://doi.org/10.1109/comst.2020.2986024>
- [47] P. Kierkegaard, "Electronic health record: Wiring europe's healthcare," *Computer Law & Security Review*, vol. 27, no. 5, pp. 503–515, Sep. 2011. [Online]. Available: <https://doi.org/10.1016/j.clsr.2011.07.013>
- [48] M. Tasnim, A. J. Patinga, H. Shahriar, and S. Sneha, "Cardiovascular health management compliance with health insurance portability and accountability act," in *2023 IEEE 47th Annual Computers, Software, and Applications Conference (COMPSAC)*. IEEE, Jun. 2023. [Online]. Available: <https://doi.org/10.1109/compsac57700.2023.00218>
- [49] D. Gupta, O. Kayode, S. Bhatt, M. Gupta, and A. S. Tosun, "Hierarchical federated learning based anomaly detection using digital twins for smart healthcare," in *2021 IEEE 7th International Conference on Collaboration and Internet Computing (CIC)*. IEEE, Dec. 2021. [Online]. Available: <https://doi.org/10.1109/cic52973.2021.00013>
- [50] M. Aledhari, R. Razzak, R. M. Parizi, and F. Saeed, "Federated learning: A survey on enabling technologies, protocols, and applications," *IEEE Access*, vol. 8, pp. 140699–140725, 2020. [Online]. Available: <https://doi.org/10.1109/access.2020.3013541>
- [51] H. Tran-Dang, S. Bhardwaj, T. Rahim, A. Musaddiq, and D.-S. Kim, "Reinforcement learning based resource management for fog computing environment: Literature review, challenges, and open issues," *Journal of*

- Communications and Networks*, vol. 24, no. 1, pp. 83–98, Feb. 2022. [Online]. Available: <https://doi.org/10.23919/jcn.2021.000041>
- [52] D. Silver, S. Singh, D. Precup, and R. S. Sutton, “Reward is enough,” *Artificial Intelligence*, vol. 299, p. 103535, Oct. 2021. [Online]. Available: <https://doi.org/10.1016/j.artint.2021.103535>
- [53] J. Fan, C. Xiao, and Y. Huang, “Gdi: Rethinking what makes reinforcement learning different from supervised learning,” 2021. [Online]. Available: <https://arxiv.org/abs/2106.06232>
- [54] M. A. Wiering and M. Van Otterlo, “Reinforcement learning,” *Adaptation, learning, and optimization*, vol. 12, no. 3, p. 729, 2012. [Online]. Available: <https://doi.org/10.1007/978-3-642-27645-3>
- [55] L. P. Kaelbling, M. L. Littman, and A. W. Moore, “Reinforcement learning: A survey,” *Journal of Artificial Intelligence Research*, vol. 4, pp. 237–285, May 1996. [Online]. Available: <https://doi.org/10.1613/jair.301>
- [56] B. Jang, M. Kim, G. Harerimana, and J. W. Kim, “Q-learning algorithms: A comprehensive classification and applications,” *IEEE Access*, vol. 7, pp. 133 653–133 667, 2019. [Online]. Available: <https://doi.org/10.1109/access.2019.2941229>
- [57] A. B. Arrieta, N. Díaz-Rodríguez, J. D. Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera, “Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI,” *Information Fusion*, vol. 58, pp. 82–115, Jun. 2020. [Online]. Available: <https://doi.org/10.1016/j.inffus.2019.12.012>
- [58] N. Burkart and M. F. Huber, “A survey on the explainability of supervised machine learning,” *Journal of Artificial Intelligence Research*, vol. 70, pp. 245–317, Jan. 2021. [Online]. Available: <https://doi.org/10.1613/jair.1.12228>
- [59] C. J. Cai, S. Winter, D. Steiner, L. Wilcox, and M. Terry, ““hello AI”: Uncovering the onboarding needs of medical practitioners for human-AI collaborative decision-making,” *Proceedings of the ACM on Human-Computer Interaction*, vol. 3, no. CSCW, pp. 1–24, Nov. 2019. [Online]. Available: <https://doi.org/10.1145/3359206>
- [60] M. Du, N. Liu, and X. Hu, “Techniques for interpretable machine learning,” *Communications of the ACM*, vol. 63, no. 1, pp. 68–77, Dec. 2019. [Online]. Available: <https://doi.org/10.1145/3359786>
- [61] K. Zhang, P. Xu, and J. Zhang, “Explainable AI in deep reinforcement learning models: A SHAP method applied in power system emergency control,” in *2020 IEEE 4th Conference on Energy Internet and Energy System Integration (EI2)*. IEEE, Oct. 2020. [Online]. Available: <https://doi.org/10.1109/ei250167.2020.9347147>
- [62] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.

- [63] F. D. Martino and F. Delmastro, "Explainable AI for clinical and remote health applications: a survey on tabular and time series data," *Artificial Intelligence Review*, vol. 56, no. 6, pp. 5261–5315, Oct. 2022. [Online]. Available: <https://doi.org/10.1007/s10462-022-10304-3>
- [64] S. Chaudhari, V. Mithal, G. Polatkan, and R. Ramanath, "An attentive survey of attention models," *ACM Transactions on Intelligent Systems and Technology*, vol. 12, no. 5, pp. 1–32, Oct. 2021. [Online]. Available: <https://doi.org/10.1145/3465055>
- [65] A. Galassi, M. Lippi, and P. Torroni, "Attention in natural language processing," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 10, pp. 4291–4308, Oct. 2021. [Online]. Available: <https://doi.org/10.1109/tnnls.2020.3019893>
- [66] S. Lu, M. Liu, L. Yin, Z. Yin, X. Liu, and W. Zheng, "The multi-modal fusion in visual question answering: a review of attention mechanisms," *PeerJ Computer Science*, vol. 9, p. e1400, May 2023. [Online]. Available: <https://doi.org/10.7717/peerj-cs.1400>
- [67] H. Alharthi, "Healthcare predictive analytics: An overview with a focus on Saudi Arabia," *Journal of Infection and Public Health*, vol. 11, no. 6, pp. 749–756, Nov. 2018. [Online]. Available: <https://doi.org/10.1016/j.jiph.2018.02.005>
- [68] C. Cheng, B. Zhou, G. Ma, D. Wu, and Y. Yuan, "Wasserstein distance based deep adversarial transfer learning for intelligent fault diagnosis with unlabeled or insufficient labeled data," *Neurocomputing*, vol. 409, pp. 35–45, Oct. 2020. [Online]. Available: <https://doi.org/10.1016/j.neucom.2020.05.040>
- [69] J. Whittlestone, K. Arulkumaran, and M. Crosby, "The societal implications of deep reinforcement learning," *Journal of Artificial Intelligence Research*, vol. 70, Mar. 2021. [Online]. Available: <https://doi.org/10.1613/jair.1.12360>
- [70] S. J. Trenfield, A. Awad, L. E. McCoubrey, M. Elbadawi, A. Goyanes, S. Gaisford, and A. W. Basit, "Advancing pharmacy and healthcare with virtual digital technologies," *Advanced Drug Delivery Reviews*, vol. 182, p. 114098, Mar. 2022. [Online]. Available: <https://doi.org/10.1016/j.addr.2021.114098>
- [71] Z. Salahuddin, H. C. Woodruff, A. Chatterjee, and P. Lambin, "Transparency of deep neural networks for medical image analysis: A review of interpretability methods," *Computers in Biology and Medicine*, vol. 140, p. 105111, Jan. 2022. [Online]. Available: <https://doi.org/10.1016/j.compbiomed.2021.105111>
- [72] T. Capel and M. Brereton, "What is human-centered about human-centered AI? a map of the research landscape," in *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. ACM, Apr. 2023. [Online]. Available: <https://doi.org/10.1145/3544548.3580959>