

Speech Collection Guideline for Speaker Recognition: Audio Collection at a Temporary Location

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



OSAC Proposed Standard

Speech Collection Guideline for Speaker Recognition: Audio Collection at a Temporary Location

Prepared by
Speaker Recognition Subcommittee
Version: 1.12
June 2018

Disclaimer:

This document has been developed by the Speaker Recognition Subcommittee of the Organization of Scientific Area Committees (OSAC) for Forensic Science through a consensus process and is *proposed* for further development through a Standard Developing Organization (SDO). This document is being made available so that the forensic science community and interested parties can consider the recommendations of the OSAC pertaining to applicable forensic science practices. The document was developed with input from experts in a broad array of forensic science disciplines as well as scientific research, measurement science, statistics, law, and policy.

This document has not been published by a SDO. Its contents are subject to change during the standards development process. All interested groups or individuals are strongly encouraged to submit comments on this proposed document during the open comment period administered by the Audio Engineering Society (www.aes.org).

1.0 Introduction

In this document, the person overseeing the collection session will be designated the “interviewer” and the individual being recorded the “subject”. The intended audience for this guideline is those interviewers called on to perform speech collection using portable equipment at a temporary location not originally designed or intended for audio recording.

The goal of speech collection is to collect an audio recording containing a combination of subject identifying information (such as their name, date of birth, etc.) and speech which can be used for future comparison against the speech of unknown speakers using unspecified speaker recognition methods. Although the specific method of speaker recognition is left undefined, automated/semi-automated computer-based methods were the primary driver for some of the specific parameters found in this guideline.

These guidelines should be viewed as providing minimum requirements for usable speech collection in an operational or field environment, and **are not intended as data collection guidelines for research applications**, speech intended for transcription, or other applications.

It is important that before implementation of this guideline the user coordinate their activities with any elements of their parent organization which will be storing and using the collected data. Issues related to audio channel (microphone, telephone, radio, etc.), desired languages and dialects, data formats, etc. should be worked out beforehand to ensure that the data is maximally useful. Similarly, the proper storage and protection of personal identifying information (if collected) should also be coordinated as required within your organization prior to data collection.

2.0 Scope

This guideline is intended to be one of a series, and covers only one scenario - the collection of speech samples for speaker recognition at a temporary, non-laboratory location. One example of this would be the collection of speech samples from subjects during some type of field activity. The field activity could be associated with a range of purposes, such as law enforcement, intelligence, military or sociological. This guideline presumes portable resources and somewhat limited time to perform the collection. It also assumes that the interviewer is fluent in the subject’s language, or that an interpreter is present who is both fluent in the interviewer’s and subject’s languages and cognizant of the goals of the interview.

This guideline does not deal with details related to the handling, transmission, storage, or preservation of collected data. Specifically:

- It also does not deal with any issues related to collection, storage or protection of personal information.

- It does not recommend how to protect speaker recognition analysts from seeing or hearing subject identifying information which could result in biased analysis results.

It is the responsibility of the guideline's user to determine what their organization's rules and policies are on these matters and to tailor their implementation of this guideline to comply with those rules and policies. It is also the guideline user's responsibility to learn how to operate the selected equipment, especially the placement orientation and distance of the microphone and the interpretation of any meters on the recording device.

This guideline also does not deal with possibly important concerns such as personnel and equipment security, power sources for equipment, etc. The alleviation of such highly situation specific concerns is the responsibility of the guideline's user or their agency.

3.0 Terms and Definitions

A/D	Analog-to-digital
AGC	Automatic Gain Control. A closed-loop regulating circuit which provides a controlled signal amplitude at its output, despite variation of the amplitude in the input signal.
Codec	Algorithm designed to encode or decode a stream of digital audio data.
Hz	Hertz. A one Hz sample rate equals one sps.
Mbyte	Megabyte (of digital storage) – this normally refers to 1024*1024 (1,048,576) bytes in digital applications though others may intend 1000*1000 or 1,000,000 bytes.
MD5	A widely used cryptographic hash function producing a 16-byte hash value, typically expressed in text format as a 32 digit hexadecimal number. MD5 is commonly used to verify data integrity.
min	Minute (of time)
MP3	A widely used audio file format which uses a lossy audio encoding algorithm defined in the MPEG-1 standard, Audio layer 3. The details of this standard are which is published as ISO/IEC 11172-3.
PCM	Pulse Code Modulation. Refers to a method of representing an analog audio waveform with a series of quantized digital sample values.
PCM-WAV	A version of the WAV file format which saves the data as uncompressed linear PCM samples with a standard RIFF header.
RIFF	The Resource Interchange File Format is a generic file container format which can be used to store audio data.
sps	Sample(s)-per-second.
USB	Universal Serial Bus. Refers to a family of standardized computer peripheral interfaces.
WAV	A specific implementation of the RIFF file format for audio data.
WMA	Windows Media Audio. This refers to both a Microsoft proprietary audio file format and the audio codecs it uses. WMA can use either lossy or lossless codecs.

4.0 Speech Collection Scenario: Audio Collection at a Temporary Location

4.1 Collection Environment

The collection environment should:

- Be an indoor space as free as possible from background noises such as air conditioners, generators, fans, or other motorized or electrical devices. Avoid locations that have music, white noise, or other audio playing in the background at any audio level. It may be necessary to turn the interference source off to fully mitigate it.
- Be a location that is not near outside traffic noise (human, animal, vehicular, or aircraft).
- Have a minimum of large, flat, hard sound-reflective surfaces which can cause reverberation and echoes. The effects of a reverberant room can be mitigated by hanging fabric (curtains, blankets, etc.) or other sound deadening materials on the walls or as dividers in the room.
- Allow the subject to be as comfortable as possible, preferably sitting, to lower cognitive/voice stress levels and to facilitate natural conversation.

4.2 Collection Equipment

Although high-quality digital audio recording equipment is preferred, recordings made with equipment meeting the minimum requirements detailed below should be useable by most speaker recognition methods. If there are multiple recording sites, the same type of recording system should be used at all collection locations to minimize recording differences due to equipment variation. Nothing in these requirements precludes the concurrent use of multiple recording devices if that is required for the intended application. An example of such a requirement would be to create a pair of recordings in which one is a high quality reference while the other is condition-matched to a specific use case.

Recording devices should fulfill the following requirements:

- The speech must be recorded digitally and saved as uncompressed PCM data with *at least* 16-bit samples at a minimum rate of 16,000 sps. The audio can be mono or stereo. Recording at a higher sample rate and bit depth is greatly preferred if that is possible. Many current devices support the recording of 16 or 24-bit samples at rates up to 48 kHz and storage of the recorded data in PCM-WAV format.¹

¹ If the recording equipment used samples at rates lower than these recommendations, the recordings may still be suitable for some applications but must be seen as information-losing and thus suboptimal.

Note: the digital re-capture of audio originally on analog tape DOES NOT fulfill these minimum requirements. The recording of the analog reproduction of a digital recording (such as the headphone output of a digital audio device) is also not acceptable.

- The audio should be saved in a standard lossless file format such as PCM-WAV, or be in a file which can be converted to a standard format without loss of fidelity. The audio should not be saved in a file format such as MP3 or WMA which use a lossy codec to compress the audio data.
- Any type of automatic gain control (AGC) on the microphone or recorder should be turned off/disabled during the recording session.
- For recordings made using laptop or other computers, it is preferred to use an external USB condenser microphone with an on-board analog-to-digital (A/D) converter. This is because the internal microphone or external microphones plugged into a “mic” port can pick up noise from internal circuitry.
- The subject’s microphone should ideally be a headset mic since:
 - it fixes the location of the mic with respect to the mouth.
 - it reduces interference from the interviewer’s speech and any background sounds.
 - the speaker more quickly forgets about its presence.Otherwise, a microphone on a stable stand or tripod which places it at an appropriate distance from the subject for the type of microphone being used is acceptable. If the recording device/microphone is directional, it should be situated to best pick up the subject’s speech. If not chosen to mimic an operational scenario, the microphone should ideally provide a flat frequency response.
- A foam or fabric breath guard (commercial or improvised) should be used with the microphone to mitigate puffing noises from the subject’s breath or burst type speech sounds such as stop consonants (i.e., /t/,/p/,/k/) if possible.
- The interviewer should have some indicator available on the recording device that shows that the audio is being recorded at an appropriate amplitude level and not too low (resulting in a noisy recording due to quantization effects) or too high (which causes clipping and thereby introduces nonlinear distortion into the audio stream). The relationship of appropriate sound levels and equipment indicator displays is not standard, therefore the indicator level providing acceptable audio should be determined by testing before the collection activity.
- There should be some method available to back-up the collected data, such as writing it to optical media, external hard drives, or USB thumb drives.

4.3 Speech Collection

After the collection environment and equipment have been arranged, the interviewer should record and audibly review an initial sample of test speech in the same recording environment and

using the same equipment as for the collection to confirm that the equipment is working properly and the audio quality meets the parameters discussed above. This may also expose other sources of noise not originally noted, such as the buzz of fluorescent lights or sounds from air handlers, which can be addressed as discussed in section 5.1. Once the setup is verified, it should be documented, ideally including the model identification and serial numbers of the equipment and a diagram or photographs of how it was connected and arranged.

Once recording begins, either the interviewer or the subject must provide a preamble with some subject identifying information along with the date, time and location of the recording session. The preamble information may be elicited with a set of fixed questions to provide identity information (such as name, address, and date of birth), or it may consist of the statement of a unique identifier, such as an identifying number which relates the subject to separately documented identifying information. If of value to the speaker recognition method used by the interviewer's organization, the identifying questions or statements could be tailored to provide specific known-text or phonetic content.

The type of speech (conversational, reading, preaching, etc.) recorded during the collection should ideally match what is expected to be in the unknown speech samples to be compared against in the future. This is also true of the language used. If the interviewer knows what type of speech and the language the subject recordings will be compared against, the collection should be designed to elicit those. These may not be possible to know beforehand, and the capture of conversational speech is discussed below as a general example.

During the recording, the interviewer should strive to elicit periods of conversational speech from the subject. Conversational speech could be elicited in multiple ways, such as:

- Asking open-ended questions or prompts. A list of possible questions is given in Appendix 1. This list is not exhaustive, and the interviewer should tailor any questions to be appropriate for the circumstances, the subject's culture, etc.
- Asking the subject to describe or interpret an image. These could be simple drawings of an object or scene, photographs of a general nature, or other images that have content understandable by the subject and which will elicit a conversational response. Some examples are provided in Appendix 2. Giving the subject a choice of several images gives them the freedom to choose one for themselves. The interviewer could ask multiple questions about an image to elicit additional speech, such as "are there any dangerous things in the image?" Alternatively, the interviewer could circle some items in the image, and ask the subject to describe them.
- Ask the subject to discuss an article from a local newspaper, news website, or social media outlet.

Note – in the last two methods, the use of paper copies of the image, drawings or newspaper articles should be avoided since their movement can add undesired noises to the recording.

It can be expected that the longer the subject speaks conversationally (presuming that fatigue does not occur), the greater chance that they will become comfortable with the collection

situation, resulting in a more “natural” speech sample. In case the subject does not engage in conversation, it is recommended that 20 or more additional questions similar to those in Appendix 1 be developed beforehand so that an adequate amount of speech can be collected.

The interviewer should avoid interjections while the subject is speaking (e.g., nodding to acknowledge the subject instead of saying "uh-huh").

The subject’s portion of the recording (including the identification segment, any answers to questions, and conversational speech) should contain a minimum of two minutes of speech and preferably up to five minutes. This is to be measured *after* the removal of speech from the interviewer, any noisy segments, and extended pauses.

After the completion of the recording session, the interviewer should document any comments on the collection (via hand written notes, verbal recordings, etc.) prior to beginning the next session.

5.0 Appendices

5.1 Appendix 1: Possible Known-Text Phrases for the Subject

Questions based on the following phrases can be used to capture subject-identifying information and elicit known text responses from the subject. They should also help encourage discussion which will fulfill the desire for the capture of conversational speech. If additional conversational speech is needed, questions should be designed to allow open-ended responses instead of short answers.

- “My name is [Subject name]”
- “I was born on [Subject date of birth]”
- “I was born in [Subject place of birth]” – this can include town, city, region, country, etc.
- “I currently live in [Current residence location] at [address]”
- “My current job/occupation is [Occupation]”
- “I am [married/single] ...”
- “My [height/weight/eye color] is ...”
- “...my wife (husband) is named [name] ...”
- “...and we have [#] children – who are [Names, ages, etc].”
- “My favorite pet’s name is [Boots, Garfield, etc]”
- “My favorite sport is [sport, ex: rugby, running, swimming, etc.]”
- “My first vehicle was a [Vehicle with details, ex: red 1970 Ford Mustang]”
- “The best phone number you can reach me at is [phone number]”

5.2 Appendix 2: Examples of Images to Elicit Speech Samples

The images below are examples of ones that could be used to elicit speech from subjects. As noted above, the images must be understandable by the subject and personally/socially acceptable to them.

(The four photographs below were acquired from the Wikimedia Commons and are attributed to their authors.)



(Geneva Rugby Cup - 20140808 - SF vs LOU. By Clement Bucco-Lechat.)



(A desert: The rainshadow region of Tirunelveli, India. By Arun Ganesh.)



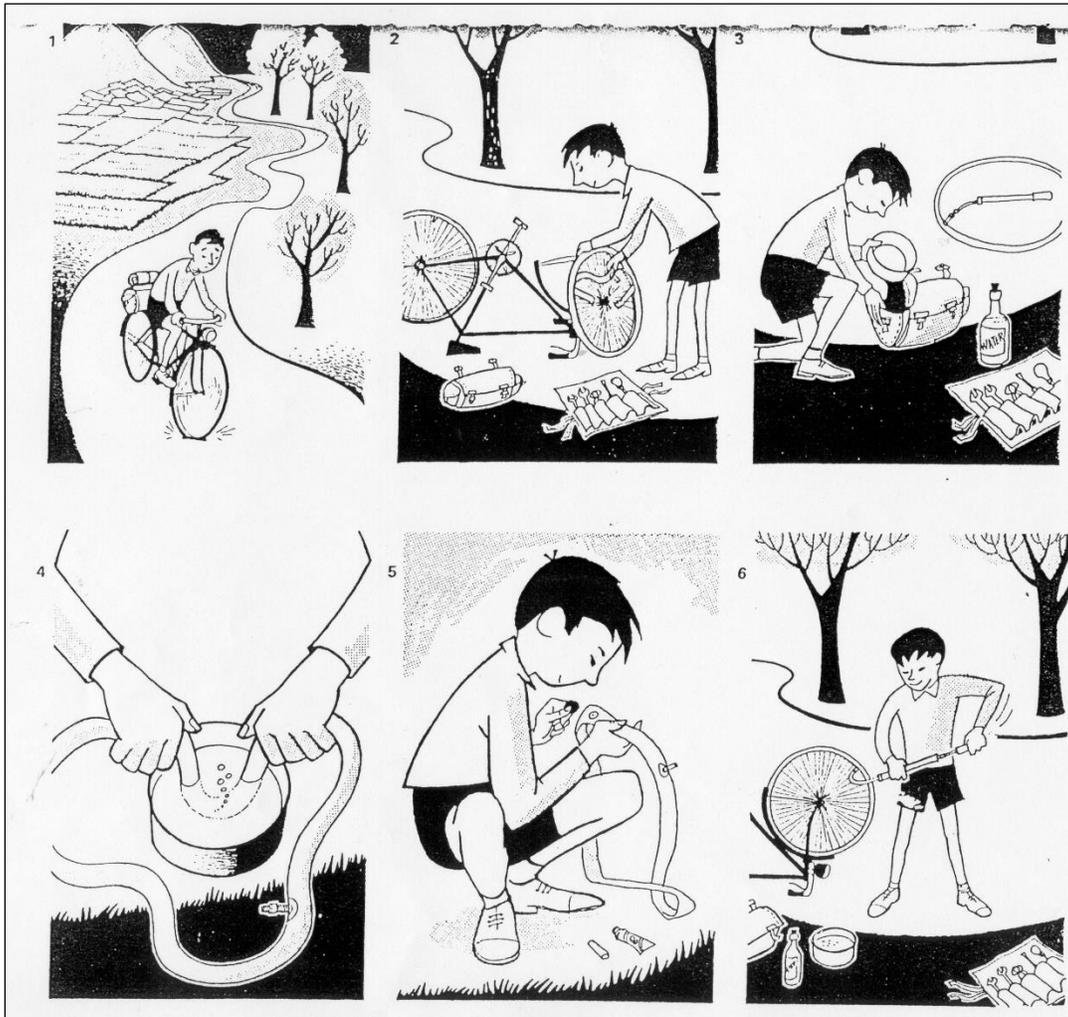
(Patagonia 4x4 Off Road Expeditions by Mil Outdoor Adventure,
El Calafate, Provincia de Santa Cruz, Patagonia Argentina. By total13.net.)



(Philippine rice terraces, Batad village. By [Adi.simionov](#).)



(Camping/Outdoors. Provided by Jos Bouten.)



(Fixing a bicycle tire. Provided by Jos Bouten.)