



# SNRi Target Training for Joint Speech Enhancement and Recognition

Yuma Koizumi, Shigeki Karita, Arun Narayanan, Sankaran Panchapagesan, Michiel Bacchiani

Google Research

{koizumiyuma, karita, arunnt, panchi, michiel}@google.com

## Abstract

Speech enhancement (SE) is used as a frontend in speech applications including automatic speech recognition (ASR) and telecommunication. A difficulty in using the SE frontend is that the appropriate noise reduction level differs depending on applications and/or noise characteristics. In this study, we propose “*signal-to-noise ratio improvement (SNRi) target training*”; the SE frontend is trained to output a signal whose SNRi is controlled by an auxiliary scalar input. In joint training with a backend, the target SNRi value is estimated by an auxiliary network. By training all networks to minimize the backend task loss, we can estimate the appropriate noise reduction level for each noisy input in a data-driven scheme. Our experiments showed that the SNRi target training enables control of the output SNRi. In addition, the proposed joint training relatively reduces word error rate by 4.0% and 5.7% compared to a Conformer-based standard ASR model and conventional SE-ASR joint training model, respectively. Furthermore, by analyzing the predicted target SNRi, we observed the jointly trained network automatically controls the target SNRi according to noise characteristics. Audio demos are available in our demo page<sup>1</sup>.

**Index Terms:** Speech enhancement, signal-to-noise ratio improvement, multi-task learning, noise robust ASR.

## 1. Introduction

Speech enhancement (SE) is the task of recovering target speech from a noisy signal [1]. Single-channel SE is an indispensable frontend in most speech tasks; for example, improving speech intelligibility for telecommunication [2–4], reducing noise for automatic speech recognition (ASR) [5–9] systems, and estimating spatial covariance matrix for multi-channel SE [10, 11].

A difficulty in using SE frontend is that the appropriate noise reduction level is different depending on the applications and/or noise characteristics. For example, maximizing signal-to-noise ratio (SNR) improvement does not necessarily lead to better ASR performance [9] and perceptual speech quality [12]. This is likely due to the distortions introduced by non-linear processing in single-channel SE such as time-frequency (TF) masking. While typical single-channel SE frontends aim to perfectly remove noise, in practice they cause artifacts in the resulting denoised speech.

One strategy to solve this problem is restricting SE performance [12–20]. Earlier SE studies limit noise suppression in non-speech TF bins by flooring [12, 13], smoothing [12, 14], and scaling [15, 16] the estimated TF mask in the short-time Fourier transform (STFT) domain. Since recent time-domain SE [21–25] cannot control noise reduction level by manipulating estimated masks, the observed and/or estimated noise signals are added to the estimated speech signal [19, 20]. This strategy is certainly effective, and therefore, it would be worth-

<sup>1</sup>google.github.io/df-conformer/snri\_target/

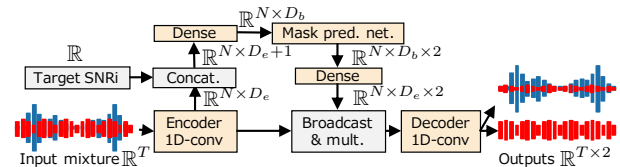


Figure 1: Network architecture of SNRi-Net trained by SNRi target training.  $N$ ,  $D_e$ , and  $D_b$  mean number of time frames, encoder/decoder basis, and bottleneck feature, respectively.

while to extend this approach to an more interpretable form, optimized in data-driven manner.

In this paper, we address the following question: *How much signal-to-noise ratio improvement (SNRi) is required in each task for a given noisy input?* To provide interpretability to the SE frontend, we propose a new framework named as *SNRi target training*. The SE frontend uses an auxiliary scalar input, which represents the target SNRi of the output signal. Instead of optimizing the enhancement to maximize the SNR, its goal is altered to produce an output signal with the specified target SNRi. We call this SE model, “*SNRi-Net*”. A block diagram of SNRi-Net is shown in Fig. 1. As a specific use case, we adopt ASR as the backend task to evaluate the merit of SNRi-Net. In joint training, the target SNRi value is estimated by an auxiliary network, which we call, *SNRi-Pred-Net*. All networks are trained to minimize ASR loss in an end-to-end manner as shown in Fig. 2. This way, we can estimate the appropriate noise reduction level for each noisy input in a data-driven manner.

Experiments show that our SNRi target training enables control of the SNRi more accurately than post-mixing of separate signals. For evaluating the proposed method as the SE frontend, our ASR system was compared with a Conformer-based standard ASR model and conventional SE-ASR joint training model. We used noisy datasets with and without reverberation, and the proposed method achieved the WER reduction by 12.5% (dry) and 4.0% (reverberant) from the standard ASR model, and 1.5% (dry) and 5.7% (reverberant) from the SE-ASR joint model, respectively. Furthermore, we observed that the predicted target SNRi was controlled dynamically according to noise characteristics. Audio demos are available in our demo page<sup>1</sup>.

## 2. Conventional Method

Let the  $T$ -sample time-domain observation  $\mathbf{x} \in \mathbb{R}^T$  be a mixture of a target speech  $\mathbf{s}$  and noise  $\mathbf{n}$ , such that,  $\mathbf{x} = \mathbf{s} + \mathbf{n}$ . The goal of standard SE is to recover  $\mathbf{s}$  from  $\mathbf{x}$ . A popular strategy in supervised SE is the time-domain mask-based method [21–25]. As an implementation example, [25] estimates masks for separating speech and noise by a mask prediction network and applies it to the representation of  $\mathbf{x}$  encoded by an encoder. The estimated signals  $\mathbf{y} \in \mathbb{R}^{T \times 2}$  are then re-synthesized using a decoder. Here,  $\mathbf{y}_{:,1}$  and  $\mathbf{y}_{:,2}$  are the estimates of speech and noise,

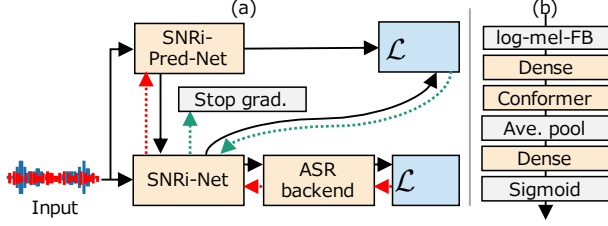


Figure 2: (a) Overview of proposed joint training of SNRi-Net and ASR backend. (b) Network architecture of SNRi-Pred-Net. Black solid lines show variable flow in forward-propagation, and red/green dotted lines show gradient flow in back-propagation.

respectively. Then, a mixture consistency projection layer [26] is applied to ensure the sum of  $\mathbf{y}_{:,1}$  and  $\mathbf{y}_{:,2}$  equals  $\mathbf{x}$ :

$$\mathbf{y}_{:,1} \leftarrow \mathbf{y}_{:,1} + \zeta \mathbf{e}, \quad \mathbf{y}_{:,2} \leftarrow \mathbf{y}_{:,2} + (1 - \zeta) \mathbf{e}, \quad (1)$$

where  $\mathbf{e} = \mathbf{x} - (\mathbf{y}_{:,1} + \mathbf{y}_{:,2})$ , and  $\zeta \in [0, 1]$  is a tunable hyperparameter. The negative thresholded SNR [24] is used as the loss:

$$\mathcal{L}^{\text{SE}} = \alpha \mathcal{L}_{s, \mathbf{y}_{:,1}}^{\text{SNR}} + (1 - \alpha) \mathcal{L}_{n, \mathbf{y}_{:,2}}^{\text{SNR}}, \quad (2)$$

$$\mathcal{L}_{a,b}^{\text{SNR}} = -10 \log_{10} (\|\mathbf{a}\|^2 / (\|\mathbf{a} - \mathbf{b}\|^2 + \tau \|\mathbf{a}\|^2)), \quad (3)$$

where  $\|\cdot\|$  is  $\ell_2$  norm and  $\tau = 10^{-3}$  is a soft threshold that clamps the loss at 30 dB [24] and  $\alpha \in [0, 1]$  is a tunable hyperparameter [24].

The SE frontend can be joined with a backend such as ASR [8, 9]. In a typical framework, the estimated speech  $\mathbf{y}_{:,1}$  is fed to the backend, then both frontend and backend are jointly fine-tuned by minimizing a loss function. In noise robust ASR tasks, the loss function  $\mathcal{L}$  can be a weighted sum of the ASR loss  $\mathcal{L}^{\text{ASR}}$  and  $\mathcal{L}^{\text{SE}}$  as

$$\mathcal{L} = \mathcal{L}^{\text{ASR}} + \gamma \mathcal{L}^{\text{SE}}, \quad (4)$$

where  $\gamma \geq 0$  is a tunable hyperparameter.

In practice,  $\mathbf{y}_{:,1}$  includes artifacts in speech. Such distortions cause degradation of the backend tasks such as ASR. To reduce distortion, several post-processing methods have been proposed [15–20]. For time-domain speech separation [19] and SE [20], a possible strategy is to add  $\mathbf{y}_{:,2}^2$ . Since SE estimates of clean speech  $\mathbf{y}_{:,1}$  and noise  $\mathbf{y}_{:,2}$ , we can control noise reduction level as

$$\mathbf{y} = \mathbf{y}_{:,1} + w \mathbf{y}_{:,2}, \quad (5)$$

where  $w \in [0, 1]$ . By assuming that the SE module perfectly separate speech and noise, SNRi of  $\mathbf{y}$  can be controlled as  $w = 10^{-\lambda/20}$ , where  $\lambda$  is a target SNRi scalar.

### 3. Proposed Method

#### 3.1. SNRi target training

In contrast the conventional SE, the goal of SNRi target training is to control SNRi of the SE output  $\mathbf{y}_{:,1}$  according to  $\lambda$ . To achieve this,  $\lambda$  is also input to the mask prediction network as an auxiliary variable. Specifically, we concatenate  $\lambda$  to the encoder

<sup>2</sup>In conventional studies [19,20],  $\mathbf{x}$  is added instead of  $\mathbf{y}_{:,2}$ , which is the same as in Eq. (5) except for the constant multiplication. To clarify the relationship with SNRi, we formulate the post-mixing using  $\mathbf{y}_{:,2}$ .

output in the feature dimension as shown in Fig. 1. We call this network as SNRi-Net. After applying a mixture consistency projection layer [26], the loss value is calculated as the squared-error between the target SNRi,  $\lambda$ , and SNRi of the output signal as:

$$\mathcal{L}^{\text{SNRi}} = |\lambda - \text{SNRi}|^2 + \beta \mathcal{L}^{\text{SAR}}, \quad (6)$$

$$\text{SNRi} = 10 \log_{10} \left( \frac{\|\mathbf{s}\|^2}{\|\mathbf{y}_{:,1} - \mathbf{s}\|^2} \right) - 10 \log_{10} \left( \frac{\|\mathbf{s}\|^2}{\|\mathbf{n}\|^2} \right), \quad (7)$$

$$\mathcal{L}^{\text{SAR}} = -10 \log_{10} (\|\mathbf{s}\|^2 / (\|\mathbf{e}_{\text{artif}}\|^2 + \tau \|\mathbf{s}\|^2)), \quad (8)$$

where  $\beta \in [0, 1]$  is a weight parameter and  $\mathcal{L}^{\text{SAR}}$  is the negative thresholded source-to-artifact ratio (SAR) which reduces artifacts in the output signal based on the SAR [27].  $\mathbf{e}_{\text{artif}}$  are the artifacts in the output signal which can be obtained by decomposing the residual noise  $\mathbf{y}_{:,1} - \mathbf{s} = \mathbf{e}_{\text{interf}} + \mathbf{e}_{\text{artif}}$  where  $\mathbf{e}_{\text{interf}}$  is the orthogonal projection of the residual noise onto the subspace spanned by both  $\mathbf{s}$  and  $\mathbf{n}$  [27].

#### 3.2. Joint training of SNRi-Net and ASR backend

SNRi-Net can also be jointly trained with a backend task, like ASR. The loss function for joint-training then becomes:

$$\mathcal{L} = \mathcal{L}^{\text{ASR}} + \eta \mathcal{L}^{\text{SNRi}}, \quad (9)$$

where,  $\eta \geq 0$  is a tunable hyperparameter. The problem of the joint network is that the appropriate target SNRi  $\lambda$  is unknown for each task and/or input. To address this, we propose using an auxiliary network called as SNRi-Pred-Net, which predicts a  $\lambda$  as  $\hat{\lambda}$  that minimizes  $\mathcal{L}^{\text{ASR}}$ . Fig. 2 (a) shows the overview of the proposed joint training framework. The architecture of SNRi-Pred-Net used in this study is shown in Fig. 2 (b).

Log-mel filterbank is used as inputs by SNRi-Pred-Net. The input features are first passed to a dense layer for computing bottleneck features, which are then processed by Conformer blocks [28]. We apply average-pooling across time to the the output of the Conformer to obtain a single, vector representation for the inputs. A second dense layer followed by a sigmoid function converts the vector to a scalar, which is the predicted target SNRi. A scaling operation is also used to restrict the predicted target SNRi to the range  $[\lambda_{\min}, \lambda_{\max}]$ . Here,  $\lambda_{\min}$  and  $\lambda_{\max}$  are hyperparameters that represents the minimum and maximum values of target SNRi  $\hat{\lambda}$ .

To enhance the noisy signal,  $\mathbf{x}$  and  $\hat{\lambda}$  are passed to SNRi-Net, and its speech output  $\mathbf{y}_{:,1}$  is passed to the ASR backend. We train the joint network by minimizing Eq. (9). That is, the ASR model, SNRi-Net and SNRi-Pred-Net are all trained to minimize ASR loss,  $\mathcal{L}^{\text{ASR}}$ . Importantly, only SNRi-Net uses the SE loss,  $\mathcal{L}^{\text{SNRi}}$ ; we stop the gradient of  $\mathcal{L}^{\text{SNRi}}$  from optimizing SNRi-Pred-Net. Since the predicted  $\lambda$  is the prediction target for SNRi-Net, if the SE loss is back-propagated to SNRi-Pred-Net, this network will be optimized to output  $\lambda$  that makes the task of SNRi-Net easy, but sub-optimal for ASR.

## 4. Experiment

#### 4.1. Experimental settings

**Dataset:** We used the same training dataset as [29] with and without reverberation. The dataset without reverberation is called ‘‘Dry’’ and with reverberation as ‘‘Reverb’’. This dataset includes 4,249 hours of clean speech which consists of 281k utterances from the LibriSpeech [30] training set and 1,916k utterances from an internal dataset. The noisy utterances

were generated using a room simulator [31], with SNR from -10 dB to 30 dB. Noise is sampled from internally collected noise snippets that simulate conditions like cafe, kitchen, and cars, and freely available noise sources from Getty [32] and YouTube Audio Library [33]. For Reverb, the room configurations have reverberation times (RT60) ranging from 0 ms to 900 ms. For Dry, RT60 is always 0 ms. We generated multiple copies of the data under different mixing conditions in order to model enough combinations of clean speech, background noise, and room-configuration. The training dataset includes 39.6M noisy samples (55,027 hours) in total.

For the validation and test datasets, we used utterances from the LibriSpeech [30] dev and test sets. The noisy utterances were generated using the same manner as the training dataset with SNR from -5 dB to 20 dB, but with noise sources that are held out from the training dataset. We generated “Dry” and “Reverb” dev and test sets using the same T60 condition as the training set. The test dataset consists of 2,620 utterances and the average SNR of the input signals was 7.5 dB.

**Models and hyper-parameters:** As SNRi-Net, we used the same architecture of DF-Conformer-8 [25] except for concatenating the target SNRi. All hyper-parameters were the same as that used in [25]; the number of Conformer layers was 8,  $D_e = 256$ ,  $D_b = 216$ , and the window and hop sizes of filterbanks were 2.5 ms and 1.25 ms. Other hyper-parameters were decided based on the WER on the validation dataset as  $\beta = 0.01$ ,  $\lambda_{\min} = 0.0$  and  $\lambda_{\max} = 20.0$ .

We use the large size Conformer Transducer ASR backend referred as Conformer(L) [28] that has 17 Conformer blocks of 512-dim 8-head dot-product attention, 512-dim 32-frame kernel 1d convolution, 2048-dim feedforward module, and a 640-dim long short-term memory (LSTM) decoder. The input feature was 80-dim filterbank features computed from a 25 ms window with a stride of 10 ms. We use SpecAugment [34] with mask parameter ( $F = 27$ ), and 10 time-masks with maximum time-mask ratio ( $p_S = 0.05$ ). All transcriptions were tokenized with a word piece model with a 1,024 vocabulary built from LibriSpeech 960h. The ASR loss  $\mathcal{L}^{\text{ASR}}$  was the recurrent neural network transducer (RNN-T) loss [35, 36].

SNRi-Pred-Net had 2-Conformer-blocks where each block has 128-dim 6-head dot-product attention, 128-dim 5-frame kernel 1d convolution, and 512-dim feedforward module. The parameter of mel-filterbank was the same as the ASR backend.

SNRi-Net and ASR backend were pre-trained for 200k steps individually; SNRi-Net was pre-trained to minimize  $\mathcal{L}^{\text{SNRi}}$  with the training dataset, and the ASR backend was pre-trained to minimize  $\mathcal{L}^{\text{ASR}}$  using both clean and noisy speech in the training dataset. While pre-training of SNRi-Net,  $\hat{\lambda}$  was randomly drawn from the uniform distribution  $\mathcal{U}(\lambda_{\min}, \lambda_{\max})$ .

The joint network was fine-tuned for an additional 100k steps to minimize Eq. (9) where  $\eta = 0.01$  was determined based on WER on the validation dataset. In the fine-tune stage, we skipped SNRi-Net and SNRi-Pred-Net with a probability of 5% as a multi-condition learning strategy [8]. In addition, we used a random  $\hat{\lambda}$  drawn from  $\mathcal{U}(\lambda_{\min}, \lambda_{\max})$  instead of the predicted one with 25% probability. All training used the Adam optimizer [37] with the same setting as [25] except for using 1/10 learning rate in the fine-tuning stage, and 128 Google TPUv3 cores with a global batch size of 512.

## 4.2. Evaluation for control accuracy of output SNRi

We compared SNRi target training with post-mixing based SNRi control in Eq. (5). We used DF-Conformer-8 [25] with

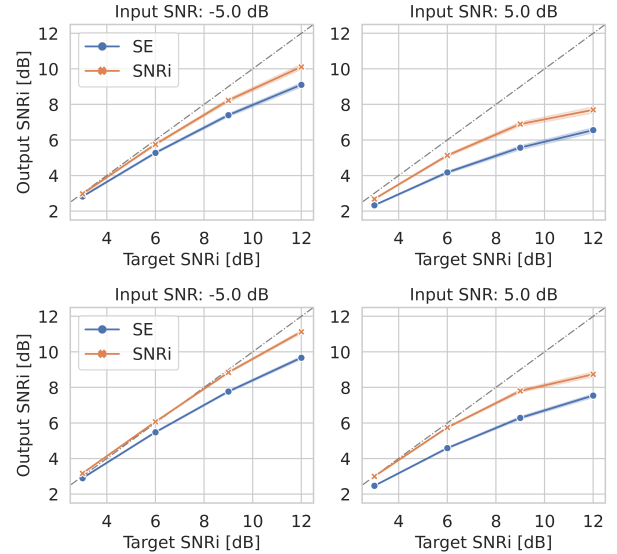


Figure 3: Results of experiment on controlling SNRi of output signal. Top and bottom figures are results on Dry and Reverb test sets, respectively. Solid lines and colored are show mean and 99 % confidence intervals, respectively. Legends SE and SNRi mean SE+post-mixing and SNRi-Net, respectively.

hyper-parameters set to the same values as [25], i.e.  $\alpha = 0.8$  and  $\zeta = 0.5$ . We tested both methods with two input SNR conditions; SNR of all test samples were adjusted to -5 dB and 5 dB. The target SNRi were 3 dB, 6 dB, 9 dB, and 12 dB.

Figure 3 shows the experimental results and audio examples of SNRi-Net are available in our demo page<sup>1</sup>. SNRi-Net trained by SNRi target training achieved significantly better control accuracy of output SNRi than supervised SE on both Dry and Reverb test sets. Although the post-mixing-based control assumes that the output signals are perfectly separated, the outputs usually contains separation errors. Such separation errors affects the output SNRi of the post-mixed signal, resulting in the under-separation problem; the output SNRi is always lower than the target SNRi. SNRi target training has succeeded to avoid the under-separation problem by directly controlling output SNRi in 3 dB and 6 dB target SNRi conditions.

Whereas in 9 dB and 12 dB target SNRi conditions, output SNRi of both methods are significantly lower than target SNRi, even though SNRi-Net was still better than the supervised method. It might be due to the performance upper-bound of the base SE network. It is necessary to accurately distinguish speech and noise to output high SNR signals, and the accuracy should correlate on the separation performance of the SE network. Therefore, to achieve higher control performance in high target SNRi conditions, it is necessary to improve the base SE network performance.

## 4.3. Evaluation as ASR frontend

We compared the ASR performance of the proposed joint-network on WER metrics with two standard ASR models and an SE-ASR joint model. The first model is Conformer(L) [28], and the second model is Conformer(L)+ which has 19 Conformer blocks and whose model size is roughly the same as the proposed method. These models give us the baseline WER of using a large ASR models without SE frontend. These models were trained for 300k steps from scratch. The third model

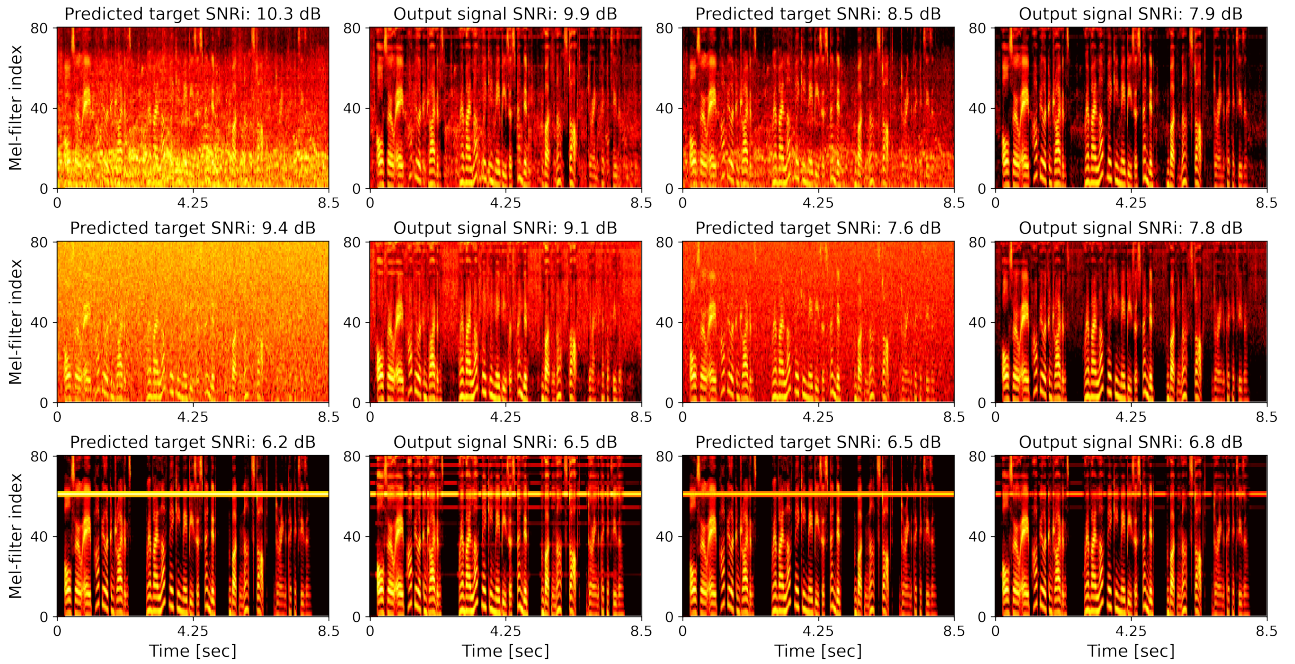


Figure 4: Log-mel spectrogram examples of input and output of joint trained model. Input SNR of first two column figures are  $-5.0$  dB, and later two figures are  $5.0$  dB. In the same input SNR figures, left figures are input signal, and right ones are output signal. Noise type of top, middle, and bottom figures are environmental noise, white noise, and  $4$  kHz sine wave.

Table 1: Word error rate on test dataset and number of trainable parameters for each model.

Network	#Params	WER [%] ( $\downarrow$ )	
		Dry	Reverb
Conformer(L)	118.8M	10.5	14.4
Conformer(L)+	132.3M	10.7	14.9
SE+Conformer(L)	127.4M	9.3	14.7
SNRi+Conformer(L)	128.3M	<b>9.2</b>	<b>13.8</b>

jointly trains the supervised SE frontend and Conformer(L) ASR model. This model gives us WER of the joint training strategy of an SE frontend and a large ASR backend. The network architecture of the SE frontend was DF-Conformer-8 used in Sec. 4.2.  $\gamma = 0.25$  was determined based on WER on the validation dataset. In addition, we skipped the SE frontend with 50% probability in the fine-tuning stage as a multi-condition learning strategy [8]. In this experiment, we excluded Eq. (5)-based joint models [19, 20] for fair comparison because these methods did not fine-tune the joint model.

Table 1 shows WER on the test datasets. The proposed method achieved the best WER performance on both datasets; it reduced WER by 12.5% and 1.5% compared to Conformer(L) and SE+Conformer(L) for `Dry`, and 4.0% and 5.7% `Reverb`, respectively. In addition, SE-Conformer(L) increased WER on `Reverb` compared to Conformer(L), whereas SNRi-Conformer(L) improved WER. From these results, the proposed method consistently improved the ASR accuracy comparing to the models which merely increase the ASR backend model size and/or jointly training with a conventional SE frontend.

Although the performance gains in terms of ASR quality from the SE+Conformer(L) are limited, we observed performance improvements in all conditions. Moreover, the model provides a novel way of controlling the SNRi, thereby providing

new insights on how to model SE in a larger system. Figure 4 shows log-mel spectrogram examples of the input and output of the proposed joint model. These examples indicate that the predicted SNRi is affected by two factors; input SNR and noise type. By comparing input SNR at  $-5$  dB and  $5$  dB cases, the predicted SNRi tends to be high when the input SNR is low. This is an intuitive result due to the fact that the larger noise makes ASR more difficult. In the case of noise types that only affects a certain frequency, such as a tonal noise, the predicted SNRi tends to be low even if the input SNR is low. As can be seen from the bottom spectrograms, if the majority of the frequency bands are clean, the speech characteristics can be analyzed even if the noise is large. These results show that the proposed method also provides interpretable insights for future development of single-channel SE frontends, e.g. tuning  $w$  in Eq. (5) using a small scale dataset without joint training [20].

## 5. Conclusion

We proposed “SNRi target training”; the SE frontend uses an auxiliary scalar input which represents target SNRi, and enhances the input so that SNRi of the output signal is close to the desired target value. In joint training with the ASR backend, the target SNRi value was also estimated to minimize the ASR loss. Experiments showed that SNRi-Net controls the SNRi more accurately than post-mixing of separated signals, and our joint ASR system achieved the best WER on noisy ASR datasets. Furthermore, by analyzing the output of the joint model, we observed the model automatically controls the target SNRi according to noise characteristics.

The limitation of the proposed method is the requirement of clean speech, in contrast to several conventional joint training strategies that only need noisy speech and its transcription [8, 19, 38]. Future work will include fine-tuning only SNRi-PredNet [19] and/or incorporating an unsupervised SE training [24].

## 6. References

- [1] D. Wang and J. Chen, "Supervised speech separation based on deep learning: An overview," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 26, no. 10, p. 1702–1726, 2018.
- [2] B. C. J. Moore, "Speech processing for the hearing-impaired: successes, failures, and implications for speech mechanisms," *Speech Commun.*, 2003.
- [3] D. L. Wang, "Time-frequency masking for speech separation and its potential for hearing aid design," trends in amplification," *Trends Amplif.*, 2008.
- [4] C. K. A. Reddy, H. Dubey, V. Gopal, R. Cutler, S. Braun, H. Gampfer, R. Aichner, and S. Srinivasan, "ICASSP 2021 deep noise suppression challenge," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2021.
- [5] M. Fujimoto, and Y. Ariki, "Robust speech recognition in additive and channel noise environments using GMM and EM algorithm," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2004.
- [6] F. Weninger, H. Erdogan, S. Watanabe, E. Vincent, J. Le Roux, J. R. Hershey, and B. Schuller, "Speech enhancement with lstm recurrent neural networks and its application to noise-robust ASR," in *Proc. Int. Conf. Latent Var. Anal. Signal Sep. (LVA/ICA)*, 2015.
- [7] H. Erdogan, J. R. Hershey, S. Watanabe, and J. Le Roux, "Phase-sensitive and recognition-boosted speech separation using deep recurrent neural networks," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2015.
- [8] K. Kinoshita, T. Ochiai, M. Delcroix, and T. Nakatani, "Improving noise robust automatic speech recognition with single-channel time-domain enhancement network," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2020.
- [9] C. Li, J. Shi, W. Zhang, A. S. Subramanian, X. Chang, N. Kamo, M. Hira, T. Hayashi, C. Boeddeker, Z. Chen, and S. Watanabe, "ESPnet-SE: end-to-end speech enhancement and separation toolkit designed for ASR integration," in *Proc. IEEE Spok. Lang. Technol. Workshops (SLT)*, 2021.
- [10] H. Erdogan, J. R. Hershey, S. Watanabe, M. I. Mandel, and J. Le Roux, "Improved MVDR beamforming using single-channel mask prediction networks," in *Proc. Interspeech*, 2016.
- [11] T. Higuchi, K. Kinoshita, N. Ito, S. Karita, and T. Nakatani, "Frame-by-frame closed-form update for mask-based adaptive MVDR beamforming," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2018.
- [12] Y. Koizumi, K. Niwa, Y. Hioaka, K. Kobayashi, and Y. Haneda, "DNN-based source enhancement to increase objective sound quality assessment score," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 26, no. 10, pp. 1780–1792, 2018.
- [13] I. Cohen, "Optimal speech enhancement under signal presence uncertainty using log-spectral amplitude estimator," *IEEE Signal Process. Lett.*, 2002.
- [14] E. Vincent, "An experimental evaluation of wiener filter smoothing techniques applied to under-determined audio source separation," in *Proc. Int. Conf. Latent Var. Anal. Signal Sep. (LVA/ICA)*, 2010.
- [15] A. Narayanan and D. Wang, "Joint noise adaptive training for robust automatic speech recognition," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2014.
- [16] Q. Wang, K. A. Lee, T. Koshinaka, K. Okabe, and H. Yamamoto, "Task-aware warping factors in mask-based speech enhancement," in *Proc. Eur. Signal Process. Conf. (EUSIPCO)*, 2021.
- [17] T. Menne, R. Schlüter, and H. Ney, "Investigation into joint optimization of single channel speech enhancement and acoustic modeling for robust ASR," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2019.
- [18] A. Pandey, C. Liu, Y. Wang, and Y. Saraf, "Dual application of speech enhancement for automatic speech recognition," in *IEEE Spok. Lang. Technol. Workshop (SLT)*, 2021.
- [19] H. Sato, T. Ochiai, M. Delcroix, K. Kinoshita, N. Kamo, and T. Moriya, "Learning to Enhance or Not: Neural network-based switching of enhanced and observed signals for overlapping speech recognition," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2022.
- [20] K. Iwamoto, T. Ochiai, M. Delcroix, R. Ikeshita, H. Sato, S. Araki, and S. Katagiri, "How Bad Are Artifacts?: Analyzing the impact of speech enhancement errors on ASR," *arXiv:2201.06685*, 2022.
- [21] Y. Luo and N. Mesgarani, "Conv-TasNet: Surpassing ideal time-frequency magnitude masking for speech separation," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [22] I. Kavalerov, S. Wisdom, H. Erdogan, B. Patton, K. Wilson, J. Le Roux, and J. R. Hershey, "Universal sound separation," in *Proc. IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA)*, 2019.
- [23] M. Pariente, S. Cornell, A. Deleforge, and E. Vincent, "Filterbank design for end-to-end speech separation," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2020.
- [24] S. Wisdom, E. Tzinis, H. Erdogan, R. Weiss, K. Wilson, and J. R. Hershey, "Unsupervised sound separation using mixture invariant training," in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, 2020.
- [25] Y. Koizumi, S. Karita, S. Wisdom, H. Erdogan, J. R. Hershey, L. Jones, and M. Bacchiani, "DF-Conformer: Integrated architecture of Conv-TasNet and Conformer using linear complexity self-attention for speech enhancement," in *IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA)*, 2021.
- [26] S. Wisdom, J. R. Hershey, K. Wilson, J. Thorpe, M. Chinen, B. Patton, and R. A. Saurous, "Differentiable consistency constraints for improved deep speech enhancement," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2020.
- [27] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR-Half-baked or well done?" in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2019.
- [28] A. Gulati, C.-C. Chiu, J. Qin, J. Yu, N. Parmar, R. Pang, S. Wang, W. Han, Y. Wu, Y. Zhang, and Z. Zhang, "Conformer: Convolution-augmented transformer for speech recognition," in *Proc. Interspeech*, 2020.
- [29] A. Narayanan, C.-C. Chiu, T. O'Malley, Q. Wang, and Y. He, "Cross-attention Conformer for context modeling in speech enhancement for ASR," in *Proc. IEEE Autom. Speech Recognit. Underst. Workshop (ASRU)*, 2021.
- [30] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. Int. Conf. on Acoust., Speech, and Signal Process. (ICASSP)*, 2015.
- [31] C. Kim, A. Misra, K. Chin, T. Hughes, A. Narayanan, T. Sainath, and M. Bacchiani, "Generation of large-scale simulated utterances in virtual rooms to train deep-neural networks for far-field speech recognition in Google Home," in *Proc. Interspeech*, 2017.
- [32] <https://www.gettyimages.com/about-music>.
- [33] <https://youtube.com/audiolibrary>.
- [34] 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," 2019.
- [35] A. Graves, "Sequence transduction with recurrent neural networks," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2012.
- [36] S. Karita, Y. Kubo, M. Bacchiani, and L. Jones, "A comparative study on neural architectures and training methods for Japanese speech recognition," in *Proc. Interspeech*, 2021.
- [37] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic optimization," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2015.
- [38] S. E. Eskimez, X. Wang, M. Tang, H. Yang, Z. Zhu, Z. Chen, H. Wang, and T. Yoshioka, "Human listening and live captioning: Multi-task training for speech enhancement," in *Proc. Interspeech*, 2021.