

Advancing a paradigm shift in evaluation of forensic evidence – Part I: The rise of forensic data science

Authors and affiliations:

Geoffrey Stewart Morrison^{1,2,*}

¹ Forensic Data Science Laboratory, Aston University, Birmingham, UK

² Forensic Evaluation Ltd, Birmingham, UK

*Corresponding author: G.S. Morrison, e-mail: geoff-morrison@forensic-evaluation.net

ORCID:

Geoffrey Stewart Morrison 0000-0001-8608-8207

Disclaimer:

All opinions expressed in the present paper are those of the author, and, unless explicitly stated otherwise, should not be construed as representing the policies or positions of any organizations with which the author is associated.

Declaration of competing interest:

The author declares that he has no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements:

This research was supported by Research England's Expanding Excellence in England Fund as part of funding for the Aston Institute for Forensic Linguistics 2019–2022.

I dedicate this paper in memory of Prof Terrance M Nearey, my former PhD supervisor, who taught me how to be a scientist. Terry passed away on December 18, 2021.

1

2 **Advancing a paradigm shift in evaluation of forensic evidence – Part I: The rise** 3 **of forensic data science**

4

5 **Abstract**

6 Widespread practice across the majority of branches of forensic science uses analytical
7 methods based on human perception and interpretive methods based on subjective
8 judgement. These methods are non-transparent and are susceptible to cognitive bias,
9 interpretation is often logically flawed, and forensic-evaluation systems are often not
10 empirically validated. We describe a paradigm shift in which existing methods are
11 replaced by methods based on relevant data, quantitative measurements, and statistical
12 models; methods that are transparent and reproducible, are intrinsically resistant to
13 cognitive bias, use the logically correct framework for interpretation of evidence (the
14 likelihood-ratio framework), and are empirically validated under casework conditions.

15 **Keywords**

16 forensic science; forensic data science; likelihood ratio; paradigm shift; validation

17 **1 Introduction**

18 The present paper is published in two parts. Part I describes an ongoing paradigm shift
19 in evaluation of forensic evidence. It describes the current state of affairs (*staus quo*),
20 the new paradigm (*quo vadis?*), obstacles to the advancement of the paradigm shift
21 (*impedimenta*), and a strategy to advance the paradigm shift (*via progredi*). It also
22 describes approaches to dealing with problems presented by the new paradigm, in
23 particular, how to calculate likelihood ratios for forensic data and how to validate
24 forensic-evaluation systems that output likelihood ratios.

25 Part of our strategy involves collaboration between forensic-data-science researchers
26 and researchers and practitioners in particular branches of forensic science who want
27 to adopt the new paradigm. In any branch of forensic science, we expect the number of
28 practitioners who initially want to adopt the new paradigm to be a very small, but it
29 will be more productive to work with a small minority on developing practical solutions
30 than to try to convince the majority without providing practical solutions. Once the
31 practical solutions are being used by the small minority, use of the new paradigm has
32 the potential to spread.

33 Part II (Basu et al., 2022) will provide an example of the application of part of our
34 strategy. It will describe building a relevant database and developing and validating
35 statistical models for forensic comparison of fired cartridge cases, a common tasks in
36 forensic firearm examination, a branch of forensic science in which the new paradigm
37 has so far made almost no progress. It focuses in particular on testing different
38 functional-data-analysis-based methods for feature extraction from 3D images of fired
39 cartridge cases. This will be a first step in attempting to advance the paradigm shift in
40 forensic firearm examination, and is intended to be an example that can be copied in
41 other branches of forensic science.

42

43 **2 A paradigm shift in evaluation of forensic evidence**

44 **2.1 Status quo**

45 Curran (2013):

46 Is forensic science the last bastion of resistance against statistics?

47 UK House of Lords Science and Technology Committee (HoL, 2019):

48 In regard to pattern comparison methods, ... “the comparison of fingerprints,
49 toolmarks, footwear, tire marks and ballistics” [are] “spot-the-difference”

50 techniques in which “there is little, if any, robust science involved in the analytical
51 or comparative processes used and as a consequence there have been questions
52 raised around the reproducibility, repeatability, accuracy and error rates of such
53 analysis.” (§155)

54 In forensic science, the process of *evaluation of strength of evidence* consists of:
55 *analysis*, i.e., extraction of information from items of interest (the evidence);¹ and
56 *interpretation*, i.e., drawing inferences with respect to the meaning of the information
57 extracted by the analysis. Items of interest may be, for example: a fingermark of
58 questioned source recovered from a crime scene and a fingerprint collected from a
59 known individual; a recording of a speaker of questioned identity on an intercepted
60 telephone call and a recording of a police interview with a speaker of known identity;
61 a fired cartridge case recovered from a crime scene and cartridge cases fired in a
62 forensic laboratory from a gun found in the possession of a suspected shooter. Forensic
63 practitioners conduct evaluations in order to assist legal-decision makers to make
64 decisions with respect to questions of legal concern such as: Do the fingermark and
65 fingerprint originate from the same finger? Is the speaker of questioned identity on the
66 intercepted recording the same as the speaker of known identity? Was the cartridge
67 case recovered from the crime scene fired from the suspect’s gun?²

68 Currently, across the majority of branches of forensic science, widespread practice is
69 that analysis is conducted using *human perception*, and interpretation is conducted
70 using *subjective judgement*. Even in branches of forensic science in which analysis is
71 conducted using instrumental measurement, interpretation is commonly based on
72 subjective judgement, e.g., by eyeballing graphical representations of the measured

¹ In the forensic-inference-and-statistics literature “evidence” is the term commonly used to refer to the items of interest (i.e., the input to the analysis) or to the information output by the analysis (i.e., the input to the interpretation). Usage is somewhat fluid, but, either way, this is evidence from the perspective of the forensic practitioner. From the perspective of the court, evidence is the oral testimony or written submission of the forensic practitioner. In wider forensic-science literature, the term “trace” is often used to refer to items of interest.

² The present paper is framed in the context of source-level comparison.

73 values. Human-perception-based analysis methods and subjective-judgement-based
74 interpretation methods are non-transparent and susceptible to cognitive bias. Across
75 the majority of branches of forensic science, even branches of forensic science in which
76 interpretation is conducted using statistical models, interpretation of evidence is often
77 logically flawed, and forensic-evaluation systems (the end-to-end combination of
78 analysis and interpretation methods) are often not empirically validated or not
79 adequately empirically validated.³

80 **2.2 Quo vadis?**

81 **2.2.1 Introduction**

82 Saks & Koehler (2005):

83 we envision a paradigm shift in the traditional forensic identification sciences in
84 which untested assumptions and semi-informed guesswork are replaced by a
85 sound scientific foundation and justifiable protocols. Although obstacles exist
86 both inside and outside forensic science, the time is ripe for the traditional forensic
87 sciences to replace antiquated assumptions of uniqueness and perfection with a
88 more defensible empirical and probabilistic foundation. (p. 895)

89 US President's Council of Advisors on Science and Technology (PCAST, 2016):

90 neither experience, nor judgment, nor good professional practice ... can substitute
91 for actual evidence of foundational validity and reliability. The frequency with
92 which a particular pattern or set of features will be observed in different samples,
93 which is an essential element in drawing conclusions, is not a matter of
94 "judgment." It is an empirical matter for which only empirical evidence is
95 relevant. (p. 6)

³ Claims made in this introductory paragraph are supported by details and references provided later in §2.2 and in §2.3.

96 Objective methods are, in general, preferable to subjective methods. Analyses that
97 depend on human judgment (rather than a quantitative measure ...) are obviously
98 more susceptible to human error, bias, and performance variability across
99 examiners. In contrast, objective, quantified methods tend to yield greater
100 accuracy, repeatability and reliability, including reducing variation in results
101 among examiners. Subjective methods can evolve into or be replaced by objective
102 methods. (p. 47)

103 A *paradigm shift* in evaluation of forensic evidence is ongoing in which methods based
104 on human perception and subjective judgement are being replaced by methods based
105 on *relevant data, quantitative measurements, and statistical models*; methods that:

- 106 1. are *transparent and reproducible*;
- 107 2. are *intrinsically resistant to cognitive bias*;
- 108 3. use the *logically correct framework for interpretation of evidence* (the
109 *likelihood-ratio framework*); and
- 110 4. are *empirically validated under casework conditions*.

111 We address each of these elements in the following four subsections (§2.2.2–§2.2.5).

112 **2.2.2 Transparency and reproducibility**

113 Methods dependent on human perception and subjective judgement are intrinsically
114 non-transparent and therefore not reproducible by others. Human introspection is often
115 mistaken, hence a forensic practitioner's explanation of how they reached their
116 conclusion may not reflect how they actually reached that conclusion (Edmond et al.,
117 2017). In contrast, procedures based on data, quantitative measurement, and statistical
118 models are transparent and reproducible: measurement (feature-extraction) and
119 statistical-modelling methods can be described in detail, and data and software tools
120 can potentially be shared with others.

121 **2.2.3 Cognitive bias**

122 There has been a great deal of concern about cognitive bias in forensic science
123 (National Research Council, 2009; Expert Working Group on Human Factors in Latent
124 Print Analysis, EWG, 2012; Found, 2015; Stoel et al., 2015; PCAST, 2016; Edmond
125 et al., 2017; Cooper & Meterko, 2019). Cognitive bias is subconscious bias, it cannot
126 be controlled by strength of will. Forensic practitioners are susceptible to cognitive bias
127 when making perceptual observations: their belief in the probability that a hypothesis
128 is true can affect their analysis of the evidence and therefore the information that feeds
129 into their interpretation. Forensic practitioners are susceptible to cognitive bias when
130 they are making subjective judgements and are exposed to information that could
131 influence their belief in the probability that a hypothesis is true but that would not
132 logically affect the probability of obtaining the evidence conditional on whether the
133 hypothesis were true. Some potentially biasing information is task-irrelevant and
134 should be withheld from practitioners, but some potentially biasing information is task-
135 relevant and practitioners employing human-perception and subjective-judgement
136 methods will need to be exposed to it at some point in the evaluation process, e.g.,
137 practitioners who visually compare known-source fingerprints and questioned-source
138 fingermarks must be exposed to both, but exposure to a higher-quality print may bias
139 their analysis of ambiguous details in a lower-quality mark. Systems in which the
140 strength-of-evidence conclusion is directly the result of subjective judgement are
141 particularly susceptible to cognitive bias.

142 Systems based on quantitative measurements and statistical models require subjective
143 judgements on decisions such as whether the data used for training the system and the
144 data used for validating the system are sufficiently representative of the relevant
145 population for the case and sufficiently reflective of the conditions of the items of
146 interest in the case so that the output of the system will be a meaningful answer to the
147 question posed in the case and so that the results of the validation will provide a
148 meaningful indication of the performance of the systems under the conditions of the

149 case. These decisions, however, are made at the beginning of the process before the
150 practitioner has analyzed the items of interest, hence the practitioner cannot know what
151 effect these decisions will have on the strength-of-evidence conclusion. The remainder
152 of the evaluation process is automated, hence not susceptible to cognitive bias.

153 **2.2.4 Likelihood-ratio framework**

154 In current practice, interpretation of evidence is often logically flawed, e.g., it is based
155 on the uniqueness or the individualization fallacy (Saks & Koehler, 2008; Cole, 2009,
156 2014), and conclusions are often expressed categorically, e.g., “identification”,
157 “inconclusive”, “exclusion” (i.e., posterior probability of 1 or 0 with respect to the
158 same-source hypothesis, with “inconclusive” meaning no conclusion rather than an
159 intermediate probability), or using some form of uncalibrated verbal posterior-
160 probability scale, e.g., “identification”, “probable identification”, “inconclusive”,
161 “probable exclusion”, “exclusion”. Jackson (2009) and Kaye (2015) review these and
162 other commonly used but logically flawed conclusions.

163 In contrast, the likelihood-ratio framework is advocated as the logically correct
164 framework for evaluation of evidence by the vast majority of experts in forensic
165 inference and statistics (including: Aitken et al., 2011; Morrison et al., 2017; Morrison,
166 Enzinger, et al., 2021; with 31, 19 and 20 authors and supporters respectively), and by
167 key organizations including: Association of Forensic Science Providers of the United
168 Kingdom and of the Republic of Ireland (2009); Royal Statistical Society (Aitken et
169 al., 2010); European Network of Forensic Science Institutes (Willis et al., 2015);
170 National Institute of Forensic Science of the Australia New Zealand Policing Advisory
171 Agency (Ballantyne et al., 2017); American Statistical Association (Kafadar et al.,
172 2019); Forensic Science Regulator for England & Wales (FSR, 2021).

173 The likelihood-ratio framework requires assessment of the probability of obtaining the
174 evidence if one hypothesis were true versus the probability of obtaining the evidence
175 if an alternative hypothesis were true. The two hypotheses must be mutually exclusive.

176 One hypothesis should represent the position of the prosecution in the case, and the
177 other the position of the defence, e.g., the fingerprint of questioned origin was
178 deposited by a finger of a particular known individual, versus the fingerprint of
179 questioned origin was deposited by a finger of some other individual selected at random
180 from the relevant population. In this source-level fingerprint-fingerprint comparison
181 example, the numerator of the likelihood ratio quantifies the *similarity* between the
182 mark and the print, and the denominator quantifies the *typicality* of the mark with
183 respect to the relevant population. For continuously-valued data, likelihood ratios can
184 be calculated as the ratio of two probability-density functions evaluated at the value of
185 the evidence. In the forensic inference and statistics literature “likelihood ratio” is
186 commonly used as a cover term for both likelihood ratios based only on sample data
187 and for Bayes’ factors based on sample data and prior distributions for model
188 parameters. The present discussion is intended to be neutral with respect to
189 likelihoodist or Bayesian approaches.

190 **2.2.5 Empirical validation**

191 Empirical validation under conditions reflecting those of the case to which a forensic-
192 evaluation system is to be applied is the only way to know how well that system
193 performs under the conditions of the case. Protocols for validating systems that output
194 likelihood ratios have been developed, including metrics and graphics appropriate for
195 representing the results of such validations (Morrison, 2011; Meuwly et al., 2017;
196 Ramos et al., 2020; Morrison, Enzinger, et al., 2021). Much of the latter development
197 has been conducted in the context of forensic voice comparison, but the results are
198 applicable across forensic science in general. The need for validation under casework
199 conditions has been emphasized by FSR (2020b), and by PCAST (2016):

200 Without appropriate estimates of accuracy, an examiner’s statement that two
201 samples are similar—or even indistinguishable—is scientifically meaningless: it
202 has no probative value, and considerable potential for prejudicial impact.

203 Nothing—not training, personal experience nor professional practices—can
204 substitute for adequate empirical demonstration of accuracy. (p. 46)

205 Despite this, practitioners in multiple branches of forensic science often claim that
206 training and experience provide sufficient warrant for their conclusions (see Mnookin
207 et al., 2011; Risinger, 2013; PCAST, 2016; Morrison & Thompson, 2017), deny or
208 obfuscate about the need for validation (see Cole, 2006; Morrison, 2014; PCAST,
209 2016; Koehler, 2017; Morrison et al., 2018), or propose lax validation protocols that
210 do not require demonstration of performance under casework conditions (see Morrison,
211 Neumann, et al., 2020, 2021).

212 **2.2.6 A Kuhnian paradigm shift**

213 The idea that evaluation of forensic evidence is undergoing a paradigm shift is not new.
214 The most famous article heralding a paradigm shift is Saks & Koehler (2005). Allowing
215 for differences in wording and level of detail, and allowing for difference due to the
216 passage of time, we believe that Saks & Koehler (2005) and the present paper describe
217 the same paradigm shift. In contrast to Saks & Koehler’s (2005) statement that they
218 intended “paradigm shift” as a metaphor, however, we view the paradigm shift in
219 evaluation of forensic evidence as a true Kuhnian paradigm shift (Kuhn, 1962) in the
220 sense that it requires rejection of existing methods and the ways of thinking that
221 underpin them, and rejection of the idea that progress can be made by incremental
222 improvements to existing methods. Instead, it requires the wholesale adoption of an
223 entire constellation of new methods and new ways of thinking.⁴

224 That a paradigm shift requires the wholesale adoption of an entire constellation of new
225 methods and new ways of thinking remains the case irrespective of whether one
226 considers the shift to be from one paradigm to another or to be from a pre-paradigm to

⁴ The only aspects in which we think the current paradigm shift in evaluation of forensic evidence deviates from Kuhn’s (1962) description of paradigm shifts relate to forensic science being an applied science which is not isolated from societal pressures.

227 a paradigm period of science. As suggested in Saks & Koehler (2005), a pre-paradigm
228 period would seem to be a more accurate description of the *status quo*, with multiple
229 traditions of evaluation of evidence used both within individual branches of forensic
230 science and across different branches of forensic science, and hence there being no
231 established widely-accepted overarching paradigm in use.

232 Some authors have used the term “paradigm shift” in relation to a single element or a
233 subset of the elements of the paradigm shift as we have outlined it above, but we believe
234 that all of these elements are required as part of the constellation. This may be viewed
235 as a radical stance, and it faces resistance, but over the last decade and a half we have
236 had substantial success in contributing to advancing this paradigm shift in forensic
237 voice comparison.⁵

238 Over many years, multiple people have suggested to us that we give the new paradigm
239 a name, rather than simply referring to it as the new paradigm. We hereby suggest that
240 the new paradigm can be referred to as *forensic data science*, on the condition that the
241 definition of forensic data science is understood to encompass the new paradigm as we
242 have described it above.

243 **2.3 Impedimenta**

244 **2.3.1 Introduction**

245 The paradigm shift in evaluation of forensic evidence is ongoing, but progress is slow
246 or stalling for multiple reasons, including the following:

- 247 1. The new paradigm has only been adopted in a few branches of forensic
248 sciences, and only by a minority of researchers and practitioners.
- 249 2. Only some elements of the new paradigm have been adopted as part of

⁵ “we” here refers to Morrison and colleagues in general.

250 incremental change.

251 3. There is misunderstanding of the new paradigm and resistance to its adoption.

252 4. Research is often not informed by practice and has no impact on practice.

253 5. It is difficult to obtain funding for evidential-forensic-science research.

254 6. There are genuine practical impediments to implementing the new paradigm.

255 We discuss each of these impediments in the following six subsections (§2.3.2–§2.3.7).

256 **2.3.2 The new paradigm has only been adopted in a few branches of forensic**
257 **sciences, and only by a minority of researchers and practitioners**

258 The new paradigm has only been adopted in a few branches of forensic sciences, and
259 only by a minority of researchers and practitioners. In the 1990s, the new paradigm
260 was widely adopted for **forensic evaluation of DNA** (Foreman et al., 2003). Although
261 the volume and importance of casework in this branch of forensic science makes it
262 influential, single-source DNA profiles are invariant and discrete, and therefore have a
263 very different structure from the continuously-valued data with within-source
264 variability that results from analyses in most other branches of forensic science. The
265 situation is more complex for low-template DNA and for DNA mixtures, but there is
266 still a difference in data structure. Interpretation of DNA profiles is also dependent on
267 well-developed theory of genetic inheritance, whereas interpretation in most branches
268 of forensic science will have to be data driven (as is the case in machine learning in
269 general including in biometric applications). The potential for transfer of new-
270 paradigm knowledge and methods from DNA to other branches of forensic science is
271 therefore limited.

272 Since around 2000, a growing number of researchers and practitioners in **forensic voice**
273 **comparison** have developed and adopted methods for calculation of likelihood ratios
274 based on acoustic measurements and statistical models, and methods for calibration

275 and validation of likelihood-ratio systems under casework conditions. This has
276 included adoption of state-of-the-art machine-learning approaches to automatic
277 speaker recognition (Lee et al., 2020; Matějka et al., 2020; Morrison, Enzinger, et al.,
278 2020; Villalba et al., 2020; Weber et al., 2022). At present, however, only a minority
279 of practitioners have adopted the new paradigm: In a survey of law-enforcement
280 agencies in INTERPOL member countries (Morrison et al., 2016), among respondents
281 who had forensic-voice-comparison capabilities, the reported rates of adoption of
282 human-supervised-automatic approaches and numeric likelihood ratios were 33% and
283 23% respectively; however, even the latter have probably not adopted all elements of
284 the new paradigm. In a survey including private practitioners (Gold & French, 2019),
285 the reported rates of adoption of human-supervised-automatic approaches and numeric
286 likelihood ratios were 41% and 13% respectively (up from 20% and 9% in Gold &
287 French, 2011); however, inconsistent with the new paradigm, most reported combining
288 human-supervised-automatic approaches with human-perception- and subjective-
289 judgement-based approaches. Data in human-supervised-automatic approaches are
290 continuously valued and have intrinsic within-source variability, a data structure shared
291 with many other branches of forensic science. Compared to DNA, new-paradigm
292 knowledge and methods from forensic voice comparison, including statistical models
293 and calibration and validation procedures, should therefore be easier to transfer to and
294 adapt for other branches of forensic science. Forensic voice comparison is, however, a
295 relatively low-volume branch of forensic science, which limits the extent to which
296 developments in forensic voice comparison are noticed and adopted by researchers and
297 practitioners in other branches of forensic science.

298 Curran (2013) lamented that only 13% of laboratories surveyed used the likelihood-
299 ratio framework for **glass evidence**, but this may be one of the highest rates of adoption
300 after DNA. In many **other branches of forensic science**, including firearm
301 examination, the rate of adoption of the likelihood-ratio framework by practitioners is
302 near zero (Bali et al., 2020; Cole & Barno, 2020).

303 **2.3.3 Only some elements of the new paradigm have been adopted as part of**
304 **incremental change**

305 Only some elements of the new paradigm have been adopted as part of incremental
306 change. Although in the short term this may be viewed as a step in the right direction,
307 in the long term it may actually impede a paradigm shift.

308 Just because it is a transition between incommensurables, the transition between
309 competing paradigms cannot be made a step at a time, ... Like the gestalt switch,
310 it must occur all at once (though not necessarily in an instant) or not at all. (Kuhn,
311 1962, p. 149)

312 Some practitioners **assign likelihood-ratio values based on subjective judgement,**
313 and the values they assign are not subject to empirical calibration or empirical
314 validation (see Risinger, 2013; Morrison & Thompson, 2017; Morrison, Ballantyne,
315 Geoghegan, 2020). Some authors emphasize the logic of the likelihood-ratio
316 framework and consider subjective assignment of likelihood ratio an acceptable end
317 goal or consider it a step in the right direction, but such incremental steps potentially
318 impede a paradigm shift which requires the abandonment of interpretation methods
319 that are entirely dependent on subjective judgement.⁶ In addition, placing an emphasis
320 on subjectivist concepts of probability is detrimental to attempts to encourage
321 practitioners to adopt methods based on relevant data, quantitative measurements, and
322 statistical models, and to adopt empirical validation under casework conditions
323 (Morrison, 2017).

324 The majority of **proposals to address cognitive bias** in forensic science (e.g., EWG,
325 2012; Stoel et al., 2015; Thompson et al., 2017; FSR, 2020a) **assume the continued**
326 **use of human-perception- and subjective-judgement-based methods.** Although this
327 may be necessary in the short term, it potentially impedes a paradigm shift to

⁶ Our stance on where in the forensic-evaluation process use of subjective judgement is acceptable is more restrictive than that of some leaders in the field, e.g., Evett et al. (2017).

328 quantitative-measurement- and statistical-model-based methods.

329 Some practitioners **make use of systems based on quantitative measurements and**
330 **statistical models, but do not empirically calibrate or validate** their system using
331 data that reflect the relevant population and the conditions for the case, and/or **rather**
332 **than directly report the output of the system, they use it as input to a subjective-**
333 **judgement process** that also considers other information including from human-
334 perception-based analyses (see Morrison & Thompson, 2017; Morrison, 2018a,
335 2018b). Such approaches are pernicious in that use of technology may give the false
336 impression of scientific validity, and reaction against this may impede a paradigm shift
337 that includes adoption of quantitative measurements and statistical models.

338 **2.3.4 There is misunderstanding of the new paradigm and resistance to its** 339 **adoption**

340 As with all Kuhnian paradigm shifts, there is misunderstanding of the new paradigm
341 and resistance to its adoption. **Some resistance is cultural.** People in general tend to
342 prefer methods which involve greater human input even when validation results
343 indicate that data-driven methods lead to better results, but over time they can come to
344 accept the latter (Swofford & Champod, 2021). The cultures of some branches of
345 forensic science seem to be especially resistant to the adoption of statistical-model-
346 based methods and of validation (see Mnookin et al., 2011; Curran, 2013; Morrison,
347 2014; Morrison & Stoel, 2014; Morrison, Neumann, et al., 2020, 2021; Swofford et al.,
348 2021).

349 There is a **belief that likelihood ratios are difficult to understand** (Bali et al. 2020;
350 Swofford et al., 2021). Commonly occurring misunderstandings have even been given
351 names, e.g., the “prosecutor’s fallacy” and the “defense attorney’s fallacy” (Thompson
352 & Schurmann, 1987). There are many examples of legal rulings in which judges have
353 misunderstood the meaning of a likelihood ratio (the England & Wales Court of Appeal
354 2010 ruling in *R v T* is an infamous example, see, e.g., Berger et al., 2011; Redmayne

355 et al., 2011; Morrison, 2012; Thompson, 2012). Results of empirical research on lay
356 understanding of expressions of strength of evidence are mixed (Eldridge, 2019;
357 Martire & Edmond, 2020).

358 Despite legal rulings and recommendations concerning the need for validation, **courts**
359 **often do not understand empirical validation and its necessity**, and accept
360 testimony based on forensic-evaluation methods that have not been validated under
361 conditions reflecting those of the case under consideration, or even that have not been
362 empirically validated at all (see Bernstein, 2013; Morrison, 2014, 2018a; Cooper, 2016;
363 Edmond, 2018).

364 **2.3.5 Research is often not informed by practice and has no impact on practice**

365 Research is often not informed by practice and has no impact on practice. **Research**
366 **that appears to be about forensic science may not actually be about solving real**
367 **forensic-science problems**. For example, research may really be about a method in
368 statistics or machine learning with an apparent forensic-science problem being used as
369 an example of the application of that method.

370 Research in forensic science is sorely needed, but it should address primarily
371 forensic science questions—not questions relating to the application of chemistry,
372 biology, statistics, or psychology. (Margot, 2011, p. 801)

373 Research that is not focused on solving real forensic science problems and is divorced
374 from forensic practice may lead to academic papers but nothing else.

375 it is critical that researchers and funding bodies understand the importance of
376 conducting research that is informed by practice and can be translated into
377 practical applications. (Roux & Weyermann, 2021, p. 2)

378 Even research that does address real forensic-science problems will fail to have impact
379 unless it involves genuine collaboration in which researchers understand the demands

380 of practice, and in which practitioners are willing to embrace research-informed change
381 (Curran, 2013).

382 **2.3.6 It is difficult to obtain funding for evidential-forensic-science research**

383 It is difficult to obtain funding for evidential-forensic-science research. **Few funding**
384 **agencies have sustained funding targeted at forensic science**, and funding agencies
385 seldom have panels of reviewers knowledgeable about evidential forensic science.
386 Applications for funding for evidential-forensic-science research made to non-
387 forensic-science-targeted calls are often rejected because reviewers do not understand
388 the epistemology or value of forensic-science research. Applications are often rejected
389 because their goals are to improve forensic science, which is an applied science, and
390 **funding-agency criteria or reviewers often do not value applied science.**

391 The larger scientific community must now come to the aid of our forensic
392 colleagues in advocating both for: (i) the research and financial support that is so
393 clearly needed to advance the field and (ii) the requirement for empirical testing
394 that is so clearly needed to advance the cause of justice. ... Forensic scientists have
395 long complained that their work is not always valued by their scientific colleagues
396 because of its applied nature; it is time for the scientific community to move
397 beyond that conceit. (Bell et al., 2018)

398 At the other extreme, when there are calls specific to forensic science, they usually
399 **focus exclusively or primarily on short-term goals related to law-enforcement**
400 **investigative applications rather than on courtroom-evidential applications**
401 (investigative and evidential applications have very different requirements), and they
402 usually focus on technology rather than on forensic inference.

403 technology-oriented development ... often overrul[es] the importance of
404 appropriate scientific reasoning to solve actual problems (Roux et al., 2021, p.
405 679)

406 **Research calls requiring deliverables with a high technology readiness level (TRL)**
407 **are not sources of funding for paradigm-shifting research.**

408 In 2018, United Kingdom Research and Innovation (UKRI) informed the UK House
409 of Lords Science and Technology Select Committee that UKRI had invested GBP 56M
410 over 10 years in forensic-science research (less than 0.1% of UKRI's total budget), but
411 on closer inspection most of the funding counted to obtain that figure was not for
412 research projects that actually focused on (or even made any contribution to) forensic
413 science: only about GBP 17M went to forensic-science focussed research, and GBP
414 15M of that went to TRL research, only GBP 2M to foundational research (HoL, 2019;
415 Morgan & Levin, 2019). HoL (2019) recommended that UKRI "urgently and
416 substantially increase the amount of dedicated funding allocated to forensic science"
417 (§187), but (as of time of writing ~2.5 years after the publication of the HoL report)
418 this has not (yet) happened.

419 **2.3.7 There are genuine practical impediments to implementing the new**
420 **paradigm**

421 There are genuine practical impediments to implementing the new paradigm. Even if
422 practitioners want to adopt the new paradigm, they will be unable to do so unless they
423 are provided with the quantitative-measurement and statistical-modelling tools and the
424 case-relevant data necessary to calculate likelihood ratios and validate systems under
425 the conditions of the cases on which they work. Practitioners will also not be able to
426 adopt the new paradigm unless they are provided with training on understanding the
427 new paradigm in general and on how to implement it for the types of cases they work
428 on.

429 **2.4 Via progredi**

430 Kuhn (1962):

431 The transfer of allegiance from paradigm to paradigm is a conversion experience

432 that cannot be forced. ... a generation is sometimes required to effect the change
433 ... Conversions will occur a few at a time until, after the last holdouts have died,
434 the whole profession will again be practicing under a single, but now a different,
435 paradigm. (pp. 150–151)

436 Kuhnian paradigm shifts are not rapid and individuals cannot be forced to embrace the
437 new paradigm, but our aim is to facilitate and thereby advance the adoption of the new
438 paradigm. Our **strategy** is to work with researchers and practitioners who want to adopt
439 the new paradigm, to work with them on addressing practical impediments to applying
440 the new paradigm in casework, i.e.:

- 441 1. to provide researchers, practitioners, and lawyers with training leading to
442 understanding of the new paradigm;
- 443 2. to collaborate with researchers and practitioners on building relevant databases
444 and on developing and validating statistical models applicable in their
445 particular branches of forensic science; and
- 446 3. to conduct research on how to present likelihood ratios and validation results
447 so as to maximize understanding by laypeople, and thereby provide guidance
448 to forensic practitioners on how to communicate forensic-evaluation results to
449 legal-decision makers.

450 Part II of the present paper will focus on element 2 of the strategy. We will focus on
451 elements 1 and 3 elsewhere. Element 2 of the strategy requires collaboration between
452 researchers with expertise in forensic data science and researchers and practitioners
453 with expertise in particular branches of forensic science. Academic publications are
454 unlikely to convince practitioners to adopt the new paradigm, but other practitioners
455 successfully applying the new paradigm are potentially convincing. In any branch of
456 forensic science, the number of practitioners who initially want to adopt the new
457 paradigm and who want to collaborate on this endeavour will almost certainly be a very

458 small minority, but it will be more productive to work with a small minority on
459 developing practical solutions than to try to convince the majority of practitioners
460 without providing practical solutions. Once the practical solutions are being used by
461 the small minority, use of the new paradigm has the potential to spread. Even then, we
462 do not expect adoption of the new paradigm to be rapid, but we do expect higher rates
463 of adoption among newer practitioners and trainees, leading to a generational shift.

464

465 **3 Calculating likelihood ratios and validating likelihood-ratio systems**

466 **3.1 Introduction**

467 A new paradigm introduces new ways of thinking, which can introduce new problems
468 to be solved, problems that may not even have been conceivable under the old
469 paradigm. In the context of source attribution, one such problem is the need to take
470 account not only of the similarity between the items of interest, but also of their
471 typicality with respect to the relevant population. The relevant population is the
472 population from which, in the context of the legal case, the item of questioned origin
473 could conceivably have come had it not come from the known source.⁷ The data used
474 to train the likelihood-ratio models (particularly the model in the denominator) must
475 be representative of the relevant population for the case, otherwise the calculated
476 likelihood ratio will not answer the question of interest for the case.

477 Typicality with respect to the relevant population can be incorporated into the
478 calculation of a likelihood ratio using either a specific-source or a common-source
479 approach (Ommen & Saunders, 2021). Another approach, similarity-score-based

⁷ The relevant population could be explicitly proposed by the defence, but in common-law systems the defence in a criminal trial is under no obligation to propose an alternative to the prosecution's proposition. A forensic practitioner usually has to decide what to adopt as the relevant population and communicate that decision to the court so that the court can later decide whether the forensic practitioner's decision was appropriate, i.e., will it result in a likelihood ratio that answers a question of interest for the court.

480 calculation of likelihood ratios does not properly incorporate typicality with respect to
481 the relevant population. We discuss each of these approaches in §3.2 below.

482 Another problem introduced by the new paradigm is how to validate the performance
483 of forensic-evaluation systems that output likelihood ratios. In a classification
484 framework, a test-pair input is either “same source” or “different source” and the
485 system being tested outputs either “same source” or “different source”. The *miss rate*
486 can be calculated as the proportion of “same source” inputs that elicited “different
487 source” outputs, the *false-alarm rate* can be calculated as the proportion of “different
488 source” inputs that elicited “same source” outputs, and the mean of the miss rate and
489 false-alarm rate can be reported as the *classification-error rate*, a single-value
490 summarizing the performance of the system. In a likelihood-ratio framework, however,
491 the system being tested does not output a discrete classification, instead, it outputs a
492 continuous likelihood-ratio value. What must be assessed in terms of performance is
493 the gradient goodness of the output: Given a same-source input, a good output would
494 be a likelihood-ratio value that is much larger than 1, a less good output would be a
495 value that is only a little larger than 1, a bad output would be a value less than 1, and a
496 worse output would be a value much less than 1. *Mutatis mutandis*, given a different-
497 source input, a good output would be a value much less than 1. In §3.3, we describe a
498 popular single-value metric for summarizing the performance of likelihood-ratio
499 systems, the *log-likelihood-ratio cost* (C_{llr} , Brümmer and du Preez, 2006), and a
500 popular graphical representation of likelihood-ratio validation results, a *Tippett plot*
501 (Meuwly, 2001). We will use C_{llr} and Tippett plots when reporting the results of the
502 empirical research described in Part II of the present paper.

503 **3.2 Calculating likelihood ratios**

504 **3.2.1 Specific-source approach**

505 A specific-source likelihood ratio answers the two-part question:

506 1. What is the probability of obtaining the measured properties of the item of
507 questioned source if it came from the specific known source?

508 versus

509 2. What is the probability of obtaining the measured properties of the item of
510 questioned source if it came not from the specific known source but from some
511 other source selected at random from the relevant population?

512 A specific-source calculation has the form given in Equation (1), in which Λ is the
513 likelihood ratio, $f(x|M)$ is a probability-density function, x_q is the feature (or feature
514 vector) extracted from the questioned-source item, and M_k and M_r are the specific-
515 known-source model and the relevant-population model respectively. The specific-
516 known-source model is trained using features extracted from multiple items sampled
517 from the specific known source, and the relevant-population model is trained using
518 features extracted from multiple items sampled from the relevant population.

519 (1)

$$520 \quad \Lambda = \frac{f(x_q|M_k)}{f(x_q|M_r)}$$

521 Equation (2) provides an example of a simple model for calculating specific-source
522 likelihood ratios. The model assumes univariate Gaussian distributions in both the
523 numerator and the denominator, and assumes that all sources have the same within-
524 source variance σ_w^2 . In Equation (2), $f(x|\mu, \sigma^2)$, is a Gaussian probability-density
525 distribution (parametrized using mean and variance), μ_k and μ_r are the specific-known-
526 source mean and the relevant-population mean respectively, and σ_w^2 and σ_b^2 are the
527 within-source variance and between-source variance respectively.

528 (2)

529
$$\Lambda = \frac{f(x_q | \mu_k, \sigma_w^2)}{f(x_q | \mu_r, \sigma_w^2 + \sigma_b^2)}$$

530 Figure 1 shows examples of the calculation of two specific-source likelihood ratios
 531 using Equation (2). The wider distribution represents the relevant-population model
 532 (denominator model), and the two peakier distributions represent two different
 533 specific-source models (numerator models). For these examples: $\mu_r = 0$, $\sigma_b^2 = 100$,
 534 and $\sigma_w^2 = 1$; for the filled circles $\mu_k = 0$ and $x_q = -1$, and for the unfilled circles $\mu_k =$
 535 20 and $x_q = 19$.⁸ The resulting likelihood-ratio values, corresponding to the filled and
 536 unfilled circles respectively, are $0.24/0.040 = 6$ and $0.24/0.0066 = 36$.

537 [insert figure about here]

538 **Figure 1.** Examples of the calculation of specific-source likelihood ratios.

539

540 3.2.2 Common-source approach

541 A common-source likelihood ratio answers the two-part question:

542 1. What is the probability of obtaining the measured properties of the items of
 543 questioned- and known-source if they both came from the same source (a
 544 source selected at random from the relevant population)?

545 versus

546 2. What is the probability of obtaining the measured properties of the items of
 547 questioned- and known-source if they each came from a different source (each
 548 a source selected at random from the relevant population)?

⁸ The values in these examples, and in other examples below, were chosen purely for illustrative purposes.

549 A common-source calculation has the form given in Equation (3), in which Λ is the
550 likelihood ratio, $f(x_q, x_k|M)$ is a joint probability-density function, x_q and x_k are the
551 features (or feature vectors) extracted from the questioned- and known-source items
552 respectively, and M_s and M_d are the same-source and different-source models
553 respectively.

554 (3)

555
$$\Lambda = \frac{f(x_q, x_k|M_s)}{f(x_q|M_d)f(x_k|M_d)}$$

556 Equation (4) provides an example of a simple model for calculating common-source
557 likelihood ratios.⁹ The model assumes univariate Gaussian distributions in both the
558 numerator and the denominator, and assumes that all sources have the same within-
559 source variance σ_w^2 . The numerator of Equation (4) integrates over all possible values
560 for source means given the between-source distribution, with the constraint that x_q and
561 x_k come from the same source. The denominator of Equation (4) integrates over all
562 possible values for source means given the between-source distribution, but does so
563 independently for x_q and for x_k . The solutions to the integrals can be expressed as
564 bivariate Gaussian distributions in which for the same-source model (the numerator
565 model) the covariances equal the between-source variance, σ_b^2 , but for the different-
566 source model (the denominator model) the covariances are zero, 0. This reflects the
567 logic that if a source is selected at random from the population and the mean of the
568 source is high, then two items selected at random from that source will both be expected
569 to have high values. Likewise, if the mean of the source is low, the values of both items
570 will be expected to be low. In general, the values of two items selected from the same

⁹ This is the univariate version of the model described Aitken & Lucy (2014) as the “multivariate normal (MVN) procedure”, and in the automatic-speaker-recognition and forensic-voice-comparison literature (e.g., Prince & Elder, 2007; Kenny, 2010; Brümmer & de Villiers, 2010; Sizov et al., 2014; Morrison, Enzinger, et al., 2020) as the “two-covariance model” for “probabilistic linear discriminant analysis (PLDA)”.

571 source are expected to be correlated. In contrast, if two sources are selected at random
 572 from the population and the mean of one is high, there is no expectation that the mean
 573 of the other source will also be high. The mean of the other source is more likely to the
 574 average, and it is equally likely to be low as to be high. The values of two items each
 575 selected from a different source are not expected to be correlated.

576 (4)

$$\begin{aligned}
 577 \quad \Lambda &= \frac{\int f(x_q|\mu_i, \sigma_w^2) f(x_k|\mu_i, \sigma_w^2) f(\mu_i|\mu_r, \sigma_b^2) d\mu_i}{\int f(x_q|\mu_i, \sigma_w^2) f(\mu_i|\mu_r, \sigma_b^2) d\mu_i \int f(x_k|\mu_j, \sigma_w^2) f(\mu_j|\mu_r, \sigma_b^2) d\mu_j} \\
 578 \quad &= \frac{f\left(\begin{bmatrix} x_q \\ x_k \end{bmatrix} \middle| \begin{bmatrix} \mu_r \\ \mu_r \end{bmatrix}, \begin{bmatrix} \sigma_w^2 + \sigma_b^2 & \sigma_b^2 \\ \sigma_b^2 & \sigma_w^2 + \sigma_b^2 \end{bmatrix}\right)}{f(x_q|\mu_r, \sigma_w^2 + \sigma_b^2) f(x_k|\mu_r, \sigma_w^2 + \sigma_b^2)} \\
 579 \quad &= \frac{f\left(\begin{bmatrix} x_q \\ x_k \end{bmatrix} \middle| \begin{bmatrix} \mu_r \\ \mu_r \end{bmatrix}, \begin{bmatrix} \sigma_w^2 + \sigma_b^2 & \sigma_b^2 \\ \sigma_b^2 & \sigma_w^2 + \sigma_b^2 \end{bmatrix}\right)}{f\left(\begin{bmatrix} x_q \\ x_k \end{bmatrix} \middle| \begin{bmatrix} \mu_r \\ \mu_r \end{bmatrix}, \begin{bmatrix} \sigma_w^2 + \sigma_b^2 & 0 \\ 0 & \sigma_w^2 + \sigma_b^2 \end{bmatrix}\right)}
 \end{aligned}$$

580 Figure 2 shows examples of the calculation of two common-source likelihood ratios
 581 using Equation (4). In Figure 2a the univariate distribution represents the relevant-
 582 population model with $\mu_r = 0$, and $\sigma_r^2 = \sigma_w^2 + \sigma_b^2$ ($\sigma_b^2 = 100$, and $\sigma_w^2 = 1$), and the
 583 pairs of circles represent pairs of x_q and x_k feature values. For the filled circles $x_q =$
 584 -1 and $x_k = 1$, and for the unfilled circles $x_q = 19$ and $x_k = 21$. Figure 2b shows the
 585 projection of the x_q and x_k feature values into the two-dimensional space of Equation
 586 (4), and shows the peakier same-source model (the numerator model) and the flatter
 587 different-source model (the denominator model). The resulting likelihood-ratio values,
 588 corresponding to the filled and unfilled circles respectively, are
 589 $(41 \times 10^{-4}) / (16 \times 10^{-4}) = 2.6$ and $(56 \times 10^{-5}) / (3.0 \times 10^{-5}) = 19$.

590 [insert figure about here]

591 **Figure 2.** Examples of the calculation of common-source likelihood ratios.

592

593 Although the example specific-source and common-source models above assume
594 univariate Gaussian distributions, both specific-source and common-source models can
595 fit more complex multivariate distributions. The latter can potentially exploit more
596 useful information resulting in better performance, but they require larger amounts of
597 data in order to estimate larger numbers of parameter values. In forensic casework, the
598 amount of case-relevant data is often limited. Solutions may include use of dimension-
599 reduction techniques, extraction of features using functional data analysis, and use of
600 embeddings from deep neural networks (DNNs) as feature vectors.

601 **3.2.3 Similarity-score approach**

602 Across multiple branches of forensic science, there have been repeated proposals to
603 calculate likelihood ratios using similarity-score-based approaches (see Morrison &
604 Enzinger, 2018). Such approaches may be motivated by situations in which a specific-
605 source approach cannot be adopted because the data only allow for a comparison
606 between one feature vector and one other feature vector (and this constraint may be
607 due to casework conditions, hence collecting more data may not be a possible solution).
608 They may also be motivated by the difficulties of trying to fit potentially complex
609 distributions in high-dimensional spaces using limited amounts of training data. Instead
610 of fitting models to feature vectors, as for the specific-source and common-source
611 models described above, models are fitted to scores which quantify degree of similarity
612 (or, inversely, degree of difference) between pairs of items. Similarity-scores are scalar
613 values that can be based, for example, on Manhattan distance, Euclidian distance, or a

614 correlation coefficient.¹⁰ Similarity-scores are calculated for pairs of items that are
615 known to come from the same source and for pairs of items known to come from
616 different sources, resulting in a set of same-source scores and a set of different-source
617 scores. The latter scores are used to train a model of the form given in Equation (5), in
618 which $\delta(x_q, x_k)$ is the function that calculates the score, and $M_{\delta,s}$ and $M_{\delta,d}$ are
619 univariate models trained on same-source scores and different-source scores
620 respectively.

621 (5)

$$622 \quad \Lambda = \frac{f(\delta(x_q, x_k)|M_{\delta,s})}{f(\delta(x_q, x_k)|M_{\delta,d})}$$

623 A similarity-score-based likelihood ratio answers the two-part question:

624 1. What is the probability of obtaining the measured degree of similarity between
625 the items of questioned- and known-source if they both came from the same
626 source (a source selected at random from the relevant population)?

627 versus

628 2. What is the probability of obtaining the measured degree of similarity between
629 the items of questioned- and known-source if they each came from a different
630 source (each a source selected at random from the relevant population)?

631 It may not be immediately obvious, but similarity-score-based approaches do not
632 properly account for typicality with respect to the relevant population and are therefore

¹⁰ Some published literature has failed to distinguish between the calculation of likelihood ratios based on similarity-scores and the calibration of uncalibrated likelihood ratios (see Morrison, 2013, for an introduction to the latter). The confusion may stem from the fact that uncalibrated log likelihood ratios have usually been called “scores” (and the calibration process has been called “score to likelihood-ratio conversion”), but these scores are not similarity-scores, they are scores that take account of both similarity and typicality with respect to the relevant population.

633 not appropriate for calculating likelihood ratios addressing questions of interest for a
634 court. We briefly demonstrate that similarity-scores do not properly account for
635 typicality with respect to the relevant population (more extensive arguments and
636 demonstrations are presented in Morrison & Enzinger, 2018; Neumann & Ausdemore,
637 2020; and Neumann et al., 2020). In Figure 2a the distance between the two filled
638 circles and the distance between the two non-filled circles is the same, 2 in each case,
639 hence the two pairs of circles will both have the same similarity-score value. The two
640 curves in Figure 3 represent Weibull distributions fitted to same-source scores and to
641 different-source scores, which were obtained by drawing Monte Carlo samples of pairs
642 of same-source feature values and pairs of different-source feature values from the
643 population distribution (consisting of within- and between-source Gaussian
644 distributions with $\mu_r = 0$, $\sigma_b^2 = 100$, and $\sigma_w^2 = 1$), and using unsigned difference as
645 the score function, $\delta(x_q, x_k) = |x_q - x_k|$. Since the two pairs of circles in Figure 2a
646 have the same score value $\delta(x_q, x_k) = 2$, when score-based calculation of likelihood
647 ratios is applied (as graphically illustrated in Figure 3) they both result in the same
648 likelihood-ratio value, which in this example is $0.18/0.058 = 3.1$. Note, however, that
649 in the feature space in Figure 2a, the probability of obtaining the pair of filled-circle
650 feature values if each came from one of two different sources each selected at random
651 from the population is higher than the probability of obtaining the pair of non-filled-
652 circle feature values if each of them came from one of two different sources each
653 selected at random from the population, because the former pair are more typical than
654 the latter pair. Selecting two sources at random, the probability that they will both be
655 in the middle of the distribution is relatively high, whereas the probability that they
656 will both be on the same tail of the distribution is very low. Hence, the likelihood-ratio
657 value associated with the unfilled circles should be higher than the likelihood-ratio
658 value associated with the filled circles, which was the case when the common-source
659 approach was used (19 versus 2.6), but was not the case when the similarity-score-
660 based approach was used (both 3.1), *quod erat demonstrandum*: similarity-score-based
661 approaches do not properly account for typicality with respect to the relevant

662 population.

663 [insert figure about here]

664 **Figure 3.** Example of the calculation of a similarity-score-based likelihood ratio.

665

666 Some published research reports have cited publications such as Morrison & Enzinger
667 (2018), Neumann & Ausdemore (2020), and Neumann et al. (2020), but have then
668 proceeded to use similarity-score-based approaches anyway. Sometimes, they
669 characterize the use of similarity-score-based approaches as entailing “some” loss of
670 information, but we consider the loss of information about typicality with respect to
671 the relevant population to be a loss of essential information, hence we take the position
672 that the use of similarity-score-based approaches are not appropriate for evaluating
673 strength of forensic evidence.

674 **3.3 Validating the performance of a likelihood-ratio system**

675 **3.3.1 Protocol**

676 Protocols for validating systems that output likelihood ratios, and metrics and graphics
677 for representing the results of such validations are described in: Morrison (2011);
678 Meuwly et al. (2017); Ramos et al. (2020); Morrison, Enzinger, et al. (2021).

679 In order to empirically validate the performance of a forensic-evaluation system that
680 outputs likelihood ratios, the tester must input to the system pairs of test items for which
681 it is known that the two items come from the same source and pairs of test items for
682 which it is known that the two items come from different sources. Each pair of items
683 must be sampled from the relevant population, and must reflect the conditions of the
684 questioned-source and known-source items in the case. If there is a mismatch in the
685 conditions of the questioned-source and known-source items, one member of each test
686 pair must reflect the conditions of the questioned-source item and the other must reflect

687 the conditions of the known-source item. If the test pairs do not represent the relevant
 688 population and/or do not reflect the conditions of the case, then the validation results
 689 will not be indicative of the expected performance of the system in the context of the
 690 case. The tester inputs a set of same-source pairs and obtains a set of same-source
 691 likelihood-ratio values as outputs, and inputs a set of different-source pairs and obtains
 692 a set of different-source likelihood-ratio values as outputs.

693 The system being tested must not be told the truth as to which input pairs are same
 694 source and which input pairs are different source. If the system automatically extracts
 695 features and passes the feature values to statistical models, the validation process can
 696 be automated, and the forensic practitioner can act as the tester.

697 3.3.2 Log-likelihood-ratio cost (C_{llr})

698 As previously stated in §3.1, given a same-source input, a good output would be a
 699 likelihood-ratio value that is much larger than 1, a less good output would be a value
 700 that is only a little larger than 1, a bad output would be a value less than 1, and a worse
 701 output would be a value much less than 1. *Mutatis mutandis*, given a different-source
 702 input, a good output would be a value much less than 1. A metric that captures this
 703 gradient goodness is C_{llr} , which is calculated as in Equation (6), in which Λ_s and Λ_d
 704 are likelihood-ratio outputs corresponding to same-source and different-source input
 705 pairs respectively, and N_s and N_d are the number of same-source and different-source
 706 input pairs respectively. C_{llr} is equivalent to the deviance statistic with equal priors.

707 (6)

$$708 \quad C_{llr} = \frac{1}{2} \left(\frac{1}{N_s} \sum_{i=1}^{N_s} \log_2 \left(1 + \frac{1}{\Lambda_{s_i}} \right) + \frac{1}{N_d} \sum_{j=1}^{N_d} \log_2 \left(1 + \Lambda_{d_j} \right) \right)$$

709 Figure 4 plots the logarithmic penalty functions for same-source likelihood-ratio
 710 outputs (the function within the leftmost summation in Equation (6)) and for different-

711 source likelihood-ratio outputs (the function within the rightmost summation in
712 Equation (6)). The x -axis of Figure 4 is scaled in \log_{10} likelihood ratios. If the value of
713 a same-source log likelihood ratio is much greater than 0 it receives a small penalty
714 value, but if its value is lower it receives a higher penalty value. If the value of a
715 different-source log likelihood ratio is much less than 0 it receives a small penalty
716 value, but if its value is higher it receives a higher penalty value. C_{llr} is calculated as
717 the mean of the penalty values with equal weight given to the set of same-source
718 penalty values and the set of different-source penalty values.

719 [insert figure about here]

720 **Figure 4.** Penalty functions used in the calculation of C_{llr} .

721

722 As with classification-error-rate values, lower C_{llr} values indicate better performance,
723 and C_{llr} values cannot be less than 0. A system that gave no useful information and
724 always responded with a likelihood ratio of 1, irrespective of the input, would have a
725 C_{llr} value of 1. A system with a C_{llr} of less than 1 is providing useful information. C_{llr}
726 values substantially greater than 1 can be produced by uncalibrated or miscalibrated
727 systems, but this can be resolved by appropriately calibrating the system.

728 C_{llr} was introduced by Brümmer & du Preez (2006) in the context of automatic speaker
729 recognition. One of the first forensic-science papers to use C_{llr} was González-
730 Rodríguez et al. (2007). Its use is recommended in the *Consensus on validation of*
731 *forensic voice comparison* (Morrison, Enzinger, et al., 2021).

732 3.3.3 Tippett plot

733 Tippett plots consist of plots of the empirical cumulative probability distributions of
734 the same-source likelihood-ratio values and of the different-source likelihood-ratio
735 values. The tradition is to plot lines joining the data points rather than to plot the data

736 points themselves. Examples are shown in Figure 5. The y -axis values corresponding
737 to the curves rising to the right give the proportion of same-source test results with log
738 likelihood-ratio values less than or equal to the corresponding value on the x -axis. The
739 y -axis values corresponding to the curves rising to the left give the proportion of
740 different-source test results with log likelihood-ratio values greater than or equal to the
741 corresponding value on the x -axis. In general, shallower curves with greater separation
742 between the two curves indicates better performance. Tippett plots give an indication
743 of the range of possible likelihood-ratio values that the system could generate under
744 the test conditions, and can also reveal problems such as bias in the output. In Figure
745 5(a) the separation between the same-source and different-source curves is small and
746 the system is clearly biased – both the same-source and different-source curves are too
747 far to the right. The C_{llr} value corresponding to these results is 1.07. Figure 5(b) shows
748 the results of a validation of a better performing system. This Tippett plot has somewhat
749 greater separation between the same-source and different-source curves and the results
750 are not obviously biased. The C_{llr} value corresponding to these results is 0.70. The
751 system that generated the results in Figure 5(b) was actually the same as the system
752 that generated the results in Figure 5(a) except that the Figure 5(b) system included a
753 calibration stage whereas the Figure 5(a) system did not. Figure 5(c) shows the results
754 of a validation of a system with substantially better performance – the same-source and
755 different-source curves have greater separation and are shallower. The C_{llr} value
756 corresponding to these results is 0.31.

757 [insert figure about here]

758 **Figure 5.** Example Tippett plots.

759

760 Tippett plots (i.e., plotting both same-source and different-source empirical cumulative
761 distributions on the same plot) were introduced by Meuwly (2001) in the context of
762 forensic voice comparison. Their use is recommended in the *Consensus on validation*

763 *of forensic voice comparison* (Morrison, Enzinger, et al., 2021).

764

765 **4 Preview of Part II**

766 Element 2 of our strategy (§2.4) is to collaborate with researchers and practitioners on
767 building relevant databases and on developing and validating statistical models
768 applicable in their particular branches of forensic science.

769 Forensic firearm examination is a branch of forensic science in which the new
770 paradigm has so far made almost no progress. A common task in forensic firearm
771 examination is forensic comparison of fired cartridge cases. In Part II of the present
772 paper we provide an example of the application of element 2 of our strategy. In Part II:

773 1. We provide an overview of current practice in forensic comparison of fired
774 cartridge cases.

775 2. We provide a brief review of published research on feature-extraction and
776 statistical-modelling methods applied to forensic comparison of fired cartridge
777 cases.

778 3. We describe the building of a database and the development and validation of
779 feature-extraction and statistical-modelling methods for calculating likelihood
780 ratios for forensic comparison of fired cartridge cases.

781 Part II focuses particularly on the problem of feature extraction, i.e., what are the best
782 features to extract and what region of the base of a fired cartridge case is it best to
783 extract the features from?

784 Part II is intended to provide an example of the application of the new paradigm that
785 can potentially be copied in other branches of forensic science. The authors of Part II
786 include researchers with expertise in machine learning, forensic data science, and

787 forensic firearm examination. The research is conducted in collaboration with forensic
788 practitioners who have contributed advice and provided the fired cartridge cases used
789 to build the database.

790 Part II describes a first attempt to develop a feature-based statistical model for
791 calculating likelihood ratios from 3D digital images of the bases of fired cartridge
792 cases. We will use the common-source approach for calculating likelihood ratios
793 (§3.2.2). An advantage of the common-source approach is that it can be applied when
794 there is a single questioned-source feature vector and a single known-source feature
795 vector, whereas the specific-source approach requires multiple known-source feature
796 vectors (e.g., each extracted from a different item sampled from the same specific
797 known source) in order to train the specific-source model. This means that the
798 common-source approach can also be applied when, rather than having a questioned-
799 source item and one or more known-source items, there are only questioned-source
800 items to be compared with one another: multiple questioned-source items cannot be
801 assumed to have come from the same source, so a specific-source model cannot be
802 trained using multiple feature vectors each obtained from a different questioned-source
803 item. The statistical modelling process will be based on the back-end of a pipeline
804 commonly used in human-supervised-automatic approaches to forensic voice
805 comparison. The features will be extracted using functional data analysis, i.e., they will
806 be values of functions fitted to the 3D-image data. In this research, we will focus
807 particularly on feature extraction, and will test several different feature sets. The
808 validation will make use of protocols, metrics, and graphics commonly used in human-
809 supervised-automatic approaches to forensic voice comparison (§3.3).

810 The database is described in Part II. The current size of the database is relatively small
811 (the COVID-19 pandemic slowed data collection), hence we have kept the statistical-
812 modelling process relatively parsimonious. We plan to continue increasing the size of
813 the database and then conduct additional research to develop and validate more
814 complex models that could potentially have better performance. After the latter stage,

815 we will collaborate with practitioners on field testing a prototype forensic-evaluation
816 system.

817

818 **5 Conclusion**

819 A paradigm shift in evaluation of forensic evidence is ongoing. The shift is away from
820 methods based on human perception and subjective judgement, to methods based on
821 relevant data, quantitative measurements, and statistical models; methods that are
822 transparent and reproducible, that are intrinsically resistant to cognitive bias, that use
823 the logically correct framework for interpretation of evidence (the likelihood-ratio
824 framework), and that are empirically validated under casework conditions. The new
825 paradigm can be called *forensic data science*.

826 Some branches of forensic science, such as forensic voice comparison, are more
827 advanced in the paradigm shift than others. Knowledge gained in advancing the
828 paradigm shift in forensic voice comparison can assist in advancing the paradigm shift
829 in other branches of forensic science. Statistical models used in forensic voice
830 comparison can even be transferred and adapted for use in other branches of forensic
831 science.

832 Our strategy for advancing the paradigm shift requires collaboration between
833 researchers with expertise in forensic data science and researchers and practitioners
834 with expertise in particular branches of forensic science. Our strategy is to work with
835 researchers and practitioners who want to adopt the new paradigm. Our strategy is to
836 work with them on addressing practical impediments to applying the new paradigm in
837 casework, i.e.:

- 838 1. to provide researchers, practitioners, and lawyers with training leading to
839 understanding of the new paradigm;

840 2. to collaborate with researchers and practitioners on building relevant databases
841 and on developing and validating statistical models applicable in their
842 particular branches of forensic science; and

843 3. to conduct research on how to present likelihood ratios and validation results
844 so as to maximize understanding by laypeople, and thereby provide guidance
845 to forensic practitioners on how to communicate forensic-evaluation results to
846 legal-decision makers.

847 In any branch of forensic science, the number of practitioners who initially want to
848 adopt the new paradigm and who want to collaborate on this endeavour will almost
849 certainly be very small, but it will be more productive to work with a small minority
850 on developing practical solutions than to try to convince the majority of practitioners
851 without providing practical solutions. Once the practical solutions are being used by
852 the small minority, use of the new paradigm has the potential to spread. Even then, we
853 do not expect adoption of the new paradigm to be rapid, but we do expect higher rates
854 of adoption among newer practitioners and trainees, leading to a generational shift.

855 In Part II of the present paper, we focus on building a relevant database and developing
856 and validating statistical models for forensic comparison of fired cartridge cases, a
857 common tasks in forensic firearm examination, a branch of forensic science in which
858 the new paradigm has so far made almost no progress. Part II focuses particularly on
859 the problem of feature extraction, i.e., what are the best features to extract and what
860 region of the base of a fired cartridge case is it best to extract the features from? This
861 is a first step in attempting to advance the paradigm shift in that branch of forensic
862 science.

863

864 **6 References**

- 865 Aitken, C.G.G., Berger, C.E.H., Buckleton, J.S., Champod, C., Curran, J.M., Dawid,
866 A.P., Evett, I.W., Gill, P., González-Rodríguez, J., Jackson, G., Kloosterman,
867 A., Lovelock, T., Lucy, D., Margot, P., McKenna, L., Meuwly, D., Neumann,
868 C., Nic Daéid, N., Nordgaard, A., Puch-Solis, R., Rasmusson, B., Redmayne,
869 M., Roberts, P., Robertson, B., Roux, C., Sjerps, M.J., Taroni, F., Tjin-A-Tsoi,
870 T., Vignaux, G.A., Willis, S.M., Zadora, G. (2011) Expressing evaluative
871 opinions: A position statement. *Science & Justice*, 51, 1–2.
872 <http://dx.doi.org/10.1016/j.scijus.2011.01.002>
- 873 Aitken, C.G.G., Lucy, D. (2004) Evaluation of trace evidence in the form of
874 multivariate data. *Applied Statistics*, 53, 109–122.
875 <http://dox.doi.org/10.1046/j.0035-9254.2003.05271.x> [Corrigendum: (2004) 53,
876 665–666. <http://dox.doi.org/10.1111/j.1467-9876.2004.02031.x>]
- 877 Aitken, C.G.G., Roberts, P., Jackson, G. (2010) *Fundamentals of Probability and*
878 *Statistical Evidence in Criminal Proceedings: Guidance for Judges, Lawyers,*
879 *Forensic Scientists and Expert Witnesses*. London, UK: Royal Statistical
880 Society. (Available from [https://rss.org.uk/news-publication/publications/law-](https://rss.org.uk/news-publication/publications/law-guides/)
881 [guides/](https://rss.org.uk/news-publication/publications/law-guides/))
- 882 Association of Forensic Science Providers (2009) Standards for the formulation of
883 evaluative forensic science expert opinion. *Science & Justice*, 49, 161–164.
884 <http://dx.doi.org/10.1016/j.scijus.2009.07.004>
- 885 Bali, A.S., Edmond, G., Ballantyne, K.N., Kemp, R.I., Martire, K.A. (2020)
886 Communicating forensic science opinion: An examination of expert reporting
887 practices. *Science & Justice*, 60, 216–224.
888 <https://doi.org/10.1016/j.scijus.2019.12.005>

- 889 Ballantyne, K., Bunford, J., Found, B., Neville, D., Taylor, D., Wevers, G., Catoggio,
890 D. (2017) *An Introductory Guide to Evaluative Reporting*. National Institute of
891 Forensic Science of the Australia New Zealand Policing Advisory Agency.
892 (Available from [http://www.anzpaa.org.au/forensic-science/our-](http://www.anzpaa.org.au/forensic-science/our-work/projects/evaluative-reporting)
893 [work/projects/evaluative-reporting](http://www.anzpaa.org.au/forensic-science/our-work/projects/evaluative-reporting))
- 894 Basu, N., Bolton-King R., Morrison, G.S. (2022) Advancing a paradigm shift in
895 evaluation of forensic evidence – Part II: Feature-extraction methods for forensic
896 comparison of fired cartridge cases. (Preprint available from
897 <http://firearms.forensic-data-science.net/>)
- 898 Bell, S., Sah, S., Albright, T.D., Gates, S.J., Denton M.B., Casadevall, A. (2018) A
899 call for more science in forensic science. *Proceedings of the National Academy*
900 *of Sciences*, 115, 4541–4544; <https://doi.org/10.1073/pnas.1712161115>
- 901 Berger, C.E.H., Buckleton, J., Champod, C., Evett, I.W., Jackson, G. (2011) Evidence
902 evaluation: A response to the Court of Appeal judgment in R v T. *Science &*
903 *Justice*, 51, 43–49. <http://dx.doi.org/10.1016/j.scijus.2011.03.005>
- 904 Bernstein, D.E. (2013) The misbegotten judicial resistance to the Daubert revolution.
905 *Notre Dame Law Review*, 89, 27–70. (Available from
906 <http://ndlawreview.org/publications/archives/volume-89/issue-1/>)
- 907 Brümmer, N., du Preez, J. (2006) Application independent evaluation of speaker
908 detection. *Computer Speech and Language*, 20, 230–275.
909 <https://doi.org/10.1016/j.csl.2005.08.001>
- 910 Brümmer, N., de Villiers, E. (2010) The speaker partitioning problem. In Proceedings
911 of Odyssey 2010: The speaker and language recognition workshop, pp. 194–
912 201. (Available from [https://www.isca-](https://www.isca-speech.org/archive_open/odyssey_2010/od10_034.html)
913 [speech.org/archive_open/odyssey_2010/od10_034.html](https://www.isca-speech.org/archive_open/odyssey_2010/od10_034.html))

- 914 Cole, S.A. (2006) Is fingerprint identification valid? Rhetorics of reliability in
915 fingerprint proponents' discourse. *Law & Policy*, 28, 109–135.
916 <https://doi.org/10.1111/j.1467-9930.2005.00219.x>
- 917 Cole, S.A. (2009) Forensics without uniqueness, conclusions without
918 individualization: The new epistemology of forensic identification. *Law,*
919 *Probability and Risk*, 8, 233–255. <https://doi.org/10.1093/lpr/mgp016>
- 920 Cole, S.A. (2014) Individualization is dead, long live individualization! Reforms of
921 reporting practices for fingerprint analysis in the United States. *Law, Probability*
922 *and Risk*, 13, 117–150. <https://doi.org/10.1093/lpr/mgt014>
- 923 Cole, S.A., Barno, M. (2020) Probabilistic reporting in criminal cases in the United
924 States: A baseline study. *Science & Justice*, 60, 406–414.
925 <https://doi.org/10.1016/j.scijus.2020.06.001>
- 926 Cooper, G.S., Meterko, V. (2019) Cognitive bias research in forensic science: A
927 systematic review. *Forensic Science International*, 297,35–46.
928 <https://doi.org/10.1016/j.forsciint.2019.01.016>
- 929 Cooper, S.L. (2016) Forensic science identification evidence: Tensions between law
930 and science. *Journal of Philosophy, Science & Law*, 16(2), 1–35.
931 <https://doi.org/10.5840/jpsl20161622>
- 932 Curran, J.M. (2013) Is forensic science the last bastion of resistance against statistics?
933 *Science & Justice*, 53, 251–252. <http://dx.doi.org/10.1016/j.scijus.2013.07.001>
- 934 Edmond, G. (2018) Re-assessing reliability. In *Forensic Science Evidence and Expert*
935 *Witness Testimony* (eds P. Roberts, M. Stockdale), pp. 71–105. Cheltenham,
936 UK: Elgar.

- 937 Edmond, G., Towler, A., Grouns, B., Ribeiro, G., Found, B., White, D., Ballantyne,
938 K., Searston, R.A., Thompson, M.B., Tangen, J.M., Kemp, R.I., Martire, K.A.
939 (2017) Thinking forensics: Cognitive science for forensic practitioners. *Science*
940 & *Justice*, 57, 144–154. <http://dx.doi.org/10.1016/j.scijus.2016.11.005>
- 941 Evett, I.W., Berger, C.E.H., Buckleton, J.S., Champod, C., Jackson, G. (2017)
942 Finding the way forward for forensic science in the US – A commentary on the
943 PCAST report. *Forensic Science International*, 278, 16–23.
944 <http://dx.doi.org/10.1016/j.forsciint.2017.06.018>
- 945 Expert Working Group on Human Factors in Latent Print Analysis (2012) *Latent*
946 *Print Examination and Human Factors: Improving the Practice through a*
947 *Systems Approach*. Gaithersburg, MD: National Institute of Standards and
948 Technology. <https://doi.org/10.6028/NIST.IR.7842>
- 949 Foreman, L.A., Champod, C., Evett, I.W., Lambert, J.A., Pope, S. (2003) Interpreting
950 DNA evidence: A review. *International Statistical Review*, 71, 473–495.
951 <http://dx.doi.org/10.1111/j.1751-5823.2003.tb00207.x>
- 952 Forensic Science Regulator (2020a) Guidance: Cognitive Bias Effects Relevant to
953 Forensic Science Examinations (FSR-G-217 Issue 2). (Available from
954 [https://www.gov.uk/government/publications/cognitive-bias-effects-relevant-to-](https://www.gov.uk/government/publications/cognitive-bias-effects-relevant-to-forensic-science-examinations)
955 [forensic-science-examinations](https://www.gov.uk/government/publications/cognitive-bias-effects-relevant-to-forensic-science-examinations))
- 956 Forensic Science Regulator (2020b). Guidance: Validation (FSR-G-201 Issue 2).
957 (Available from [https://www.gov.uk/government/publications/forensic-science-](https://www.gov.uk/government/publications/forensic-science-providers-validation)
958 [providers-validation](https://www.gov.uk/government/publications/forensic-science-providers-validation))
- 959 Forensic Science Regulator (2021) Codes of Practice and Conduct: Development of
960 Evaluative Opinions (FSR-C-118 Issue 1). (Available from
961 [https://www.gov.uk/government/publications/development-of-evaluative-](https://www.gov.uk/government/publications/development-of-evaluative-opinions)
962 [opinions](https://www.gov.uk/government/publications/development-of-evaluative-opinions))

- 963 Found, B. (2015) Deciphering the human condition: The rise of cognitive forensics.
964 *Australian Journal of Forensic Sciences*, 47, 386–401.
965 <http://dx.doi.org/10.1080/00450618.2014.965204>
- 966 Gold, E., French, J.P. (2011) International practices in forensic speaker comparison.
967 *International Journal of Speech, Language and the Law*, 18, 143–152.
968 <http://dx.doi.org/10.1558/ijssl.v18i2.293>
- 969 Gold, E., French, J.P. (2019) International practices in forensic speaker comparison:
970 Second survey. *International Journal of Speech, Language and the Law*, 26, 1–
971 20. <https://doi.org/10.1558/ijssl.38028>
- 972 Eldridge, H. (2019). Juror comprehension of forensic expert testimony: A literature
973 review and gap analysis. *Forensic Science International: Synergy*, 1, 24–34.
974 <https://doi.org/10.1016/j.fsisyn.2019.03.001>
- 975 House of Lords Science and Technology Select Committee (2019) Forensic Science
976 and the Criminal Justice System: A Blueprint for Change. (Available from
977 <https://publications.parliament.uk/pa/ld201719/ldselect/ldsctech/333/333.pdf>)
- 978 Jackson, G. (2009) Understanding forensic science opinions. In *Handbook of*
979 *Forensic Science* (eds J. Fraser, R. Williams), pp. 419–445. Cullompton, UK:
980 Willan. <https://doi.org/10.4324/9781843927327>
- 981 Kafadar, K., Stern, H., Cuellar, M., Curran, J., Lancaster, M., Neumann, C.,
982 Saunders, C., Weir, B., Zabell, S. (2019) American Statistical Association
983 Position on Statistical Statements for Forensic Evidence. (Available from
984 <https://www.amstat.org/asa/files/pdfs/POL-ForensicScience.pdf>)

- 985 Kenny, P. (2010) Bayesian speaker verification with heavy tailed priors. In
986 Proceedings of Odyssey 2010: The Speaker and Language Recognition
987 Workshop, paper 014. (Available from [https://www.isca-](https://www.isca-speech.org/archive_open/odyssey_2010/od10_014.html)
988 [speech.org/archive_open/odyssey_2010/od10_014.html](https://www.isca-speech.org/archive_open/odyssey_2010/od10_014.html))
- 989 Kaye, D.H. (2015) Presenting forensic identification findings: The current situation.
990 In *Communicating the Results of Forensic Science Examinations: Final*
991 *Technical Report for NIST Award Number 70NANB12H014* (eds C. Neumann,
992 A. Ranadive, D.H. Kaye), pp. 12–30. (Available from
993 <https://ssrn.com/abstract=2690899>)
- 994 Koehler, J.J. (2017) Forensics or fauxrensic? Ascertaining accuracy in the forensic
995 sciences. *Arizona State Law Journal*, 49(4), 1369–1416. (Available from
996 [https://arizonastatelawjournal.org/2018/02/07/forensics-or-fauxrensic-](https://arizonastatelawjournal.org/2018/02/07/forensics-or-fauxrensic-ascertaining-accuracy-in-the-forensic-sciences/)
997 [ascertaining-accuracy-in-the-forensic-sciences/](https://arizonastatelawjournal.org/2018/02/07/forensics-or-fauxrensic-ascertaining-accuracy-in-the-forensic-sciences/))
- 998 Kuhn, T.S. (1962) *The Structure of Scientific Revolutions*. Chicago IL: University of
999 Chicago Press.
- 1000 Lee, K.A., Yamamoto, H., Okabe, K., Wang, Q., Guo, L., Koshinaka, T., Zhang, J.,
1001 Shinoda, K. (2020) NEC-TT System for mixed-bandwidth and multi-domain
1002 speaker recognition. *Computer Speech, Language*, 61, article 101033.
1003 <https://doi.org/10.1016/j.csl.2019.101033>
- 1004 Margot, P. (2011) Commentary on the need for a research culture in the forensic
1005 sciences. *UCLA Law Review*, 58, 795–801. (Available from
1006 [https://www.uclalawreview.org/commentary-on-the-need-for-a-research-culture-](https://www.uclalawreview.org/commentary-on-the-need-for-a-research-culture-in-the-forensic-sciences-3-2/)
1007 [in-the-forensic-sciences-3-2/](https://www.uclalawreview.org/commentary-on-the-need-for-a-research-culture-in-the-forensic-sciences-3-2/))

- 1008 Matějka, P., Plchot, O., Glembek, O., Burget, L., Rohdin, J., Zeinali, H., Mošner, L.,
1009 Silnova, A., Novotný, O., Diez, M., Černocký, J.H. (2020) 13 years of speaker
1010 recognition research at BUT, with longitudinal analysis of NIST SRE. *Computer*
1011 *Speech & Language*, 63, article 101035.
1012 <https://doi.org/10.1016/j.csl.2019.101035>
- 1013 Meuwly, D. (2001) Reconnaissance de locuteurs en sciences forensiques: l'apport
1014 d'une approche automatique. Doctoral dissertation, University of Lausanne.
1015 (Available from
1016 <https://www.unil.ch/files/live/sites/esc/files/shared/These.Meuwly.pdf>)
- 1017 Meuwly, D., Ramos, D., Haraksim, R. (2017) A guideline for the validation of
1018 likelihood ratio methods used for forensic evidence evaluation. *Forensic Science*
1019 *International*, 276, 142–153. <http://dx.doi.org/10.1016/j.forsciint.2016.03.048>
- 1020 Mnookin, J.L., Cole, S.A., Dror, I.E., Fisher, B.A.J., Houck, M.M., Inman, K., Kaye,
1021 D.H., Koehler, J.J., Langenburg, G., Risinger, D.M., Rudin, N., Siegel, J.,
1022 Stoney, D.A. (2011) The need for a research culture in the forensic sciences.
1023 *UCLA Law Review*, 58, 725–777. (Available from
1024 [https://www.uclalawreview.org/the-need-for-a-research-culture-in-the-forensic-](https://www.uclalawreview.org/the-need-for-a-research-culture-in-the-forensic-sciences-2/)
1025 [sciences-2/](https://www.uclalawreview.org/the-need-for-a-research-culture-in-the-forensic-sciences-2/))
- 1026 Morgan, R.M., Levin, E.A. (2019) A crisis for the future of forensic science: Lessons
1027 from the UK of the importance of epistemology for funding research and
1028 development. *Forensic Science International: Synergy*, 1, 243–252.
1029 <https://doi.org/10.1016/j.fsisyn.2019.09.002>
- 1030 Morrison, G.S. (2011) Measuring the validity and reliability of forensic likelihood-
1031 ratio systems. *Science & Justice*, 51, 91–98.
1032 <http://dx.doi.org/10.1016/j.scijus.2011.03.002>

- 1033 Morrison, G.S. (2012) The likelihood-ratio framework and forensic evidence in court:
1034 A response to R v T. *International Journal of Evidence and Proof*, 16, 1–29.
1035 <http://dx.doi.org/10.1350/ijep.2012.16.1.390>
- 1036 Morrison, G.S. (2013). Tutorial on logistic-regression calibration and fusion:
1037 Converting a score to a likelihood ratio. *Australian Journal of Forensic*
1038 *Sciences*, 45, 173–197. <http://dx.doi.org/10.1080/00450618.2012.733025>
- 1039 Morrison, G.S. (2014) Distinguishing between forensic science and forensic
1040 pseudoscience: Testing of validity and reliability, and approaches to forensic
1041 voice comparison. *Science & Justice*, 54, 245–256.
1042 <http://dx.doi.org/10.1016/j.scijus.2013.07.004>
- 1043 Morrison, G.S. (2017) What should a forensic practitioner’s likelihood ratio be? II.
1044 *Science & Justice*, 57, 472–476. <http://dx.doi.org/10.1016/j.scijus.2017.08.004>
- 1045 Morrison, G.S. (2018a) Admissibility of forensic voice comparison testimony in
1046 England and Wales. *Criminal Law Review*, 2018(1), 20–33. (Available from
1047 http://geoff-morrison.net/#Admissibility_EW_2018)
- 1048 Morrison, G.S. (2018b) The impact in forensic voice comparison of lack of
1049 calibration and of mismatched conditions between the known-speaker recording
1050 and the relevant-population sample recordings. *Forensic Science International*,
1051 283, e1–e7. <http://dx.doi.org/10.1016/j.forsciint.2017.12.024>
- 1052 Morrison, G.S., Ballantyne, K., Geoghegan, P.H. (2018) A response to Marquis et al
1053 (2017) What is the error margin of your signature analysis? *Forensic Science*
1054 *International*, 287, e11–e12. <https://doi.org/10.1016/j.forsciint.2018.03.009>

- 1055 Morrison, G.S., Enzinger, E. (2018) Score based procedures for the calculation of
1056 forensic likelihood ratios – Scores should take account of both similarity and
1057 typicality. *Science & Justice*, 58, 47–58.
1058 <http://dx.doi.org/10.1016/j.scijus.2017.06.005>
- 1059 Morrison, G.S., Enzinger, E., Hughes, V., Jessen, M., Meuwly, D., Neumann, C.,
1060 Planting, S., Thompson, W.C., van der Vloed, D., Ypma, R.J.F., Zhang, C.,
1061 Anonymous, A., Anonymous, B. (2021) Consensus on validation of forensic
1062 voice comparison. *Science & Justice*, 61, 229–309.
1063 <https://doi.org/10.1016/j.scijus.2021.02.002>
- 1064 Morrison, G.S., Enzinger, E., Ramos, D., González-Rodríguez, J., Lozano-Díez, A.
1065 (2020) Statistical models in forensic voice comparison. In *Handbook of Forensic*
1066 *Statistics* (eds D. Banks, K. Kafadar, D.H. Kaye, M. Tackett), pp. 451–497.
1067 Boca Raton, FL: CRC. <https://doi.org/10.1201/9780367527709>
- 1068 Morrison, G.S., Kaye, D.H., Balding, D.J., Taylor, D., Dawid, P., Aitken, C.G.G.,
1069 Gittelsohn, S., Zadora, G., Robertson, B., Willis, S.M., Pope, S., Neil, M.,
1070 Martire, K.A., Hepler, A., Gill, R.D., Jamieson, A., de Zoete, J., Ostrum, R.B.,
1071 Caliebe, A. (2017) A comment on the PCAST report: Skip the “match”/“non-
1072 match” stage. *Forensic Science International*, 272, e7–e9.
1073 <http://dx.doi.org/10.1016/j.forsciint.2016.10.018>
- 1074 Morrison, G.S., Neumann, C., Geoghegan, P.H. (2020) Vacuous standards –
1075 subversion of the OSAC standards-development process. *Forensic Science*
1076 *International: Synergy*, 2, 206–209. <https://doi.org/10.1016/j.fsisyn.2020.06.005>

- 1077 Morrison, G.S., Neumann, C., Geoghegan, P.H., Edmond, G., Grant, T., Ostrum,
1078 R.B., Roberts, P., Saks, M., Syndercombe Court, D., Thompson, W.C., Zabell,
1079 S. (2021) Reply to Response to Vacuous standards – subversion of the OSAC
1080 standards-development process. *Forensic Science International: Synergy*, 3,
1081 article 100149. <https://doi.org/10.1016/j.fsisyn.2021.100149>
- 1082 Morrison, G.S., Sahito, F.H., Jardine, G., Djokic, D., Clavet, S., Berghs, S., Goemans
1083 Dorny, C. (2016) INTERPOL survey of the use of speaker identification by law
1084 enforcement agencies. *Forensic Science International*, 263, 92–100.
1085 <http://dx.doi.org/10.1016/j.forsciint.2016.03.044>
- 1086 Morrison, G.S., Stoel, R.D. (2014) Forensic strength of evidence statements should
1087 preferably be likelihood ratios calculated using relevant data, quantitative
1088 measurements, and statistical models – a response to Lennard (2013) Fingerprint
1089 identification: How far have we come? *Australian Journal of Forensic Sciences*,
1090 46, 282–292. <http://dx.doi.org/10.1080/00450618.2013.833648>
- 1091 Morrison, G.S., Thompson, W.C. (2017) Assessing the admissibility of a new
1092 generation of forensic voice comparison testimony. *Columbia Science and
1093 Technology Law Review*, 18, 326–434. <https://doi.org/10.7916/stlr.v18i2.4022>
- 1094 National Research Council of the National Academies (2009) *Strengthening Forensic
1095 Science in the United States: A Path Forward*. Washington, DC: National
1096 Academies Press. <https://doi.org/10.17226/12589>
- 1097 Neumann, C., Ausdemore, M. (2020) Defence against the modern arts: The curse of
1098 statistics – Part II: ‘Score-based likelihood ratios’. *Law, Probability and Risk*,
1099 19, 21–42. <http://dx.doi.org/10.1093/lpr/mgaa006>

- 1100 Neumann, C., Hendricks, J., Ausdemore, M. (2020) Statistical support for
1101 conclusions in fingerprint examinations. In *Handbook of Forensic Statistics* (eds
1102 D. Banks, K. Kafadar, D.H. Kaye, M. Tackett), pp. 277–324. Boca Raton, FL:
1103 CRC. <https://doi.org/10.1201/9780367527709>
- 1104 Ommen, D.M., Saunders, C.P. (2021) A problem in forensic science highlighting the
1105 differences between the Bayes factor and likelihood ratio. *Statistical Science*, 36,
1106 344–359. <https://doi.org/10.1214/20-STS805>
- 1107 President’s Council of Advisors on Science and Technology (2016) Forensic Science
1108 in Criminal Courts: Ensuring Scientific Validity of Feature-Comparison
1109 Methods. (Available from
1110 <https://obamawhitehouse.archives.gov/administration/eop/ostp/pcast/docsreports>
1111 /)
- 1112 Prince, S.J.D., Elder, J.H. (2007). Probabilistic linear discriminant analysis for
1113 inferences about identity. In Proceedings of the IEEE 11th International
1114 Conference on Computer Vision, pp. 1–8.
1115 <https://doi.org/10.1109/ICCV.2007.4409052>
- 1116 Ramos, D., Meuwly, D., Haraksim, R., Berger, C.E.H. (2020) Validation of forensic
1117 automatic likelihood ratio methods. In *Handbook of Forensic Statistics* (eds D.
1118 Banks, K. Kafadar, D.H. Kaye, M. Tackett), pp. 143–163. Boca Raton, FL:
1119 CRC. <https://doi.org/10.1201/9780367527709>
- 1120 Redmayne, M., Roberts, P., Aitken, C.G.G., Jackson, G. (2011). Forensic science
1121 evidence in question. *Criminal Law Review*, 2011(5), 347–356. (Available from
1122 <https://rke.abertay.ac.uk/en/publications/forensic-science-evidence-in-question>)
- 1123 Risinger, D.M. (2013) Reservations about likelihood ratios (and some other aspects
1124 of forensic ‘Bayesianism’). *Law, Probability and Risk*, 12, 63–73,
1125 <http://dx.doi.org/10.1093/lpr/mgs011>

- 1126 Roux, C., Weyermann, C. (2021) From research integrity to research relevance to
1127 advance forensic science. *Forensic Sciences Research*.
1128 <https://doi.org/10.1080/20961790.2021.1977480>
- 1129 Roux, C., Willis, S., Weyermann, C. (2021) Shifting forensic science focus from
1130 means to purpose: A path forward for the discipline? *Science & Justice*, 61,
1131 678–686. <https://doi.org/10.1016/j.scijus.2021.08.005>
- 1132 Saks, M.J., Koehler, J.J. (2005) The coming paradigm shift in forensic identification
1133 science. *Science*, 309, 892–895. <https://doi.org/10.1126/science.1111565>
- 1134 Saks, M.J., Koehler, J.J. (2008) The individualization fallacy in forensic science.
1135 *Vanderbilt Law Review*, 61 199–219. (Available from
1136 <https://ssrn.com/abstract=1432516>)
- 1137 Stoel, R.D., Berger, C.E.H., Kerkhoff, W., Mattijssen, E.J.A.T., Dror, E.I. (2015)
1138 Minimizing contextual bias in forensic casework. In *Forensic Science and the*
1139 *Administration of Justice: Critical Issues and Directions* (eds K.J. Strom, M.J.
1140 Hickman), pp. 67–86. Thousand Oaks CA: Sage.
1141 <http://dx.doi.org/10.4135/9781483368740.n5>
- 1142 Swofford, H., Champod, C. (2021) Implementation of algorithms in pattern &
1143 impression evidence: A responsible and practical roadmap. *Forensic Science*
1144 *International: Synergy*, 3, article 100142.
1145 <https://doi.org/10.1016/j.fsisyn.2021.100142>
- 1146 Swofford, H., Cole, S., King, V. (2021) Mt. Everest – we are going to lose many: A
1147 survey of fingerprint examiners’ attitudes towards probabilistic reporting. *Law,*
1148 *Probability and Risk*, 19, 255–291. <https://doi.org/10.1093/lpr/mgab003>
- 1149 Thompson, W.C. (2012) Discussion paper: Hard cases make bad law – reactions to R
1150 v T. *Law, Probability and Risk*, 11, 347–359. <https://doi.org/10.1093/lpr/mgs020>

- 1151 Thompson, W.C., Black J., Jain A., Kadane J. (2017) Forensic Science Assessments:
1152 A Quality and Gap Analysis- Latent Fingerprint Examination. Washington, DC:
1153 American Association for the Advancement of Science.
1154 <https://doi.org/10.1126/srhrl.aag2874>
- 1155 Thompson, W.C., Schumann, E.L. (1987) Interpretation of statistical evidence in
1156 criminal trials: The prosecutor's fallacy and the defense attorney's fallacy. *Law*
1157 *and Human Behavior*, 11(3), 167–187. <http://www.jstor.org/stable/1393631>
- 1158 Villalba, J., Chen, N., Snyder, D., García-Romero, D., McCree, A., Sell, G.,
1159 Borgstrom, J., García-Perera, L.P., Richardson, F., Dehak R., Torres-
1160 Carrasquillo, P.A., Dehak, N. (2020) State-of-the-art speaker recognition with
1161 neural network embeddings in NIST SRE18 and Speakers in the Wild
1162 evaluations. *Computer Speech & Language*, 60, article 101026.
1163 <https://doi.org/10.1016/j.csl.2019.101026>
- 1164 Weber, P., Enzinger, E., Labrador-Serrano, B., Lozano-Díez, A., Ramos, D.,
1165 González-Rodríguez, J., Morrison, G.S. (2022 in press) Forensic voice
1166 comparison – Human-supervised-automatic approach. In *Encyclopedia of*
1167 *Forensic Sciences* (3rd ed). Elsevier.
- 1168 Willis, S.M., McKenna, L., McDermott, S., O'Donnell, G., Barrett, A., Rasmusson,
1169 A., Nordgaard, A., Berger, C.E.H., Sjerps, M.J., Lucena-Molina, J.J., Zadora,
1170 G., Aitken, C.G.G., Lunt, L., Champod, C., Biedermann, A., Hicks, T.N.,
1171 Taroni, F. (2015) ENFSI Guideline for Evaluative Reporting in Forensic
1172 Science. (Available from [http://enfsi.eu/wp-](http://enfsi.eu/wp-content/uploads/2016/09/m1_guideline.pdf)
1173 [content/uploads/2016/09/m1_guideline.pdf](http://enfsi.eu/wp-content/uploads/2016/09/m1_guideline.pdf))









