

1 Manufacturing big data ecosystem: A 2 systematic literature review

3
4 Yesheng Cui^{a,*}, Sami Kara^a, Ka C. Chan^{a,b}

5 ^a *Sustainable Manufacturing and Life Cycle Engineering Research Group, School of*
6 *Mechanical and Manufacturing Engineering, The University of New South Wales,*
7 *Sydney, NSW 2052, Australia*

8 ^b *School of Management and Enterprise, Faculty of Business, Education, Law and*
9 *Arts, University of Southern Queensland, Springfield, QLD 4305, Australia*

10

11

12 **Abstract**

13 Advanced manufacturing is one of the core national strategies in the US (AMP),
14 Germany (Industry 4.0) and China (Made-in China 2025). The emergence of the
15 concept of Cyber Physical System (CPS) and big data imperatively enable
16 manufacturing to become smarter and more competitive among nations. Many
17 researchers have proposed new solutions with big data enabling tools for
18 manufacturing applications in three directions: product, production and business.
19 Big data has been a fast-changing research area with many new opportunities for
20 applications in manufacturing. This paper presents a systematic literature review
21 of the state-of-the-art of big data in manufacturing. Six key drivers of big data
22 applications in manufacturing have been identified. The key drivers are system
23 integration, data, prediction, sustainability, resource sharing and hardware.
24 Based on the requirements of manufacturing, nine essential components of big
25 data ecosystem are captured. They are data ingestion, storage, computing,
26 analytics, visualization, management, workflow, infrastructure and security.
27 Several research domains are identified that are driven by available capabilities
28 of big data ecosystem. Five future directions of big data applications in
29 manufacturing are presented from modelling and simulation to realtime big data
30 analytics and cybersecurity.

31

1 Introduction

Smart manufacturing is critical to national economies by providing jobs, improving innovation and advancing sustainability [1]. Several national strategies were initiated to boost their competitiveness of manufacturing, such as: ‘Industry 4.0’ in Germany [2], ‘Advanced Manufacturing Partnership (AMP)’ program in the United States [3], ‘Made in China 2025’ and so on. These initiatives provide massive potential to envision the future of manufacturing: Smart manufacturing, which is defined by National Institute of Standards and Technology (NIST) as a completely integrated, collaborative manufacturing system that respond in real time to meet changing demands and conditions in the factory, in the supply network and in customer needs[4].

Manufacturing industry uses a wide range of software and automation systems to increase efficiency and productivity from shop floors to enterprise layers such as CNC machines, Programmable Logic Controllers (PLC), Supervisory Control And Data Acquisition System (SCADA) [5], Manufacturing Executive System (MES) [6], product design and development (CAx: CAD, CAPP, CAM, CAE [7]), Product Lifecycle Management (PLM) [8], Enterprise Resource Planning system (ERP) [9], Operating and Maintenance (O&M) [10], Energy Management System (EMS) [11], Supply Chain Management (SCM) [12], Customer Relationship Management (CRM) [13] etc. However, smart manufacturing cannot be realised with traditional manufacturing software and technologies due to two main challenges. First, these systems and software cannot be fully integrated and collaborative since they are developed by multiple vendors using different interfaces or protocols. Second, manufacturers cannot perceive and respond to the real-time changes on time from the factory, supply chain and market since the traditional manufacturing software lack sensory data to notice the changes inside and outside the systems.

Digital thread and digital twin are two recent concepts proposed by integrating disparate systems over the product lifecycle [14] and building up the real-time relationship between the physical space and the cyberspace in manufacturing [15] respectively. Cyber-Physical systems (CPS) are physical and engineered systems, which are monitored, controlled, coordinated and integrated with computing and communicating core [16]. Internet of Things (IoT) is data-accessing and data-processing technologies on the cyberspace to perceive the real-time changes of physical space with sensory tools [17][18]. Digital thread and digital twin can be enabled by using IoT and CPS. With the practices of these new concepts and technologies, massive data will be generated from the systems, which have to go through the processes of data collection, storage, aggregation, analysis and exchange to provide timely information to manufacturers. Being empowered with cloud computing [19], data science [20] and Artificial Intelligence (AI) [21], big data focuses on addressing the big data issues among the processes, which traditional manufacturing tools cannot. Big data could be the enabler of the concepts of digital thread and digital twin.

Smart manufacturing gains actionable knowledge in real-time with the fusion of big data and manufacturing knowledge. As big data are collected and analysed to extract timely information, manufacturing industry may still not know which approach to use, and their impacts without the domain knowledge [22]. The actionable knowledge is created when manufacturers get timely information from big data and apply manufacturing knowledge in a specific application. Some examples found in the reviewed literature are: identifying the reasons of faults from the production process by analysing real-time big processing data and manufacturing knowledge [23], predicting maintenance intervals by utilising knowledge discovery and hundreds of machine data attributes [24], making real-time scheduling and cost-effective decisions in MES system with streaming shop floor data and existing manufacturing systems [25].

Similar to big data, manufacturing industry faces the same challenges associated with 5Vs of big data (Volume, Velocity, Variety, Veracity and Value) [26]. IDC reports that manufacturing has the largest

1 share of data (3584 Exabyte) in 2018 and will have 30% annual growth rate of data from 2018 to 2025
2 [27]. Of the reported data types, structured data, such as tabular data in relational databases or
3 spreadsheets, accounts for only 5% of all the data generated [23]; while the rest is made up of semi-
4 structured and unstructured data with formats JSON, XML, image, video, and audio, etc. Issues of
5 velocity, variety and veracity can be explained that the same type of data come from different devices
6 with various sampling frequencies, formats, precisions, which leads to inconsistent data and makes
7 challenging to extract the value-added insight to manufacturers. Amir et al. illustrate that the
8 limitations of the traditional methods (relational database management systems (RDBMS) and on-
9 premise software) could not handle big data [28]. As more manufacturing enterprises generate big
10 data, the issues of big data will become pressing.

11 Tapping into the capabilities of big data tools presents enormous opportunities for smart
12 manufacturing. A large number of big data tools are developed by the big players in the Internet
13 industry such as Google, Yahoo, Facebook for their own applications in search engines, social media,
14 and business analytics [29][30][31], such as Apache Hadoop [30], Apache Spark, Apache Flume,
15 Apache Flink, Apache Storm, NoSQL and NewSQL databases [32], Apache Hive, Apache Pig,
16 Apache Zookeeper etc. As these tools are enterprise-ready to use in manufacturing, many researchers
17 reported big data based solutions, which use a set of big data tools to address problems to enable
18 smart manufacturing. In 2012, a Hadoop-based sensor data management framework was proposed for
19 cloud manufacturing [33]. In 2014, Tao proposed the architecture of cloud manufacturing system as
20 well as the investigation of applying cloud computing and IoT technology in manufacturing [34]. In
21 2015, the Cloud-Based Design Manufacturing paradigm (CBDM) was proposed by comparing other
22 design methods [35]. In 2017, Nagorny etc. conclude that big manufacturing data and big data
23 analytics would provide a vast potential in smart manufacturing [26]. In 2017, Wu presented a fog
24 computing-based framework to monitor machine health in cyber-manufacturing. In 2018, Tao
25 proposed a data-driven conceptual framework to interoperate with ERP, MES, CRM, PLM systems in
26 manufacturing [36]. Big data tools complement and provide additional functionalities to address
27 manufacturing big data issues which could not have been solved by traditional approaches.

28 Three challenging issues have to be addressed in the research of big data in manufacturing. Firstly,
29 big data tools from Internet industry do not consider the differences between Internet and
30 manufacturing. Most manufacturing data is standardized which is supported by various industrial
31 vendors and associations such as manufacturers of CNC machines, meters and sensors, controllers
32 and software companies. The manufacturers use different hardware interfaces, communication
33 protocols, manufacturing machine readable languages, and semantical definitions. Whereas, most
34 data in the internet is based on natural languages and easier to be exchanged without the difficulties
35 associated with multiple interfaces and protocols. The differences between the two industries have not
36 been taken into account in developing big data tools. Secondly, big data tools are massive, diverse,
37 with many overlap functions. It is challenging to design big data based solution by selecting suitable
38 big data tools. However, designing big data solutions not only depends on big data tools but are
39 closely associated with the specific manufacturing applications and scenarios. This paper categorizes
40 similar types of big data tools and identifies the differences as the preparation to achieve this purpose.
41 Thirdly, many manufacturing systems have dedicated aims, complicated and sophisticated functions,
42 and are closely specialised with application scenarios. Social media uses big data tools to collect and
43 store time series data from billions of customers who follow each other online; and to report trending
44 events. For manufacturing, time series data from multiple data sources can be collected, integrated,
45 and analysed to explain the states of manufacturing entities. For example, the sensor data collected
46 from a CNC machine reflects the state of the machine and can be used to develop simulation models
47 or prediction models for preventive maintenance. As many solutions are proposed, a systematic
48 literature review is required to identify the data issues in manufacturing, capabilities of big data tools,
49 essential components to design big data based solutions in manufacturing and the potential research
50 directions of big data in manufacturing.

51 The rest of this paper is presented as follows: Section 2 presents the methodology to systematically
52 review the state of the art of big data research in manufacturing; Section 3 presents the outcomes of
53 the systematic review; Section 4 discusses several critical issues of big data ecosystem in

1 manufacturing including critical drivers, system requirements, essential components, research
2 innovation and future directions; and finally, Section 5 presents the conclusion of this systematic
3 review.

4

Nomenclature

BDA	Big Data Analytic	NoSQL	Not Only Structured Query Language
BI	Business Intelligence	NIST	National Institute of Standards and Technology, USA
CAD	Computer-Aided Design	OLAP	On-Line Analytic Processing
CAPP	Computer-Aided Process Planning	OLTP	On-Line Transaction Processing
CAM	Computer-Aided Manufacturing	O & M	Operation and Maintenance
CAE	Computer-Aided Engineering	OPC-UA	OPC Unified Architecture
CNC	Computer Numerical Control	OWL	Web Ontology Language
CRM	Customer Relationship Management	PDM	Product Data Management
DSS	Decision Support System	PLM	Product Lifecycle Management
EOL	End-Of-Life	QMS	Quality Management System
ERP	Enterprise Resource Planning	RDF	Resource Description Framework
HTML	Hypertext Markup Language	RDBMS	Relational Database Management System
IIoT	Industry Internet Of Things	SCADA	Supervisory Control And Data Acquisition
JSON	JavaScript Object Notation	SCM	Supply Chain Management
KM	Knowledge Management	STEP	Standard for Exchange of Product data
MES	Manufacturing Execution System	STEP-NC	STEP for Numerical Control
MOM	Manufacturing Operations Management	XML	Extensible Markup Language
MES	Manufacturing Executions System		

5

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

2 Methodology

This paper presents a systematic literature review (SLR) on the current state of research associated with big data technologies in manufacturing [37]. To apply big data technologies in manufacturing successfully, it is essential to systematically review the literature of big data technologies in manufacturing from the following three perspectives: manufacturing data, big data technologies and data applications in manufacturing.

Firstly, manufacturing data is the foundation to conduct data-driven manufacturing. It is impossible to propose one big data based solution to fit all manufacturing circumstances since different applications have different data issues (data types, data formats and data sources) and require specific tools to address. Therefore, systematically analysing manufacturing data could provide a useful guideline to select appropriate big data enabling technologies.

Secondly, the big data tools in the big data ecosystem have to be identified the similarities and differences. The 5Vs characters of big data are widely recognised as challenges, such as volume (TB/PB level of data size), velocity (ingesting or processing big data in streams or batches, in real time or non-real time), variety (dealing with complex big data formats, schemas, semantic models and information), value (analysing data to deliver added-value to some events), and veracity (validate data consistency and trustworthy) [38]. In general, these big data technologies are intended to address some Vs of big data. Hence, their capabilities need to be effectively classified and analysed to know which Vs are addressed.

Thirdly, gaps of data applications in manufacturing could be identified by systematically reviewing the capabilities of the traditional manufacturing systems and big data analysis. Since much traditional manufacturing software has been widely used in enterprises, big data could integrate and collaborate the software and systems as well as providing timely information. The massive amount of data generated by these applications can be fed back to the big data ecosystems for analytics and innovative applications such as prediction, optimization, monitoring, simulation and visualisation, etc. Therefore, these application gaps would be the future research directions in academia and the demands in the industry.

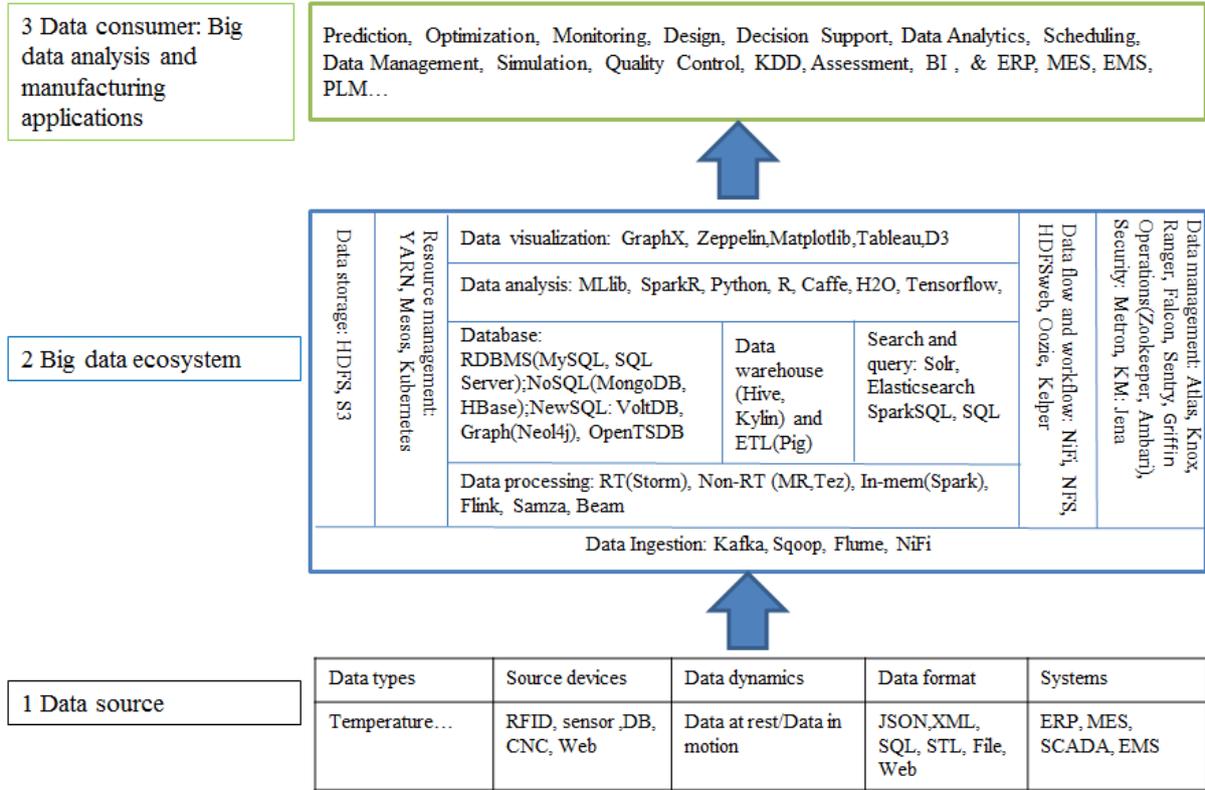
In summary, knowing the data requirements of manufacturing applications, understanding the capabilities of big data tools, and identifying the gaps will help define future research directions and generate new ideas for innovative applications. This systematic literature review presents a holistic overview of big data in manufacturing to study the possible use cases for manufacturing. As shown in Figure 1, the conceptual framework of this systematic literature review includes three layers: data source, big data ecosystem, and data consumers.

The first layer, at the bottom of Figure 1, is the data source, which consists of five aspects:

1. Data types refer to the meaning of data such as temperature, humidity from the physical space, and log, email, operational data from cyberspace;
2. Source devices to collect data sources, which include sensors, controllers, actuators, software systems. The data type and source devices have a close relationship with the first four characteristics of big data (Volume, Velocity, Variety, Veracity) [39]. For example, in order to know the temperature of a production line, a temperature sensor is selected with the determined sampling rate (speeds of data generation), the sizes of data accumulating by time, formats and quality of data from the sensor.
3. Data dynamics describe the states of data. Data-at-rest refers to the inactive data stored in spreadsheets, databases and data warehouses; while data-in-motion refers to the active data generated by sensors, equipment or machines, and fed into the big data ecosystem in real-time.
4. Data formats are structures of the data. Data is exchangeable among various systems with consistent data formats and languages.

The system aspect of data source refers to the system where data is originated. There is a diverse range of manufacturing systems used in different applications such as product design, manufacturing

1 pyramid, product lifecycle management, supply chain management, logistics. The second layer, the
 2 middle layer shown in Figure 1, is the big data ecosystem comprising all the big data software. This
 3 layer plays the role of connecting the data sources from the layer below and the big data analytics
 4 applications at the layer above.
 5



6
7 Figure 1. Conceptual framework of systematic literature review

8 The big data ecosystem is a set of complex and interrelated components to process and analyse big
 9 data [40]. Also, the ecosystem needs to store data from various data sources for data integration and
 10 analytics as well as other applications. Therefore, data storage layer in the ecosystem includes
 11 database and file system technologies to store big data. The ecosystem consists of the following
 12 components and tools:

- 13 • Data collection and ingestion: log data collection (Flume), bulk data collection from a
 14 relational database (Sqoop), distributed messaging system (Kafka), dataflow (NiFi);
- 15 • Computing engines: batch processing (MapReduce), iterative/near real-time processing
 16 (Spark, Flink), real-time processing/streaming (Storm, Flink) [41];
- 17 • Database: Relational database (RDBMS) has a standard schema but without scalable
 18 capabilities (MySQL, Oracle DB, SQL server, ProgresSQL); NoSQL database does not have
 19 a standard schema and has scalable capability (four types NoSQL: Column-based: HBase;
 20 Document-based: MongoDB; Key-value-based: Redis; Graph-based: Neo4j); NewSQL
 21 database is scalable relational database (VoltDB) [42][43], search engine (Solr,
 22 Elasticsearch) ;
- 23 • Data analysis (BDA): Machine Learning (MLlib, Caffe, Tensorflow, Python), statistic
 24 (SparkR, R), OLAP,
- 25 • Data visualization (Zeppelin, Matplotlib, Tableau, D3 [44], GraphX);
- 26 • Workflow which is a scheduler of the jobs of various big data tools and dataflow which
 27 manages data transfer and data transformation among different big data tools: Oozie, Kepler,
 28 Apache NiFi;
- 29 • Data management and KM: Apache Falcon, Apache Atlas, Apache Sentry, Apache Hive,
 30 Operation (Zookeeper, Ambari), Apache Griffin, Apache Ranger, Apache Jena;

- Big data infrastructure (BDI): computing resources (general purpose computing and HPC), cluster management (YARN, Mesos) [37], network communication (Software-Defined Networks (SDN) [45], InfiniBand [46], 5G [47]) etc.;
- Big data security: Apache Metron [44], Apache Knox;

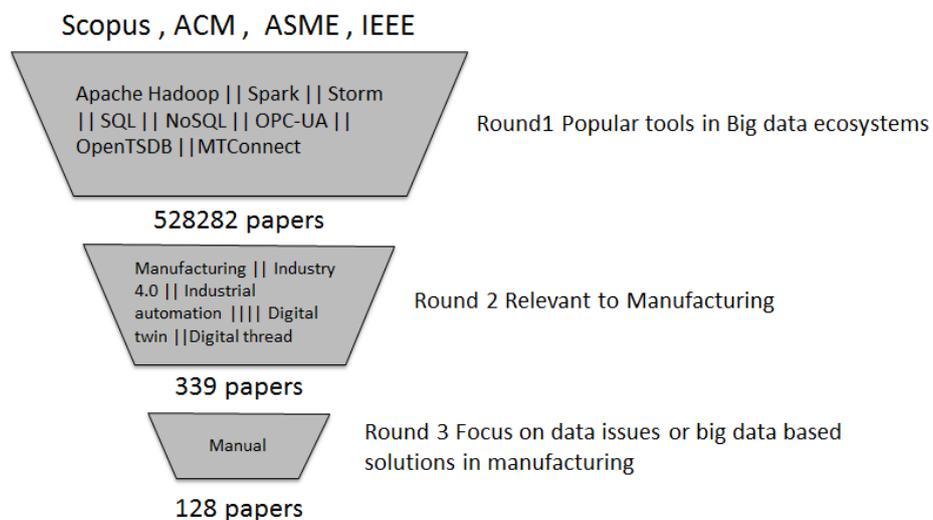
Although Hadoop is a big part of the big data ecosystem with many big data tools [48], it lacks functions such as data flow, data management and security [41].

Finally, the top layer represents the way data is used and data users. It includes the applications of big data analytics [26] and manufacturing applications [36]. The databases of traditional manufacturing software maybe the data sources of the ecosystem as well, such as SCADA or EPR.

2.1 Literature identification

Pertinent articles are identified within the scope of manufacturing data, big data technologies and big data based solution in manufacturing. Figure 2 illustrates the research method in this literature review. First, four citation databases are chosen due to their comprehensive coverage and high relevance to the scope: Scopus, IEEE Xplore, ASME Digital Collection, and ACM Digital Library. Second, several big data technologies and popular manufacturing data collection tools are selected as keywords. Hadoop is chosen since it is the earliest big data technology which is well studied and used. The underlying technologies of big data are computing and storage. There are a few big data computing engines. Three common ones are selected: Spark, Storm and Flink. For storage, because there are hundreds of available databases, it is not appropriate to limit to specific databases. All the databases could be categorized into three types: SQL, NoSQL and NewSQL. We select SQL and NoSQL because NewSQL could be recognized as SQL database with better features than traditional SQL. The features of NewSQL will be discussed in the following Sections. Time series data widely exist in manufacturing such as machines, sensors, controllers. Time-Series database (OpenTSDB) is selected since it seems more suitable for manufacturing data. Two widely adopted manufacturing data collection tools are selected: OPC-UA [49] and MTConnect[14].

Third, to further focus on the nature of this research paper, the papers were filtered by using the abstract “Manufacturing”, “Industry 4.0”, “Industrial automation”, “Smart manufacturing”, “Digital twin” and “Digital thread”. Four, manual review is implemented to select the papers about manufacturing data issues, or big data based solutions and applications in manufacturing. For example: some papers of other industrial sectors are found since they merely mentioned manufacturing in the abstract, such as Oil and Gas, Healthcare, Energy and Agriculture. Hence, the 339 articles are reviewed and 128 relevant articles are selected. The search strings of four databases are listed in Table 1 in Appendices.



34

35

Figure 2. Process of Literature review methodology

3 Results

3.1 Manufacturing systems

In 2016, NIST reported three dimensions of concerns in smart manufacturing systems (SMS): product, production and business. Many traditional manufacturing systems and software can be categorized into one of the dimensions [49] (Figure 3). Business dimension is presented in the upper rectangle block with dash lines, which includes suppliers, customers and manufacturing enterprises (SCM, CRM, BI, asset management). Product dimension is located at the bottom rectangle block with solid lines. It includes objects and activities from product design to end-of-life of the product (CAx: CAD, CAM, CAPP and CAE, PLM). Production dimension is the triangle block, which includes an entire production system (ERP, MOM/MES, SCADA/DCS/HMI, O&M, Safety, quality management). Industry Internet of Things (IIoT) and RFID technology (grey blue) are widely used in manufacturing, logistics [50], in-use and product End-of-Life (EOL) [51].

Afterwards, all the reviewed articles are classified into four categories: Product, Production, Business and ICT (Information Communication Technology). The first three categories focus on engineering functions and business, ICT architecture underpins all three dimensions to provide the ICT infrastructure and digitalization to manufacturing, which includes several topics: CPS, Cloud manufacturing (CM), ICT, Data analytics/Data management(DM), KM. Table 2 in Appendices illustrates the distribution of reviewed articles by these four categories.

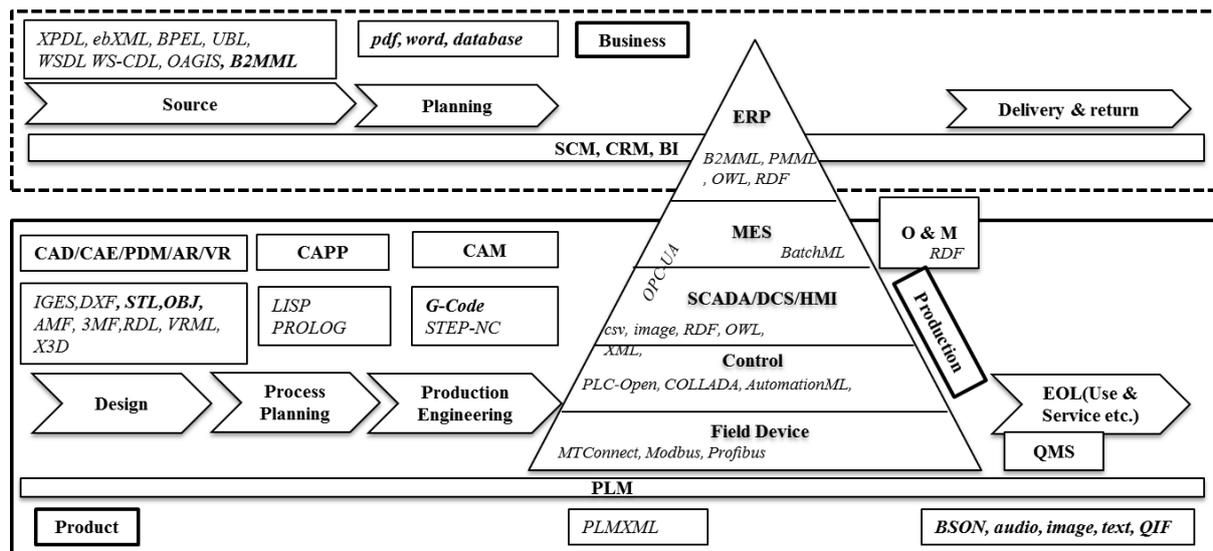


Figure 3. Smart Manufacturing Systems and various data formats

3.2 Data source

3.2.1 Data format

Based on the three dimensions of SMS in NIST report, the standards of data formats, computer languages are listed with the italic style in every manufacturing system in Figure 3. The bold black text represents some of the data formats found in the reviewed articles, whereas the black text is not found in the review but mentioned in the NIST report [49].

Table 1 demonstrates the complete data formats found from the reviewed articles and NIST report. To discuss the formats conveniently in the following chapter, all the data formats are categorized into three groups:

- Structured data: data that is presented in tables and can be stored in a relational database;

- Semi-structured data: data that has a self-described structure and is not presented in tables, such as XML, JSON, HTML [52];
- Unstructured data: data that does not have a self-described structure, such as document, image, audio, video, text and e-mail.

Figure 3 and Table 1 show that one challenging issue to realize that smart manufacturing requires different data formats from various manufacturing systems. In order to make these systems collaborative and integrated, the transformation of these data formats is an essential function to the manufacturing big data solutions. It also illustrates that some data formats in the specific manufacturing systems are missing in the proposed big data solutions. It requires solutions to fill the gap to realize data exchange among the systems. Big data tools can address the variety issue of manufacturing data such as NiFi.

Categories	Systems	Structured	Semi-structured	Unstructured
Product	CAD/CAE/CAPP/CAM	[53]	XML [54], G-code [55], STL [56], (aml, obj, UML, AutomationML) [57], IGES, DXF, AMF, RDL	○
	PLM	[58]	(XML, B2MML) [59], PMML[60], (RDF, SPARQL, STEP, QIF, STL) [61], PLMXML [60]	pdf [61]
Production	ERP	[62]	(XML, HTML, SCUFL) [63], PMML	e-mail [64]
	MOM/MES	[65][66]	JSON [67], RDF [68], AutomationML [69], BatchML	○
	SCADA/DCS/HMI	[70]	(RDF, OWL, XML) [71], SPDML [72]	Image [73]
	IIoT/CNC/Robot	○	(XML, UML, AutomationML) [69], OPC-UA, PLC Open, COLLADA) [74], BSON [75], JSON [76]	image [77][78],
	O&M	[79][80]	cad [80], RDF [81], XML [82]	Image [83], video [82]
	QMS	[84]	BSON [85], (XML, QIF) [86]	image [87], (audio, document) [88]
	Safety	[89]	○	○
Business	SCM//BI/AM	[90]	(XML, JSON, RDF) [91], (WSDL, EPL) [92], XPDL, ebXML, BPEL, UBL, WS-CDL, OAGIS,	document [93]
ICT	CPS/CM/ICT	[94][95]	(JSON, STEP, JT Open) [96], HTML [97], (AutomationML, PLCopen) [98], (XML, XSD, RDF) [99], (UML, SysML, STEP, B2MML) [5], EDDL [100], [101]	e-mail [64]
	Data analytics/DM	[102]	BSON [102], JSON [103], Parquet [104]	(image, video, document) [105]
	KM	○	(OWL, UML, RDF, SWRL) [23], (JSON, PMML, AMPL) [106], (STEP-NC, G-code, XML, DMIS, QIF) [107]	○

○: Not found

Text: Not found in reviewed articles (mentioned in the NIST report)

Table 1. Data formats in reviewed articles and NIST report

3.2.2 Data issues

Data issues are fundamental challenges to smart manufacturing, which extract actionable information from good quality of data. In order to prepare the suitable data for smart applications, amount of cost

1 and time is consumed to address the data issues. For example, data scientists spend over 90% of their
2 time on data preparation before analysing data for innovative tasks such as machine learning, AI [108].
3 Enterprises spend billions of dollars on their data warehousing systems, which can only use well and
4 pre-defined methods (ETL) to process product structure data and produce business reports in non-real
5 time. Therefore, understanding and addressing data issues is critical to design big data-based solutions.
6 From Table 2, we have identified 11 common issues of manufacturing data. Data issues are generally
7 related to the Big data 5Vs features[38]: volumes and variety (large scale); velocity (inconsistent
8 sampling frequencies or timestamp, batching and streaming data); veracity (missing value, imbalance,
9 data outlier, noisy, drifting, asynchronization, data correlation); variety, veracity and value (data
10 model, data format exchange and data integration). The researchers of these kind of literature address
11 these issues without using big data tools. Traditional software cannot address these data issues if the
12 data is generated in large scale systems with large number of devices. Some big data tools are
13 discussed to provide potential solutions in Sections 4.2.7 and 4.4.2.

Issues	5Vs	Description	References
Large scale	Volume, Variety	Large volume dataset with massive features.	[109]
Inconsistent sampling frequencies or timestamp	Velocity	Sensors use various sampling frequency and timestamps; Unnecessary high sampling frequency affect the real time performance, it is related with "Smart data" topic.	[110][111][112] [113]
Batching data and streaming data	Velocity	Data modification is sensitive to time or not, some literature also uses the terms: data-at-rest and data-in-motion [13];	[114]
Missing value	Veracity	Record is empty when equipment is an anomaly; Product passes some portion of machines; Sensor is an anomaly or losing communication with sensor; Too costly to capture data by installing sensor or building models.	[115] [109] [110] [113]
Imbalance	Veracity	Small probability data in a very large dataset with most of normal data.	[109]
Data outlier	Veracity	Data is out of range of measurement device.	[110][84]
Noise or an anomaly data	Veracity	Data is out of the similar clustering dataset; Noise data is possibly generated by replacing the missing data.	[115] [109]
Drifting data	Veracity	Process drift caused by vulnerability to external environment; Sensor drift caused by modification in measuring device or calibration.	[110] [113]
Asynchronization	Veracity	Several data producers (sensor and machine) use non-central time server in manufacturing enterprise.	[113]
Data correlation	Veracity	Nature and structure of data caused by redundant sensor arrangement; Data correlation is to improve data quality with process variables. Inconsistent simulation and collection of data from CNC machine.	[110] [116] [107]
Data model, data format exchange and data integration	Variety, Veracity, Value	Merge data from multiple data sources into a single view; Exchange different data formats from various systems: OPC-UA and IoT, OPC-UA and AutomationML, MTConnect with QIF, MTConnect and IEEE 1451 wired smart transducer Integrate data models to explain the relationships of data; Data information must be available for information sharing. OWL is the enabling technology on resource description.	[115] [117],[98] [86], [118], [113] [99][118]

14 Table 2. Data issues in manufacturing

15

3.3 Big data ecosystem

Through this SLR, new big-data based research innovation could be identified by closely tracking the attention of various big data tools in manufacturing and other research domains. Figure 4 shows the distributions of big data tools over the years from the literature, including manufacturing engineering and others. Fig 4(a) shows the numbers of the term Hadoop occurred in all literature increased from 2008, reached its peak in 2016 and slightly decreased in 2017. The number of Spark articles had been higher than Hadoop since 2016. There are four factors to compare Spark and Hadoop: volume, velocity, fault-tolerant and data analysis. MapReduce and Spark Streaming are the computation engines of Hadoop and Spark, respectively. MapReduce executes batch processing by reading and writing data on disk multiple times. Spark Streaming executes micro-batch processing in memory. It results in the differences between Hadoop and Spark because disk can persistently store a larger volume of data with slower velocity than memory, which temporarily stores limited volume of data. By comparing data analysis, Spark has built-in tools (MLlib) and support third-party tools (Mahout, H2O), whereas Hadoop is only supported by the third-party tool: Mahout [41]. Spark supports iterative computation with GraphX, which is the graph processing engine. Therefore, Hadoop is suitable for the applications which need planned extraction of non-real time and critical information from a larger volume of data and guarantee without loss of data, such as ERP, Production planning. Spark can be used to provide near real-time monitoring and analytics by processing streaming data such as monitoring process and product quality, MES, SCADA, predictive maintenance. The number of Storm articles has kept increasing since 2012; however, its total number is smaller than Hadoop and Spark since it only executes streaming with limited data analysis supported by SAMOA, which is a version of Mahout [119]. Fig 4(b) presents the growing number of literatures about these three tools in the manufacturing research literature since their first inception.

From Fig 4(c) and Fig 4(d), NoSQL database, as the highest frequently mentioned database, is compared with NewSQL and time series databases (OpenTSDB) over the last eight years. Fig 4(e) also demonstrates that Kafka and OPC-UA are closely related to manufacturing as seen by their increasing patterns. Because OPC-UA protocol provides open connection to monitoring and automation systems as well as communications between MES and SCADA systems, these systems are prevalent in manufacturing [120]. Kafka is also used widely in applications at the shop floor level of manufacturing such as processing, machine and sensor data, due to its message streaming capability [121]. Apache NiFi is entirely new to manufacturing. Only four articles were found in the industry in general and no article in manufacturing industry [122][123][124].

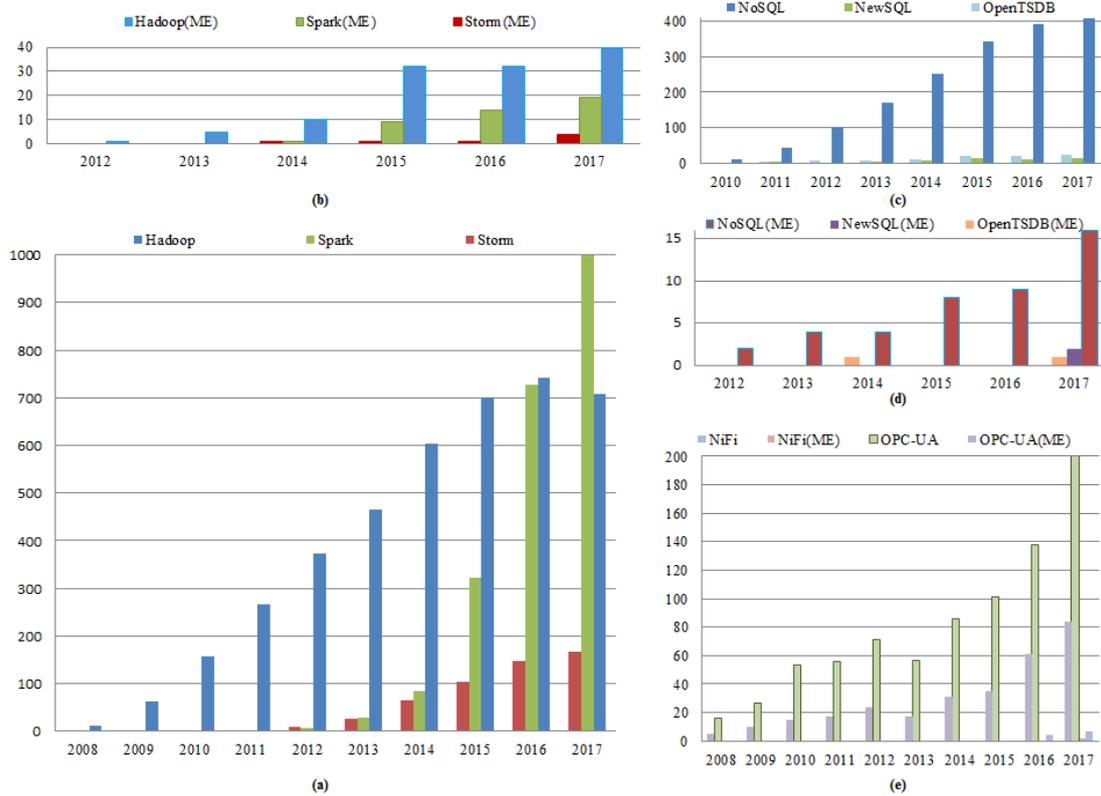


Figure 4. Chronological distribution of big data tools

3.4 Applications of big data in manufacturing

Data applications are essential to realize Smart manufacturing. Identified scopes of data applications could provide a clear guideline to design manufacturing big data platforms. Among the reviewed 128 articles, 78 articles are big-data based applications, which are categorized into 17 applications. The percentages of the applications are presented in Figure 5. Monitoring (25%), prediction (23.8%), ICT framework (11.9%) and data analytics (9.5%) are the four most frequently used big-data applications in manufacturing. Through this statistic, it can be shown that the researches of big-data based solutions focus on monitoring, prediction, data analytics and propose ICT solutions in manufacturing in Table 3.

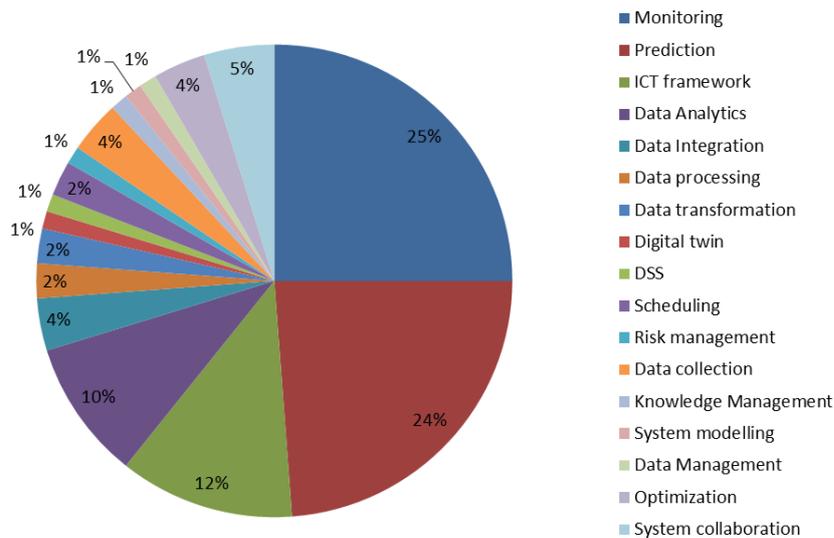


Figure 5. Percentage allocation of Big data applications (Based on our literature investigation)

Field	Systems	Big data framework	Big data computing and storage	Big data analytics
Product	CAD/CAE/CAPP/CAM		[57], [56], [125],	Regressions [53], ANN [53],
	PLM	[58],		
Production	ERP			MP [63], regression [62], K-means [126]
	MOM/MES	[69],[116],	[66],	regression [127], Distance, Regression, Self-organizing map, principal component analysis [128],
	SCADA/DCS/HMI	[129], [71],	[104],[130],[131],[104],[132],[112]	Classification [133], OPL [106], KM [134], GA [73],
	O&M	[39],[135]	[136],[81],[82],[80]	logistic regression, naïve Bayes, and a decision tree [109], regression [137], LSM [138], SVM [83], Anomaly detection [139], DTW [140], RF [141], K-means, Markov [142], KD [24],
	QMS			SPC [87], ANN, decision tree, random forest, SVM [143],[84]
	Safety		[144],	graph analytics [89],
Business	SCM/CRM/BI/AM	[93], [90],	[145]	
ICT architecture	IIoT/CPS/CM/ICT	[64],[146],[76],[147]	[68],[148],[95],[149],[150],	
	Data analytics/DM	[151],[33]	[152],[94],	Stream data analytics [103],
	KM			KM [106], Semantic data integration [23],[153],

1 Table 3. Allocations of proposed solutions of reviewed literatures (Based on our literature
2 investigation)

3 4 Discussion

4 Four fundamental questions about the relationship between big data ecosystem and smart
5 manufacturing are six drivers and requirements for big data application in manufacturing, seven
6 essential components of the big data ecosystem, harnessing big data capabilities for research
7 innovation in manufacturing, and future directions of big data application in manufacturing. They
8 illustrate the driving factors of big data applications in manufacturing. These are summarized in Table
9 3 and Table 4 for ease of reading to interested readers.

Topics	Sub-topics	Application systems	Enabling tools	References
6 drivers for big data in smart manufacturing	System integration	Product design, AM, ERP, MES, BI, SCM, PLM,	Kepler, Hadoop, OPC-UA, RESTful API	[49][57][56][125][63][62][90][153][97][55][128][154][155]
	Data	Predictive maintenance, KM, Production planning, Safety, Anomaly detection, industrial process control, model prediction, QC, shop floor scheduling,	Cassandra, MongoDB, Blueflood, OpenTSDB, DalmatinerDB and InfluxDB, Storm, Spark, Flink, Apache Hive,	[121][62][73][156][157][82][138][158][24][89][94][153][159][75][160][88]
	Prediction	3D printing, product performance, production planning, energy consumption, MES, QC, SCM	Random Forest, Bayesian Network, statistic	[39][109][63][62][93][136][137][138][141][24][87][143][90][55][161][85][88][25][162][79]
	Sustainability	PLM, Maintenance		[163][59]

	Resource sharing and networking	SCM, ERP, Data integration	Public cloud, Private cloud, Hybrid cloud, hypervisor, container,	[19][63][62][144][95][152][97][101][164]
	Low cost hardware	SCM, Industrial automation	RFID, IoT, Robotic	[65][165][166][167]
9 essential components of big data ecosystem	Data ingestion	ERP, MES, SCADA, O&M, QC, PLM, Data management	Sqoop, Flume, Kafka,	[39][116][58][71][138][139][140][151][85][79][65]
	Storage	ERP, SCM, SCADA, OLAP, OLTP	Redis, HBase, Cassandra, MongoDB, Neo4j, HDFS, VoltDB, Clustrix, NuoDB	[30][32][43][168][169]
	Computation	ERP, SCM, PLM, MES, SCADA, O&M, QC, IoT	MapReduce, Spark, Flink, Storm,	[170][171]
	Analytics	DSS, CRM, machine vision, QC, O&M	MLlib, Scikit-Learn, CNTK, Caffe, Kylin, CaffeOnSpark, Hive, CaffeOnSpark	[109][73][82][87][25][105][172][173]
	Visualization	MES, SCADA, O&M, QC, Security,	Zeppelin, Tableau, D3.js, Matplotlib, QlikView	[126][76]
	Workflow and dataflow	Business	Oozie, Kepler, InfoSphere, Wings/Pegasus, NiFi	[116][123][63][172]
	Data management	ICT, KM, SCADA, SCM,	Apache Falcon, Apache Atlas, Apache Sentry, Apache Griffin, Jena	[121][23][174][175][176][177][178][179]
	Infrastructure and deployment model	ICT	HPC(AWS),	[180][181]
	Cybersecurity	SCADA, ICT		[182][183]
5 future directions	Modelling and simulation	VR/AR, PHM, PLM, CAx	Cloud computing, Quantum computing	[184][185]
	Connectivity and interoperability	File formats of PLM (PLM XML) and 3D printing (AMF, 3MF)	NiFi,	[186]
	Standardized big data platform design			[5]
	Real time big data analytics	SCADA, MES, Data warehousing	Spark, Storm, Flink, Beam, Spark R, MLlib, GraphX, SparkSQL	[26][187]
	Cybersecurity	SCM, SCADA, Security, Safety	Apache Metron, Apache Ranger, Apache Knox	[49][188]

Table 4. Summary of Discussion section

4.1 Question 1: What are the drivers and requirements for big data applications in smart manufacturing?

Identifying drivers for big data applications is essential to implement feasible smart manufacturing initiatives. Kusiak discusses the future developments in manufacturing and identified six drivers for smart manufacturing theoretically [4], which are manufacturing technology and processes, material, data, predictive engineering, sustainability, resource sharing and networking. Through reviewing the proposed big-data based solutions, six drivers are identified: data, prediction, sustainability, resource sharing, system integration, and low-cost hardware. Four common drivers from Kusiak's six drivers are verified with big-data based solutions.

4.1.1 Driver 1: System integration

System integration with big data technologies is a crucial enabler of smart manufacturing to integrate and cooperate manufacturing systems to timely adapt dynamic demands from production and supply chain [49]. Integration of production systems demonstrates the significant improvement in production efficiency and productivity since the 1980s [4]. In the context of smart manufacturing, it needs to

1 further expand the scope of system integration from production to product and business domains.
2 Figure 3 shows that various manufacturing systems use different networks and protocols, which are
3 challenging due to the necessity of implementing data and information exchange among these systems.
4 Big data technologies can integrate these independent systems by using “cloud” as a common place to
5 collect data, extract and exchange the required information on the cloud. With the fusion of IoT, BDA
6 not only can integrate manufacturing systems but integrating physical and cyber worlds closer [36].
7 From the system integration in Table 4, many benefits of system integration are demonstrated with the
8 proposed big-data based solutions.

9 Integrating product design and additive manufacturing with big data provides many benefits. With the
10 integration of product design and additive manufacturing, product costs can be estimated by analysing
11 more features of product model with DBA than the traditional approach [55]; using Hadoop clusters
12 demonstrates faster velocity of converting a huge 3D model to G-code for 3D printer than traditional
13 methods [56]. Spark and Cassandra demonstrate capabilities of computing and storing a large volume
14 of streaming data from the application of real time monitoring 3D printing, which traditional
15 manufacturing systems cannot offer [125].

16 As to production, big-data based solutions drive the integration of ERP and MES. The big data
17 scientific workflow management system (Kepler) demonstrates the efficient scheduling capability for
18 smart manufacturing [63]. Another paper concludes that the critical cycle time can be predicted with
19 Hadoop for production planning [62].

20 In the business domain, several researchers focus on designing big data architecture to integrate BI,
21 SCM, ERP, MES and PLM systems. Although a business intelligence architecture is proposed to
22 advance the integration of business information from various existing systems such as ERP, CAx,
23 SCM, PDM/PLM, the specific technical framework is not provided [90]. It demonstrates that it is
24 necessary to have a capability to store data and information with various data formats, structures and
25 models to integrate various systems. Another big-data based solution is proposed to integrate supply
26 chain and production planning by retrieving and integrating SCM, ERP and MES data [97]. However,
27 the framework focuses on designing business functionality without providing support to process the
28 collected big data. A cloud manufacturing collaboration system is proposed to achieve better
29 performance and functionality on production, resource planning with the Hadoop ecosystem. It is
30 concluded that information integration from the web and other data sources is a critical issue to
31 implement system collaboration [153]. Processing and analysing data from MES and SCADA with
32 MapReduce and BDA can detect anomaly minutes beforehand in a large-scale production [128]. In
33 the industry, Bosch presented a conceptual, analytic platform with a data integration method to
34 integrate various data sources [154].

35 System integration with BDA requires standard data format and standard interfaces. New standards of
36 file formats need to be developed in order to fill the deficiencies of the existing standards to transfer
37 consistent content [49]. Additive Manufacturing File (AMF) is a new XML based standard format to
38 replace STL by providing many new features in additive manufacturing such as materials, material
39 properties, colours. [49]. Because RDBMS is not suitable to store data with flexible data models
40 (unstructured or semi-structured data), big data technologies are thus mainly used to process
41 unstructured data such as Hadoop, NoSQL databases. The standard interface is of importance to
42 seamlessly integrate systems in manufacturing. Because OPC-UA technology provides dedicated
43 interfaces to production equipment such as PLC, it received increasing attention in the industry [128].
44 RESTful API demonstrated to be more efficient to connect web data by comparing with the traditional
45 method SOAP [155].

46 **4.1.2 Driver 2: Data**

47 Timely comprehensive data with enabling big data tools is the key driver to smart manufacturing. The
48 data row in Table 4 presents many manufacturing applications could be implemented with more
49 comprehensive data and big data tools, which were challenging with traditional tools. Firstly, a large
50 volume of data collecting from various sources provides sufficient data for big data analytics. Data
51 from various industrial equipment, IoT devices, web and smartphones, is called data-in-motion [159],

1 which is continuously generated and ingested into the systems to provide real-time response from the
2 physical world, such as dynamic shop floor scheduling [75], predictive maintenance [141][136][24],
3 anomaly detection [81], diagnosis [82], prognosis [161] and systems collaboration [153]. As another
4 significant data source for big data ecosystem, databases in manufacturing is data-at-rest, which
5 represents static, historical data [159]. This data is mainly used to predict the long-term performance
6 in production planning [62], global manufacturing network design [94], critical event detection in
7 safety [89]. Both historical batching data and real-time streaming data are integrated to train models
8 and monitor real time condition information such as anomaly detection of machines' energy
9 consumption data [139].

10 Secondly, big data technologies make the management of manufacturing big data feasible.
11 Traditionally, RDBMS is mainly designed to store structured data with limited scalability. However,
12 NoSQL databases present better performance for handling semi-structured (JSON, XML) and
13 unstructured data (audio, video, and email) with unlimited scalability. For example, column NoSQL
14 database Cassandra was used to store event data of automation controller [131], document NoSQL
15 database MongoDB was used to store machine data [85]. Time-series databases (TSDB) begin to
16 receive increasing attention by providing dedicated applications for sensor data. A comprehensive
17 evaluation was implemented to several TSDBs: Blueflood, OpenTSDB, DalmatinerDB and InfluxDB
18 [82]. The collected data needs cleaning before usage in order to resolve issues such as noisy and
19 incorrect format, as shown in Table 2. Streaming (Flink, Storm), micro-batching (Spark) and batching
20 (MapReduce) data processing technologies provide the capabilities to clean and calculate big volume
21 of manufacturing data. The following big-data based solutions are proposed and implemented in
22 various manufacturing applications: complex event processing (CEP) with Storm[138], anomaly
23 detection with Flink [121], industrial process control with Spark [73], model prediction [137] and
24 quality control with MapReduce [88].

25 Lastly, past and new knowledge can emerge from the generated big data by harvesting big data
26 technologies. Knowledge of predictive maintenance can be extracted with an Apache Hive-based
27 platform [24]. Knowledge of intelligent applications of a smart factory is managed with Hadoop and
28 OWL technologies [134].

29 **4.1.3 Driver 3: Prediction**

30 Prediction enables manufacturing to change from reaction to prevention. Because of big data and
31 increasing applications of data analysis in manufacturing, it is feasible to predict the behaviours of
32 various manufacturing systems accurately.

33 Prediction attracts many researchers' attention to manufacturing. From Figure 5, the prediction is the
34 common BDA application in manufacturing. In the product domain, the costs of 3D printed products
35 can be predicted with the proposed big-data based solution and three machine learning algorithms
36 [55]. In the production domain, machine learning is used for the prediction of product performance
37 degradation [63], cycle times for production planning [62], energy consumption and KPI values on
38 MES system [25], and production efficiency [162]. A number of papers focus on the prediction of
39 product quality with different machine learning algorithms such as Bayesian Network [88] and
40 statistic analytics [87]. ANN is identified as the highest prediction accuracy by comparing with a
41 decision tree, random forest and support vector machine [143]. Some applications need a trade-off
42 between acquiring higher accuracy and shorter calculation time since it takes more time to calculate to
43 get higher accuracy, such as real-time quality control. In terms of the full consideration of shorter
44 calculation time and higher accuracy, Random forest is outperformed Naive Bayesian, Multi-Layer
45 Perceptron and Logistic Regression [85] [109]. Many papers publish the solutions of predictive
46 maintenance on machines by using machine learning algorithms to recommend scheduling of
47 proactive measures before outages occurred [137][138][24][39][79][161][136][141]. In the business
48 domain, some potential applications are found such as predicting user behaviour with using principal
49 component analysis and Hadoop [128], proactive inventories, location and throughput times on
50 logistics with a big data based platform [90].

4.1.4 Driver 4: Sustainability

Sustainability with big data technologies plays a vital role in smart manufacturing. Sustainable manufacturing considers the four factors: material, manufacturing processes, energy and pollutants. Big data technologies can provide a data-driven solution to analyse the big volume of data for these four factors. For example, product design could be guided by analysing the End-of-Life data of products; strategic decision making by analysing marketing, production and supply network data from CRM, ERP, MES and SCM systems; energy consumption and pollutant influence would be monitored through IoT sensors and RFID tags. Although several conceptual big-data based frameworks are proposed as shown in Table 4 [58][59], the performance of sustainability with big data technologies has not been evaluated. Further research is required to discuss the effectiveness of big data technologies on sustainability in manufacturing.

4.1.5 Driver 5: Resource sharing and networking

Manufacturing could benefit by sharing virtual and physical resources with the supply chain network. The issue of information silo results in losses of productivity and economy in manufacturing.

- **Sharing information**

Sharing information is beneficial to manufacturing systems of product, production and business. Helu et al. discuss that digital thread can improve product design and manufacturing processes by sharing data of product lifecycle systems [14]. For instance, some information on idle equipment among enterprises can be shared with big data tools in order to reduce holding costs such as machines, 3D printers, equipment, transport and warehouses [101]. Sharing information about SCM and PLM systems could meet the new business requirements such as find best supplier of a specific raw material [164]. Integrated data from various systems could identify potential production problems and improve work efficiency by selecting suitable maintenance time [152]. Collection and sharing data of supply chains with big data could be helpful to interactions among customers, manufacturers, and suppliers [97].

Addressing information silo is challenging because data cannot be easily shared among traditional manufacturing systems by different protocols. Some practical solutions are implemented, such as sharing data in various systems, which leads to the issue of data redundancy. Manufacturers also have to take massive maintenance effort on synchronizing the shared data of every system. Because the big data platform could be sat on the cloud, it has to establish only one connection to synchronize data of the platform and the system. Through big data tools, data is much easier collected from various systems to the data lake of the big data platform, synchronized, managed and shared[121].

- **Sharing BDA Infrastructure and software**

Sharing big-data infrastructure and software brings economic benefits to manufacturing. According to NIST definition, there are four deployment models of cloud computing: private cloud, community cloud, public cloud and hybrid cloud [19]. From reviewed papers, 18 solutions use private cloud while only two solutions use public cloud. The main reason for adopting a private cloud is to provide better privacy and security than public cloud [62][63][144]. However, there is a significant investment on the hardware of private cloud clusters. Wang et al. demonstrate that manufacturing can gain economic benefits from the public cloud with three aspects: pay-as-you-go service model, reducing maintenance fee on data centres and operating cost [189]. No security function was identified in the 18 proposed solutions, which means the security and privacy of the private cloud solutions do not outperform public cloud. Therefore, using the public cloud in comparison to private cloud may bring significant economic benefits to manufacturing.

Virtualization technologies such as hypervisor and container, provide faster deployments, high efficiency to share software packages in enterprises [95]. Virtualization can be used to quickly test and validate the proposed big data solutions with minimum influence on the other systems.

1 **4.1.6 Driver 6: Low cost hardware**

2 Low-cost hardware makes smart manufacturing more accessible. Low-cost actuator and IoT sensor
3 reduce the wiring cost to collect data and improve automation at the factory floor. Wang et al. present
4 a collaboration mechanism to deploy large scale robotics with big data technologies at a small factory
5 [65]. Big-data based solution with smart sensors elevates the constraints of time and geolocation to
6 monitor manufacturing processes [165]. Low cost RFID tags make traceability of enormous resources
7 more feasible at the supply chain level. Industrial Internet of Things hub is proposed to realize the
8 smart connection of various resources in a manufacturing facility with RFID tags and sensors [166].
9 Cost competitive NC machines and 3D printers can be flexibly deployed to demand sides such as
10 design prototype product or serving customers. Big data technologies can secure critical data
11 transmission between design departments and manufacturing equipment. Furthermore, low-cost data
12 processing and storage technologies provide a cost-effective approach to manage large volumes of
13 data from massive data sources in manufacturing [167]. Hence, cheap hardware makes it more
14 feasible for manufacturers to monitor and respond to timely changes from production, to supply chain
15 networks.

16 **4.2 Question 2: What are the essential components of big data ecosystem to** 17 **better serve smart manufacturing?**

18 As big data applications enable smart manufacturing, several essential components of the big data
19 ecosystem should be utilized to build up BDA platform for smart manufacturing, including data
20 ingestion, storage, computing, analytics, visualization, workflow and dataflow, data management,
21 infrastructure and security.

22 **4.2.1 Data ingestion**

23 Data ingestion or inception is of necessity to manufacturing in order to bring big volume data into its
24 BDA platform. There are two types of big data: data at rest and data in motion. Apache Sqoop is used
25 to transfer bundle of data from a relational database (MySQL, SQL server) to Hadoop in several
26 applications such as ERP and MES [116] [65], SCADA [71], O&M [39][79] and Data Management
27 [121][151].

28 Streaming data is data continuously generated from manufacturing systems and devices. Apache
29 Flume is mainly used to collect large amounts of logs from controllers, sensors, equipment and
30 actuators [65][71] [121][116]. Kafka is a general-purpose messaging system to collect streaming data
31 and publish it to data consumers who subscribe to the topic of data. Kafka is applied in some cases of
32 SCADA [140], O&M[139]. Apache Storm is a real-time data processor, which is used to collect and
33 ingest streaming data to data consumers straightway (O&M [138][139], Quality control [85],
34 PLM[58]). Although some data collection tools of manufacturing are widely used in production
35 systems such as MTCConnect and OPC-UA, big data ingestion tools can complement their limitations
36 as discussed in Section 4.4.2. Therefore, data ingestion is an essential component of the big-data
37 ecosystem to collect batching data and streaming data.

38 **4.2.2 Storage**

39 Data storage is critical to big data applications in smart manufacturing. As the manufacturing industry
40 increasingly benefits from the use of big data, it is of importance to store more data [4]. Various
41 applications require different storage technologies to provide different features, which are file system
42 and databases. Although both technologies can store structured, semi-structured and unstructured data,
43 they have some differences that file system is suitable to store data-at-rest or unstructured data such as
44 files, search and compile files manually. Database is suitable to store data-in-motion, semi-structured
45 or structured data, faster query data automatically.

46 There are three types of databases: RDBMS, NoSQL and NewSQL databases. RDBMS database has
47 been used in manufacturing for general purpose applications for decades such as SQL Server, Oracle,

1 MySQL. Whereas, RDBMS is unable to address the challenges of big data's 3Vs (Volume, Velocity
2 and Variety) for storage and query [32]. NoSQL and NewSQL databases can provide almost
3 unlimited scalability and faster query capability for industrial big data. Figure 4 illustrates that the
4 increased use of a NoSQL database is already happening in manufacturing. NoSQL is suitable for one
5 kind of OLTP application, which does not require consistent data all the time, but simple query and
6 frequent updates to data [168]. Hence, manufacturing could use NoSQL databases for real-time big
7 data analytics such as quality monitoring and prediction. Moreover, manufacturing could benefit from
8 using the four types of NoSQL data models (Key-value, Wide column, Document, Graph) to easily
9 manage semi-structured data (XML and JSON) [43]. The widely used NoSQL databases are Redis,
10 HBase, Cassandra, MongoDB and Neo4j. Since NoSQL databases were not designed to meet data
11 consistency, they are not suitable for some OLTP applications, where data consistency needs to be
12 guaranteed anytime [169]. Therefore, NoSQL should not be used in an environment of many
13 operations and controllers that are present since the control data may be inconsistent. The NewSQL
14 solution is used to provide relational query (SQL) and data consistency all the time for an OLTP
15 application. Manufacturing could use NewSQL for the scenarios requiring consistent data all the time
16 (finance data in ERP, inventory data in SCM, control signal in SCADA systems). Some examples of
17 NewSQL are VoltDB, Clustrix and NuoDB [43]. However, NewSQL focuses on relational data,
18 which may not fully support unstructured data for big data analytics.

19 Unlike database solutions, in which the data is structured as a data model for data consumers'
20 demands, file systems do not need a data model to store unstructured data. Hadoop Distributed File
21 System (HDFS) could store petabytes of data with a redundantly low-cost method [30]. However,
22 querying data in HDFS is much slower than the speed in databases. HDFS could be utilized as a
23 central data storage to retain all the data in manufacturing. Therefore, selection of the correct storage
24 solutions is required for big data applications in manufacturing.

25 **4.2.3 Computation**

26 Computation is the foundation of implementing big data applications. Three types of computation
27 engines are available to manufacturing: batching, micro-batching and streaming. MapReduce in
28 Hadoop provides the batching method to process a petabyte level of big data by less memory usage
29 and cannot provide real-time analytics [170]. Spark is a micro-batching processing engine, which
30 provides near real-time computation with more memory resources than MapReduce. Flink and Storm
31 are real-time streaming engines to process small volumes of data [171]. The differences and
32 application scenarios of these computation engines have been discussed in Section 3.3. Based on the
33 analysis of these computation engines and outcomes of Figure 4, there would be more Spark-based
34 big data solutions on the factory floor from device to production planning in the near future. The
35 proposed solutions could use Spark to replace Hadoop to get faster outcomes.

36 **4.2.4 Analytics**

37 Big data analytics is intended to extract information from collected big data. Big data analytics
38 includes two types of analysis: 1) data mining and machine learning algorithms (clustering, regression,
39 Bayesian networks, artificial neural networks (ANN), deep learning. 2) On-Line Analytic Processing
40 (OLAP).

41 Big data analytics tools of the first type are identified as MLlib and Scikit-learn for machine learning
42 [25][109], Spark R for high-level statistical analysis [73], Tensorflow and CNTK for deep learning
43 [105]. Analytics tools have many potential use cases in manufacturing, such as machine vision for
44 robotics[105], speech recognition for alarm and security, image processing for quality control [87]and
45 O&M[82]. Manufacturers can save time to develop algorithms from scratch by using these tools in
46 their big data platforms, such as image recognition for product quality control[87]. Manufacturers also
47 benefits from new tools such as CaffeOnSpark [105], which is a deep learning framework widely used
48 for autonomous driving in the automotive industry[105]. Because Caffe does not work on Spark
49 clusters, it is challenging to meet the two strengths of better data analytics algorithm and faster big
50 data computation at the same time. CaffeOnSpark could address the issue to work on Spark and
51 Hadoop clusters for Caffe applications.

1 OLAP is an approach to analyse large multidimensional datasets for complex business analytics, such
2 as BI reporting, Decision Support System (DSS), and CRM in manufacturing. Analytics tools of
3 OLAP are Apache Hive [172] and Apache Kylin [173]. Hive is one utility of Hadoop ecosystem
4 which is used as a data warehouse. Kylin can query a large volume of data faster than Hive.

5 By analysing the collected big data, data analytics tools can efficiently extract timely information to
6 manufacturers to make decisions. It is challenging to the decision makers to apply their experience
7 and knowledge on the new circumstances. Their experience is acquired under the previous
8 circumstance, which may be different from the current one and their knowledge may be out of date.
9 With the given data and big data analytics tools, manufacturers can analyse historic data, discover
10 new knowledge, build actionable intelligence to make data-driven decisions. It would happen in some
11 areas of different systems feed with new data such as production planning with streaming real-time
12 IoT data.

13 **4.2.5 Visualization**

14 Manufacturing systems require various visualization methods to present analytics results, such as
15 interactive dashboard, reporting, graph, document. Most of the proposed solutions did not provide a
16 tool to construct visualization. As the Python language is widely used in big data analytics, it should
17 be convenient to use Python plotting library (Matplotlib) to present big data analytic results by data
18 scientists in manufacturing. However, Matplotlib works on a command line interface (CLI), which is
19 not user-friendly to business users with less programming experience. The reviewed solutions could
20 benefit from visualization tools such as Zeppelin, Tableau, and D3.js[76]. Zeppelin provides
21 interpreters with big data tools (HBase, Cassandra, HDFS, Spark, Flink etc.) and supports multi-
22 agents. Manufacturers could use Zeppelin in their existing big data platform with less developing
23 effort. D3.js provides more complex visualization templates than Zeppelin. D3.js is a JavaScript
24 library for producing dynamic, interactive data visualizations in the web browser. It is suitable for
25 applications, which require dynamic monitoring and control such as MES, SCADA, O&M, Quality
26 control, safety and security systems[126]. D3.js has less support to produce a report, which is not
27 suitable for business analyses. Tableau and QlikView are commercial data visualization software
28 focusing on BI and support many common databases. . Big data and visualisation can help bring data
29 together and show the value of big data in meaningful ways. AI helps to make automatic decisions,
30 and visualisation helps to make manual decisions. As a result, AI and human can collaboratively
31 make data-driven decisions based on the actionable intelligence generated by big data or mined by
32 machine learning.

33 **4.2.6 Workflow and dataflow**

34 Workflow and dataflow components provide efficient approaches to manage workflows for business
35 and processes to data management. Two workflow tools are found in the review papers: Oozie[116]
36 and Kelper [63]. Oozie is specialized in managing Hadoop jobs, which would be suitable for a
37 manufacturer with Hadoop ready platform [174] . Kepler provides a convenient method to share
38 workflows with other business users in cloud manufacturing [63]. However, Kepler does not support
39 horizontal scalability to store large volume data. Another two workflow tools Wings and IBM
40 InfoSphere have not been found through this review [175]. Wings/Pegasus provides a more intelligent
41 method to construct and execute workflow for users by optimising the workflow automatically. For
42 parallel executing of workflows, Wings has to work with a resource management framework
43 (Pegasus), which may be a constraint to manufacturers with Hadoop-based platform since Hadoop
44 uses Yarn as its cluster resource manager. IBM InfoSphere is a commercial enterprise-ready
45 framework, which is suitable for inexperienced users. It can work with Hadoop and other IBM
46 software together. InfoSphere is not open source, which is unable to build add-on functions. For the
47 dataflow component, researchers had to construct their dataflow component from scratch. Usually,
48 there is a long learning curve for beginners to go through trial-error processes. Apache NiFi is an
49 efficient tool to construct dataflow and data integration with an interactive GUI. One paper uses NiFi
50 to collect streaming data in process industry[123]. NiFi is initially developed and used by the National
51 Security Agency of the United States (NSA), which has been verified in the real application

1 environment. Through this study, there is limited work reported in this area. Smart manufacturing
2 needs the workflow and dataflow components to get data automatically processed among various
3 manufacturing systems. If assisted by AI and optimization, they could significantly improve the
4 efficiency of workflow and dataflow in manufacturing such as production planning and scheduling.

5 **4.2.7 Data management**

6 Data management would make big data platform more feasible for manufacturers. Data management
7 focuses on data governance, metadata management, data modelling, data quality management, master
8 data management, data integration and knowledge management. Although data management is highly
9 related to the traditional IT area, it provides a holistic management method to meet the requirements
10 of enterprise compliance, data policy, data lifecycle.

11 Some data management tools are available to the big data ecosystem of manufacturing:

- 12 • Data lifecycle management is based on enterprises' policy to manage data lifecycle from
13 creation, storage, obsolescence to delete. Apache Falcon includes the functions of data
14 retention (persistence), data replication for disaster recovery, aggregation and archive on
15 Hadoop[176]. Falcon is beneficial to manufacturing research. For example, each acoustic
16 experiment of aerospace engine collects hundreds of gigabytes of audio data through over a
17 hundred sensors. The traditional method uses a single disk to store them, which is costly to
18 manage with high risk of disk failure or data loss. Falcon could manage the data lifecycle with
19 a friendly user interface.
- 20 • Data governance provides an enterprise policy approach to manage data availability, usability,
21 and integrity. Apache Atlas is the Hadoop ecosystem tool for audit, lineage and service level
22 agreement (SLA), which is used to apply agile enterprise compliance through consistent
23 metadata management across the big data ecosystem [177].
- 24 • Data authorization manages users' privilege to access sensitive data. Apache Sentry is utilized
25 to authorize data and metadata on Hadoop clusters based on the roles of people in
26 manufacturing [121].
- 27 • Data quality management is crucial to the outcomes of data analyses. Table 2 presents that
28 data quality is a challenging data issue to manufacturing. Apache Griffin maintains data
29 quality by automating data profiling and validation on Hadoop and Spark[177]. Apache
30 Griffin improves data quality by pre-processing data automatically from various data sources,
31 which reduces the amount of data analyst's time to prepare data for analysing.
- 32 • Data integration has two approaches to address data silos issue [178]: Data Warehouse (DW)
33 and Data Lake (DL). DW is the traditional approach that integrates data from various data
34 sources into a central data store with a predefined extract-transform-load (ETL)
35 method(schema-on-read) [176]. As data volume increases, distributed DW solution (Hive)
36 and associated ETL tool (Apache Pig) is available to manufacturing for business applications
37 such as BI reporting, DSS. However, the predefined ETL method is not flexible and
38 expensive to build before this finally used by data consumers. Data Lake is a new solution of
39 data integration, which is defined as central storage to store any data (sizes, types, rates) with
40 the raw formats in an enterprise [190]. DL is more suitable to data consumers for an ad-hoc
41 query of the data, which is undefined until issuing the query(schema-on-read)[179]. HDFS is
42 a popular DL tool to store extensive unstructured data in manufacturing (video, audio,
43 image)[105]. However, DL probably becomes a "data swamp" without practical data
44 management tools (Falcon, Sentry, Atlas etc.)[176].
- 45 • Semantic KM includes a series of operations for including, creating, classifying, sharing,
46 using information and knowledge in manufacturing [180]. Many technologies of the semantic
47 web are utilized to manage knowledge in manufacturing such as application of ontology web
48 language (OWL) in SCADA [134], Resource Description Framework (RDF), RDF query
49 language (SPARQL) in SCM[91], RDF database(Jena) in KM[23].

50 Data management tools include broad areas to address the veracity and value of big data issues in
51 manufacturing. The tools are still developing since most of the tools are Hadoop-based. As

1 manufacturing begins to focus on velocity of big data computation, some new data management tools
2 based on faster computation engines would be developed in the near future such as Spark and Flink.
3 The big data management tools are still manually operated by data stewards, which would be
4 challenging to manufacturers when various massive data is ingested into the big data platform, and
5 different roles of users apply to use them. New methods of big data management could address this
6 issue by using algorithms such as rule-based, machine learning-based or hybrid of both.

7 **4.2.8 Infrastructure and deployment model**

8 The infrastructure and deployment model are the foundation of big data applications in manufacturing.
9 There are two types of cloud computing infrastructures:

- 10 1. General purpose commodity computer cluster (Hadoop, Spark) to process completely parallel
11 computing problems such as computing massive sensor data separately or responding millions
12 of users' requests (Facebook);
- 13 2. High-Performance Computing (HPC) cluster provides faster computing speed with dedicated
14 hardware.

15 Hence, HPC outperforms others in processing highly dependent data computing such as complex
16 modelling and simulation workload in manufacturing. However, the disadvantages of HPC are that its
17 investment is enormous, and its utilization rate is low. Some solutions address its low utilization by
18 moving HPC to cloud [182]. Some IAAS providers offer public cloud HPC (AWS) and hybrid cloud
19 HPC service (Microsoft Azure).

20 Manufacturing enterprises have different perspectives on economic, security, the privacy of big data
21 platforms. Four deployment models (Public cloud, Private cloud, Hybrid cloud, Community cloud
22 [19]) are available to satisfy the requirements [183]. For example: in terms of better privacy and
23 reducing waste of computing resource, some platforms could adopt a hybrid cloud model, which puts
24 sensitive data on private cloud and process insensitive data on public cloud. However, Hybrid and
25 Community cloud are not found by this review. To meet diverse requirements of manufacturing
26 enterprises, some deployment technologies, including OpenStack and Docker [191], are used to
27 quickly deploy an agile software environment with different micro-service frameworks and
28 programming languages.

29 All deployment models and infrastructure should be taken into consideration in the big data
30 ecosystem of manufacturing.

31 **4.2.9 Cybersecurity**

32 Cybersecurity is an essential component of big data platform to protect data assets in manufacturing.
33 As manufacturing becomes data-driven, the standards and tools are required to secure data in the IT
34 architecture of manufacturers completely. Although some security standards have been provided in
35 manufacturing systems such as SCADA [192], the standards cannot fully address the challenges of
36 SCADA on the Internet [193]. Traditional control systems are vulnerable to unauthorized attacks from
37 the Internet since they are designed as close systems with few capabilities of cybersecurity [194].
38 Because big data platform integrates the physical space and the cyber space closely, the risk of
39 cybersecurity could quickly escalate to the physical system in manufacturing. If unauthorized people
40 manage the critical equipment or information, it will bring irreversible disaster to manufacturing such
41 as economic loss, personal safety.

42 **4.3 Question 3: How can we harness the capabilities available in the big data** 43 **ecosystem to drive research innovations in manufacturing?**

44 Big data ecosystem is the comprehension of massive functional components with various enabling
45 tools. Capabilities of the big data ecosystem are not only about computing and storing big data, but
46 also the advantages of its systematic platform and potentials of big data analytics. Hence, according to
47 proposed solutions of reviewed literature and big data capabilities, the maturity of big data ecosystem
48 application is categorized into three stages:

- 1 Stage 1: proposing a big data framework and platform;
- 2 Stage 2: harvesting cloud computing capacity for big data computing and storage;
- 3 Stage 3: analysing big data with various algorithms for the applications (prediction, fault detection,
- 4 optimisation etc.).

5 Table 3 presents the allocation of proposed solutions by these three stages. MES, SCADA and O&M
6 have been studied within the three stages. However, big-data based solution is missing in some
7 manufacturing systems. For example, no general big data framework has been identified in CAx
8 (Stage 1); no practices of Stage 2 and Stage 3 are found in PLM, SCM, BI and AM. Available
9 capabilities of big data ecosystems could drive research innovation in various big data applications
10 (Figure 5) of the immature manufacturing areas.

11 Manufacturing could benefit from new development and deployment methods of the big data
12 ecosystem. Because many tools are fast changing with highly frequent updates, it makes the design of
13 the big data platform dramatically challenging. Most of the proposed solutions in manufacturing
14 constructed the platform from scratch by installing and testing every tool step by step. Researchers
15 spent much time preparing the software but not focusing on programming and big data analytics.
16 Several popular vendors offer their Hadoop distributions to mitigate the issue, such as Cloudera
17 (CDH), Hortonworks (HDP, HDF), MapR, IBM (Infosphere BigInsights), Microsoft (HD Insight),
18 Pivotal HD [195]. Because the tools and versions were tested in the distributions, they are ready to be
19 used by researchers and enterprises in manufacturing. Moreover, the packages could be easily
20 deployed and shared with virtualization technology (VM and container) [191]. However, developers
21 need to evaluate several conditions to select the suitable distribution such as open source, pricing,
22 customer support, sizes of community and so on [195].

23 **4.4 Question 4: What are the future directions of big data applications in** 24 **Manufacturing?**

25 Manufacturing could benefit from fast developing big data technology with new and matured tools.
26 Hence, the full potential of big data has not been discovered in manufacturing. Some potential
27 directions are proposed for the future work of interested researchers:

28 **4.4.1 Modelling and simulation**

29 Modelling and simulation will naturally play an important role in extracting value from data once the
30 volume of data is available [4]. Digital twin and digital thread are the two essential methods with two
31 important tasks: collecting real-time data from manufacturing input and output devices such as CNC
32 machine and sensors and building up the real-time simulation model with the collected data. These
33 two tasks require considerable computing resources and data storage to process streaming data from
34 manufacturing devices. Researchers commonly use Matlab or some scientific software to develop the
35 selected algorithm to simulate the models [196]. As one simulation model may be implemented on a
36 single computer, many models are required to design, and manually implemented on computers. It is a
37 challenging and onerous task to manufacturers to convert the Matlab simulation models to executive
38 programs on the selected cloud computing engine. However, from the reviewed articles, simulation is
39 rare in the big data based solutions. The ongoing practice is that of a digital twin, which provide a
40 real-time, bi-directional management between a physical object and its digital object [36]. It is
41 essential to simulate the digital twin model with parameters across different domains (e.g. product,
42 process and logistics) in order to predict its dynamic performance, such as Virtual/Augmented
43 Reality[4]. For example, simulation of the product model with predefined parameters includes FEA
44 analysis (ANSYS). However, the parameters from other domains (processes and logistics) are not
45 accounted for in the FEA simulation. One potential research direction is customized products by
46 integrating End-Of-Life data from PLM into CAx software. It could improve the product's impact
47 during its EOL phase, such as Prognostics and health management (PHM) [184]). Another direction is
48 that simulation of digital twin models with real-time data could be used for predictive maintenance.
49 Moreover, simulation can provide integral data to data scientists or business users to design machine

1 learning algorithms. Because Matlab simulation manager is available to simulate multiple models on
2 the public cloud (Azure, AWS), more case studies of methods above could be implemented to verify
3 the performance of big data tools for digital twin and digital thread. Since Matlab is commercial
4 software, unlike most of the software in the big data ecosystem that is open source and free to use,
5 users have to purchase the specific licenses to run the models on the cloud. AnyLogic is also a very
6 popular simulation software which provides cloud-based simulation recently. Researchers could use
7 AnyLogic Cloud to deploy and verify their digital twin models.

8 Simulations with these large volumes of data require enormous computing, storage and
9 communication resources, and cloud computing is well placed to satiate these. Another disruptive
10 technology is Quantum computing, which can provide unlimited resources to process and store the
11 data [197]. However, the technology is still under development and still has limitations, such as fault
12 tolerance and error correlation[185].

13 Future research directions could be:

- 14 • Develop big data tools to convert simulation models from scientific software to implement on
15 public cloud or private cloud;
- 16 • Using a simulation method to generate testing and training data for machine learning at
17 planning and decision-making processes;
- 18 • Using general-purpose cloud computing cluster to simulate FEA at product design.

19 **4.4.2 Connectivity and interoperability**

20 Big data ecosystem complements of existing manufacturing approaches to have connectivity and
21 interoperability. With the aim of systems integration and collaboration for smart manufacturing, the
22 basis is that data and information have to be timely collected, correctly formatted, analysed and
23 exchanged among the systems. MTConnect and OPC-UA are both data collection approaches in
24 manufacturing. MTConnect focuses on device level and control level by monitoring CNC machine
25 tools using a predefined consistent data model, data format and definition, which matches different
26 vendors' machine tools. OPC-UA focuses on SCADA, MES and ERP by using a generic data model
27 which can flexibly match more industrial devices with additional configuration effort. There are three
28 challenging issues of both approaches. One issue is the performance limitation of OPC-UA which the
29 CPU of OPC-UA server is identified as the main bottleneck in production [198]. Since both
30 approaches are implemented on edge devices at the factory floor, the capabilities of data processing
31 cannot be flexibly scaled up. The servers of both approaches have to be replaced with better
32 performing hardware, while more data sources are connecting to the systems such as IoT devices.
33 Another issue is interoperability of the systems using both approaches since there is no common
34 ontology on the top of both information models. The last issue is the capabilities of data analysis and
35 information exchange between both systems and other systems such as SCM, PLM, and CRM.
36 Although some solutions are proposed to exchange both data formats and other formats (such as
37 MTConnect and IEEE 1451 [118], OPC-UA and IoT [117], OPC-UA and AutomationML [98]), more
38 solutions are required for data exchange of massive data formats between these two systems and
39 others systems as shown in Table 1.

40 Four existing big data ingestion tools (Sqoop, Flume, Kafka and Storm) and a new tool (NiFi) could
41 complement the weaknesses of both approaches. All five big data tools can scale out the capability of
42 data processing by adding new hardware to the clusters without replacing the old ones. There are
43 some differences between the five tools. The first difference is that Sqoop and Flume are two tools of
44 Hadoop ecosystem; while Kafka and Storm are not dedicated to Hadoop. Secondly, data consumers
45 require to pull data with Kafka, whereas Flume pushes data to consumers. Thirdly, Kafka provides
46 better fault-tolerance and scalable than others. Kafka provides event duplication, which means other
47 nodes continuously make data available when one node is failed. Compared with Kafka, Storm,
48 Flume and NiFi do not provide event duplication, which is theoretically not suitable for application of
49 critical missions such as safety, security, and finance in manufacturing. Storm, Flume and NiFi
50 require less developing effort to work with Hadoop.

1 Interoperability has two main issues: data format and data quality. With the aim of interoperability
2 among systems on BDA platform, data needs to be correctly formatted with good quality before
3 exchange happens. Table 1 demonstrated various data formats in manufacturing systems. In terms of
4 data formats, there are two types of data: data with schema (XML with given XSD/DTD) and data
5 without schema such as XML without given XSD/DTD, JSON and unstructured data (document,
6 report, email). Since data with schema define the schemas in standard schema files (XSD/DTD: PLM
7 XML), data is can be exchanged automatically by mapping elements of each other's schemas. It is
8 challenging to automatically transform data without schema to an intended schema since there is no
9 mapping between both sides. One solution may be using natural language processing (NLP) to
10 process human readable documents. BDA platform also has to support new file formats such as AMF
11 and 3MF[199].

12 Issues of data quality in Table 2 includes missing value, noise data or anomaly data, uncertainty, data
13 outlier, data correlation, timing and synchronization. Data transformation can address these issues to
14 extract data to correct timely information. No generic transformation tool is identified in the reviewed
15 papers which can be used as "One size fit all" tool for all the applications in manufacturing. Because
16 manufacturers use various manufacturing systems with different data characterises, there are many
17 combinations of data issues which require massive data transformation tools to address. It is
18 challenging for manufacturers to develop customized big data transformation tool for every specific
19 scenario. Without addressing these data issues, correct information cannot be extracted and exchanged
20 among systems. Systems integration and collaboration cannot be achieved for smart manufacturing.

21 Future research directions in this area could be:

- 22 • Review the availability and feasibility of data collection tools in various manufacturing
23 scenarios;
- 24 • Develop a generic data transformation solution with big data technologies to exchange data of
25 manufacturing systems;

26 **4.4.3 Standardized big data platform design**

27 Standardisation of big data platform design improves the feasibility of enterprise-ready solutions in
28 manufacturing. Because some essential components are missing in the proposed solutions, it is likely
29 to increase the difficulty to apply them to manufacturing. Missing components would be insufficient
30 to apply big data applications to smart manufacturing. However, there is no standard approach to
31 design big data platform in manufacturing. The reason may be that different profiles of manufacturing
32 enterprises have varieties of system requirements for big data applications.

33 A systematic assessment method is required to analyse the limitations and strengths of the proposed
34 big data solution in the various manufacturing systems. Data issues are not entirely assessed, such as
35 latency of data transmission among clusters, data quality and data format exchange. Therefore, it is
36 essential to provide a standard approach to design big data platform with related assessment method to
37 manufacturing.

38 **4.4.4 Real time big data analytics**

39 Big data analytics go deeper from batch analysing to real-time streaming analysing in manufacturing.
40 On the one hand, streaming big data analysis is considered as a high research requirement in
41 manufacturing [187]. One the other hand, enabling technologies are changing from non-real time
42 analytics to real time. Figure 4 illustrates that the streaming computation engine Spark is becoming a
43 popular tool than the traditional batching engine (MapReduce of Hadoop). However, it also shows
44 that Hadoop still receives more focus than Spark in manufacturing from Figure 4. Batch processing is
45 not able to provide real-time analytics response such as real time monitoring, dynamic scheduling and
46 planning on systems of workshop floor (SCADA, MES). Micro-batching and streaming engines
47 (Spark, Storm, Flink) can provide real-time big data analytics of streaming data [26]. Spark could take
48 advantage of batching and streaming to replace MapReduce engine.

1 Moreover, compared to Storm and Flink, Spark has more powerful analytics tools such as SparkSQL,
2 Spark R, GraphX, and MLlib. However, Storm and Flink outperform Spark in real-time concerns.
3 Hence, with the objective of streaming and analytics, it is necessary to use Storm or Flink with Spark
4 together. The issue is that it requires more development effort to work on several computation engines.
5 Apache Beam provides a uniform abstraction layer to run these real-time engines at the execution
6 layer[44]. Although it has not been used in manufacturing, researchers in manufacturing could focus
7 on analytics logics without spending time on learning various usages of engines. Another fact of big
8 data analytics forwarding to real-time is that data warehousing tools (Pig, Hive) of BI used
9 MapReduce to batch processing large datasets, now that their latest versions support Spark engine.

10 **4.4.5 Cybersecurity in manufacturing**

11 Cybersecurity will continuously challenge manufacturing since security standards are still not
12 available in some system such as SCM [49]. Recently, NIST published a framework to improve
13 cybersecurity on critical infrastructures [188]. It is envisioned that there would be more developments
14 in cybersecurity tools based on the new standard. No big data tool of cybersecurity was found in all
15 reviewed papers, which is likely a promising research direction. Manufacturing could benefit from the
16 following new security tools. Firstly, Apache Metron is an enterprise-ready real-time big data security
17 tool, which is used by Telstra Company [200]. Secondly, Apache Ranger provides security
18 administration management on Hadoop clusters. Thirdly, Apache Knox provides gateway service to
19 access Hadoop clusters. Future direction in this area is to explore the capability of big data
20 cybersecurity tools on critical systems in manufacturing such as safety, security and SCADA.
21 Vulnerability in both the physical and cyberspace of the manufacturing systems must be identified
22 and protected. The potential risk and damages warrant high priority of future research in this direction.

23 **5 Conclusions and future work**

24 This paper systematically reviews the state of art of big data research in manufacturing to evaluate the
25 capabilities of big data ecosystem and requirements of smart manufacturing. Six key drivers of big
26 data ecosystem are identified for smart manufacturing, which are system integration, data, prediction,
27 sustainability, resource sharing and hardware. Afterwards, the nine essential components of big data
28 ecosystem are presented to design a feasible big data solution to manufacturing enterprises. These are
29 data ingestion, storage, computing, analytics, visualization, management, workflow, infrastructure and
30 security. The evaluation reveals that there is no enterprise-ready big data solution in the reviewed
31 literature.

32 It is important to note that some research areas have received less attention from the manufacturing
33 community such as PLM, CAx, ERP and SCM. Many big data utilities are applicable to these areas,
34 which could drive research innovation.

35 Regarding future work, there are five promising directions: modelling and simulation, connectivity
36 and interoperability, standardized big data platform design, real-time big data analytics and
37 cybersecurity.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

6 References

- [1] S. de Treville, M. Ketokivi, V. Singhal, Competitive manufacturing in a high-cost environment: Introduction to the special issue, *J. Oper. Manag.* 49–51 (2017) 1–5. doi:10.1016/j.jom.2017.02.001.
- [2] H. (Deutsche P.A. Henning, Kagermann(National Academy of Science and Engineering). Wolfgang, Wahlster (German Research Center for Artificial Intelligence). Johannes, Recommendations for implementing the strategic initiative INDUSTRIE 4.0, Final Rep. Ind. 4.0 WG. (2013) 82.
- [3] Y. Liao, F. Deschamps, E. de F.R. Loures, L.F.P. Ramos, Past, present and future of Industry 4.0 - a systematic literature review and research agenda proposal, *Int. J. Prod. Res.* 55 (2017) 3609–3629. doi:10.1080/00207543.2017.1308576.
- [4] A. Kusiak, Smart manufacturing, *Int. J. Prod. Res.* 56 (2018) 508–517. doi:10.1080/00207543.2017.1351644.
- [5] B. (Serm) Kulvatunyou, N. Ivezic, V. Srinivasan, On Architecting and Composing Engineering Information Services to Enable Smart Manufacturing, *J. Comput. Inf. Sci. Eng.* 16 (2016) 031002. doi:10.1115/1.4033725.
- [6] B.W. Jeon, J. Um, S.C. Yoon, S. Suk-Hwan, An architecture design for smart manufacturing execution system, *Comput. Aided. Des. Appl.* 4360 (2016) 1–14. doi:10.1080/16864360.2016.1257189.
- [7] X. Xu, From cloud computing to cloud manufacturing, *Robot. Comput. Integr. Manuf.* 28 (2012) 75–86. doi:10.1016/j.rcim.2011.07.002.
- [8] J. Li, F. Tao, Y. Cheng, L. Zhao, Big Data in product lifecycle management, *Int. J. Adv. Manuf. Technol.* 81 (2015) 667–684. doi:10.1007/s00170-015-7151-x.
- [9] J. Wang, Q. Chang, G. Xiao, N. Wang, S. Li, Data driven production modeling and simulation of complex automobile general assembly plant, *Comput. Ind.* 62 (2011) 765–775. doi:10.1016/j.compind.2011.05.004.
- [10] S. Yang, B. Bagheri, H.-A. Kao, J. Lee, A Unified Framework and Platform for Designing of Cloud-Based Machine Health Monitoring and Manufacturing Systems, *J. Manuf. Sci. Eng.* 137 (2015) 040914. doi:10.1115/1.4030669.
- [11] H. Sequeira, P. Carreira, T. Goldschmidt, P. Vorst, Energy cloud: Real-time cloud-native energy management system to monitor and analyze energy consumption in multiple industrial sites, *Proc. - 2014 IEEE/ACM 7th Int. Conf. Util. Cloud Comput. UCC 2014.* (2014) 529–534. doi:10.1109/UCC.2014.79.
- [12] R.Y. Zhong, S.T. Newman, G.Q. Huang, S. Lan, Big Data for supply chain management in the service and manufacturing sectors: Challenges, opportunities, and future perspectives, *Comput. Ind. Eng.* 101 (2016) 572–591. doi:10.1016/j.cie.2016.07.013.
- [13] H.M. Chen, R. Schutz, R. Kazman, F. Matthes, Amazon in the air: Innovating with big data at Lufthansa, *Proc. Annu. Hawaii Int. Conf. Syst. Sci.* 2016-March (2016) 5096–5105. doi:10.1109/HICSS.2016.631.
- [14] M. Helu, T. Hedberg, A. Barnard Feeney, Reference architecture to integrate heterogeneous manufacturing systems for the digital thread, *CIRP J. Manuf. Sci. Technol.* 19 (2017) 191–195. doi:10.1016/j.cirpj.2017.04.002.
- [15] F. Tao, F. Sui, A. Liu, Q. Qi, M. Zhang, B. Song, Z. Guo, S.C.Y. Lu, A.Y.C. Nee, Digital twin-driven product design framework, *Int. J. Prod. Res.* 7543 (2018) 1–19. doi:10.1080/00207543.2018.1443229.

- 1 [16] L. Monostori, B. Kádár, T. Bauernhansl, S. Kondoh, S. Kumara, G. Reinhart, O. Sauer, G.
2 Schuh, W. Sihn, K. Ueda, Cyber-physical systems in manufacturing, *CIRP Ann. - Manuf.*
3 *Technol.* 65 (2016) 621–641. doi:10.1016/j.cirp.2016.06.005.
- 4 [17] F. Bonomi, R. Milito, P. Natarajan, J. Zhu, Big Data and Internet of Things: A Roadmap for
5 Smart Environments, 2014. doi:10.1007/978-3-319-05029-4.
- 6 [18] J. Gubbi, R. Buyya, S. Marusic, M. Palaniswami, Internet of Things (IoT): A vision,
7 architectural elements, and future directions, *Futur. Gener. Comput. Syst.* 29 (2013) 1645–
8 1660. doi:10.1016/j.future.2013.01.010.
- 9 [19] P. Mell, T. Grance, The NIST definition of cloud computing, *NIST Spec. Publ.* 145 (2011) 7.
10 doi:10.1136/emj.2010.096966.
- 11 [20] B.V. Dhar, Data Science and Prediction, *Commun. ACM.* 56 (2013) 64–73.
- 12 [21] A.J.C. Trappey, C. V. Trappey, U. Hareesh Govindarajan, A.C. Chuang, J.J. Sun, A review of
13 essential standards and patent landscapes for the Internet of Things: A key enabler for Industry
14 4.0, *Adv. Eng. Informatics.* 33 (2017) 208–229. doi:10.1016/j.aei.2016.11.007.
- 15 [22] R. Kosara, C. Healey, Visualization viewpoints: Data, Information and Knowledge in
16 Visualization, *Comput. Graph.* (2003).
17 http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=1210860.
- 18 [23] S. Wang, J. Wan, D. Li, C. Liu, Knowledge reasoning with semantic data for real-time data
19 processing in smart factory, *Sensors (Switzerland).* 18 (2018) 1–10. doi:10.3390/s18020471.
- 20 [24] L. Spendla, M. Kebisek, P. Tanuska, L. Hrcka, Concept of predictive maintenance of
21 production systems in accordance with industry 4.0, *SAMI 2017 - IEEE 15th Int. Symp. Appl.*
22 *Mach. Intell. Informatics, Proc.* (2017) 405–410. doi:10.1109/SAMI.2017.7880343.
- 23 [25] O. Morariu, C. Morariu, T. Borangiu, S. Răileanu, Manufacturing Systems at Scale with Big
24 Data Streaming and Online Machine Learning, *Stud. Comput. Intell.* 762 (2018) 253–264.
25 doi:10.1007/978-3-319-73751-5_19.
- 26 [26] K. Nagorny, P. Lima-Monteiro, J. Barata, A.W. Colombo, Big Data Analysis in Smart
27 Manufacturing: A Review, *Int. J. Commun. Netw. Syst. Sci.* 10 (2017) 31–58.
28 doi:10.4236/ijcns.2017.103003.
- 29 [27] D. Reinsel, J. Gantz, J. Rydning, Data Age 2025: The Digitization of the World From Edge to
30 Core, 2018. [https://www.seagate.com/files/www-content/our-story/trends/files/idc-seagate-](https://www.seagate.com/files/www-content/our-story/trends/files/idc-seagate-dataage-whitepaper.pdf)
31 [dataage-whitepaper.pdf](https://www.seagate.com/files/www-content/our-story/trends/files/idc-seagate-dataage-whitepaper.pdf).
- 32 [28] A. Gandomi, M. Haider, Beyond the hype: Big data concepts, methods, and analytics, *Int. J.*
33 *Inf. Manage.* 35 (2015) 137–144. doi:10.1016/j.ijinfomgt.2014.10.007.
- 34 [29] J. Dean, S. Ghemawat, MapReduce : Simplified Data Processing on Large Clusters, (n.d.) 1–
35 13.
- 36 [30] K. Shvachko, H. Kuang, S. Radia, R. Chansler, The Hadoop Distributed File System, in: 2010
37 IEEE 26th Symp. Mass Storage Syst. Technol., 2010: pp. 1–10.
38 doi:10.1109/MSST.2010.5496972.
- 39 [31] J. Duda, Business Intelligence and NoSQL Databases, *Inf. Syst. Manag.* 1 (2012) 25–37.
- 40 [32] R. Cattell, Scalable SQL and NoSQL Data Stores, *Acm Sigmod Rec.* 39 (2010) 12–27.
- 41 [33] Y. Bao, L. Ren, L. Zhang, X. Zhang, Y. Luo, Massive sensor data management framework in
42 cloud manufacturing based on Hadoop, *IEEE Int. Conf. Ind. Informatics.* (2012) 397–401.
43 doi:10.1109/INDIN.2012.6301192.
- 44 [34] F. Tao, Y. Cheng, L. Da Xu, L. Zhang, B.H. Li, CCIoT-CMfg: Cloud computing and internet
45 of things-based cloud manufacturing service system, *IEEE Trans. Ind. Informatics.* 10 (2014)
46 1435–1442. doi:10.1109/TII.2014.2306383.

- 1 [35] D. Wu, D.W. Rosen, L. Wang, D. Schaefer, Cloud-based design and manufacturing: A new
2 paradigm in digital manufacturing and design innovation, *CAD Comput. Aided Des.* 59 (2015)
3 1–14. doi:10.1016/j.cad.2014.07.006.
- 4 [36] F. Tao, Q. Qi, A. Liu, A. Kusiak, Data-driven smart manufacturing, *J. Manuf. Syst.* (2018).
5 doi:10.1016/j.jmsy.2018.01.006.
- 6 [37] M. Soualhia, F. Khomh, S. Tahar, Task Scheduling in Big Data Platforms: A Systematic
7 Literature Review, *J. Syst. Softw.* 134 (2017) 170–189. doi:10.1016/j.jss.2017.09.001.
- 8 [38] Y. Demchenko, P. Grosso, C. De Laat, P. Membrey, Addressing big data issues in Scientific
9 Data Infrastructure, *Proc. 2013 Int. Conf. Collab. Technol. Syst. CTS 2013.* (2013) 48–55.
10 doi:10.1109/CTS.2013.6567203.
- 11 [39] J. Moyne, J. Samantaray, M. Armacost, Big Data Capabilities Applied to Semiconductor
12 Manufacturing Advanced Process Control, *IEEE Trans. Semicond. Manuf.* 29 (2016) 283–291.
13 doi:10.1109/TSM.2016.2574130.
- 14 [40] Y. Demchenko, C. De Laat, P. Membrey, Defining architecture components of the Big Data
15 Ecosystem, *2014 Int. Conf. Collab. Technol. Syst. CTS 2014.* (2014) 104–112.
16 doi:10.1109/CTS.2014.6867550.
- 17 [41] S. Landset, T.M. Khoshgoftaar, A.N. Richter, T. Hasanin, A survey of open source tools for
18 machine learning with big data in the Hadoop ecosystem, *J. Big Data.* 2 (2015) 24.
19 doi:10.1186/s40537-015-0032-1.
- 20 [42] S. Binani, A. Gutti, S. Upadhyay, SQL vs. NoSQL vs. NewSQL-A Comparative Study,
21 *Commun. Appl. Electron.* 6 (2016) 43–46.
22 <http://www.caeaccess.org/archives/volume6/number1/binani-2016-cae-652418.pdf>.
- 23 [43] K. Grolinger, W.A. Higashino, A. Tiwari, M.A.M. Capretz, Data management in cloud
24 environments: NoSQL and NewSQL data stores, *J. Cloud Comput.* 2 (2013).
25 doi:10.1186/2192-113X-2-22.
- 26 [44] M. Gökalp, K. Kayabay, M. Zaki, A. Koçyiğit, Big-Data Analytics Architecture for Businesses:
27 a comprehensive review on new open-source big-data tools, (2017) 1–35.
28 doi:10.13140/RG.2.2.30306.84165.
- 29 [45] J. Wan, S. Tang, Z. Shu, D. Li, S. Wang, M. Imran, A. V. Vasilakos, Software-Defined
30 Industrial Internet of Things in the Context of Industry 4.0, *IEEE Sens. J.* 16 (2016) 7373–
31 7380. doi:10.1109/JSEN.2016.2565621.
- 32 [46] M.J. Koop, W. Huang, K. Gopalakrishnan, D.K. Panda, Performance analysis and evaluation
33 of PCIe 2.0 and quad-data rate InfiniBand, *Proc. - Symp. High Perform. Interconnects, Hot*
34 *Interconnects.* (2008) 85–92. doi:10.1109/HOTI.2008.26.
- 35 [47] J. Cheng, L. Da Xu, W. Chen, F. Tao, C.-L. Lin, Industrial IoT in 5G Environment towards
36 Smart Manufacturing, *J. Ind. Inf. Integr.* (2018). doi:10.1016/j.jii.2018.04.001.
- 37 [48] P. Zikopoulos, C. Eaton, *Understanding Big Data: Analytics for Enterprise Class Hadoop and*
38 *Streaming Data: Analytics for Enterprise Class Hadoop and Streaming Data*, 2011.
39 <https://www.immagic.com/eLibrary/ARCHIVES/EBOOKS/I111025E.pdf>.
- 40 [49] Y. Lu, K. Morris, S. Frechette, Current Standards Landscape for Smart Manufacturing
41 Systems, *Natl. Inst. Stand. Technol. NISTIR.* 8107 (2016) 39. doi:10.6028/NIST.IR.8107.
- 42 [50] S. Evdokimov, *RFID and the Internet of Things: Technology, Applications, and Security*
43 *Challenges*, 2010. doi:10.1561/0200000020.
- 44 [51] D. Kiritsis, V.K. Nguyen, J. Stark, How closed-loop PLM improves Knowledge Management
45 over the complete product lifecycle and enables the factory of the future, *Int. J. Prod. Lifecycle*
46 *Manag.* 3 (2008) 54. doi:10.1504/IJPLM.2008.019970.
- 47 [52] P. Buneman, Semistructured data, *Proc. 16th Symp. Princ. Database Syst.* (1997) 117–121.

- 1 doi:10.1145/263661.263675.
- 2 [53] Y. Xu, G. Chen, J. Zheng, An integrated solution—KAGFM for mass customization in
3 customer-oriented product design under cloud manufacturing environment, *Int. J. Adv. Manuf.*
4 *Technol.* 84 (2016) 85–101. doi:10.1007/s00170-015-8074-2.
- 5 [54] S. Chen, C. Yin, X. Li, Implementation of MTCConnect in Machine Monitoring System for
6 CNCs, *Proc. - 2017 5th Int. Conf. Enterp. Syst. Ind. Digit. by Enterp. Syst. ES 2017.* (2017)
7 70–75. doi:10.1109/ES.2017.19.
- 8 [55] S.L. Chan, Y. Lu, Y. Wang, Data-driven cost estimation for additive manufacturing in
9 cybermanufacturing, *J. Manuf. Syst.* 46 (2018) 115–126. doi:10.1016/j.jmsy.2017.12.001.
- 10 [56] S.K. and K.C. and H.Y. and S.O. Yang, G-code conversion from 3D model data for 3D
11 printers on Hadoop systems, in: *2017 4th Int. Conf. Comput. Appl. Inf. Process. Technol.*,
12 2017: pp. 1–4. doi:10.1109/CAIPT.2017.8320709.
- 13 [57] S. Sierla, V. Kyrki, P. Aarnio, V. Vyatkin, Automatic assembly planning based on digital
14 product descriptions, *Comput. Ind.* 97 (2018) 34–46. doi:10.1016/j.compind.2018.01.013.
- 15 [58] Y. Zhang, S. Ren, Y. Liu, T. Sakao, D. Huisingh, A framework for Big Data driven product
16 lifecycle management, *J. Clean. Prod.* 159 (2017) 229–240. doi:10.1016/j.jclepro.2017.04.172.
- 17 [59] Y. Zhang, S. Ren, Y. Liu, S. Si, A big data analytics architecture for cleaner manufacturing
18 and maintenance processes of complex products, *J. Clean. Prod.* (2016).
19 doi:10.1016/j.jclepro.2016.07.123.
- 20 [60] D. Ramanujan, W.Z. Bernstein, M.A. Totorikaguena, C.F. Ilvig, K.B. Ørskov, Generating
21 Contextual Design for Environment Principles in Sustainable Manufacturing Using Visual
22 Analytics, *J. Manuf. Sci. Eng.* 141 (2018) 021016. doi:10.1115/1.4041835.
- 23 [61] T. Hedberg, A.B. Feeney, M. Helu, J.A. Camelio, Toward a Lifecycle Information Framework
24 and Technology in Manufacturing, *J. Comput. Inf. Sci. Eng.* 17 (2017) 021010.
25 doi:10.1115/1.4034132.
- 26 [62] J.W. Wang, J.Y. Yang, J. Zhang, X.X. Wang, W. (Chris) Zhang, Big data driven cycle time
27 parallel prediction for production planning in wafer manufacturing, *Enterp. Inf. Syst.* 12 (2018)
28 1–19. doi:10.1080/17517575.2018.1450998.
- 29 [63] X. Li, J. Song, B. Huang, A scientific workflow management system architecture and its
30 scheduling based on cloud service platform for manufacturing big data analytics, *Int. J. Adv.*
31 *Manuf. Technol.* 84 (2015) 119–131. doi:10.1007/s00170-015-7804-9.
- 32 [64] M.Y. Santos, J. Oliveira e Sá, C. Andrade, F. Vale Lima, E. Costa, C. Costa, B. Martinho, J.
33 Galvão, A Big Data system supporting Bosch Braga Industry 4.0 strategy, *Int. J. Inf. Manage.*
34 37 (2017) 750–760. doi:10.1016/j.ijinfomgt.2017.07.012.
- 35 [65] S. Wang, C. Zhang, C. Liu, D. Li, H. Tang, Cloud-assisted interaction and negotiation of
36 industrial robots for the smart factory, *Comput. Electr. Eng.* 63 (2017) 66–78.
37 doi:10.1016/j.compeleceng.2017.05.025.
- 38 [66] N.M. Khushairi, N.A. Emran, M.M. Mohd Yusof, Database performance tuning methods for
39 manufacturing execution system, *World Appl. Sci. J.* 30 (2014) 91–99.
40 doi:10.5829/idosi.wasj.2014.30.icmrp.14.
- 41 [67] B.W. Jeon, J. Um, S.C. Yoon, S. Suk-hwan, An architecture design for smart manufacturing
42 execution system, 4360 (2017). doi:10.1080/16864360.2016.1257189.
- 43 [68] M. Jirkovský, V. and Obitko, Enabling semantics within industry 4.0, in: *Int. Conf. Ind. Appl.*
44 *Holonic Multi-Agent Syst.*, 2017: pp. 39–52. doi:10.1007/978-3-642-40090-2.
- 45 [69] O. Sauer, Developments and trends in shopfloor-related ICT systems, *IEEE Int. Conf. Ind. Eng.*
46 *Eng. Manag.* 2015-Janua (2014) 1352–1356. doi:10.1109/IEEM.2014.7058859.

- 1 [70] B.E.I. Systron, D. Division, Managing Configuration Control in an Automotive Sensor Mass
2 Customization Manufacturing Product Line, 2006 World Autom. Congr. (2006) 1--16.
- 3 [71] J.-J. Kim, D.-W. Lee, D.-B. Ko, S.-I. Jeong, J.-M. Park, An autonomic computing based on big
4 data platform for high-reliable smart factory, *J. Eng. Appl. Sci.* 12 (2017) 2662–2666.
5 [https://www.scopus.com/inward/record.uri?eid=2-s2.0-](https://www.scopus.com/inward/record.uri?eid=2-s2.0-85041023417&partnerID=40&md5=003f099909f142b4879dffa79001fe1c)
6 [85041023417&partnerID=40&md5=003f099909f142b4879dffa79001fe1c](https://www.scopus.com/inward/record.uri?eid=2-s2.0-85041023417&partnerID=40&md5=003f099909f142b4879dffa79001fe1c).
- 7 [72] P. Reboredo, M. Keinert, Integration of discrete manufacturing field devices data and services
8 based on OPC UA, *IECON Proc. (Industrial Electron. Conf.)* (2013) 4476–4481.
9 doi:10.1109/IECON.2013.6699856.
- 10 [73] A.R. Khan, H. Schioler, M. Kulahci, T. Knudsen, Big data analytics for industrial process
11 control, 2017 22nd IEEE Int. Conf. Emerg. Technol. Fact. Autom. (2017) 1–8.
12 doi:10.1109/ETFA.2017.8247658.
- 13 [74] X. Ye, S.H. Hong, An AutomationML/OPC UA-based Industry 4.0 Solution for a
14 Manufacturing System, *IEEE Int. Conf. Emerg. Technol. Fact. Autom. ETFA. 2018-Septe*
15 *(2018) 543–550.* doi:10.1109/ETFA.2018.8502637.
- 16 [75] D. Mourtzis, E. Vlachou, A cloud-based cyber-physical system for adaptive shop-floor
17 scheduling and condition-based maintenance, *J. Manuf. Syst.* 47 (2018) 179–198.
18 doi:10.1016/j.jmsy.2018.05.008.
- 19 [76] C. Toro, I. Barandiaran, J. Posada, A perspective on knowledge based and intelligent systems
20 implementation in industrie 4.0, *Procedia Comput. Sci.* 60 (2015) 362–370.
21 doi:10.1016/j.procs.2015.08.143.
- 22 [77] E. Poormohammady, J.H. Reelfs, M. Stoffers, K. Wehrle, A. Papageorgiou, Dynamic
23 algorithm selection for the logic of tasks in IoT stream processing systems, 2017 13th Int.
24 Conf. Netw. Serv. Manag. CNSM 2017. 2018-Janua (2018) 1–5.
25 doi:10.23919/CNSM.2017.8256009.
- 26 [78] D. Mourtzis, E. Vlachou, N. Milas, Industrial Big Data as a result of IoT adoption in
27 Manufacturing, *Procedia CIRP.* 55 (2016) 290–295. doi:10.1016/j.procir.2016.07.038.
- 28 [79] J. Moyne, J. Samantaray, M. Armacost, Big data emergence in semiconductor manufacturing
29 advanced process control, *Adv. Semicond. Manuf. Conf. (ASMC), 2015 26th Annu. SEMI.*
30 *(2015) 130–135.* doi:10.1109/ASMC.2015.7164483.
- 31 [80] D. Mourtzis, A. Vlachou, V. Zogopoulos, Cloud-Based Augmented Reality Remote
32 Maintenance Through Shop-Floor Monitoring: A Product-Service System Approach, *J. Manuf.*
33 *Sci. Eng.* 139 (2017) 061011. doi:10.1115/1.4035721.
- 34 [81] Y. Busnel, N. Riveei, A. Gal, Y. Busnel, N. Riveei, A. Gal, A. Gal, FlinkMan : Anomaly
35 Detection in Manufacturing Equipment with Apache Flink : Grand Challenge, (2017).
36 doi:10.1145/3093742.3095099.
- 37 [82] I. Yen, S. Zhang, F. Bastani, A Framework for IoT-Based Monitoring and Diagnosis of
38 Manufacturing Systems, in: *Proc. - 11th IEEE Int. Symp. Serv. Syst. Eng. SOSE 2017, 2017:*
39 *pp. 1–8.* doi:10.1109/SOSE.2017.26.
- 40 [83] S. Kang, W.T.K. Chien, J.G. Yang, A study for big-data (Hadoop) application in
41 semiconductor manufacturing, *IEEE Int. Conf. Ind. Eng. Eng. Manag.* 2016-Decem (2016)
42 1893–1897. doi:10.1109/IEEM.2016.7798207.
- 43 [84] X. Li, Z. Tu, Q. Jia, X. Man, H. Wang, X. Zhang, Deep-level quality management based on
44 big data analytics with case study, *Proc. - 2017 Chinese Autom. Congr. CAC 2017. 2017-*
45 *Janua (2017) 4921–4926.* doi:10.1109/CAC.2017.8243651.
- 46 [85] M. Syafrudin, N.L. Fitriyani, D. Li, G. Alfian, J. Rhee, Y.S. Kang, An open source-based real-
47 time data processing architecture framework for manufacturing sustainability, *Sustain.* 9
48 (2017). doi:10.3390/su9112139.

- 1 [86] S. Venkatesh, Web-Enabled Real-Time Quality Feedback for Factory Systems Using
2 MTConnect, in: ASME 2012 Int. Des. Eng. Tech. Conf. Comput. Inf. Eng. Conf., 2012: pp.
3 403–409.
- 4 [87] H.-M. Hou, J.-F. Kung, Y.-B. Hsu, Y. Yamazaki, K. Maruyama, Y. Toyoshima, C. Chen,
5 Prediction of ppm level electrical failure by using physical variation analysis, Proc.SPIE.
6 (2016) 9778. doi:10.1117/12.2229410.
- 7 [88] M.-K. Zheng, X.-G. Ming, X.-Y. Zhang, G.-M. Li, MapReduce Based Parallel Bayesian
8 Network for Manufacturing Quality Control, Chinese J. Mech. Eng. 30 (2017) 1216–1226.
9 doi:10.1007/s10033-017-0179-0.
- 10 [89] S.A. Jacobs, A. Dagnino, Large-Scale Industrial Alarm Reduction and Critical Events Mining
11 Using Graph Analytics on Spark, in: Proc. - 2016 IEEE 2nd Int. Conf. Big Data Comput. Serv.
12 Appl. BigDataService 2016, 2016: pp. 66–71. doi:10.1109/BigDataService.2016.21.
- 13 [90] H. Kemper, H. Baars, H. Lasi, An integrated business intelligence framework: Closing the gap
14 between IT support for management and for production, in: Bus. Intell. Perform. Manag.,
15 Springer London, 2013: pp. 13–27. doi:10.1007/978-1-4471-4866-1.
- 16 [91] B.R. Ferrer, W.M. Mohammed, J.L.M. Lastra, A solution for processing supply chain events
17 within ontology-based descriptions, (2016) 4877–4883.
- 18 [92] M.J.A.G. Izaguirre, A. Lobov, J.L.M. Lastra, OPC-UA and DPWS interoperability for factory
19 floor monitoring using complex event processing, IEEE Int. Conf. Ind. Informatics. (2011)
20 205–210. doi:10.1109/INDIN.2011.6034874.
- 21 [93] J. Campos, P. Sharma, U.G. Gabiria, E. Jantunen, D. Baglee, A Big Data Analytical
22 Architecture for the Asset Management, Procedia CIRP. 64 (2017) 369–374.
23 doi:10.1016/j.procir.2017.03.019.
- 24 [94] P. Gölzer, L. Simon, P. Cato, M. Amberg, Designing global manufacturing networks using Big
25 Data, Procedia CIRP. 33 (2015) 191–196. doi:10.1016/j.procir.2015.06.035.
- 26 [95] M. Zimmermann, U. Breitenbucher, M. Falkenthal, F. Leymann, K. Saatkamp, Standards-
27 Based Function Shipping - How to Use TOSCA for Shipping and Executing Data Analytics
28 Software in Remote Manufacturing Environments, Proc. - 2017 IEEE 21st Int. Enterp. Distrib.
29 Object Comput. Conf. EDOC 2017. 2017-Janua (2017) 50–60. doi:10.1109/EDOC.2017.16.
- 30 [96] L.K. B, C. Gr, K. Jan, E. Hoos, C. Kiefer, C. Weber, S. Silcher, B. Mitschang, The stuttgart IT
31 architecture for manufacturing an architecture for the data-driven factory, Int. Conf. Enterp. Inf.
32 Syst. (2016) 53–80. doi:10.1007/978-3-319-62386-3.
- 33 [97] W.M. Mohammed, B.R. Ferrer, L. Jose, M. Lastra, D. Aleixo, C. Agostinho, Configuring and
34 visualizing the data resources in a cloud-based data collection framework, 2017 Int. Conf. Eng.
35 Technol. Innov. Eng. Technol. Innov. Manag. Beyond 2020 New Challenges, New
36 Approaches, ICE/ITMC 2017 - Proc. 2018-Janua (2018) 1201–1208.
37 doi:10.1109/ICE.2017.8280017.
- 38 [98] X. Ye, T.Y. Park, S.H. Hong, Y. Ding, A. Xu, Implementation of a Production-Control System
39 Using Integrated Automation ML and OPC UA, 2018 Work. Metrol. Ind. 4.0 IoT, MetroInd
40 4.0 IoT 2018 - Proc. (2018) 242–247. doi:10.1109/METROI4.2018.8428310.
- 41 [99] Y. Cao, S. Wang, L. Kang, C. Li, L. Guo, Study on machining service modes and resource
42 selection strategies in cloud manufacturing, Int. J. Adv. Manuf. Technol. 81 (2015) 597–613.
43 doi:10.1007/s00170-015-7222-z.
- 44 [100] S.K. Panda, T. Schroder, L. Wisniewski, C. Diedrich, PlugProduce Integration of Components
45 into OPC UA based data-space, IEEE Int. Conf. Emerg. Technol. Fact. Autom. ETFA. 2018-
46 Septe (2018) 1095–1100. doi:10.1109/ETFA.2018.8502663.
- 47 [101] F. Tao, L. Zhang, Y. Liu, Y. Cheng, L. Wang, X. Xu, Manufacturing Service Management in
48 Cloud Manufacturing: Overview and Future Research Directions, J. Manuf. Sci. Eng. 137

- 1 (2015) 040912. doi:10.1115/1.4030510.
- 2 [102] H. Yang, M. Park, M. Cho, M. Song, S. Kim, A system architecture for manufacturing process
3 analysis based on big data and process mining techniques, *Big Data (Big Data)*, 2014 IEEE Int.
4 Conf. (2014) 1024–1029. doi:10.1109/BigData.2014.7004336.
- 5 [103] J.C.C. Tseng, J.Y. Gu, P.F. Wang, C.Y. Chen, C.F. Li, V.S. Tseng, A scalable complex event
6 analytical system with incremental episode mining over data streams, 2016 IEEE Congr. Evol.
7 Comput. CEC 2016. (2016) 648–655. doi:10.1109/CEC.2016.7743854.
- 8 [104] B. Suryajaya, C.C. Chen, M.H. Hung, Y.Y. Liu, J.X. Liu, Y.C. Lin, A fast large-size
9 production data transformation scheme for supporting smart manufacturing in semiconductor
10 industry, *IEEE Int. Conf. Autom. Sci. Eng. 2017-Augus (2018)* 275–281.
11 doi:10.1109/COASE.2017.8256114.
- 12 [105] A. Luckow, M. Cook, N. Ashcraft, E. Weill, E. Djerekarov, B. Vorster, Deep Learning in the
13 Automotive Industry: Applications and Tools, *Big Data Int. Conf. Big Data. (2016)* 3759–
14 3768. doi:10.1109/BigData.2016.7841045.
- 15 [106] A. Brodsky, G. Shao, M. Krishnamoorthy, A. Narayanan, D. Menasce, R. Ak, Analysis and
16 optimization in smart manufacturing based on a reusable knowledge base for process
17 performance models, *Proc. - 2015 IEEE Int. Conf. Big Data, IEEE Big Data 2015. (2015)*
18 1418–1427. doi:10.1109/BigData.2015.7363902.
- 19 [107] S.C. Feng, W.Z. Bernstein, T. Hedberg, A. Barnard Feeney, Toward Knowledge Management
20 for Smart Manufacturing, *J. Comput. Inf. Sci. Eng. 17 (2017)* 031016. doi:10.1115/1.4037178.
- 21 [108] M. Stonebraker, I.F. Ilyas, Data Integration: The Current Status and the Way Forward., *IEEE*
22 *Data Eng. Bull. 41 (2018)* 3–9.
- 23 [109] D. Zhang, B. Xu, J. Wood, Predict failures in production lines: A two-stage approach with
24 clustering and supervised learning, *Proc. - 2016 IEEE Int. Conf. Big Data, Big Data 2016.*
25 (2016) 2070–2074. doi:10.1109/BigData.2016.7840832.
- 26 [110] V. Jirkovský, M. Obitko, P. Novák, P. Kadera, Big data analysis for sensor time-series in
27 automation, 19th IEEE Int. Conf. Emerg. Technol. Fact. Autom. ETFA 2014. (2014).
28 doi:10.1109/ETFA.2014.7005183.
- 29 [111] J. Bakakeu, J. Fuchs, T. Javied, M. Brossog, J. Franke, H. Klos, W. Eberlein, S. Tolksdorf, J.
30 Peschke, L. Jahn, Multi-Objective Design Space Exploration for the Integration of Advanced
31 Analytics in Cyber-Physical Production Systems, *IEEE Int. Conf. Ind. Eng. Eng. Manag.*
32 2019-Decem (2019) 1866–1873. doi:10.1109/IEEM.2018.8607483.
- 33 [112] R. Lynn, W. Louhichi, M. Parto, E. Wescoat, T. Kurfess, RAPIDLY DEPLOYABLE
34 MTCONNECT-BASED MACHINE TOOL MONITORING SYSTEMS, in: *Proc. ASME*
35 *12TH Int. Manuf. Sci. Eng. Conf. - 2017, VOL 3, 2017.*
- 36 [113] D. Libes, S. Shin, J. Woo, Considerations and recommendations for data availability for data
37 analytics for manufacturing, in: *3rd IEEE Int. Conf. Big Data, IEEE Big Data, 2015: pp. 68–75.*
38 doi:10.1109/BigData.2015.7363743.
- 39 [114] C. Zhao, L. Zhang, X.Z.L. Zhang, Cloud manufacturing resource management based on
40 metadata, *ASME 2015 Int. Manuf. Sci. Eng. Conf. MSEC 2015. 2 (2015)* 1–8.
41 doi:10.1115/MSEC20159388.
- 42 [115] Y. Cheng, W. Shang, L. Zhu, D. Zhang, D. Feng, Items analysis of postal supervision, 2016
43 IEEE/ACIS 15th Int. Conf. Comput. Inf. Sci. ICIS 2016 - Proc. (2016) 3–5.
44 doi:10.1109/ICIS.2016.7550949.
- 45 [116] P. Gaj, K. Andrzej, P. Stera, Ontology-Based Integrated Monitoring of Hadoop Clusters in
46 Industrial Environments with OPC UA and RESTful Web Services, *Commun. Comput. Inf.*
47 *Sci. 522 (2015)* 162–171. doi:10.1007/978-3-319-19419-6.

- 1 [117] H. Derhamy, J. Ronnholm, J. Delsing, J. Eliasson, J. Van Deventer, Protocol interoperability
2 of OPC UA in service oriented architectures, Proc. - 2017 IEEE 15th Int. Conf. Ind.
3 Informatics, INDIN 2017. (2017) 44–50. doi:10.1109/INDIN.2017.8104744.
- 4 [118] K.B. Lee, E.Y. Song, P.S. Gu, Integration of MTCConnect and Standard-Based Sensor
5 Networks for Manufacturing Equipment Monitoring, in: ASME 2012 Int. Manuf. Sci. Eng.
6 Conf. MSEC, 2012: pp. 4–8.
- 7 [119] A.N. Richter, T.M. Khoshgoftaar, S. Landset, T. Hasanin, A Multi-dimensional Comparison of
8 Toolkits for Machine Learning with Big Data, Proc. - 2015 IEEE 16th Int. Conf. Inf. Reuse
9 Integr. IRI 2015. (2015) 1–8. doi:10.1109/IRI.2015.12.
- 10 [120] A. Jos, Integration of Sensors , Controllers and Instruments, Sensors. (2017).
11 doi:10.3390/s17071512.
- 12 [121] A. Luckow, K. Kennedy, F. Manhardt, E. Djerekarov, B. Vorster, A. Apon, Automotive big
13 data: Applications, workloads and infrastructures, Proc. - 2015 IEEE Int. Conf. Big Data,
14 IEEE Big Data 2015. (2015) 1201–1210. doi:10.1109/BigData.2015.7363874.
- 15 [122] C. Mathis, Data Lakes, Datenbank-Spektrum. (2017). doi:10.1007/s13222-017-0272-7.
- 16 [123] M. Sarnovsky, P. Bednar, M. Smatana, Data integration in scalable data analytics platform for
17 process industries, in: 2017 IEEE 21st Int. Conf. Intell. Eng. Syst., 2017: pp. 187–192.
18 doi:10.1109/INES.2017.8118553.
- 19 [124] M. Sarnovsky, P. Bednar, M. Smatana, Big Data Processing and Analytics Platform
20 Architecture for Process Industry Factories, Big Data Cogn. Comput. 2 (2018) 3.
21 doi:10.3390/bdcc2010003.
- 22 [125] G.L. Ooi, Y.-H. Wang, P.S. Tan, Z. Zhang, Y. Gao, J.K. Chow, Y. Wu, Q. Yuan,
23 Customizable and scalable geotechnical laboratory testing and field monitoring with new
24 sensing and big data technologies, ICSMGE 2017 - 19th Int. Conf. Soil Mech. Geotech. Eng.
25 (2017) 471–474.
- 26 [126] H. Liang, L. Feng, Z. Chun, Application of the Big Data Technology for Massive Data of the
27 Whole Life Cycle of EMU, in: F. Xhafa, S. Patnaik, Z. Yu (Eds.), Recent Dev. Intell. Syst.
28 Interact. Appl., Springer International Publishing, Cham, 2017: pp. 219–224.
- 29 [127] L. Zheng, L. Tang, T. Li, B. Duan, M. Lei, P. Wang, C. Zeng, L. Li, Y. Jiang, W. Xue, J. Li, C.
30 Shen, W. Zhou, H. Li, Applying data mining techniques to address critical process
31 optimization needs in advanced manufacturing, Proc. 20th ACM SIGKDD Int. Conf. Knowl.
32 Discov. Data Min. - KDD '14. (2014) 1739–1748. doi:10.1145/2623330.2623347.
- 33 [128] S. Windmann, A. Maier, O. Niggemann, C. Frey, A. Bernardi, Y. Gu, H. Pfrommer, T. Steckel,
34 M. Krüger, R. Kraus, Big data analysis of manufacturing processes, J. Phys. Conf. Ser. 659
35 (2015). doi:10.1088/1742-6596/659/1/012055.
- 36 [129] R. Bohlin, L. Lindkvist, J. Hagmar, J.S. Carlson, K. Bengtsson, DATA FLOW AND
37 COMMUNICATION FRAMEWORK SUPPORTING DIGITAL TWIN FOR GEOMETRY
38 ASSURANCE Robert, in: Proc. ASME 2017 Int. Mech. Eng. Congr. Expo., 2017: pp. 1–7.
- 39 [130] K.E. Harper, J. Zheng, S.A. Jacobs, A. Dagnino, A. Jansen, T. Goldschmidt, A. Marinakis,
40 Industrial analytics pipelines, Proc. - 2015 IEEE 1st Int. Conf. Big Data Comput. Serv. Appl.
41 BigDataService 2015. (2015) 242–248. doi:10.1109/BigDataService.2015.38.
- 42 [131] T. Goldschmidt, M.K. Murugaiah, C. Sonntag, B. Schlich, S. Biallas, P. Weber, Cloud-Based
43 Control : A Multi-Tenant , Horizontally Scalable Soft-PLC, in: 2015 IEEE 8th Int. Conf.
44 Cloud Comput., 2015. doi:10.1109/CLOUD.2015.124.
- 45 [132] H.-S. Park, J.-H. Kim, C.-H. Choi, B.-R. Jung, K.-H. Lee, S.-Y. Chi, W.-S. Cho, In-Memory
46 Data Grid System for Real-Time Processing of Machine Sensor Data in a Smart Factory
47 Environment, Proc. 2015 Int. Conf. Big Data Appl. Serv. - BigDAS '15. (2015) 92–97.
48 doi:10.1145/2837060.2837073.

- 1 [133] N. Ramakrishnanus, R. Ghosh, Distributed Dynamic Elastic Nets : A Scalable Approach for
2 Regularization in Dynamic Manufacturing Environments, 2015 IEEE Int. Conf. Big Data.
3 (2015) 2752–2761.
- 4 [134] S. Wang, J. Ouyang, D. Li, C. Liu, An Integrated Industrial Ethernet Solution for the
5 Implementation of Smart Factory, IEEE Access. (2017) 25455–25462.
6 doi:10.1109/ACCESS.2017.2770180.
- 7 [135] S. Division, C. Hsing, N. Village, N. City, N. County, M. Availability, Developing a cloud
8 virtual maintenance system for machine tools management, 2015 11th Int. Conf. Heterog.
9 Netw. Qual. Reliab. Secur. Robustness. (2015) 358–364.
- 10 [136] J. Wan, S. Tang, D. Li, S. Wang, C. Liu, H. Abbas, S. Member, A. V Vasilakos, A
11 Manufacturing Big Data Solution for Active Preventive Maintenance, in: IEEE Trans. Ind.
12 INFORMATICS, 2017: pp. 2039–2047.
- 13 [137] J.-H. Ku, A Study on Prediction Model of Equipment Failure Through Analysis of Big Data
14 Based on RHadoop, Wirel. Pers. Commun. 98 (2017) 3163–3176. doi:10.1007/s11277-017-
15 4151-1.
- 16 [138] W. Lee, J. Cho, L. Lee, S. Korea, Time Series Abnormal Data Detection for Smart Factory, Int.
17 J. Control Autom. 11 (2018) 91–98.
- 18 [139] H. Chen, X. Fei, S. Wang, X. Lu, G. Jin, W. Li, X. Wu, Energy Consumption Data Based
19 Machine Anomaly Detection, 2014 Second Int. Conf. Adv. Cloud Big Data. (2014) 136–142.
20 doi:10.1109/CBD.2014.24.
- 21 [140] N. Stojanovic, M. Dinic, L. Stojanovic, A data-driven approach for multivariate contextualized
22 anomaly detection: Industry use case, Proc. - 2017 IEEE Int. Conf. Big Data, Big Data 2017.
23 2018-Janua (2018) 1560–1569. doi:10.1109/BigData.2017.8258090.
- 24 [141] M. Canizo, E. Onieva, A. Conde, S. Charramendieta, S. Trujillo, Real-time predictive
25 maintenance for wind turbines using Big Data frameworks, IEEE Int. Conf. Progn. Heal.
26 Manag. (2017) 1–8. doi:10.1109/ICPHM.2017.7998308.
- 27 [142] T. Zaarour, N. Pavlopoulou, S. Hasan, U. ul Hassan, E. Curry, Grand challenge: Automatic
28 anomaly detection over sliding windows, Proc. 11th ACM Int. Conf. Distrib. Event-Based Syst.
29 - DEBS '17. (2017) 310–314. doi:10.1145/3093742.3095105.
- 30 [143] J.H. Lee, S. Do Noh, H.-J. Kim, Y.-S. Kang, Implementation of cyber-physical production
31 systems for quality prediction and operation control in metal casting, Sensors (Switzerland). 18
32 (2018). doi:10.3390/s18051428.
- 33 [144] A. Stojadinović, Industry Paper : Dynamic Monitoring for Improving Worker Safety at the
34 Workplace : Use Case from a Manufacturing Shop Floor, (2015) 205–216.
- 35 [145] H. Haskamp, F. Orth, J. Wermann, A.W. Colombo, Implementing an OPC UA interface for
36 legacy PLC-based automation systems using the Azure cloud: An ICPS-architecture with a
37 retrofitted RFID system, Proc. - 2018 IEEE Ind. Cyber-Physical Syst. ICPS 2018. (2018) 115–
38 121. doi:10.1109/ICPHYS.2018.8387646.
- 39 [146] D. Cemernek, H. Gursch, R. Kern, Big data as a promoter of industry 4.0: Lessons of the
40 semiconductor industry, Proc. - 2017 IEEE 15th Int. Conf. Ind. Informatics, INDIN 2017.
41 (2017) 239–244. doi:10.1109/INDIN.2017.8104778.
- 42 [147] C. Ellwein, O. Riedel, O. Meyer, D. Schel, Rent'n'Produce: A Secure Cloud Manufacturing
43 Platform for Small and Medium Enterprises, 2018 IEEE Int. Conf. Eng. Technol. Innov.
44 ICE/ITMC 2018 - Proc. (2018) 1–6. doi:10.1109/ICE.2018.8436332.
- 45 [148] N. Ferry, G. Terrazas, P. Kalweit, A. Solberg, S. Ratchev, D. Weinelt, Towards a big data
46 platform for managing machine generated data in the cloud, Proc. - 2017 IEEE 15th Int. Conf.
47 Ind. Informatics, INDIN 2017. (2017) 263–270. doi:10.1109/INDIN.2017.8104782.

- 1 [149] A. Angrish, B. Starly, Y.S. Lee, P.H. Cohen, A flexible data schema and system architecture
2 for the virtualization of manufacturing machines (VMM), *J. Manuf. Syst.* 45 (2017) 236–247.
3 doi:10.1016/j.jmsy.2017.10.003.
- 4 [150] R.S. Peres, A.D. Rocha, A. Coelho, J. Barata Oliveira, A Highly Flexible, Distributed Data
5 Analysis Framework for Industry 4.0 Manufacturing Systems, in: T. Borangiu, D. Trentesaux,
6 A. Thomas, P. Leitão, J.B. Oliveira (Eds.), *Serv. Orientat. Holonic Multi-Agent Manuf.*,
7 Springer International Publishing, Cham, 2017: pp. 373–381.
- 8 [151] M.Y. Santos, J.O. e Sá, C. Costa, J. Galvão, C. Andrade, B. Martinho, F.V. Lima, E. Costa, A
9 Big Data Analytics Architecture for Industry 4.0, *WorldCIST 2017 Recent Adv. Inf. Syst.*
10 *Technol. 0* (2017). doi:10.1007/978-3-319-56538-5.
- 11 [152] P. Tanuska, L. Spendla, M. Kebisek, Data Integration for Incidents Analysis in Manufacturing
12 Infrastructure, in: *Comput. Conf. 2017*, 2017: pp. 340–345.
- 13 [153] H. Lin, J.A. Harding, C. Chen, A Hyperconnected Manufacturing Collaboration System Using
14 the Semantic Web and Hadoop Ecosystem System, *Procedia CIRP.* 52 (2016) 18–23.
15 doi:10.1016/j.procir.2016.07.075.
- 16 [154] C. Gröger, Building an Industry 4.0 Analytics Platform, *Datenbank-Spektrum.* (2018).
17 doi:10.1007/s13222-018-0273-1.
- 18 [155] M.C. Domenech, L.P. Rauta, M.D. Lopes, P.H. Da Silva, R.C. Da Silva, B.W. Mezger, M.S.
19 Wangham, Providing a smart industrial environment with the web of things and cloud
20 computing, *Proc. - 2016 IEEE Int. Conf. Serv. Comput. SCC 2016.* (2016) 641–648.
21 doi:10.1109/SCC.2016.89.
- 22 [156] A.V. Wan, J., Tang, S., Li, D., Wang, S., Liu, C., Abbas, H. and Vasilakos, A Manufacturing
23 Big Data Solution for Active Preventive Maintenance, *IEEE Trans. Ind. Informatics.* 13 (2017)
24 2039–2047. doi:10.1201/b15906-13.
- 25 [157] N. Rivetti, Y. Busnel, A. Gal, Grand challenge: Flinkman - Anomaly detection in
26 manufacturing equipment with apache flink, *Proc. 11th ACM Int. Conf. Distrib. Event-Based*
27 *Syst. - DEBS '17.* (2017) 274–279. doi:10.1145/3093742.3095099.
- 28 [158] M. Canizo, E. Onieva, A. Conde, S. Charramendieta, S. Trujillo, Real-time Predictive
29 Maintenance for Wind Turbines Using Big Data Frameworks, (2017) 1–8.
- 30 [159] M.R. Brule, Big data in EP: Real-time adaptive analytics and data-flow architecture, *Soc. Pet.*
31 *Eng. - SPE Digit. Energy Conf. Exhib. 2013.* (2013) 305–311. doi:10.2118/163721-MS.
- 32 [160] D. Wu, J. Terpenney, L. Zhang, R. Gao, T. Kurfess, Fog-enabled architecture for data-driven
33 cyber-manufacturing systems, *ASME 2016 11th Int. Manuf. Sci. Eng. Conf. MSEC 2016.* 2
34 (2016) V002T04A032. doi:10.1115/MSEC2016-8559.
- 35 [161] R. Gao, L. Wang, R. Teti, D. Dornfeld, S. Kumara, M. Mori, M. Helu, Cloud-enabled
36 prognosis for manufacturing, *CIRP Ann. - Manuf. Technol.* 64 (2015) 749–772.
37 doi:10.1016/j.cirp.2015.05.011.
- 38 [162] Y. Wu, S. Wang, Streaming Analytics Processing in Manufacturing Performance Monitoring
39 and Prediction, *2017 IEEE Int. Conf. Big Data (Big Data).* (2017) 3285–3289.
- 40 [163] Y. Zhang, S. Ren, Y. Liu, T. Sakao, D. Huisingsh, A framework for Big Data driven product
41 lifecycle management, *J. Clean. Prod.* 159 (2017) 229–240. doi:10.1016/j.jclepro.2017.04.172.
- 42 [164] M. Naeem, N. Moalla, Y. Ouzrout, A. Bouaras, An ontology based digital preservation system
43 for enterprise collaboration, *2014 IEEE/ACS 11th Int. Conf. Comput. Syst. Appl.* (2014) 691–
44 698. doi:10.1109/AICCSA.2014.7073267.
- 45 [165] Q.P. He, J. Wang, D. Shah, N. Vahdat, Statistical Process Monitoring for IoT-Enabled
46 Cybermanufacturing: Opportunities and Challenges, *IFAC-PapersOnLine.* 50 (2017) 14946–
47 14951. doi:10.1016/j.ifacol.2017.08.2546.

- 1 [166] F. Tao, J. Cheng, Q. Qi, IIHub: An Industrial Internet-of-Things Hub Toward Smart
2 Manufacturing Based on Cyber-Physical System, *Ieee Trans. Ind. Informatics*. 14 (2018)
3 2271–2280. doi:10.1109/TII.2017.2759178.
- 4 [167] S. Ramírez-Gallego, A. Fernández, S. García, M. Chen, F. Herrera, Big Data: Tutorial and
5 guidelines on information and process fusion for analytics algorithms with MapReduce, *Inf.*
6 *Fusion*. 42 (2018) 51–61. doi:10.1016/j.inffus.2017.10.001.
- 7 [168] S. Michael, SQL databases v. NoSQL databases, *Commun. ACM*. 53 (2010) 10–11.
8 doi:10.1145/1721654.1721659.
- 9 [169] I. Kovačević, I. Mekterović, Alternative Business Intelligence Engines, *Inf. Commun. Technol.*
10 *Electron. Microelectron.* (2017) 1617–1622. doi:10.23919/MIPRO.2017.7973638.
- 11 [170] P. Kannan, Beyond Hadoop MapReduce Apache Tez and Apache Spark, (2015).
- 12 [171] S. Chintapalli, D. Dagit, B. Evans, R. Farivar, T. Graves, M. Holderbaugh, Z. Liu, K.
13 Nusbaum, K. Patil, B.J. Peng, P. Poulosky, Benchmarking streaming computation engines:
14 Storm, flink and spark streaming, *Proc. - 2016 IEEE 30th Int. Parallel Distrib. Process. Symp.*
15 *IPDPS 2016*. (2016) 1789–1792. doi:10.1109/IPDPSW.2016.138.
- 16 [172] J. Abadi, D., Agrawal, R., Ailamaki, A., Balazinska, M., Bernstein, P.A., Carey, M.J.,
17 Chaudhuri, S., Dean, J., Doan, A., Franklin, M.J. and Gehrke, The Beckman Report on
18 Database Research, *Commun. ACM*. 59 (2016) 92–99.
- 19 [173] S.V. Ranawade, Online Analytical Processing on Hadoop using Apache Kylin, 12 (2017) 1–5.
- 20 [174] M. Islam, A.K. Huang, M. Battisha, M. Chiang, S. Srinivasan, C. Peters, A. Neumann, A.
21 Abdelnur, Oozie: Towards a Scalable Workflow Management System for Hadoop, *Proc. 1st*
22 *ACM SIGMOD Work. Scalable Work. Exec. Engines Technol. - SWEET '12*. (2012) 1–10.
23 doi:10.1145/2443416.2443420.
- 24 [175] K. Sundaravarathan, P. Martin, D. Rope, M. McRoberts, C. Statchuk, MEWSE: Multi-engine
25 Workflow Submission and Execution on Apache YARN, *Proc. 26th Annu. Int. Conf. Comput.*
26 *Sci. Softw. Eng.* (2016) 194–200. <http://dl.acm.org/citation.cfm?id=3049877.3049897>.
- 27 [176] I. Suriarachchi, B. Plale, Crossing Analytics Systems : A Case for Integrated Provenance in
28 Data Lakes, in: 2016 IEEE 12th Int. Conf. e-Science Crossing, IEEE, Baltimore, MD, USA,
29 2016: pp. 349–354. doi:10.1109/eScience.2016.7870919.
- 30 [177] M.O. Gökalp, K. Kayabay, M. Zaki, A. Koçyiğit, P.E. Eren, A. Neely, Big-Data Analytics
31 Architecture for Businesses : a comprehensive review on new open-source big-data tools,
32 (2017).
- 33 [178] B. Stein, A. Morrison, The enterprise data lake: Better integration and deeper analytics, PWC
34 Technol. Forecast Rethink. Integr. (2014). [http://www.pwc.com/us/en/technology-](http://www.pwc.com/us/en/technology-forecast/2014/cloud-computing/assets/pdf/pwc-technology-forecast-data-lakes.pdf)
35 [forecast/2014/cloud-computing/assets/pdf/pwc-technology-forecast-data-lakes.pdf](http://www.pwc.com/us/en/technology-forecast/2014/cloud-computing/assets/pdf/pwc-technology-forecast-data-lakes.pdf).
- 36 [179] N. Miloslavskaya, A. Tolstoy, Big Data, Fast Data and Data Lake Concepts, *Procedia Comput.*
37 *Sci.* 88 (2016) 300–305. doi:10.1016/j.procs.2016.07.439.
- 38 [180] R. Bose, V. Sugumaran, Application of Knowledge Management Technology in Customer
39 Relationship Management, 10 (2003) 3–17. doi:10.1002/kpm.163.
- 40 [181] B.R. Ferrer, W.M. Mohammed, J.L.M. Lastra, A solution for processing supply chain events
41 within ontology-based descriptions, *IECON Proc. (Industrial Electron. Conf.)* (2016) 4877–
42 4883. doi:10.1109/IECON.2016.7793020.
- 43 [182] M. Mercier, D. Glessner, Y. Georgiou, O. Richard, M. Mercier, D. Glessner, Y. Georgiou, O.
44 Richard, B. Data, Big Data and HPC collocation : Using HPC idle resources for Big Data
45 Analytics, in: *Big Data (Big Data)*, 2017 IEEE Int. Conf., 2017: pp. 347--352.
- 46 [183] G. Adamson, L. Wang, M. Holm, P. Moore, Cloud manufacturing – a critical review of recent
47 development and future trends, *Int. J. Comput. Integr. Manuf.* (2017) 1–34.

- 1 doi:10.1080/0951192X.2015.1031704.
- 2 [184] J. Lee, H.D. Ardakani, S. Yang, B. Bagheri, Industrial Big Data Analytics and Cyber-physical
3 Systems for Future Maintenance & Service Innovation, *Procedia CIRP*. 38 (2015) 3–7.
4 doi:10.1016/j.procir.2015.08.026.
- 5 [185] E. Knill, Quantum computing with realistically noisy devices, *Nature*. 434 (2005) 39–44.
6 doi:10.1038/nature03350.
- 7 [186] Deutsche Kommission Elektrotechnik DKE; DIN e.V., German Standardization Roadmap
8 Industry 4.0, (2016) 523.
- 9 [187] W. and S. Kagermann, H., Riemensperger, F., Hoke, D., Helbig, J., Stocksmeier, D., Wahlster,
10 Recommendations for the Strategic Initiative Coordination and editing, 2015.
- 11 [188] NIST, Framework for Improving Critical Infrastructure Cybersecurity, 2018.
12 doi:10.1109/JPROC.2011.2165269.
- 13 [189] P. Wang, R.X. Gao, Z. Fan, Cloud Computing for Cloud Manufacturing: Benefits and
14 Limitations, *J. Manuf. Sci. Eng.* 137 (2015) 044002. doi:10.1115/1.4030209.
- 15 [190] H. Fang, Managing Data Lakes in Big Data Era, in: *Cyber Technol. Autom. Control. Intell.*
16 *Syst. (CYBER)*, 2015 IEEE Int. Conf., 2015: pp. 820–824.
- 17 [191] D. Bernstein, Containers and cloud: From LXC to docker to kubernetes, *IEEE Cloud Comput.*
18 1 (2014) 81–84. doi:10.1109/MCC.2014.51.
- 19 [192] K. Stouffer, J. Falco, K. Kent, Guide to supervisory control and data acquisition (SCADA) and
20 industrial control systems security, NIST Spec. Publ. SP800-82. (2006) 800--82.
21 doi:10.6028/NIST.SP.800.82.
- 22 [193] S. Nazir, S. Patel, D. Patel, Assessing and augmenting SCADA cyber security: A survey of
23 techniques, *Comput. Secur.* 70 (2017) 436–454. doi:10.1016/j.cose.2017.06.010.
- 24 [194] H. Shih, W and Ludwig, The biggest challenges of data-driven manufacturing, *Harv. Bus. Rev.*
25 (2016). <https://hbr.org/2016/05/the-biggest-challenges-of-data-driven-manufacturing> (accessed
26 June 6, 2019).
- 27 [195] Allae Erraissi, Abdessamad Belangour, Abderrahim Tragha, Digging into Hadoop-based Big
28 Data Architectures, *Int. J. Comput. Sci. Issues IJCSI*. 14 (2017) 52–59.
29 doi:10.20943/01201706.5259.
- 30 [196] H. Fleischmann, S. Spreng, J. Kohl, D. Kisskalt, J. Franke, Distributed condition monitoring
31 systems in electric drives manufacturing, 2016 6th Int. Electr. Drives Prod. Conf. EDPC 2016
32 - Proc. (2016) 52–57. doi:10.1109/EDPC.2016.7851314.
- 33 [197] J. Manyika, M. Chui, J. Bughin, R. Dobbs, P. Bisson, Marrs, Disruptive technologies:
34 Advances that will transform life, business, and the global economy, McKinsey Glob. Insitute.
35 (2013) 163.
36 http://www.mckinsey.com/insights/business_technology/disruptive_technologies%5Cnhttp://www.chrysalixevc.com/pdfs/mckinsey_may2013.pdf.
37
- 38 [198] A. Burger, H. Koziolk, J. Rückert, M. Platenius-Mohr, G. Stomberg, Bottleneck
39 Identification and Performance Modeling of OPC UA Communication Models, (2019) 231–
40 242. doi:10.1145/3297663.3309670.
- 41 [199] Deutsche Kommission Elektrotechnik DKE; DIN e.V., German Standardization Roadmap
42 Industry 4.0, 2016.
- 43 [200] D. Loshin, Data Integration, *Bus. Intell.* (2013) 189–210. doi:10.1016/B978-0-12-385889-
44 4.00013-2.
- 45 [201] J.M. Gutierrez-Guerrero, J.A. Holgado-Terriza, IMMAS an industrial meta-model for
46 automation system using OPC UA, *Elektron. Ir Elektrotechnika*. 23 (2017) 3–11.

1 doi:10.5755/j01.eie.23.3.18324.

2 [202] Y.H. Wu, S. De Wang, L.J. Chen, C.J. Yu, Streaming analytics processing in manufacturing
3 performance monitoring and prediction, Proc. - 2017 IEEE Int. Conf. Big Data, Big Data 2017.
4 2018-Janua (2018) 3285–3288. doi:10.1109/BigData.2017.8258312.

5 [203] A. Stojadinović, N. Stojanović, L. Stojanović, Dynamic monitoring for improving worker
6 safety at the workplace, in: Proc. 9th ACM Int. Conf. Distrib. Event-Based Syst. - DEBS '15,
7 2015: pp. 205–216. doi:10.1145/2675743.2771881.

8 [204] S. Scholze, K. Nagorny, R. Siafaka, K. Krone, An Approach for Cloud-Based Situational
9 Analysis for Factories Providing Real-Time Reconfiguration Services, in: 2017: pp. 118–127.
10 doi:10.1007/978-3-319-65151-4.

11 [205] L. Angrisani, G. Ianniello, I. Elettrica, N. Federico, Cloud based system for measurement data
12 management in large scale electronic production, (n.d.).

13 [206] G. Hesse, B. Reissaus, C. Matthies, M. Lorenz, M. Kraus, M. Uflacker, Senska – Towards an
14 enterprise streaming benchmark, 2018. doi:10.1007/978-3-319-72401-0_3.

15 [207] L. Banica, A. Hagi, 4 - Using big data analytics to improve decision-making in apparel
16 supply chains A2 - Choi, Tsan-Ming BT - Information Systems for the Fashion and Apparel
17 Industry, in: Woodhead Publ. Ser. Text., Woodhead Publishing, 2016: pp. 63–95.
18 doi:https://doi.org/10.1016/B978-0-08-100571-2.00004-X.

19 [208] S. Saeidlou, M. Saadat, E.A. Sharifi, D. Guiovanni, An ontology-based intelligent data query
20 system in manufacturing networks, Prod. Manuf. Res. 3277 (2017) 1–18.
21 doi:10.1080/21693277.2017.1374887.

22

23 **7 APPENDICES**

	Search strings	Number of papers
Scopus	TITLE-ABS-KEY (("hadoop") OR ("spark") OR ("Storm") OR ("Flink") OR ("SQL") OR ("nosyl") OR ("time-series database") OR ("opc-ua") OR ("MTconnect")) AND (TITLE-ABS-KEY ("manufacturing") OR TITLE-ABS-KEY ("Industry 4.0") OR TITLE-ABS-KEY ("Smart manufacturing") OR TITLE-ABS-KEY ("Digital twin") OR TITLE-ABS-KEY ("Digital thread")) AND NOT (TITLE-ABS-KEY ("construction") OR TITLE-ABS-KEY ("healthcare") OR TITLE-ABS-KEY ("oil") OR TITLE-ABS-KEY ("energy") OR TITLE-ABS-KEY ("Agriculture"))	228
IEEE Xplore	(("Author Keywords":hadoop OR "Author Keywords":spark OR "Author Keywords":storm OR "Author Keywords":flink OR "Author Keywords":sql OR "Author Keywords":nosql OR "Author Keywords": "time-series database" OR "Author Keywords": "opc-ua" OR "Author Keywords": "mtconnect") AND ("Abstract":manufacturing OR "Abstract": "Industrial 4.0" OR "Abstract": "industrial automation" OR "Abstract": "smart manufacturing" OR "Abstract": "digital twin" OR "Abstract": "digital thread")	63
ASME	(hadoop, spark, storm, flink, sql, nosql, "time series database", "opc-ua", mtconnect) AND (manufacturing, "industry 4.0", "industrial automation", "smart manufacturing", "digital twin", "digital thread")	25

ACM ((keywords.author.keyword:(hadoop, spark, storm, flink) AND content.ftsec:(apache)) OR keywords.author.keyword:(sql, nosql, "time series database", opc-ua, mtconnect)) AND recordAbstract:(manufacturing, "industrial automation", "smart manufacturing", "digital twin", "digital thread")

22

1 Table 1. Search strings of four citation databases

Categories	Systems	Reference
Product	CAD/CAE/CAPP/CAM	[55][53][56][125]
	PLM	[58][59]
Production	ERP	[126][62][62][97][94]
	MOM/MES	[25][65][127][128][69][66][116][6]
	SCADA/DCS/HMI	[137][73][71][39][134][201][133][131][130][129]
	IIoT/CNC/Robot	[166][202][78][77][140]
	O&M	[138][24][83][139][161][156][157][142][82]
	QMS	[85][88][109][87]
	Safety	[89][203]
Business	SCM/CRM/BI/AM	[93][90][181]
ICT architecture	CPS/CM/ICT	[64][146][68][160][153][99][132][101][143][75][148][95][149][204][96][76]
	Data analytics/DM	[105][121][205][102][152][206][151][207][208]
	KM	[164][23][106]

2 Table 2. Category distribution of reviewed articles