

1 Realistic solutions for practising forensic scientists – a response to Morrison (2023)

3 1. Introduction

4 In Kirchhübel et al. (2023) [1], we address the issue of method validation for the Auditory
5 Phonetic and Acoustic (AuPhA) approach to forensic voice comparison (FVC). We illustrate
6 one possible solution, which centres around competency testing of the human expert
7 implementing AuPhA. We received a response to our article by Morrison (2023) [2] which
8 conveys the idea that there is but a single acceptable way of meaningfully demonstrating
9 validity of FVC methods. This position not only contradicts a range of relevant literature
10 (including Morrison et al. (2021) [3] – henceforth the “*Consensus*”), but it also defies common
11 sense. We take this opportunity to respond to Morrison (2023) [2].

13 2. Clarifying the issue of validation

14 It appears that the main issue that Morrison (2023) [2] raises relates to the number of recordings
15 needed to validate FVC methods. This is reflected in the title of Morrison (2023) [2] which is
16 “*A single test pair does not a method validation make: A response to Kirchhübel et al. (2023)*”.
17 Morrison (2023) [2] recognises that it would be challenging in practical terms to validate
18 AuPhA according to the recommendations made in the *Consensus*. However, it is not the case
19 that it would just be challenging; it would in fact be impossible. An illustration may help here.
20 The *Consensus* recommends that validation data should be sufficiently representative of the
21 case at hand and contain a sufficient number of speakers. Importantly, the *Consensus* does not
22 specify a number. This in itself indicates that it would be difficult to agree a specific number to
23 apply to all cases, and this implies that there is a need for flexibility. Despite this, Morrison
24 (2023) [2: 328] states that ‘*It may be that a practically achievable validation set consists of*
25 *pairs of recordings from only upper tens of speakers to lower hundreds of speakers.*’ Let’s
26 assume for the purposes of exposition that the validation set contains 100 same-speaker FVC
27 tests and 100 different-speaker FVC tests. Using the AuPhA approach exemplified in
28 Kirchhübel et al. (2023) [1], it would take at least 10 hours to complete one FVC test. We are
29 therefore looking at 2,000 hours, or in other words, 50 full-time working weeks (net analysis
30 time) to complete all the tests as part of the proposed validation exercise. And then this of
31 course would need to be repeated for a new type of case in which the recordings are judged to
32 be insufficiently comparable to those contained in the validation exercise (e.g., different
33 recording types or speaker demographics involved).

34

35 To overcome these time and resource impracticalities, Morrison (2023) [2] puts forward the
36 solution of simply discarding AuPhA for FVC and only using human-supervised automatic
37 methods. However, we consider this proposal to be detrimental, in part, because not all
38 casework recordings meet the quality requirements in order to apply an automatic system. Even
39 for those cases where the recordings do meet the quality requirements to use an automatic
40 system, there are practical hurdles to overcome. These hurdles include, for example, obtaining
41 case-relevant recordings to form suitable validation datasets, ensuring the competence of the
42 humans supervising the automatic systems, and navigating its reception within a legal
43 environment. Therefore, solely relying on automatic systems would mean that a huge number
44 of FVC cases go unanalysed and forensic speech science would become largely ineffectual as
45 a forensic discipline. We are very much in favour of incorporating automatic systems into FVC,
46 but this is not with a view to displace AuPhA. There is widespread recognition within the
47 forensic speech science community that the two methods – AuPhA analysis and automatic
48 speaker recognition – complement one another. As such, practitioners who regularly use
49 automatic systems in FVC casework use them in conjunction with an AuPhA analysis (e.g.,
50 Jessen (2018) [4]; van der Vloed and Cambier-Langeveld (2023) [5]).

51

52 Rather than discarding AuPhA, a much more productive solution, in our view, is to find a way
53 of validating AuPhA which is both meaningful and practically feasible. The approach contained
54 within Kirchhübel et al. (2023) [1] exhibits both of these characteristics by carefully selecting
55 data which are reflective of casework conditions and by including a comprehensive external
56 review process.

57

58 In any case, Morrison (2023) [2] has misrepresented the proposal contained in Kirchhübel et
59 al. (2023) [1] (particularly in the title of Morrison (2023) [2]) by suggesting that we have
60 proposed a validation protocol that only includes an analysis of a single pair of recordings.
61 When reading Kirchhübel et al. (2023) [1] in full, it becomes clear that what we are proposing
62 is a portfolio approach to validation which involves numerous FVC and speech analysis
63 exercises covering different accents, speaker demographics, recording characteristics, sample
64 durations, etc. The single 1:1 comparison is presented as an example of one component within
65 this approach. It is not intended as the sole contribution. In fact, we explicitly list “blind-
66 grouping” (e.g., Cambier-Langeveld et al. (2014) [6]) as another possible component which
67 would provide opportunity to include a larger number of speech samples into a validation

68 exercise. We propose that this validation approach is an ongoing and cumulative process
69 throughout the careers of practitioners. Morrison (2023) [2] therefore fails to acknowledge our
70 full proposal and appears to have taken a selective approach to his response.

71

72 **3. Further Clarifications**

73 There are further clarifications to make in response to Morrison (2023) [2]. The main ones are
74 set out below.

75

76 3.1 Scope of the *Consensus* (Morrison et al. (2021) [3])

77 Morrison et al. (2021) [3] present a '*Consensus*' that describes recommendations for how
78 method validation can be achieved using an automatic approach to FVC. Morrison (2023) [2]
79 points out that the *Consensus* had 13 authors and an additional 7 supporters. We acknowledge
80 the efforts in producing this work and we do not take issue with its contents (as indicated by
81 its citation in Kirchhübel et al. (2023) [1]). However, it is important to highlight that the
82 *Consensus* has a restricted scope, i.e., it applies to "*validation of forensic-voice comparison*
83 *systems that are based on relevant data, quantitative measurements, and statistical models, and*
84 *that output numeric likelihood ratios*". That is, it largely applies to forensic application of
85 automatic speaker recognition. However, Morrison (2023) [2] appears to suggest that, "*with*
86 *minor wording changes*", the scope of the *Consensus* can be extended to cover other approaches
87 to FVC. The AuPhA method as described in Kirchhübel et al. (2023) [1] does not fall within
88 the original scope because it is not dependent on statistical models, nor does it necessarily result
89 in numeric likelihood ratios. Further, in Section 2 of this response, we have already explained
90 why the *Consensus* cannot be extended to AuPhA, and "*minor wording changes*" would
91 therefore make no difference. In any case, we do not consider it to be appropriate to attempt to
92 extend a consensus, which by its very nature is based on explicit agreement from all
93 contributors, to matters outside its scope.

94

95 3.2 Admissibility of automatic systems

96 Footnote 3 of Morrison (2023) [2] states:

97 "*Kirchhübel et al. (2023) claims that the auditory-acoustic-phonetic approach "is the only*
98 *admissible approach in UK jurisdictions for voice comparison analysis."* The implication
99 *that the human-supervised-automatic approach would not in-principle be admissible is*
100 *incorrect."*

101

102 We offer two responses in relation to this comment. First, Morrison (2023) [2] misrepresents
103 our statement by omitting some words from the beginning of the relevant sentence. It actually
104 reads:

105 *“At the time of writing, AuPhA analysis is the only admissible approach in UK jurisdictions*
106 *for voice comparison analysis.”*

107

108 Second, we certainly did not intend for this to imply a general inadmissibility of automatic
109 methods for FVC in UK jurisdictions. The key legal authority on this point is R v Slade
110 [2015] EWCA Crim 71 where the court was presented with voice comparison analysis by
111 way of an automatic speaker recognition system. The court voiced a number of concerns
112 about this evidence and ultimately declined to admit it in that case, although the court did not
113 make a definitive ruling as to whether automatic speaker recognition evidence can ever be
114 admissible. On reflection, we accept that we could have presented a clearer account in our
115 original article. We are keen to reiterate that we are very much in favour of incorporating
116 automatic methods in FVC.

117

118 3.3 Expression of conclusions

119 Footnote 4 of Morrison (2023) [2] questions whether we have actually followed the logic of
120 the likelihood ratio framework:

121 *“Kirchhübel et al. [22] states that the practitioner and the reviewer “expressed their*
122 *conclusions with reference to the scale that is recommended by the UK [sic] Association of*
123 *Forensic Science Providers [25] and ENFSI [26] (however, their conclusions were not*
124 *derived from a numerical likelihood ratio).” Both those scales, however, are intended to*
125 *provide verbal expression corresponding to numerical ranges of likelihood ratios, and the*
126 *ENFSI Guideline states that even if the numerator and denominator of a likelihood ratio are*
127 *“informed by subjective probabilities using expert knowledge. These probability assignments*
128 *shall still be expressed by a number between 0 and 1 rather than by an undefined qualifier*
129 *(such as frequent, rare, etc.).” There is no evidence in Kirchhübel et al. [22] that the*
130 *practitioner or reviewer actually followed the logic of the likelihood-ratio framework.”*

131

132 For the avoidance of any doubt, the practitioner and the reviewer in Kirchhübel et al. (2023)
133 [1] did follow the logic of the likelihood ratio framework (which, of course, is not dependent
134 on computing or estimating numerical probabilities). Instead of computing probabilities with
135 reference to databases, probabilities were derived from knowledge and judgment, with

136 appropriate caution. Instead of assigning numerical values to these subjective probabilities,
137 the practitioner and reviewer made use of verbal expressions. This is because the practitioner
138 and reviewer also followed legal and regulatory authority that is directly relevant to UK
139 jurisdictions when it comes to the expression of expert opinions which are the product of
140 knowledge and experience, and not based on explicit databases. From R v Atkins and Atkins
141 [2009] EWCA Crim 1876, it is clear that experts should not be using numerical estimates in
142 these circumstances. This is because of the danger of misleading the jury to believe that
143 quantitative processes have been applied when in fact they have not.

144

145 In addition, in guidance published in 2021, the UK Forensic Science Regulator (FSR) makes
146 clear that numerical values should only be expressed if they are based on appropriate
147 statistical data. In the absence of appropriate statistical data, the results should still be
148 expressed probabilistically, but without numerical values. At paragraphs 7.2.7 and 8.5.13 in
149 the FSR 2021 guidance [7] it states:

150 *“In an ideal situation, the scientist would have access to relevant, high quality*
151 *datasets to inform their evaluation. However, the nature of forensic science, where*
152 *every case has different circumstances, means that large and directly relevant*
153 *datasets are very rarely available. A LR [likelihood ratio] determined by the expert*
154 *using experience alone can be no more wrong than the expert’s opinion on which it is*
155 *based and that opinion is admissible evidence. However, it is less easily tested than a*
156 *LR based on relevant data, quantitative measurements, and statistical models’ (7.2.7)*
157 *‘Where no (published or unpublished) structured data are available, a qualitative*
158 *evaluation shall be reported, based on the expert’s qualitative evaluation of the*
159 *probability of the observations under each proposition...” (8.5.13)*

160

161 **4. Conclusion**

162 Morrison (2023) [2] paints a picture of forensic speech science being a field that has had years
163 of opportunities to “progress” by way of analysis methods, and that these opportunities have
164 been blocked by reluctant FVC practitioners who “do not think like forensic scientists”.
165 However, Morrison (2023) [2] fails to accommodate the realities of FVC casework. By
166 constraining ourselves to a human-supervised automatic method to FVC (simply because such
167 a method could allow for hundreds of validation tests) would mean that very little FVC
168 casework could be carried out. It is of little help to let a doctrinaire view of validation stifle
169 other FVC methods, the scientific validity of which can be demonstrated through more flexible

170 means. In saying this, we are thinking like practising forensic scientists with experience of
171 working within the criminal justice system. The point of Kirchhübel et al. (2023) [1] was to
172 find real solutions to real problems while advocating for quality in FVC.

173

174 **References**

175 [1] Kirchhübel, C., Brown, G., & Foulkes, P. (2023). What does method validation look like
176 for forensic voice comparison by a human expert? *Science & Justice*, 63, 251-257. DOI:
177 10.1016/j.scijus.2023.01.004.

178 [2] Morrison, G. S. (2023). A single test pair does not a method validation make: A response
179 to Kirchhübel et al. (2023). *Science & Justice*, 63, 327-329. DOI:
180 10.1016/j.scijus.2023.03.001.

181 [3] Morrison, G. S., Enzinger, E., Hughes, V., Jessen, M., Meuwly, D., Neumann, C.,
182 Planting, S., Thompson, W. C., van der Vloed, D., Ypma, R. J. F., Zhang, C., Anonymous, A.,
183 & Anonymous, B. (2021). Consensus on validation of forensic voice comparison. *Science &*
184 *Justice*, 61(3), 299-309. DOI: 10.1016/j.scijus.2021.02.002.

185 [4] Jessen, M. (2018). Forensic voice comparison. In *Handbook of Communication in the*
186 *Legal Sphere*, edited by Jacqueline Visconti, Berlin, Boston: De Gruyter Mouton, pp. 219-
187 255. DOI: 10.1515/9781614514664-012.

188 [5] van der Vloed, D., and Cambier-Langeveld, T. (2023). How we use automatic speaker
189 comparison in forensic practice. *International Journal of Speech, Language and the Law*,
190 29(2), 201–224. DOI: 10.1558/ijssl.23955.

191 [6] Cambier-Langeveld, T., van Rossum, M. and Vermeulen, J. (2014). Whose voice is that?
192 Challenges in forensic phonetics. In: J. Caspers, Y. Chen, W. Heeren, J. Pacilly, N.O. Schiller
193 & E. van Zanten (eds), *Above and Beyond the Segments. Experimental linguistics and*
194 *phonetics*. Amsterdam: John Benjamins Publishing Company, 14-27.

195 [7] Forensic Science Regulator: Development of Evaluative Opinions. FSR-C-118 (Issue 1),
196 2021. URL: [https://www.gov.uk/government/publications/development-of-evaluative-](https://www.gov.uk/government/publications/development-of-evaluative-opinions)
197 [opinions](https://www.gov.uk/government/publications/development-of-evaluative-opinions).