

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

ARTICLE INFO

Article History:

Received 6 May 2019

Revised 17 September 2019

Accepted 17 October 2019

Available online 29 October 2019

Keywords:

NIST SRE

Speaker recognition

Speaker Recognition Evaluation

Speaker verification

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.

© 2019 Published by Elsevier Ltd.

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 methods in order to reduce the total overhead necessary to conduct research and to benchmark the current state of the art and identify the most promising research directions (Dorr et al., 2016). There are four basic components that make up *evaluation-driven research*: planning, design, assessment, and a workshop. The planning component involves identifying research goals for the technology (e.g., to be able to improve performance of the fundamental underlying technology or to be robust to certain conditions), obtaining data that supports the evaluation goals, creating and documenting the evaluation plan, as well as identifying and notifying interested researchers and organizations. The design component involves deciding the tasks, metrics, and measurement methods that will make up the evaluation, and analyzing the available data to create necessary data sets (e.g., typically some data is provided to researchers in advance of the assessment period to assist in research, and other data is used as test data for the assessment). During the assessment component, either the technology developers or the

*Corresponding author.

E-mail address: craig.greenberg@nist.gov (C.S. Greenberg).

evaluator runs the systems with the specified test data, and the evaluator analyzes the systems' performances. At the workshop, results and lessons learned are shared and future research goals are identified, which support the planning of future evaluations.

In 1996, NIST conducted its first evaluation of technology for automatically recognizing speakers by their voices. Over the following two decades¹, NIST conducted 15 Speaker Recognition Evaluations (SREs), in addition to an evaluation held in 2018 and evaluations planned for 2019 and 2020. During that time, speaker recognition technology has evolved substantially, and the SRE series has as well. What started as an evaluation of approximately 10 systems completing 4000 trials has expanded into a series that commonly includes hundreds of systems completing millions of trials. This has been necessary, as the 1996 evaluation would be grossly insufficient for the research needs in 2019, and the 2019 evaluation would have been impossible in 1996—specifically, the 1996 SRE data set is too small and the data too easy to analyze performance of modern state-of-the-art systems, and the amount of data and challenging data conditions planned for SRE19 would have overwhelmed state of the art systems in 1996.

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

In this article, we present an overview of the NIST ITR/IAD/MIG approach to evaluating speaker recognition technology over the past two decades and provide insights into what evaluations may look like moving into the next decade. The aim is to provide a review of the *evaluation-driven research* methodology employed by the SRE series that is accessible by newcomers to the field of Speaker Recognition Evaluation. We discuss some of the key considerations necessary when conducting speaker recognition technology evaluations, and how NIST has addressed evaluating speaker recognition in general and for specific, specialized tasks. A brief survey of past SREs and results from recent evaluations are also provided, as well as a brief overview of plans for the 2019 and 2020 evaluations. We conclude the article with some general projections about how future evaluations may look as research directions have dramatically evolved since the inaugural 1996 SRE.

2. Considerations in evaluating speaker recognition technology

There is a great deal that could be said about the considerations necessary when running large-scale research-focused evaluations of speaker recognition technology.² Indeed, NIST has published several lengthy articles covering various aspects of this topic (Doddington et al., 2000; Przybocki and Martin, 2004; Martin et al., 2005; Przybocki et al., 2006). While still more could be said 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 performance,³ 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 (Larcher et al., 2014).

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 (Doddington et al., 2000). 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 speaker⁴. An analysis of differences among various framings of the problem can be found in Doddington et al. (2000), 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 followed by a clustering task, where the recording is segmented into chunks

¹ Over 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*.

² During 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."

³ <https://www.nytimes.com/2003/11/13/technology/what-s-next-is-that-you-son-voice-authentication-trips-up-the-experts.html>.

⁴ Those in the machine learning community will recognize detection as a binary classification task.

of speech and the segments are clustered by speaker. As another example, in the 2010 and 2012 evaluations, an alternate task involved human-in-the-loop speaker recognition, also known as human assisted speaker recognition (HASR). This was a spontaneous text-independent speaker detection task, however humans were permitted to listen to the speech and otherwise interact with it in ways forbidden in the traditional SREs (Greenberg et al., 2010).

2.2. Data

“Data is the new oil.” (The Economist 2017b). “The data economy is the new economy.” (The Economist 2017a). While data is becoming recognized as increasingly important by society, it has always been the single most critical element of evaluation driven research. If the data is too easy, the systems will not be challenged and the evaluation is of limited value. If the data is too difficult, the systems will balk and error analysis will prove mostly fruitless. If there is not enough data, the results will lack significance. If there is too much data⁵, participants lacking the necessary compute resources will be unable to participate, the logistics of the evaluation will be burdensome, and the analysis can become impractically complex. Finally, the data must capture the desired conditions to support the specific evaluation goals and not be otherwise idiosyncratic in some detrimental way.

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

Since its founding in 1992, the Linguistic Data Consortium (LDC) at the University of Pennsylvania has been the primary collector and provider of data used in the SRE series. Data collections are jointly designed by the LDC and NIST, the collections are implemented by the LDC, and the data and annotations are provided to NIST. The collection is then analyzed and processed by NIST prior to splitting the data into appropriate sets for system development and evaluation. Collecting data and finding a split of the data that provides sufficient (but not excessive) amounts for system development while also allowing the necessary data for the evaluation has become increasingly difficult. The difficulty lies in the need to collect more data and that the data collected meet some specified properties. That is, precisely measuring system performance of better performing systems requires (1) more data to obtain significant results, and (2) data that is more challenging for the systems in useful ways (from a research perspective), which can prove difficult to collect.

One of the challenges of transitioning research systems into production environments is that performance “in the lab” varies substantially from performance “in the field.” This has been attributed entirely to differences in the nature of the data in these two contexts. As a result, there has been an increasing move toward access to more “realistic” data in technology evaluation settings. In the SRE series, this move has recently involved the collection of telephone recordings not routed through the Philadelphia⁶ public switched telephone network (PSTN), as well as including voice over internet protocol (VOIP) and audio from video (AFV) recordings. As this transition to increasingly “real” data progresses, there is a resulting loss of the carefully controlled data collection parameters, simultaneously increasing the importance and challenge of being able to measure various properties of the recordings necessary for understanding what aspects of the data are challenging for current systems. The tradeoffs can be even more nuanced. For example, selectively drawing from a real data source in a manner that eases data labeling often results in data that does not have carefully controlled independent variables and still does not sufficiently represent the data source.

2.3. Measurement and analysis

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

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

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

⁵ The 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.

⁶ The LDC is located in Philadelphia, Pennsylvania, United States.

NIST has primarily measured system performance using a *detection cost function* (DCF), which is a weighted linear combination of one or more sets of false reject (aka miss) and false alarm rates observed in the evaluation trials, as the main SRE performance metric. Alternate functions over error rates have also been utilized in NIST SREs, including a function sweeping over all observable error rates (Brümmer and Du Preez, 2006). Although popular among speaker recognition technology researchers due to its easy interpretability, NIST has typically not been a proponent of using the equal error rate (EER) as an SRE performance metric because of its inability to weight false alarm and false reject (miss) errors differently. NIST has found that in nearly all contexts, the applications of speaker recognition technology tend to strongly favor either few false alarms or few misses, making the equal error rate a counterproductive choice of operating point to focus attention. Instead, the SREs have focused attention on the low false positive region of the operating range, which is most appropriate for contexts where a high rate of false alarms is problematic (Przybocki and Martin, 2004), such as biometric authentication applications.

Simply measuring the performance of multiple systems on a fixed, well-chosen data set using a single, meaningful measurement is inherently valuable (Doddington and Schalk, 1981; Pallett, 2003). Doing this regularly allows tracking performance progress over time. Implicit in this process is the need to understand how performance varies under different conditions present in the data, e.g., environmental noise or speaker vocal effort, as this suggests immediate research directions to improve technology performance. Analysis of SRE results have been a driver of researcher efforts as well as many data collections. Past analyses have included differences in speaker environment, vocal effort, speech modality (e.g., reading, interviews, phone conversation among strangers, phone conversations among friends), speaker aging, language, sensor, speaker demographics, and channel. NIST has also conducted analysis of the progress of speaker recognition technology over time.

As more dimensions of variation are added to the data set, more careful analysis is necessary. In order to understand how the co-occurrence of independent variables impact system performance, more data are needed, and data sets must have a sufficient number of trials to support a meaningful analysis. Further, once a relationship between an independent variable and performance has been established, a question is raised about what to do when some values of the independent variables have disproportionate representation in the evaluation data set. Recent SREs have separately measured performance across several such variables and then applied a balanced weighting to measure performance, which has also been proposed at various points in the past (Leeuwen, 2009). This approach has advantages and disadvantages, though the realized impact of this decision on SRE analysis has not been thoroughly explored.

An important, if under-recognized, aspect of analysis is how information is displayed. Numbers have relatively little meaning outside their proper context. An effective visualization method enables the interpretation process. *Detection Error Tradeoff* (DET) curves, a method that visually depicts the error rates at different operating points on a normal deviate scale, were introduced in 1997 by NIST for SRE (Martin et al., 1997). A DET curve's general shape, distance from origin, slope, "steppiness" (or quantization), and relative distance to other DET curves are all meaningful and relatively easy to interpret, making them popular in speaker recognition as well as various other detection tasks (Croft and Lafferty, 2013; Rose et al., 2009).

3. NIST Speaker Recognition Evaluations: a brief history

The first SRE was held in 1996⁷. Since then, NIST has conducted more than 15 evaluations of speaker recognition technology, including a Human Assisted Speaker Recognition Evaluation (Greenberg et al., 2010), which encouraged participation from human experts and humans collaborating with automatic systems, as well as several online challenges, which distributed speaker representations (a.k.a. embeddings) to participants rather than audio recordings to reduce the barrier for participation (Greenberg et al., 2014). 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 (Doddington et al., 2000).

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 (Przybocki and Martin, 1999; Martin and Przybocki, 2000). 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 (Martin and Przybocki, 2001). SRE00 also included data from the Spanish AHUMADA corpus (Ortega-Garcia et al., 2000), making 2000 the first year that 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 ideolectic features⁸. A Federal Bureau of Investigation (FBI) forensic database was included in the 2002 evaluation (Nakasone and Beck, 2001).

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 (Przybocki et al., 2007) included recordings in multiple languages spoken by bilingual speakers as well as room microphone recordings, allowing for cross-language and cross-channel

⁷ NIST 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.

⁸ This emphasis on higher-level features in speaker recognition was further pursued in a SuperSid workshop following the 2002 SRE (Reynolds et al., 2003).

trials. This was extended in 2008 (Martin and Greenberg 2009), by including face-to-face interview data as well. The 2010 SRE (SRE10) (Martin and Greenberg 2010) 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 room reverberation. The 2012 SRE (SRE12) (Greenberg et al. 2013) 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 speakers⁹.

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 distribution¹⁰. 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 dealing with these issues at the time, and they proved overly distracting, filling email threads and workshop discussions. To help limit these occurrences, NIST began collaborating with a team at MIT Lincoln Laboratory¹¹ to detect anomalous data and to gauge expected performance prior to the evaluation. This collaboration has been successful and has had tremendous positive impact, especially with respect to reducing data related distractions¹².

In 2016, NIST developed and began using baseline speaker recognition systems (Sadjadi et al., 2017) to explicitly test the impact of various evaluation design decisions on system performance measurement. The use of NIST developed baseline systems has also improved NIST’s ability to more precisely understand how speaker recognition technology performance has changed over time. Past evaluations have relied on researchers to voluntarily run “mothballed” systems, i.e., systems used in prior evaluations, to help assess how much a change in performance between evaluations is due to system changes and how much is due to the changes in the data. Having a collection of baseline speaker recognition systems, each utilizing the state-of-the-art approach from a past evaluation, has allowed NIST to better quantify the source of changes in performance. Additionally, evaluation participants have reported that the baseline systems’ results have proven useful for debugging their research systems.

As a result of the many advances in information technology in recent years, NIST has been able to substantially improve evaluation logistics. In the past, participants needed to register for the evaluation by mail, fax, or email, and then NIST would mail them hard drives and/or optical media containing the evaluation data. Special care would be taken so that the data would be expected to arrive at all participating sites around the world at approximately the same time. The necessary logistics were burdensome and subject to human error. NIST now manages the evaluation logistics through a custom built online web platform¹³, that allows sites to register for the evaluation, create formal evaluation teams composed of individual participant sites, sign all necessary documents, download data, upload system output, receive the evaluation results, keys, and analysis, as well as upload and share system descriptions and workshop presentations. This change has had tremendous value for the evaluation participants as well as for NIST, substantially reducing the effort needed for, and increasing the speed of completion of, the necessary evaluation administration.

4.2. Evaluation design

Prior to each evaluation, participants receive data for use in building their speaker recognition systems. It has been the common practice of SRE participants to split the provided data into training and development sets. Current evaluations have specified training and development sets within the provided data. This was in part by popular demand, but it also facilitated the introduction of *fixed* and *open* system training conditions in the evaluation series. The *fixed* training condition limits system training and development to a predetermined common set of corpora to facilitate meaningful system comparisons in terms of core speaker recognition algorithms and/or techniques. The *open* training condition allows participants to use any other proprietary and/or

⁹ This turned out to be a major logistical challenge.

¹⁰ In the 2008 SRE, the data collection finished only two weeks before the evaluation began!

¹¹ MIT 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.

¹² As performance improves, the impact of any errors in data labeling or analysis increases, further adding value to the success of this effort.

¹³ After first being developed for SRE, the web platform has been used for many different technology evaluations at NIST.

publicly available data in addition to the corpora provided in the fixed condition to demonstrate the gains that could be achieved with unconstrained amounts of data. Previously, training data was always unconstrained, though only data that was or would become publicly available was permitted for use.

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

An ongoing trend in the SRE series has been the fusion of several speaker recognition systems to create a single “fusion” submission to an evaluation. While it remains interesting to see how much this approach can improve performance, there is a growing sense that the resultant fused systems complicate the error analysis and are impractical to deploy. Therefore, current evaluations have encouraged sites to also report results on their best “single” systems¹⁴.

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 (Jones et al., 2017) 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 (Tracey and Strassel, 2018) 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 s 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 s 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, therefore 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.

4.4. SRE16 and 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).

These evaluations explored the impact of several factors on system performance, most notably channel/domain (Fig. 1), duration (Fig. 2), and language (Fig. 3). They also found that the effective use of the provided unlabeled development data and choice of calibration data substantially impacted system performance, particularly for the data from the AfV domain. Approaches based on recent advances in neural networks, found to be less successful in SRE16¹⁵, were dominant in SRE18 due to the availability of large amounts of training data from a large number of speakers, the use of data augmentation in system development, the introduction of a new end-to-end speaker embedding extraction paradigm, and the use of more complex models. While fusion systems continued to maintain some of the performance advantages seen in SRE16, SRE18 witnessed strong single system results that were nearly as good as the best fused systems (Fig. 4). We include a figure comparing SRE16 systems with SRE18 systems (Fig. 5). The interested reader can find additional results for SRE16 and SRE18 in Sadjadi et al. (2017) and Sadjadi et al. (2019) respectively.

4.5. SRE19 and SRE20

Plans for the 2019 (SRE19) and 2020 (SRE20) Speaker Recognition Evaluations were publicized at the SRE18 participant workshop in December 2018. Acknowledging the observed performance challenges presented by the AfV data in SRE18 and the growing interest of the speaker recognition research community in applying speaker recognition to more realistic multimedia

¹⁴ While 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.

¹⁵ This is believed to be due to the language and domain mismatch presented in the 2016 evaluation.

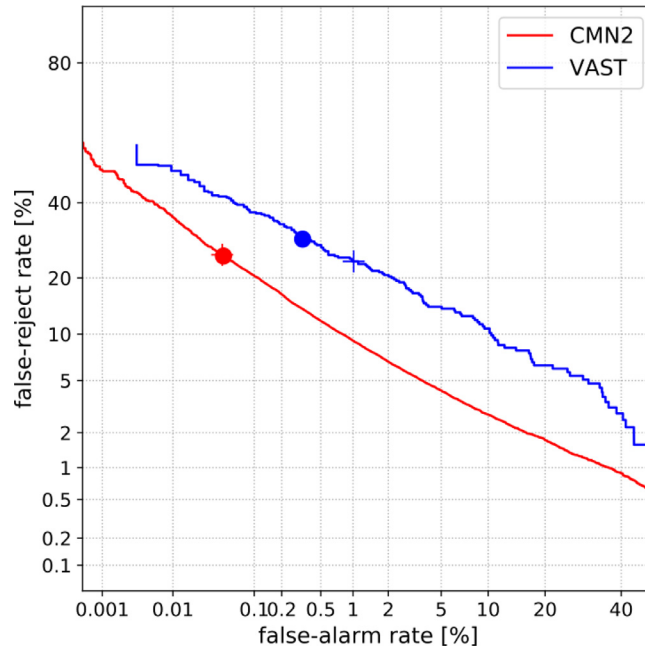


Fig. 1. DET curves for a leading system's performance on CTS data (CMN2) and AfV data (VAST) in SRE18. The circles denote the operating point that minimizes the detection cost function and the cross hairs denote the operating point selected by the system. Systems performed consistently better on CTS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

applications, both SRE19 and SRE20 have the goal of further exploring speaker recognition technology for audio from amateur video data. In addition to exploiting the audio from video data, these evaluations will provide participants the opportunity to explore the possibility of fusing face recognition with speaker recognition.

SRE19 will serve as a special evaluation allowing more in depth analysis and exploration into each of the data domains used in SRE18. There will be two components to SRE19: the SRE19 CTS Challenge and the SRE19 Audio-visual (AV) evaluation. The SRE19 CTS challenge will be conducted entirely online in a manner similar to the NIST 2014 and 2015 i-vector challenges (Greenberg et al., 2014; Tong et al., 2016), however actual audio recordings will be used as the source data instead of feature embeddings. Unexposed CTS data from the CMN2 corpus will be used to support the SRE19 CTS challenge. System performance scores will be

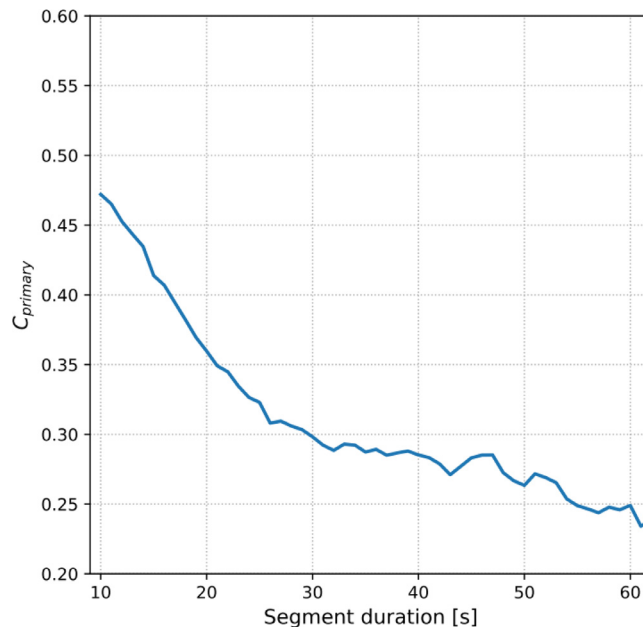


Fig. 2. Performance as a function of the speech duration in a test recording for a deep learning based system submission in SRE18. Systems performed consistently better as the speech duration increased, as anticipated.

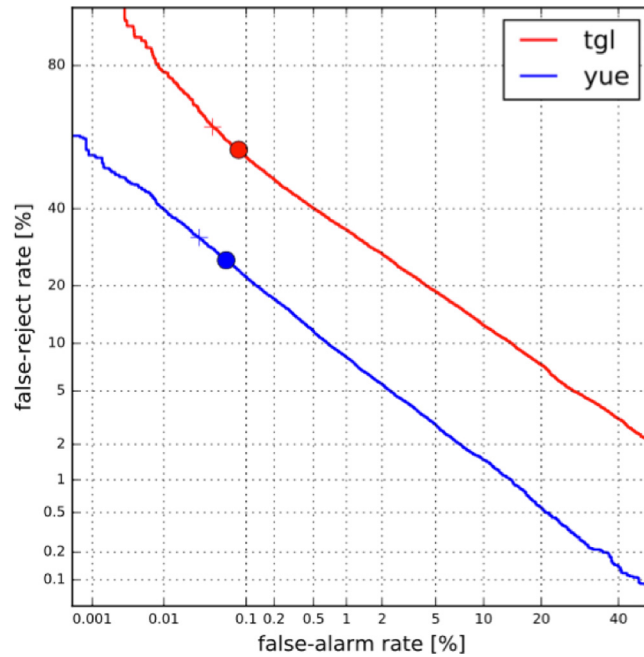


Fig. 3. DET curves for a leading system's performance on Tagalog speech (tgl) and Cantonese speech (yue) in SRE16. The circles denote the operating point that minimizes the detection cost function and the cross hairs denote the operating point selected by the system. System performances were consistently 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.

made available throughout the entire evaluation period instead of at the end, and multiple submissions will be allowed, enabling participants to explore how low they can drive error rates on the traditional CTS data domain.

The SRE19 AV evaluation will be conducted in the same manner as the traditional SREs, with training and development data released in early summer 2019, evaluation data released in late summer 2019, evaluation results submitted in October 2019, and a post-evaluation workshop held in December 2019¹⁶. Unexposed multimedia data from the VAST corpus will be used to support the SRE19 evaluation which will feature two core evaluation tracks: audio only and audio-visual fusion. An optional visual only track will also be available for participants.

The plans for SRE20 are based on the availability of a data corpus currently being collected by the LDC from multilingual speakers in both the CTS and AfV data domains. This corpus is designed to allow for explorations into cross-domain enroll-test

SRE18 Fixed Submissions

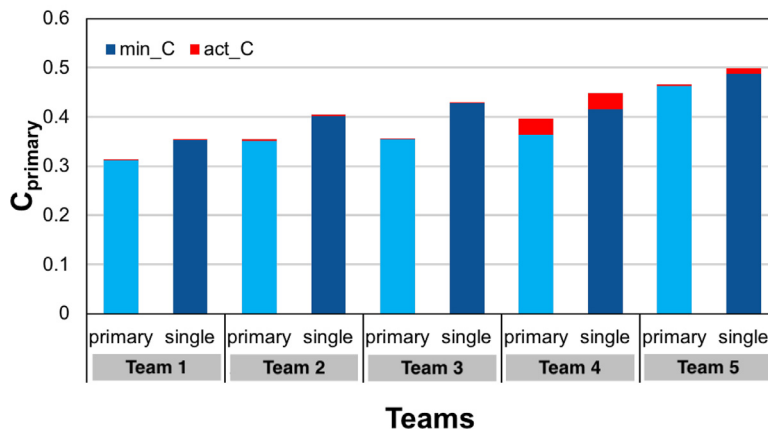


Fig. 4. A comparison of system performance for fused (primary) and the best single systems from five teams in SRE18. The detection cost is displayed at both the minimum operating (min_C) and the actual operation point (act_C). The observed differences between the fused system and single systems within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

¹⁶ The SRE19 workshop will be co-located with the 2019 IEEE Automatic Speech Recognition and Understanding (ASRU) Workshop in Sentosa, Singapore.

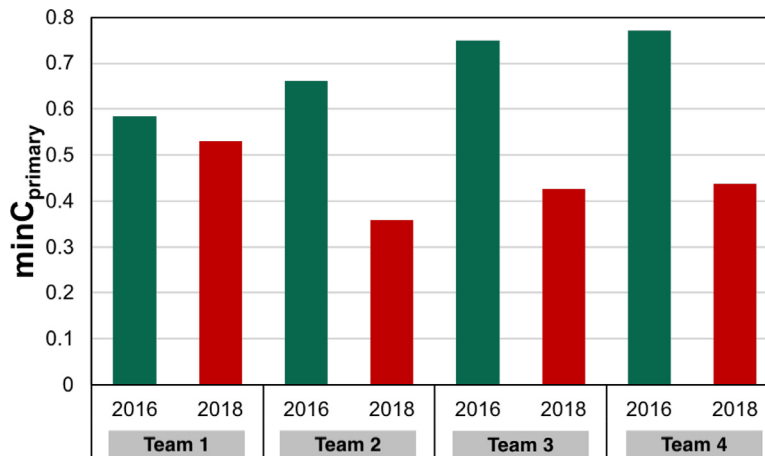


Fig. 5. A comparison of system performance for the SRE16 and SRE18 systems submitted by four teams that participated in both evaluations. A data set drawn from the *Call My Net* corpus Jones et al. (2017) was used to measure the performance of these 2016 and 2018 systems. Substantial improvements can be seen between the systems submitted in 2016 and those submitted in 2018.

trials (e.g. enroll on CTS data and test on AFV data for a single target speaker). The corpus is also designed to provide image data to support multimodal fusion explorations similar to SRE19. Continuing with the SRE16 and SRE18 data paradigms, this corpus is being collected outside of North America and will feature non-English data.

5. The future of NIST Speaker Recognition Evaluations

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

5.1. Evaluation tasks

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

Continuing with the core speaker detection task will also allow NIST to have a continued focus on the technological challenges presented by data domain and channel mismatches as new domains/channels become of interest to the speaker recognition research community. And as multimedia applications become more relevant to the speaker recognition community, like realtime group discussion transcription applications that use visual data to help with speaker identification, tasks involving the fusion of audio and video data such as those introduced in SRE19 are also anticipated to continue to be considered in future SREs.

5.2. Evaluation data

While the core SRE task will remain the same moving into the future, the data used to evaluate that task will continue to evolve in order to support exploration in more challenging domains and channels. Conversational telephony speech (CTS) data will remain a focus of the SRE series moving forward, and NIST maintains the goal of including recordings from different languages, from microphones with varying distances from the speaker, and different communication networks and technologies. It is anticipated that NIST will continue to partner with LDC to collect data for future evaluations. The collaboration has provided NIST with the largest amount of control over desired data collection parameters and data properties, which will become more important as more challenging data properties are introduced to the SRE series.

In addition to evolving CTS data characteristics, a continued progression towards data that mimics more realistic modern application conditions is also a possible focus area for future SREs (e.g., multimedia data, virtual assistant enabled devices, etc.). Recent SREs have leveraged publicly available speaker recognition data sources using “found data”,¹⁷ and this trend may continue in the future as long as these sources remain available for public research use.

¹⁷ VoxCeleb (Nagrani et al. 2017; Chung et al. 2018) and SITW (McLaren et al. 2016) corpora were allowable under the SRE18 fixed training condition.

Disclaimer

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

Acknowledgments

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

References

- Brümmer, N., Du Preez, J., 2006. Application-independent evaluation of speaker detection. *Comput. Speech Lang.* 20 (2–3), 230–275.
- Chung, J.S., Nagrani, A., Zisserman, A., 2018. VoxCeleb2: deep speaker recognition. In: *Proceedings of the INTERSPEECH 2018*, pp. 1086–1090.
- Croft, W.B., Lafferty, J., 2013. *Language Modeling for Information Retrieval*, Vol. 13. Springer Science & Business Media.
- Doddington, G.R., Przybocki, M.A., Martin, A.F., Reynolds, D.A., 2000. The NIST speaker recognition evaluation—overview, methodology, systems, results, perspective. *Speech Commun.* 31 (2–3), 225–254.
- Doddington, G.R., Schalk, T.B., 1981. Speech recognition: turning theory to practice: new ICs have brought the requisite computer power to speech technology; an evaluation of equipment shows where it stands today. *IEEE Spectr.* 18 (9), 26–32.
- Dorr, B.J., Fontana, P.C., Greenberg, C.S., Le Bras, M., Przybocki, M., 2016. Evaluation-driven research in data science: leveraging cross-field methodologies. In: *Proceedings of the IEEE International Conference on Big Data (Big Data)*, 2016, pp. 2853–2862.
- Greenberg, C.S., Bansé, D., Doddington, G.R., Garcia-Romero, D., Godfrey, J.J., Kinnunen, T., Martin, A.F., McCree, A., Przybocki, M., Reynolds, D.A., 2014. The NIST 2014 speaker recognition i-vector machine learning challenge. In: *Proceedings of the Odyssey 2014: The Speaker and Language Recognition Workshop*, pp. 224–230.
- Greenberg, C.S., Martin, A.F., Brandschain, L., Campbell, J.P., Cieri, C., Doddington, G.R., Godfrey, J.J., 2010. Human assisted speaker recognition in NIST SRE10. In: *Proceedings of the Odyssey 2010: The Speaker and Language Recognition Workshop*, pp. 180–185.
- Greenberg, C.S., Stanford, V.M., Martin, A.F., Yadagiri, M., Doddington, G.R., Godfrey, J.J., Hernandez-Cordero, J., 2013. The 2012 nist speaker recognition evaluation. In: *Proceedings of the INTERSPEECH 2013*, pp. 1971–1975.
- Jones, K., Strassel, S.M., Walker, K., Graff, D., Wright, J., 2017. Call my net corpus: a multilingual corpus for evaluation of speaker recognition technology. In: *Proceedings of the INTERSPEECH 2017*, pp. 2621–2624.
- Larcher, A., Lee, K.A., Ma, B., Li, H., 2014. Text-dependent speaker verification: classifiers, databases and RSR2015. *Speech Commun.* 60, 56–77.
- Leeuwen, D.A.V., 2009. Overall performance metrics for multi-condition speaker recognition evaluations. In: *Proc. INTERSPEECH 2009*, pp. 908–911.
- Martin, A., Doddington, G., Kamm, T., Ordowski, M., Przybocki, M., 1997. The DET curve in assessment of detection task performance. Technical Report. National Institute of Standards and Technology, Gaithersburg, MD.
- Martin, A., Przybocki, M., 2000. The NIST 1999 speaker recognition evaluation—an overview. *Dig. Signal Process.* 10 (1–3), 1–18.
- Martin, A., Przybocki, M., Campbell, J.P., 2005. The NIST speaker recognition evaluation program. *Biometric Systems*. Springer, pp. 241–262.
- Martin, A.F., Greenberg, C.S., 2009. NIST 2008 speaker recognition evaluation: performance across telephone and room microphone channels. In: *Proceedings of the INTERSPEECH 2009*, pp. 2579–2582.
- Martin, A.F., Greenberg, C.S., 2010. The NIST 2010 speaker recognition evaluation. In: *Proceedings of the INTERSPEECH 2010*, pp. 2726–2729.
- Martin, A.F., Przybocki, M.A., 2001. The NIST speaker recognition evaluations: 1996–2001. In: *Proceedings of the 2001: A Speaker Odyssey—The Speaker Recognition Workshop*, pp. 39–43.
- McLaren, M., Ferrer, L., Castan, D., Lawson, A., 2016. The speakers in the wild (SITW) speaker recognition database. In: *Proceedings of the INTERSPEECH 2016*, pp. 812–822.
- Nagrani, A., Chung, J.S., Zisserman, A., 2017. VoxCeleb: a large-scale speaker identification dataset. In: *Proceedings of the INTERSPEECH 2017*, pp. 2616–2620.
- Nakasone, H., Beck, S., 2001. Forensic automatic speaker recognition. In: *Proceedings of the Odyssey 2001: The Speaker and Language Recognition Workshop*.
- Ortega-García, J., Gonzalez-Rodriguez, J., Marrero-Aguilar, V., 2000. AHUMADA: a large speech corpus in spanish for speaker characterization and identification. *Speech Commun.* 31 (2–3), 255–264.
- Pallett, D.S., 2003. A look at NIST's benchmark ASR tests: past, present, and future. In: *Proceedings of the IEEE ASRU Workshop 2003*, pp. 483–488.
- Przybocki, M., Martin, A., Le, A., 2007. NIST speaker recognition evaluations utilizing the mixer corpora – 2004, 2005, 2006. *IEEE Trans. Audio Speech Lang. Process.* 15 (7), 1951–1959.
- Przybocki, M.A., Martin, A.F., 1999. The 1999 NIST speaker recognition evaluation, using summed two-channel telephone data for speaker detection and speaker tracking. In: *Proceedings of the EUROSPEECH 1999*, pp. 2215–2218.
- Przybocki, M.A., Martin, A.F., 2004. NIST speaker recognition evaluation chronicles. In: *Proceedings of the Odyssey 2004: The Speaker and Language Recognition Workshop*, pp. 15–22.
- Przybocki, M.A., Martin, A.F., Le, A.N., 2006. NIST speaker recognition evaluation chronicles—Part 2. In: *Proceedings of the Odyssey 2006: The Speaker and Language Recognition Workshop*, pp. 1–6.
- Reynolds, D., Andrews, W., Campbell, J., Navratil, J., Peskin, B., Adami, A., Jin, Q., Klusacek, D., Abramson, J., Mihaescu, R., et al., 2003. The SuperSID project: exploiting high-level information for high-accuracy speaker recognition. In: *Proceedings of the IEEE ICASSP 2003*, Vol. 4, pp. IV–784.
- Rose, T., Fiscus, J., Over, P., Garofolo, J., Michel, M., 2009. The TRECVID 2008 event detection evaluation. In: *Proceedings of the IEEE Workshop on Applications of Computer Vision (WACV)*, 2009, pp. 1–8.
- Sadjadi, S.O., Greenberg, C.S., Reynolds, D.A., Mason, L., 2019. The 2018 nist speaker recognition evaluation. In: *Proceedings of the INTERSPEECH 2019*, pp. 1483–1487.
- Sadjadi, S.O., Kheyrikhah, T., Tong, A., Greenberg, C.S., Reynolds, D.A., 2017. The 2016 nist speaker recognition evaluation. In: *Proceedings of the INTERSPEECH 2017*, pp. 1353–1357.
- The Economist, 2017a. Data is giving rise to a new economy. <https://www.economist.com/briefing/2017/05/06/data-is-giving-rise-to-a-new-economy>. Accessed: 2019-05-05.
- The Economist, 2017b. 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>. Accessed: 2019-05-05.
- Tong, A., Greenberg, C., Martin, A., Banse, D., Howard, J., Zhao, H., Doddington, G., Garcia-Romero, D., McCree, A., Reynolds, D., Singer, E., Hernandez-Cordero, J., Mason, L., 2016. Summary of the 2015 NIST language recognition i-vector machine learning challenge. In: *Proceedings of the Odyssey 2016: The Speaker and Language Recognition Workshop*, pp. 297–302.
- Tracey, J., Strassel, S., 2018. VAST: a corpus of video annotation for speech technologies. In: *Proceedings of the LREC 2018*, pp. 4318–4321.