

Speaker-Conditional Chain Model for Speech Separation and Extraction

Jing Shi^{1,2}, Jiaming Xu¹, Yusuke Fujita³, Shinji Watanabe², Bo Xu¹

¹Institute of Automation, Chinese Academy of Sciences (CASIA)

²Center for Language and Speech Processing, Johns Hopkins University

³Hitachi, Ltd. Research & Development Group

shijing2014@ia.ac.cn

Abstract

Speech separation has been extensively explored to tackle the cocktail party problem. However, these studies are still far from having enough generalization capabilities for real scenarios. In this work, we raise a common strategy named Speaker-Conditional Chain Model to process complex speech recordings. In the proposed method, our model first infers the identities of variable numbers of speakers from the observation based on a sequence-to-sequence model. Then, it takes the information from the inferred speakers as conditions to extract their speech sources. With the predicted speaker information from whole observation, our model is helpful to solve the problem of conventional speech separation and speaker extraction for multi-round long recordings. The experiments from standard fully-overlapped speech separation benchmarks show comparable results with prior studies, while our proposed model gets better adaptability for multi-round long recordings.

Index Terms: speech separation, speaker extraction, cocktail party problem

1. Introduction

Human interactions are often in a broad range of complex auditory scenes, consisting of several speech sources from different speakers and various noises. This complexity poses challenges for many speech technologies, because they usually assume one or zero speaker to be active at the same time [1]. To tackle these challenging scenes, many techniques have been studied.

Speech separation aims at isolating individual speaker's voices from a recording with overlapped speech [2–8]. With the separation results, both the speech intelligibility for human listening and speech recognition accuracy could be improved [9]. Different from the separation task, speaker extraction makes use of additional information to distinguish a target speaker from other participating speakers [10–13]. Besides, speech denoising [14, 15] and speaker diarization [16, 17] tasks have also been studied for solving the problem of complex acoustic scenes.

Although many works have been proposed towards each task mentioned above, the processing of natural recordings is still challenging. Overall, these tasks are designed to accomplish one particular problem, which has assumptions that do not hold in complex speech recordings. For instance, speech separation was heavily explored with pre-segmented audio samples with a length of several seconds (less than 10 seconds), which makes it difficult to form reasonable results for long recordings. Because most existing separation methods only output a fixed number of speech sources with agnostic order, and it is unable to process the variable number of speakers and the relation of the orders between different segments. Similarly, the speaker diarization bypassed the overlapped part before. Recently, the emergence of EEND approaches [16, 17] could fix the prob-

lem of overlapped speech parts to some extent. However, the diarization results seem an intermediate product without the extraction of each speaker, especially for the overlapped parts.

To address these limitations, we believe that integrating speaker information (used in aim speaker extraction, speaker diarization) into speaker-independent tasks (e.g., speech separation, speech denoising and even speech recognition) will help broaden the application of these techniques towards real scenes. To be specific, we reconstruct the speech separation/extraction task with the strategy over probabilistic chain rule by importing the conditional probability based on speaker information. In practice, our model automatically infers the information of speakers' identities and then takes it as condition to extract speech sources. The speaker information here is some learned hidden representation related to the speaker's identity, which makes it also suitable for open speaker tasks. We believe this design actually better meets the expectation about an intelligent front-end speech processing pipeline. Because users usually want to get the information about not only the extracted clean speech sources but also which ones speak what.

In this work, we propose our Speaker-Conditional Chain Model (SCCM) to separate the speech sources of different speakers with overlapped speech. Meanwhile, the proposed method can handle a long recording with multiple rounds of utterances spoken by different speakers. Based on this model, we verified its effectiveness in getting both the identity information of each speaker and the extracted speech sources of them.

The contributions of this paper span the following aspects: (1) we built a common chain model for the processing of speech with one or more speakers. Through the inference-to-extraction pipeline, our model solves the problem about the variable and even unknown number of speakers; (2) with the same architecture, our model shows a comparative performance with the base model, while we could additionally offer accurate speaker identity information for further downstream usage; (3) we proved the effectiveness of this design for both short overlapped segments and long recordings with multi-round conversations, (4) we analyze the advantages and drawbacks of this model. Our demo video and Supplementary Material are available at <https://shincling.github.io/>.

2. Related work

2.1. Speech separation

As the core part of the cocktail party problem [18], speech separation gains much attention recently. The common design of this task is to disentangle fully overlapped speech signals from a given short mixture (less than 10 seconds) with a fixed number of speakers. Under this design, from spectrogram-based methods [4–6, 19, 20] to time-domain methods [21–23], speaker-agnostic separation approaches have been intensively studied. However, with the steady improvement in performance, most

existing approaches might overfit the fully overlapped audio data, which is far from the natural situation with less than 20% overlap ratio in conversations [24]. Besides, most existing separation models should know the number of speakers in advance and could only tackle the data with the same number of speakers [25]. These constraints further limit their application to real scenes, while our proposed SCCM can provide a solution to the above sparse overlap and unknown speaker number issues. A similar idea with recurrent selective attention networks [26] has been proposed before to tackle the variable number of speakers in separation. However, this model performs with residual spectrograms without leveraging the time-domain methods. And their uPIT [19] based training is hard to process a long recording, due to the speaker tracing problem raised when chunking the long recording into short segments.

2.2. Speaker extraction

Another task related to our model is the speaker extraction [10–13]. The idea of speaker extraction is to provide a reference from a speaker, and then use such reference to direct the attention to the specified speaker. The reference may be taken from different characteristic binding with the specific speaker, such as voiceprint, location, onset/offset information, and even visual representation [27]. The speaker extraction technique is particularly useful when the system is expected to respond to a specific target speaker. However, for a meeting or conversation with multiple speakers, the demand for additional references makes it inconvenient. In our work, the reference could be directly inferred from the original recordings, which shows an advantage when the complete analysis of each speaker is needed.

3. Speaker-conditional chain model

This section describes our Speaker-Conditional Chain Model (SCCM). As illustrated in Figure 1, the chain here refers to a pipeline through two sequential components: speaker inference and speech extraction. These models are integrated based on a joint probability formulation, which will be described in Section 3.1. Speaker identities play an important role in our strategy. The speaker inference module aims to predict the possible speaker identities and the corresponding embedding vectors. The speech extraction module takes each embedding from the speaker inference module as the query to disentangle the corresponding source audio from the input recording.

This design will bring several advantages. First, the possible speakers are inferred by a sequence-to-sequence model with an end-of-sequence label, which easily handles variable and unknown numbers of speakers. Second, the inference part is based on a self-attention network, which utilizes the full context information in a recording to form a speaker embedding. This avoids the calculation inefficiency problem in some clustering-based models [4, 5, 9], which needs an iterative k-means algorithm in each frame. Third, the information about each speaker will make it suitable for our model to some further applications in speaker diarization or speaker tracking.

3.1. Problem setting and formulation

Assume there is a training dataset with a set of speaker identities \mathcal{Y} with $|\mathcal{Y}| = N$ known distinct speakers in total. In a T -length segment of waveform observation $O \in \mathbb{R}^T$, there are I different speakers $Y = (y_1, \dots, y_i, \dots, y_I)$. Each speaker y_i ¹

¹Although $y_i \in \mathcal{Y}$ during training, potentially $y_i \notin \mathcal{Y}$ during inference in the open speaker task, where the system could still provide a meaningful speaker embedding vector for downstream applications.

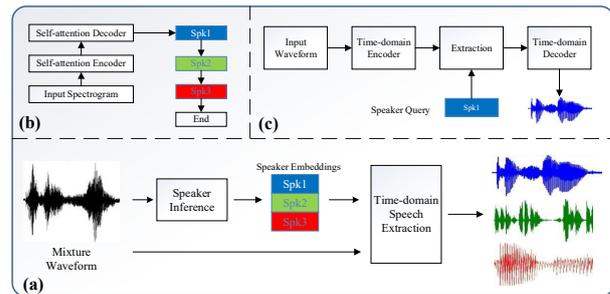


Figure 1: The framework of the proposed Speaker-Conditional Chain Model (SCCM). (a) shows the whole strategy of our proposed SCCM; (b) is the module of **speaker inference**, which predicts the speaker identities and corresponding embeddings. (c) refers to the time-domain **speech extraction** module. This module takes the each inferred information from (b) respectively to conduct a conditional extraction.

has the corresponding speech source $s_i \in \mathbb{R}^T$ to form the set of sources $S = (s_1, \dots, s_i, \dots, s_I)$. The basic formulation of our strategy is to estimate the joint probability of speaker labels and corresponding sources, i.e., $p(S, Y|O)$. This is factorized with speaker inference probability $p(Y|O)$ and speech extraction probability $p(S|Y, O)$ as follows:

$$p(S, Y|O) = p(S|Y, O)p(Y|O). \quad (1)$$

We further factorize each probability distribution based on the probabilistic chain rule.

The speaker inference probability $p(Y|O)$ in Eq. (1) recursively predicts variable numbers of speaker identities as follows:

$$p(Y|O) = \prod_i p(y_i|O, y_{i-1}, \dots, y_1). \quad (2)$$

We adopt a *sequence-to-sequence model* based on self-attention transformer [28], as illustrated in Figure 1(b). The network architecture of $p(Y|O)$ will be discussed in Section 3.2.

The speech extraction probability $p(S|Y, O)$ in Eq. (1) is also factorized by using the probabilistic chain rule and the conditional independence assumption, as follows:

$$p(S|Y, O) = \prod_i p(s_i|y_i, y_{i-1}, \dots, y_1, O) = \prod_i p(s_i|y_i, O). \quad (3)$$

As illustrated in Figure 1(c), our speech extraction module takes the speaker identity y_i , which is predicted from the speaker inference module $p(Y|O)$ in Eq. (2), to conduct a conditional extraction. Every speaker information here serves as the condition to guide the following extraction. For multi-round long recordings, the speaker information will be formed as global information from the whole observation to track the specific speaker. The network architecture of $p(s_i|y_i, O)$ will be discussed in Section 3.3.

3.2. Speaker inference module

In the speaker inference part, we build a model to simulate the probability $p(Y|O)$ in Eq. (1) and (2). We adopt a self-attention based transformer [28] architecture as the encoder-decoder structure. In this part, we take the observation spectrogram (Short-Time Fourier Transform (STFT) coefficients) as an input. The reason we do not use the time-domain approach here is to avoid excessive computation complexities which may

consume too much GPU memory to train the model, especially with inputs of long recordings.

In detail, for a given spectrogram \mathbf{X} containing \tilde{T} frames and F frequency bins, it is viewed as a sequence of frames. For the encoder part, we use the Transformer Encoder as follows:

$$\mathbf{E}_0 = \text{Linear}^{(F \rightarrow D)}(\mathbf{X}) \in \mathbb{R}^{D \times \tilde{T}}, \quad (4)$$

$$\mathbf{E}_m = \text{Encoder}(\mathbf{E}_{m-1}) \in \mathbb{R}^{D \times \tilde{T}} \quad (1 \leq m \leq M), \quad (5)$$

where, $\text{Linear}^{(F \rightarrow D)}()$ is a linear projection that maps F -dimensional vector to D -dimensional vector for each column of the input matrix. $\text{Encoder}()$ is the Transformer Encoder block that contains multi-head self-attention layer, position-wise feed-forward layer, and residual connections. By stacking the encoder M times, $\mathbf{E}_M \in \mathbb{R}^{D \times \tilde{T}}$ is an output of the encoder part.

For the decoder part, the neural network outputs probability distribution \mathbf{z}_i for the i -th speaker, calculated as follows:

$$\mathbf{j}_i = \text{Linear}^{(1 \rightarrow D)}(\mathbf{i}) \in \mathbb{R}^D, \quad (6)$$

$$\mathbf{h}_i = \text{Decoder}(\mathbf{E}_M, \mathbf{h}_{i-1}, \mathbf{j}_i) \in \mathbb{R}^D, \quad (7)$$

$$\mathbf{z}_i = \text{Softmax}(\text{Linear}^{(D \rightarrow \tilde{N})}(\mathbf{h}_i)) \in \mathbb{R}^{\tilde{N}+1}, \quad (8)$$

where \mathbf{j}_i is the positional encoding in each step to predict the speaker. $\text{Decoder}()$ is the Transformer Decoder block, which takes the states from the output of encoder and the hidden state from the previous step to output the speakers embedding \mathbf{h}_i at this step. Finally, a linear projection with a softmax produces a $(\tilde{N} + 1)$ -dimensional vector \mathbf{z}_i as the network output, where \mathbf{z}_i is the i -th predicted probability distribution over the union of speaker set \mathcal{Y} and the additional end-of-sequence label $\langle \text{EOS} \rangle$, i.e., $y^* \in \{\mathcal{Y}, \langle \text{EOS} \rangle\}$.

3.3. Speech extraction module

For the speech extraction module, each speaker channel $p(s_i|y_i, O)$ will be processed independently, as formed in Eq. (3). This part takes each inferred speaker embedding \mathbf{h}_i predicted in Eq. (7) instead of identity y_i , and the raw waveform O as input to produce the corresponding clean signal \hat{s}_i :

$$\hat{s}_i = \text{Extractor}(O, \mathbf{h}_i) \in \mathbb{R}^T, \quad (9)$$

where, $\text{Extractor}()$ takes a similar architecture with time-domain speech separation methods from the Conv-TasNet [22]. The difference lies in that we will output one channel towards each speaker embedding rather than separate several sources together. To be specific, at the end of the separator module in [22], we will concatenate the \mathbf{h}_i with each frame of the output features. Then, a single channel $1 \times 1 - \text{conv}$ operation is conducted towards this speaker, rather than multi-channel (as the number of speakers in this mixture). Besides this simple fusion approach, we have tested several different methods to integrate the condition vector \mathbf{h}_i into the model. For example, to concatenate it at the beginning of the separator, or use the similar method in [9] with FiLM [29] in each block in TasNet’s separator. However, we found both of the other methods cause severe overfitting.

3.4. Training targets

Our whole model is end-to-end, with the loss \mathcal{L} , which corresponds to optimize the joint probability $p(S, Y|O)$ in Eq. (1). \mathcal{L} is calculated from both the cross-entropy loss \mathcal{L}_c , which corresponds to deal with speaker inference $p(Y|O)$ in Section 3.2,

and the source reconstruction loss (SI-SNR) \mathcal{L}_r , which corresponds to deal with speech extraction $p(S|Y, O)$ in a non-probabilistic manner in Section 3.3. One critical problem in training SCCM is to decide the order of the inferred speakers. For one possible permutation ρ , the speakers list Y and the speech sources S will be re-ordered synchronously as follows:

$$Y_\rho = (y_1^*, y_2^*, \dots, y_I^*), \forall y_i^* \xrightarrow{\rho} y \in Y, \quad (10)$$

$$S_\rho = (s_1^*, s_2^*, \dots, s_I^*), \forall s_i^* \xrightarrow{\rho} s \in S. \quad (11)$$

Some former works have shown that the seq2seq structure helps to improve the accuracy in the inference module by setting a fixed order in training [30]. We compared several options to use a random fixed order or use the order defined by the energy in the spectrogram (observed well in [31]). But we found the order decided by the model itself gets better performance in practice. Therefore, we take the best permutation θ with least reconstruction error in the extraction part as the order to train the inference part as follows:

$$\theta = \underset{\rho \sim \text{Perms}}{\text{argmin}} \mathcal{L}_r(\hat{\mathbf{S}}, S_\rho), \quad (12)$$

$$\mathcal{L} = \mathcal{L}_r(\hat{\mathbf{S}}, S_\theta) + \alpha \mathcal{L}_c(\mathbf{Z}, Y_\theta), \quad (13)$$

where we use $\alpha = 50$ in all our experiments.

4. Experiments

As a generalized framework to tackle the problem of extracting speech sources of all speakers, we tested the effectiveness of SCCM with different tasks. Besides the signal reconstruction quality (e.g., SDRi, SI-SNRi) used in speech separation task, we also verified the performance over speaker identification and speech recognition. In our experiments, all data are re-sampled to 8 kHz. For the speaker inference module, the magnitude spectra are used as the input feature, computed from STFT with 32 ms window length, 8 ms hop size, and the sine window. More detailed configuration of the proposed architecture could be seen in Section A.1 of our [Supplementary Material](#)².

4.1. Speech separation for overlapped speech

First, we evaluated our method on fully-overlapped speech mixtures from the Wall Street Journal (WSJ0) corpus. The WSJ0-2mix and 3mix datasets are the benchmarks designed for speech separation in [4]. In the validation set, we used the so-called Closed Conditions (CC) in [4, 5], where the speakers are all from the training set. As a contrast, for the evaluation set, we use Open Condition (OC), which provides unknown speakers. For the separation performance, we compare our results with the TasNet, which is our base model described in Section 3.3, without changing any hyper-parameter. Table 1 listed the speech separation performance over the different training sets.

Table 1 shows that our SCCM got slightly worse performance than the base model in OC with the same architecture and training dataset. However, unlike the fixed-speaker-number speech separation method, SCCM could be trained and tested in the variable number of speakers with a single model thanks to our speaker-conditional strategy with the sequence-to-sequence model. As we expect, the training with both WSJ0-2mix and WSJ0-3mix datasets got better performance than the training with each dataset in close condition. Although we did not achieve obvious improvement in the OC case, with the careful tuning based on the cascading technique (the similar meth-

²<https://drive.google.com/open?id=1aQJy465dLHaWPdMqG-BgJAgYEG70q7as>

Table 1: *Speech separation performance (SI-SNRi) for the benchmark datasets with overlapped speech.*

Models	Training Dataset	SI-SNRi CC		SI-SNRi OC	
		WSJ0-2mix	WSJ0-3mix	WSJ0-2mix	WSJ0-3mix
TasNet	WSJ0-2mix	-	-	14.6	-
	WSJ0-3mix	-	-	-	11.6
SCCM	WSJ0-2mix	15.4	-	14.5	-
	WSJ0-3mix	-	11.9	-	11.4
SCCM+	both	16.4	12.1	14.3	11.3
	both	17.7	13.4	15.4	12.5

Table 2: *Speaker inference performance of SCCM.*

Training Dataset	F1 scores in Validset (CC)		Speaker counting accuracy in Testset (OC)	
	WSJ0-2mix	WSJ0-3mix	WSJ0-2mix	WSJ0-3mix
WSJ0-2mix	89.2	-	99.7	-
WSJ0-3mix	-	-	-	98.9
both	90.4	75.5	96.8	94.5

ods used in [19]), the separation performance gets a notable improvement, which also exceeds the base model. For the SCCM+ model, we use the extracted speech source, along with the raw observation, as input to go through another extraction module (TasNet). With this cascading method, the details of the extracted source get further optimized, which may fix the ambiguity caused by the independence assumption in Eq. (3).

Also, as the former node in the chain, the ability to predict the correct speakers or get the distinct and informative embeddings is quite crucial. Table 2 shows the performance of the speaker inference module, as discussed in Section 3.2. For the CC, micro-F1 is calculated to evaluate the correctness of the predicted speakers. For the OC, we use the speaker counting accuracy to measure the speaker inference module, which guarantees the success of the subsequent speech extraction module. From the results, we could see that the speaker inference module in SCCM could reasonably infer the correct speaker identity in CC and the correct number of speakers in OC.

It should be mentioned that the number of speakers in training data (N in Section 3.1) with WSJ0-2mix and 3mix is 101, much smaller than the number in a standard speaker recognition task (e.g., 1,211 in VoxCeleb1 [32]). We infer that this limited number somewhat limits the performance of the speaker inference part and the following extraction module, especially for the open condition. Besides, compared with the state-of-the-art speaker recognition methods, our model takes the overlapped speech as input, which also brings more complexity.

4.2. Extraction performance for multi-round recordings

As mentioned before, the natural conversions in real scenes usually get multi-round utterances from several speakers. And the ratio of overlapped speech is less than 20% in general. For the conventional speech separation methods, there exists a problem with the consistent order of several speakers in different parts in a relatively long recording, especially when the dominant speaker changes [9]. To validate this, we extend each mixture in the standard WSJ0-mix to multiple rounds. In detail (seen in Algorithm 1 and Section A.2 in [Supplementary Material](#)), we take the list of the original mixtures from WSJ0-2mix and sample several additional utterances from the provided speakers. After getting the sources from different speakers, the long recording will be formed by concatenating the sources one by one. The beginning of the following source gets a random shift around the end of the former one, making it similar to a natural conversation with an overlap-ratio around 15%.

Without any change in our model, we could directly train our SCCM on the synthetic multi-round data. It should be mentioned that our speaker inference module takes the whole spectrogram as an input. In contrast, the speech extraction module takes a random segment with 4 seconds from the long recording

Table 3: *Extraction performance for multi-round recordings.*

	Valid SI-SNRi		Test SI-SNRi	
TasNet	14.2		11.5	
SCCM	17.5		13.7	
	<5dB	>5dB	<5dB	>5dB
TasNet	17.0%	83.0%	33.6%	66.4%
SCCM	12.6%	87.4%	26.8%	73.2%

Table 4: *WERs for utterance-wise evaluation over the single-channel LibriCSS dataset with clean mixtures. OS: 0% overlap with short inter-utterance silence (0.1-0.5 s). OL: 0% overlap with long inter-utterance silence (2.9-3.0 s).*

System	Overlap ratio in %					
	OS	OL	10	20	30	40
No separation	2.7	3.0	11.9	20.4	30.2	43.0
Single-channel SCCM	9.5	9.4	6.5	9.3	11.9	13.9

to avoid the problem with out-of-memory. Table 3 shows the performance difference compared with the base model. Both valid set and test set are fixed with four rounds of conversations with an average length of 10 seconds. As we expect, the results show that SCCM stays more stable than the baseline model with multi-round recordings. To further understand the model, we observed the attention status of the Decoder in Eq. (7). We find the attention of the inference reflects the speaker’s activities at different parts within a recording. More details and visualization could be viewed in Section A.3 in the [Supplementary Material](#).

4.3. Speech recognition in continuous speech separation

To further validate the downstream application, we conducted the speech recognition in the recently proposed continuous speech separation dataset [33]. LibriCSS is derived from LibriSpeech [34] by concatenating the corpus utterances to simulate conversations. In line with the utterance-wise evaluation in LibriCSS, we directly use our trained model from the former multi-round task to test the recognition performance. The original raw recordings in LibriCSS are from far-field scenes with noise and reverberation, which is inconsistent with ours. So we use the single-channel clean mixtures and convert to 8 kHz to separate them. Moreover, we use the trained model from the Espnet’s [35] LibriSpeech recipe to recognize each utterance. Table 4 shows the WER results in this dataset.

We observed that (1) the results show a similar tendency with the provided baseline model in LibriCSS [33]. (2) With the increase of overlap ratio, the performance on the original clean mixture becomes much worse, while our model stays a low level of WER. (3) Because the training data of our model comes from the situation of multi speakers, the performance on the no-overlapped segments becomes worse. And we think this could be avoided by adding some single speaker’s segments in the training set.

5. Conclusions

We introduced the Speaker-conditional chain model as a common framework to process audio recordings with multiple speakers. Our model could be applied to tackle the separation problem towards fully-overlapped speech with variable and unknown number of speakers. Meanwhile, multi-round long audio recordings in natural scenes can also be modeled and extracted effectively using this method. Experimental results showed the effectiveness and good adaptability of the proposed model. Our following work will extend this model to the real scenes with noisy and reverberant multi-channel recordings. We would also like to explore the factors to improve the generalization ability of this approach, like the introduction of more speakers or changes in the network and training objectives.

6. References

- [1] R. Haeb-Umbach, S. Watanabe, T. Nakatani, M. Bacchiani, B. Hoffmeister, M. L. Seltzer, H. Zen, and M. Souden, "Speech processing for digital home assistants: Combining signal processing with deep-learning techniques," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 111–124, 2019.
- [2] P. Huang, M. Kim, M. Hasegawajohnson, and P. Smaragdis, "Deep learning for monaural speech separation," in *ICASSP*, 2014, pp. 1562–1566.
- [3] G. Wang, C. Hsu, and J. Chien, "Discriminative deep recurrent neural networks for monaural speech separation," in *ICASSP*, 2016, pp. 2544–2548.
- [4] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, "Deep clustering: Discriminative embeddings for segmentation and separation," in *ICASSP*, 2016, pp. 31–35.
- [5] Y. Isik, J. L. Roux, Z. Chen, S. Watanabe, and J. R. Hershey, "Single-channel multi-speaker separation using deep clustering," in *INTERSPEECH*, 2016.
- [6] D. Yu, M. Kolbæk, Z.-H. Tan, and J. Jensen, "Permutation invariant training of deep models for speaker-independent multi-talker speech separation," in *ICASSP*, 2017, pp. 241–245.
- [7] Z. Chen, Y. Luo, and N. Mesgarani, "Deep attractor network for single-microphone speaker separation," in *ICASSP*, 2017, pp. 246–250.
- [8] L. Drude, T. von Neumann, and R. Haeb-Umbach, "Deep attractor networks for speaker re-identification and blind source separation," in *ICASSP*, 2018, pp. 11–15.
- [9] N. Zeghidour and D. Grangier, "Wavesplit: End-to-end speech separation by speaker clustering," *arXiv preprint arXiv:2002.08933*, 2020.
- [10] M. Delcroix, K. Zmolikova, K. Kinoshita, A. Ogawa, and T. Nakatani, "Single channel target speaker extraction and recognition with speaker beam," in *ICASSP*, 2018, pp. 5554–5558.
- [11] Q. Wang, H. Muckenhirn, K. Wilson, P. Sridhar, Z. Wu, J. R. Hershey, R. A. Saurous, R. J. Weiss, Y. Jia, and I. L. Moreno, "Voice-Filter: Targeted Voice Separation by Speaker-Conditioned Spectrogram Masking," in *INTERSPEECH*, 2019, pp. 2728–2732.
- [12] J. Xu, J. Shi, G. Liu, X. Chen, and B. Xu, "Modeling attention and memory for auditory selection in a cocktail party environment," in *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI)*, 2018, pp. 2564–2571.
- [13] C. Xu, W. Rao, E. S. Chng, and H. Li, "Spex: Multi-scale time domain speaker extraction network," *arXiv preprint arXiv:2004.08326*, 2020.
- [14] C. Donahue, B. Li, and R. Prabhavalkar, "Exploring speech enhancement with generative adversarial networks for robust speech recognition," in *ICASSP*, 2018, pp. 5024–5028.
- [15] D. Rethage, J. Pons, and X. Serra, "A wavenet for speech denoising," in *ICASSP*, 2018, pp. 5069–5073.
- [16] Y. Fujita, N. Kanda, S. Horiguchi, K. Nagamatsu, and S. Watanabe, "End-to-end neural speaker diarization with permutation-free objectives," in *INTERSPEECH*, 2019, pp. 4300–4304.
- [17] Y. Fujita, N. Kanda, S. Horiguchi, Y. Xue, K. Nagamatsu, and S. Watanabe, "End-to-end neural speaker diarization with self-attention," in *ASRU*, 2019.
- [18] E. C. Cherry, "Some experiments on the recognition of speech, with one and with two ears," *Journal of the Acoustical Society of America*, vol. 25, no. 5, pp. 975–979, 1953.
- [19] M. Kolbaek, D. Yu, Z. H. Tan, J. Jensen, M. Kolbaek, D. Yu, Z. H. Tan, and J. Jensen, "Multitalker speech separation with utterance-level permutation invariant training of deep recurrent neural networks," *IEEE/ACM Transactions on Audio Speech and Language Processing*, vol. 25, no. 10, pp. 1901–1913, 2017.
- [20] Y. Luo, Z. Chen, and N. Mesgarani, "Speaker-independent speech separation with deep attractor network," *IEEE/ACM Transactions on Audio Speech and Language Processing*, vol. 26, no. 4, pp. 787–796, 2018.
- [21] Y. Luo and N. Mesgarani, "Real-time single-channel dereverberation and separation with time-domain audio separation network," in *INTERSPEECH*, 2018, pp. 342–346.
- [22] —, "Tasnet: time-domain audio separation network for real-time, single-channel speech separation," in *ICASSP*, 2018, pp. 696–700.
- [23] Y. Luo, Z. Chen, and T. Yoshioka, "Dual-path rnn: efficient long sequence modeling for time-domain single-channel speech separation," in *ICASSP*, 2020, pp. 46–50.
- [24] Ö. Çetin and E. Shriberg, "Analysis of overlaps in meetings by dialog factors, hot spots, speakers, and collection site: Insights for automatic speech recognition," in *Ninth international conference on spoken language processing*, 2006.
- [25] J. Shi, J. Xu, G. Liu, and B. Xu, "Listen, think and listen again: Capturing top-down auditory attention for speaker-independent speech separation," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI)*, 2018.
- [26] K. Kinoshita, L. Drude, M. Delcroix, and T. Nakatani, "Listening to each speaker one by one with recurrent selective hearing networks," in *ICASSP*, 2018, pp. 5064–5068.
- [27] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, and M. Rubinstein, "Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation," *Acm Transactions on Graphics*, vol. 37, no. 4, pp. 1–11, 2018.
- [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [29] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "Film: Visual reasoning with a general conditioning layer," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [30] J. Shi, J. Xu, and B. Xu, "Which ones are speaking? speaker-inferred model for multi-talker speech separation," *INTERSPEECH*, pp. 4609–4613, 2019.
- [31] C. Weng, D. Yu, M. L. Seltzer, and J. Droppo, "Deep neural networks for single-channel multi-talker speech recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 10, pp. 1670–1679, 2015.
- [32] A. Nagrani, J. S. Chung, and A. Zisserman, "Voxceleb: a large-scale speaker identification dataset," in *INTERSPEECH*, 2017.
- [33] Z. Chen, T. Yoshioka, L. Lu, T. Zhou, Z. Meng, Y. Luo, J. Wu, and J. Li, "Continuous speech separation: dataset and analysis," *arXiv preprint arXiv:2001.11482*, 2020.
- [34] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *ICASSP*, 2015, pp. 5206–5210.
- [35] S. Watanabe, T. Hori, S. Karita, T. Hayashi, J. Nishitoba, Y. Unno, N. E. Y. Soplin, J. Heymann, M. Wiesner, N. Chen *et al.*, "Espnet: End-to-end speech processing toolkit," in *INTERSPEECH*, 2018.
- [36] Y. Luo and N. Mesgarani, "Conv-tasnet: Surpassing ideal time-frequency magnitude masking for speech separation," *IEEE/ACM transactions on audio, speech, and language processing*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [37] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *Computer Science*, 2014.
- [38] Y. Kim, C. Denton, L. Hoang, and A. M. Rush, "Structured attention networks," in *5th International Conference on Learning Representations ICLR, Toulon, France, April 24-26, 2017*.

Algorithm 1: Multi-round recordings simulation.

Input: \mathcal{Y} // Speaker lists set in WSJ0-mix
 N_{spk} // number of speakers per mixture
 k_{min}, k_{max} // Min&Max number of rounds per mixture
 β // random shift range
 R // SRN range

Output: $\mathbb{O} \leftarrow \{o\}$ // Simulated list of mixtures

```
1 forall  $Y \in \mathcal{Y}$  do
2    $o \leftarrow \phi$  // initial mixture signal
3    $t \leftarrow 0$  // beginning of one mixture
4   for  $k \in [k_{min}, k_{max}]$  do
5     for  $y \in Y$  do
6       // each speaker in one mixture
7       Sample one audio  $s$  towards speaker  $y$ 
8       Sample SNR  $r$  from the given range  $R$ 
9        $s = s \times 10^{\frac{r}{20}}$  // scale with SNR
10       $o[t : ].add(s)$  // extend the mixture around the end
11       $t = \text{length}(o)$ 
12       $t = t + \text{random}(-\beta, +\beta)$ 
13    $\mathbb{O}.append(o)$ 
```

A. Supplementary Material

A.1. Model details

For the inference module, we used self-attention based encoder-decoder architecture to predict several possible speakers. For both the encoder and decoder, we used one encoder blocks with 512 attention units containing eight heads ($M = 1, d_{model} = 512, H = 8$). The size of dimension used in key and value is 64 ($d_k = 64, d_v = 64$). We used 2048 internal units ($d_{ff} = 2048$) in a position-wise feed-forward layer. And, we used the Adam optimizer with the learning rate decayed by a factor of 2×10^{-1} after every 20 epochs. We tested several different configuration in the model architecture, we found that the large number of layers (above 4) resulted in unconvrgent training. And the configuration with $M = 2$ shows similar results with $M = 1$.

Different from the original transformer model, we did not feed the output embeddings offset by one position to the next step in decoder. Instead, position i is embedded with a linear layer to \mathbf{j}_i (as shown in Eq. (6)) to serve as input at each step. This is to ensure the decoding process can be done without knowing the order of the true speakers, and the order will be decided after the following extraction module by choosing the best permutation with the \mathcal{L}_r .

For the extraction module, we used the original configuration from Conv-TasNet [22] with $N = 256, L = 20, B = 256, H = 512, P = 3, X = 8, R = 4$. Also, we noticed the update of the base model in extraction could further improve the performance like the same tendency in [23, 36]. In this paper, we mainly focus the relative performance over the original TasNet.

For the training strategy, we set a large ratio α in Eq. (13) to balance the \mathcal{L}_c and \mathcal{L}_r , which get a large difference in their ranges. To be specific, with training continues, the cross-entropy criterion \mathcal{L}_c tends to a small positive number close to zero, while the non-probabilistic \mathcal{L}_r changes from positive to almost -20 because of the negative SI-SNR loss definition. Therefore, we set $\alpha = 50$ to keep a reasonable balance between these two factors. Besides, in practice, we found that the extraction

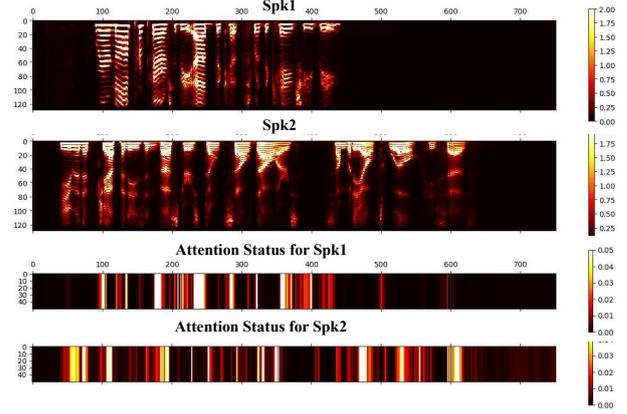


Figure 2: Visualization of one sample of the learned attention status in speaker inference module for overlapped speech in WSJ0-2mix.

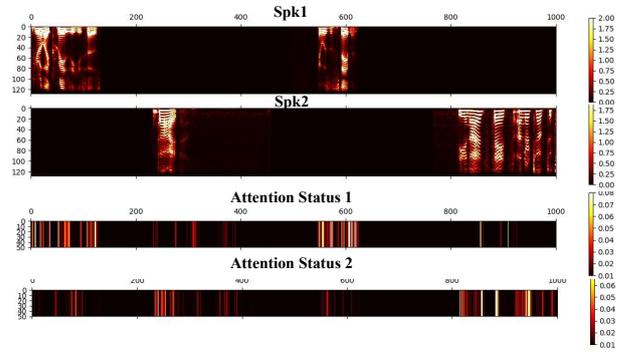


Figure 3: Visualization of one sample of the learned attention status in 2 rounds of utterances.

module takes much more time to converge than the speaker inference module. To avoid the overfitting, the speaker inference module is early-stopped based on the \mathcal{L}_c in validation set, which the extraction module will continue until converged.

A.2. Simulation of WSJ0-mix multi-round recordings

For the multi-round mixtures mentioned in Section 4.2, we simulated them by Algorithm 1. The algorithm is to simply simulate the natural conversations with several parts of overlapped part.

A.3. Attention status

Attention mechanisms have become an integral part of compelling sequence modeling and transduction models in various tasks, allowing modeling of dependencies without regard to their distance in the input or output sequences [28, 37, 38]. For the speech related tasks, the vocal characteristics from one specific speaker stay stable in a short segment and a long conversation. Based on these, we use the self-attention based model in our inference part to utilize the relation between different frames from the same speaker. Therefore, the attention status could be used to check the specific process to find the possible speakers. As shown in Figure 2, we visualized one example

from WSJ0-2mix test set about the real spectrograms of the two speakers and the corresponding attention status towards them. The attention status is from the multi-head self-attention block in the decoder, and we added the weights from each head to form the attention status $\in \mathbb{R}^{1 \times \tilde{T}}$.

As we expect, the attention status shows significant consistency with the real spectrogram. In particular, the attention tends to focus on the frame with larger difference. This is to say, if one speaker gets dominant in some frames, then the attention of this one tends to place emphasis on these dominant frames. Similarly, the attention from multi-round mixture also shows the consistency for one speaker in the whole audio, which could be taken as the implicit speech activity outputted by speaker diarization task.