# BREAKING THE TRADE-OFF IN PERSONALIZED SPEECH ENHANCEMENT WITH CROSS-TASK KNOWLEDGE DISTILLATION

*Hassan Taherian*[1,2*], *Sefik Emre Eskimez*[1], and *Takuya Yoshioka*[1]

[1]Microsoft, Redmond, WA, USA        [2]The Ohio State University, Columbus, OH, USA

taherian.1@osu.edu, {seeskime, tayoshio}@microsoft.com

## ABSTRACT

Personalized speech enhancement (PSE) models achieve promising results compared with unconditional speech enhancement models due to their ability to remove interfering speech in addition to background noise. Unlike unconditional speech enhancement, causal PSE models may occasionally remove the target speech by mistake. The PSE models also tend to leak interfering speech when the target speaker is silent for an extended period. We show that existing PSE methods suffer from a trade-off between speech over-suppression and interference leakage by addressing one problem at the expense of the other. We propose a new PSE model training framework using cross-task knowledge distillation to mitigate this trade-off. Specifically, we utilize a personalized voice activity detector (pVAD) during training to exclude the non-target speech frames that are wrongly identified as containing the target speaker with hard or soft classification. This prevents the PSE model from being too aggressive while still allowing the model to learn to suppress the input speech when it is likely to be spoken by interfering speakers. Comprehensive evaluation results are presented, covering various PSE usage scenarios.

***Index Terms***— personalized speech enhancement, target speech extraction, knowledge distillation.

## 1. INTRODUCTION

Remote meetings have become part of our daily lives in the rapidly emerging hybrid work era. Causal and real-time speech enhancement (SE) algorithms are now integrated into most teleconferencing services to attenuate background noise. Meanwhile, personalized speech enhancement (PSE) is gaining increased attention from the research community. PSE utilizes additional cues such as a speaker embedding vector of a target speaker to enhance only the speaker's signal even when interfering speech and background noise are both present [1, 2, 3]. The PSE task may be regarded as a combination of speech separation, enhancement, and speaker verification tasks.

Despite the advantage over SE, the current causal PSE methods face two major problems, i.e., speech over-suppression and interference leakage. Speech over-suppression refers to the problem of the target speaker's voice being identified as an interfering speaker and wrongly removed. This problem is worse for the same-gender mixtures due to voice characteristic similarities between the target and interfering speakers [4]. As reported by prior studies, speech over-suppression negatively impacts automatic speech recognition (ASR) accuracy [5] and human communication experiences [1].

The second problem, or interference leakage, means that the PSE models often fail to remove interfering speakers when the target

speaker is not present at all or for a sustained period. This problem has yet to be fully investigated, as most prior works assumed the target speaker to be actively speaking. In practical scenarios such as video conferencing, the target speaker can be inactive or silent for a long time. A naive solution for reducing the interference leakage would be to add inactive target speaker (ITS) samples in the training data [6] and train the PSE model to generate zero signals for the ITS samples. However, precisely identifying ITS frames is challenging due to the causality constraint and the model size limitation. Forcing the PSE model to generate zero signals for all the ITS frames regardless of their difficulty levels results in increased speech over-suppression. Previously proposed PSE models suffered from a trade-off between the speech over-suppression and the interference leakage by addressing only one problem at the expense of the other.
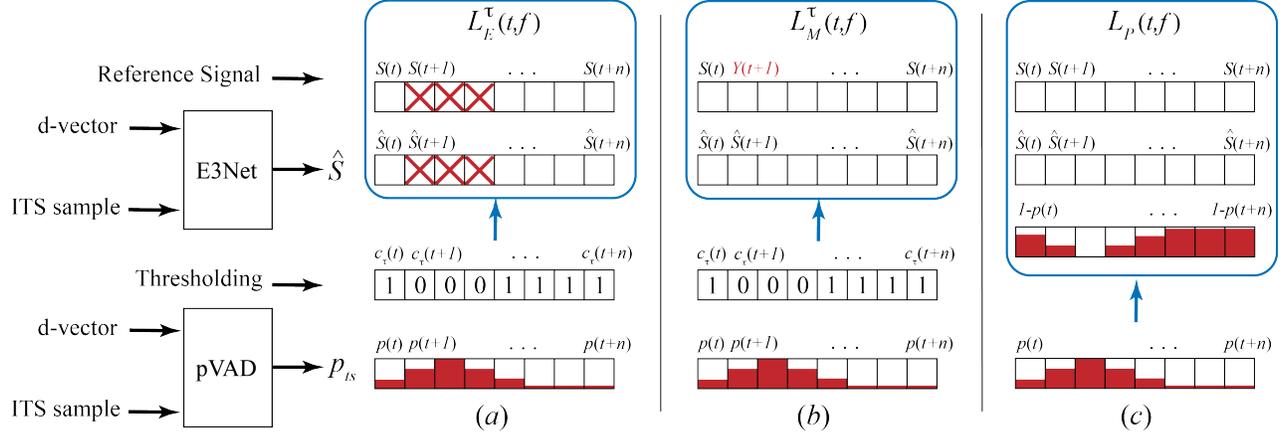
We propose a cross-task knowledge distillation approach to reduce both speech over-suppression and interference leakage and thus overcome the trade-off between these two problems. Specifically, we utilize a causal personalized voice activity detector (pVAD) to identify the frames in the ITS training samples that are wrongly classified as the target speaker (note that the ITS samples contain no target speakers). We then modify the PSE loss function based on the pVAD outputs to adjust the contribution of each frame. With the modified PSE loss, we exclude or de-emphasize the misclassified frames as these frames are difficult to handle, and including them during training can exacerbate the speech over-suppression. We show the effectiveness of our proposed training method in different scenarios.

## 2. RELATED WORK

Several studies developed causal PSE models utilizing a speaker embedding vector to extract the target speaker's voice. [1, 2, 3, 7, 8]. Giri et al. proposed a perceptually motivated PSE model with low complexity [2]. In [1], two real-time PSE models were proposed and evaluated in various scenarios. Thakker et al. introduced an efficient real-time PSE model with low computational cost [7]. [8] employed a multi-stage and multi-loss framework to train a full band PSE. In [9], a dual-stage PSE network is proposed where the target speech magnitude is estimated in the first stage, and the clean phase information is retrieved in the second stage. These studies paid limited attention to the trade-off problem between speech over-suppression and interference leakage.

Wang et al. proposed an asymmetric loss [5] for a target speaker extraction system for speech recognition to mitigate the speech over-suppression. The asymmetric loss penalizes the time-frequency bins where the target speaker's voice is over-suppressed. While it reduces speech over-suppression, the asymmetric loss significantly increases the interference leakage. [1] proposed a PSE model with ASR-based multi-task training to alleviate the speech over-suppression problem.

---

**Fig. 1**: Schematic diagram of E3Net training with cross-task knowledge distillation. (a) Misclassified frames are excluded from PSE loss. (b) Noisy signal $Y$ is used as the reference signal for misclassified frames. (c) Active target speaker probabilities are used as weights in PSE loss. In (a), crossed-out frames are excluded from the loss computation.

A few recent studies attempted to address the interference leakage problem to handle the case where the target speaker is inactive [6, 10, 11]. [11] and [10] proposed time-domain speaker extraction models with a modified signal-to-noise (SNR) ratio loss for the inactive target speaker scenario. [6] trained a target speaker extraction model with ITS samples using a modified SNR loss that preserves the input signal amplitude at the system's output. To reduce interference leakage, they utilized an extra speaker verification module to detect if the extracted speech belonged to the target speaker. However, this approach increases the computational cost and is unsuitable for real-time processing. In this paper, we address both the speech over-suppression and interference leakage problems in causal PSE without increasing the inference cost.

## 3. SYSTEM DESCRIPTION

### 3.1. Baseline PSE and Problem

We build our PSE models based on the end-to-end enhancement network (E3Net) architecture of [7] while the proposed approach is applicable to other model architectures. E3Net uses a learnable encoder and decoder. The encoded features are concatenated with a speaker embedding vector (d-vector) and fed into a stack of long short-term memory (LSTM) blocks. Each LSTM block consists of two fully connected layers, an LSTM layer with residual connection, and layer normalization modules. On top of the last LSTM block, it has a fully connected layer for generating feature masks. They are multiplied with the encoded features and transformed into a waveform with the decoder to estimate the target speaker audio. The model is trained to minimize a power-law compressed phase-aware (PLCPA) loss function, which is defined as [1]:

$$
\begin{aligned}
\mathcal{L}_{S,\hat{S}}(t,f) = \alpha \; & \left| |S(t,f)|^p - |\hat{S}(t,f)|^p \right|^2 + \\
(1-\alpha) \; & \left| |S(t,f)|^p e^{j\varphi(S(t,f))} - |\hat{S}(t,f)|^p e^{j\varphi(\hat{S}(t,f))} \right|^2,
\end{aligned}
\tag{1}
$$

where $\hat{S}(t,f)$ and $S(t,f)$ are the estimated and clean speech signals, respectively, at time $t$ and frequency $f$ in the short-time Fourier transform domain. Operator $\varphi$ calculates the argument of a complex number. The loss to be minimized is obtained by averaging $\mathcal{L}_{S,\hat{S}}$ over all time and frequency units. See [7] for further details.

E3Net is causal and was shown to achieve good accuracy with a low computational cost. In [7], as with other prior models, the E3Net model was trained by using a dataset that always contained target speech signals. However, this training scheme promotes an interference leakage behavior when the input signal does not contain the target speaker at all or for a long time.

### 3.2. PSE Training with Cross-task Knowledge Distillation

A naive approach to address the interference leakage issue would be to include ITS samples during training. An ITS sample is a noisy signal where the target speaker corresponding to the provided d-vector is completely inactive or silent. For the ITS samples, a PSE model is supposed to generate zero signals. Training the PSE model with the ITS samples helps it learn to remove the interfering speech signals when the target speaker is inactive. However, this simple approach tends to make the trained model so aggressive that it frequently attenuates the target speaker's speech signal.

Our hypothesis about the root cause of the increased speech over-suppression when the ITS samples are used during training is as follows. Due to the causality constraint and the limited model capacity for real-time operation, there will be some frames that are difficult to identify as target speech or interference. Forcing the PSE model to generate zero signals in these frames will encourage the model to occasionally zero out the signal gain even when the target speaker is present, worsening the speech over-suppression problem.

To circumvent this problem and break the trade-off between the speech over-suppression and the interference leakage, we propose a cross-task distillation approach. Specifically, we use a separately trained causal personalized voice activity detection (pVAD) model to detect the challenging frames in the ITS training samples and exclude them from the PSE loss calculation. The pVAD model performs two-way classification for each frame to produce the posterior probability of each frame being spoken by the target speaker or not. Our pVAD model is based on a modified E3Net architecture. Instead of a learnable encoder, the pVAD model adopts 40-dimensional log Mel-filterbank energies as input by using the same window and hope size as the E3Net PSE model. The masking and decoder layers are replaced with a softmax layer for classification. The interference-only time frames of the ITS samples that the pVAD model misclassifies as containing the target speech are regarded as the challenging

frames and removed from the PSE loss computation either with hard or soft decisions, as described below.

We examine three methods to modify the PSE loss for handling the misclassified frames by pVAD. In the first method, we explicitly exclude the misclassified frames from the PSE loss:

$$\mathcal{L}_E^\tau(t,f) = c_\tau(t)\mathcal{L}_{S,\hat{S}}(t,f), \quad (2)$$

where

$$c_\tau(t) = \begin{cases} 1 & \text{if } p_{ts}(t) < \tau \\ 0 & \text{otherwise} \end{cases}, \quad (3)$$

with $p_{ts}$ and $\tau$ being the frame-wise target speaker posterior probability generated by the pVAD model and a threshold, respectively. Fig. 1 depicts the diagram of the proposed loss functions. For frame $t$, the target speaker is considered active by the pVAD model when $c_\tau(t) = 0$ and inactive when $c_\tau(t) = 1$. We use Eq. (2) as the PSE loss function to exclude the contribution of the misclassified frames in the ITS samples. Alternatively, the second method leaks the misclassified frames by using the noisy signal $Y$ as the reference signal:

$$\mathcal{L}_M^\tau(t,f) = \begin{cases} \mathcal{L}_{S,\hat{S}}(t,f) & \text{if } c_\tau(t) = 1 \\ \mathcal{L}_{Y,\hat{S}}(t,f) & \text{otherwise} \end{cases}. \quad (4)$$

This method alleviates target speaker over-suppression at the cost of slightly increased interference leakage. In the third method, we adjust the contribution of each frame by using the active target speaker probability as a weight in the PSE loss calculation:

$$\mathcal{L}_P(t,f) = (1 - p_{ts}(t))\mathcal{L}_{S,\hat{S}}(t,f). \quad (5)$$

Intuitively, Eq. (5) reduces the loss contribution from misclassified frames and instead emphasizes the frames that are correctly predicted as an inactive target speaker. Unlike the previous methods, Eq. (5) does not require the threshold. Note that the proposed cross-task knowledge distillation scheme is applied only to the ITS training samples and that we use Eq. (1) as the loss function for the training samples containing the target speech.

## 4. EXPERIMENTAL RESULTS

We conducted a comprehensive test for the proposed training method by using datasets covering various conditions and performance metrics to evaluate different aspects of PSE systems.

### 4.1. Datasets

The evaluation was carried out based on simulated datasets. Room impulse responses (RIRs) were generated by using the image method with reverberation time (T60) between 0.15 and 0.6 seconds. In our simulation, we assumed the target speaker to be closer to the microphone than the interfering speaker, which seems a reasonable assumption for telecommunication applications. The target speaker's distance to the microphone was in the range of (0, 1.3] m, while the interfering speaker's distance was greater than 2 meters.

We generated 2,000 and 50 hours of audio for the training and validation datasets, respectively, based on the clean speech data of the Deep Noise Suppression challange [12]. The clean speech signals were corrupted by the simulated RIRs and the noises from the AudioSet and Freesound datasets [13, 14] with signal-to-noise ratios (SNRs) in the range of [0, 15] dB. Half of the training and validation utterances contained the target and interfering speakers as well as

noise with signal-to-interference ratios (SIRs) between 0 and 10 dB. The other half contained samples comprising the target speaker and noise only. The sampling rate was 16 kHz. The d-vectors had 128 dimensions and were extracted with a pre-trained Res2Net model (see [15] for the details). For training PNS models with ITS samples, we also randomly replaced the clean target speech with a zero signal in 15% of the above training data to simulate an inactive target speaker scenario. We refer to the training datasets without and with the ITS samples as Base and Base/ITS, respectively.

The voice cloning toolkit (VCTK) corpus was used to create the test sets. The VCTK dataset contains clean utterances of 109 speakers with different English accents. We set aside 30 utterances of each speaker for d-vector extraction. To simulate a teleconference session, we concatenated the noisy reverberant mixtures generated from the same speaker's utterances to create a single long audio file for each speaker. The average audio file duration was 27.5 minutes. The following three test sets were created to evaluate PSE models in different scenarios. TS1: the target speech signal is corrupted by both interfering speech and background noise; TS2: the target signal is corrupted by background noise; and TS3: the target speaker is inactive for the whole session and the audio file includes only interfering speakers and noise.

### 4.2. Evaluation Metrics

We used the word error rate (WER), deletion error rate (DEL), short-time objective intelligibility (STOI) [16], and DNSMOS [17] for performance measurement. DNSMOS is a neural network-based mean opinion score (MOS) estimator which was shown to be highly correlated with subjective quality ratings. To directly measure the target speech over-suppression (TSOS) at the signal level, in addition to DEL, we also used the TSOS metric proposed in [1]. For each time frame, it is defined as

$$\mathcal{TSOS}(t) = \begin{cases} 1 & \text{if } \sum_f \mathcal{L}_{OS}(t,f) > \gamma \sum_f |S(t,f)|^p \\ 0 & \text{otherwise} \end{cases}, \quad (6)$$

where $\mathcal{L}_{OS}$ represents the following over-suppression index:

$$\mathcal{L}_{OS}(t,f) = \left| \text{ReLU}(|S(t,f)|^p - |\hat{S}(t,f)|^p) \right|^2. \quad (7)$$

ReLU(.) is the rectified linear unit function, and $\gamma$ is a threshold value set at 0.1. Note that Eq. (7) is a special version of the asymmetric loss proposed in [5]. Since reference clean utterances occasionally contained modest non-speech sounds, we applied forced alignment to ignore the time frames with no speech activity. To make it easy to interpret the resultant numbers, we counted the segments where the frame-level TSOS values continued to be one for one second or longer. Finally, to measure the interference leakage in the TS3 scenario, we calculated the energy difference between the input and residual signals, i.e.,

$$\Delta N = 10\log|Y|^2 - 10\log|\hat{S}|^2. \quad (8)$$

### 4.3. Implementation Details

Following [7], we used an E3Net model consisting of 4 LSTM blocks and an encoder-decoder pair with 2,048 filters. The dimensions of the LSTM and fully connected layers of each LSTM block were 256 and 1024, respectively. We set the window size to 20 ms and the hop size to 10 ms. During the training, we generated mixtures on the fly by randomly selecting reverberated target and interfering speech signals and noise samples. We also applied a

**Table 1**: Comparison of different PSE training methods for TS1, TS2, and TS3 scenarios. TS1 includes the target, interfering speaker, and noise, while TS2 includes the target speaker and noise. TS3 includes only interfering speakers and noise. System U1 means unprocessed audio, B* are baseline PNS models trained with Eq. (1), and S* systems are based on the proposed training methods. All models used E3Net.

| System | Train data | Loss | TS1 | | | | | TS2 | | | | | TS3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | WER↓ | DEL↓ | DNSMOS↑ | STOI↑ | TSOS↓ | WER↓ | DEL↓ | DNSMOS↑ | STOI↑ | TSOS↓ | $\Delta N$↑ |
| U1 | – | – | 43.0 | 3.99 | 2.92 | 78.9 | 0.00 | 13.4 | 2.00 | 2.98 | 85.0 | 0.00 | 0.0 |
| B1 | Base | $\mathcal{L}_{S,\hat{S}}$ | 31.9 | 4.56 | 3.56 | 88.8 | 1.35 | 16.8 | 2.55 | 3.80 | 93.4 | 0.45 | 46.5 |
| B2 | Base/ITS | $\mathcal{L}_{S,\hat{S}}$ | 35.9 | 8.28 | 3.49 | 85.5 | 3.95 | 20.3 | 6.17 | 3.70 | 89.7 | 2.54 | 148.3 |
| S1 | | $\mathcal{L}_E^{0.5}$ | 34.7 | 4.46 | 3.52 | 88.6 | 1.66 | 17.8 | 2.53 | 3.75 | 93.3 | 0.37 | 148.5 |
| S2 | | $\mathcal{L}_E^{0.25}$ | 34.5 | 4.65 | 3.53 | 88.6 | 1.98 | 17.2 | 2.68 | 3.76 | 93.2 | 0.72 | 148.4 |
| S3 | Base/ITS | $\mathcal{L}_E^{0.1}$ | 35.1 | 4.90 | 3.50 | 88.2 | 2.06 | 17.8 | 3.01 | 3.74 | 92.8 | 1.15 | 148.3 |
| S4 | | $\mathcal{L}_P$ | 35.7 | 5.73 | 3.50 | 87.5 | 1.80 | 18.4 | 3.48 | 3.72 | 92.4 | 0.48 | 148.6 |
| S5 | | $\mathcal{L}_M^{0.5}$ | 35.6 | 4.73 | 3.48 | 87.8 | 1.19 | 18.1 | 2.65 | 3.70 | 92.9 | 0.32 | 145.6 |

**Table 2**: Experimental results with asymmetric loss function.

| System | Train data | Loss | TS1 | | | | | TS2 | | | | | TS3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | WER↓ | DEL↓ | DNSMOS↑ | STOI↑ | TSOS↓ | WER↓ | DEL↓ | DNSMOS↑ | STOI↑ | TSOS↓ | $\Delta N$↑ |
| B1a | Base | $\mathcal{L}_{S,\hat{S}} + \mathcal{L}_{OS}$ | 32.3 | 3.70 | 3.59 | 90.3 | 0.13 | 16.5 | 2.28 | 3.79 | 94.0 | 0.03 | 32.4 |
| B2a | Base/ITS | $\mathcal{L}_{S,\hat{S}} + \mathcal{L}_{OS}$ | 35.7 | 4.24 | 3.55 | 89.3 | 0.90 | 18.1 | 2.73 | 3.75 | 93.4 | 0.56 | 145.6 |
| S4a | Base/ITS | $\mathcal{L}_P + \mathcal{L}_{OS}$ | 34.7 | 5.56 | 3.46 | 87.1 | 1.81 | 17.9 | 2.60 | 3.70 | 93.1 | 0.27 | 148.6 |

signal-domain variant of SpecAugment [18] to input mixtures. The PLCPA loss parameters were set as $p = 0.3$ and $\alpha = 0.5$. The value of threshold $\tau$ was set at 0.5 by default for $\mathcal{L}_E^\tau$ and $\mathcal{L}_M^\tau$ losses.

Our pVAD model, used for the cross-task knowledge distillation, was based on E3Net with 40-dimensional log mel-filterbank input, three LSTM blocks, and a two-way softmax output layer. The model was trained with a binary cross-entropy loss on the Base/ITS training dataset. The ground-truth target speaker activity labels were generated by applying a DNN-based VAD model [19] to the underlying clean speech signals.

### 4.4. Results and Discussions

Table 1 shows the experimental results. Two baseline E3Net models were built based on PLCPA loss $\mathcal{L}_{S,\hat{S}}$ without the proposed cross-task scheme. One was trained on the Base training dataset (B1), and the other used the Base/ITS dataset (B2). We can observe that including the ITS samples during training significantly reduced interference leakage in the TS3 scenario. The average noise energy of TS3 was 148.8 dB, which means that the B2 model removed the noise and interference signals almost completely when the target speaker was silent. However, B2 considerably increased speech over-suppression in TS1 and TS2 compared with B1. For example, DEL and TSOS were increased by 3.62 percentage points and 2.09 seconds, respectively, in the T2 scenario. This shows the previous training scheme using PLCPA loss suffers from the trade-off between the speech over-suppression and the interference leakage.

Model S1 trained with the proposed loss of $\mathcal{L}_E^{0.5}$ yielded the DEL and TSOS values that are close to the results of B1 for both TS1 and TS2 while achieving almost the same $\Delta N$ value as B3 for TS3. This means that excluding the ITS frames misclassified by the pVAD model led to decreasing the interference leakage without incurring increases in speech over-suppression. While S1 modestly increased the WER compared with B1, especially for the TS1 scenario (31.89% → 34.67%), which might be due to increased processing artifacts, the gain (46.52 dB → 148.49 dB) in TS3 was much more significant.

The effect of the threshold in $\mathcal{L}_E^\tau$ was examined by changing the $\tau$ value to 0.25 and 0.1 (see S2 and S3). By decreasing $\tau$, more frames would be considered as being misclassified and would be excluded from the loss function. The results show that decreasing the pVAD threshold value did not lead to further speech over-suppression reduction. Instead of using hard pVAD decisions, training the PNS model with soft decisions using $\mathcal{L}_P$ loss resulted in a similar performance for all metrics (S5) without threshold adjustment. Finally, using $\mathcal{L}_M^{0.5}$ as the loss function further reduced the speech over-suppression measured by TSOS at the slight expense of the interference leakage. All results show that the proposed cross-task knowledge distillation training method based on pVAD improved the PNS performance in both speech over-suppression and leakage interference.

Table 2 shows the results of the PNS models obtained by combining the asymmetric loss of (7) with PLCPA or the proposed loss. As we can see in B1a, this improved the baseline system with respect to the speech over-suppression at the cost of increased interference leakage in TS3. Adding the ITS training samples mitigated interference leakage significantly while showing a reasonable performance with respect to the speech over-suppression (B2a). The proposed method, denoted as S4a, further reduced the interference leakage amount while the TSOS results were mixed compared with B2a (i.e., improvement was observed for TS2 while the TSOS was degraded in TS1).

### 5. CONCLUSION

In this work, we introduced a new causal PSE model training method to reduce interference leakage when the target speaker is inactive without over-suppressing the target speech. We used a pVAD task for cross-task knowledge distillation to achieve this goal. Specifically, we used misclassification patterns of a pVAD model to identify challenging frames of ITS training samples and excluded or de-emphasized them from the PNS model loss calculation. The experimental results showed the effectiveness of the proposed method.

# 6. REFERENCES

[1] S. E. Eskimez, T. Yoshioka, H. Wang, X. Wang, Z. Chen, and X. Huang, "Personalized speech enhancement: New models and comprehensive evaluation," in *Proc. ICASSP*, 2022, pp. 356–360.

[2] R. Giri, S. Venkataramani, J.-M. Valin, U. Isik, and A. Krishnaswamy, "Personalized PercepNet: Real-time, low-complexity target voice separation and enhancement," in *Proc. Interspeech*, 2021, pp. 1124–1128.

[3] H. Taherian, S. E. Eskimez, T. Yoshioka, H. Wang, Z. Chen, and X. Huang, "One model to enhance them all: Array geometry agnostic multi-channel personalized speech enhancement," in *Proc. ICASSP*, 2022, pp. 271–275.

[4] M. Delcroix, T. Ochiai, K. Zmolikova, K. Kinoshita, N. Tawara, T. Nakatani, and S. Araki, "Improving speaker discrimination of target speech extraction with time-domain speakerbeam," in *Proc. ICASSP*, 2020, pp. 691–695.

[5] Q. Wang, I. L. Moreno, M. Saglam, K. Wilson, A. Chiao, R. Liu, Y. He, W. Li, J. Pelecanos, M. Nika, and A. Gruenstein, "VoiceFilter-Lite: Streaming targeted voice separation for on-device speech recognition," in *Proc. Interspeech*, 2020, pp. 2677–2681.

[6] M. Delcroix, K. Kinoshita, T. Ochiai, K. Zmolikova, H. Sato, and T. Nakatani, "Listen only to me! How well can target speech extraction handle false alarms?" in *Proc. Interspeech*, 2022, pp. 216–220.

[7] M. Thakker, S. E. Eskimez, T. Yoshioka, and H. Wang, "Fast real-time personalized speech enhancement: End-to-end enhancement network (E3Net) and knowledge distillation," in *Proc. Interspeech*, 2022, pp. 991–995.

[8] L. Chen, C. Xu, X. Zhang, X. Ren, X. Zheng, C. Zhang, L. Guo, and B. Yu, "Multi-stage and multi-loss training for fullband non-personalized and personalized speech enhancement," in *Proc. ICASSP*, 2022, pp. 9296–9300.

[9] Y. Ju, W. Rao, X. Yan, Y. Fu, S. Lv, L. Cheng, Y. Wang, L. Xie, and S. Shang, "TEA-PSE: Tencent-Ethereal-Audio-Lab personalized speech enhancement system for ICASSP 2022 DNS challenge," in *Proc. ICASSP*, 2022, pp. 9291–9295.

[10] M. Borsdorf, C. Xu, H. Li, and T. Schultz, "Universal speaker extraction in the presence and absence of target speakers for speech of one and two talkers," in *Proc. Interspeech*, 2021, pp. 1469–1473.

[11] Z. Zhang, B. He, and Z. Zhang, "X-TaSNet: Robust and accurate time-domain speaker extraction network," in *Proc. Interspeech*, 2020, pp. 1421–1425.

[12] C. K. Reddy, H. Dubey, K. Koishida, A. Nair, V. Gopal, R. Cutler, S. Braun, H. Gamper, R. Aichner, and S. Srinivasan, "INTERSPEECH 2021 deep noise suppression challenge," in *Proc. Interspeech 2021*, 2021, pp. 2796–2800.

[13] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in *Proc. ICASSP*, 2017, pp. 776–780.

[14] E. Fonseca, J. Pons Puig, X. Favory, F. Font Corbera, D. Bogdanov, A. Ferraro, S. Oramas, A. Porter, and X. Serra, "Freesound datasets: a platform for the creation of open audio datasets," in *Proc. ISMIR*, 2017, pp. 486–493.

[15] T. Zhou, Y. Zhao, and J. Wu, "ResNeXt and Res2Net structures for speaker verification," in *Proc. SLT*, 2021, pp. 301–307.

[16] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An algorithm for intelligibility prediction of time–frequency weighted noisy speech," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, pp. 2125–2136, 2011.

[17] H. Gamper, C. Reddy, R. Cutler, I. Tashev, and J. Gehrke, "Intrusive and non-intrusive perceptual speech quality assessment using a convolutional neural network," in *Proc. WASPAA*, 2019.

[18] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "SpecAugment: A simple data augmentation method for automatic speech recognition," in *Proc. Interspeech*, 2019, pp. 2613–2617.

[19] S. Braun and I. Tashev, "On training targets for noise-robust voice activity detection," in *Proc. EUSIPCO*, 2021, pp. 421–425.