

The FFSVC 2020 Evaluation Plan

Xiaoyi Qin¹, Ming Li¹, Hui Bu⁴, Rohan Kumar Das², Wei Rao², Shrikanth Narayanan³, Haizhou Li²

¹Data Science Research Center, Duke Kunshan University, Kunshan, China

²Department of Electrical & Computer Engineering, National University of Singapore, Singapore

³Signal Analysis and Interpretation Lab, University of Southern California, Los Angeles, USA

⁴AI Shell Foundation, Beijing, China

1. Introduction

Speaker verification is a key technology in speech processing and biometric authentication, which has broad impact on our daily lives, e.g. security, customer service, mobile devices, smart speakers. Recently, speech based human computer interaction has become more and more popular in far-field smart home and smart city applications, e.g. mobile devices, smart speakers, smart TVs, automobiles. Due to the usage of deep learning methods, the performances of speaker verification in telephone channel and close-talking microphone channel have been enhanced dramatically. However, there are still some open research questions that can be further explored for speaker verification in the far-field and complex environments, including but not limited to

- Far-field text-dependent speaker verification for wake up control
- Far-field text-independent speaker verification with complex environments
- Far-field speaker verification with cross-channel enrollment and test
- Far-field speaker verification with single multi-channel microphone array
- Far-field speaker verification with multiple distributed microphone arrays
- Far-field speaker verification with front-end speech enhancement methods
- Far-field speaker verification with end-to-end modeling using data augmentation
- Far-field speaker verification with front-end and back-end joint modeling
- Far-field speaker verification with transfer learning and domain adaptation

The Far-Field Speaker Verification Challenge 2020 (FFSVC20) is designed to boost the speaker verification research with special focus on far-field distributed microphone arrays under noisy conditions in real scenarios. The objectives of this challenge are to: 1) benchmark the current speech verification technology under this challenging condition, 2) promote the development of new ideas and technologies in speaker verification, 3) provide an open, free, and large scale speech database to the community that exhibits the far-field characteristics in real scenes.

The challenge has three tasks in different scenes.

- Task 1: Far-Field Text-Dependent Speaker Verification from single microphone array

Table 1: The details of the FFSVC20 challenge data

Utterance ID	Content	Noise
001-030	ni hao, mi ya (text-dependent)	F - TV/Office + electric fan T - electric fan
091-	text independent	S - clean

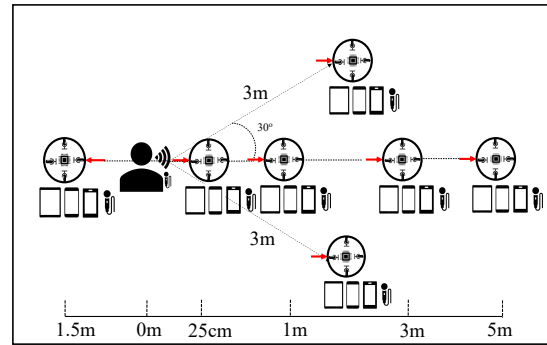


Figure 1: The setup of the recording environment

- Task 2: Far-Field Text-Independent Speaker Verification from single microphone array
- Task 3: Far-Field Text-Dependent Speaker Verification from distributed microphone arrays

All three tasks follow the cross-channel setup. The recordings of close-talking cellphone will be selected as enrollment and the recordings of far-field microphone array will be used for test.

2. Database

2.1. The FFSVC20 Challenge Database

2.1.1. The DMASH Database

The Distributed Microphone Arrays in Smart Home (DMASH) database is recorded in real smart home scenarios with two different rooms. The detail information of room sizes will be released after the challenge. Figure 1 shows the recording environment setup of DMASH, includes one close-talking microphone; one iPhone, one Android phone, one iPad, one microphone and one circular microphone array (named PCM) placed at 25cm, 1m, 3m, 5m, left 3m, right 3m and -1.5m distances, respectively.

2.1.2. FFSVC20 Challenge Data Description

The FFSVC20 challenge database is part of the DMASH Database. The recording devices include one close-talking microphone (48kHz, 16 bit), one iPhone (48kHz, 16 bit) at 25cm distance and 6 circular microphone arrays (16kHz, 16bit, 16 microphones, 5cm radius). The language is Mandarin. Text content include *ni hao mi ya* as text dependent utterances as well as other text independent ones.

The data collection setup of the challenge database is shown in Figure 2. Red arrow points to channel 0 of microphone arrays. Each speaker visits 3 times with 7-15 days gap.

The first letter of the file name denotes the visit index. *F* stands for the speaker’s first visit, *S* denotes the speaker’s second visit and *T* means the speaker’s third visit.

2.1.3. Named structure

Here is an example of speaker files.

```
T0003/
  003MIC/
  003I0.25M/
  003PCM5M/
  003PCML3M/
  003PCM3M/
  .....
  T0003-003PCM3M_recorded14_0308_normal.wav
  .....
```

The corresponding structure is as follows,

```
<visit><spk_id>/
  <spk_id><device_and_distance>/
    <visit><spk_id>_<spk_id><device_distance>
      _<channel_id>_<utt_id>_<speed>.wav
```

In PCM (microphone array), *recorded 2* stand for channel 0, *recorded 6* denotes channel 4 and so on (in total there are 18 channels in each PCM, *recorded 0, 1* is empty).

Here we provide three examples.

- *F0148_148I0.25M_1_0218_normal.wav* means this audio is the utterance 218 in the first visit of speaker 148, the recorded device is iPhone at a distance of 25cm. 1 stands for channel 1, which is meaningless since iPhone at 25cm distance only contains one channel.
- *S0183_183MIC_Tr2_0138_normal.wav* means this audio is the utterance 138 in the second visit of speaker 183, the recorded device is a close-talking MIC, *Tr2* stands for close-talking. In this challenge, we just provide one close-talking microphone.
- *T0003-003PCML3M_recorded14_0308_normal.wav* means this audio is the utterance 308 in the third visit of speaker 3, the recorded device is a microphone array (this audio is channel 12 in this array), located at 3m distance in front of the speaker, on the left side with a 30 degrees angle.

In this challenge database, we provide three randomly selected microphone arrays out of the total six arrays in the training and development set; for each microphone array, we only provide 4 channels’ data (channel 0,4,8,12, denoted by *recorded 2,6,10,14*) due to the large size of the whole database.

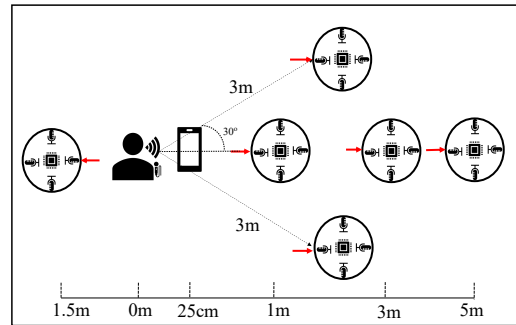


Figure 2: The setup of the FFSVC20 challenge data

2.2. The SLR-85 HI-MIA database

The original HI-MIA database includes two sub databases, which are the AISHELL- wakeup1 with 254 speakers and the AISHELL-2019B-eval with 86 speakers. The content of utterances covers two wake-up words, ‘*ni hao, mi ya*’ in chinese and ‘*Hi, Mia*’ in English.

During the recording process, seven recording devices (one close-talking microphone and six 16-channel circular microphone arrays) were set in a real smart home environment. The 16-channel circular microphone array records signals in the 16kHz, 16 bit format, and the close-talking microphone records waveforms in the 44.1kHz, 16 bit format.

¹The SLR-85 HI-MIA open source database is the 2019 AISHELL Speaker Verification Challenge database which is a subset of the original HI-MIA database. It contains one close-talking microphone and three microphone arrays located at 1m, 3m and 5m distance right in front of speaker. The SLR-85 HI-MIA database covers all the 254 speakers, but only includes the mandarin utterances.

For more details, please refer to the latest version (v3) of [1].

3. Task Description

3.1. Task 1:Far-Field Text-Dependent Speaker Verification from single microphone array

3.1.1. Training data

The training data includes 120 speakers and each speaker has 3 visits. In each visit, there are multiple (‘*ni hao, mi ya*’) text-dependent utterances as well as multiple text-independent utterances. The recording from five recording devices for each utterance are provided for training. These five recording devices include one close-talk microphone, one 25cm distance cellphone, and three randomly selected microphone arrays (4 channels per array).

Any publicly open and freely accessible speech database shared on openslr.org before Feb 1st 2020 (including SLR-85 HI-MIA) can be used in this challenge.

3.1.2. Development Data

The Development data includes 35 speakers and each speaker has 3 visits. In each visit, there are multiple (‘*ni hao, mi ya*’) text-dependent utterances as well as multiple text-independent utterances. The recording from five recording devices for each utterance are provided. These five recording devices include

¹<http://openslr.org/85/>

one close-talk microphone, one 25cm distance cellphone, and three randomly selected microphone arrays (4 channels per array).

3.1.3. Evaluation Data

The evaluation data includes 80 speakers and each speaker has 3 visits. In each visit, there are multiple ('*ni hao, mi ya*') utterances, The recording from two recording devices for each utterance are provided. These two recording devices include one 25cm distance cellphone and one randomly selected microphone arrays (4 channels per array).

The recording from 25cm distance cellphone will be selected as enrollment and recording from single far-field microphone array will be used for test. For any true trial, the enrollment and the testing utterances are from different visits of the same speaker.

3.2. Task 2: Far-Field Text-Independent Speaker Verification from single microphone array

3.2.1. Training data

The same as the training data for task 1.

3.2.2. Development Data

The same as the development data for task 1.

3.2.3. Evaluation Data

The evaluation data includes 80 speakers and each speaker has 3 visits. In each visit, there are multiple text-independent utterances, The recording from two recording devices for each utterance are provided. These two recording devices include one 25cm distance cellphone and one randomly selected microphone arrays (4 channels per array).

The recording from 25cm distance cellphone will be selected as enrollment and recording from single far-field microphone array will be used for test. For any true trial, the enrollment and the testing utterances are from different visits of the same speaker.

3.3. Task 3: Far-Field Text-Dependent Speaker Verification from distributed microphone arrays

3.3.1. Training data

The same as the training data for task 1.

3.3.2. Development Data

The same as the development data for task 1.

3.3.3. Evaluation Data

The evaluation data includes 80 speakers and each speaker has 3 visits. In each visit, there are multiple ('*ni hao, mi ya*') utterances. For each utterance, its corresponding recordings from one 25cm distance cellphone and 2-4 randomly selected microphone arrays are provided. For each microphone array, the selected four microphones are equally distributed along the circle with a random start channel index to simulate the scenarios with unknown array orientation angles. (e.g. channel 0, 4, 8, 12; channel 1, 5, 9, 13; channel 6, 10, 14, 2, etc.)

Recording from 25cm distance cellphone will be selected as enrollment and the recordings from 2-4 randomly selected far-field microphone arrays will be used for test. For any true

trial, the enrollment and the testing utterances are from different visits of the same speaker.

There is no overlapping among the speakers in the training data, development data, evaluation data in task 1, task 2, and task 3.

4. Evaluation Rules

4.1. Evaluation Results

Before the mid-term deadline of score submission, each team have 5 times to submit the result. We sincerely suggest each team to test the system performance on the development set due to limited opportunities for submission. The evaluation set and the development set are from the same database, the only difference is that the development set only has 35 speakers, while the evaluation set for each task has 80 speakers.

After the mid-term deadline and before the final deadline, each team has another 5 times to submit the final score file.

For the results on the evaluation set, we will release the results calculated based on a fixed 30% of the trials in the leaderboard. So ranking on the leaderboard is not the final ranking. We will announce the official result in the Interspeech 2020 FFSVC special session. We encourage each team to explore more novel ideas, not just for the first place.

4.2. The trials

The trial file consists of three segments: enrollment audio ID, test audio ID and label. Label denotes that the trial is target or non-target. The 25cm distance iPhone signal is selected as the enrollment data, and the microphone array audio is considered as the testing data.

4.3. Performance Measures

In this challenge, we will use several metric to evaluate the system performance. The primary metric we adopt is the min C_{det} cost value. In addition, Equal Error Rate (EER) and C_{llr} will be provided to participant as auxiliary metrics.

4.3.1. Primary measures

The primary metric is based on the following detection cost function which is the same function as used in the NIST 2010 SRE, but with modified parameters [2]. It is a weighted sum of miss and false alarm error probabilities in the form:

$$C_{det} = C_{miss} \times P_{miss} \times P_{tar} + C_{fa} \times P_{fa} \times (1 - P_{tar}) \quad (1)$$

We assume a prior target probability, P_{tar} of 0.01 and equal costs between misses and false alarms. The model parameters are 1.0 for both C_{miss} and C_{fa} . The C_{det} will be normalized by P_{tar} the same way as in [2].

We use the minimum C_{det} as our primary metric.

4.3.2. Alternate Performance Metrics

The EER and C_{llr} will be provided as alternate performance metrics and a brief description of C_{llr} is provided as follows:

To analyze how well a system performs and is calibrated across different operating points, a log-likelihood ratio based cost metric, C_{llr} , is suggested. Assuming trials scores are rep-

resented as LLRs, then C_{llr} can be calculated as [3],

$$C_{llr} = \frac{1}{2 \times \log(2)} \times \left(\frac{\sum \log(1 + 1/s)}{N_{tar}} + \frac{\sum \log(1 + s)}{N_{non}} \right) \quad (2)$$

where s is the likelihood ratio for a trial, and N_{tar} and N_{non} represent the number of target and non-target trials, respectively.

5. Registration and Submission

5.1. Registration

First, please create an account if you do not have one. We kindly request you to associate your account to an institutional e-mail. The organizing committee reserves the right to revoke your access to the challenge sites otherwise, please read the evaluation plan carefully. If you are part of a team, at least one person in your team will need an account to participate. Make sure to set the name of your team in the user's profile, or it will not be visible on the leaderboard.

Participants can register in one or more tasks. If your team participates in multiple tasks, we kindly request you to use the same user account to participate in all tasks.

Please note that any deliberate attempts to bypass the submission limit (for instance, by creating multiple accounts and using them to submit) will lead to automatic disqualification.

In case of any issues, the final interpretation right belongs to the organizing committee.

5.2. Submission

5.2.1. Score submission

Participants are required to submit at least one valid score file for each participating task to the FFSVC 2020 platform. The score files need to follow the following rules:

`<TeamName>.<Task>.<SystemNumber>.txt`

The score files should be in UTF-8 format with one line per trial. Each line must include two space-delimited fields:

`<enrol_data><space><test_data><space><score>`

We will provided a sample in the challenge website.

5.2.2. System description submission

Each registered team is required to submit a technical system description report. Please submit this report using the Interspeech 2020 paper template. All reports must be a minimum of 2 pages (including references). Reports must be written in English. The system description does not need to repeat the content of the evaluation plan, such as the introduction of database, evaluation metric, etc. The system description must include the following items:

- a complete description of the system components, including front-end (e.g., speech activity detection, features, normalization, front-end speech enhancement) and back-end (e.g., background models, i-vector/embedding extractor, Probabilistic Linear Discriminant Analysis (PLDA), speaker features fusion) modules along with their configurations (i.e., filterbank configuration, dimensionality and type of the acoustic feature parameters, as well as the acoustic model and the backend model configurations).

ters, as well as the acoustic model and the backend model configurations).

- a complete description of the data partitions used to train the various models.
- performance of the submission systems on the development dataset (calculated based on the provided tools) and the evaluation dataset (calculated by the challenge platform). Teams are encouraged to quantify the contribution of their major system components that they believe resulted in significant performance gains.²
- novel ideas, strategies and methods are strongly recommended to be shared.
- a report of the model size, CPU (single threaded) and GPU execution times as well as the amount of memory used to process a single trial (i.e., the time and memory used for creating a speaker model from enrollment data as well as processing a test segment to compute the score).

5.2.3. Paper submission

The organizing committee highly encourages the participating teams to submit a paper to the INTERSPEECH 2020 Far-Field Speaker Verification Challenge special session. However, paper contributions within the scope are also welcome if the authors do not intend to participate in the Challenge itself. In any case, please submit your paper until 30 March 2020 (and final results by 15 June 2020) using the standard style info and respecting length limits, and submit to the Interspeech 2020 paper submission system. Important: as topic you should choose only this Special Session (FFSVC 2020). The papers will undergo the normal review process similar to the regular session papers.

6. Schedule

- Feb 1st: Releasing the training and development data as well as the evaluation plan
- March 1st: Releasing the evaluation data and launching the leaderboard (30% of the trials)
- March 15th: Challenge registration deadline
- March 23th: Mid-term deadline of score submission (up to 5 chances)
- March 30th: Interspeech 2020 paper submission deadline
- June 15th: Final deadline of score submission (another 5 chances)
- Interspeech 2020 Camera Ready Paper deadline: System description submission deadline
- Interspeech 2020 special session: Official results announcement

7. References

- [1] X. Qin, H. Bu, and M. Li, "HI-MIA : A Far-field Text-Dependent Speaker Verification Database and the Baselines," in *Proc. of ICASSP*, 2020, arXiv:1912.01231.
- [2] M. K. Nandwana, J. van Hout, M. McLaren, C. Richey, A. Lawson, and M. A. Barrios, "The voices from a distance challenge 2019 evaluation plan," *arXiv:1902.10828*, 2019.

²https://www.nist.gov/system/files/documents/2019/07/22/2019_nist_speaker_recognition_challenge_v8.pdf

- [3] N. Brummer and J. Preez, "Application-independent evaluation of speaker detection," *Computer Speech & Language*, vol. 20, pp. 230–275, 04 2006.