

Two Decades of Speaker Recognition Evaluation at the National Institute of Standards and Technology

Craig S. Greenberg^a, Lisa P. Mason^b, Seyed Omid Sadjadi^a,
Douglas A. Reynolds^c

^a*NIST ITL/IAD/Multimodal Information Group, MD, USA*

^b*U.S. Department of Defense, MD, USA*

^c*MIT Lincoln Laboratory, MA, USA*

Abstract

The National Institute of Standards and Technology has been conducting Speaker Recognition Evaluations (SREs) for over 20 years. This article provides an overview of the practice of evaluating speaker recognition technology as it has evolved during this time. Focus is given to the current state of speaker recognition evaluation. Highlights from past SREs and future plans are also discussed.

Keywords: NIST SRE, Speaker Recognition, Speaker Recognition Evaluation, Speaker Verification

1. Introduction

The *Information Technology Laboratory* (ITL) at the National Institute of Standards and Technology (NIST) conducts three major activities: 1) fundamental research in mathematics, statistics, and Information Technology (IT); 2) applied IT research and development; and 3) standards development and technology transfer. Part of the ITL, the NIST Speech Group was founded in the mid-1980s to conduct these activities in service of speech-related technologies, and toward that end, held its first evaluation of automatic speech recognition technology in 1987. Since that time, the Speech Group has evolved into the Multimodal Information Group (MIG) at NIST (formerly National Bureau of Standards) and has been conducting *evaluation-driven research* of speech, text, images, video, and multimedia technologies.

Evaluation-driven research is a method of community-focused technology research that utilizes a set of common tasks, data, metrics, and measurement

Preprint submitted to Computer Speech and Language

September 17, 2019

15 methods in order to reduce the total overhead necessary to conduct research
16 and to benchmark the current state of the art and identify the most promis-
17 ing research directions [1]. There are four basic components that make up
18 *evaluation-driven research*: planning, design, assessment, and a workshop.
19 The planning component involves identifying research goals for the technol-
20 ogy (e.g., to be able to improve performance of the fundamental underlying
21 technology or to be robust to certain conditions), obtaining data that sup-
22 ports the evaluation goals, creating and documenting the evaluation plan,
23 as well as identifying and notifying interested researchers and organizations.
24 The design component involves deciding the tasks, metrics, and measure-
25 ment methods that will make up the evaluation, and analyzing the available
26 data to create necessary data sets (e.g., typically some data is provided to
27 researchers in advance of the assessment period to assist in research, and
28 other data is used as test data for the assessment). During the assessment
29 component, either the technology developers or the evaluator runs the sys-
30 tems with the specified test data, and the evaluator analyzes the systems’
31 performances. At the workshop, results and lessons learned are shared and
32 future research goals are identified, which support the planning of future
33 evaluations.

34 In 1996, NIST conducted its first evaluation of technology for automati-
35 cally recognizing speakers by their voices. Over the following two decades¹,
36 NIST conducted 15 Speaker Recognition Evaluations (SREs), in addition
37 to an evaluation held in 2018 and evaluations planned for 2019 and 2020.
38 During that time, speaker recognition technology has evolved substantially,
39 and the SRE series has as well. What started as an evaluation of approx-
40 imately 10 systems completing 4,000 trials has expanded into a series that
41 commonly includes hundreds of systems completing millions of trials. This
42 has been necessary, as the 1996 evaluation would be grossly insufficient for
43 the research needs in 2019, and the 2019 evaluation would have been impos-
44 sible in 1996—specifically, the 1996 SRE data set is too small and the data
45 too easy to analyze performance of modern state-of-the-art systems, and the
46 amount of data and challenging data conditions planned for SRE19 would

¹Over the 20+ years of running the Speaker Recognition Evaluation series, NIST has received support from other U.S. government agencies, such as Department of Defense, Department of Justice and Intelligence Advanced Research Projects Activity (IARPA), to build a forum for the advancement of speaker recognition technology through *evaluation-driven research*.

47 have overwhelmed state of the art systems in 1996.

48 Despite the substantial changes the SRE series has undergone over time,
49 certain elements have remained constant. For example, the goals of the eval-
50 uation series have always been to drive the technology forward, to benchmark
51 the current state of the art, and to identify the most promising research di-
52 rections. The evaluations have also remained open to all researchers working
53 on the general problem of text-independent speaker recognition, and have
54 consistently been designed to focus on core technology issues and to be sim-
55 ple and accessible to those wishing to participate. The requirement that
56 submitted systems must be fully automatic and humans may not listen to,
57 or otherwise interact with, the evaluation data has also been maintained for
58 the entire SRE series.

59 In this article, we present an overview of the NIST ITL/IAD/MIG ap-
60 proach to evaluating speaker recognition technology over the past two decades
61 and provide insights into what evaluations may look like moving into the
62 next decade. The aim is to provide a review of the *evaluation-driven re-*
63 *search* methodology employed by the SRE series that is accessible by new-
64 comers to the field of speaker recognition evaluation. We discuss some of the
65 key considerations necessary when conducting speaker recognition technol-
66 ogy evaluations, and how NIST has addressed evaluating speaker recognition
67 in general and for specific, specialized tasks. A brief survey of past SREs and
68 results from recent evaluations is also provided, as well as a brief overview
69 of plans for the 2019 and 2020 evaluations. We conclude the article with
70 some general projections about how future evaluations may look as research
71 directions have dramatically evolved since the inaugural 1996 SRE.

72 2. Considerations in Evaluating Speaker Recognition Technology

73 There is a great deal that could be said about the considerations necessary
74 when running large-scale research-focused evaluations of speaker recognition
75 technology.² Indeed, NIST has published several lengthy articles covering
76 various aspects of this topic [2, 3, 4, 5]. While still more could be said

²During an informal conversation with a speech researcher, who at that time had recently worked with NIST on creating an evaluation of speaker recognition technology for the IARPA BEST program, he remarked to one of the authors that despite having been a long-time SRE participant, he was shocked by “how much actually had to be taken into account when conducting a speaker recognition evaluation.”

and some material bears repeating, in the interest of focusing this article, we will limit the discussion to three main considerations: task, data, and metrics. It should be noted, however, that driving all decisions must be a set of underlying goals, framed in large part by the current maturity of the technology and the needs of the researchers, system developers, and end-users.

2.1. Task

Speakers are multifarious. Put differently, speech is a behavior, and it varies wildly both within and across individuals. As a result, speaking fixed phrases, reading, and spontaneous text-independent speech are substantially different from one another, and the performances of speaker recognition systems (and the approaches taken by these systems) in these contexts are substantially different as well [6].

Spontaneous text-independent speaker recognition has been recognized as the most general setting for speaker recognition and progress in this area seems most likely to impact other settings [2]. For this reason, NIST has chosen to make spontaneous text-independent speech the focus of the SREs. Even within this setting there are several ways of presenting the task. For example, it could be framed as an *identification task*, where the system must associate each recording with one of a fixed set of speakers (or possibly none of them); a *clustering task*, where systems must partition the speech into an unspecified number of speaker clusters; or a *detection task*, where two recordings are compared, and the task is to say whether the recordings are spoken by the same speaker³. An analysis of differences among various framings of the problem can be found in [2], and an argument is given in favor of detection, particularly in technology oriented evaluations. Since the goal of NIST SREs is to drive progress by focusing on the core technology, the evaluations are technology-oriented, and, as a result, the NIST SREs have been focused on spontaneous text-independent speaker detection.

While all the evaluations have had this primary task in common, several evaluations have included one or more alternate tasks. For example, *speaker diarization*, labeling a recording based on who spoke when, has been included in several past evaluations. This might be viewed as a segmentation task

³Those in the machine learning community will recognize detection as a binary classification task.

110 followed by a clustering task, where the recording is segmented into chunks
111 of speech and the segments are clustered by speaker. As another example,
112 in the 2010 and 2012 evaluations, an alternate task involved human-in-the-
113 loop speaker recognition, also known as **human assisted speaker recognition**
114 **(HASR)**. This was a spontaneous text-independent speech speaker detection
115 task, however humans were permitted to listen to the speech and otherwise
116 interact with it in ways forbidden in the traditional SREs [7].

117 2.2. Data

118 “*Data is the new oil.*”[8] “*The data economy is the new economy.*”[9]
119 While data is becoming recognized as increasingly important by society, data
120 has always been the single most critical element of evaluation driven research.
121 **If the data is too easy, the systems will not be challenged and the evaluation**
122 **is of limited value. If the data is too difficult, the systems will balk and**
123 **error analysis will prove mostly fruitless.** If there is not enough data, the
124 results will lack significance. If there is too much data⁴, participants
125 lacking the necessary compute resources will be unable to participate, the
126 logistics of the evaluation will be burdensome, and the analysis can become
127 impractically complex. **Finally, the data must capture the desired conditions**
128 **to support the specific evaluation goals and not be otherwise idiosyncratic in**
129 **some detrimental way.**

130 Past SRE data collection goals have included collection of recordings **in**
131 **different languages, using different microphones with varying distances from**
132 **the speaker, high and low vocal efforts, noisy environments, the utilization**
133 **of different communication networks and technologies, and collections with**
134 **targeted speaker demographics.** Originally, data was collected by offering
135 study participants a handful of free long distance phone calls in exchange for
136 the conversations being recorded. Due to the reduction in cost of making
137 long distance phone calls, this model of data collection has been abandoned,
138 instead favoring paying participants to make phone calls or be interviewed,
139 as well as using “found” data, e.g., recordings from the internet.

140 Since its founding in 1992, the Linguistic Data Consortium (LDC) at the
141 University of Pennsylvania has been the primary collector and provider of
142 data used in the SRE series. Data collections are jointly designed by the

⁴The idea of too much data is in conflict with Bob Mercer’s widely-used comment at Arden House Conference “*There’s no data like more data.*” Like all general truths, there are limits to its application.

143 LDC and NIST, the collections are implemented by the LDC, and the data
144 and annotations are provided to NIST. The collection is then analyzed and
145 processed by NIST prior to splitting the data into appropriate sets for system
146 development and evaluation. Collecting data and finding a split of the data
147 that provides sufficient (but not excessive) amounts for system development
148 while also allowing the necessary data for the evaluation has become increas-
149 ingly difficult. The difficulty lies in the need to collect more data and that
150 the data collected meet some specified properties. That is, precisely measur-
151 ing system performance of better performing systems requires 1) more data
152 to obtain significant results, and 2) data that is more challenging for the sys-
153 tems in useful ways (from a research perspective), which can prove difficult
154 to collect.

155 One of the challenges of transitioning research systems into production
156 environments is that performance “*in the lab*” varies substantially from per-
157 formance “*in the field*.” This has been attributed entirely to differences in the
158 nature of the data in these two contexts. As a result, there has been an in-
159 creasing move toward access to more “realistic” data in technology evaluation
160 settings. In the SRE series, this move has recently involved the collection of
161 telephone recordings not routed through the Philadelphia⁵ public switched
162 telephone network (PSTN), as well as including voice over internet protocol
163 (VOIP) and audio from video (AfV) recordings. As this transition to increas-
164 ingly “real” data progresses, there is a resulting loss of the carefully controlled
165 data collection parameters, simultaneously increasing the importance and
166 challenge of being able to measure various properties of the recordings nec-
167 essary for understanding what aspects of the data are challenging for current
168 systems. The tradeoffs can be even more nuanced. For example, selectively
169 drawing from a real data source in a manner that eases data labeling often
170 results in data that does not have carefully controlled independent variables
171 and still does not sufficiently represent the data source.

172 2.3. Measurement & Analysis

173 Measurement is a foundational requirement of science and engineering.
174 Without the ability to measure, it is not possible to distinguish between
175 change and progress. It is difficult to overstate the fundamental importance
176 of measurement.

⁵The LDC is located in Philadelphia, Pennsylvania, United States.

177 Equally important is what is being measured and how. SREs have always
178 measured system performance using some function of error rate. This seems
179 a bleak and arbitrary choice over focusing on success rate. However, there
180 are advantages to focusing explicitly on errors. When the goal is to improve
181 system performance, focusing on errors is intuitive and naturally leads to
182 areas to direct future effort. It is also worth mentioning that the impact of
183 halving the error rate is more apparent than a relatively small increase in
184 success rate, which will be the case when system performance is well above
185 chance.

186 As mentioned in section 2.1, the task in NIST SREs is detection, and there
187 are two types of errors in detection tasks. Sometimes referred to as type I
188 and type II errors in the statistics and machine learning communities, in the
189 speaker recognition community these errors are often called misses (short for
190 missed detections), false negatives or false rejects (when the speakers are in
191 fact the same) and false alarms, false positives or false accepts (when the
192 speakers are in fact not the same). Each evaluation consists of a series of
193 trials, and a trial consists of one or more recordings of a target speaker for
194 enrollment (or model creation) and a recording of a speaker whose identity is
195 unknown to the system (i.e., may or may not be the target speaker) for testing
196 purposes. Each system submitted to the evaluation must output a real-valued
197 response for every trial, where a greater value indicates greater confidence
198 that the enrollment and test recordings both contain speech spoken by the
199 target speaker.

200 NIST has primarily measured system performance using a *detection cost*
201 *function* (DCF), which is a weighted linear combination of one or more sets of
202 false reject (aka miss) and false alarm rates observed in the evaluation trials,
203 as the main SRE performance metric. Alternate functions over error rates
204 have also been utilized in NIST SREs, including a function sweeping over
205 all observable error rates [10]. Although popular among speaker recognition
206 technology researchers due to its easy interpretability, NIST has typically
207 not been a proponent of using the equal error rate (EER) as an SRE perfor-
208 mance metric because of its inability to weight false alarm and false reject
209 (miss) errors differently. NIST has found that in nearly all contexts, the
210 applications of speaker recognition technology tend to strongly favor either
211 few false alarms or few misses, making the equal error rate an counterpro-
212 ductive choice of operating point to focus attention. Instead, the SREs have
213 focused attention on the low false positive region of the operating range,
214 which is most appropriate for contexts where a high rate of false alarms is

215 problematic [3], such as biometric authentication applications.

216 Simply measuring the performance of multiple systems on a fixed, well-
217 chosen data set using a single, meaningful measurement is inherently valu-
218 able [11, 12]. Doing this regularly allows tracking performance progress over
219 time. Implicit in this process is the need to understand how performance
220 varies under different conditions present in the data, e.g., environmental
221 noise or speaker vocal effort, as this suggests immediate research directions
222 to improve technology performance. Analysis of SRE results have been a
223 driver of researcher efforts as well as many data collections. Past analyses
224 have included differences in speaker environment, vocal effort, speech modal-
225 ity (e.g., reading, interviews, phone conversation among strangers, phone
226 conversations among friends), speaker aging, language, sensor, speaker de-
227 mographics, and channel. NIST has also conducted analysis of the progress
228 of speaker recognition technology over time.

229 As more dimensions of variation are added to the data set, more care-
230 ful analysis is necessary. In order to understand how the co-occurrence of
231 independent variables impact system performance, more data are needed,
232 and data sets must have a sufficient number of trials to support a mean-
233 ingful analysis. Further, once a relationship between an independent variable
234 and performance has been established, a question is raised about what to do
235 when some values of the independent variables have disproportionate repre-
236 sentation in the evaluation data set. Recent SREs have separately measured
237 performance across several such variables and then applied a balanced weight-
238 ing to measure performance, which has also been proposed at various points
239 in the past [13]. This approach has advantages and disadvantages, though
240 the realized impact of this decision on SRE analysis has not been thoroughly
241 explored.

242 An important, if under-recognized, aspect of analysis is how information
243 is displayed. Numbers have relatively little meaning outside their proper con-
244 text. An effective visualization method enables the interpretation process.
245 *Detection Error Tradeoff* (DET) curves, a method that visually depicts the
246 error rates at different operating points on a normal deviate scale, were intro-
247 duced in 1997 by NIST for SRE [14]. A DET curve’s general shape, distance
248 from origin, slope, “steppiness” (or quantization), and relative distance to
249 other DET curves are all meaningful and relatively easy to interpret, mak-
250 ing them popular in speaker recognition as well as various other detection
251 tasks [15, 16].

3. NIST Speaker Recognition Evaluations: A Brief History

The first SRE was held in 1996⁶. Since then, NIST has conducted more than 15 evaluations of speaker recognition technology, including a human assisted speaker recognition evaluation [7], which encouraged participation from human experts and humans collaborating with automatic systems, as well as several online challenges, which distributed embeddings to participants rather than audio recordings to reduce the barrier for participation [17]. Rather than detail each evaluation, we offer a brief summary of the early evaluations and include citations to detailed descriptions for the interested reader.

In the 1996 and 1997 evaluations, the effect of multiple-session training was explored and handset variation was featured as a prominent technical challenge. While handset variation remained a formidable challenge, the 1998 evaluation focused on matched-source training and test data [2].

The 1999 evaluation introduced two new tasks utilizing recordings with multiple speakers: multi-speaker detection, determining which speaker spoke when, and speaker tracking, performing speaker detection as a function of time [18, 19]. The test recordings for both of these tasks consisted of a recording of a telephone call mixed into a single track. The 2000 SRE (SRE00) added a speaker segmentation task, in which no specified target speakers are given and the number of different speakers may or may not be known [20]. SRE00 also included data from the Spanish *AHUMADA* corpus [21], making 2000 the first year SRE made use of non-English data.

In 2001, the SREs began including cellular data and provided automated transcripts produced by a then state-of-the-art automatic speech recognizer as part of an effort to encourage research into idiolectic features⁷. A Federal Bureau of Investigation (FBI) forensic database was included in the 2002 evaluation [23].

In 2004, NIST introduced an unsupervised adaptation mode, where the systems may optionally update the speaker model after each trial involving that model. The 2005 and 2006 evaluations [24] included recordings in

⁶NIST was involved in a limited 1992 speaker identification evaluation for a DARPA program and another small speaker identification evaluation in 1995, though it is difficult to find reference to these events elsewhere in the literature.

⁷This emphasis on higher-level features in speaker recognition was further pursued in a SuperSid workshop following the 2002 SRE [22].

multiple languages spoken by bilingual speakers as well as room microphone recordings, allowing for cross-language and cross-channel trials. This was extended in 2008 [25], by including face-to-face interview data as well. The 2010 SRE (SRE10) [26] explored several new areas, including high and low vocal effort and speaker aging, and featured a new decision cost function metric stressing even lower false positive rates. A human-assisted speaker recognition evaluation was included as part of SRE10 as well. While not part of the SRE series, in 2011 NIST conducted an evaluation of speaker recognition featuring a broad range of test conditions as part of the IARPA BEST program, most notably added noise and reverb. The 2012 SRE (SRE12) [27] explored the performance impact of allowing multiple models to be considered in a given trial by defining model speakers beforehand and distinguishing between “known” and “unknown” test speakers⁸.

4. The Current State of NIST Speaker Recognition Evaluations

The 2016 Speaker Recognition Evaluation (SRE16) was not only the 20th anniversary of the SRE series, but was also the first evaluation to begin introducing a variety of changes that distinguish the current SREs from the past. These changes span all aspects of the evaluation. We highlight several of them in the contexts of evaluation administration, evaluation design, and data collection. We also offer some highlights from the most recent SREs.

4.1. Evaluation Administration

Several early SREs were impacted by delays in data collection, giving a limited amount of time to analyze, process, and organize the data sets prior to distribution⁹. This was seen as detrimental, and NIST decided to not host an SRE in 2014, which would have maintained the then biannual schedule, to allow additional time to collect and organize the data. The series resumed its biannual schedule in 2016 with SRE16.

Early SREs also included a relatively small amount of data with undesirable characteristics, e.g., a trial lacking speech, a mislabeled recording, too little data to support a more fine-grained analysis. Despite their trivial impact on performance measurement, much effort and attention went toward

⁸This turned out to be a major logistical challenge.

⁹In the 2008 SRE, the data collection finished only two weeks before the evaluation began!

314 dealing with these issues at the time, and they proved overly distracting, fill-
315 ing email threads and workshop discussions. To help limit these occurrences,
316 NIST began collaborating with a team at MIT Lincoln Laboratory¹⁰ to detect
317 anomalous data and to gauge expected performance prior to the evaluation.
318 This collaboration has been successful and has had tremendous positive im-
319 pact, especially with respect to reducing data related distractions¹¹.

320 In 2016, NIST developed and began using baseline speaker recognition
321 systems [28] to explicitly test the impact of various evaluation design deci-
322 sions on system performance measurement. The use of NIST developed base-
323 line systems has also improved NIST’s ability to more precisely understand
324 how speaker recognition technology performance has changed over time. Past
325 evaluations have relied on researchers to voluntarily run “mothballed” sys-
326 tems, i.e., systems used in prior evaluations, to help assess how much a change
327 in performance between evaluations is due to system changes and how much
328 is due to the changes in the data. Having a collection of baseline speaker
329 recognition systems, each utilizing the state-of-the-art approach from a past
330 evaluation, has allowed NIST to better quantify the source of changes in
331 performance. Additionally, evaluation participants have reported that the
332 baseline systems’ results have proven useful for debugging their research sys-
333 tems.

334 As a result of the many advances in information technology in recent
335 years, NIST has been able to substantially improve evaluation logistics. In
336 the past, participants needed to register for the evaluation by mail, fax, or
337 email, and then NIST would mail them hard drives and/or optical media
338 containing the evaluation data. Special care would be taken so that the data
339 would be expected to arrive at all participating sites around the world at
340 approximately the same time. The necessary logistics were burdensome and
341 subject to human error. NIST now manages the evaluation logistics through
342 a custom built online web platform¹², that allows sites to register for the
343 evaluation, create formal evaluation teams composed of individual partici-
344 pant sites, sign all necessary documents, download data, upload system out-

¹⁰MIT Lincoln Laboratory also has a team that participates in the evaluations. There is no overlap in staff between these two teams and they do not collaborate on the evaluations.

¹¹As performance improves, the impact of any errors in data labeling or analysis increases, further adding value to the success of this effort.

¹²After first being developed for SRE, the web platform has been used for many different technology evaluations at NIST.

345 put, receive the evaluation results, keys, and analysis, as well as upload and
346 share system descriptions and workshop presentations. This change has had
347 tremendous value for the evaluation participants as well as for NIST, substan-
348 tially reducing the effort needed for, and increasing the speed of completion
349 of, the necessary evaluation administrivia.

350 4.2. Evaluation Design

351 Prior to each evaluation, participants receive data for use in building
352 their speaker recognition systems. It has been the common practice of SRE
353 participants to split the provided data into training and development sets.
354 Current evaluations have specified training and development sets within the
355 provided data. This was in part by popular demand, but it also facilitated the
356 introduction of *fixed* and *open system* training conditions in the evaluation
357 series. The *fixed* training condition limits system training and development
358 to a predetermined common set of corpora to facilitate meaningful system
359 comparisons in terms of core speaker recognition algorithms and/or tech-
360 niques. The *open* training condition allows participants to use any other
361 proprietary and/or publicly available data in addition to the corpora pro-
362 vided in the fixed condition to demonstrate the gains that could be achieved
363 with unconstrained amounts of data. Previously, training data was always
364 unconstrained, though only data that was or would become publicly available
365 was permitted for use.

366 Current SREs have also begun distributing data without speaker labels
367 for use in system development, motivated by the availability of unlabeled
368 data from the data source that can be useful for system adaptation. Typi-
369 cally, researchers have applied a clustering algorithm on this data, intending
370 to cluster recordings based on speaker, and then model the characteristics of
371 the various channels in the data source from the resultant clusters. Interest-
372 ingly, it has been found that a perfect, or oracle, clustering of this data by
373 speaker when using this method does not necessarily lead to optimal speaker
374 recognition performance.

375 An ongoing trend in the SRE series has been the fusion of several speaker
376 recognition systems to create a single “fusion” submission to an evaluation.
377 While it remains interesting to see how much this approach can improve
378 performance, there is a growing sense that the resultant fused systems com-
379 plicate the error analysis and are impractical to deploy. Therefore, current
380 evaluations have encouraged sites to also report results on their best “single”

systems¹³.

4.3. Evaluation Data

The data emphasis in every SRE has always been conversational telephony speech (CTS) recorded over public switched telephone networks (PSTN), though other varieties of speech data have been explored. This emphasis remains in the most recent evaluations, though two new data domains have also been introduced: voice over Internet Protocol (VOIP) and audio from video (AfV). Both the PSTN and VOIP CTS data used for the latest evaluations were extracted from *Call My Net* (CMN) 1 and 2 [29] corpora collected outside of North America, which was a new emphasis for the SREs. On the other hand, the AfV data was extracted from the *Video Annotation for Speech Technologies* (VAST) corpus [30] which was collected from amateur online video blogs (Vlogs) spoken in English, representing more modern data sources.

One factor affecting performance is the amount of speech available to the system. Current SREs explore this variability to a greater extent than in the past. It was previously common to have evaluation recordings either contain approximately 10 seconds of speech or approximately 180 or more seconds of speech for CTS data. Current evaluations now include additional segment durations spanning between 10 and 60 seconds of speech for CTS data, as well as segments potentially containing less or much more speech in the case of AfV data.

Practically speaking, recruiting subjects and collecting speech in a way that is balanced from an experimental design standpoint has always been difficult. This challenge has only grown as the number of data sources and independent variables being explored has increased. One approach is to discard data from any subject that completes only a portion of their intended recordings and then remove other subjects as well to maintain the desired balance. Large amounts of data can be discarded using this approach, so NIST has instead favored accounting for any imbalances during analysis. As mentioned in Section 2.3, current evaluations have also begun re-balancing data as part of computing the performance metric.

¹³While the definition of a “single” system is somewhat subjective, the aim is to encourage more intuitively cohesive and simplified systems versus a score level fusion of a large basket of slightly modified systems.

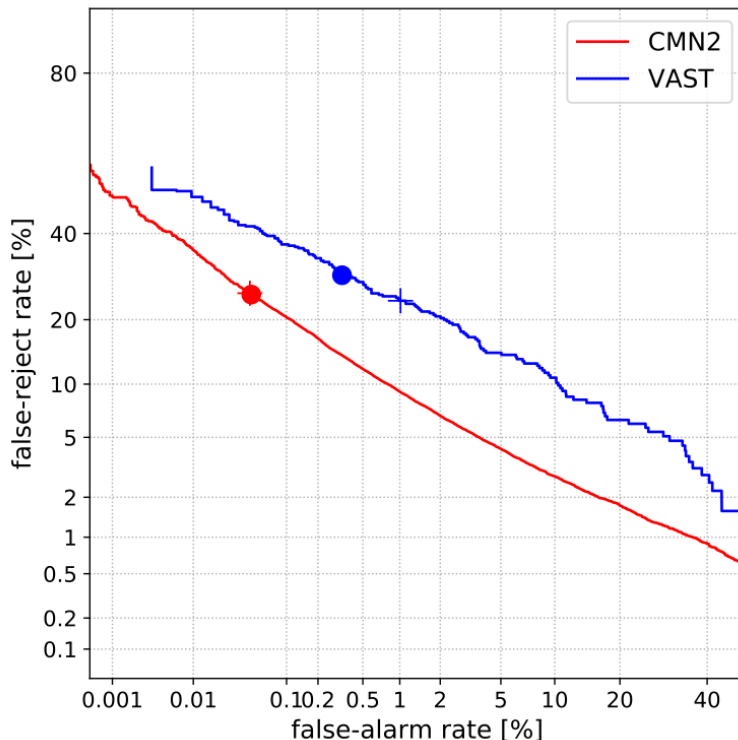


Figure 1: DET curves for a leading system’s performance on CTS data (CMN2) and AfV data (VAST) in SRE18. The circles denote the operating point that minimizes the detection cost function and the cross hairs denote the operating point selected by the system. Systems performed consistently better on CTS data than AfV data in SRE18.

4.4. SRE16 & SRE18 Participation and Performance

The 2018 Speaker Recognition Evaluation (SRE18), held in September of 2018, was the latest in the series of formal NIST evaluations to support research and innovation for text-independent speaker recognition. SRE18 was organized in a manner similar to the 2016 SRE (SRE16), held in September of 2016, and included all of the above mentioned changes.

In SRE18, a total of 48 teams from 78 academic and industrial sites participated. A total of 129 valid system submissions were made, with 120 for the fixed training condition and 9 for the open training condition. The participation in SRE16 was similar, with 66 teams from 34 countries submitting 121 valid submissions (103 for fixed training condition and 18 for open training).

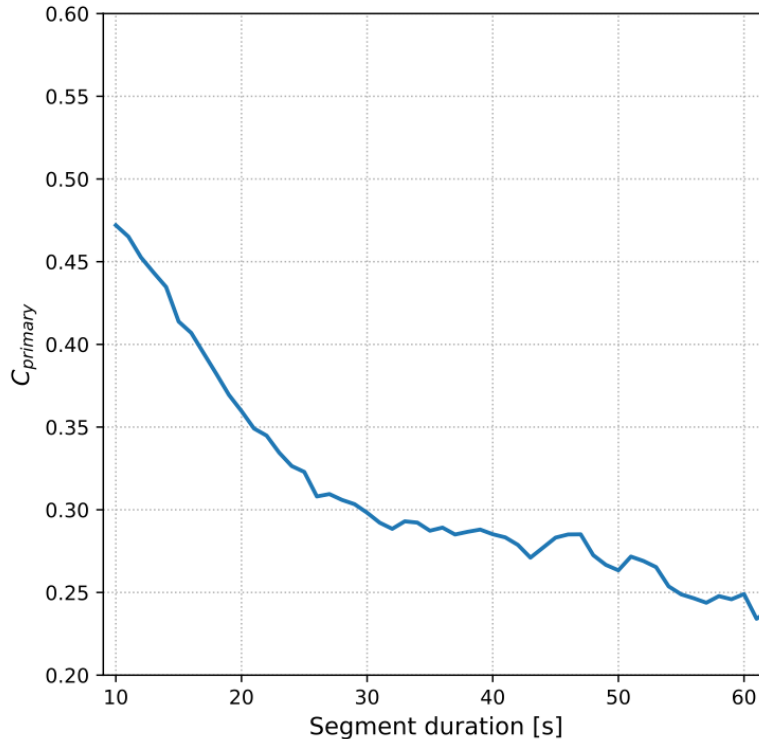


Figure 2: Performance as a function of the speech duration in a test recording for a deep learning based system submission in SRE18. Systems performed consistently better as the speech duration increased as anticipated.

425 These evaluations explored the impact of several factors on system per-
 426 formance, most notably channel/domain (Fig. 1), duration (Fig. 2), and
 427 language (Fig. 3). They also found that the effective use of the provided
 428 unlabeled development data and choice of calibration data substantially im-
 429 pacted system performance, particularly for the data from the AfV domain.
 430 Approaches based on recent advances in neural networks, found to be less suc-
 431 cessful in SRE16¹⁴, were dominant in SRE18 due to the availability of large
 432 amounts of training data from a large number of speakers, the use of data
 433 augmentation in system development, and the use of more complex models.

¹⁴This is believed to be due to the language and domain mismatch presented in the 2016 evaluation.

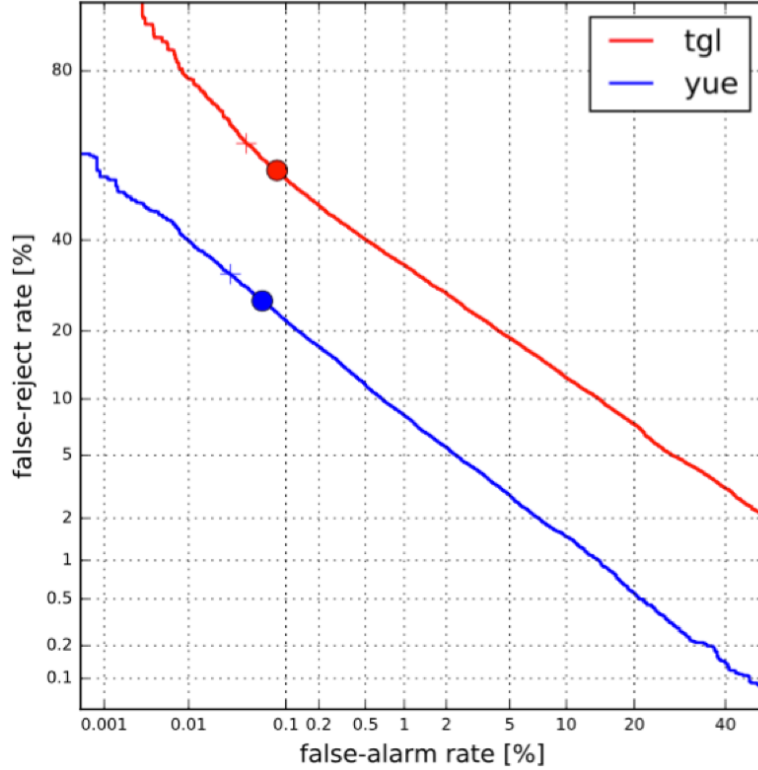


Figure 3: DET curves for a leading system’s performance on Tagalog speech (tgl) and Cantonese speech (yue) in SRE16. The circles denote the operating point that minimizes the detection cost function and the cross hairs denote the operating point selected by the system. System performances were consistently better on Cantonese speech than Tagalog speech, though there were channel differences between the Tagalog and Cantonese recordings that may have led to the observed performance differences.

434 While fusion systems continued to maintain some of the performance advan-
 435 tages seen in SRE16, SRE18 witnessed strong single system results that were
 436 nearly as good as the best fused systems (Fig. 4). We include a figure com-
 437 paring SRE16 systems with SRE18 systems (Fig. 5). The interested reader
 438 can find additional results for SRE16 and SRE18 in [28] and [31] respectively.

439 4.5. SRE19 & SRE20

440 Plans for the 2019 (SRE19) and 2020 (SRE20) Speaker Recognition Eval-
 441 uations were publicized at the SRE18 participant workshop in December
 442 2018. Acknowledging the observed performance challenges presented by the

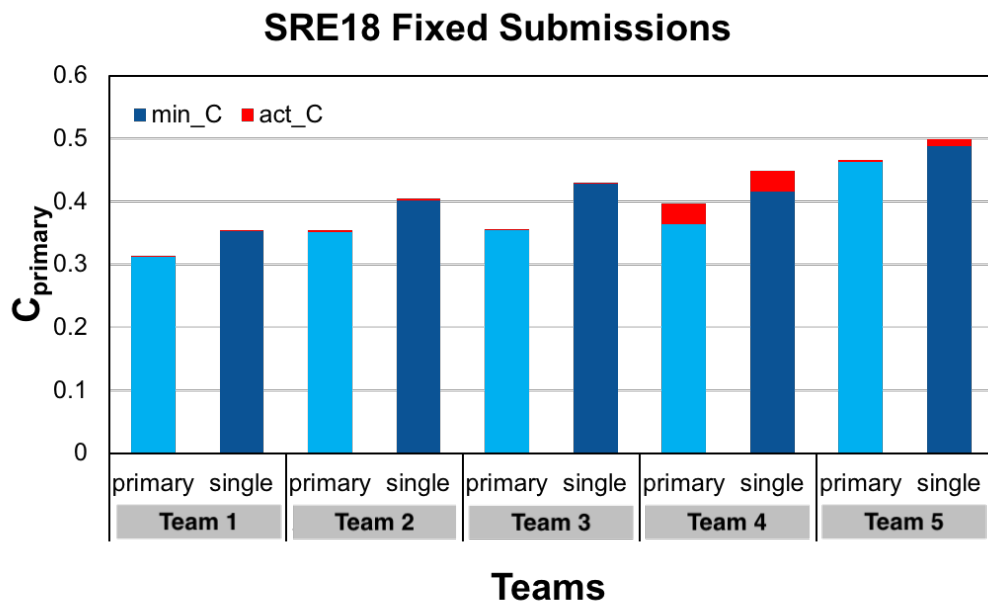


Figure 4: A comparison of system performance for fused (primary) and the best single systems from five teams in SRE18. The detection cost is displayed at both the minimum operating (min_C) and the actual operation point (act_C). The observed differences between the fused system and single systems within teams is relatively small. Further, the best single system in the evaluation was competitive with the best fused systems in the evaluation.

443 AfV data in SRE18 and the growing interest of the speaker recognition re-
 444 search community in applying speaker recognition to more realistic multime-
 445 dia applications, both SRE19 and SRE20 have the goal of further exploring
 446 speaker recognition technology for audio from amateur video data. In addi-
 447 tion to exploiting the audio from video data, these evaluations will provide
 448 participants the opportunity to explore the possibility of fusing face recogni-
 449 tion with speaker recognition.

450 SRE19 will serve as a special evaluation allowing more in depth analysis
 451 and exploration into each of the data domains used in SRE18. There will be
 452 two components to SRE19: the SRE19 CTS Challenge and the SRE19 Audio-
 453 visual (AV) evaluation. The SRE19 CTS challenge will be conducted entirely
 454 online in a manner similar to the NIST 2014 and 2015 i-vector challenges [17,
 455 32], however actual audio recordings will be used as the source data instead

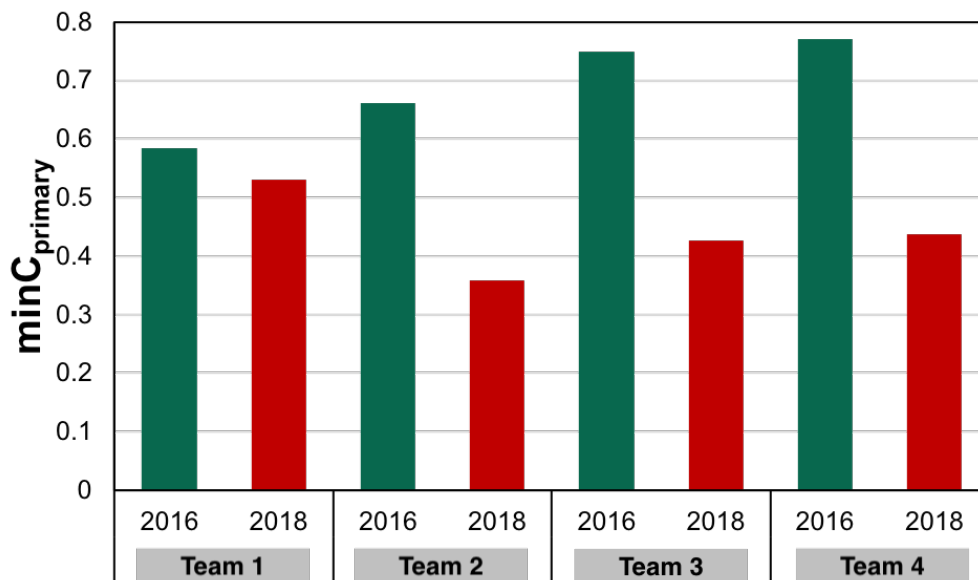


Figure 5: A comparison of system performance for the SRE16 and SRE18 systems submitted by four teams that participated in both evaluations. A data set drawn from the *Call My Net* corpus [29] was used to measure the performance of these 2016 and 2018 systems. Substantial improvements can be seen between the systems submitted in 2016 and those submitted in 2018.

456 of feature embeddings. Unexposed CTS data from the CMN2 corpus will
 457 be used to support the SRE19 CTS challenge. System performance scores
 458 will be made available throughout the entire evaluation period instead of
 459 at the end, and multiple submissions will be allowed, enabling participants
 460 to explore how low they can drive error rates on the traditional CTS data
 461 domain.

462 The SRE19 AV evaluation will be conducted in the same manner as the
 463 traditional SREs, with training and development data released in early sum-
 464 mer 2019, evaluation data released in late summer 2019, evaluation results
 465 submitted in October 2019, and a post-evaluation workshop held in December
 466 2019¹⁵. Unexposed multimedia data from the VAST corpus will be used to
 467 support the SRE19 evaluation which will feature two core evaluation tracks:
 468 audio only and audio+visual fusion. An optional visual only track will also

¹⁵The SRE19 workshop will be co-located with the 2019 IEEE Automatic Speech Recognition and Understanding (ASRU) Workshop in Sentosa, Singapore.

469 be available for participants.

470 The plans for SRE20 are based on the availability of a data corpus cur-
471 rently being collected by the LDC from multilingual speakers in both the
472 CTS and AfV data domains. This corpus is designed to allow for explo-
473 rations into cross-domain enroll-test trials (e.g. enroll on CTS data and test
474 on AfV data for a single target speaker). The corpus is also designed to pro-
475 vide image data to support multimodal fusion explorations similar to SRE19.
476 Continuing with the SRE16 and SRE18 data paradigms, this corpus is being
477 collected outside of North America and will feature non-English data.

478 5. The Future of NIST Speaker Recognition Evaluations

479 Pending the availability of sufficient and appropriate data, it is expected
480 that the NIST SREs will continue after 2020 and resume a bi-annual schedule
481 in 2022 with a focus on challenging data domains and channels. NIST will
482 also continue to explore ways to collaborate with organizers of other speaker
483 recognition technology evaluations, where feasible, to ensure maximal com-
484 munity benefit. As the SRE series moves into its next decade, we highlight
485 some of the projected trends for the future in the contexts of evaluation tasks
486 and evaluation data.

487 5.1. Evaluation Tasks

488 The one constant throughout the SRE series from its inception has been
489 a focus on speaker detection for spontaneous text-independent speech. The
490 consistency of this task has allowed NIST to drive core speaker recognition
491 technology forward and track the technological advancements over the last
492 two decades. Moving into the next decade of speaker recognition evaluation,
493 NIST maintains the same goal of driving speaker recognition progress by
494 focusing on the core technology and anticipates maintaining a core focus on
495 spontaneous text-independent speaker detection.

496 Continuing with the core speaker detection task will also allow NIST to
497 have a continued focus on the technological challenges presented by data do-
498 main and channel mismatches as new domains/channels become of interest
499 to the speaker recognition research community. And as multimedia appli-
500 cations become more relevant to the speaker recognition community, like
501 realtime group discussion transcription applications that use visual data to
502 help with speaker identification, tasks involving the fusion of audio and video

503 data such as those introduced in SRE19 are also anticipated to continue to
504 be considered in future SREs.

505 5.2. *Evaluation Data*

506 While the core SRE task will remain the same moving into the future,
507 the data used to evaluate that task will continue to evolve in order to sup-
508 port exploration in more challenging domains and channels. Conversational
509 telephony speech (CTS) data will remain a focus of the SRE series moving
510 forward, and NIST maintains the goal of including recordings from different
511 languages, from microphones with varying distances from the speaker, and
512 different communication networks and technologies. It is anticipated that
513 NIST will continue to partner with LDC to collect data for future evalu-
514 ations. The collaboration has provided NIST with the largest amount of
515 control over desired data collection parameters and data properties, which
516 will become more important as more challenging data properties are intro-
517 duced to the SRE series.

518 In addition to evolving CTS data characteristics, a continued progression
519 towards data that mimics more realistic modern application conditions is also
520 a possible focus area for future SREs (e.g., multimedia data, virtual assistant
521 enabled devices, etc.). Recent SREs have leveraged publicly available speaker
522 recognition data sources using “found data”,¹⁶ and this trend may continue
523 in the future as long as these sources remain available for public research use.

524 6. Acknowledgements

525 A multitude of people worldwide have contributed to the success of the
526 NIST SREs over the last two decades as sponsors, evaluation designers, and
527 evaluation participants. While it is not feasible to list all their names, the
528 authors would like to acknowledge their contributions covered in this paper.

529 7. Disclaimer

530 The results presented in this paper are not to be construed or represented
531 as endorsements of any participants system, methods, or commercial product,
532 or as official findings on the part of NIST or the U.S. Government.

¹⁶VoxCeleb [33, 34] and SITW [35] corpora were allowable under the SRE18 fixed train-
ing condition.

References

- [1] B. J. Dorr, P. C. Fontana, C. S. Greenberg, M. Le Bras, M. Przybocki, Evaluation-driven research in data science: Leveraging cross-field methodologies, in: Proc. IEEE International Conference on Big Data (Big Data), 2016, pp. 2853–2862.
- [2] G. R. Doddington, M. A. Przybocki, A. F. Martin, D. A. Reynolds, The NIST speaker recognition evaluation—Overview, methodology, systems, results, perspective, *Speech Communication* 31 (2000) 225–254.
- [3] M. A. Przybocki, A. F. Martin, NIST Speaker Recognition Evaluation Chronicles, in: Proc. Odyssey 2004: The Speaker and Language Recognition Workshop, pp. 15–22.
- [4] A. Martin, M. Przybocki, J. P. Campbell, The NIST speaker recognition evaluation program, in: *Biometric Systems*, Springer, 2005, pp. 241–262.
- [5] M. A. Przybocki, A. F. Martin, A. N. Le, NIST Speaker Recognition Evaluation Chronicles-Part 2, in: Proc. Odyssey 2006: The Speaker and Language Recognition Workshop, pp. 1–6.
- [6] A. Larcher, K. A. Lee, B. Ma, H. Li, Text-dependent speaker verification: Classifiers, databases and RSR2015, *Speech Communication* 60 (2014) 56–77.
- [7] C. S. Greenberg, A. F. Martin, L. Brandschain, J. P. Campbell, C. Cieri, G. R. Doddington, J. J. Godfrey, Human Assisted Speaker Recognition in NIST SRE10, in: Proc. Odyssey 2010: The Speaker and Language Recognition Workshop, pp. 180–185.
- [8] The Economist, The world’s most valuable resource is no longer oil, but data, <https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data>, 2017. Accessed: 2019-05-05.
- [9] The Economist, Data is giving rise to a new economy, <https://www.economist.com/briefing/2017/05/06/data-is-giving-rise-to-a-new-economy>, 2017. Accessed: 2019-05-05.
- [10] N. Brümmer, J. Du Preez, Application-independent evaluation of speaker detection, *Computer Speech & Language* 20 (2006) 230–275.

- 565 [11] G. R. Doddington, T. B. Schalk, Speech recognition: Turning theory to
566 practice: New ICs have brought the requisite computer power to speech
567 technology; an evaluation of equipment shows where it stands today,
568 IEEE Spectrum 18 (1981) 26–32.
- 569 [12] D. S. Pallett, A look at NIST’s benchmark ASR tests: past, present,
570 and future, in: Proc. IEEE ASRU Workshop 2003, pp. 483–488.
- 571 [13] D. A. v. Leeuwen, Overall performance metrics for multi-condition
572 Speaker Recognition Evaluations, in: Proc. INTERSPEECH 2009, pp.
573 908–911.
- 574 [14] A. Martin, G. Doddington, T. Kamm, M. Ordowski, M. Przybocki,
575 The DET curve in assessment of detection task performance, Technical
576 Report, National Institute of Standards and Technology, Gaithersburg,
577 MD, 1997.
- 578 [15] W. B. Croft, J. Lafferty, Language modeling for information retrieval,
579 volume 13, Springer Science & Business Media, 2013.
- 580 [16] T. Rose, J. Fiscus, P. Over, J. Garofolo, M. Michel, The TRECVID 2008
581 event detection evaluation, in: Proc. IEEE Workshop on Applications
582 of Computer Vision (WACV), 2009, pp. 1–8.
- 583 [17] C. S. Greenberg, D. Bansé, G. R. Doddington, D. Garcia-Romero, J. J.
584 Godfrey, T. Kinnunen, A. F. Martin, A. McCree, M. Przybocki, D. A.
585 Reynolds, The NIST 2014 Speaker Recognition i-vector Machine Learn-
586 ing Challenge, in: Proc. Odyssey 2014: The Speaker and Language
587 Recognition Workshop, pp. 224–230.
- 588 [18] M. A. Przybocki, A. F. Martin, The 1999 NIST speaker recognition eval-
589 uation, using summed two-channel telephone data for speaker detection
590 and speaker tracking, in: Proc. EUROSPEECH 1999, pp. 2215–2218.
- 591 [19] A. Martin, M. Przybocki, The NIST 1999 Speaker Recognition Evalu-
592 ationan overview, Digital Signal Processing 10 (2000) 1–18.
- 593 [20] A. F. Martin, M. A. Przybocki, The NIST Speaker Recognition Evalu-
594 ations: 1996-2001, in: 2001: A Speaker Odyssey-The Speaker Recogni-
595 tion Workshop, pp. 39–43.

- 596 [21] J. Ortega-Garcia, J. Gonzalez-Rodriguez, V. Marrero-Aguiar, AHU-
597 MADA: A large speech corpus in Spanish for speaker characterization
598 and identification, *Speech Communication* 31 (2000) 255–264.
- 599 [22] D. Reynolds, W. Andrews, J. Campbell, J. Navratil, B. Peskin,
600 A. Adami, Q. Jin, D. Klusacek, J. Abramson, R. Mihaescu, et al., The
601 SuperSID project: Exploiting high-level information for high-accuracy
602 speaker recognition, in: *Proc. IEEE ICASSP 2003*, volume 4, pp. IV–
603 784.
- 604 [23] H. Nakasone, S. Beck, Forensic Automatic Speaker Recognition, in:
605 *Proc. Odyssey 2001: The Speaker and Language Recognition Workshop*.
- 606 [24] M. Przybocki, A. Martin, A. Le, NIST Speaker Recognition Evaluations
607 Utilizing the Mixer corpora – 2004, 2005, 2006, *IEEE Transactions on*
608 *Audio, Speech, and Language Processing* 15 (2007) 1951–1959.
- 609 [25] A. F. Martin, C. S. Greenberg, NIST 2008 Speaker Recognition Eval-
610 uation: Performance across telephone and room microphone channels,
611 in: *Proc. INTERSPEECH 2009*, pp. 2579–2582.
- 612 [26] A. F. Martin, C. S. Greenberg, The NIST 2010 Speaker Recognition
613 Evaluation, in: *Proc. INTERSPEECH 2010*, pp. 2726–2729.
- 614 [27] C. S. Greenberg, V. M. Stanford, A. F. Martin, M. Yadagiri, G. R. Dod-
615 dington, J. J. Godfrey, J. Hernandez-Cordero, The 2012 NIST Speaker
616 Recognition Evaluation., in: *Proc. INTERSPEECH 2013*, pp. 1971–
617 1975.
- 618 [28] S. O. Sadjadi, T. Kheyrkhah, A. Tong, C. S. Greenberg, D. A. Reynolds,
619 The 2016 NIST Speaker Recognition Evaluation., in: *Proc. INTER-*
620 *SPEECH 2017*, pp. 1353–1357.
- 621 [29] K. Jones, S. M. Strassel, K. Walker, D. Graff, J. Wright, Call My Net
622 corpus: A Multilingual Corpus for Evaluation of Speaker Recognition
623 Technology, in: *Proc. INTERSPEECH 2017*, pp. 2621–2624.
- 624 [30] J. Tracey, S. Strassel, VAST: A corpus of video annotation for speech
625 technologies, in: *Proc. LREC 2018*, pp. 4318–4321.

- 626 [31] S. O. Sadjadi, C. S. Greenberg, D. A. Reynolds, L. Mason, The 2018
627 NIST Speaker Recognition Evaluation, in: Proc. INTERSPEECH 2019.
- 628 [32] A. Tong, C. Greenberg, A. Martin, D. Banse, J. Howard, H. Zhao,
629 G. Doddington, D. Garcia-Romero, A. McCree, D. Reynolds, E. Singer,
630 J. Hernandez-Cordero, L. Mason, Summary of the 2015 NIST Language
631 Recognition i-vector Machine Learning Challenge, in: Proc. Odyssey
632 2016: The Speaker and Language Recognition Workshop, pp. 297–302.
- 633 [33] A. Nagrani, J. S. Chung, A. Zisserman, VoxCeleb: a large-scale speaker
634 identification dataset, in: Proc. INTERSPEECH 2017, pp. 2616–2620.
- 635 [34] J. S. Chung, A. Nagrani, A. Zisserman, VoxCeleb2: Deep Speaker
636 Recognition, in: Proc. INTERSPEECH 2018, pp. 1086–1090.
- 637 [35] M. McLaren, L. Ferrer, D. Castan, A. Lawson, The Speakers in the
638 Wild (SITW) Speaker Recognition Database, in: Proc. INTERSPEECH
639 2016, pp. 812–822.