

**Diagnostic information produces better calibrated judgments about forensic
comparison evidence than likelihood ratios**


Gianni Ribeiro*, Blake Malcolm McKimmie, Jason Marcus Tangen

School of Psychology, The University of Queensland

Author Note

Gianni Ribeiro  <https://orcid.org/0000-0002-2594-8311>

Blake Malcolm McKimmie  <https://orcid.org/0000-0002-2750-6111>

Jason Marcus Tangen  <https://orcid.org/0000-0002-0649-2566>

*Correspondence regarding this article should be addressed to Gianni Ribeiro,
School of Psychology, The University of Queensland, Saint Lucia QLD 4072. Email:
g.ribeiro@uq.edu.au

This article is scheduled for publication in the *Journal for Applied Research in Memory and Cognition* (accepted on 25/07/22).

© 2022, American Psychological Association. This paper is not the copy of record and may not exactly replicate the final, authoritative version of the article. Please do not copy or cite without authors' permission. The final article will be available, upon publication, via its DOI: 10.1037/mac0000062

DIAGNOSTIC INFORMATION VERSUS LIKELIHOOD RATIOS

Abstract

Forensic expert testimony is slowly starting to reflect the uncertain nature of forensic science, but the way experts should express the uncertainty of their decisions is under debate. Here we compare the likelihood model approach to a diagnostic approach — which provides information about performance and error rates — to determine which produces a more calibrated understanding and evaluation of the evidence. In Experiment 1 (N = 738), participants were more sensitive to differences in evidence strength when the evidence was expressed as diagnostic information than as a likelihood ratio, as predicted. In Experiment 2 (N = 499), however, when provided with both diagnostic information and a likelihood ratio, participants tended to discount the presence of the likelihood ratio in favour of the diagnostic information, which we did not predict. Together, these results suggest that providing fact-finders with diagnostic information might aid their understanding and evaluation of forensic evidence.

Keywords: forensic evidence; likelihood ratio; expert testimony; juror decision-making

DIAGNOSTIC INFORMATION VERSUS LIKELIHOOD RATIOS

General Audience Summary

Forensic experts were once allowed to testify that two samples (e.g., fingerprints) “match” a defendant to the exclusion of all other people and that their decisions are infallible.

However, numerous wrongful identifications and convictions have shown that these claims of infallibility simply are not true. These days, forensic testimony is slowly starting to reflect the uncertain nature of forensic decisions. One approach is known as the likelihood ratio, which conveys the probability of one hypothesis (that the fingerprints come from the same person) against another (that the fingerprints come from two different people). However, research has shown that jurors struggle to understand and interpret likelihood ratios. In this paper, we present a different approach — the diagnostic information approach — which presents jurors with information about the past performance of forensic decisions (how often forensic examiners’ decisions are correct and how often they are incorrect).

The first experiment we find that presenting expert testimony in this way results in better calibrated judgments compared to likelihood ratios. In the second experiment, we find that presenting both simultaneously does not lead to improved understanding. In conclusion, the diagnostic approach to presenting forensic testimony might be better understood, interpreted, and applied by jurors than likelihood ratios.

Diagnostic information produces better calibrated judgments about forensic comparison evidence than likelihood ratios

To date, the National Registry of Exonerations in the United States has documented 710 individuals who have been wrongfully convicted of crimes they did not commit on the basis of false or misleading forensic evidence, making it the fourth leading cause of wrongful convictions in their database (National Registry of Exonerations, 2022, January 31). While experts were once permitted to testify at trial that two prints, marks, or samples unequivocally “matched” the defendant, after numerous scathing reports about the state of forensic science (see, for example, the 2009 report from the National Academy of Science and the 2016 report from the President’s Council of Advisors on Science and Technology) the “match” terminology is no longer used their testimony is slowly starting to reflect the uncertain (and human) nature of their decisions. The NAS report (2009) noted that the terminology that forensic examiners use when testifying in court can have a profound effect on how the fact-finder perceives and evaluates that evidence (p. 21) and that terminology should be standardised and presented to the court in lay terms, so that testimony is understandable and interpretable by all trial participants (p. 186). By making these improvements to communicating forensic decisions, the NAS report noted that the number of wrongful convictions should be reduced, ensuring that innocent people will not be incarcerated for crimes they did not commit whilst true perpetrators of crime will not continue to commit crimes (pp. 4-5).

The likelihood ratio approach to communication forensic evidence

One approach to communicate the strength of forensic decisions is the likelihood ratio, which conveys the probability of an observation given two competing source hypotheses: (1) that the samples come from the same source, and (2) that the samples come from different sources. There is substantial support for the likelihood ratio approach (see Morrison’s introduction to the special issue on likelihood ratios in *Science & Justice*,

2016). However, others are more skeptical, suggesting that the focus on precise, large likelihood ratios — often in the order of millions or even billions — is futile when the probability of a false positive occurring due to human error (such as mislabelling samples or accidental contamination) is far greater (Jackson et al., 2015; Spellman, 2017; Eldridge, 2019). For example, in 2011 British man Adam Scott was charged with the rape of a woman in Manchester after private forensics company LGC concluded that his DNA matched the sample taken from the victim (Peachey, 2012). However, Scott claimed, and his phone records showed, that he was at home in Plymouth more than 300 kilometres away from the crime scene and had never been to Manchester in his life (Hall, 2017). After five months on remand in prison, an investigation found that Scott's DNA sample, which was taken for an unrelated incident, had contaminated the crime scene sample due to poor handling of the samples by LGC staff (Peachey, 2012).

Unfortunately, this is not an isolated incident of human error. A review of over 3,500 cases from Houston Crime Lab revealed that there were major issues in 32% of the DNA cases reviewed, including four inmates on death row (Bromwich, 2007). In 2008, the Detroit Crime Laboratory's firearms unit was shut down following a review which found that the unit complied with only 42% of the essential standards and found serious errors in 19 of a total 200 randomly selected firearms cases (Bunkley, 2008). In 2010, a review of North Carolina's state crime laboratory concluded that exculpatory evidence was either withheld or misrepresented in over 230 cases from 1987-2003 (Swecker & Wolf, 2010). In 2013, New York City's medical examiner's office undertook a review of over 800 rape cases over a period of 10 handled by a particular examiner who had consistently poor reviews and underperformance during her employment (Goldstein, 2013; Scott, 2013). Here in Australia, a miscode in the computer software used to analyse DNA samples affected the probability statistics in up to 60 Queensland cases. Further, the developers of

the program revealed that Queensland did not purchase the manual that accompanies the program, which may have resulted in human error (Murray, 2015).

While any approach to quantify the strength of evidence is a step in the right direction, if the ultimate fact-finder — the jury — is unable to properly understand, evaluate, and apply the information appropriately it could do more harm than good. A significant body of research in the forensic sciences as focused on how evidence is expressed in court, with a particular focus on random match probabilities (for a comprehensive review, see Eldridge, 2019). The overwhelming majority of studies find that, despite being numerically equivalent, fact-finders come to different ultimate decisions based on how the evidence is expressed to them (e.g., single-target probability versus multi-target frequency expressions in Koehler, 2001).

Research also suggests that fact-finders struggle to understand and interpret evidence when expressed as likelihood ratios (Martire et al., 2013; 2014).

Recommendations have been made to express likelihood ratios verbally rather than numerically (De Kinder & Olsson, 2011; Association of Forensic Science Providers, 2009), perhaps because many forensic sciences lack the data necessary to compute a numerical likelihood ratio (de Keijser & Elffers, 2012). However, we need to consider whether the recipients of likelihood ratio expressions, the jury, actually interpret these numerical and verbal expressions the way that they were intended (Campbell, 2011). Research regarding probabilities and statistics in general suggests that people have difficulty understanding and converting numerical and verbal expressions (Gigerenzer & Edwards, 2003; Kahneman et al., 1982). In addition, individuals may interpret the same verbal expressions of probability differently, as the meanings attributed to words may vary between people and contexts (Budescu et al., 2009; Wallsten & Budescu, 1995; Shaw & Dear, 1990). In Martire and colleagues (2013) study, the authors manipulated whether the expert's opinion was expressed as either a verbal likelihood ratio or a numeric likelihood ratio and also

varied the strength of the evidence (low, moderate, or high). They found that when the evidence strength was moderate or high, numerical and verbal expressions resulted in equal amounts of belief change. However, when evidence strength was low, there was significantly greater belief change for the numerical expression than the verbal expression. Interestingly, in the low-strength numerical condition, belief change was in the intended direction (i.e., towards guilt), however in the low-strength verbal condition, belief change was in the opposite direction than intended (i.e., towards innocence). This finding clearly contradicts the recommendation to express likelihood ratios verbally rather than numerically. Even when both verbal and numerical likelihood ratio expressions are presented in table format amongst the complete range of possible expressions, Martire and colleagues (2014) found that belief change was in the opposite direction (i.e. towards innocence) than what the expert had intended.

Given the evidence from Martire and colleagues' (2013, 2014) studies that fact-finders struggle to understand and interpret likelihood ratios, perhaps there is an alternative approach that may help fact-finders better understand and evaluate forensic evidence. Some have suggested that informing jurors about the accuracy and error rates of forensic decisions may help them to understand and weigh the evidence more appropriately (Mnookin, 2010; Garrett & Mitchell, 2013; Howes, 2015).

The diagnostic information approach to communicating forensic evidence

Like decisions made about forensic comparison evidence, many medical decisions are also binary in nature, involve human judgment, and have varying degrees of error. Similarly, the ground truth of any *particular* decision is unknowable, so the best indicator is experts' past performance on similar tasks. Tomlinson, Marshall, and Ellis (2008) conducted an experiment to determine the accuracy of a home pregnancy test (*Answer*TM) by administering the test to 120 pregnant women and 120 women who were not pregnant. The results are summarised in Figure 1A below using natural frequencies which are

demonstrably more concrete and easier to grasp than probabilistic information (see, for example, Gigerenzer & Hoffrage, 1995; Hoffrage & Gigerenzer, 1998; Akl et al., 2011).

While a woman who has just taken the *Answer*TM pregnancy test cannot know, on the basis of this information alone, with absolute certainty whether or not she is pregnant, she can use this information to guide her interpretation of the result of her own test.

Figure 1

A) Results from Tomlinson, Marshall, & Ellis (2008). B) Results from Tangen, Thompson, & McCarthy (2011).

A			B		
	Test said "pregnant"	Test said "not pregnant"		Expert said "match"	Expert said "no match"
Pregnant	98	22	120	409	35
Not pregnant	2	118	120	3	441
	100	140		412	476
			Same source		444
			Different source		444

The same diagnostic information approach can be applied to forensic evidence. A juror cannot know, with absolute certainty, whether or not a particular fingerprint found at a crime scene belongs to a defendant or not. But with information about the accuracy of fingerprint examiners' decisions from studies like Tangen and colleagues (2011) summarised in Figure 1B above, a juror can better evaluate a forensic examiner's testimony in order to guide their interpretation of whether or not the fingerprint is likely to have come from the defendant or some other person. Taking a diagnostic information approach to forensic evidence would provide information, such as performance and error rates on previous, similar decisions, to help jurors reason about the present case. This process is often referred to as group-to-individual inference (see Faigman, Monahan, & Slobogin, 2014). Jurors may use this general performance information to help them weigh the evidence in a particular case and adjust their degree of belief accordingly.

Should Forensic Examiners Present Both Likelihood Ratios and Diagnostic Information Simultaneously?

Authoritative bodies, such as PCAST (2016) and the NRC (1992) have suggested that error rate estimates be provided alongside probabilistic testimony such as a random match probability or likelihood ratio. However, because error rates derived from proficiency tests provide the accuracy of an examiner in making a binary decision, while a random match probability provides the strength of the evidence without requiring a binary decision, jurors may struggle to understand and combine these two types of information. In fact, studies have demonstrated that people have difficulty combining probabilities (see, for example, the conjunction fallacy in Tversky & Kahneman, 1983). This has been empirically demonstrated by Koehler and colleagues (1995), who found that the presence of a laboratory error rate (0.1% or 2%) had no impact on verdicts. Further, as previously discussed, because the chance of human error is far greater than the chance of a random match (or the inverse, a likelihood ratio), the error rate would effectively swamp the other statistic, rendering it futile.

In two preregistered experiments, we compared the likelihood ratio approach to a diagnostic information approach, which included information about accuracy and errors of fingerprint examiners' decisions. The aim of Experiment 1 was to determine whether diagnostic information would result in a more calibrated understanding and evaluation of the fingerprint evidence than a likelihood ratio, whereas the aim of Experiment 2 was to determine how presenting both a likelihood ratio and diagnostic information together would affect participants' evaluation of the evidence.

Experiment 1

The aim of Experiment 1 was to directly compare the likelihood ratio approach to the diagnostic information approach using a 2 (Evidence Expression: Likelihood Ratio vs Diagnostic Information) by 5 (Evidence Strength: 5 levels) between-groups design.

Participants

This experiment was preregistered prior to data collection (see: <http://osf.io/5xmk8>). Both experiments reported in this paper received ethical approval from The University of Queensland's Low & Negligible Risk Ethics Sub-Committee (Clearance Number: 2018000922). We recruited 740 participants for this experiment, where sample size was determined a priori as having 95% power to detect a small-to-medium effect ($d = 0.4$) at an alpha level of .05. We recruited participants residing in Australia or England from the *Prolific* online testing platform who received £1.25 for their participation.

Following our preregistered exclusion plan, we excluded two participants for failing, or not responding to, the attention check question. The final sample ($N = 738$) consisted of 501 women, 231 men, three participants who identified as another gender, and three who did not report their gender with an age range of 18 – 88 years ($M = 37.59$, $SD = 12.32$). Most (99.5%) participants reported living in England, however two participants did not report whether they lived in England or Australia. Other participant demographic information, such as highest level of education and prior jury service, can be viewed in the supplemental materials.

Procedure and Materials

Ethical clearance for this study Participants completed the experiment on their own devices. They first read a short case vignette detailing the facts of a murder and were told that a latent fingerprint was retrieved from the victim's windowsill. Participants who were randomly assigned to the Likelihood Ratio conditions received information about the strength of evidence in a likelihood ratio in one of five strengths: 5, 55, 550, 5,500, or 550,000. These values were selected as they fall in the middle of each of the categories in the Association of Forensic Science Providers' (2009) *Standards for Numerical and Verbal Expression of Likelihood Ratios*. Below is an example of the highest strength likelihood ratio condition:

A man with the initials B.M. is alleged to have murdered the victim in the case you just read about.

During the defendant's trial, a forensic fingerprint examiner testified about the fingerprint retrieved from the victim's window sill. The forensic examiner compared the fingerprint retrieved from the victim's window sill with a fingerprint obtained from the defendant, B.M.

When assessing the significance of any similarity or differences between two fingerprints, the likelihood of obtaining that similarity or difference is considered against two alternative propositions: (Hypothesis 1) the two fingerprints originated from the same person; (Hypothesis 2) the two fingerprints did not originate from the same person.

The forensic examiner testified that the correspondence between the fingerprint retrieved from the victim's window sill and the fingerprint obtained from the defendant, B.M., is 5,500,000 times more likely if the two fingerprints originated from the same person (Hypothesis 1) than if the two fingerprints originated from different people (Hypothesis 2).

Whereas, participants randomly assigned to the Diagnostic Information conditions received information about how fingerprint examiners perform. This information was adapted from Edmond, Tangen, and Thompson's (2014) *Guide to Interpreting Forensic Testimony*. This information was presented in one of five strengths: 50%, 63%, 75%, 88%, 99%. We chose these strengths as they represent equal intervals between chance and (near) perfect performance. Below is an example of the highest strength diagnostic information condition:

A man with the initials B.M. is alleged to have murdered the victim in the case you just read about.

During the defendant's trial, a forensic fingerprint examiner testified about the fingerprint retrieved from the victim's window sill. The examiner testified that the fingerprint retrieved from the victim's window sill matched the defendant, B.M.'s, fingerprint.

A decision about whether a fingerprint found at a crime scene matches a person's fingerprint or not is based on decisions by a human examiner, not a computer.

In laboratory-based experiments, a group of fingerprint examiners were tested on whether they could tell the difference between two fingerprints from the same person and two fingerprints from different people. They repeated this test 200 times.

For the 100 pairs of fingerprints that matched, the examiners correctly said "match" for 99, but incorrectly said "no match" for 1.

For the 100 pairs of fingerprints that did not match, the examiners correctly said "no match" for 99, but incorrectly said "match" for 1.

After reading the case vignette and the evidence information, participants were asked to respond to three key dependent variables: source likelihood (On a scale of 0-100, how likely is it that the DNA evidence retrieved from underneath the victim's fingernails belongs to B.M.?), guilt likelihood (On a scale of 0-100, how likely is it that B.M. committed this murder?), and a dichotomous verdict judgment (not guilty vs. guilty). As source and guilt likelihood measures were consistent across Experiments 1 and 2, these are reported here, whereas results for the dichotomous guilt judgment can be found in the Supplementary Materials. Finally, participants responded to one attention check question, three manipulation check questions, and five demographic questions (see the full set of experimental materials on the Open Science Framework).

Hypotheses

We predicted a main effect of how the evidence was expressed, such that participants would be more likely to conclude that the fingerprint came from the defendant, that the defendant committed the murder, and that the defendant was guilty of murder when the evidence was expressed as a likelihood ratio compared to diagnostic information (H1). We also predicted a main effect of Strength, such that participants would be more likely to conclude that the fingerprint came from the defendant, that the defendant committed the murder, and that the defendant was guilty of murder as the strength of the evidence increases (H2). Finally, we predicted an Expression by Strength interaction, such that Strength will have a greater effect on the dependent variables when the evidence is expressed as Diagnostic Information compared to Likelihood Ratio (H3).

Results

Overview of Analyses

To test our predictions, in line with our preregistration we conducted 2 (Evidence Expression: Likelihood Ratio vs Diagnostic Information) by 5 (Evidence Strength: 5 levels) between-groups factorial ANOVAs on both source and guilt likelihood. As the data violated

the assumptions of normality and homogeneity of variance, we also conducted robust ANOVAs in *R* using 20% trimmed means (Wilcox, 2012). Unless otherwise stated, the results of both approaches were the same, thus we will report the original ANOVA results here for ease of interpretation.

Manipulation Checks

92.4% of participants who received evidence expressed as a Likelihood Ratio correctly identified the strength of the likelihood ratio that they received. Similarly, 99.5% of participants who received evidence expressed as diagnostic information correctly identified the strength of the diagnostic information that they received. As per the preregistered analysis plan, all participants were retained in the main analysis.

Source Likelihood

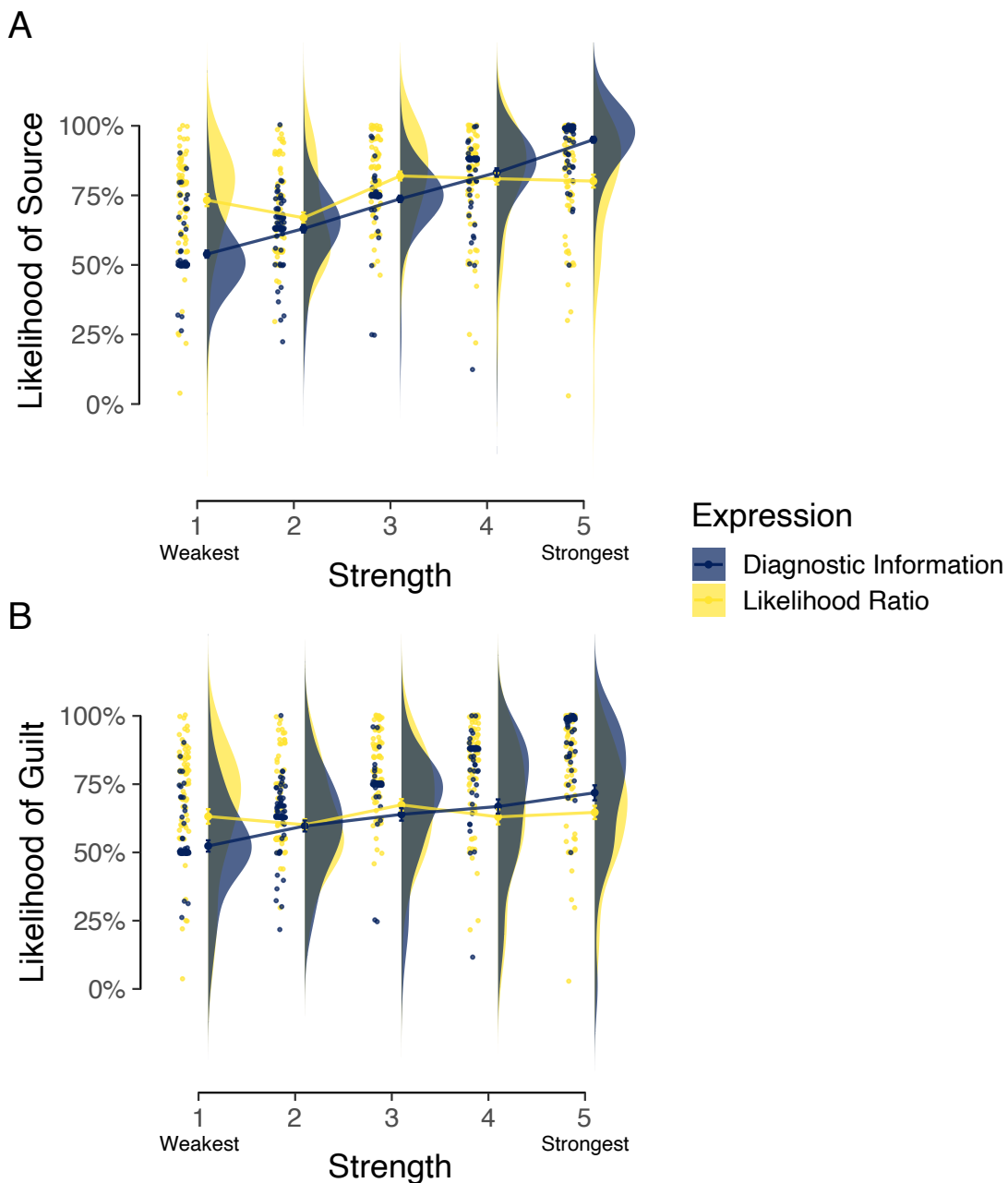
As illustrated in Figure 2 below, consistent with H1 there was a significant main effect of Expression, $F(1,726) = 7.144$, $p = .008$, $\eta^2_p = .010$. Participants who received the evidence expressed as a Likelihood Ratio ($M = 76.65$, $SD = 18.28$) were more likely to conclude that the defendant was the source of the fingerprint than participants who received the evidence expressed as Diagnostic Information ($M = 73.89$, $SD = 18.18$). Furthermore, consistent with H2 we found a significant main effect of Evidence Strength, $F(4,726) = 76.99$, $p < .001$, $\eta^2_p = .298$. Polynomial planned contrasts revealed a significant linear trend, $C = 20.59$, $p < .001$, showing that participants were more likely to conclude that the defendant was the source of the fingerprint as the strength of the evidence increased.

Finally, there was a significant Expression \times Strength interaction, $F(4,726) = 27.79$, $p < .001$, $\eta^2_p = .133$, where the strength of the expression influenced the likelihood that the defendant was the source of the fingerprint for both the Likelihood Ratio, $F(4,726) = 14.20$, $p < .001$, and Diagnostic Information expressions, $F(4,726) = 89.92$, $p < .001$. A planned linear contrast revealed that the positive linear trend was greater for the

Diagnostic expression, $C = 102.46$, $p < .001$, compared to the Likelihood Ratio expression, $C = 27.76$, $p < .001$, demonstrating that participants are more sensitive to the strength of the evidence when expressed diagnostically rather than as a likelihood ratio, which was consistent with H3.

Figure 2

Panels A and B depict participants' source and guilt likelihood ratings, respectively. The raincloud plots depict a half violin plot of participants' mean ratings overlaid with jittered data points from each of the participants who were randomly assigned to one of five evidence strengths (1 = weakest, 5 = strongest) along with the standard error of the mean per condition.



Guilt Likelihood

There was no main effect of Expression, $F(1,731) = 0.25$, $p = .619$, $\eta^2_p = <.001$, which we did not predict (H1). That is, participants were equally likely to conclude that the defendant committed the murder when the evidence was expressed as a Likelihood Ratio ($M = 63.67$, $SD = 20.85$) or using Diagnostic Information ($M = 62.87$, $SD = 21.08$).

Consistent with H2, there was a main effect of Evidence Strength, $F(4,731) = 6.57$, $p = <.001$, $\eta^2_p = .035$. Polynomial planned contrasts revealed a significant linear trend, $C = 8.21$, $p < .001$, where participants were more likely to conclude that the defendant committed the murder as the strength of the evidence increased.

Finally, there was a significant Expression \times Strength interaction, $F(4,731) = 4.24$, $p = .002$, $\eta^2_p = .023$. The strength of the evidence influenced the likelihood that the defendant committed the crime when expressed as Diagnostic information, $F(4,722) = 9.51$, $p < .001$, but not as a Likelihood Ratio, $F(4,722) = 1.22$, $p = .299$. A planned linear contrast revealed that there was a positive linear trend for the Diagnostic expression, $C = 46.07$, $p < .001$, but not the Likelihood Ratio expression, $C = 5.82$, $p = .439$, which was consistent with H3.

Experiment 2

In Experiment 1, we found that participants were more sensitive to the strength of the evidence when it was expressed as diagnostic information rather than as a likelihood ratio. What if participants were presented with both a likelihood ratio and diagnostic information simultaneously? As mentioned previously, providing a large likelihood ratio may not be all that useful given that the chance of a human error occurring is far greater (Jackson et al., 2015; Spellman, 2017; Eldridge, 2019). Thus, ideally, participants would disregard the likelihood ratio in favour of the diagnostic information. However, studies have demonstrated that people have difficulty combining probabilities (see, for example, the conjunction fallacy in Tversky & Kahneman, 1983). Further, as likelihood ratios may

provide the “aura of precision” (Tribe, 1971) and large likelihood ratios may be particularly compelling, it may be more likely that participants would rely more on the likelihood ratio than the diagnostic information. The aim of Experiment 2 was to determine how participants would evaluate the fingerprint evidence when presented with both a likelihood ratio as well as diagnostic information using a 2 (Diagnostic Information: Low vs High) by 2 (Likelihood Ratio: Absent vs Present) + 1 (Likelihood Ratio Only) between-groups design.

Participants

This experiment was preregistered prior to data collection (see: <http://osf.io/n27p8>). We recruited 500 participants for this experiment, where sample size was determined a priori as having 95% power to detect a small-to-medium effect ($d = 0.4$) at an alpha level of .05. We recruited participants residing in Australia or England from the *Prolific* online testing platform who agreed to receive £1.25 for their participation.

Following our preregistered exclusion plan, we excluded one participant for failing the attention check question. The final sample ($N = 499$) consisted of 312 women, 185 men, and two participants who did not report their gender, with an age range of 18 – 72 years ($M = 35.43$, $SD = 13.10$). 99.4% of participants reported living in England, 0.4% of participants reported living in Australia, and one participant did not report whether they lived in England or Australia. Other participant demographic information, such as highest level of education and prior jury service, can be viewed in the supplemental materials.

Procedure and Materials

Participants completed the experiment on their own devices. Participants first read a short case vignette (same as Experiment 1). Participants were then randomly assigned to one of five experimental conditions. In the Likelihood Ratio Present condition, participants received a likelihood ratio value of 5,500,000 — which falls in the “extremely strong” category of the Association of Forensic Science Providers’ (2009) *Standards for Numerical and Verbal Expression of Likelihood Ratios* — whereas, participants in the

Likelihood Ratio Absent condition received no likelihood ratio information at all.

Participants in the Diagnostic Information condition received information about fingerprint examiners' performance, presented similarly to that in Experiment 1. In the Low condition, this information conveyed an overall accuracy of 63%, whereas in the High condition, this information conveyed an overall accuracy of 99%. Finally, participants in the Likelihood Ratio Only condition received only the likelihood ratio information, but no diagnostic information.

After reading the scenario and the evidence information, participants were asked to respond to three key dependent variables: source likelihood (On a scale of 0-100, how likely is it that the fingerprint retrieved from the victim's window sill belongs to B.M.?), guilt likelihood (On a scale of 0-100, how likely is it that B.M. committed this murder?), and error likelihood (On a scale of 0-100, how likely is it that there could have been an error during the analysis and comparison of B.M.'s fingerprint to the fingerprint retrieved from the victim's window sill?). As source and guilt likelihood measures were consistent across Experiments 1 and 2, these are reported here, whereas results for the error likelihood measure can be found in the Supplementary Materials. Finally, as in Experiment 1, participants responded to one attention check question, three manipulation check questions, and five demographic questions (see full experimental materials on the Open Science Framework).

Hypotheses

We predicted a main effect of Likelihood Ratio, such that participants would be more likely to conclude that the evidence came from the defendant and that the defendant committed the murder, but less likely to conclude that an error could have occurred, when a Likelihood Ratio was Present compared to Absent (H1). We also predicted a main effect of Diagnostic Information, such that participants would be more likely to conclude that the evidence came from the defendant and that the defendant committed the murder, but less

likely to conclude that an error could have occurred, when Diagnostic Information was High compared to Low (H2). Further, we predicted a Likelihood Ratio by Diagnostic Information interaction, such that we expect Diagnostic Information to have a greater effect on the dependent variables when Likelihood Ratio is Absent compared to Present (H3). Finally, if participants are not sensitive to the diagnostic information and are instead focused solely on the likelihood ratio information, then the floating control condition (Likelihood Ratio Only) should only differ significantly from the conditions where a Likelihood Ratio is absent and only diagnostic information is presented, but not when a Likelihood Ratio is present alongside diagnostic information (H4).

Results

Overview of Analyses

To test our predictions, we conducted 2 (Diagnostic Information: Low vs High) \times 2 (Likelihood Ratio: Absent vs Present) between-groups factorial ANOVAs on each of our dependent measures: source, guilt, and error likelihood. As the data violated the assumptions of normality and homogeneity of variance, we also conducted robust ANOVAs in *R* using 20% trimmed means (Wilcox 2012). Unless otherwise stated, the pattern of results of both approaches were the same, thus we will report the original ANOVA results here for ease of interpretation. We also conducted a one-way between-groups ANOVA on the three Likelihood Ratio conditions: Likelihood Ratio with high and low diagnostic information, and Likelihood Ratio only.

Manipulation Checks

In the Likelihood Ratio Absent condition, 87.88% of participants correctly identified that they did not receive a likelihood ratio, whereas 88.33% of participants in the Present condition correctly identified that they received a likelihood ratio of 5,500,000. For the diagnostic information manipulation, 94% of participants in the Low diagnostic condition and 97.9% of participants in the High diagnostic condition correctly identified the level of

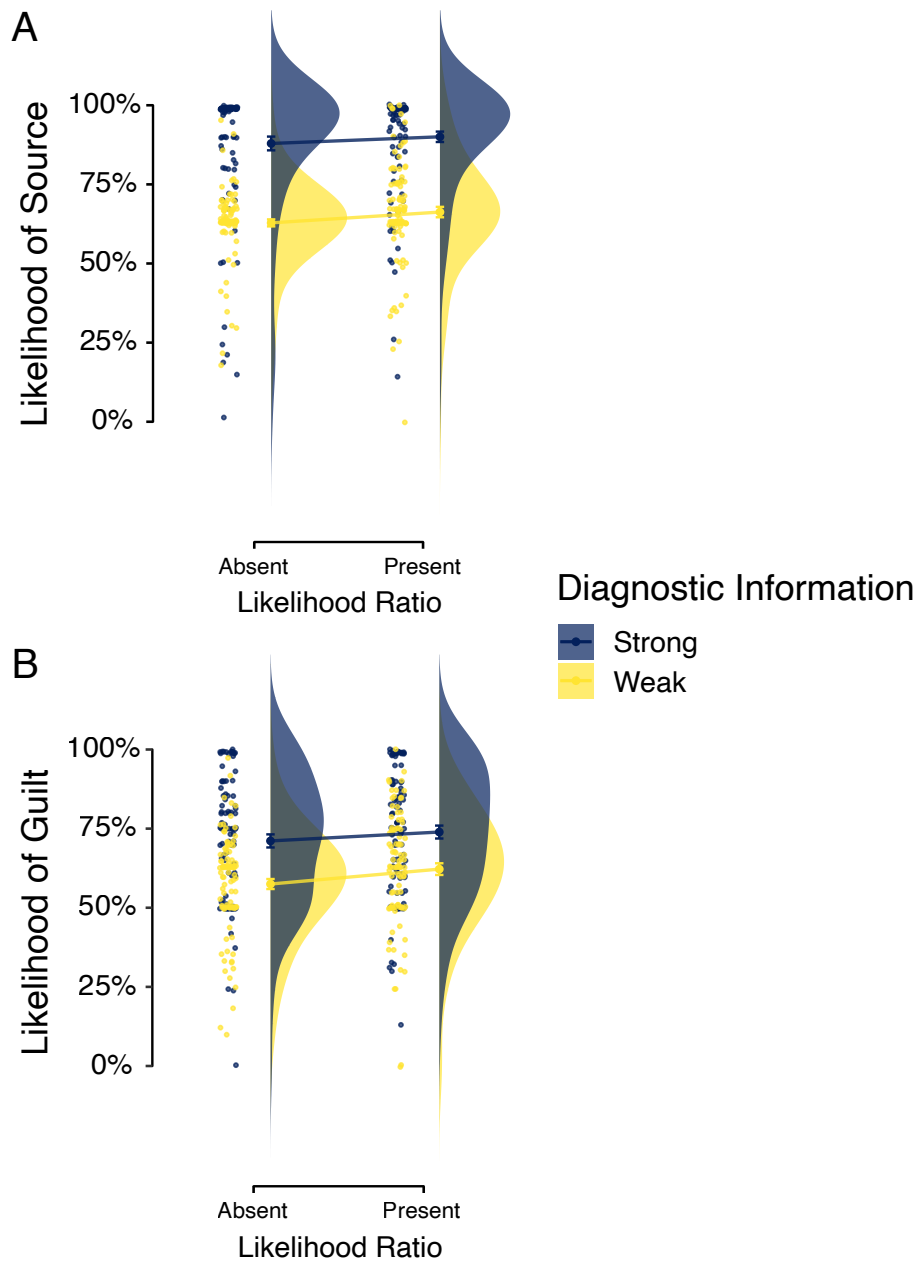
diagnostic information they were presented in the vignette. As per the preregistered analysis plan, all participants were retained in the main analysis.

Source Likelihood

As illustrated in Figure 3 below, there was no main effect of Likelihood Ratio, $F(1, 395) = 2.71, p = .10, \eta^2_p = .007$, therefore participants were equally likely to conclude that the defendant was the source of the fingerprint when a Likelihood Ratio was Present ($M = 78.26, SD = 20.27$) or Absent ($M = 75.19, SD = 21.14$), which we did not predict (H1). However, consistent with H2, there was a main effect of Diagnostic Information, $F(1, 395) = 212.54, p < .001, \eta^2_p = .350$. Participants were more likely to conclude that the defendant was the source of the fingerprint when Diagnostic Information was High ($M = 89.00, SD = 19.10$) compared to Low ($M = 64.53, SD = 14.05$). Contrary to our prediction (H3), there was no significant Likelihood Ratio \times Diagnostic Information interaction, $F(1, 395) = .155, p = .694, \eta^2_p < .001$, indicating that the effect of diagnostic information on participants' judgments of whether the defendant was the source of the fingerprint did not differ depending on whether participants received a likelihood ratio or not.

Figure 3

Panels A and B depict participants' source and guilt likelihood ratings, respectively. The raincloud plots depict a half violin plot of participants' mean ratings overlaid with jittered data points from each of the participants along with the standard error of the mean per condition.



Guilt Likelihood

Consistent with H1, there was a significant main effect of Likelihood Ratio, $F(1, 386) = 3.94$, $p = .048$, $\eta^2_p = .010$, therefore participants were slightly more likely to conclude that

the defendant committed the murder when a Likelihood Ratio was Present ($M = 68.03$, $SD = 20.03$) or Absent ($M = 64.20$, $SD = 19.33$). However, this effect was not significant with a robust ANOVA, $Q = 346$, $p = .065$. Furthermore, consistent with H2, there was a main effect of Diagnostic Information, $F(1, 386) = 44.86$, $p < .001$, $\eta^2_p = .104$, where participants were more likely to conclude that the defendant committed the murder when Diagnostic Information was High ($M = 72.52$, $SD = 20.19$) compared to Low ($M = 59.84$, $SD = 17.18$). However, contrary to our prediction (H3), there was no significant Likelihood Ratio \times Diagnostic Information interaction, $F(1, 386) = .254$, $p = .615$, $\eta^2_p < .001$, indicating that the effect of diagnostic information on participants' judgments of whether the defendant committed the murder did not differ depending on whether participants received a likelihood ratio or not.

Comparing Likelihood Ratio Conditions

Finally, we compared the floating control condition, Likelihood Ratio Only, to the other conditions in the experimental design. A one-way ANOVA comparing all five experimental conditions revealed a significant difference between conditions on the measure of source likelihood, $F(4, 498) = 55.79$, $p < .001$. As per H4, we hypothesised that the Likelihood Ratio Only condition would only differ from conditions where a Likelihood Ratio was absent (i.e., diagnostic information only) compared to present. However, contrary to our hypothesis, post hoc multiple comparisons using Tukey's range test revealed that participants in the Likelihood Ratio Only condition ($M = 83.67$, $SD = 17.93$) were significantly more likely to conclude that the defendant was the source of the evidence than participants who received a likelihood ratio with weak diagnostic information ($M = 66.25$, $SD = 16.38$), $p < .001$, and participants who only received weak diagnostic information ($M = 62.83$, $SD = 11.12$), $p < .001$, but not participants who received a likelihood ratio with strong diagnostic information ($M = 90.03$, $SD = 16.49$), $p = .063$, or

participants who received only strong diagnostic information ($M = 83.67$, $SD = 17.93$), $p = .396$.

A similar pattern of results was found for the measure of guilt likelihood, $F(4, 487) = 12.83$, $p < .001$. Post hoc multiple comparisons using Tukey's range test revealed that participants in the Likelihood Ratio Only condition ($M = 70.36$, $SD = 20.35$) were significantly more likely to conclude that the defendant was the source of the evidence than participants who received a likelihood ratio with weak diagnostic information ($M = 62.20$, $SD = 18.28$), $p = .024$, and participants who only received weak diagnostic information ($M = 57.50$, $SD = 15.76$), $p < .001$, but not participants who received a likelihood ratio with strong diagnostic information ($M = 73.92$, $SD = 20.08$), $p = .687$, or participants who received only strong diagnostic information ($M = 71.11$, $SD = 20.30$), $p = .999$. Together, these findings demonstrate that participants' conclusions about the likelihood that the fingerprint belongs to the defendant and the likelihood that the defendant committed the crime is driven largely by the diagnostic information and not the likelihood ratio.

Discussion

There is wide support for forensic examiners adopting likelihood ratios to communicate the uncertainty of forensic decisions (Morrison, 2016), however research has shown that mock jurors struggle to understand and interpret likelihood ratios (Martire et al., 2013; 2014). The aim of the current paper was to compare the use of likelihood ratios to a diagnostic information approach, where jurors would be presented with information regarding forensic examiner's past performance — the number of correct and incorrect decisions — in natural frequencies, which are demonstrably better understood than probabilistic information (see, for example, Gigerenzer & Hoffrage, 1995; Hoffrage & Gigerenzer, 1998; Akl et al., 2011).

Experiment 1 compared information about the strength of evidence expressed either as a likelihood ratio (which communicates the strength of the observations between the two fingerprints) or as diagnostic information (which communicates the strength of fingerprint examiners' conclusions through information about examiners' performance). We predicted that participants would be more sensitive to the strength of the evidence when it was presented diagnostically, including information about accuracy and error, rather than as a likelihood ratio. This hypothesis was supported, as participants' judgments about whether the defendant was the source of the evidence and committed the murder increased linearly as the strength of the evidence increased in the diagnostic expression condition, but not the likelihood ratio expression condition. Participants' insensitivity to the strength of the evidence when presented as a likelihood ratio is consistent with Martire and colleagues' (2013; 2014) findings that jurors struggle to interpret the strength of a likelihood ratio, even when it is presented in the context of a table where the full range of possible likelihood ratio values were presented. Our findings suggest that expressing the strength of forensic evidence diagnostically, rather than as a likelihood ratio, may result in greater sensitivity to the strength of the evidence by fact-finders and aid their interpretation. It is important to note that we do not intend to claim that the likelihood ratios and diagnostic information presented for each strength are equivalent (e.g., a likelihood ratio of 5 is not equivalent to an overall accuracy of 50% in the lowest strength condition), hence why we conducted analyses *within* each expression type rather than *between* expression types.

In Experiment 2, we wanted to determine whether providing participants with a likelihood ratio in addition to the diagnostic information would help or hinder their judgments. As the chance of an error having occurred outweighs the strong likelihood ratio, participants should ideally disregard the likelihood ratio in favour of the diagnostic information. However, previous research suggests that people combine probabilities poorly

(e.g., Tversky & Kahneman, 1983) and, as the strong likelihood ratio may seem overly compelling and precise, participants may favour the likelihood ratio over the diagnostic information. Thus, we predicted that participants would only be sensitive to the diagnostic information when it was presented alone rather than presented alongside a likelihood ratio, which the data did not support. The only effect was a medium-to-large sized main effect of diagnostic information, where participants were more likely to conclude that the defendant was the source of the evidence and committed the murder when the diagnostic information was high compared to low, regardless of whether the diagnostic information was accompanied by a likelihood ratio or not. Thus, there appears to be no added benefit to providing participants with a likelihood ratio in addition to the diagnostic information, but it does not appear to hinder their ability to interpret the diagnostic information.

These results appear to contradict the existing literature which suggests that people have difficulty combining probabilities. For example, Koehler and colleagues (1995) found that the presence of a random match probability, whether accompanied by an error rate or not, increased participants' guilty verdicts and that the error rate information had no impact on participants' verdicts. Whereas, the results of our experiment are the opposite: the presence of a likelihood ratio had no effect, but the error rate information did. One possible explanation for this difference is that Koehler and colleagues' (1995) study presented error rate information very briefly; participants were simply told that "the probability that a human error occurred which may have led to an incorrect report was about 1 in 1,000" (p. 219) or "1 in 50." Whereas, in our experiment, participants were told about a test that fingerprint examiners have taken involving pairs of known matching and non-matching prints and are provided with the results of that test, indicating their overall accuracy and error rate. It may be that providing participants with this added information about how the error rate was derived makes it easier to understand and, in turn, leads to increased sensitivity to the strength of the evidence.

Further, Koehler and colleagues' (1995) study used a random match probability as the evidence statistic, compared to a likelihood ratio in our experiment. While a likelihood ratio is the inverse of a random match probability, it is possible that participants may interpret these statistics differently. Previous research seems to support this argument, as studies have found that people give more weight to the evidence when the random match probability is low compared to high (e.g., Goodman, 1992; Smith et al, 1996) whereas people do not appear to be sensitive to the value of the likelihood ratio (Martire et al., 2013, 2014). It is also worth that we did not assess participants' jury eligibility in their country of residence, thus we could not exclude any jury ineligible participants which may affect the ecological validity of our findings.

Taken together, the results of both experiments suggest that people are more sensitive to the strength of evidence when presented with diagnostic information rather than a likelihood ratio and, when presented with both, rely primarily on the diagnostic information to inform their judgments. If this is the case, then forensic examiners should strongly consider adopting a diagnostic approach to communicate the strength of their evidence.

References

- Akl, E. A., Oxman, A. D., Herrin, J., Vist, G. E., Terrenato, I., Sperati, F., Costiniuk, C., Blank, D., & Schünemann, H. (2011, March 16). Using different statistical formats for presenting health information. Retrieved from: https://www.cochrane.org/CD006776/COMMUN_using-different-statistical-formats-for-presenting-health-information
- Association of Forensic Science Providers. (2009). Standards for the formulation of evaluative forensic science expert opinion. *Science & Justice*, 49(3), 161-164. doi:10.1016/j.scijus.2009.07.004
- Edmond, G., Thompson, M. B., & Tangen, J. M. (2014). A guide to interpreting forensic testimony: Scientific approaches to fingerprint evidence. *Law, Probability, and Risk*, 13(1), 1-25. doi: 10.1093/lpr/mgt011
- Faigman, D. L., Monahan, J., & Slobogin, C. (2014). Group to individual (G2i) inference in scientific expert testimony. *The University of Chicago Law Review*, 81(2), 417-480.
- Gigerenzer, G., & Hoffrage, U. (1995). How to improve bayesian reasoning without instruction: Frequency formats. *Psychological Review*, 102(4), 684–704. doi: 10.1037/0033-295X.102.4.684
- Hall, L. (2017, February 5). *DNA profiling is not infallible and can lead to innocent people being wrongly convicted: Report*. <https://www.smh.com.au/national/nsw/dna-profiling-is-not-infallible-and-can-lead-to-innocent-people-being-wrongly-convicted-report-20170202-gu3rwe.html>
- Hoffrage, U., & Gigerenzer, G. (1998). Using natural frequencies to improve diagnostic inferences. *Academic Medicine*, 73(5), 538–540. doi: 10.1097/00001888-199805000-00024
- Koehler, J. (2011). If the shoe fits they might acquit: The value of forensic science testimony. *Journal of Empirical Legal Studies*, 8, 21–48.

- Martire, K. A., Kemp, R. I., Sayle, M., & Newell, B. R. (2014). On the interpretation of likelihood ratios in forensic science evidence: presentation formats and the weak evidence effect. *Forensic Science International*, 240(61). doi:10.1016/j.forsciint.2014.04.005
- Martire, K. A., Kemp, R. I., Watkins, I., Sayle, M. A., & Newell, B. R. (2013). The expression and interpretation of uncertain forensic science evidence: Verbal equivalence, evidence strength, and the weak evidence effect. *Law and Human Behavior*, 37(3), 197-207. doi:10.1037/lhb0000027
- Morrison, G. S. (2016). Special issue on measuring and reporting the precision of forensic likelihood ratios: Introduction to the debate. *Science & Justice*, 56(5), 371-373. doi: 10.1016/j.scijus.2016.05.002
- National Registry of Exonerations. (2022, January 31). % exonerations by contributing factor. Retrieved from: <https://www.law.umich.edu/special/exoneration/Pages/ExonerationsContribFactorsByCrime.aspx>
- National Academy of Science (2009). *Strengthening forensic science in the United states: A path forward*. Washington, DC: The National Academies Press. Retrieved from: <https://www.ncjrs.gov/pdffiles1/nij/grants/228091.pdf>
- Peachey, P. (2012, October 2). *Rape accused Adam Scott was victim of forensics error, regulator finds*. <https://www.independent.co.uk/news/uk/crime/rape-accused-adam-scott-was-victim-of-forensics-error-regulator-finds-8193163.html>
- President's Council of Advisor's on Science and Technology (2016). *Forensic science in criminal courts: Ensuring scientific validity of feature-comparison methods*. Retrieved from: https://obamawhitehouse.archives.gov/sites/default/files/microsites/ostp/PCAST/pcast_forensic_science_report_final.pdf
- Ribeiro, G., McKimmie, B. M., & Tangen, J. M. (2022). *Comparing the use of likelihood ratios with diagnostic information*. Retrieved from: <https://osf.io/5xmk8/>

Ribeiro, G., McKimmie, B. M., & Tangen, J. M. (2022). *How do jurors evaluate likelihood ratios alongside diagnostic information?* Retrieved from: <https://osf.io/n27p8/>

Spellman, B. A. (2018). Communicating Forensic Evidence: Lessons from Psychological Science. *Seton Hall Law Review*, 48(3), 827-840.

Tribe, L. H. (1971). Trial by mathematics: Precision and ritual in the legal process. *Harvard Law Review*, 84(6), 1329-1393. doi: 10.2307/1339610

Wilcox, R. R. (2012). *Introduction to robust estimation and hypothesis testing* (3rd Ed.). Boston, MA: Academic Press.