

# Speaker Diarization with Lexical Information

Tae Jin Park<sup>1</sup>, Kyu J. Han<sup>2</sup>, Jing Huang<sup>2</sup>, Xiaodong He<sup>2</sup>, Bowen Zhou<sup>2</sup>, Panayiotis Georgiou<sup>1</sup> and Shrikanth Narayanan<sup>1</sup>

<sup>1</sup>University of Southern California

<sup>2</sup>JD AI Research

taejinpa@usc.edu

## Abstract

This work presents a novel approach for speaker diarization to leverage lexical information provided by automatic speech recognition. We propose a speaker diarization system that can incorporate word-level speaker turn probabilities with speaker embeddings into a speaker clustering process to improve the overall diarization accuracy. To integrate lexical and acoustic information in a comprehensive way during clustering, we introduce an adjacency matrix integration for spectral clustering. Since words and word boundary information for word-level speaker turn probability estimation are provided by a speech recognition system, our proposed method works without any human intervention for manual transcriptions. We show that the proposed method improves diarization performance on various evaluation datasets compared to the baseline diarization system using acoustic information only in speaker embeddings.

**Index Terms:** speaker diarization, automatic speech recognition, lexical information, adjacency matrix integration, spectral clustering.

## 1. Introduction

Speaker diarization is a process of partitioning a given multi-speaker audio signal in terms of “who spoke when”, generally consisting of two sub-processes: *speaker segmentation* of cutting the given audio into homogeneous speech segments in terms of speaker characteristics and *speaker clustering* of grouping all the segments from one speaker into the same cluster and assigning them with the same speaker label. Speaker diarization plays a critical role in speech applications like automatic speech recognition (ASR) or behavioral analytics [1, 2, 3, 4].

Speaker diarization has long been considered a pre-processing step in the context of ASR. This is mostly because, considering research setups where oracle results for speech activity detection or segmentation are given, grouping speech portions from the same speakers in a multi-speaker audio signal can benefit ASR systems. It can enable speaker-specific feature transformation, e.g., fMLLR [5] or total variability factor analysis for i-vectors [6]. However, such oracle results would not be available for speaker diarization in practice. Also, performing speaker diarization before ASR on production systems in the wild without proper post-processing would degrade recognition accuracy significantly since it is likely to determine speaker change points in the middle of words, not between words, and result in word cuts or deletions. These practical issues were also pointed out in [7]. In addition, we recently showed in [8] that lexical cues in words or utterances can help diarization accuracy improve when combined with acoustic features. All of these suggest that it would make more sense and practical for speaker diarization to be considered as a *post-processing* step

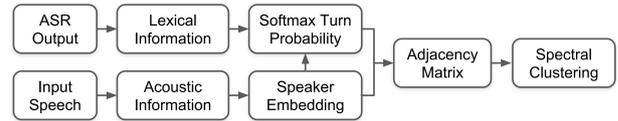


Figure 1: The data flow of the proposed system.

and take advantage of utilizing ASR outputs within the ASR pipeline. With this regard, in this paper, we assume that there are available ASR outputs in a text form for speaker diarization and propose a system to incorporate such lexical information into the diarization process.

There have been a handful of works to employ ASR outputs to enhance speaker diarization systems, but mostly limited to speaker segmentation. ASR outputs are used in [9] for determining potential speaker change points. In [10], the lexical information provided by an ASR system is utilized to train a character-level language model and improve speaker segmentation performance. In our previous work [8], we exploited lexical information, from either reference transcripts or ASR outputs, along with acoustic information to enhance speaker segmentation in estimating speaker turns and showed the overall improvement in speaker diarization accuracy.

In this paper we extend the exploitation of lexical information provided by an ASR system to a *speaker clustering* process in speaker diarization. The challenge of employing lexical information to speaker clustering is multifaceted and requires practical design choices. In our proposal, we use *word-level speaker turn probabilities* as lexical representation and combine them with acoustic vectors of *speaker embedding* when performing *spectral clustering* [11]. In order to integrate lexical and acoustic representations in the spectral clustering framework, we create *adjacency matrices representing lexical and acoustic affinities between speech segments respectively* and combine them later with a per-element max operation. It is shown that the proposed speaker diarization system improves a baseline performance on two evaluation datasets.

The rest of the paper is organized as follows. In Section 2, we explain the data flow of our proposed speaker diarization system. In Sections 3 and 4, we detail how we process acoustic and lexical information, respectively. In Section 5, we describe the integration of the two sets of information in the framework of spectral clustering. Experimental results are discussed in Section 6 and we conclude the work with some remarks in Section 7.

## 2. Proposed speaker diarization system

The overall data flow of our proposed speaker diarization system is depicted in Fig 1. In the proposed system, there are two streams of information: lexical and acoustic. On the lexical information side, we receive the automated transcripts with the corresponding time stamps for word boundaries from an avail-

able ASR system. This text information is passed to the speaker turn probability estimator to compute word-level speaker turn probabilities. On the acoustic information side, we perform a common diarization task. MFCCs are extracted from the input speech signal after speech activity detection (SAD). Following SAD, we uniformly segment the SAD outputs. These uniform segments are relayed to the speaker embedding extractor that provides speaker embedding vectors. We use the publicly available Kaldi ASPIRE SAD Model<sup>1</sup> [12] for SAD in our proposed diarization pipeline.

After processing the two streams of information, we create two adjacency matrices which hold lexical as well as acoustic affinities between speech segments, respectively, and combine them with a per-element max operation to generate the combined affinity matrix that contains lexical and acoustic information in a comprehensive space. With the integrated adjacency matrix, we finally obtain speaker labels using a spectral clustering algorithm. Each module in Fig. 1 is explained in more details in the following sections.

### 3. Acoustic information stream: Speaker embedding extractor

We employ the x-vector model<sup>2</sup> [13] as our speaker embedding generator that has been showing the state-of-the-art performances for speaker verification and diarization tasks. To perform the general diarization task with acoustic information in the proposed system pipeline, we use 0.5 second window, 0.25 second shift and 0.5 second minimum window size for 23-dimensional MFCCs. The performance improvement of speaker embedding is out of the scope of this paper.

### 4. Lexical information stream: Speaker turn probability estimator

While acoustic speaker characteristics can be used for speaker turn detection tasks [14], our proposal of word-level speaker turn probability estimation comes behind the reasoning that lexical data can also provide an ample amount of information for similar tasks. It is likely for words in a given context (i.e., utterance) to have different probabilities on whether speaker turns change at the time of being spoken. For example, the words in the phrase “how are you” are very likely to be spoken by a single speaker rather than by multiple speakers. This means that each word in this phrase “how are you” would likely have lower speaker turn probabilities than the word right after the phrase would have. In addition to lexical information, we also fuse a speaker embedding vector per each word to increase the accuracy of the turn probability estimation.

To estimate speaker turn probability, we train bi-directional three-layer gated recurrent units (GRUs) [15] with 2,048 hidden units on the Fisher [16] and Switchboard [17] corpora using the force-aligned texts. The actual inputs to the proposed speaker turn probability estimator would be the decoder outputs of the ASR. The words and the corresponding word boundaries are used to generate word embedding and speaker embedding vectors respectively, as follows:

- **Speaker embedding vector (S):** With the given start and end time stamps of a word from ASR, we retrieve the speaker embedding vector using the speaker embedding extractor described in Section 3. The x-vector speaker embedding is 128-dimensional.

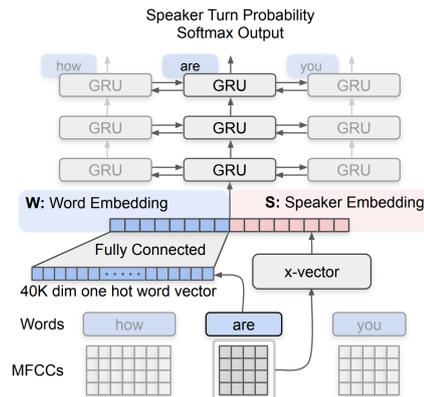


Figure 2: An illustration of the proposed speaker turn probability estimator.

- **Word embedding vector (W):** We map the same word input to a 40K-dimensional one-hot vector, which is fully connected to the word embedding layer shown in Fig. 2. The dimension of the embedding layer is set to 256.

These two vectors are appended to make a 384-dimensional vector for every word and fed to the GRU layer. The softmax layer has one node and, during inference, outputs speaker turn probability. The parameters of the speaker turn probability estimator are trained with the cross entropy loss. The ASR system used in this paper for decoding is the Kaldi ASPIRE recipe<sup>3</sup> [12] that is publicly available.

## 5. Adjacency matrix integration

### 5.1. Adjacency matrix calculation

The biggest challenge of integrating speaker turn probabilities (from lexical information) and speaker embedding vectors (from acoustic information) in the spectral clustering framework is the heterogeneity of the information sources for these representations. To tackle this challenge, we first create two independent adjacency matrices that contain lexical and acoustic affinities between speech segments, respectively, and then combine them with a per-element max operation to handle the information from the two different sources in the common space used for spectral clustering. For each adjacency matrix, we employ undirected graphs to represent the corresponding affinities between the segments.

- **Adjacency matrix using speaker embeddings**

- 1) Initially compute the cosine similarity  $p_{i,j}$  between speaker embedding vectors for segments  $s_i$  and  $s_j$  to form the adjacency matrix  $\mathbf{P}$ , where  $1 \leq i, j \leq M$  and  $M$  is the total number of segments in a given audio signal.
- 2) For every  $i$ -th row of  $\mathbf{P}$ , update  $p_{i,j}$  as follows:

$$p_{i,j} = \begin{cases} 1 & \text{if } p_{i,j} \leq W(r) \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

where  $W(r)$  is the cosine similarity value that is at  $r$ -percentile in each row and  $r$  is optimized on the dev set. This operation converts  $\mathbf{P}$  to a discrete-valued affinity matrix through  $N$  nearest neighbor connections.

- 3) Note that  $\mathbf{P}$  is asymmetric and can be seen an adjacency matrix for a directed graph where each node represents a

<sup>1</sup><http://kaldi-asr.org/models/m4>

<sup>2</sup><http://kaldi-asr.org/models/m6>

<sup>3</sup><http://kaldi-asr.org/models/m1>

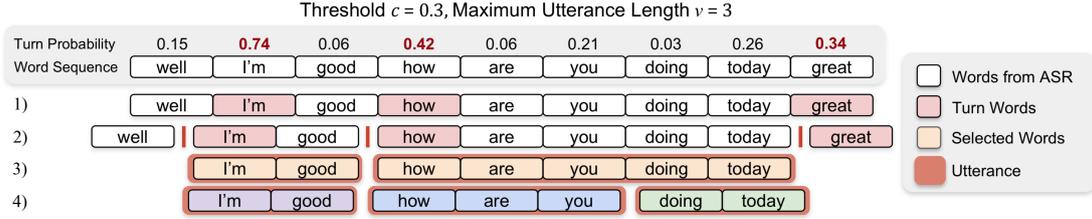


Figure 3: An example of the word sequence processing for the adjacency matrix calculation using the speaker turn probabilities.

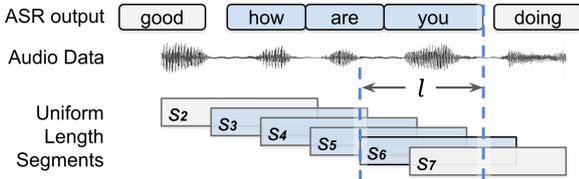


Figure 4: An example of the speech segment selection process using the utterance boundary information.

speech segment in our case. As spectral clustering finds the minimum cuts on an *undirected* graph in theory [11], we choose an undirected version of  $\mathbf{P}$ ,  $\mathbf{P}_{ud}$ , as the adjacency matrix for speaker embeddings by averaging  $\mathbf{P}$  and  $\mathbf{P}^T$  as below:

$$\mathbf{P}_{ud} = \frac{1}{2}(\mathbf{P} + \mathbf{P}^T) \quad (2)$$

The pictorial representation of  $\mathbf{P}_{ud}$  is given in the left side of Fig. 5.

#### • Adjacency matrix using speaker turn probabilities

The following steps 1) to 4) match to the numbered illustrations in Fig. 3, where  $c = 0.3$  and  $\nu = 3$  are given as example parameters.

- 1) For a given threshold  $c$ , pick all the *turn words* that have speaker turn probabilities greater than  $c$  in the word sequence provided by ASR. The colored boxes in Fig. 3-1) indicate the turn words. The threshold  $c$  is determined by the eigengap heuristic that we will discuss in Section 5.2.
- 2) Break the word sequence at every point where the turn word starts as in Fig. 3-2). The given word sequence is broken into multiple utterances as a result.
- 3) Pick all the utterances that have more than one word because a duration spanning one word may be too short to carry any speaker-specific information. In Fig. 3-3), the words “well” and “great” are thus not considered for further processing. Additionally, we always arrange seven back channel words (“yes”, “oh”, “okay”, “yeah”, “uh-huh”, “mhm”, “[laughter]”) as independent utterances regardless of their turn probabilities.
- 4) To mitigate the effect of any miss detection by the speaker turn probability estimator, we perform over-segmentation on the utterances by limiting the max utterance length to  $\nu$ . In Fig. 3-4), the resulting utterances are depicted with different colors. Maximum utterance length  $\nu$  is optimized on the dev set in the range of 2 to 9.
- 5) Find all the speech segments that fall into the boundary of each utterance. Fig. 4 explains how speech segments within the boundary of the example utterance “how are you” are selected. If a segment partly falls into the utterance boundary and its overlap ( $l$  in Fig. 4) is greater than 50% of the segment length, the segment is considered to fall into the utterance boundary.

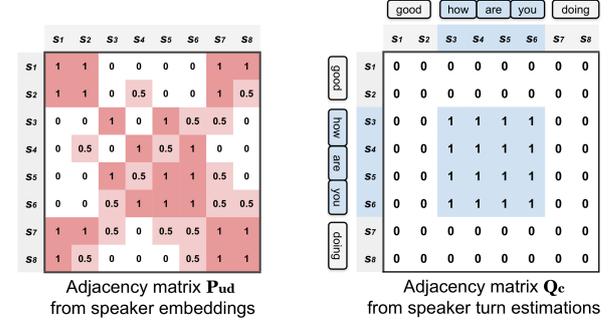


Figure 5: Examples of the two adjacency matrices.

- 6) Let  $s_m$  be the first segment and  $s_n$  be the last segment falling into the utterance boundary (e.g., segments  $s_3$  and  $s_6$ , respectively, in Fig. 4). For the elements  $q_{i,j}$  in an adjacency matrix  $\mathbf{Q}_c$  (with the threshold  $c$ ) being initialized with zeros, we do the following operation for all the utterances:

$$q_{i,j} = \begin{cases} 1 & \text{if } m \leq i, j \leq n \\ q_{i,j} & \text{otherwise} \end{cases} \quad (3)$$

The right side of Fig. 5 shows an example of  $\mathbf{Q}_c$  by the utterance “how are you” in Fig. 4.

#### • Combining adjacency matrices

We combine the two adjacency matrices:

$$\mathbf{A}_c = \max(\mathbf{P}_{ud}, \mathbf{Q}_c) = \max\left(\frac{1}{2}(\mathbf{P} + \mathbf{P}^T), \mathbf{Q}_c\right) \quad (4)$$

where  $\max$  is a per-element max operation.

#### 5.2. Eigengap analysis

In spectral clustering, the Laplacian matrix is employed to get the spectrum of an adjacency matrix. In this work, we employ the unnormalized graph Laplacian matrix  $\mathbf{L}_c$  [11] as below:

$$\mathbf{L}_c = \mathbf{D}_c - \mathbf{A}_c \quad (5)$$

where  $\mathbf{D}_c = \text{diag}\{d_1, d_2, \dots, d_M\}$ ,  $d_i = \sum_{k=1}^M a_{ik}$  and  $a_{ij}$  is the element in the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column of the adjacency matrix  $\mathbf{A}_c$ . We calculate eigenvalues from  $\mathbf{L}_c$  and set up an eigengap vector  $\mathbf{e}_c$ :

$$\mathbf{e}_c = [\lambda_2 - \lambda_1, \lambda_3 - \lambda_2, \dots, \lambda_M - \lambda_{M-1}] \quad (6)$$

where  $\lambda_1$  is the smallest eigenvalue and  $\lambda_M$  is the largest eigenvalue. The resulting adjacency matrix  $\mathbf{A}_c$  is passed to the spectral clustering algorithm, for which we use the implementation in [18].

The number of clusters (in our case, number of speakers) is estimated by finding the  $\arg \max$  value of the eigengap vector  $\mathbf{e}_c$  as in the following equation:

$$\widehat{n}_s = \arg \max_n(\mathbf{e}_c) \quad (7)$$

Table 1: DER (%) on the RT03-CTS dataset.

Number of Speakers		Unknown			
Dataset Split(Quantity)		Dev(14)		Eval(58)	
Error Type		DER	SER	DER	SER
Quan <i>et al.</i> [19] System SAD		-	-	12.3	3.76
<b>Baseline</b>	<b>M1</b>	4.00	1.03	6.97	2.90
<b>Proposed</b>	<b>M2 W</b>	3.97	1.00	5.19	1.93
	<b>M3 W+S</b>	3.79	0.82	5.11	1.85

where  $\widehat{n}_s$  refers to the estimated number of speakers.

## 6. Experimental Results and Discussion

### 6.1. Datasets

Evaluation datasets for our proposed system are limited to English speech corpora since the proposed system relies on the English ASR system at the moment. We report the diarization performance on the following corpora:

- **RT03-CTS** (LDC2007S10): We use the 14-vs-58 dev and eval split provided by the authors in [19]. All the parameters appear on this paper are optimized on the RT03 dev set.
- **CH American English Speech (CHAES)** (LDC97S42) [20]: Note that the CHAES corpus is different from the commonly-used multilingual dataset “NIST SRE 2000 CALLHOME (LDC2001S97)” which is the superset of the CHAES corpus. Within the CHAES corpus, we use two different subsets. (1) **CH-Eval**: Evaluation set from CHAES. (2) **CH-109**: 109 conversations from the CHAES corpus that have 2 speakers only. The CH-109 subset is popularly used such as in [21] when evaluating diarization systems focusing only on the 2 speaker cases.

### 6.2. Evaluation setups

We evaluate the proposed system (M3) with the baseline system configuration (M1) on the two evaluation datasets (CH-109 and CH-Eval) as well as the RT03-CTS dataset. To evaluate the systems in terms of diarization error rate (DER) and speaker error rate (SER), we use the *md-eval* software presented in [22]. The gap between DER and SER originates from the false alarms and missed detections that are caused by SAD. The systems compared in the tables above are configured in the following manners:

- **M1**: This baseline system configuration only exploits  $\mathbf{P}_{ud}$  as  $\mathbf{A}_c$  for spectral clustering (i.e.,  $\mathbf{A}_c = \mathbf{P}_{ud}$ ). This is the general speaker diarization system utilizing acoustic information only in speaker embeddings. The results of this system would contrast how much lexical information can contribute to the speaker clustering process to enhance the overall speaker diarization accuracy in M2 and M3.
- **M2**: This configuration for the proposed system excludes the speaker embedding part for the speaker turn probability estimator in Fig. 2 to show the contribution of lexical information in the speaker turn probability estimation process.
- **M3**: This is the full-blown configuration, as explained throughout this paper.

### 6.3. Evaluation results

The performance of our proposed system is compared to previously published results [19, 21] on the same dataset. However, it should be noted that results in [19] and our proposed system are based on system SAD that is bound to give higher DER than

Table 2: DER (%) on the CHAES dataset.

Number of Speakers		Unknown		Known	
Dataset(Quantity)		CH-Eval(20)		CH-109(109)	
Error Type		DER	SER	DER	SER
Quan <i>et al.</i> [19] System SAD		12.54	5.97	12.48	6.03
Zajc <i>et al.</i> [21] Oracle SAD		-	-	-	7.84
<b>Baseline</b>	<b>M1</b>	7.00	2.94	6.42	2.13
<b>Proposed</b>	<b>M2 W</b>	7.04	2.97	5.96	1.67
	<b>M3 W+S</b>	6.97	2.9	6.03	1.73

the systems based on oracle SAD. On the other hand, the system in [21] uses oracle SAD which makes DER equal to SER.

- **Table1 (RT03-CTS)**: The M3 system improves the performance over M2, but the relative improvements are minimal as compared to the improvements of M2 over M1. This shows that most of the performance gain by the proposed speaker diarization system comes from employing lexical information to the speaker clustering process.
- **Table2 (CH-Eval, CH-109)**: This table compares our proposed speaker diarization system with the recently published results [19, 21] on the CHAES datasets. For a fair comparison, we applied the eigengap analysis based speaker number estimation in Eq. (6) only to the CH-Eval dataset while fixing the number of speakers to 2 in the CH-109 dataset (since CH-109 is the chosen set of the CHAES conversations with only 2 speakers). It is shown in the table that our proposed system (M3) outperforms the previously published results in [19, 21] on both CH-Eval and CH-109. It is worthwhile to mention that the proposed system did not gain the noticeable improvement in the CH-Eval dataset as compared to the baseline configuration (M1). As for the CH-109 dataset, on the other hand, M3 seems to provide a noticeable jump in SER over M1. Given our observation that in the CH-109 evaluation most of the performance improvement from M1 to M3 was from the worst 10 sessions that the baseline system performed poorly on, we presume that the proposed system improves the clustering results on such challenging data.

### 6.4. Discussion

The experimental results show that the baseline system outperforms the previously published results due to the performance of ASPIRE SAD [12] and x-vector [13]. However, our proposed system still improves the competitive baseline system by 36% for RT03-Eval and 19% for CH-109 in terms of SER.

## 7. Conclusions

In this paper, we proposed the speaker diarization system to exploit lexical information from ASR to the speaker clustering process to improve the overall DER. The experimental results showed that the proposed system provides meaningful improvements on both of the CHAES and RT03-CTS datasets outperforming the baseline system which is already competitive against the previously published state-of-the-art results. This supports our claim that lexical information can improve diarization results by incorporating turn probability and word boundaries. Further studies should target the optimal approaches of integrating the adjacency matrices by employing improved search techniques which can improve not only the clustering performance but also the processing speed.

## 8. Acknowledgements

This research was supported in part by NSF, NIH, and DOD.

## 9. References

- [1] D. Liu, D. Kiecza, A. Srivastava, and F. Kubala, "Online speaker adaptation and tracking for real-time speech recognition," in *Ninth European Conference on Speech Communication and Technology*, 2005.
- [2] T. Hain, L. Burget, J. Dines, P. N. Garner, F. Grézl, A. El Hanani, M. Huijbregts, M. Karafiat, M. Lincoln, and V. Wan, "Transcribing meetings with the AMIDA systems," *IEEE Trans. Audio Speech Lang. Process.*, vol. 20, no. 2, pp. 486–498, 2012.
- [3] P. Georgiou, M. Black, and S. Narayanan, "Behavioral signal processing for understanding (distressed) dyadic interactions: Some recent developments," in *Joint ACM Workshop on Human Gesture and Behavior Understanding*, 2011, pp. 7–12.
- [4] S. Narayanan and P. Georgiou, "Behavioral signal processing: Deriving human behavioral informatics from speech and language," *Proceedings of the IEEE*, vol. PP, no. 99, pp. 1–31, 2013.
- [5] M. Gales, "Maximum likelihood linear transformations for HMM-based speech recognition," *Comp. Speech and Lang.*, vol. 12, pp. 75–98, 1997.
- [6] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2011.
- [7] D. Dimitriadis and P. Fousek, "Developing on-line speaker diarization system," in *Interspeech*, 2017, pp. 2739–2743.
- [8] T. J. Park and P. Georgiou, "Multimodal speaker segmentation and diarization using lexical and acoustic cues via sequence to sequence neural networks," in *Interspeech*, 2018.
- [9] P. Cerva, J. Silovsky, J. Zdansky, J. Nouza, and L. Seps, "Speaker-adaptive speech recognition using speaker diarization for improved transcription of large spoken archives," *Speech Communication*, vol. 55, no. 10, pp. 1033–1046, 2013.
- [10] M. À. India Massana, J. A. Rodríguez Fonollosa, and F. J. Hernandez Pericás, "LSTM neural network-based speaker segmentation using acoustic and language modelling," in *INTERSPEECH 2017: 20-24 August 2017: Stockholm*. International Speech Communication Association (ISCA), 2017, pp. 2834–2838.
- [11] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [12] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz *et al.*, "The kaldi speech recognition toolkit," *IEEE Signal Processing Society*, Tech. Rep., 2011.
- [13] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, "X-vectors: Robust dnn embeddings for speaker recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018.
- [14] H. Bredin, "Tristounet: Triplet loss for speaker turn embedding," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 5430–5434.
- [15] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [16] C. Cieri, D. Miller, and K. Walker, "Fisher English training speech parts 1 and 2," *Philadelphia: Linguistic Data Consortium*, 2004.
- [17] J. J. Godfrey and E. Holliman, "Switchboard-1 release 2," *Linguistic Data Consortium, Philadelphia*, vol. 926, p. 927, 1997.
- [18] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [19] Q. Wang, C. Downey, L. Wan, P. A. Mansfield, and I. L. Moreno, "Speaker diarization with LSTM," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5239–5243.
- [20] L. D. Consortium *et al.*, "CALLHOME American English speech," 1997, 1997.
- [21] Z. Zajc, M. Hrz, and L. Mller, "Speaker diarization using convolutional neural network for statistics accumulation refinement," in *Proc. Interspeech 2017*, 2017, pp. 3562–3566. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2017-51>
- [22] J. G. Fiscus, J. Ajot, M. Michel, and J. S. Garofolo, "The Rich Transcription 2006 spring meeting recognition evaluation," in *International Workshop on Machine Learning for Multimodal Interaction*. Springer, 2006, pp. 309–322.