

L-SPEX: LOCALIZED TARGET SPEAKER EXTRACTION

Meng Ge^{1,2}, Chenglin Xu^{3,*}, Longbiao Wang^{1,*}, Eng Siong Chng⁴, Jianwu Dang^{1,5}, Haizhou Li^{2,6}

¹ College of Intelligence and Computing, Tianjin University, Tianjin, China

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

³ Kuaishou Technology, Beijing, China

⁴ School of Computer Science and Engineering, Nanyang Technological University, Singapore

⁵ Japan Advanced Institute of Science and Technology, Ishikawa, Japan

⁶ The Chinese University of Hong Kong, Shenzhen, China

ABSTRACT

Speaker extraction aims to extract the target speaker’s voice from a multi-talker speech mixture given an auxiliary reference utterance. Recent studies show that speaker extraction benefits from the location or direction of the target speaker. However, these studies assume that the target speaker’s location is known in advance or detected by an extra visual cue, e.g., face image or video. In this paper, we propose an end-to-end localized target speaker extraction on pure speech cues, that is called L-SpEx. Specifically, we design a speaker localizer driven by the target speaker’s embedding to extract the spatial features, including direction-of-arrival (DOA) of the target speaker and beamforming output. Then, the spatial cues and target speaker’s embedding are both used to form a top-down auditory attention to the target speaker. Experiments on the multi-channel reverberant dataset called MC-Libri2Mix show that our L-SpEx approach significantly outperforms the baseline system.

Index Terms— Speaker extraction, speaker localizer, beamforming, DOA estimation, speaker embedding

1. INTRODUCTION

Human has the ability to selectively listen to a particular speaker through various stimuli in a multi-talker scenario, that is called selective auditory attention in *cocktail party* problem [1]. Ever since the theory is proposed, researchers

never stop seeking the engineering solution to confer human’s selective attention capability on machines, as there is high demand for various real-world applications, such as speech recognition [2, 3] and speaker verification [4, 5].

With the advent of deep learning in recent years, blind speech separation methods have been widely studied to solve the cocktail party problem by applying neural networks [6–9] and beamforming [10–13]. The neural network seeks the regular patterns (i.e., masks) between the time-frequency representation of the target speech and mixture speech, while beamforming incorporates the spatial statistics (i.e., spatial covariance matrix) obtained from the estimated masks to compute beamformer’s weights and filter the desired voice. However, speech separation always requires that the number of speakers is known as a prior, and assumes the label permutation is unchanged during training, which greatly limits its scope of real-world applications.

Unlike blind speech separation, speaker extraction only extracts the target speech from a mixture speech driven by spectral [14, 15] or spatial cues [16–18] of the target speaker. The spectral cue is always represented by a speaker embedding from an enrolled reference utterance, while the spatial cue is usually transformed into spectrum-like features derived from the target speaker’s location. For example, Chen et al. [16] introduced a location-based angle feature to guide separation network, which was the cosine distance between the steering vector and inter-channel phase difference (IPD) for each speaker in the mixture. Gu et al. [17] suggested that the beamforming output could be regarded as an alternative way of spatial cues, as beamforming aimed to summarize the signals from the target speaker’s direction and suppress non-target signals. However, these studies often require the speaker location is known in advance or detected using an extra visual cue.

To address this issue, we propose an end-to-end localized target speaker extraction on pure speech cues, that is called L-SpEx. We design a target speaker localizer driven by an enrolled utterance of the target speaker to extract the target

This work is supported by A*STAR under its RIE2020 Advanced Manufacturing and Engineering Domain (AME) Programmatic Grant (Grant No. A1687b0033); National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG-100E-2018-006); Human-Robot Interaction Phase 1 (Grant No. 192 25 00054), National Research Foundation (NRF) Singapore under the National Robotics Programme; National Natural Science Foundation (61771333), Tianjin Municipal Science and Technology Project (18ZXZNGX00330). This research is also funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy (University Allowance, EXC 2077, University of Bremen). * Corresponding author.

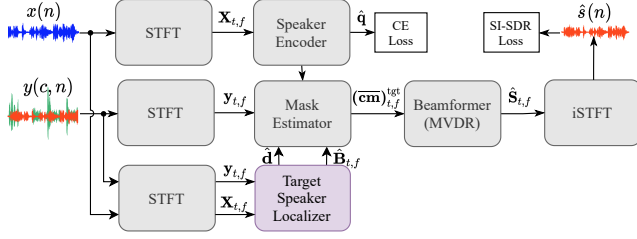


Fig. 1. The diagram of the proposed L-SpEx. The target speaker localizer is illustrated in Fig. 2.

speaker's DOA and beamforming output, simultaneously. By doing so, the extracted spatial cues and the enrolled utterance can be further used to guide the network to learn which direction and speaker to be extracted.

2. L-SPEX ARCHITECTURE

As illustrated in Fig. 1, the difference between the proposed L-SpEx system and other speaker extraction systems lies in the target speaker localizer (Fig. 2).

2.1. Target speaker localizer

The target speaker localizer learns to encode the spatial cues related to the target speaker's direction from the multi-channel mixture signal $y(c, n)$, with reference to a reference utterance $x(n)$ by the target speaker. We have,

$$y(c, n) = s(c, n) + \sum_{i=1}^I b_i(c, n) \quad (1)$$

where c denotes the channel index and n denotes discrete time index. $s(c, n)$ represents the target signal and $b_i(c, n)$ represents the interference signal corresponding to speaker i .

Formally, let $\mathbf{Y}_{t,f,c} \in \mathbb{C}$ be the STFT coefficient of the c -th channel mixture signal $y(c, n)$ at time-frequency bin (t, f) , and let $\mathbf{X}_{t,f} \in \mathbb{C}$ be the STFT coefficient of an enrolled single-channel utterance $x(n)$ of the corresponding target speaker. As shown in Fig. 2, we first employ a network to estimate a complex-valued mask, as opposed to a real-valued mask, for speaker location estimation. This is motivated by Sharath's work [19], which shows that speaker location estimation relies strongly on both phase-differences and magnitude-differences between the microphones. Thus, the complex-valued masks $(\mathbf{cm})_{t,f}^{\text{tgt}}$ is estimated as follow:

$$\begin{aligned} (\mathbf{cm})_{t,f}^{\text{tgt}} &= \{\text{Re}[(\mathbf{cm})_{t,f}^{\text{tgt}}], \text{Im}[(\mathbf{cm})_{t,f}^{\text{tgt}}]\} \\ &= \text{CMaskEst}(\mathbf{y}_{t,f}, \text{Enc}_{\text{speaker}}(\mathbf{X}_{t,f})) \end{aligned} \quad (2)$$

where $\mathbf{y}_{t,f} = \{\mathbf{Y}_{t,f,c}\}_{c=1}^C \in \mathbb{C}^C$ is the spatial vector of the signals obtained from all C -microphones for each time-frequency bin (t, f) . $\text{CMaskEst}(\cdot)$ and $\text{Enc}_{\text{speaker}}(\cdot)$ represent a complex mask estimator and speaker encoder, and both structures are built with several BLSTM layers as shown in

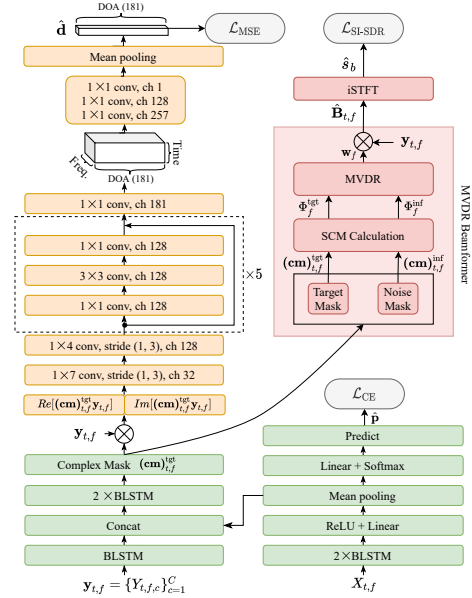


Fig. 2. The network structure of target speaker localizer.

Fig. 2. $\text{Re}[\cdot]$ and $\text{Im}[\cdot]$ denotes the real and imaginary part of a complex tensor, respectively.

After the masks are estimated, we use $(\mathbf{cm})_{t,f}^{\text{tgt}} \mathbf{y}_{t,f}$ as the input of DOA estimator to predict target speaker's DOA. We employ two CNN layers with kernel size of 1×7 and 1×4 on the masked input, followed by residual network blocks with a number of 5. Then, a 1×1 convolutional layer with 181 output channels (corresponding to 181 azimuth directions) projects the features to the DOA space. Finally, three 1×1 convolutional layers and a mean pooling operation summarizes the time and frequency to obtain the 181-dimension DOA vector:

$$\hat{\mathbf{d}} = \text{DOAEst}((\mathbf{cm})_{t,f}^{\text{tgt}} \mathbf{y}_{t,f}) \quad (3)$$

where $\text{DOAEst}(\cdot)$ denotes the DOA estimator.

Motivated by the studies [17, 20], beamforming output is also regarded as a direction-related spatial cue, as beamforming has the ability to summarize the signals from target speaker's direction. Therefore, the estimated masks are also used to compute the cross-channel spatial covariance matrices (SCMs) Φ_f^j and then obtain the beamforming output $\hat{\mathbf{B}}_{t,f}$,

$$\Phi_f^j = \frac{1}{\sum_{t=1}^T (\mathbf{cm})_{t,f}^j} \sum_{t=1}^T (\mathbf{cm})_{t,f}^j \mathbf{y}_{t,f} \mathbf{y}_{t,f}^H \quad (4)$$

$$\hat{\mathbf{B}}_{t,f} = \mathbf{w}_f^H \mathbf{y}_{t,f} = \left[\frac{(\Phi_f^{\text{inf}})^{-1} \Phi_f^{\text{tgt}}}{\text{tr}((\Phi_f^{\text{inf}})^{-1} \Phi_f^{\text{tgt}}) \mathbf{u}} \right]^H \mathbf{y}_{t,f} \quad (5)$$

Here, $j \in \{\text{tgt}, \text{inf}\}$ and $\mathbf{m}_{t,f}^{\text{inf}} = 1 - \mathbf{m}_{t,f}^{\text{tgt}}$. $\mathbf{u} \in \mathbb{R}^C$ is a vector denoting the reference microphone, and $\text{tr}(\cdot)$ denotes the trace operation. $\mathbf{w}_f = \{\mathbf{W}_{t,f,c}\}_{c=1}^C \in \mathbb{C}^C$ is corresponding to time-invariant beamformer coefficients. C denotes the number of channels and H denotes the conjugate transpose.

2.2. Localized target speaker extraction

Given the estimated DOA likelihood coding $\hat{\mathbf{d}}$ and beamforming output $\hat{\mathbf{B}}_{t,f}$ of the target speaker in Sec. 2.1, we can obtain two direction-related spatial features, i.e., DF_{angle} and DF_{beam} . The former [16] is derived from the estimated target speaker's angle $\hat{\theta} = \text{argmax}(\hat{\mathbf{d}})$, the latter is derived from beamforming output $\hat{\mathbf{B}}_{t,f}$, which are defined as follows

$$\text{DF}_{\text{angle}}(t, f) = \frac{1}{P} \sum_{l,r \in \Omega} \cos(\mathbf{o}_{l,r} - \frac{\pi f_s f \Delta_{l,r} \cos \hat{\theta}}{(N_{\text{FFT}} - 1)v}), \quad (6)$$

$$\text{DF}_{\text{beam}}(t, f) = \sqrt{\text{Re}[\hat{\mathbf{B}}_{t,f}]^2 + \text{Im}[\hat{\mathbf{B}}_{t,f}]^2} \quad (7)$$

where Ω contains P microphone pairs, and $\mathbf{o}_{l,r} = \angle \mathbf{Y}_{t,f,l} - \angle \mathbf{Y}_{t,f,r}$ represents the observed inter-channel phase difference (IPD) between left channel l and right channel r . N_{FFT} is the number of FFT bins, v is the sound velocity and f_s is the sampling rate. Note that f ranges from 0 to $(N_{\text{FFT}} - 1)$. $\Delta_{l,r}$ denotes the distance between the microphone pair (l, r) .

The above two direction-related spatial features have an ability to inform the extraction network of target speaker's direction, while the speaker embedding obtained from a speaker encoder can guide the network to attend to the target speaker. Thus, we use both spatial features as the inputs of the mask estimator to predict better masks for the target speaker extraction. The better masks $(\overline{\mathbf{cm}})_{t,f}^{\text{tgt}}$ are calculated as

$$\mathbf{y}_{t,f}^{\text{new}} = \text{Concat}[\mathbf{y}_{t,f}, \text{DF}_{\text{beam}}, \text{DF}_{\text{angle}}], \quad (8)$$

$$(\overline{\mathbf{cm}})_{t,f}^{\text{tgt}} = \overline{\text{CMaskEst}\{\mathbf{y}_{t,f}^{\text{new}}, \overline{\text{Enc}}_{\text{speaker}}(\mathbf{X}_{t,f})\}} \quad (9)$$

where $\text{Concat}[\cdot]$ is the concatenation operation. $\overline{\text{CMaskEst}}$ and $\overline{\text{Enc}}_{\text{speaker}}$ denote the complex mask estimator and speaker encoder in extraction network for speaker extraction.

Finally, the masks are used to obtain the extracted STFT spectrum $\hat{\mathbf{S}}_{t,f} \in \mathbb{C}$ by using MVDR formula, i.e., Eq. (4) and (5), and further estimate the target signal \hat{s} via iSTFT,

$$\hat{s} = \text{iSTFT}(\hat{\mathbf{S}}_{t,f}) = \text{iSTFT}(\text{MVDR}((\overline{\mathbf{cm}})_{t,f}^{\text{tgt}}, \mathbf{y}_{t,f})) \quad (10)$$

2.3. End-to-end training

We first pretrain the target speaker localizer by a multi-task learning strategy, and then optimize the whole network. The loss function of target speaker localizer is defined as follow:

$$\mathcal{L}_{\text{localizer}} = \mathcal{L}_{\text{SI-SDR}}(\hat{s}_b, s) + \alpha \mathcal{L}_{\text{CE}}(\hat{\mathbf{p}}, \mathbf{p}) + \beta \mathcal{L}_{\text{MSE}}(\hat{\mathbf{d}}, \mathbf{d})$$

$$\mathcal{L}_{\text{SI-SDR}}(\hat{s}_b, s) = -20 \log_{10} \frac{||(\hat{s}_b^T s / s^T s) \cdot s||}{||(\hat{s}_b^T s / s^T s) \cdot s - \hat{s}_b||}$$

$$\mathcal{L}_{\text{CE}}(\hat{\mathbf{p}}, \mathbf{p}) = - \sum_{i=1}^N \mathbf{p}_i \log(\hat{\mathbf{p}}_i)$$

$$\mathcal{L}_{\text{MSE}}(\hat{\mathbf{d}}, \mathbf{d}) = \sum_{i=1}^M ||\hat{\mathbf{d}}_i - \mathbf{d}_i||$$

where α and β are scaling factors. $\mathcal{L}_{\text{SI-SDR}}$ aims to minimize the signal reconstruction error. $\hat{s}_b = \text{iSTFT}(\hat{\mathbf{B}}_{t,f})$ and s are

the estimated signal and the target clean signal of reference microphone, respectively. \mathcal{L}_{CE} is the cross-entropy loss for speaker classification. N denotes the number of speakers. \mathbf{p}_i is the true class label for speaker i , and $\hat{\mathbf{p}}_i$ represents the predicted probability corresponding to i -th speaker. \mathcal{L}_{MSE} is the mean squared error for target DOA estimation. M is the number of azimuth directions, here $M = 181$. $\hat{\mathbf{d}}_i$ and \mathbf{d}_i are the predicted and ground-truth DOA coding of the target speaker. Based on the likelihood-based coding in [21], the desired ground-truth values \mathbf{d}_i are defined as follows:

$$\mathbf{d}_i = \begin{cases} e^{-d(\theta_i, \theta)^2 / \sigma^2}, & \text{if } \theta \text{ exists} \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

where θ is true target speaker's angle and θ_i is one of 181 azimuth directions. σ is the parameter to control the width of the Gaussian curves. $d(\cdot, \cdot)$ is the azimuth angular distance.

After target speaker localizer is pretrained, we optimize the whole L-SpEx network by using the following loss,

$$\mathcal{L}_{\text{L-SpEx}} = \mathcal{L}_{\text{SI-SDR}}(\hat{s}, s) + \gamma \mathcal{L}_{\text{CE}}(\hat{\mathbf{q}}, \mathbf{p}), \quad (12)$$

where \hat{s} is the extracted signal via Eq. 10, and $\hat{\mathbf{q}}$ is predicted probability from the speaker encoder $\overline{\text{Enc}}_{\text{speaker}}$. γ is also a scale factor to balance the two objectives.

3. EXPERIMENTS AND DISCUSSION

To facilitate the evaluation, we introduce a multi-channel reverberated version of the Libri2Mix¹ dataset which we refer to as MC-Libri2Mix. The original Libri2Mix is a clean 2-talker mixture corpus generated from the LibriSpeech corpus by mixing two randomly selected utterances. Libri2Mix contains training data with 64,700 utterances (270 hours, 1,172 speakers), development data with 3,000 utterances (11 hours, 40 speakers), and test data with 3,000 utterances (11 hours, 40 speakers). The average duration of the utterances is 14.8s.

The room impulse responses (RIRs) in MC-Libri2Mix is simulated using pyroomacoustics² package. For the room configurations, the length and width of each room are randomly drawn in the range [5, 10]m, and the height is selected in the range [3, 4]m. The reverberation time (RT_{60}) of the reverberant data ranges from 200ms to 600ms. In MC-Libri2Mix, we consider a linear array with four microphones, where the microphone-to-microphone distance is 5cm [20]. Target speakers are placed in the frontal plane and are at least 15° apart from each other. The distance of speaker and the center of microphone is from 0.75m to 2m.

Unlike in blind speech separation, we set that the speakers in each 2-talker mixed speech acted as the target speaker in turn, and the corresponding auxiliary reference speech is randomly selected from original LibriSpeech corpus. In practice, the training set (127,056 examples, 1,172 speakers) and

¹<https://github.com/JorisCos/LibriMix>

²<https://github.com/LCAV/pyroomacoustics>

Table 1. SDR (dB) and SI-SDR (dB) in a comparative study on the MC-Libri2Mix dataset under open condition. “m” and “cm” represent the real-value mask and complex-value mask, respectively. “Pretrained Speaker Localizer” indicates the extracted signal \hat{s}_b that is derived from the output spectrum $\hat{\mathbf{B}}_{t,f}$ of the pretrained target speaker localizer in Fig. 2.

ID	Methods	Mask Type	Spatial Cues		E2E Train	SDR	SI-SDR
			DF _{beam}	DF _{angle}			
1	Unprocessed	-	-	-	-	0.46	0.07
2	Mask MVDR (m)	m	✗	✗	-	8.03	6.36
3	Mask MVDR (cm)	cm	✗	✗	-	8.02	6.26
4	Pretrained Speaker Localizer	cm	✗	✗	-	7.44	5.80
5	L-SpEx	cm	✓	✓	✗	8.96	7.17
6		cm	✓	✓	✗	9.41	7.29
7		cm	✓	✓	✓	9.68	7.45

development set (2,344 examples, 1,172 speakers) are randomly selected from the training data of MC-Libri2Mix. The test set contains 6,000 examples from the test data of MC-Libri2Mix. The details of configuration and data simulation can be found at <https://github.com/gemengtju/L-SpEx.git>

3.1. Experimental setup

We train all systems for 70 epochs on the 4-channel mixture segments and their corresponding reference utterances. The learning rate is initialized to $1e^{-4}$ and decays by 0.5 if the accuracy of validation set is not improved in 2 consecutive epochs. Early stopping is applied if no best model is found in the validation set for 5 consecutive epochs. Adam is used as the optimizer. For feature extraction, STFT is performed with a 8k Hz sampling rate and a 25ms window length with a 10ms stride, and the feature dimension is 257. For mask estimation network, we use three BLSTM layers with 512 cells, and the dimension of speaker embedding is 256. The speaker embedding is inserted between first and second BLSTM layer. For the loss configuration, we used $\alpha = \gamma = 0.5, \beta = 10$ to balance the loss. The parameter σ in Gaussian curve is 6.

3.2. Results and analysis

We compare L-SpEx with the mask-based MVDR baseline systems on MC-Libri2Mix in terms of SDR and SI-SDR. From Table 1, we conclude: 1) Our L-SpEx approach achieves the best performance under the open condition. Compared to the “Mask MVDR (cm)” system, L-SpEx leads to 20.7% and 19.0% relative improvement in terms of SDR and SI-SDR measure, respectively. The improvements mainly come from the augmented spatial features through the proposed target speaker localizer module. 2) The proposed L-SpEx with DF_{beam} achieves 0.94 dB and 0.89 dB performance gain over the “Mask MVDR (cm)” baseline in terms of SDR and SI-SDR. This result proves that the spatial feature derived from beamforming outputs can let network attend to the signal from the target direction and ignore non-target signals. 3) The result of L-SpEx system with DF_{beam} and DF_{angle} further show the effectiveness of spatial feature. This shows that DF_{angle} and DF_{beam} represents different aspects of spatial cues of target speaker, and the two features

Table 2. SDR (dB) and SI-SDR (dB) in a comparative study of different angle distance under open condition. The percentage in the table head indicates the ratio of each condition.

ID	Methods	< 45° (34.6%)		45°-90° (36.5%)		> 90° (28.9%)	
		SDR	SI-SDR	SDR	SI-SDR	SDR	SI-SDR
1	Unprocessed	0.46	0.06	0.43	0.06	0.48	0.08
2	Mask MVDR (m)	7.52	5.92	8.35	6.66	8.23	6.51
3	Mask MVDR (cm)	7.49	5.82	8.37	6.57	8.22	6.42
4	Pretrained Speaker Localizer	6.95	5.36	7.72	6.09	7.65	6.00
5	L-SpEx	8.27	6.55	9.37	7.56	9.26	7.45
6		8.64	6.63	9.88	7.68	9.75	7.59
7		8.97	7.06	9.78	7.66	9.75	7.66

Table 3. SDR (dB) and SI-SDR (dB) in a comparative study of different and same gender mixture under open condition.

ID	Methods	Diff. Gender (25.2%)		Same Gender (74.8%)	
		SDR	SI-SDR	SDR	SI-SDR
1	Unprocessed	0.35	-0.05	0.49	0.11
2	Mask MVDR (m)	9.91	8.22	7.39	5.74
3	Mask MVDR (cm)	9.82	8.02	7.42	5.67
4	Pretrained Speaker Localizer	9.02	7.39	6.90	5.28
5	L-SpEx	10.45	8.76	8.46	6.64
6		11.10	9.00	8.85	6.71
7		11.11	9.20	8.94	6.86

complement each other. 4) The result evaluated on the output of pre-trained speaker localizer is worse than the mask-based MVDR systems. The reason is that the DOA estimation and beamforming have something in common, but they are still two separate tasks and one mask is hard to optimize.

We further report the speaker extraction performance on different angle distance mixture speech in Table 2. From Table 2, we find that the mixture speech with smaller angle difference (i.e., < 45°) is more difficult to extract target speaker. It is observed that our L-SpEx system achieves the SI-SDR from 5.82 dB to 7.06 dB, which is even higher than the performance of “Mask MVDR (cm)” on large angle distance (i.e., > 45°). Furthermore, we report the extraction performance with different and same gender mixture speech, separately in Table 3. We observe that separating same-gender mixture is a more challenging task, as the spectral cues between each speaker are similar. When we incorporate the direction-related spatial information (i.e., DF_{beam} and DF_{angle}) to increase the discrimination between speakers, the extraction process starts to be easier. Specifically, our L-SpEx system improve 1.19 dB SI-SDR performance compared with the spectral-only “Mask MVDR (cm)” system.

4. CONCLUSIONS

In this paper, we proposed an end-to-end localized target speaker extraction on pure speech cues, that is called L-SpEx, to eliminate the assumption of known target angle in existing studies. We took the advantages of a target speaker’s enrolled utterance to design speaker localizer for estimating direction-related spatial cues. Experiments showed that the extracted spatial cues and the enrolled utterance input let extraction network works better, as target speech can be extracted based on speaker and direction aspects.

5. REFERENCES

- [1] E Colin Cherry, “Some experiments on the recognition of speech, with one and with two ears,” *The Journal of the acoustical society of America*, vol. 25, no. 5, pp. 975–979, 1953.
- [2] Hangting Chen, Pengyuan Zhang, Qian Shi, and Zuozhen Liu, “Improved guided source separation integrated with a strong back-end for the chime-6 dinner party scenario.,” in *Interspeech*, 2020, pp. 334–338.
- [3] Jing Shi, Xuankai Chang, Pengcheng Guo, Shinji Watanabe, Yusuke Fujita, Jiaming Xu, Bo Xu, and Lei Xie, “Sequence to multi-sequence learning via conditional chain mapping for mixture signals,” *arXiv preprint arXiv:2006.14150*, 2020.
- [4] Wei Rao, Chenglin Xu, Eng Siong Chng, and Haizhou Li, “Target speaker extraction for multi-talker speaker verification,” in *Interspeech*, 2019, pp. 1273–1277.
- [5] Chenglin Xu, Wei Rao, Jibin Wu, and Haizhou Li, “Target speaker verification with selective auditory attention for single and multi-talker speech,” *IEEE/ACM Transactions on audio, speech, and language processing*, vol. 29, pp. 2696–2709, 2021.
- [6] John R Hershey, Zhuo Chen, Jonathan Le Roux, and Shinji Watanabe, “Deep clustering: Discriminative embeddings for segmentation and separation,” in *ICASSP*, 2016, pp. 31–35.
- [7] Dong Yu, Morten Kolbæk, Zhenghua Tan, and Jesper Jensen, “Permutation invariant training of deep models for speaker-independent multi-talker speech separation,” in *ICASSP*, 2017, pp. 241–245.
- [8] Yi Luo and Nima 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.
- [9] Yi Luo, Zhuo Chen, and Takuya Yoshioka, “Dual-path rnn: efficient long sequence modeling for time-domain single-channel speech separation,” in *ICASSP*, 2020, pp. 46–50.
- [10] Jian Wu, Zhuo Chen, Jinyu Li, Takuya Yoshioka, Zhili Tan, Ed Lin, Yi Luo, and Lei Xie, “An end-to-end architecture of online multi-channel speech separation,” *arXiv preprint arXiv:2009.03141*, 2020.
- [11] Tsubasa Ochiai, Marc Delcroix, Rintaro Ikeshita, Keisuke Kinoshita, Tomohiro Nakatani, and Shoko Araki, “Beam-tasnet: Time-domain audio separation network meets frequency-domain beamformer,” in *ICASSP*, 2020, pp. 6384–6388.
- [12] Xuankai Chang, Wangyou Zhang, Yanmin Qian, Jonathan Le Roux, and Shinji Watanabe, “Mimo-speech: End-to-end multi-channel multi-speaker speech recognition,” in *ASRU 2019*, 2019, pp. 237–244.
- [13] Wangyou Zhang, Aswin Shanmugam Subramanian, Xuankai Chang, Shinji Watanabe, and Yanmin Qian, “End-to-end far-field speech recognition with unified dereverberation and beamforming,” in *Interspeech*, 2020.
- [14] Chenglin Xu, Wei Rao, Chng Eng Siong, and Haizhou Li, “SpEx: Multi-scale time domain speaker extraction network,” *IEEE/ACM Transaction on Audio, Speech, and Language Processing*, vol. 28, pp. 1370–1384, 2020.
- [15] Meng Ge, Chenglin Xu, Longbiao Wang, Eng Siong Chng, Jianwu Dang, and Haizhou Li, “Spex+: A complete time domain speaker extraction network,” in *Proc. Interspeech*, 2020, pp. 1406–1410.
- [16] Zhuo Chen, Xiong Xiao, Takuya Yoshioka, Hakan Erdogan, Jinyu Li, and Yifan Gong, “Multi-channel overlapped speech recognition with location guided speech extraction network,” in *SLT 2018*, 2018, pp. 558–565.
- [17] Rongzhi Gu, Lianwu Chen, Shi-Xiong Zhang, Jimeng Zheng, Yong Xu, Meng Yu, Dan Su, Yuexian Zou, and Dong Yu, “Neural spatial filter: Target speaker speech separation assisted with directional information.,” in *Interspeech*, 2019, pp. 4290–4294.
- [18] Rongzhi Gu, Shi-Xiong Zhang, Yong Xu, Lianwu Chen, Yuexian Zou, and Dong Yu, “Multi-modal multi-channel target speech separation,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 3, pp. 530–541, 2020.
- [19] Sharath Adavanne, Archontis Politis, Joonas Nikunen, and Tuomas Virtanen, “Sound event localization and detection of overlapping sources using convolutional recurrent neural networks,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 1, pp. 34–48, 2018.
- [20] Zhong-Qiu Wang and DeLiang Wang, “Combining spectral and spatial features for deep learning based blind speaker separation,” *IEEE/ACM Transactions on audio, speech, and language processing*, vol. 27, no. 2, pp. 457–468, 2018.
- [21] Weipeng He, Petr Motlicek, and Jean-Marc Odobez, “Neural network adaptation and data augmentation for multi-speaker direction-of-arrival estimation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 1303–1317, 2021.