# A single test pair does not a method validation make:

# A response to Kirchhübel et al. (2023)

## Author and affiliations:

Geoffrey Stewart Morrison *

Forensic Data Science Laboratory, Aston University, Birmingham, UK

Forensic Evaluation Ltd, Birmingham, UK

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

## ORCID:

Geoffrey Stewart Morrison        0000-0001-8608-8207

## Disclaimer:

All opinions expressed in the present response 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:**

**Acknowledgements:**

**Abstract:**

In terms of development of methods that are transparent and reproducible, that are intrinsically resistant to cognitive bias, that use the logically correct framework for interpretation of evidence (the likelihood-ratio framework), and that are empirically validated under casework conditions, forensic voice comparison may be one of the most progressive branches of forensic science. Forensic voice comparison is not, however, a monolithic branch of forensic science. The aforementioned progress has been made in the context of human-supervised-automatic methods. Unfortunately, many or most forensic-voice-comparison practitioners, who use auditory-acoustic-phonetic methods, have not joined in this progress. There are calls going back to the late 1960s for forensic-voice-comparison methods to be meaningfully validated under casework conditions, but many or most forensic-voice-comparison practitioners have still not heeded those calls. A recent example appears in Kirchhübel et al. (2023)

https://doi.org/10.1016/j.scijus.2023.01.004, which proposes a validation protocol in which the number of test pairs is only one. This does not constitute meaningful validation. If the cost and time necessary to conduct meaningful validations of auditory-acoustic-phonetic methods are such that it is not practical to conduct meaningful validations, then forensic voice comparison should not be performed using auditory-acoustic-phonetic methods.

**Keywords:**

Forensic voice comparison; Validation

**Letter to Editor:**

1   A keynote presentation at the European Academy of Forensic Science Conference 2022,
2   published as Morrison [1], described a paradigm shift in evaluation of forensic
3   evidence. According to Morrison [1], in the new paradigm, the forensic-data-science
4   paradigm, the methods used adhere to the following principles: they are transparent
5   and reproducible; they are intrinsically resistant to cognitive bias; they use the logically
6   correct framework for interpretation of evidence (the likelihood-ratio framework); and
7   they are empirically validated under casework conditions. With respect to the last
8   principle, Morrison [1] noted that "practitioners in multiple branches of forensic
9   science often claim that training and experience provide sufficient warrant for their
10  conclusions (see Mnookin et al. [2]; Risinger [3]; President's Council of Advisors on
11  Science and Technology [4]; Morrison & Thompson [5]), and deny or obfuscate about
12  the need for validation (see Cole [6]; Morrison [7]; President's Council of Advisors on
13  Science and Technology [4]; Koehler [8]; Morrison et al. [9]), or propose lax validation
14  protocols that do not require demonstration of performance under casework conditions
15  (see Morrison et al. [10], [11])." Morrison [1] also noted, however, that "protocols for

16   validating systems that output likelihood ratios have been developed, including metrics

17   and graphics appropriate for representing the results of such validations (Meuwly [12];

18   Brümmer & du Preez [13]; Morrison [14]; Meuwly et al. [15]; Ramos et al. [16];

19   Morrison et al. [17]). Much of the latter work has been conducted in the context of

20   forensic voice comparison, but the results are applicable across forensic science in

21   general."

22   Between 2016 and 2019, in the context of a virtual special issue of *Speech*

23   *Communication*,[1] training and test data reflecting the conditions of a real forensic-

24   voice-comparison case were released, and, using those data, multiple laboratories

25   validated multiple different human-supervised-automatic forensic-voice-comparison

26   systems. The test data consisted of 111 same-speaker pairs of recordings (from 61

27   unique speakers), and 9720 different-speaker pairs of recordings (from 3660 unique

28   pairs of speakers) (Morrison & Enzinger [18]). Detailed results were published in a

29   series of papers in the virtual special issue, and a summary of the results was published

30   in the virtual special issue's conclusion article (Morrison & Enzinger [19]). The data

31   and answer key have since been made available for others to use. Weber et al. [20] used

32   the data to conduct a benchmark validation of another human-supervised-automatic

33   forensic-voice-comparison system, and Basu et al. [21] used a subset of the data to

34   assess lay listeners' speaker-identification abilities and compare them with the

35   performance of the latter human-supervised-automatic forensic-voice-comparison

36   system. In 2019 and 2020, a group of researchers and practitioners collaborated on

37   developing a *Consensus on validation of forensic voice comparison*. The published

38   version, Morrison et al. [17], had 13 authors and an additional 7 supporters. In order to

39   be able to write a progressive document, the scope of the *Consensus* was restricted to

40   "validation of forensic-voice-comparison systems that are based on relevant data,

41   quantitative measurements, and statistical models, and that output numeric likelihood

42   ratios"; however, with minor wording changes the *Consensus* would be applicable to

---

[1]   https://www.sciencedirect.com/journal/speech-communication/special-issue/10KTJHC7HNM

43  validating methods for assigning likelihood ratios in other approaches to forensic voice

44  comparison or to validating methods for assigning likelihood ratios which address

45  source-level hypotheses in other branches of forensic science. Since validation is a

46  black-box exercise, the details of the method being validated are not relevant.

47  The foregoing may seem to suggest that forensic voice comparison is one of the most

48  progressive branches of forensic science, but forensic voice comparison is not a

49  monolithic branch of forensic science. Almost a decade ago, Morrison [7] wrote about

50  and reiterated calls going back to the late 1960s for forensic-voice-comparison methods

51  to be meaningfully validated under casework conditions. Unfortunately many or most

52  forensic-voice-comparison practitioners have still not heeded those calls. A recent

53  example appears in Kirchhübel et al. [22], which proposes a validation protocol in

54  which the number of test pairs is only one.[2]

55  As previously mentioned, method validation is a black-box exercise. It may be possible

56  to perform the same task using different methods, e.g., it is possible to perform forensic

57  voice comparison using different methods that fall under the human-supervised-

58  automatic approach (the aforementioned virtual special issue of *Speech*

59  *Communication* provides examples) or different methods that fall under the auditory-

60  acoustic-phonetic approach.[3]  What validation does is test how well a particular method

61  performs the task, e.g., the task of assessing the likelihood of obtaining the speech of

---

[2]  Kirchhübel et al. [22] frames the task that the practitioner performed as a "proficiency test", but claims that this serves the purpose of method validation – the title of the article is: *What does method validation look like for forensic voice comparison by a human expert?*

[3]  See Morrison & Zhang [23] for a recent description of different approaches to forensic voice comparison, including human-supervised-automatic and auditory-acoustic-phonetic approaches. Kirchhübel et al. [22] claims that the auditory-acoustic-phonetic approach "is the only admissible approach in UK jurisdictions for voice comparison analysis." The implication that the human-supervised-automatic approach would not in-principle be admissible is incorrect. For a review of relevant case law in England & Wales and in Northern Ireland, and a discussion of admissibility of forensic voice comparison in light of England & Wales Criminal Practice Directions 19A ([2015] EWCA Crim 1567 V 19A), see Morrison [24].

62   interest on the questioned- and known-speaker recordings if they were both produced

63   by the same speaker versus the likelihood of obtaining the speech of interest on the

64   questioned- and known-speaker recordings if they were produced by two different

65   speakers from the relevant population.[4] Black-box testing is not concerned with how

66   a method performs the task, only with how well the method performs the task. It

67   therefore does not matter whether a method for performing forensic voice comparison

68   is a human-supervised-automatic method or an auditory-acoustic-phonetic method;

69   different methods for performing the same task can, and should, be validated using the

70   same validation protocol. The statement in Kirchhübel et al. [22] that "it would not be

71   possible to simply adopt the recommendations made in [the *Consensus*] for the

72   [auditory-acoustic-phonetic] approach" is therefore incorrect unless there is a practical

73   impediment to conducting validation of auditory-acoustic-phonetic methods according

74   to the recommendations of the *Consensus*. The amount of time and amount of human

75   effort that it takes to compare each pair of recordings in the validation set will be much

76   much greater for an auditory-acoustic-phonetic method than for a human-supervised-

77   automatic method. Validating an auditory-acoustic-phonetic method will therefore be

78   practically much more difficult than validating a human-supervised-automatic method.

79   This does not, however, excuse auditory-acoustic-phonetic methods from the

80   requirement applicable to all methods that they be meaningfully validated.

81   As stated in the *Consensus*, a necessary condition for validation to be meaningful is the

---

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

82    following:[5]

83        2.5.2. Validation data [pairs of same-speaker recordings and pairs of different-
84        speaker recordings] should be sufficiently representative of the relevant
85        population for the case, and sufficiently reflective of the conditions of the
86        questioned-speaker and known-speaker recordings in the case, that the results of
87        validating the system using those data will be informative as to the expected
88        performance of the system when it is applied in the case.

89        2.5.3. One of the criteria for the validation data to be sufficient is that the number
90        of speakers included be sufficient. Because of sampling variability, small
91        validation sets can give results that are not representative of the case conditions.

92    The number of speakers is a constraint on how many unique same-speaker and unique
93    different-speaker pairs can be constructed. Factorial combinations allow for many
94    more unique different-speaker pairs to be constructed, but the number of unique same-
95    speaker pairs that can be constructed will be limited to the number of speakers from
96    whom pairs of recordings are available. The *Consensus* does not recommend a specific
97    value for what would constitute a sufficient number of speakers, but instead states that:

98        2.6.1. The decision as to whether the calibration data and the validation data are
99        sufficiently representative of the relevant population for the case and sufficiently
100        reflective of the conditions of the questioned-speaker and known-speaker
101        recordings in the case will be the result of a subjective judgment made by the
102        forensic practitioner.

103        2.6.7. The forensic practitioner should communicate to the court a clear
104        description of the calibration data and the validation data used.

105        2.6.8. A description of the calibration and validation data is a prerequisite for a

---

[5] The following quotations use the original paragraph numbering from Morrison et al. [17].

106    second forensic practitioner to be able to conduct an independent review so as to
107    be able to opine on whether the data are sufficient.

108    2.6.9. A description of the calibration and validation data is also a prerequisite for
109    the court to be able to decide to either accept or reject the first forensic
110    practitioner's decision about the sufficiency of the data.

111  Larger numbers of speakers would be better, but the number of speakers included in
112  the validation set will be constrained by the cost and time required to obtain data that
113  are sufficiently representative of the relevant population and sufficiently reflective of
114  the questioned-speaker and known-speaker recordings' conditions. In may be that a
115  practically achievable validation set consists of pairs of recordings from only upper
116  tens of speakers to lower hundreds of speakers. Whether a validation set of this size
117  would be sufficiently representative of the relevant population and sufficiently
118  reflective of the recording conditions of a case is a matter of judgement, and, ultimately,
119  of acceptance by the court. It is very clear, however, that a validation set consisting of
120  a single test pair, as proposed in Kirchhübel et al. [22], would not be sufficiently
121  representative of the relevant population nor sufficiently reflective of the recording
122  conditions for the validation results to be informative as to the expected performance
123  of the system when it is applied in the case. A validation consisting of a single test pair
124  is not meaningful.[6]

125  Practitioners of forensic voice comparison should think like forensic scientists not like
126  phoneticians. A practitioner thinking like a forensic scientist will use whatever method
127  they believe will best perform the task of forensic voice comparison. Their choice of
128  method should be informed by prior validation studies. The decision as to whether the
129  performance of the method is sufficiently good in the context of the particular case

---

[6] In Kirchhübel et al. [22], the single test pair was selected from a set of nine potential pairs. The validation set presented
to the practitioner in Kirchhübel et al. [22], however, consisted of that single test pair. If the validation set had consisted
of nine test pairs, we would argue that that would also have been too small.

130   must be based on a validation of the method using data that are representative of the
131   relevant population for the case and reflective of the conditions of the questioned-
132   speaker and known-speaker recordings in the case. The validation may have been
133   conducted ahead of time (anticipatory validation) and a judgement made that the
134   conditions of the case are sufficiently similar to the conditions under which the existing
135   validation was conducted; or, if such a validation does not already exist, a new
136   validation should be conducted using data that are judged to be sufficiently similar to
137   the conditions of the case (case-by-case validation). In contrast, a practitioner thinking
138   like a phonetician will persist in only using auditory-phonetic or acoustic-phonetic
139   methods, even when other methods (such as human-supervised-automatic methods)
140   have been demonstrated to result in superior performance, and even when the
141   performance of the auditory-phonetic or acoustic-phonetic methods have not been
142   meaningfully demonstrated at all. If the cost and time necessary to conduct meaningful
143   validations of auditory-acoustic-phonetic methods are such that it is not practical to
144   conduct meaningful validations, then forensic voice comparison should not be
145   performed using auditory-acoustic-phonetic methods.

## References

[1]   Morrison G.S. (2022). Advancing a paradigm shift in evaluation of forensic
      evidence: The rise of forensic data science. *Forensic Science International:
      Synergy*, 5, 100270. https://doi.org/10.1016/j.fsisyn.2022.100270

[2]   Mnookin J.L., Cole S.A., Dror I.E., Fisher B.A.J., Houck M.M., Inman K.,
      Kaye D.H., Koehler J.J., Langenburg G., Risinger D.M., Rudin N., Siegel J.,
      Stoney D.A. (2011). The need for a research culture in the forensic sciences.
      *UCLA Law Review*, 58, 725–777. https://www.uclalawreview.org/the-need-for-
      a-research-culture-in-the-forensic-sciences-2/

[3]     Risinger D.M. (2013). Reservations about likelihood ratios (and some other aspects of forensic 'Bayesianism'). *Law, Probability and Risk*, 12, 63–73, http://dx.doi.org/10.1093/lpr/mgs011

[4]     President's Council of Advisors on Science and Technology (2016). *Forensic Science in Criminal Courts: Ensuring Scientific Validity of Feature-Comparison Methods*. https://obamawhitehouse.archives.gov/administration/eop/ostp/pcast/docsreports/

[5]     Morrison G.S., Thompson W.C. (2017). Assessing the admissibility of a new generation of forensic voice comparison testimony. *Columbia Science and Technology Law Review*, 18, 326–434. https://doi.org/10.7916/stlr.v18i2.4022

[6]     Cole S.A. (2006). Is fingerprint identification valid? Rhetorics of reliability in fingerprint proponents' discourse. *Law & Policy*, 28, 109–135. https://doi.org/10.1111/j.1467-9930.2005.00219.x

[7]     Morrison G.S. (2014). Distinguishing between forensic science and forensic pseudoscience: Testing of validity and reliability, and approaches to forensic voice comparison. *Science & Justice*, 54, 245–256. http://dx.doi.org/10.1016/j.scijus.2013.07.004

[8]     Koehler J.J. (2017). Forensics or fauxrensics? Ascertaining accuracy in the forensic sciences. *Arizona State Law Journal*, 49(4), 1369–1416. https://arizonastatelawjournal.org/2018/02/07/forensics-or-fauxrensics-ascertaining-accuracy-in-the-forensic-sciences/

[9]     Morrison G.S., Ballantyne K., Geoghegan P.H. (2018). A response to Marquis et al (2017) What is the error margin of your signature analysis? *Forensic Science International*, 287, e11–e12. https://doi.org/10.1016/j.forsciint.2018.03.009

[10] Morrison G.S., Neumann C., Geoghegan P.H. (2020). Vacuous standards – subversion of the OSAC standards-development process. *Forensic Science International: Synergy*, 2, 206–209. https://doi.org/10.1016/j.fsisyn.2020.06.005

[11] Morrison G.S., Neumann C., Geoghegan P.H., Edmond G., Grant T., Ostrum R.B., Roberts P., Saks M., Syndercombe Court D., Thompson W.C., Zabell S. (2021). Reply to Response to Vacuous standards – subversion of the OSAC standards-development process. *Forensic Science International: Synergy*, 3, 100149. https://doi.org/10.1016/j.fsisyn.2021.100149

[12] Meuwly D. (2001). *Reconnaissance de locuteurs en sciences forensiques: l'apport d'une approche automatique*. Doctoral dissertation, University of Lausanne. https://www.unil.ch/files/live/sites/esc/files/shared/These.Meuwly.pdf

[13] Brümmer N., du Preez J. (2006). Application independent evaluation of speaker detection. *Computer Speech and Language*, 20, 230–275. https://doi.org/10.1016/j.csl.2005.08.001

[14] Morrison G.S. (2011). Measuring the validity and reliability of forensic likelihood-ratio systems. *Science & Justice*, 51, 91–98. http://dx.doi.org/10.1016/j.scijus.2011.03.002

[15] Meuwly D., Ramos D., Haraksim R. (2017). A guideline for the validation of likelihood ratio methods used for forensic evidence evaluation. *Forensic Science International*, 276, 142–153. http://dx.doi.org/10.1016/j.forsciint.2016.03.048

[16] Ramos D., Meuwly D., Haraksim R., Berger C.E.H. (2020). Validation of forensic automatic likelihood ratio methods. In Banks D., Kafadar K., Kaye D.H., Tackett M. (Eds.), *Handbook of Forensic Statistics* (pp. 143–163). Boca

Raton, FL: CRC. https://doi.org/10.1201/9780367527709

[17]   Morrison G.S., Enzinger E., Hughes V., Jessen M., Meuwly D., Neumann C., Planting S., Thompson W.C., van der Vloed D., Ypma R.J.F., Zhang C., Anonymous A., Anonymous B. (2021). Consensus on validation of forensic voice comparison. *Science & Justice*, 61, 229–309. https://doi.org/10.1016/j.scijus.2021.02.002

[18]   Morrison G.S., Enzinger E. (2016). Multi-laboratory evaluation of forensic voice comparison systems under conditions reflecting those of a real forensic case (forensic_eval_01) - Introduction. *Speech Communication*, 85, 119–126. https://doi.org/10.1016/j.specom.2016.07.006

[19]   Morrison G.S., Enzinger E. (2019). Multi-laboratory evaluation of forensic voice comparison systems under conditions reflecting those of a real forensic case (forensic_eval_01) - Conclusion. *Speech Communication*, 112, 37–39. https://doi.org/10.1016/j.specom.2019.06.007

[20]   Weber P., Enzinger E., Labrador B., Lozano-Díez A., Ramos D., González-Rodríguez J., Morrison G.S. (2022). Validation of the alpha version of the $E^3$ Forensic Speech Science System ($E^3FS^3$) core software tools. *Forensic Science International: Synergy*, 4, 100223. https://doi.org/10.1016/j.fsisyn.2022.100223

[21]   Basu N., Bali A.S., Weber P., Rosas-Aguilar C., Edmond G., Martire K.A., Morrison G.S. (2022). Speaker identification in courtroom contexts – Part I: Individual listeners compared to forensic voice comparison based on automatic-speaker-recognition technology. *Forensic Science International*, 341, 111499. https://doi.org/10.1016/j.forsciint.2022.111499

[22]   Kirchhübel C., Brown G., Foulkes P. (2023). What does method validation look like for forensic voice comparison by a human expert? *Science & Justice*,

63, 251–257. https://doi.org/10.1016/j.scijus.2023.01.004

[23]   Morrison G.S., Zhang C. (2023). Forensic voice comparison: Overview. In Houck M., Wilson L., Eldridge H., Lewis S., Lothridge K., Reedy P. (Eds.), *Encyclopedia of Forensic Sciences* (3rd Ed.), vol. 2, pp. 737–750. Elsevier. https://doi.org/10.1016/B978-0-12-823677-2.00130-6

[24]   Morrison G.S. (2018). Admissibility of forensic voice comparison testimony in England and Wales. *Criminal Law Review*, 2018(1), 20–33. [Preprint available at http://geoff-morrison.net/#Admissibility_EW_2018]

[25]   Association of Forensic Science Providers (2009). Standards for the formulation of evaluative forensic science expert opinion. *Science & Justice*, 49, 161–164. http://dx.doi.org/10.1016/j.scijus.2009.07.004

[26]   Willis S.M., McKenna L., McDermott S., O'Donell G., Barrett A., Rasmusson A., Nordgaard A., Berger C.E.H., Sjerps M.J., Lucena-Molina J.J., Zadora G., Aitken C.G.G., Lunt L., Champod C., Biedermann A., Hicks T.N., Taroni F. (2015). *ENFSI Guideline for Evaluative Reporting in Forensic Science*. http://enfsi.eu/wp-content/uploads/2016/09/m1_guideline.pdf