



Deposited via The University of York.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/236015/>

Version: Accepted Version

Article:

HUGHES, VINCENT, LLAMAS, CARMEN and Kettig, Thomas (2026) The effects of expert conclusions and additional forensic evidence on listener evaluations in a voice comparison task. *International Journal of Speech, Language and the Law*. ISSN: 1748-8885

<https://doi.org/10.3138/ijsl-2025-0022>

Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

1 The effects of expert conclusions and additional forensic evidence on listener
2 evaluations in a voice comparison task

3
4 Abstract

5
6 Little is known about how lay people understand and evaluate forensic evidence. Yet in
7 many jurisdictions, lay people make up juries that are responsible for determining the
8 innocence or guilt of the accused. In this paper, we examine how lay people's
9 evaluations of voice comparison evidence are affected by likelihood ratio (LR)-based
10 conclusions of different magnitudes, expressed in the form of numerical and verbal
11 statements, as well as the suggestion of the existence of additional forensic evidence in
12 a mock case. We use a novel "jury-of-the-future" game to elicit participant judgments
13 about whether pairs of voices belong to the same or different speakers. We then assess
14 how those judgments are affected by expert conclusions and additional forensic
15 evidence. Results show that verbal LRs shifted participants' responses in a consistent
16 way relative to the strength of the conclusion presented (i.e. people thought the voices
17 were more likely to be the same if the LR supported the prosecution). This effect was
18 stronger for voices from the North East of England rather than Standard Southern
19 British English (sometimes referred to as Received Pronunciation). Less consistent
20 patterns were found for numerical LRs, suggesting that participants struggled to
21 interpret numerical values, particularly those of the highest magnitude. When primed
22 with the suggestion of additional forensic evidence (DNA, footprint, fingerprint),
23 participant responses shifted towards the voices being more likely to be the same, for
24 both same- and different-speaker pairs, with the effect being strongest for footprint and
25 fingerprint evidence. We discuss the implications of these findings for the presentation
26 of conclusions in forensic voice comparison cases, and for forensic science more
27 generally.

28
29 Keywords: evidence evaluation, voice comparison, likelihood ratio, lay listeners,
30 additional evidence

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

47 1. Introduction

48
49 Forensic voice comparison (FVC) typically involves the analysis of a recording of the
50 voice of an unknown criminal and a recording of the voice of a known suspect (usually
51 in a police interview). A forensic expert is tasked with comparing the speech features in
52 the two recordings to assess the similarity between the voices and their typicality
53 relative to a wider relevant population. In doing so, the analyst evaluates the relative
54 strength of the evidence considering the prosecution and defence propositions. This
55 conclusion is presented to an end-user in the form of a written report and, in some
56 cases, oral testimony delivered to a trier-of-fact in court; in many jurisdictions around the
57 world, this is a jury made up of lay people. It is then the responsibility of the trier-of-fact
58 to evaluate the specific piece of FVC evidence in combination with the other evidence in
59 the case to assess the *ultimate issue* (Lynch and McNally 2003) of the defendant's
60 innocence or guilt.

61
62 Communicating forensic science to end-users, particularly jurors in court, remains a
63 significant challenge. By definition, forensic expertise is requested for evidence which is
64 “likely to be outside the experience and knowledge of a judge or jury.”¹ Therefore,
65 effective communication is essential to ensure that end-users can understand the
66 analytical process, develop trust in the methodology being employed, and assign the
67 appropriate weight to the evidence. Ineffective communication potentially risks
68 miscarriages of justice, irrespective of the accuracy of the scientific work underpinning
69 the analysis. Yet voice is an unusual form of forensic evidence in that non-experts can
70 listen to a pair of recordings and arrive at their own view about whether they contain the
71 voice of the same speaker or not. Voice is therefore an accessible form of evidence to
72 lay people in a way that other forensic evidence is not (e.g. DNA, forensic chemistry,
73 and even fingerprints). This is potentially problematic since recordings are often played
74 to jurors in criminal trials, in addition to being presented with expert evidence in the form
75 of a written report and/or oral testimony.

76
77 Despite its complexity and importance, little is known about the decisions made by
78 juries and the processes that underlie those decisions. In part, this is because it is not
79 possible to access jury deliberations in real cases, or to question jurors after the fact. In
80 this paper, we investigate the jury decision-making process in the context of FVC
81 evidence. Specifically, we consider the extent to which lay people's evaluations of
82 whether two voices are the same or different are affected by expert conclusions of
83 various forms and magnitudes, and the suggestion of additional forensic evidence in a
84 case.

85 86 1.1 Conclusion frameworks for forensic voice comparison evidence

87
88 There is now considerable consensus within the forensic science community that the
89 likelihood ratio (LR) is the appropriate framework for evaluating and reporting the
90 strength of expert evidence (Aitken et al. 2011, Berger et al. 2011, Redmayne et al.
91 2011). Despite this, there is still variability in terms of the frameworks used in practice to

¹ <https://www.cps.gov.uk/legal-guidance/expert-evidence> (last viewed: 05/06/25)

92 express conclusions, with many forensic disciplines and jurisdictions around the world
93 not using the LR framework at all. Where the LR is used, there is considerable
94 variability in how those LRs are computed and then expressed. The LR can be a
95 numerical value computed in an entirely quantitative way using statistical models and
96 data that are representative of the conditions of the case. However, for many
97 disciplines, principally where analysis is based on expert human judgment, it is not
98 possible to use data and statistical models to compute a numerical LR (Evetts 1998).
99 This may be because representative data do not exist, the analytic methods do not
100 provide easily quantifiable data, or because generating large-scale data is extremely
101 labour intensive and time consuming. In such disciplines, experts may still use a
102 number to represent the LR, using more subjective estimates to produce the values for
103 the numerator and denominator. More commonly, experts use a verbal statement which
104 expresses the extent to which the evidence supports one proposition over the other
105 based on categorical labels (e.g. limited/moderate/strong support). Numerical LRs
106 derived from an entirely data-driven approach may also be converted into a verbal
107 equivalent using an established scale (see Champod and Evett 2000, Association for
108 Forensic Science Providers 2009).

109
110 Within the field of FVC, there has long been debate about the appropriate framework to
111 use for evaluating the strength of evidence. In many ways FVC has been ahead of
112 many other forensic sciences in these discussions, with the first work outlining issues
113 with posterior probability published in the 1990s and early 2000s (see Broeders 1999,
114 Champod and Evett 2000, Rose 2002). In the late 2000s, many FVC experts adopted
115 what is known as the *UK Position Statement* (French and Harrison 2007) which involves
116 making separate judgments about the similarity (termed *consistency* in French and
117 Harrison 2007) and typicality of the voices, but without a single expression of the overall
118 strength of evidence. This approach was criticised by Rose and Morrison (2009) for
119 ultimately falling short of using the LR framework. Since then, most experts have moved
120 to a verbal LR scale with conclusions expressed in the form of a support statement for
121 one proposition over another. However, there is considerable demand for experts to use
122 numerical LRs, and to compute those LRs using representative data and statistical
123 models (Morrison 2014, Morrison 2022). Indeed, this will become more common in
124 casework as the use of automatic speaker recognition (ASR) systems (which output
125 numerical LRs) increases.

126 127 1.2 Jury understanding of expert evidence

128
129 Previous work has considered the extent to which potential jurors understand forensic
130 conclusions expressed in different ways. Some have argued against the use of
131 numerical LRs, with Aitken et al. (2011: 1) stating that verbal expressions are “the most
132 appropriate basis for communication”. Arscott et al. (2017) examined how verbal
133 statements are interpreted by three groups with differing levels of experience of the
134 criminal justice system: lay people, legal professionals, and those with forensic/
135 investigative knowledge. All three groups increased their evaluations of strength of
136 evidence in line with the strength of the verbal expression. However, participants
137 generally struggled to interpret differences at the highest end of the scale.

138

139 Further empirical work has reported other issues with verbal LR statements. Martire et
140 al. (2013) presented lay people with three verbal LR statements of “weak or limited”,
141 “moderately strong” and “very strong” support. These categories correspond to
142 numerical LRs of 4.5, 450, and 495,000, following the Association of Forensic Science
143 Providers (2009) scale. They report three key findings. Firstly, people respond to shifts
144 in the weight of evidence in the expected direction. That is, changes in belief towards
145 guilt are greatest for *very strong* support and smallest for *weak or limited* support.
146 Secondly, despite this, the magnitude of peoples’ belief change is generally smaller
147 than would be expected from the magnitude indicated by the verbal statement. This
148 potentially undervalues the strength of evidence. Thirdly, results demonstrate the *weak*
149 *evidence effect* (see Fernbach et al. 2011), whereby belief shifts away from guilt (i.e.
150 towards the defence proposition) when presented with a conclusion of “weak or limited”
151 evidence in support of the prosecution proposition. Similar patterns are replicated in
152 Mullen et al. (2014) who found little correlation between verbal statements and
153 perceived strength of evidence irrespective of whether participants were presented with
154 fibre, footwear, glass, or paint evidence. This leads to the conclusion that verbal
155 statements are “unlikely to be understood properly by lay people and it would appear
156 that they are actually misunderstood” (p.154).

157

158 Martire et al. (2014) analysed the understanding of LRs comparing both verbal
159 and numerical conclusions, and the use of a visual scale. Participants were presented
160 with a small amount of information about a burglary and asked their view on the guilt of
161 the accused. Subsequently, they were presented with high (*strong*; LR = 5500 times
162 more likely) and low (*limited*; LR = 5.5 times more likely) magnitude LRs relating to
163 fingerprint evidence and asked again about their view on guilt. The difference between
164 the two responses was used as a measure of belief change. Results showed that
165 numerical LRs “produce belief-change and implicit LRs most commensurate with the
166 intentions of the expert” and “most resistant to weak evidence effect” (p. 61).
167 Meanwhile, the correlation was much weaker for verbal LRs, which were also most
168 affected by the *weak evidence effect*. A similar pattern was reported for the combination
169 of verbal and numerical conclusions.

170

171 1.3 Factors affecting jury decision making

172

173 There are a variety of factors which potentially affect jury decision making. Jurors bring
174 with them a unique set of life experiences and skills, which may vary according to their
175 level of education or knowledge of the judicial process. This in turn may affect the level
176 of understanding about forensic science, statistics and probability. Jurors also have
177 inherently different levels of trust in legal institutions, which in some cases may come
178 from direct interactions with the judicial system. Jurors also necessarily have a unique
179 set of implicit and explicit biases, based on ethnicity, gender, sexuality, or other
180 personal characteristics, which may affect the extent to which they are likely to consider
181 someone innocent or guilty, or the extent to which they trust the evidence from a
182 specific witness or expert (see Curley et al. 2022). As highlighted by Bornstein and
183 Greene (2011: 64-65), the reason why individual jurors may arrive at different overall

184 decisions is that they “filter the evidence through their own experiences, expectations,
185 values, and beliefs”.

186
187 Some of those biases may affect the way in which jurors perceive a defendant, and the
188 voice has been shown to be a cue to those perceptions (see Frumkin and Stone 2020).
189 Paver et al. (2025) found that lay people generally associate certain regional accents of
190 the UK as more likely to be linked with criminal activity (a finding which replicates Dixon
191 et al. 2002 amongst others). Criminality ratings on a Likert scale for Birmingham and
192 Glasgow were generally lower than for other regional accents, such as Liverpool,
193 London and Newcastle. Standard Southern British English (also referred to as RP)
194 speakers were generally rated as less likely to behave in criminal ways, and those
195 speakers were also more associated with status traits such as intelligence, education
196 and wealth. Interestingly, ratings of criminal activity varied as a function of crime type,
197 with the regional accents rated as more likely to be involved in crimes such as physical
198 assault and shoplifting, but not sexual assault.

199
200 Other evidence in the case may also affect juror judgment (see Munro et al. 2024 for
201 exploration about how information outside of the trial itself can affect juror decisions).
202 For example, knowing that there is a lot of forensic evidence in the case may make
203 jurors attach more weight to weaker forensic evidence, or even to dismiss contradictory
204 evidence. Of course, other evidence in a case is not task-irrelevant for jurors. After all, it
205 is their job to make a decision about innocence and guilt based on all of the evidence
206 presented. However, Carlson and Russo (2001) refer to the concept of *predecisional*
207 *distortion* whereby jurors attach more weight to new evidence in a trial that supports the
208 verdict which is favoured by previous evidence. Other effects are also possible such as
209 recency bias, whereby more weight is attached to the most recent piece of evidence,
210 and anchoring bias, whereby a specific piece of earlier evidence is used as a basis for
211 assessing the weight of other evidence (Tindale and Berryman 2023).

212 213 1.4 This study

214
215 The aim of the present study is to evaluate the effects of expert conclusions and the
216 suggestion of additional forensic evidence in a mock case on lay peoples’ responses
217 in an FVC task. More specifically, we want to know whether these sources of external
218 information (i.e. factors not related to the voices themselves) can change listeners’
219 evaluations of voice sameness, as well as the direction and magnitude of any change,
220 and how this interacts with the accent of the speaker. To do this, we asked lay people to
221 listen to pairs of voice recordings of speakers of Standard Southern British English,
222 Newcastle English and Middlesbrough English (two localities in the North East of
223 England) and then to assess the likelihood that each pair contained the same or
224 different speakers. Data collection was performed using a bespoke jury-based game
225 allowing us to test the following research questions:

- 226
227 ● To what extent are listener judgments in an FVC task affected by a conclusion
228 provided by a forensic expert?

- 229
- Do listeners better understand verbal or numerical LR conclusions? And how
- 230 does this interact with the magnitude of the evidence?
- Are listeners more likely to think voices belong to the same speaker (i.e. more
- 231 support for guilt) when primed with the suggestion of other forensic evidence in
- 232 the case?
- 233
- 234

235 While there is a body of existing research on lay peoples' understanding of expert

236 conclusion frameworks, previous work has not assessed this issue with regard to voice

237 evidence. More crucially, previous work has not assessed the understanding of

238 conclusion frameworks when lay people are in a position to evaluate the evidence

239 themselves, as is the case for voice evidence, but not for other forensics, such as DNA.

240

241 2. Methods

242 2.1 Stimuli

243

244

245 The voice recordings used as stimuli in this study were taken from two corpora. The first

246 was the Dynamic Variability in Speech (DyViS; Nolan et al. 2009) corpus, which

247 contains adult (18-25 years) male speakers of Standard Southern British English

248 (SSBE). The speakers were all students at the University of Cambridge at the time of

249 the recordings. The second corpus was collected as part of The Use and Utility of

250 Localised Speech Forms in Determining Identity project (TUULS; Llamas et al. 2016-

251 19). The speakers used in our study were all adult (18-65 years) male speakers from

252 Newcastle or Middlesbrough in the North East of England. The three accents in our

253 study were chosen to test for effects of familiarity on listeners' FVC performance. SSBE

254 is an accent with high familiarity for all UK-based listeners. Newcastle English is also

255 likely to be familiar to listeners, given its high exposure in the media. Middlesbrough

256 English is segmentally similar to Newcastle English, but with lower levels of familiarity

257 (for more, see Hughes et al. 2025). A total of 45 SSBE speakers, 15 Middlesbrough

258 speakers, and 15 Newcastle speakers were used in the present study.

259

260 Stimuli were created from longer recordings which were representative of the conditions

261 of typical FVC casework. Both DyViS and TUULS elicited speech through a mock police

262 interview in which speakers had to lie about their involvement in a crime. These police

263 interviews were all high-quality studio recordings. For the SSBE speakers, DyViS

264 contains a landline telephone recording of a separate conversation with a supposed

265 accomplice where each speaker discussed the mock crime again with an 'accomplice'.

266 No 'accomplice task' speech was available from TUULS, so voice samples were taken

267 from sociolinguistic interviews with the Middlesbrough and Newcastle speakers.

268 Landline telephone-like quality was mimicked by applying a band-pass filter (300-

269 3400Hz, with 100Hz smoothing at the edges of the band-pass range) and adding 25dB

270 of white noise. From each recording, edited samples of around 10 seconds were

271 created to use as stimuli. Identifying information and semantic content about the

272 background of the speakers was avoided in creating the edits. Norming of the samples

273 ensured that the content taken from all of the TUULS and DyViS tasks were similar in

274 style and no stimulus used in the listening task contained any overtly criminal or
275 particularly guilty-sounding content.

276
277 Speech samples were paired to create same-speaker and different-speaker
278 comparisons. In each case, the telephone sample was used as the nominal ‘criminal’
279 sample (i.e. the ‘unknown’ voice) and the high-quality interview sample was used as the
280 nominal ‘suspect’ sample (i.e. the ‘known’ voice). This arrangement of samples
281 replicates the most common situation in FVC casework. In total, the set of stimuli used
282 in this study consisted of 45 same-speaker comparisons and 75 different-speaker
283 comparisons. Most comparisons were done within accent group (i.e. SSBE, Newcastle
284 or Middlesbrough), with the exception of 30 different-speaker comparisons involving a
285 Newcastle voice and a Middlesbrough voice. This allowed us to test the sensitivity of
286 listeners to regional accent variation.

287 288 2.2 Listener task

289
290 We asked our listeners to conduct something similar in design to an FVC task. This
291 involved first listening to the nominal ‘criminal’ (telephone) sample and providing a
292 judgment about the *typicality* (i.e. the distinctiveness) of the voice relative to other
293 speakers of the same accent. In each case, the listeners were told the regional accent
294 of this speaker. Listeners provided their typicality judgments on a 0 to 100 scale. They
295 were then presented with the nominal ‘suspect’ (interview) sample and asked to provide
296 a judgment about the *similarity* between the two voices on a 0 to 100 scale. Listeners
297 also provided a judgment about whether they thought the voices belonged to the same
298 speaker or not – that is, their *sameness* – again using a 0 to 100 scale. In Hughes et al.
299 (2022), we computed LR-like scores from the similarity and typicality ratings to compare
300 the performance of humans and automatic systems directly. However, in the present
301 study, we focus on the results of the *sameness* ratings only, as this is more directly
302 related to the assignment of guilt in a jury context.

303
304 Every listener was presented with a total of 24 comparisons, divided into three blocks of
305 eight comparisons. The blocks had been curated to ensure that they consisted of three
306 same-speaker comparisons and five different-speaker comparisons, split equally across
307 the accent groups. Both the order of the blocks and the order of the comparisons within
308 blocks were randomised across listeners. The experiment was a between-subjects
309 design: each listener responded only to a subset (24) of the total number of
310 comparisons in the study (120). This ensured that, outside of the context of same-
311 speaker comparisons, listeners never heard the same voice or sample more than once
312 throughout the experiment. Restrictions were not placed on the number of times a
313 sample could be played.

314 315 2.3 Data collection

316
317 Data were collected using two different platforms. The first was a jury-based game, and
318 the second was a Qualtrics study with no gamification.

319

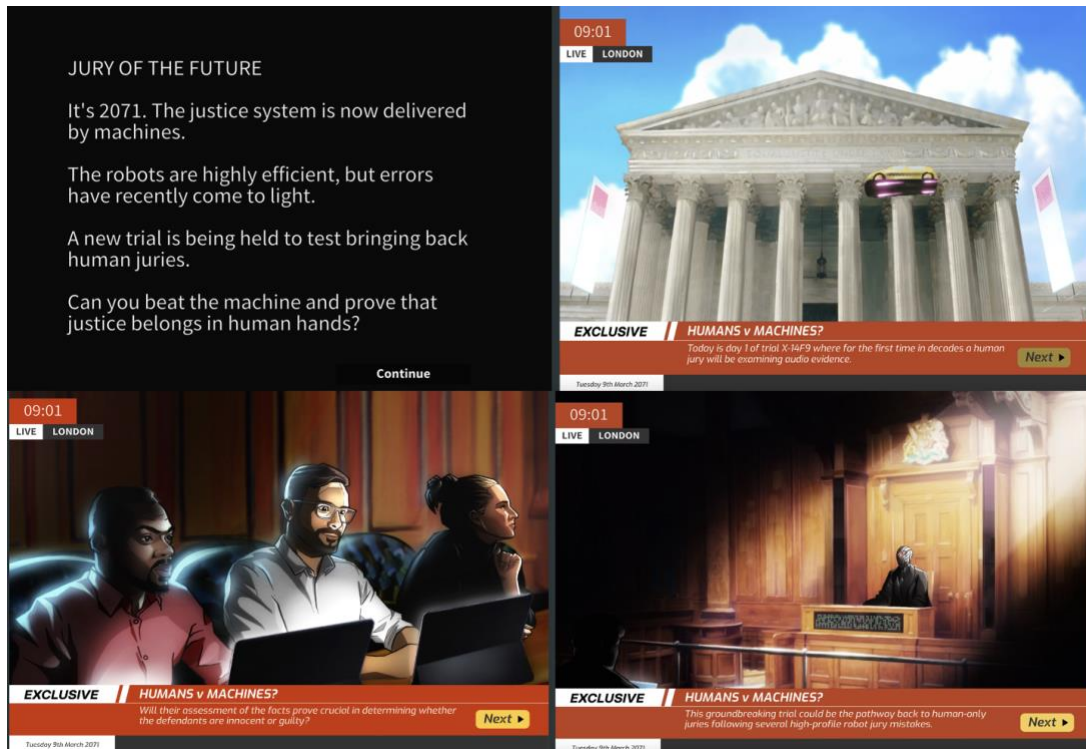
320 2.3.1 Jury-based game platform

321
322 Gamification was used to immerse the participants in a ‘jury of the future’ environment.
323 The aim here was not to replicate a real-world jury scenario. Rather, we wanted to build
324 a context-driven narrative where it was possible for participants to respond to multiple
325 pairs of voices (rather than a real jury scenario involving a single one-to-one
326 comparison of voices). We also intended to increase the stakes of participant responses
327 beyond what might typically be expected in a standard online survey. Providing the jury
328 narrative was also used to help ensure that participants took notice of the expert
329 conclusions and additional evidence presented to them.

330
331 The levels of the game were aligned with the three eight-comparison blocks of
332 comparisons described in Section 2.2. In the first level, participants were presented with
333 a standard, beige Qualtrics-style interface to provide responses to the first eight
334 comparisons. The first level was presented to participants as a ‘Level 1: Tutorial’. At the
335 end of the level, participants were told how many pairs of the eight they got ‘correct’
336 (defined as providing a *sameness* scale rating greater than 50 for same-speaker pairs
337 and below 50 for different-speaker pairs).

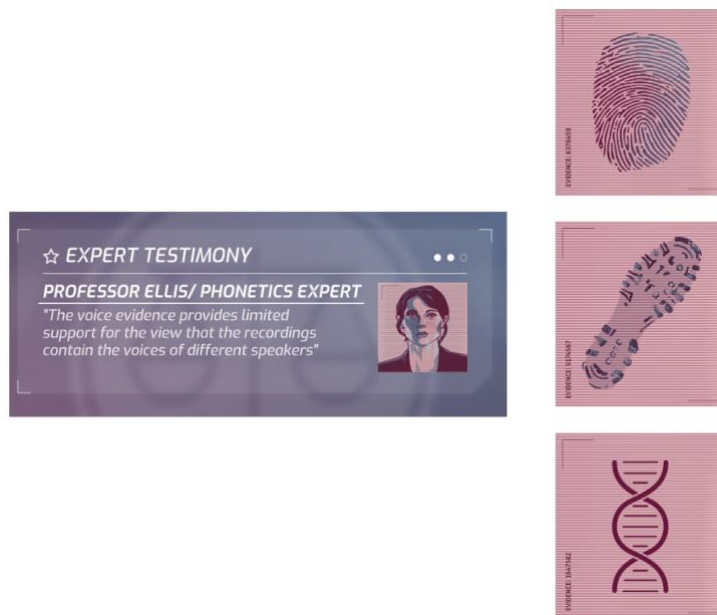
338
339 Between levels one and two, participants were introduced to the game with an
340 illustrated video mimicking a news report and court scenes describing their role on a
341 jury of the future in which they were responsible for the fate of human jurors over
342 machines (see Figure 1). Participants were then presented with a more futuristic,
343 ‘technical’ looking version of the interface from level one and responded to the next
344 block of eight comparisons. The game established that the interface was intended to
345 represent a tablet device for the jury to provide their responses. The format of the
346 questions remained the same as in level one, and participants were also given feedback
347 about their accuracy after each eight-comparison level.

348



349
 350 Figure 1: Screenshots of the illustrated scenes within the game-based data elicitation
 351 tool, presented to participants between levels one and two
 352

353 In the third level, participants were split randomly across two conditions: (i) 'expert
 354 evidence' or (ii) 'additional evidence'. In the 'expert evidence' condition, participants
 355 were presented with an LR-based conclusion from an FVC expert (see Figure 2, left).
 356 The conclusion took the form of either a verbal statement or a numerical value. In both
 357 cases, there were three equally spaced categories (as in Martire et al. 2014), with the
 358 verbal statement and numerical values aligned following the Champod and Evett (2000)
 359 scale. The three categories were 'limited' (LR=10x), 'moderately strong' (LR=1,000x),
 360 and 'very strong' (LR=100,000x) evidence in support of either the same-speaker view or
 361 the different-speaker view. Throughout this level, the conclusion provided by the expert
 362 was always consistent with the ground truth, but the magnitude of the evidential support
 363 was randomised across comparisons. The decision to align conclusions with the ground
 364 truth was intended to maximise the extent to which the participants utilised the expert
 365 conclusion, as should be expected in a real jury trial. We considered that, for example,
 366 providing a strong conclusion in support of the same-speaker view for pairs of voices
 367 that were very clearly different would potentially confuse participants and make them
 368 more likely to entirely ignore the expert conclusions in the rest of the level.
 369



370
 371 Figure 2: Screenshots of the additional information provided to participants in level three
 372 in the form of ‘expert evidence’ (left) and ‘additional evidence’ in the case (right)
 373

374 In the ‘additional evidence’ condition, we primed participants with the suggestion of
 375 additional forensic evidence in the case. Participants were told that there was forensic
 376 evidence in the form of DNA, fingerprints, or footprints and were provided with an image
 377 to represent each evidence type (see Figure 2, right). They were provided with no
 378 indication about the strength of evidence or whether it supported the same-speaker
 379 (guilty) or different-speaker (innocent) proposition. Given this, we expected the effects
 380 to be relatively subtle, but that, in principle, the priming could shift listener judgments
 381 more towards same-speaker rather than different-speakers; i.e. in the absence of other
 382 information the assumption would be that the additional evidence was incriminating.
 383

384 Comparison pairs were counterbalanced across all three levels of the game, meaning
 385 that we have listener responses for all 120 comparisons across all levels. The median
 386 number of listener responses in levels 1 and 2 was 100. The split in level three meant
 387 that fewer listener responses were available for analysis in each condition. In the ‘expert
 388 evidence’ condition, the median number of listener responses per comparison was 20
 389 for the verbal condition and 7 for the numerical condition. For the ‘additional evidence’
 390 level, the median number of responses per comparison was also 7.
 391

392 2.3.2 Qualtrics platform

393
 394 A separate set of listeners provided responses to the same comparisons via Qualtrics.
 395 Again, listeners judged 24 pairs of comparisons each, arranged into three blocks of
 396 eight, with the same distribution of accent comparisons, same- and different-speakers,
 397 and the same process of randomisation within and between blocks as in the game. We
 398 also added a ten-second break between blocks of eight comparisons, to mimic the time
 399 course of the game. The median number of responses to each of the 120 comparisons
 400 within the Qualtrics data was 20. The Qualtrics data were intended to provide a direct

401 point of comparison for assessing the effects of the gamified data elicitation method.
402 Our analysis in this paper is conducted by comparing data from the two versions of level
403 three of the game (i.e. 'expert evidence' and 'additional evidence') with the third block of
404 comparisons from Qualtrics (i.e. those from the same point in the chronology of the
405 experiment).

406 407 2.4 Listeners

408
409 Participants were recruited via Prolific and sampling was conducted in the same way for
410 both the game and Qualtrics versions of the experiment. After excluding three
411 participants for technical issues and two participants for taking more than 71 minutes to
412 complete the experiment, we report here data from a total of 1,804 participants: 1,503
413 who completed the game version and 301 who completed the Qualtrics version. Within
414 the game, 1,205 participants completed the 'expert evidence' level (302 with numeric
415 advice, 903 with verbal advice) and 298 participants completed the 'additional evidence'
416 level. Participants were aged between 18 and 79 years (median=36 years) and based in
417 the United Kingdom, with a gender split of 879 women, 902 men, and 23 non-binary or
418 other. No additional controls over the regional background of participants were
419 implemented. The game version of the experiment took participants an average of 14.5
420 minutes to complete (median=12.9 minutes, range=3.9 to 69.5 minutes), counting only
421 time spent considering audio files (not counting pauses taken between pairs or
422 animations presented between levels). The Qualtrics version of the experiment,
423 inclusive of pauses taken between pairs and levels, took participants an average of 18.8
424 minutes to complete (median=17.3 minutes, range=8.8 to 70.8 minutes). Analysis of the
425 data showed that the time taken to complete the task did not affect listener accuracy.

426 427 2.5 Analysis

428
429 Responses from the 'expert evidence' and 'additional evidence' levels were separately
430 compared with the final block of eight comparisons from the Qualtrics data set. This
431 provides a point of direct comparison involving the same pairs of voice samples, but
432 avoids the order effects observed across the experiment as a whole. Analysis was
433 conducted using the distribution of raw 0 to 100 responses to the question about the
434 *sameness* of the two samples in each comparison. No normalisation of responses was
435 applied because each listener only conducted eight comparisons within the blocks that
436 we considered in this study. Statistical analysis was conducted using linear mixed
437 effects models, with the response to the sameness question as the dependent variable.²
438 In all models, the comparison type (same- or different-speaker) was included as an
439 independent variable, and listener and specific comparison pair were included as
440 random intercepts. This accounts for the fact that there will be random variation in
441 responses across listeners and across specific pairs of speakers. For the analysis of the

² All statistical models took the form: $\text{SamenessResponse} \sim \text{ExperimentalCondition} * \text{GroundTruth} + (1|\text{ComparisonPair}) + (1|\text{Listener})$. In the 'expert evidence' level, ExperimentalCondition was the levels of the expert conclusion including Qualtrics baseline data, with separate models fitted for verbal and numerical conclusions. In the 'additional evidence' level, ExperimentalCondition was initially responses to all of the additional evidence types pooled together and the Qualtrics baseline data. A separate model was fitted where ExperimentalCondition also broke down results by evidence type.

442 'expert evidence' condition, separate mixed effects models were fitted for the verbal and
443 numerical conclusions, and the magnitude of the evidence was included as an
444 independent variable and analysed in interaction with the comparison type (same- or
445 different-speaker). In the 'additional evidence' condition, we report here two models:
446 one in which all evidence types are amalgamated, with a binary Qualtrics vs. game
447 factor analysed in interaction with the comparison type, and another in which the
448 evidence type was included as an independent variable and analysed in interaction with
449 the comparison type. In all cases, the Qualtrics level was dummy-coded as the baseline
450 against which other levels were analysed; same- and different-speaker levels were
451 relevelled as necessary to provide a baseline for significance analysis. Comparison type
452 was included as an interaction in all models to account for the fact that the effects of the
453 'expert evidence' and 'additional evidence' may be different depending on whether the
454 comparison involved the same or different speakers (note, however, that listeners did
455 not have explicit access to the ground truth). In addition to this statistical analysis,
456 results were also considered in terms of speaker accent.

457

458 3. Results

459

460 3.1 'Expert evidence' level

461

462 In this section, we report the results from the 'expert evidence' level, presenting the data
463 from the verbal LR condition first, followed by the numerical LR condition.

464

465 3.1.1 Verbal LR conclusions

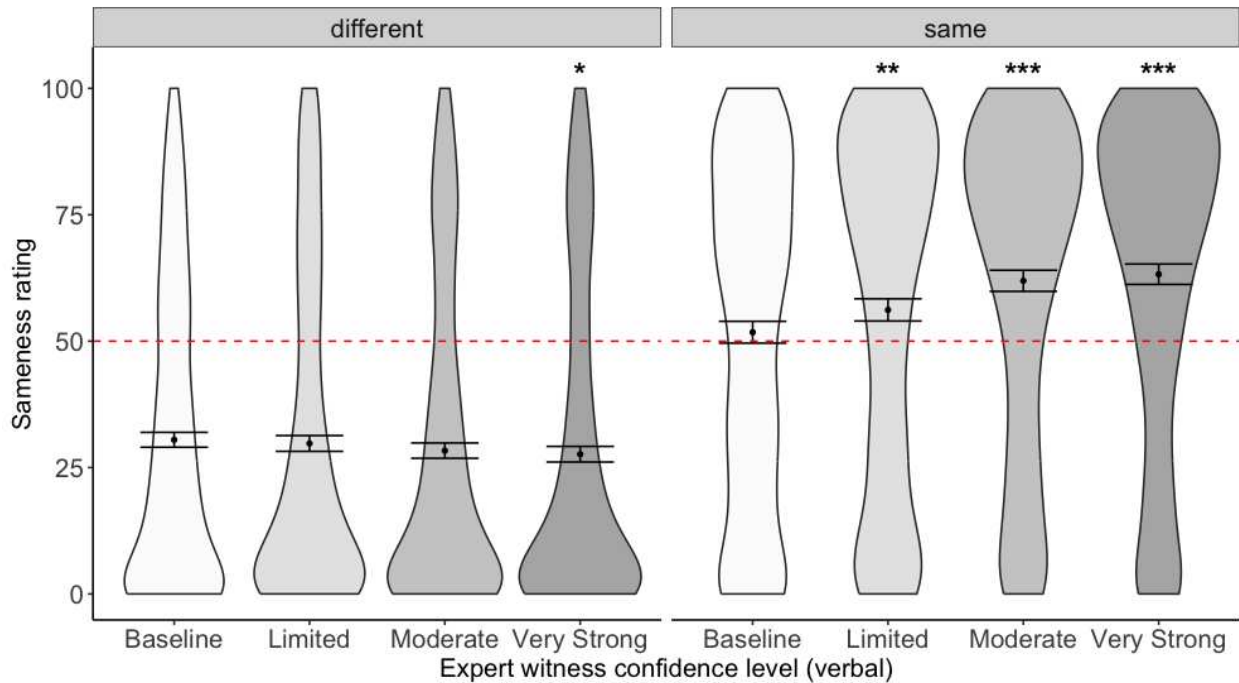
466

467 Figure 3 shows the distributions of responses to the sameness question on the 0 to 100
468 scale, as a function of the magnitude of the verbal LR conclusion presented to listeners.
469 Qualtrics responses from the third block of comparisons (i.e. comparisons 17 to 24) are
470 provided as a baseline where listeners were given no expert conclusion. Responses are
471 distributed in an expected way, with listeners judging different-speaker pairs to be less
472 likely to be the same and same-speaker pairs to be more likely to be the same when
473 provided with an expert conclusion. The shift in listener sameness responses also
474 aligns closely with the magnitude of the evidence presented, such that the smallest
475 effects were found when the expert conclusion was 'limited', and the largest effects
476 were found when the expert conclusion was 'very strong'. Smaller effects were found for
477 different-speaker comparisons than for same-speaker pairs. For different-speaker pairs,
478 statistical modeling demonstrated only a statistically significant difference between the
479 Qualtrics baseline and the 'very strong' conclusion conditions ($\beta=-2.589$, $p=0.048$); the
480 differences between Qualtrics and the 'limited' ($\beta=-1.085$, $p=0.405$) and 'moderately
481 strong' ($\beta=-2.274$, $p=0.082$) conditions did not reach a $p<0.05$ significance level. By
482 contrast, for same-speaker comparisons, all three verbal LR magnitudes produced
483 statistically significant responses from the Qualtrics baseline ('limited': $\beta=4.228$,
484 $p=0.005$, 'moderately strong': $\beta=10.236$, $p<0.001$, 'very strong': $\beta=11.545$, $p<0.001$).

485

486 Interestingly, slightly different patterns were found when analysing responses as a
487 function of the regional accents involved in the comparisons. For the Newcastle and

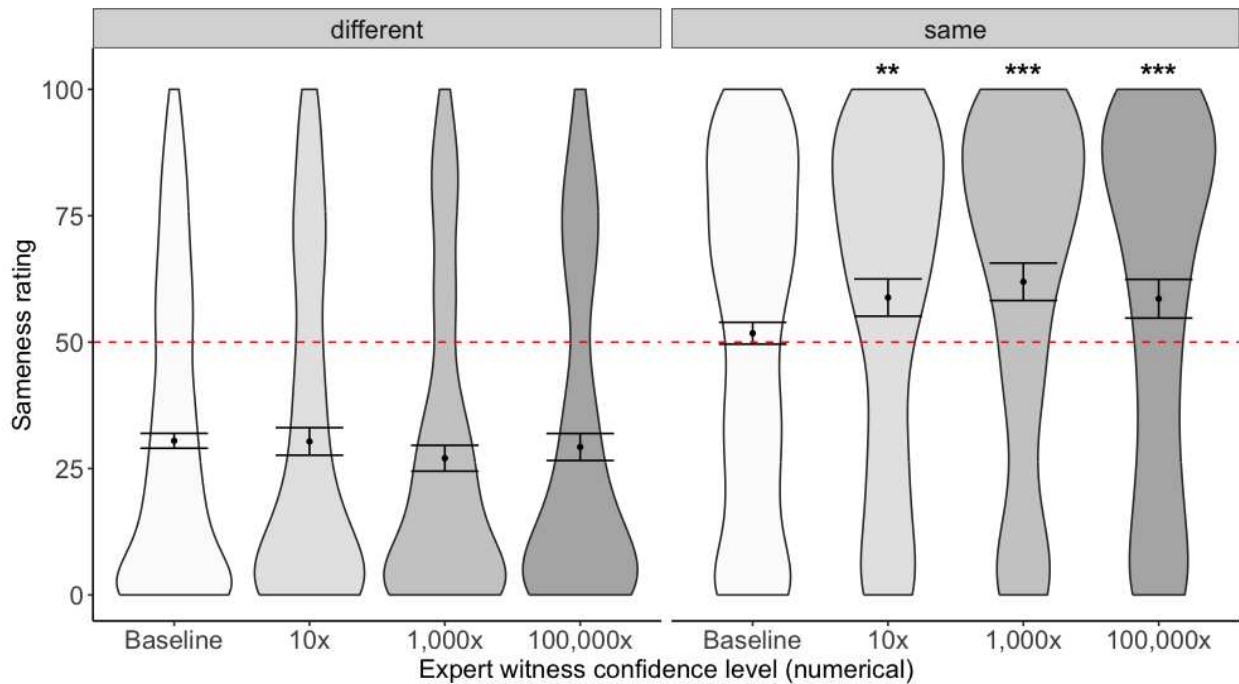
488 SSBE voices, no increases in mean sameness responses were found for same-speaker
 489 pairs between the 'moderately strong' and 'very strong' conditions. For Middlesbrough
 490 voices, the difference between the 'moderately strong' and 'very strong' conditions was
 491 statistically significant. In this way, the responses to the Middlesbrough voices appear to
 492 drive the increase in overall sameness ratings for same-speaker comparisons in Figure
 493 3. The absolute mean values for responses to the Middlesbrough voices were also
 494 lower than the responses to the Newcastle or SSBE voices across all same-speaker
 495 comparisons and all magnitudes of verbal LR conclusion.
 496
 497



498
 499 Figure 3: Distributions of sameness responses to same- (right) and different-speaker
 500 (left) pairs according to the magnitude of verbal LR conclusions (limited, moderately
 501 strong, very strong) and with responses from Qualtrics as a baseline. Whiskers
 502 represent 95% confidence intervals around the mean (significant differences from the
 503 baseline are marked: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$)
 504

505 3.1.2 Numerical LR conclusions

506
 507 Less consistent patterns were found in the responses to numerical LR conclusions.
 508 Figure 4 displays the distributions of sameness responses as a function of the
 509 magnitude of the numerical LR conclusion provided to listeners, with the Qualtrics
 510 responses given as a baseline.
 511



513

514 Figure 4: Distributions of sameness responses to same (right) and different-speaker
 515 (left) pairs according to the magnitude of numerical LR conclusions (10x, 1,000x,
 516 100,000x) and with responses from Qualtrics as a baseline. Whiskers represent 95%
 517 confidence intervals around the mean (significant differences from the baseline are
 518 marked: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$)
 519

520

521 For different-speaker comparisons, the largest shift from the Qualtrics results was found
 522 for the 1,000x condition ($\beta = -3.378$, $p = 0.054$); that is, the numerical equivalent of
 523 'moderately strong' conclusions, rather than the 100,000x 'very strong' conclusions
 524 ($\beta = -0.904$, $p = 0.606$). The 1,000x condition reached close to but just shy of statistical
 525 significance, and no other significant differences were found between the Qualtrics
 526 responses and the numerical LR conditions. For same-speaker comparisons, the
 527 distributions of sameness responses in the 10x ($\beta = 6.202$, $p = 0.003$), 1,000x ($\beta = 9.383$,
 528 $p < 0.001$), and 100,000x ($\beta = 7.921$, $p < 0.001$) numerical LR conditions were significantly
 529 higher than those from Qualtrics. However, again, the biggest effect was found for the
 530 1,000x 'moderately strong' condition, rather than the 100,000x 'very strong' condition.
 531 Unlike the verbal conclusions, no marked differences were found when analysing the
 532 results according to the regional accents involved in the comparisons.

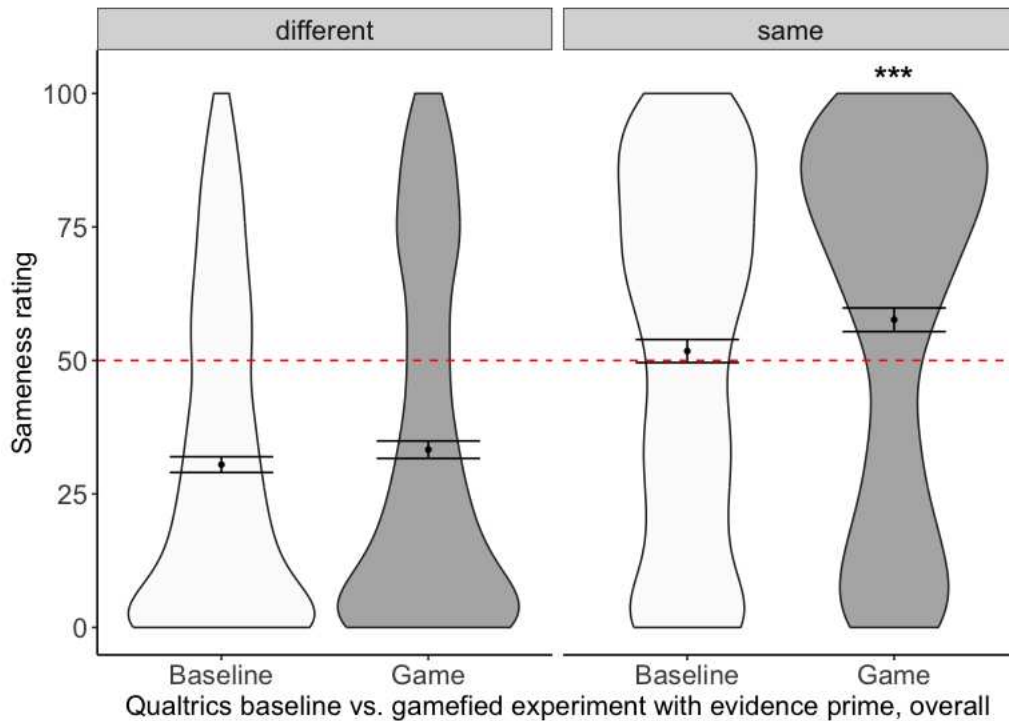
533

534 3.2 'Additional evidence' level

535

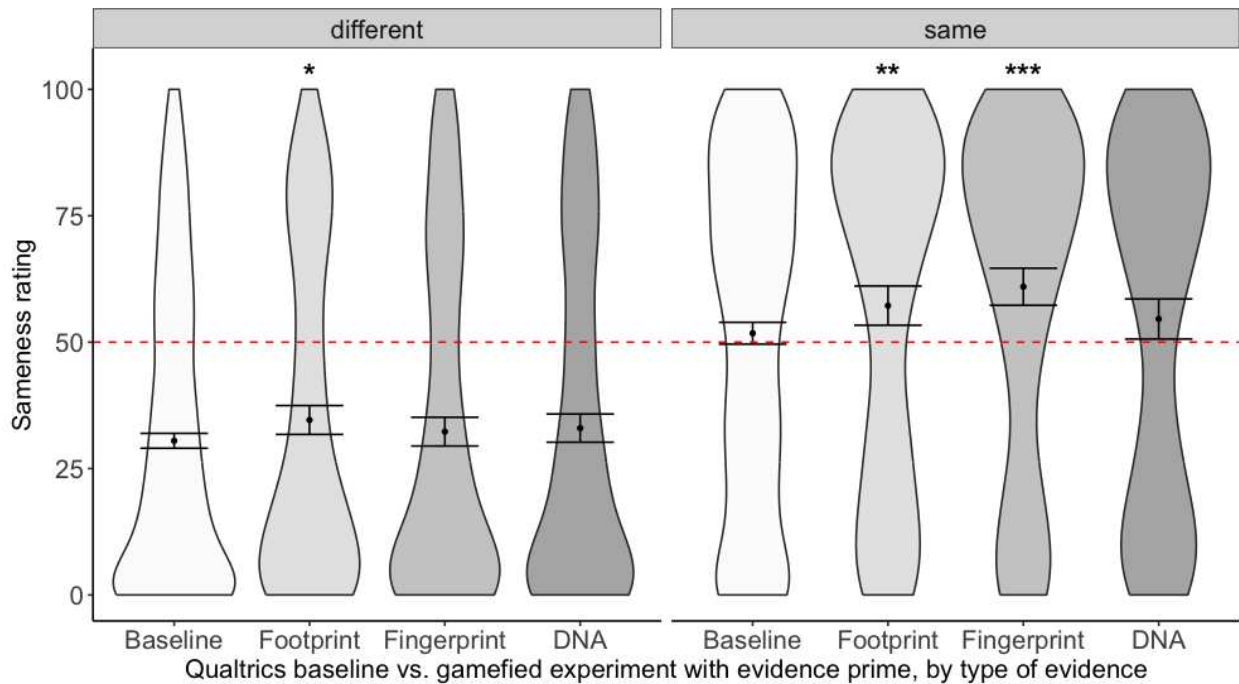
536 Figure 5 displays the overall distribution of sameness responses on the 0 to 100 scale
 537 from the 'additional evidence' level compared with the Qualtrics baseline. Sameness
 538 responses were, on average, higher in the game compared with Qualtrics, meaning that
 539 in both same- and different-speaker conditions, listeners were more likely to think a pair
 540 of voices were the same when primed with additional forensic evidence. This effect was
 541 smaller for different-speaker pairs, where the difference was not statistically significant

541 ($\beta=2.77, p=0.065$), than for same-speaker pairs, where the difference was statistically
 542 significant ($\beta=-5.84, p<0.001$).
 543



544
 545 Figure 5: Distributions of sameness responses to same (right) and different-speaker
 546 (left) pairs from Qualtrics responses as a baseline versus the ‘additional evidence’ game
 547 level (significant differences from the baseline are marked: * $p < 0.05$, ** $p < 0.01$,
 548 *** $p < 0.001$)
 549

550 Variation in listener responses was also found according to the type of ‘additional
 551 evidence’ presented (see Figure 6). Systematic effects were found for footprint evidence
 552 which produced significantly higher sameness ratings compared with the Qualtrics
 553 baseline, for both same- ($\beta=6.57, p=0.002$) and different-speaker ($\beta=3.89, p=0.030$)
 554 comparisons. When presented with fingerprint evidence, listeners provided higher
 555 sameness ratings for same-speaker pairs compared with the Qualtrics baseline
 556 ($\beta=7.74, p<0.001$), with an effect stronger than for footprint evidence. However, no
 557 effect was found for fingerprint evidence in different-speaker comparisons ($\beta=1.18,$
 558 $p=0.512$). For DNA evidence, no significant effects were found for either same-
 559 ($\beta=3.131, p=0.139$) or different-speaker ($\beta=3.24, p=0.068$) comparisons; the
 560 distributions of sameness ratings provided by listeners were therefore largely
 561 unchanged irrespective of whether or not listeners were presented with additional DNA
 562 evidence.
 563



565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

Figure 6: Distributions of sameness responses to same (right) and different-speaker (left) pairs from each of the individual 'additional evidence' types (footprint, fingerprint, DNA) with responses from Qualtrics as a baseline (white) (significant differences from the baseline are marked: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$)

4. Discussion

The present study has demonstrated a range of factors which affect the responses provided by lay listeners in a forensic voice comparison task. In contrast with previous work (see e.g. Martire et al. 2013, 2014), we found more consistent and predictable effects of expert conclusions when those conclusions were provided in the form of a verbal statement rather than a numerical value. When given a consistent-with-fact verbal conclusion, listener sameness responses on average moved in a stepwise fashion in line with the magnitude of the conclusion provided (shifting higher on the sameness scale for same-speaker comparisons and lower on the sameness scale for different-speaker comparisons). There was less systematicity in listener responses when presented with numerical LR conclusions, with no statistically significant effects found for different-speaker pairs compared with the Qualtrics baseline. While same-speaker responses to numerical LRs were all significantly higher than the Qualtrics baseline, the most extreme deviation was found for moderately strong (1,000x) evidence rather than very strong (100,000x) evidence. This finding indicates that, for our study, listeners interpreted the verbal LRs more clearly than numerical LRs. This is especially true of the highest magnitude numerical LRs, which listeners struggled to interpret. Neither the verbal nor the numerical LR responses demonstrated patterns consistent with the *weak evidence effect* reported in previous work (e.g. Fernbach et al. 2011).

593 A possible explanation for the differences between our results and previous work is the
594 fact that the expert conclusions in our study interact with stimuli that listeners respond to
595 as part of the voice comparison task. The fact that listeners were able to combine their
596 own evaluations of the voices with the expert evidence may have provided a boost
597 particularly in the verbal, *limited* evidence conditions. This may in turn have alleviated
598 the *weak evidence effect*. In principle, listeners were also able to complete the task
599 without using the expert conclusion at all, if they so wished. While it is not possible to
600 assess the extent to which listeners utilised the expert evidence directly, the
601 distributions of responses suggest that the expert conclusions were not attended to as
602 consistently for the numerical LRs as for the verbal LRs. This may reflect a general
603 difficulty in interpreting the numerical form, something which participants in previous
604 studies could not avoid as they were responding only to the conclusion itself without any
605 additional stimuli.

606
607 Bigger effects of expert conclusions were found for same-speaker comparisons
608 compared to different-speaker comparisons. We believe this reflects a general pattern,
609 demonstrated in our other work (Hughes et al. 2025), that listener performance is better
610 (i.e. easier) with different-speaker pairs than same-speaker pairs. This is because there
611 is more scope for differences between speakers compared with differences within
612 speakers. Thus, the expert conclusions in our study were most useful for listeners when
613 evaluating same-speaker pairs (without access to the ground truth).

614
615 Our results also show that the benefit of ‘very strong’ same-speaker conclusions is only
616 found in responses to the Middlesbrough voices, at least for the verbal LR conclusions.
617 This finding is consistent with the fact that people are generally much less familiar with
618 the Middlesbrough accent, compared with SSBE or Newcastle (see Hughes et al.
619 2025). For SSBE and Newcastle, sameness responses were generally higher than for
620 Middlesbrough and appear to have reached a peak by the ‘moderately strong’ LR
621 condition from which they don’t increase further. For Middlesbrough, there was still
622 scope for higher overall sameness ratings in the ‘very strong’ LR condition, indicating
623 that listeners struggled most with these comparisons.

624
625 The suggestion of additional forensic evidence also affected listener responses in
626 various ways, despite us not providing any information about the direction of the
627 evidence (i.e. whether it supported the same-speaker or different-speaker propositions).
628 Overall, listeners evaluated the pairs of voices as being more likely to be the same
629 when provided with additional evidence, although this effect was only significant for
630 same-speaker pairs. As with the expert conclusion condition, the fact that different-
631 speaker pairs were less affected by extra evidence may be due to listeners performing
632 at ceiling already, since different-speaker pairs are generally easier than same-speaker
633 pairs. In terms of the different types of extra evidence presented, the strongest effects
634 were found for footprint and fingerprint evidence, with no significant effects found for
635 DNA evidence. The explanation for our findings is not entirely clear but may be due to
636 the choice of images used to represent each evidence type. It is possible that the
637 footprint and fingerprint images were more clearly symbolic of criminal activity, leading
638 listeners to shift their sameness responses towards ‘guilt’. The double helix image used

639 to represent DNA evidence may not have had the same associations with criminality for
640 listeners.

641 642 5. Conclusions

643
644 The results presented in this paper have demonstrated some of the factors that can
645 affect lay peoples' judgments of voices in a forensic voice comparison task. Both the
646 format and magnitude of expert conclusions can have substantial impact on listeners'
647 judgments of how likely two voice recordings are to be from the same person. The
648 suggestion of additional forensic evidence is also capable of shifting listener
649 evaluations. Our study involved collecting data via a jury-based game. The aim was to
650 make the experiment immersive and to increase the stakes for listeners. In this way, we
651 wanted to capture, as far as possible within a research context, some of the elements of
652 real jury decision-making. Our results demonstrate that even relatively subtle
653 manipulations of key variables can have substantial effects on peoples' evaluations,
654 to the extent that people may hear voices as being more similar when provided with
655 additional pieces of information. Future work should further examine the interaction
656 between various sources of case-internal and -external information on jury decision
657 making, and how such information and biases accumulate across multiple pieces of
658 evidence in real legal cases.

659 660 Acknowledgements

661
662 This work was supported by a UK Arts and Humanities Research Council Early Career
663 Grant (AH/T012978/1) entitled *Humans and Machines: Novel Methods for Assessing*
664 *Speaker Recognition Performance*. Thanks also to Joe Cutting and Dan Slawson who
665 helped us develop our game and build the infrastructure to run it as an experiment.

666 667 References

- 668
669 Aitken, C. G. G. et al. (2011) Expressing evaluative opinion: a position statement.
670 *Science and Justice* 51(1): 1–2.
- 671
672 Arscott, E., Morgan, R., Meakin, G. and French, J. (2017) Understanding forensic expert
673 evaluative evidence: a study of the perception of verbal expressions of the strength of
674 evidence. *Science and Justice* 57(3): 221–227.
- 675
676 Association of Forensic Science Providers (2009) Standards for the formulation of
677 evaluative forensic science expert opinion. *Science and Justice* 49(3): 161–164
- 678
679 Berger, C. E. H., Buckleton, J., Champod, C., Evett, I. W., Jackson, G. (2011) Evidence
680 evaluation: A response to the Court of Appeal judgment in R v T. *Science and Justice*
681 51: 43–49.
- 682
683 Bornstein, B. H. and Greene, E. (2011) Jury decision making: implications for and from
684 Psychology. *Current Directions in Psychological Science* 20(1): 63–67.

685
686 Broeders, A. P. A. (1999) Some observations on the use of probability scales in forensic
687 identification. *Forensic Linguistics* 6(2): 228–241.
688
689 Carlson, K. A. and Russo, J. E. (2001) Biased interpretation of evidence by mock jurors.
690 *Journal of Experimental Psychology: Applied* 7(2): 91—103.
691
692 Champod, C. and Evett, I. W. (2000) Commentary on A. P. A. Broeders (1999) ‘Some
693 observations on the use of probability scales in forensic identification’, *Forensic*
694 *Linguistics* 6(2): 228–41. *International Journal of Speech, Language and the Law* 7(2):
695 239–243.
696
697 Corretge, R. (2023) Praat Vocal Toolkit. <https://www.praatvocaltoolkit.com/index.html>
698 (Last viewed November 10, 2023).
699
700 Curley, L. J., Murray, J., MacLean, R., Munro, J. Lages, M., Frumkin, L. A., Laybourn, P.
701 and Brown, D. (2022) Verdict spotting: investigating the effects of juror bias, evidence
702 anchors and verdict system in jurors. *Psychiatry, Psychology and Law* 29(3): 323–344.
703
704 Dixon, J. A., Mahoney, B. and Cocks, R. (2002) Accents of guilt? Effects of regional
705 accent, race, and crime type on attributions of guilt. *Journal of Language and Social*
706 *Psychology* 21(2): 162—168.
707
708 Evett, I. W. (1998) Towards a uniform framework for reporting opinions in forensic
709 science casework. *Science and Justice* 38: 198–202.
710
711 Fernbach, P. M., Darlow, A. and Sloman, S. A. (2011) When good evidence goes bad:
712 the weak evidence effect in judgment and decision-making. *Cognition* 119: 459–467.
713
714 French, J. P. and Harrison, P. (2007) Position statement concerning the use of
715 impressionistic likelihood terms in forensic speaker comparison cases. *International*
716 *Journal of Speech, Language and the Law* 14(1): 137–144.
717
718 Fumkin, L. A. and Stone, A. (2020) Not all eyewitnesses are equal: accent status, race
719 and age interact to influence evaluations of testimony. *Journal of Ethnicity in Criminal*
720 *Justice* 18(2): 123—145.
721
722 Hughes, V., Llamas, C. and Kettig, T. (2022) Eliciting and evaluating likelihood ratios for
723 speaker recognition by human listeners under forensically realistic channel-mismatched
724 conditions. *Proceedings of Interspeech*. Incheon, Korea. pp. 5238–52412.
725
726 Hughes, V., Llamas, C. and Kettig, T. (2025) Human and automatic voice recognition
727 with regionally variable speech samples. *Speech Communication* 103253.
728
729 Llamas, C., Watt, D. and French, J. P. (2016-19) *The Use and Utility of Localised*
730 *Speech Forms in Determining Identity: Forensic and Sociophonetic Perspectives*.
731 ESRC-funded project (ES/M010783/1).

732
733 Lynch, M. and R. McNally (2003) 'Science', 'common sense', and DNA evidence: a
734 legal controversy about the public understanding of science. *Public Understanding*
735 *of Science* 12: 83–103.
736
737 Martire, K. A., Kemp, R. I., Watkins, I., Sayle, M. A. and Newell, B. R. (2013) The
738 expression and interpretation of uncertain forensic science evidence: Verbal
739 equivalence, evidence strength, and the weak evidence effect. *Law and Human*
740 *Behaviour* 37(3): 197–207.
741
742 Martire, K. A., Kemp, R. I., Sayle, M. and Newell, B. R. (2014). On the interpretation of
743 likelihood ratios in forensic science evidence: presentation formats and the weak
744 evidence effect. *Forensic Science International* 240: 61–68.
745
746 Morrison, G. S. (2014) Distinguishing between forensic science and forensic
747 pseudoscience: testing of validity and reliability, and approaches to forensic voice
748 comparison. *Science and Justice* 54(3): 245–256.
749
750 Morrison, G. S. (2022) Advancing a paradigm shift in evaluation of forensic evidence:
751 the rise of forensic data science. *Forensic Science International: Synergy* 5: 100270.
752
753 Mullen, C., Spence, D., Moxey, L. and Jamieson, A. (2014) Perception problems of the
754 verbal scale. *Science and Justice* 54(2): 154–158.
755
756 Munro, J., Motson, F., Turner, J., Frumkin, L. A. and Curley L. J. (2024) Double
757 jeopardy: The effects of retrial knowledge on juror decisions. *Journal of Criminal*
758 *Psychology* 14(4): 444–455.
759
760 Nolan, F., McDougall, K., de Jong, G. and Hudson, T. (2009) The DyViS database:
761 style-controlled recordings of 100 homogeneous speakers for forensic phonetic
762 research. *International Journal of Speech, Language and the Law* 16(1): 31–57.
763
764 Paver, A., Wright, D., Braber, N. and Pautz, N. (2025) Stereotyped accent judgements
765 in forensic contexts: listener perceptions of social traits and types of behaviour.
766 *Frontiers in Communication* 9: 1462013.
767
768 Redmayne, M., Roberts, P., Aitken, C. G. G. and Jackson, G. (2011) Forensic science
769 evidence in question. *Criminal Law Review* 5: 347–356.
770
771 Rose, P. (2002) *Forensic Speaker Identification*. Taylor and Francis: New York.
772
773 Rose, P. and Morrison, G. S. (2009) A response to the UK Position Statement on
774 forensic speaker comparison. *International Journal of Speech, Language and the Law*
775 16(1): 139–163.
776

777 Tindale, R. S. and Berryman, K. (2023) Jury decision-making. In *Oxford Research*
778 *Encyclopedia of Psychology*. Oxford: OUP.