# DEEP AUDIO ZOOMING: BEAMWIDTH-CONTROLLABLE NEURAL BEAMFORMER

*Meng Yu, Dong Yu*

Tencent AI Lab, Bellevue, WA, USA
{raymondmyu, dyu}@global.tencent.com

## ABSTRACT

Audio zooming, a signal processing technique, enables selective focusing and enhancement of sound signals from a specified region, attenuating others. While traditional beamforming and neural beamforming techniques, centered on creating a directional array, necessitate the designation of a singular target direction, they often overlook the concept of a field of view (FOV), that defines an angular area. In this paper, we proposed a simple yet effective FOV feature, amalgamating all directional attributes within the user-defined field. In conjunction, we've introduced a counter FOV feature capturing directional aspects outside the desired field. Such advancements ensure refined sound capture, particularly emphasizing the FOV's boundaries, and guarantee the enhanced capture of all desired sound sources inside the user-defined field. The results from the experiment demonstrate the efficacy of the introduced angular FOV feature and its seamless incorporation into a low-power subband model suited for real-time applications.

***Index Terms***— audio zooming, field of view, beamwidth-controllable, neural beamformer
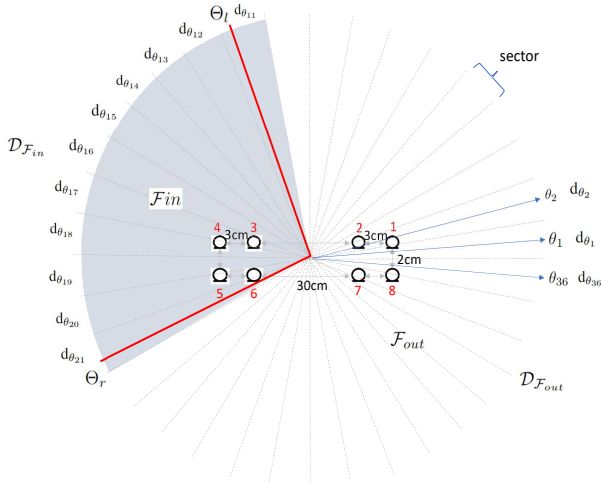
## 1. INTRODUCTION

Audio zooming technology refers to a signal processing technique that allows a user to selectively focus on and enhance the audio signals originating from a specific region of interest in a sound field, while attenuating signals from other directions. It finds applications in various fields, including teleconferencing, surveillance, broadcasting and video filming. With the audio zooming feature, people can adjust the video recording's focus using a pinch motion on the screen. When you zoom in, the audio from your focal point becomes more pronounced. Conversely, as you zoom out, the ambient sounds become more evident and are no longer diminished. Beamforming is the most closely related field of study to this concept, that amplifies the sound from the target direction while suppressing all other sounds [1, 2]. Most beamforming techniques share the same objective of improving the sound quality from a single direction. However, they do not take into account the concept of a field of view (FOV). The width of the enhanced beam in beamforming is determined by several factors, including the microphone configuration,

signal frequency, and type of beamformer used. Duong et al. created the audio zoom effect by weighted mixing the beamforming enhanced target sound source and the original microphone signal [3]. Nair et al. [4] proposed a beamforming method that depends on the sound spectral covariance matrices, taking into account the desired audio signals both within and outside the FOV. With sound reflections, audio waves from a source outside the FOV can reach the microphone after bouncing off surfaces from within the FOV. Under these circumstances, their audio zooming technique would continue to amplify these reflected sound signals.

An array of microphones carry spatial information about the source of a sound. The effectiveness of established spatial features, such as inter-channel phase difference (IPD), has been demonstrated at the input stage for speech separation methods that use time-frequency (T-F) masking [5–7]. Additionally, to further improve the extraction of source signals from a specific direction, meticulously crafted directional features indicating the dominant directional source in each T-F bin have been introduced in prior works such as [5, 8–11]. These directional features are contingent upon a specific direction of interest. Drawing inspiration from the multi-look/zone neural beamformer [12, 13], we introduced a FOV feature that consolidates all directional features within the desired field of interest. Simultaneously, we integrate a counter FOV feature to signify all the directional features outside the angular field. Employing both these features refines sound capture, especially around the FOV's edges. When this feature is input into the neural network, it ensures the capture of all sound sources or speakers within the designated angular field, while attenuating those outside this zone. The rest of the paper is organized as follows. In Section 2, we first recap the directional feature mentioned in the recent literature, and then present the new audio zooming FOV feature and its usage in the neural network model. We describe our experimental setups and evaluate the effectiveness of the proposed method in Section 3. We conclude this work in Section 4.

## 2. DEEP AUDIO ZOOMING FEATURE AND MODEL

Previous work in [5, 9, 10, 14] have proposed to leverage a proper designed directional feature of the target speaker to perform the target speaker separation. Through a short-time-Fourier-transform (STFT), the $M$ microphone signals

**Fig. 1**. The audio zooming 2D FOV feature extraction.

$y_m$ is transformed to its complex spectrum $Y_m$, where $m = 1, 2 \ldots, M$. IPD is computed by the phase difference between channels of complex spectrograms as $\text{IPD}^{(m)}(t, f) = \angle \mathbf{Y}^{m_1}(t, f) - \angle \mathbf{Y}^{m_2}(t, f)$, where $m_1$ and $m_2$ are two microphones of the $m$-th microphone pair out of $\overline{M}$ selected microphone pairs. A directional feature is incorporated as a target speaker bias. This feature was originally introduced in [5], which computes the averaged cosine distance between the target speaker steering vector and IPD on all selected microphone pairs as
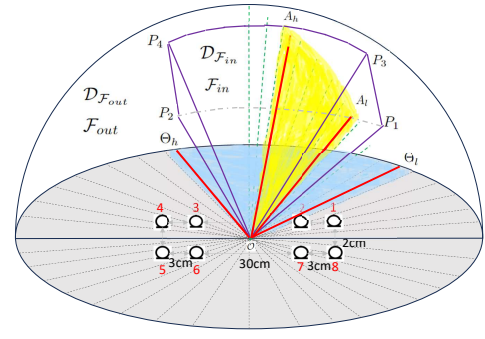
$$\mathrm{d}_\theta(t, f) = \sum_{m=1}^{\overline{M}} \left\langle \mathbf{e}^{\angle \mathbf{v}_\theta^{(m)}(f)}, \mathbf{e}^{\text{IPD}^{(m)}(t,f)} \right\rangle \quad (1)$$

where $\angle \mathbf{v}_\theta^{(m)}(f) := 2\pi f \Delta^{(m)} \cos \theta^{(m)}/c$ is phase of the steering vector from the direction $\theta$ at frequency $f$ with respect to $m$-th microphone pair, $\Delta^{(m)}$ is the distance between the $m$-th microphone pair, $c$ is the sound velocity, and vector $\mathbf{e}^{(\cdot)} := [\cos(\cdot), \sin(\cdot)]^T$. If the T-F bin $(t, f)$ is dominated by the source from $\theta$, then $d_\theta(t, f)$ will be close to 1, otherwise it deviates towards -1. As a result, $\mathrm{d}_\theta(t, f)$ indicates if a speaker from a desired direction $\theta$ dominates in each T-F bin. Such location-based input features (LBI) are found in [5, 8–10].

## 2.1. Audio Zooming FOV Feature

Beyond the LBI approach, when a scenario involves multiple speakers that need to be distinctly separated, the location-based training (LBT) method is utilized [15, 16]. For specific application scenarios, DNNs have been trained to provide distinct outputs for every designated direction, including those devoid of source activity [17, 18]. When we evenly divide the space into several sectors and sample a direction at the bisector of each sector, we establish the idea of multi-look directions. The idea of "multi-look direction" has been applied to speech separation [12, 19–21] and multi-channel acoustic model [22–24], respectively, where a small number of spatial look directions cover all possible target speaker directions.

A set of $K$ directions in the horizontal plane is sampled. For example, in Fig. 1, $K = 36$ and the horizontal space is
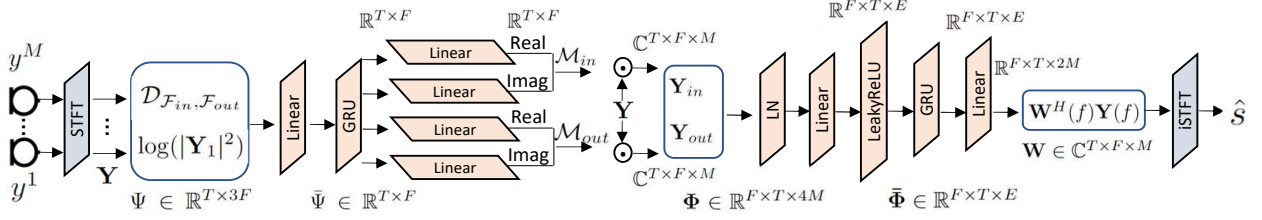
evenly partitioned to 36 sectors. The azimuths of look directions $\theta_{1,2,\ldots,36} = \{5°, 15°, \ldots, 345°, 355°\}$, corresponding to the sectors $[0°, 10°]$, $[10°, 20°]$, $\ldots$, $[340°, 350°]$, and $[350°, 360°]$, respectively. The look directions $\theta_{1,2,\ldots,K}$ result in $K$ directional feature vectors $\mathrm{d}_{\theta_k}, k = 1, 2, \ldots, K$. Obviously, $\mathrm{d}_{\theta_k}(t, f)$ is determined by the actual source direction $\theta$ and look direction $\theta_k$. Consequently, for those T-F bins influenced predominantly by the source near the look direction $\theta_k$, the value of $\mathrm{d}_{\theta_k}$ will be substantial.

For a user-specified FOV delineated by the angular boundaries $[\Theta_l, \Theta_h]$, we first determine the corresponding sector of these boundaries. The sectors on and within the range $[\Theta_l, \Theta_h]$ are represented by its index set $\mathcal{F}in$, i.e. the gray area in Fig.1. On the other hand, sectors that fall outside these boundaries are represented by $\mathcal{F}_{out}$. Each of the directional feature $\mathrm{d}_{\theta_k}$ belongs to $\mathbb{R}^{T \times F}$, where $T$ and $F$ are the total number of frames and frequency bands of the complex spectrogram, respectively. By concatenating the values of $\mathrm{d}_{\theta_k}$ where $k \in \mathcal{F}_{in}$, the resulting dimension becomes considerably large. Moreover, this dimension fluctuates with variations in the field size, making it unsuitable for model training. We introduce a method to consolidate the high-dimensional directional feature sets by employing a max operation across the sectors, as detailed in

$$\mathcal{D}_{\mathcal{F}_{in}}(t, f) = \max_{\theta_k \in \mathcal{F}_{in}} \mathrm{d}_{\theta_k}(t, f). \quad (2)$$

The underlying idea is that when a speaker resides within the field $\mathcal{F}_{in}$, there will inevitably be a sampled direction, $\theta_{k_0}, k_0 \in \mathcal{F}_{in}$, nearer to this speaker than other sampled directions. This results in $\mathrm{d}_{\theta_{k_0}}$ being the most prominent at the T-F bins where this speaker has dominance. Consequently, employing the max operation for $\mathrm{d}_{\theta_k}$ within $\mathcal{F}_{in}$ can effectively capture the prominence of all active speakers within the targeted field. Meanwhile, following the max operation, the dimension of the resulting feature vector $\mathcal{D}_{\mathcal{F}_{in}}$ remains consistent with that of any original directional feature $\mathrm{d}_{\theta_k}$. At the same time, we incorporate a counter FOV feature to represent all directional attributes beyond the angular field. Using both these features sharpens the sound capture, particularly near the boundaries of the FOV. This enhances the model's ability to differentiate between sound sources within and outside the targeted field. The counter FOV feature $\mathcal{D}_{\mathcal{F}_{out}}(t, f)$ is computed similarly as Eq. 2, with $\theta_k \in \mathcal{F}_{out}$. There are two methods to further consolidate the two FOV feature vectors.



**Fig. 2**. The audio zooming 3D FOV feature extraction.

**Fig. 3**. Deep audio zooming model diagram. Trained on 4-second chunks with the Adam optimizer and a batch size 16 for 30 epochs. The initial learning rate is set to 1e-4 with a gradient norm clipped with max norm 10.

The first approach involves concatenation, represented as

$$\mathcal{D}_{\mathcal{F}_{in},\mathcal{F}_{out}} = [\mathcal{D}_{\mathcal{F}_{in}}, \mathcal{D}_{\mathcal{F}_{out}}] \in \mathbb{R}^{T \times 2F}. \qquad (3)$$

The second is through post-processing, as calculated in Eq. 4,

$$\mathcal{D}_{\mathcal{F}_{in},\mathcal{F}_{out}}(t,f) = \begin{cases} -1 & \text{if } \mathcal{D}_{\mathcal{F}_{in}}(t,f) \leq \mathcal{D}_{\mathcal{F}_{out}}(t,f) \\ \mathcal{D}_{\mathcal{F}_{in}}(t,f) & \text{else} \end{cases}$$

$$(4)$$

The post-processing aims to remove the components of speakers outside the desired region, thereby enhancing the prominence of the speakers within the region.

To adapt the aforementioned feature extraction to a 3D space [25], we merely modify the directional feature $d_\theta$ in Eq. 1 such that $\angle \mathbf{v}_\theta^{(m)}(f) := 2\pi f \Delta^{(m)} \cos \theta^{(m)} \cos \alpha^{(m)}/c$, with $\alpha$ representing the elevation angle in the vertical dimension. Concurrently, the definitions of $\mathcal{F}_{in}$ and $\mathcal{F}_{out}$ will also account for the preferred elevation span. As illustrated in Fig. 2, two red lines at angles $\Theta_l$ and $\Theta_h$ delineate the desired angular boundary on the horizontal plane. Simultaneously, two red lines at angles $A_u$ and $A_d$ set the boundary on the vertical plane. When considering these angular boundaries in a 3D space collectively, the sectors within the field, denoted as $\mathcal{F}_{in}$, shape the pyramid defined by $\overline{\mathcal{O}P_1P_2P_3P_4}$.
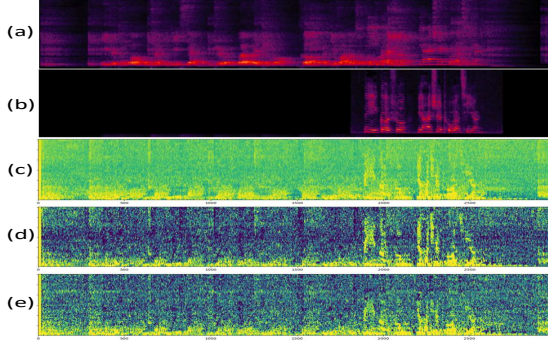
### 2.2. Subband Model
We formulate the problem as amplifying speech signals from a specified angular range $[\Theta_l, \Theta_h]$ and $[A_l, A_h]$. The diagram of the deep audio zooming model, which is similar yet more optimized than the model presented in [11, 13], can be seen in Fig. 3. By computing STFT (512 window size and 256 hop size) on $M = 8$ microphone channels, a reference channel, e.g. the first channel complex spectrogram $\mathbf{Y}_1$, is used to compute logarithm power spectrum (LPS) by $\log(|\mathbf{Y}_1|^2) \in \mathbb{R}^{T \times F}$. The spectral feature LPS, combined with the introduced audio zooming FOV feature $\mathcal{D}_{\mathcal{F}_{in},\mathcal{F}_{out}}$ in Eq. 3, are utilized as the input for the neural network. Subsequent to processing through a GRU layer, the model deduces $T \times F$ dimensional complex-valued masks: $\mathcal{M}_{in}(t,f)$ and $\mathcal{M}_{out}(t,f)$ for source estimations within and outside the FOV, respectively, via four linear layers. Noise signals are associated with the sources outside the FOV since noise reduction is integral to our objective. The two masks are multiplied with each microphone spectrum $\mathbf{Y}^m$ to estimate the target signals and interfering signals at each microphone $m$. $M$ channels of real and imaginary parts of $\mathbf{Y}_{in}$ and $\mathbf{Y}_{out}$ are then concatenated, represented as $\mathbf{\Phi} \in \mathbb{R}^{F \times T \times 4M}$. After layer

normalization [26] and a linear projection with leakyRELU activation, a RNN layer together with a linear layer receive the subband embedding of individual frequency bands $\bar{\mathbf{\Phi}}(f) \in \mathbb{R}^{T \times E}$, where $f = 1, 2, \ldots F$ and $E = 32$ is the dimension of the subband embedding, to estimate the enhancement filters at the corresponding frequency as $\mathbf{W}(f) \in \mathbb{C}^{T \times M}$. We compute the output signal $\hat{\mathbf{S}}$ by $\mathbf{W}^H(f)\mathbf{Y}(f)$. The total number of parameters is 860K and the computation cost is 184 MMACs/s. Same as the one in [27], the loss function is calculated using the combination of negative SISDR [28] in time domain and $l_1$ norm on the spectrum magnitude difference.

## 3. EXPERIMENTS AND RESULTS

We simulate the 8-channel dataset using AISHELL-2 [29]. The geometry of the microphone array, which is mounted on the ceiling, is depicted in Fig. 1. We generate 10K room impulse responses with random room characteristics and reverberation time ranging from 0.3s to 1.3s using gpuRIR [30]. In addition, environmental noises are added with signal-to-noise-ratio ranging from 10 to 40dB. 95K, 2.5K, and 100 utterances are generated for training, validation and testing, respectively. A total of up to five speakers are randomly sampled. For each training utterances we randomly sample the angular range $[\Theta_l, \Theta_h]$ and $[A_l, A_h]$, ensuring the model encounters a variety of speaker counts both within and outside the designated target field.

We begin by examining the efficacy of the introduced audio zooming FOV feature, $\mathcal{D}_{\mathcal{F}_{in},\mathcal{F}_{out}}$, in its two representations Eq. 3 and 4. The horizontal and vertical resolutions, which correspond to the sector width, are configured at $20°$ and $10°$, respectively. The microphone pairs chosen to compute IPD are (1, 4), (2, 6), (1, 7), (2, 7), (4, 6), and (3, 7). For every test sample, the beam-center aligns with the angular direction of one speaker. Meanwhile, the beam-width undergoes random sampling on both the horizontal and vertical planes. Consequently, the target output signal always contains at least one speaker. Fig. 4 displays the extracted feature for a target area with an active speaker. Comparing the raw audio zooming feature from Eq. 2 with the post-processed feature in Eq. 4, it's evident that the post-processing refines the feature pattern. This refinement results in the attenuation of T-F bins dominated by sound sources outside the target region, thereby accentuating the speakers within the desired area and enhancing their distinctiveness. The findings in Table 1 indicate that concatenating the FOV feature from both inside and outside the target region ("our best") is more effective than using the

**Fig. 4**. Audio zooming FOV feature extraction. a) microphone signal, b) speaker in the target region, c) original feature $\mathcal{D}_{\mathcal{F}in}$ in Eq. 2, d) post processed feature $\mathcal{D}_{\mathcal{F}in,\mathcal{F}out}$ in Eq. 4. e) feature in Eq. 4 with low resolution.



**Fig. 5**. Deep audio zooming real demo.

post-processed version. This discrepancy arises because the "winner-take-all" approach in post-processing doesn't cater well to bins where speakers overlap in the T-F domain. As a result, the post-processed feature might lead to attenuation of the target speaker in certain T-F bins.

The calculation of the FOV feature relies on space partitioning, essentially referring to the sampling resolution of the look directions. A higher resolution leads to a more precise raw directional feature (as per Eq. 1). Consequently, post-aggregation through the max operation detailed in Eq. 2, the resultant audio zooming feature is sharpened, allowing for superior differentiation between sources within and outside the target region. The row (e) in Fig. 4 displays the feature extracted at the horizontal resolution $60°$ and vertical resolution $15°$. A more refined resolution demonstrates markedly enhanced feature prominence in row (d).

In our subband model, the pivotal component is the second GRU layer, which processes data for each frequency band separately. The model produces $M$ channel complex weights, and the final output is derived from the weighted summation of all microphone signals, functioning similar to beamformer filtering. For a comparative analysis, we crafted a full-band GRU model. This was achieved by eliminating the layers and operations post the layer normalization as seen in Fig. 3, augmenting the initial GRU with additional layers, and expanding its size. This ensures that, in terms of computational expense, the full-band model aligns with the subband model. The ultimate output is procured via $\mathcal{M}_{in}\mathbf{Y}_1$. Results from this model comparison are detailed in Table 1. The subband model ("our best") distinctly surpasses the fullband model.

Our subsequent investigation delves into the performance implications of utilizing a significantly smaller number of microphones, such as three microphones, akin to those found on

**Table 1**. *Audio zooming experimental comparisons*

| Category | Method | PESQ [31] | SDR[dB] [32] |
|---|---|---|---|
| - | no processing | 2.01 | 4.94 |
| feature | $\mathcal{D}_{\mathcal{F}in}$ Eq. 2 | 2.66 | 10.19 |
| | $\mathcal{D}_{\mathcal{F}in,\mathcal{F}out}$ Eq. 4 | 2.73 | 11.32 |
| resolution | h.:$60°$, v.: $15°$ | 2.56 | 10.55 |
| model | fullband model | 2.44 | 10.32 |
| # of channels | 3-mic | 2.61 | 9.94 |
| baseline | oracle MVDR | 2.31 | 6.09 |
| | $d_\theta$ Eq. 1 | 2.43 | 9.10 |
| our best | $\mathcal{D}_{\mathcal{F}in,\mathcal{F}out}$ Eq. 3 h.:$20°$, v.: $10°$ subband 8-mic | **2.78** | **11.50** |

a phone. In this 3-channel array setup, microphones 1, 2, and 7 are utilized. The pairs (1,2), (1, 7), and (2, 7) are employed to compute the IPD. From Table 1, it reveals that the 3-microphone model's performance is almost on par with the 8-microphone model, albeit with slight discrepancies.

The oracle MVDR, derived from the Ideal Ratio Mask (IRM), serves as our baseline. For each test utterance, we compute the IRM using the target and microphone signals. This IRM then helps compute the noise covariance matrix. We use the ground-truth beam-center as the target MVDR direction. This approach positions the oracle MVDR as the upper limit for the audio zooming method in [4]. Additionally, the directional feature from Eq. 1 is inputted into our subband model for a direction-specific LBI enhancement, rather than a field-wide enhancement. The details of the comparison can be found in Table 1. The live demonstration under an actual scenario was carried out in an office setting, as depicted in Fig. 5. In this scene, five individuals are engaged in conversation; three of them are seated at a table, while the other two stand beside a wall at some distance from the table. The device's processing capabilities allow for user-directed control. For instance, it can focus solely on the individuals seated at the table or target a specific area away from the table. When no speaker is present in the target region, there's an attenuation of about 49.0dB. Moreover, the DNSMOS score [33] for test samples rose from 1.98 to 2.49.

## 4. CONCLUSION

In this study, we introduced a straightforward yet potent angular region feature that allows the neural network model to perform audio zooming efficiently. The methods of random field sampling during training and converting field boundaries into look directions are both essential components for achieving beamwidth-adjustable neural beamforming. The experimental results highlight the promise of our pioneering deep learning-based audio zooming approach. In subsequent studies, we plan to delve into how to more effectively utilize different microphone pairs during the feature computation.

# 5. REFERENCES

[1] B. D. Van and K. M. Buckley, "Beamforming: A versatile approach to spatial filtering," *IEEE ASSP magazine, vol. 5, no. 2, pp. 4-24*, 1988.

[2] J. Benesty, J. Chen, and Y. Huang, "Microphone array signal processing," *Springer Science & Business Media, vol. 1*, 2008.

[3] N. Duong, P. Berthet, S. Zabre, M. Kerdranvat, A. Ozerov, and L. Chevallier, "Audio zoom for smartphones based on multiple adaptive beamformers," *13th International Conference on Latent Variable Analysis and Signal Separation*, 2017.

[4] A. Nair, A. Reiter, C. Zheng, and S. Nayar, "Audio zoom for smartphones based on multiple adaptive beamformers," *ACM International Conference on Multimedia*, 2019.

[5] Z. Chen, X Xiao, T. Yoshioka, H. Erdogan, J. Li, and Y. Gong, "Multi-channel overlapped speech recognition with location guided speech extraction network," in *IEEE Spoken Language Technology Workshop (SLT)*, 2018, pp. 558–565.

[6] Q. Wang, J. L. Roux, and J. R. Hershey, "Multi-channel deep clustering: Discriminative spectral and spatial embeddings for speaker-independent speech separation," in *ICASSP*, 2018.

[7] L. Chen, M. Yu, D. Su, and D. Yu, "Multi-band PIT and model integration for improved multi-channel speech separation," in *ICASSP*. IEEE, 2019.

[8] Z. Wang and D. 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, 2019.

[9] R. Gu, L. Chen, S. X. Zhang, J. Zheng, Y. Xu, M. Yu, D. Su, Y. Zou, and D. Yu, "Neural spatial filter: Target speaker speech separation assisted with directional information," in *Proc. Interspeech*, 2019.

[10] F. Bahmaninezhad, J. Wu, R. Gu, S. X. Zhang, Y. Xu, M. Yu, and D. Yu, "A comprehensive study of speech separation: spectrogram vs waveform separation," *Proc. Interspeech*, 2019.

[11] Y. Xu, Z. Zhang, M. Yu, S.-X. Zhang, and D. Yu, "Generalized spatio-temporal rnn beamformer for target speech separation," in *Proc. 29th Eur. Signal Process. Conf.*, 2021.

[12] M. Yu, X. Ji, B. Wu, D. Su, and D. Yu, "End-to-end multi-look keyword spotting," *Interspeech*, 2020.

[13] Y. Xu, V. Kothapally, M. Yu, S. Zhang, and D. Yu, "Zoneformer: On-device neural beamformer for in-car multi-zone speech separation, enhancement and echo cancellation," *Interspeech*, 2023.

[14] R. Gu, S.-X. Zhang, Y. Xu, L. Chen, Y. Zou, , and D. Yu, "Multi-modal multi-channel target speech separation," *IEEE Journal of Selcted Topics in Signal Processing,*, 2020.

[15] H. Taherian, K. Tan, and D. Wang, "Location-based training for multichannel talker-independent speaker separation," in *ICASSP*, 2022.

[16] H. Taherian, K. Tan, and D. Wang, "Multi-channel talker-independent speaker separation through location-based training," in *IEEE/ACM Trans. Audio, Speech, Language Process., vol. 30, pp. 2791–2800*, 2022.

[17] S. E. Chazan, H. Hammer, G. Hazan, J. Goldberger, and S. Gannot, "Multi-microphone speaker separation based on deep DOA estimation," in *in Proc. 27th Eur. Signal Process. Conf.*, 2019.

[18] A. Bohlender, A. Spriet, W. Tirry, and N. Madhu, "Neural networks using full-band and subband spatial features for mask based source separation," in *in Proc. 27th Eur. Signal Process. Conf.*, 2021.

[19] X. Ji, M. Yu, J. Chen, J. Zheng, D. Su, and D. Yu, "Integration of multi-look beamformers for multi-channel keyword spotting," in *ICASSP 2020*, pp. 7464–7468.

[20] Z. Chen, J. Li, X. Xiao, T. Yoshioka, H. Wang, Z. Wang, and Y. Gong, "Cracking the cocktail party problem by multi-beam deep attractor network," in *ASRU 2017*.

[21] Z. Chen, T. Yoshioka, X. Xiao, J. Li, M. L. Seltzer, and Y. Gong, "Efficient integration of fixed beamformers and speech separation networks for multi-channel far-field speech separation," .

[22] T. N. Sainath, R. J. Weiss, K. W. Wilson, A. Senior, and O. Vinyals, "Learning the speech front-end with raw waveform CLDNNs," in *Proc. Interspeech*, 2015.

[23] T. N. Sainath, R. J. Weiss, K. W. Wilson, A. Narayanan, and M. Bacchiani, "Factored spatial and spectral multichannel raw waveform cldnns," in *ICASSP*. IEEE, 2016.

[24] T. N. Sainath, R. J. Weiss, K. W. Wilson, B. Li, A. Narayanan, E. Variani, M. Bacchiani, I. Shafran, A. Senior, K. Chin, A. Misra, and C. Kim, "Multichannel signal processing with deep neural networks for automatic speech recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2017.

[25] R. Gu, S-X Zhang, M. Yu, and D. Yu, "3D spatial features for multi-channel target speech separation," *ASRU*, 2021.

[26] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," in *arXiv preprint arXiv:1607.06450*, 2016.

[27] M. Yu, Y. Xu, C. Zhang, S.-X. Zhang, and D. Yu, "NeuralEcho: A self-attentive recurrent neural network for unified acoustic echo suppression and speech enhancement," in *arXiv preprint arXiv:2205.10401*, 2022.

[28] J. L. Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR - half-baked or well done?," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019.

[29] J. Du, X. Na, X. Liu, and H. Bu, "AISHELL-2: Transforming mandarin asr research into industrial scale," *arXiv:1808.10583*, 2018.

[30] D. Diaz-Guerra, A. Miguel, and J.R. Beltran, "gpuRIR: A python library for room impulse response simulation with gpu acceleration," in *Multimed Tools Appl*, 2020.

[31] ITU-T, "Recommendation p.862: Perceptual evaluation of speech quality (pesq), an objective method for endto-end speech quality assessment of narrowband telephone networks and speech codecs," 2001.

[32] E. Vincent, R. Gribonval, and C. Fevotte, "Performance measurement in blind audio source separation," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 14, no. 4, pp. 1462–1469, 2006.

[33] C. Reddy, V. Gopal, and R. Cutler, "Dnsmos p.835: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors," in *ICASSP*. IEEE, 2022.