

OSAC 2023-N-0023 Standard Guide to Forensic Speaker Recognition Landscape

*Speaker Recognition Subcommittee
Digital/Multimedia Scientific Area Committee
Organization of Scientific Area Committees (OSAC) for Forensic Science*

DRAFT





DRAFT OSAC Proposed Standard

OSAC 2023-N-0023 Standard Guide to the Forensic Speaker Recognition Landscape

Prepared by
Speaker Recognition Subcommittee
Version 1.0
September 2023

Disclaimer

This OSAC Proposed Standard was written by the [insert subcommittee or other unit name] of the Organization of Scientific Area Committees (OSAC) for Forensic Science following a process that includes an [open comment period](#). This Proposed Standard will be submitted to a standards developing organization and is subject to change.

There may be references in an OSAC Proposed Standard to other publications under development by OSAC. The information in the Proposed Standard, and underlying concepts and methodologies, may be used by the forensic-science community before the completion of such companion publications.

Any identification of commercial equipment, instruments, or materials in the Proposed Standard is not a recommendation or endorsement by the U.S. Government and does not imply that the equipment, instruments, or materials are necessarily the best available for the purpose.

Acknowledgements

The Speaker Recognition Subcommittee thanks the following contributing authors who have provided substantial language to the document on multiple occasions: Sandra Ferrari Disner, Michael Jessen, and Finnian Kelly; as well as the following reviewers who provided valuable and feedback on content: Anil Alexander, John Hansen, Larry Kincaid, David Marks, Stephanie Domitrovich, and Jennifer Stathakis. The Subcommittee also thanks Erica Gold who has been the editor through the document's many revisions as well as lead author.

DRAFT

Table of Contents

1. Scope	Page 5
2. Referenced Documents	Page 5
3. Terminology	Page 5
4. Summary of Practice	Page 5
5. Significance of Use	Page 6
6. Procedure	Page 8
7. Annex: Bibliography	Page 11

DRAFT

1 **1. Scope**

2 Forensic speaker recognition, also referred to or covered by the terms forensic speech science,
3 forensic phonetics, and speaker identification, aims to determine whether speakers are likely
4 to be the same person or different people from at least two recordings (e.g., known and
5 questioned recordings). This document provides a landscape of the methods used for analysis
6 in the field of speaker recognition as well as the commonly used interpretation frameworks –
7 Conclusion Frameworks. This document also establishes that the wider speaker recognition
8 community has rejected previously held beliefs regarding the scientific validity of
9 voiceprinting. This document is intended to serve as a general overview and reference (as it is
10 currently practiced in the field) for forensic speaker recognition.

11
12 **2. Referenced Documents**

13 2.1 *OSAC Standards.*

14 2.2 Organization of Scientific Area Committees (OSAC) for Forensic Science Speaker
15 Recognition Subcommittee, “Essential scientific literature for human-supervised automatic
16 approaches to forensic speaker recognition.”¹

17
18 **3. Terminology**

19 For purposes of this document, the following definitions and acronyms apply.

20
21 3.1 Definition of terms specific to this standard.

22
23 3.2 **Automatic Speaker Recognition (ASR)**, *n.* as used in this guideline, ASR requires the
24 use of specialized software to compare speech samples, producing a numerical score that is
25 evaluated from the perspective of the same-speaker origin as well as the different-speaker
26 origin.

27
28 **4. Summary of Practice**

29
30 4.1 The main objective of forensic speaker recognition is to determine whether speakers are
31 likely to be the same or different. Forensic practitioners are typically presented with a
32 minimum of two audio recordings and asked to carry out an analysis of those audio recordings.
33 The methods used for analysis in forensic speaker recognition have evolved far past previous
34 traditions of voiceprinting, which has been rejected and discredited by the speaker recognition
35 community. Methods of analysis that are currently in practice include: auditory phonetic
36 analysis, acoustic phonetic analysis, semi-automatic acoustic analysis, automatic speaker
37 recognition, human-assisted speaker recognition, and combined human and automatic speaker
38 recognition analysis.

39

¹ Prepared by Scientific Literature Working Group, Forensic Speaker Recognition Subcommittee, “Essential scientific literature for human-supervised automatic approaches to forensic speaker recognition,” Organization of Scientific Area Committees (OSAC), Online, Available: <https://www.nist.gov/document/essentialscientificliteratureforhuman>

40 Indeed, just as there are multiple methods for analysis being implemented in speaker
41 recognition, there are also a number of different conclusion frameworks that have also been
42 adopted. Interpretation frameworks (or Conclusion Frameworks) that are currently being
43 utilized by the speaker recognition community include: the binary decision, probability scales,
44 likelihood ratios (both verbal and numerical), the UK Position Statement, and support
45 statements.

46

47 5. Significance of Use

48

49 5.1 Introduction: Forensic speaker recognition involves the comparison of at least two speech
50 samples (typically from a questioned and known recording to determine whether speakers are
51 likely to be the same or different. It is common for forensic speaker recognition to be referred
52 to by a few other terms, largely dictated by the field in which the subject is researched and
53 practiced. Within the forensic speech science and forensic phonetics communities, the task is
54 often referred to as forensic speaker comparison or forensic voice comparison. Within the
55 engineering communities, forensic speaker recognition is also sometimes referred to as
56 forensic speaker identification. For clarification purposes, the task of comparing speech
57 samples is referred to in this document as forensic speaker recognition. This document is
58 intended to provide a general overview of forensic speaker recognition. For more detailed
59 information about the topics discussed in this document, please see the OSAC Speaker
60 Recognition Essential Literature document.²

61

62 5.2 The objective of the forensic practitioner carrying out forensic speaker recognition is to
63 provide the trier(s) of fact with an informed opinion regarding the probability of obtaining the
64 evidence (under the hypothesis that the samples came from the same person, versus under the
65 hypothesis that two different speakers produced each sample). This objective can be reached
66 by practitioners using a variety of methods (i.e., auditory phonetic and acoustic phonetic
67 analysis, semi-automatic acoustic analysis, automatic speaker recognition, human-assisted
68 speaker recognition, or combined human and automatic speaker recognition analysis). While
69 questioned recordings often involve an array of different sounds and speech, the task of
70 forensic speaker recognition is wholly concerned with the speech (and sounds) produced by
71 individuals. Those sounds that cannot be attributed to a person are outside the scope of forensic
72 speaker recognition, and fall more into the general area of audio or acoustic forensics. The aim
73 of this document is not to promote or suggest any one method of analysis or interpretation
74 framework over another, but rather to provide a general landscape of the methods used within
75 the speaker recognition community.

76

77 5.3 *Voiceprinting*: While there are many ways to conduct forensic speaker recognition, it is
78 important to note here that the method known as "voiceprinting" is not supported by the
79 scientific community, and has been discredited. The term "voiceprint²" was coined by the
80 author of an article (Kersta 1962) which appeared over a half-century ago. The name chosen
81 for that methodology quite transparently implied parallels between a (never-proven) "theory
82 of invariant speech" and the relative invariance of fingerprints.

83

84

² Gray and Kopp (1944) also used the term voiceprint with the same definition, however, they used the term with a space between the words voice and print. For all intents and purposes this document uses the term voiceprint without a space.

85 5.3.1 The so-called "voiceprints" were the product of sound spectrography, a technology
86 carried forward even to the present day, which is still of great utility to speech scientists.
87 The most notable scientific failing of the voiceprint method was that it did not provide
88 examiners with the vocal output (i.e., the audio). This inevitably obscured the phonetic
89 nature of the patterns of acoustic energy and reduced the analysis to a simple pattern-
90 matching exercise. Nevertheless, that exercise was initially heralded with outsized claims
91 of success. The article that introduced the voiceprinting method in 1962 reported that
92 phonetically naïve examiners³ were able to identify a target voice with 99% accuracy,
93 even from a pool of a dozen speakers. Not surprisingly, this new methodology soon caught
94 the attention of law enforcement, and was presented as evidence in a number of criminal
95 prosecutions, in the US and elsewhere.

96
97 5.3.2 However, the scientific community remained skeptical. Well-known phoneticians
98 such as Peter Ladefoged and Harry Hollien reported that mere pattern matching (which is
99 all that the young voiceprint examiners were asked to do) was incapable of yielding the
100 astonishing results reported in the 1962 article. Due to the variability present in speech
101 productions from sample to sample, spectrographic template matching is not effective and
102 it is inconsistent in speaker recognition work. In time, phoneticians began to provide
103 expert testimony against the admissibility of voiceprint evidence, and in consequence, a
104 number of lower-court convictions were eventually overturned. In response to these
105 criticisms, another academic linguist, Oscar Tosi, initiated a more rigorous, and
106 procedurally transparent, voiceprint study (Tosi et al. 1972). This yielded less vertiginous,
107 but more scientifically reliable results – 6% false identification errors and 13% false
108 elimination errors, under laboratory conditions. Still, given the high stakes of introducing
109 a still largely unsupported procedure into courts of law, a report issued by the National
110 Research Council⁴ concluded that the voiceprint method lacked an adequate scientific
111 basis for estimating reliability in many practical situations, pointing out in addition that
112 laboratory evaluations of the voiceprint method showed increasing errors as the conditions
113 for evaluation moved toward real-life situations, such as poor signal-to-noise ratios and
114 dissimilar recording conditions.

115
116 5.3.3 The Federal Rules of Evidence, adopted in 1975, further challenged the voiceprint
117 methodology by shifting the standards for admissibility in favor of practitioners whose
118 "scientific, technical, and other specialized knowledge" can help the trier of fact "to
119 understand the evidence or to determine a fact in issue"⁵ Phonetically untrained voiceprint
120 examiners, who sought to identify speakers simply by looking at pictures of their voice
121 signals, were left at a marked disadvantage.

122
123

³ To drive home his point, Kersta used high school students, who had been given only one week of training, as his examiners. The difficulty of the task was augmented by a forced-choice design; "not sure" was not an option.

⁴ National Research Council. 1979. *On the Theory and Practice of Voice Identification*. Washington, DC: The National Academies Press. <https://doi.org/10.17226/19814>.

⁵ *Federal Rules of Evidence* 702

124 5.3.4 For further overviews of the voiceprint approach one can consult (this is not
125 intended to be an exhaustive list, but merely a few selected references):

126
127 5.3.4.1 Hollien (1990)

128 5.3.4.2 Hollien (2002)

129 5.3.4.3 Rose (2002)

130 5.4 While voiceprinting has been discredited from the speaker recognition community, it has
131 also been declared inadmissible in court (for further information see U.S. v. Angleton 2003).

132
133 5.5 Additional methods that are properly applied and under certain circumstances are
134 recognized as appropriate by the community are, in turn, described in detail below. It will
135 become apparent that these methods have historically developed in parallel to one another, as
136 they have grown out of different disciplines. However, it is not uncommon to see some cross-
137 over between the various methods used in forensic speaker recognition. This will be explained
138 further in the sections that follow.

139 140 6. Procedure

141 6.1 There is no one, single method that is used by all practitioners of forensic speaker
142 recognition, and it is sometimes the case that some of these methods are combined when
143 undertaking analysis. The methods most commonly employed in forensic speaker recognition
144 are: auditory phonetic analysis, acoustic phonetic analysis, auditory phonetic + acoustic
145 phonetic analysis, semi-automatic acoustic speaker recognition, automatic speaker
146 recognition, human-assisted automatic speaker recognition, and a combination of auditory
147 phonetic + acoustic phonetic analysis and (human-assisted) automatic speaker recognition. All
148 seven approaches to speaker recognition are detailed below.

149
150 6.2 *Auditory Phonetic Analysis (AuPA)*: AuPA is defined as the process by which “the expert
151 listens analytically to the speech samples and attends to aspects of speech at the segmental and
152 suprasegmental levels” (Gold and French 2011). AuPA is very important in the identification
153 of language varieties, such as regional accent, foreign accent or in the detection of linguistic
154 correlates of various social factors. Age, sex, and gender also fall into the category of
155 characteristics most commonly judged auditorily. All of these can be classified as “group-level
156 characteristics,” in contrast to “individual-level characteristics” (Hughes and Rhodes 2018).
157 Group-level speaker characteristics are crucial in speaker profiling, but they also have their
158 established place in speaker recognition. They can be important in defining the relevant
159 population. They are also particularly powerful as evidence speaking against speaker identity:
160 if, for example, the known and the unknown voices use two different regional varieties, it is
161 likely that they are different individuals, given that bi-dialectalism is relatively rare. Group-
162 level characteristics can also provide important information that supports inclusion, within the
163 context of the case. They are also particularly helpful in excluding speakers. Part of the

164 tradition of AuPA has been to narrow down dialect to a point that it achieves the status of a
165 rare combination of linguistic parameters only used by a few individuals (and ideally, however,
166 rather unrealistic, just one individual). This can occur when a speaker uses only some features
167 of a dialect (or other language variety) to the exclusion of others, or if features from various
168 language varieties are combined. Discussions of these aspects related to language variety are
169 provided in Jessen (2010; 2021) and Hughes & Rhodes (2018).

170
171 6.2.1 AuPA is also used for the description and interpretation of various individual-level
172 speaker characteristics. According to the survey by Gold and French (2011), voice quality
173 is a particularly important one. Voice quality can be measured acoustically (Keating et al.
174 2015), but given the acoustic limitations typically encountered in forensic casework, most
175 of these methods suffer from information loss or lack of applicability. Auditory analysis,
176 instead, offers more robustness, though auditory analysis is also more subjective. Auditory
177 voice quality assessment in forensics often builds upon the classificatory framework of
178 the phonetician John Laver (1980), particularly in Europe. A description of that
179 framework and how it is adapted to forensics is found in San Segundo et al. (2019), and a
180 complete definition of voice quality can be found in McIntyre et al. (2021). Another
181 classification framework for AuPA-based analysis has been developed for disfluency
182 patterns, which include silent pauses, breathing pauses, filled pauses (utterance such as
183 uh, and um) or sentence interruptions (McDougall and Duckworth 2018, de Boer and
184 Heeren 2019, Hughes et al. 2016). Further speaker characteristics observed auditorily are
185 listed in Gold and French (2011).

186
187 6.3 *Acoustic Phonetic Analysis (AcPA)*: AcPA is the method by which “the expert analyzes
188 and quantifies physical parameters of the speech signal using computer software. As with
189 AuPA, this is labor intensive, involving a high degree of human input and judgment” (Gold
190 and French 2011). AcPA traditionally has its strongest focus on speaker characteristics that
191 have an anatomical motivation and that have been known since the 1950s to vary between
192 women, men, and children, but also between individuals within these larger speaker categories.
193 This applies to fundamental frequency and formant frequencies.

194
195 6.3.1 Fundamental frequency (f_0), which is the frequency of the vibration patterns of the
196 vocal folds, depends on the size of the larynx (especially vocal fold length), but it is to a
197 degree controllable for linguistic purposes. As a way of disregarding locally determined
198 linguistic factors, a common method is to average f_0 (mean, mode, or median) across long
199 utterances or the entire recording (Hudson et al. 2007). In this process, further irrelevant
200 factors that have a strong influence on f_0 must be controlled as much as possible; this is
201 particularly important for the f_0 -raising effect of vocal effort (that is, speaking loudly)
202 (Jessen et al. 2005). Speakers can also differ habitually in terms of how “melodically”
203 they speak (scale from speaking monotonously to highly modulated). Standard deviation
204 of f_0 across long passages or the entire recording is a way of capturing these habitual
205 speaker differences. Since mean and standard deviation of f_0 are to some extent correlated,
206 variability is sometimes expressed by the coefficient of variation (standard deviation
207 divided by mean), by means of which the correlation almost disappears (Jessen et al.
208 2005).

209 6.3.2 Vowel formant frequencies, which are characteristic patterns of amplitude peaks in
210 the speech spectrum, are associated with the length of the speaker's vocal tract and other
211 anatomical features. However, formant frequencies – especially the first formant (F1) and
212 the second formant (F2) – are also crucial carriers of linguistic information; they are the
213 main correlates of vowel distinctions in a language, and transitions between successive
214 vowels serve to distinguish any intervening consonants. There is thus a clear need to
215 control for these linguistic factors. One way, which is analogous to the processing of f0,
216 is to average all the formants across long stretches of speech. This method is referred to
217 as long-term formant analysis (Nolan and Grigoras 2005). Another method is to measure
218 formant frequencies separately for different vowels. This is the traditional way formants
219 are measured in phonetics. Beyond anatomical restrictions, there are degrees of freedom
220 in transitioning from one target sound to the next. Hence, measuring formant dynamics is
221 a third way of capturing formant information in forensic speaker recognition. It has been
222 shown in many studies that when formant measurements are not limited to targets but,
223 rather, when one takes into account the entire dynamics of the formant movements,
224 speaker recognition capability is improved (see McDougall 2006 and Morrison 2009 for
225 some of the early studies).

226
227 6.3.3 There are other speaker characteristics that can be based upon acoustic phonetic
228 analysis (Gold and French 2011). For example, it is possible to measure the spectral
229 energy distribution in fricatives or nasals (Kavanagh 2012) or to make temporal
230 measurements in the domain of rhythm and timing (Dellwo et al 2015; Plug et. al. 2021).
231 But most actual AcPA casework utilizes f0 and formants.

232
233 6.4 *Auditory Phonetic + Acoustic Phonetic Analysis (AuPA+AcPA)*: AuPA+AcPA is the
234 combination of both auditory and acoustic analysis in speaker recognition as detailed in §6.2
235 and §6.3. The combination of AuPA and AcPA has also been referred to as an “auditory-
236 acoustic-phonetic” method that is carried out by forensic practitioners (Morrison et al. 2016).”
237 The term “auditory-acoustic-phonetic by forensic practitioners (qualitative opinion)” reflects
238 how the phonetic data are traditionally interpreted by many practitioners of AuPA+AcPA:
239 though there can be quantification on the feature level (e.g., formant frequencies in Hz; values
240 on a scale of perceived voice qualities), the results are most commonly interpreted
241 qualitatively, e.g., as distances of the values of the unknown and known speaker that is visible
242 in a plot of formant values, or as an experience-based judgment of how frequently a certain
243 speaker characteristic occurs in a relevant population (Morrison et al. 2016). Such a qualitative
244 expression of AuPA+AcPA can approach a quantitative likelihood ratio-based method if there
245 are data available of the relevant population, for example of the f0 values of male speakers
246 (Gold and French 2019). But for full expression of quantitative likelihood ratios, the statistical
247 methodology has to be present, as well as all the necessary data, such as non-contemporary
248 same-speaker and different-speaker data.

249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281

Annex

Bibliography

This is not meant to be an all-inclusive list as the group recognizes other publications on this subject may exist. At the time this document was drafted, these were some of the publications available for reference. Additionally, any mention of a particular software tool or vendor as part of this bibliography is purely incidental, and any inclusion does not imply endorsement by the authors of this document.

American Nurses Association. (2021). Nursing: Scope and standards of practice (4th ed.).
Nursesbooks.org.

Cambridge University Press. (n.d.). Meanings & definitions. Cambridge Dictionary.
<https://dictionary.cambridge.org/>

National Commission on Forensic Science. (n.d.). Views document on definitions.
<https://www.justice.gov/archives/ncfs/page/file/477836/download>

National Institute of Justice. (2020). National best practices for sexual assault kits: A
multidisciplinary approach. Office of Justice Programs, U.S. Department of Justice.

Organization of Scientific Area Committees (OSAC) for Forensic Science, Crime Scene
Investigation Subcommittee. (2021). Guiding Principles for Scene Investigation and
Reconstruction (OSAC 2021-N-0015). OSAC, National Institute of Standards and
Technology, US Department of Commerce.
https://www.nist.gov/system/files/documents/2021/09/02/OSAC%202021-N0015%20Guiding%20Principles%20for%20CSI_FINAL%20OSAC%20PROPOSED%20FOR%20REGISTRY.pdf

Spellman, B. A., Eldridge, H., & Bieber, P. (2021). Challenges to reasoning in forensic science
decisions. *Forensic Science International. Synergy*, 4, 100200.
<https://doiorg.mutex.gmu.edu/10.1016/j.fsisyn.2021.100200>