Mask scalar prediction for improving robust automatic speech recognition

Arun Narayanan, James Walker, Sankaran Panchapagesan, Nathan Howard, Yuma Koizumi

Google LLC, U.S.A.

{arunnt, jswalker, panchi, ndhoward, koizumiyuma}@google.com

Abstract

Using neural network based acoustic frontends for improving robustness of streaming automatic speech recognition (ASR) systems is challenging because of the causality constraints and the resulting distortion that the frontend processing introduces in speech. Time-frequency masking based approaches have been shown to work well, but they need additional hyperparameters to scale the mask to limit speech distortion. Such mask scalars are typically hand-tuned and chosen conservatively. In this work, we present a technique to predict mask scalars using an ASR-based loss in an end-to-end fashion, with minimal increase in the overall model size and complexity. We evaluate the approach on two robust ASR tasks: multichannel enhancement in the presence of speech and non-speech noise, and acoustic echo cancellation (AEC). Results show that the presented algorithm consistently improves word error rate (WER) without the need for any additional tuning over strong baselines that use hand-tuned hyper-parameters: up to 16% for multichannel enhancement in noisy conditions, and up to 7% for AEC.

Index Terms: speech recognition, time-frequency masking, speech enhancement, acoustic echo cancellation

1. Introduction

As performance of automatic speech recognition (ASR) systems improved over the years [1, 2, 3], the number of applications that use speech as a standard modality of input has increased. With varying uses and high user expectations, an important goal is to ensure that the performance does not deteriorate significantly in harsh acoustic conditions - something current ASR models still struggle with [4]. Such conditions can be the result of significant environmental noise or competing speech. With increasing focus on building large scale, general purpose, multidomain [5] and multilingual ASR models [6, 7], addressing background noise together with variations in domain and language in the same model can lead to additional complexity in training and maintaining ASR models. While data augmentation strategies like SpecAug [8] and multi-condition training [9] help to an extent, it is often advantageous to address environmental noise, device echo and competing speech using dedicated modules. A number of techniques have therefore been proposed in the literature [10].

Time-frequency masking based techniques are commonly used to build acoustic frontends for ASR [11]. In a single-channel setting, a time-frequency mask, either in the complex spectral domain [12] or the feature domain [11], is used to estimate clean speech, which is then passed on to the backend ASR model. Alternatively, frontend processing can be done directly in the time-domain [13, 14]. Typically, such frontends introduce speech distortion. This is partly due to the mismatch

The authors thank Tom O'Malley, Joe Caroselli, and Alex Park.

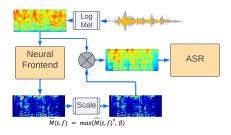


Figure 1: A block diagram of an acoustic frontend for ASR.

in the training criterion, which optimizes some measure of distance between clean and noisy speech, and the final goal, which is ASR. The constraints of streaming, which limits the model's architecture to be unidirectional and causal, makes the task even more challenging. If multiple microphones are available, the time-frequency mask can be used to compute beamforming filters to remove noise. This limits the amount of distortion, and usually yields better results [15]. Nonetheless, for short queries in a streaming setting, the gains are limited even with multiple microphones [16].

Several techniques have been proposed to reduce speech distortion when using an acoustic frontend for ASR. Menne et al. [17] propose estimating the per-frame parameters of a Wiener filter and the noise mask, and joint training with a hybrid ASR model. In [18], the input to the ASR model is a weighted sum of noisy and enhanced speech features. The weights are computed by a neural network, and jointly trained with ASR. Sato et al. [19] recommend skipping enhancement in less noisy conditions, which they later extend using a learned weighted sum of noisy and enhanced speech [20]. Adding back a scaled version of noisy data is also adopted in [21] to address distortions. Koizumi et al. [22] propose using a signal-to-noise ratio improvement (SNRi) target when doing enhancement in the waveform domain. SNRi is predicted by a separate network and optimized using an ASR loss. While most algorithms jointly optimize the the enhancement model and the ASR model [17, 18, 22], there are also approaches that freeze the ASR model [23]. We will also use this strategy, since we assume that the ASR model is trained to cover several use cases, and cannot easily be jointly optimized with the enhancement model.

When using a time-frequency mask, one way to limit distortion is to post-process the mask using a mask scalar and a mask floor before using it to enhance the features [24, 25]. The mask scalar exponentially scales the mask, which reduces speech distortion, but retains residual noise. The mask floor limits noise attenuation. One drawback of this approach is that these hyperparameters have to be hand-tuned. Moreover, a single value is chosen in the end, irrespective of the noise condition. In practice, the best value depends on the amount of noise; e.g., in clean conditions, the mask scalar should be close to 0.0 so that the enhanced features are close to the original "noisy" features, unlike in noisier conditions. Furthermore, we hypothesize that

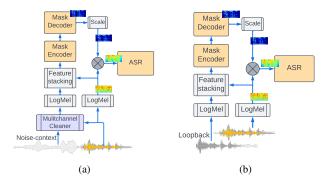


Figure 2: (a) Cleanformer based enhancement frontend for ASR; (b) Neural acoustic echo cancellation frontend for ASR.

it may even be preferable to allow the mask scalar to be dynamic, so that its values change from frame-to-frame depending on the level of noise. In this work, therefore, we propose using a *mask scalar net* to predict their optimal values. Similar to [22], mask scalar prediction is optimized exclusively using an ASR loss. Our results show that, on two separate robustness tasks – speech enhancement and acoustic echo cancellation – using predicted mask scalars provide significant improvements in WER with almost negligible increase model size.

The rest of the paper is organized as follows. Sec. 2 provides a detailed description of the proposed models. Experimental settings are described in Sec. 3 and results in Sec. 4. Finally, conclusions and future work are presented in Sec. 5.

2. System

2.1. Masking-based ASR frontend

A masking-based acoustic frontend for ASR is shown in Fig. 1. The frontend operates in the log Mel magnitude spectral domain, to match the features used by ASR. This not only simplifies joint training with ASR, which is needed for the best performance [24], but also circumvents the need to estimate phase, in contrast to complex spectral or waveform domain models.

Given a noisy Mel spectrogram, \mathbf{Y} , the goal is to estimate the clean Mel spectrogram, \mathbf{X} . This is done via an estimate of the ideal ratio mask, \mathbf{M} , which is the ratio of speech to mixture Mel magnitudes at each time-frequency bin, assuming that speech and noise are uncorrelated: $\mathbf{M}(t,c) = \frac{\mathbf{X}(t,c)}{\mathbf{X}(t,c)+\mathbf{N}(t,c)}$. Here, \mathbf{N} is the noise Mel spectrogram, and t, c are time and frequency indices. We assume that $\mathbf{Y} \approx \mathbf{X} + \mathbf{N}$, so that $\mathbf{M} \in [0,1]$, which simplifies estimation [11]. The mask is estimated using a neural net and post-processed to limit distortions:

$$\overline{\mathbf{M}}(t,c) = \max(\widehat{\mathbf{M}}(t,c)^{\alpha}, \beta), \tag{1}$$

 $\widehat{\mathbf{M}}$ and $\overline{\mathbf{M}}$ are the estimated and post-processed masks, respectively. α and β are the hyper-parameters for post-processing: $\alpha \in [0,1]$ exponentially scales the mask, and β floors the mask. The estimated clean Mel spectrogram, $\widehat{\mathbf{X}} = \mathbf{Y} \odot \overline{\mathbf{M}}$, where \odot stands for point-wise multiplication. By using a value of $\alpha < 1$ and $\beta > 0$, we can trade-off speech distortion and residual noise. $\widehat{\mathbf{X}}$ is log compressed, mean-variance normalized, and, optionally, stacked and subsampled [26] to create input ASR features.

We look at two specific frontends in this work: a multichannel enhancement frontend to remove background noise, like a T.V. playing in the background, and an acoustic echo cancellation frontend for removing device echo during playback.

2.1.1. Speech Enhancement frontend

The multichannel enhancement frontend we consider in this work is the recently proposed *Cleanformer* [27], as shown in Fig. 2a. Cleanformer uses Multichannel Cleaner (McCleaner) [28], which estimates an enhanced waveform from multichannel noisy audio and the noise context (a short segment of noise preceding the input to be recognized). McCleaner achieves excellent noise reduction, but also introduces distortion, making it less suited for ASR [27]. As input, Cleanformer concatenates the log Mel features computed from the output of McCleaner and one of the unprocessed microphone channels. It estimates M using a neural net, which consists of a Conformer encoder [29] and a fully connected layer (FCLayer) as the mask decoder. The FCLayer uses a sigmoid activation function:

$$\widehat{\mathbf{M}}(t,\cdot) = \operatorname{sigmoid}(\operatorname{FCLayer}(\mathbf{e}_t; \theta_{\mathbf{M}})). \tag{2}$$

 \mathbf{e}_t is the output of the encoder t and $\theta_{\mathbf{M}}$ the parameters of the FCLayer. Eq. 1 is then used for post-processing $\widehat{\mathbf{M}}$.

The model is trained using a combination of direct mask loss, $\ell_{\mathbf{M}}$, and an ASR loss, ℓ_{ASR} , [23], using a Conformer-based ASR model [1]:

$$\ell_{\mathbf{M}} = \|\mathbf{M} - \widehat{\mathbf{M}}\|_{1} + \|\mathbf{M} - \widehat{\mathbf{M}}\|_{2}^{2}$$

$$\ell_{ASR} = \|\mathbf{E}_{ASR}(\mathbf{f}_{X}) - \mathbf{E}_{ASR}(\mathbf{f}_{\widehat{X}})\|_{2}^{2},$$

$$\ell = \ell_{\mathbf{M}} + \lambda_{ASR} \times \ell_{ASR}.$$
(3)

Here, $\|\mathbf{A}\|_p = (\sum_{i,j} |\mathbf{A}(i,j)|^p)^{\frac{1}{p}}$ is the entry-wise matrix p-norm and ℓ is the total loss. ℓ_{ASR} is the squared-error between the output of a pre-trained ASR encoder when using clean features, \mathbf{f}_X , and the estimate of the clean features, $\mathbf{f}_{\widehat{X}}$, computed from $\widehat{\mathbf{X}}$. λ_{ASR} weights ℓ_{ASR} . Note that ℓ_{ASR} only affects the parameters of the enhancement model; ASR model parameters are kept frozen during training.

2.1.2. Acoustic Echo Cancellation

AEC follows an architecture that is very similar to the Cleanformer, as shown in Fig. 2b. The main difference is the inputs to the AEC frontend: log Mel features computed on the loopback signal concatenated with those computed on the microphone input. For AEC, we assume that the input is single channel; the goal is to remove the device echo from the microphone input. Other than the inputs, the rest of the architecture closely follows the enhancement frontend; a Conformer encoder followed by an FCLayer to compute M. As before, $\widehat{\mathbf{M}}$ is post-processed to $\overline{\mathbf{M}}$, and used for estimating clean features for ASR. The model is also trained on the combination of $\ell_{\mathbf{M}}$ and ℓ_{ASR} .

2.2. Mask-scalar prediction

Both acoustic frontends use hyperparameters α and β to limit speech distortion. The optimal values for α and β depend on the input. In clean conditions, $\alpha=0$ ensures that the input features are unchanged. But in noisy conditions, a larger value of α works better, as we show in Sec. 4.1. Since the best value is input dependent, we propose learning it directly from the data.

In this work, we predict $\alpha(t)$, a per-frame exponential scalar for the mask. The architecture we use is shown in Fig. 3. Compared to the models described in the previous sections, the model has an additional Mask Scalar Net, which is an FCLayer that maps the encoder output, \mathbf{e}_t , to $\alpha(t)$:

$$\alpha(t) = \operatorname{sigmoid}(\operatorname{FCLayer}(\operatorname{StopGrads}(\mathbf{e}_t); \theta_{\alpha})).$$
 (4)

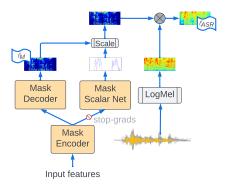


Figure 3: Mask scalar prediction.

 θ_{α} are the parameters of this FCLayer. Importantly, θ_{α} is learned exclusively using ℓ_{ASR} . The mask decoder and mask scalar nets share the encoder, which limits the number of additional parameters. However, gradients are not back propagated from the mask scalar net to the encoder (StopGrads). This ensures that mask scalar prediction is driven entirely based on ASR performance without being influenced by $\ell_{\mathbf{M}}$. Given $\alpha(t)$, the post-processed mask is computed as:

$$\overline{\mathbf{M}}(t,\cdot) = \max(\widehat{\mathbf{M}}(t,\cdot)^{\alpha(t)}, \beta). \tag{5}$$

 $\overline{\mathbf{M}}$ is finally used to obtain the inputs to the ASR model.

3. Experimental Settings

3.1. Data

The target speech data for training both frontends is created using Librispeech [30], Librivox [31] and internal, vendor-collected datasets. The training set has 50.8k hours of speech.

3.1.1. Speech Enhancement

The training data is generated using a room simulator. Noise is sampled from internally collected noise-only data that simulate conditions like cafe, kitchen, and cars, as well as from Getty¹ and YouTube Audio Library². The signal to noise ratio (SNR) is sampled from the range -10 dB to 30 dB. Simulated rooms have reverberation times (T60) from 0 msec to 900 msec. The room impulse responses (RIR) correspond to a 3 microphone array, spaced along an equilateral triangle, 66 millimeters apart. We also use multi-talker training sets, using Librivox and vendor-collected utterances as background speech, and at SNRs (defined as target to background-speech) between 1 dB and 10 dB. Multiple copies of the training data are generated, using randomly selected room configurations and SNRs for each target utterance. SpecAug [8] is also used during training.

The test sets consist of: 1) Fully simulated sets based on the test-clean subset of Librispeech; we create a reverberation-only set ("Reverb"), and sets at -5 dB and 0 dB using environmental noise or speech background. 2) Rerecorded vI, constructed by recording speech and noise separately in a room using a 3-mic array, and mixed at 0 dB and 6 dB SNR; we use pink noise, which is different from the noise types seen in training. These sets have ~ 13 k utterances each. 3) Rerecorded v2 records target speech and noise, played simultaneously, using a 5-mic array unseen during training; this set is recorded in 3 conditions: With no additional noise other than the ambient noise in the room

(Reverb), with a movie playing in the background ("Movie") and with babble noise played out from a different location in the room ("Babble"). These sets have 500 utterances each.

Each utterance in the training and test sets have \sim 6 seconds of noise context that is used by McCleaner. In multi-talker conditions, the noise context helps identify background speech, since only the interfering speaker is active in the noise context.

3.1.2. Acoustic Echo Cancellation

For AEC, we construct training data using synthetic and real echoes [23]. To simulate synthetic echoes, both target and loopback signals are convolved with synthetic RIRs, with the hypothetical loudspeaker close to the microphone. The synthetic datasets are created using Librispeech and Librivox as the source data; the same sets are also used as loopback along with waveforms from Getty and YouTube Audio Library. The signal-to-echo ratio (SER) is set to be between -20 dB and 5 dB. To create datasets with real echoes, we rerecorded, at various loudness levels, Librispeech and an internal dataset that was collected for training text-to-speech (TTS) models using Google Home devices. The recorded echoes are added to reverberant target speech at SERs ranging from -20 dB to 5 dB.

Real echoes, recorded using a held-out subset of the utterances from the internal TTS-based dataset, are mixed with the test-clean subset of Librispeech to create test sets. Test sets are constructed at SNRs from -10 dB to 5 dB in 5 dB increments.

3.2. Architecture

We use 128-dimensional log Mel features, stacked across 4 contiguous frames as input to the frontends and ASR. The window size is 32 msec, at 10 msec hops. The features are also subsampled by a factor of 3. All models use the ideal ratio mask corresponding to the reverberant speech as the target.

The encoder for all our frontends consists of a 4-layer conformer model. Each layer has 256 units; a convolutional block with kernel size 15, and 1024-dimensional feed-forward nets. We use causal masked attention with a left context of 31 frames. The FCLayer in the mask decoder projects the 256-dimensional encoder output to a 512-dimensional mask. The mask scalar net is a 256×1 dimensional FCLayer that projects the encoder output to a 1-dimensional scalar. In total, the frontends have $\sim\!6.5$ million parameters.

 λ_{ASR} in Eq. 3 is increased from 0.0 to 100.0 linearly from 20k steps to 200k steps, and kept fixed after that. When using the mask scalar net, we use a fixed value of $\alpha=0.5$ till 200k steps, after which the model optimizes it using the ASR loss. Training α from scratch resulted in divergence. θ_{α} is initialized by sampling from a Gaussian distribution with a standard deviation of 0.01. Together with the sigmoid activation used by the FCLayer, this ensures that the initial values for the predicted $\alpha(t)$ during training after 200k steps are close to 0.5 to avoid any sudden jumps that can cause training instability.

The ASR model used for evaluation is an LSTM-based multidomain recurrent neural net transducer [32], trained on ~400k hours of English speech, covering domains like near-field and far-field VoiceSearch, YouTube, and Telephony. The utterances for VoiceSearch and Telephony are anonymized and hand-transcribed. We augment training data using SpecAug [8] or simulated noise [9], making the model robust to moderate noise levels. We use the Lingvo toolkit for training [33].

 $^{^{1} \}verb|https://www.gettyimages.com/about-music|$

²https://youtube.com/audiolibrary

Table 1: Enhancement performance, in terms of ASR WER, using various models. Cleanformer uses $\alpha = 0.5, \beta = 0.01$. MSP stands for mask scalar prediction. E1, E2 use the model in Fig. 3; E1 without StopGrads, E2 with StopGrads. E1 and E2 use $\beta = 0.01$.

	Librispeech				Rerecorded v1		Rerecorded v2			
Model		Environment		Speech		Pink		Reverb	Movie	Babble
	Reverb	-5 dB	0 dB	-5 dB	0 dB	0 dB	6 dB	Keverb	wiovie	Dannie
Baseline	7.2	36.5	22.5	65.3	44.8	60.7	28.4	6.4	80.2	61.3
Cleanformer TasNet	7.3	17.5	12.6	26.3	20.0	29.3	15.3	6.3	65.5	50.5
Cleanformer	7.3	13.7	10.6	19.8	15.9	19.1	10.1	5.8	39.1	31.9
+ MSP E1 [this work]	7.3	13.4	10.5	19.2	15.5	17.9	9.6	5.8	33.0	27.6
+ MSP E2 [this work]	7.3	13.3	10.4	19.1	15.5	17.5	9.4	5.7	33.0	28.5

4. Results

4.1. Effect of Mask Scalar on ASR

Table 2: Performance of Cleanformer as a function of mask scalar α ; mask floor $\beta = 0.01$.

Condition	Base- line	1e-6	0.25	$\begin{array}{c} \alpha \\ 0.5 \end{array}$	0.75	1.0
Reverb	7.2	7.3	7.3	7.3	7.4	7.3
0 dB Env.	22.5	22.5	12.6	10.6	10.6	11.1
0 dB Sp.	44.8	44.8	20.5	15.9	15.8	16.7

In Tab. 2, we show how a fixed α affects ASR performance, when using Cleanformer for enhancement. Results are shown in Reverb and 0 dB SNR conditions. $\alpha=0.5$ and 0.75 work well across conditions. When $\alpha<0.5$, the amount of residual noise after enhancement deteriorates performance. And when α is closer to 1.0, the resulting speech distortion worsens WER by 4.7% in environmental noise and 5.0% in multi-talker conditions. Clearly, the choice of α significantly affects WERs. When not predicting $\alpha(t)$ using a mask scalar net, we conservatively set $\alpha=0.5$ for the remaining experiments.

4.2. Speech enhancement

Enhancement results are presented in Tab. 1. For comparison, the table also shows results using Cleanformer TasNet, which is a waveform version of Cleanformer [13, 34]: It uses a learnable TasEncoder layer to convert waveforms to features, followed by a conformer encoder for mask estimation, and a TasDecoder and overlap-add operation to convert back to time-domain. During training, scale invariant SNR loss and ASR loss are used. This model is also causal, and has $\sim\!1.6\mathrm{M}$ parameters operating on 5 msec windows with 2.5 msec overlap, which is typical. Even though the model has fewer parameters, it requires more computation than the remaining log Mel models.

Cleanformer TasNet provides large gains over the noisy baseline; e.g., for the Rerecorded v1 set at 0 dB, it improves WER by 51.7%. On the same set, Cleanformer outperforms Cleanformer TasNet by 34.8%, and the baseline by 68.5%. Enhancing directly in the feature space and controlling for speech distortions, likely help Cleanformer outperform Cleanformer TasNet. We compare two strategies for mask scalar prediction (MSP). Cleanformer + MSP E1 uses the model architecture in Fig. 3, but does not include StopGrads operation. This already provides significant improvements, e.g. a 6.3% relative reduction in WER at 0 dB on Rerecorded v1 set. The relative gains are generally larger on the mismatched, rerecorded test sets compared to simulated Librispeech test sets. Cleanformer + MSP E2, which uses StopGrads provides small gains over Cleanformer + MSP E1. Compared to Cleanformer, it improves WER by 8.4% at 0 dB in Rerecoded v1 set. The gains are more significant in harder conditions, e.g. Rerecorded v2 Movie noise, where Cleanformer + MSP E2 improves over Cleanformer by 15.6%. Overall, all enhancement frontends significantly improve performance over the baseline; using a predicted mask scalar consistently outperforms other enhancement approaches across all conditions, including the typical approach of using a hand-tuned mask scalar.

4.3. AEC

Table 3: AEC performance, in terms of ASR WER.

Model	SER						
Model	-10 dB	-5 dB	0 dB	5 dB			
Baseline	80.5	72.7	58.0	36.1			
LMel-NAEC [23]	29.8	21.8	16.9	14.7			
Mask-NAEC	26.2	17.4	12.5	10.0			
+ MSP E2 [this work]	24.3	16.4	12.0	9.8			

Tab. 3 shows AEC performance in terms of WER. We present results using the algorithm presented in [23] (LMel-NAEC), which estimates log Mel features directly, using regression and ASR losses. The overall architecture is identical to the proposed models, except that LMel-NAEC directly predicts the log Mel spectra as opposed to a mask. Predicting the log Mel spectra makes it challenging to limit speech distortion using post-processing techniques like the ones we use with estimated masks. Compared to LMel-NAEC, predicting the mask and using a fixed α and β during inference (Mask-NAEC) already reduces the WER significantly. Note that Mask-NAEC is similar to the approach in [35], but it additionally introduces ASR loss into AEC training [23]. Compared to LMel-NAEC, Mask-NAEC improves WER by 12.0% at -10 dB and 32.2% at 5 dB. Clearly, using the post-processed mask even with a fixed mask scalar and floor significantly reduces distortion, especially at higher SNRs. Finally, using the proposed mask scalar net further improves WER by 7.2% at -10 dB and 2.1% at 5 dB. As with the enhancement experiments, the advantage of using a predicted mask scalar increases as the SNR goes down. Compared to LMel-NAEC, the proposed model improves WER by 18.3% at -10 dB and 33.6% at 5 dB.

5. Conclusions

Limiting speech distortions introduced by acoustic frontends for ASR is important to get the best performance out of them. In this work, we propose using a neural net to predict exponential mask scalars in an end-to-end fashion. Our results show that predicted mask scalars reduce WER on multiple robustness tasks, with almost no increase in model size. The current work only focused on predicting a per-frame exponential scalar. Future work will explore predicting the mask floor, and more fine-grained time-frequency level prediction.

6. References

- B. Li, A. Gulati, J. Yu, T. N. Sainath, C.-C. Chiu, A. Narayanan, S.-Y. Chang *et al.*, "A better and faster end-to-end model for streaming ASR," in *Proc. IEEE ICASSP*, 2021.
- [2] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 29, pp. 3451–3460, 2021.
- [3] J. Li, "Recent advances in end-to-end automatic speech recognition," arXiv preprint arXiv:2111.01690, 2021.
- [4] J. Barker, S. Watanabe, E. Vincent, and J. Trmal, "The fifth CHiME speech separation and recognition challenge: Dataset, task and baselines," in *Proc. Interspeech*, 2018.
- [5] A. Narayanan, A. Misra, K. C. Sim, G. Pundak, A. Tripathi, M. Elfeky, P. Haghani, T. Strohman, and M. Bacchiani, "Toward domain-invariant speech recognition via large scale training," in *Proc. IEEE SLT Workshop*, 2018.
- [6] V. Pratap, A. Sriram, P. Tomasello, A. Hannun, V. Liptchinsky, G. Synnaeve, and R. Collobert, "Massively multilingual ASR: 50 languages, 1 model, 1 billion parameters," *Proc. Interspeech*, 2020.
- [7] Y. Zhang, D. S. Park, W. Han, J. Qin, A. Gulati, J. Shor, A. Jansen, Y. Xu, Y. Huang, S. Wang *et al.*, "BigSSL: Exploring the frontier of large-scale semi-supervised learning for automatic speech recognition," *arXiv preprint arXiv:2109.13226*, 2021.
- [8] 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," *Proc. Interspeech*, 2019.
- [9] 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.
- [10] Z. Zhang, J. Geiger, J. Pohjalainen, A. E.-D. Mousa, W. Jin, and B. Schuller, "Deep learning for environmentally robust speech recognition: An overview of recent developments," *ACM Trans*actions on Intelligent Systems and Technology (TIST), vol. 9, no. 5, pp. 1–28, 2018.
- [11] A. Narayanan and D. L. Wang, "Ideal ratio mask estimation using deep neural networks for robust speech recognition," in *Proc. IEEE ICASSP*, 2013.
- [12] Z.-Q. Wang, P. Wang, and D. Wang, "Complex spectral mapping for single-and multi-channel speech enhancement and robust ASR," *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 28, pp. 1778–1787, 2020.
- [13] Y. Luo and N. Mesgarani, "TasNet: time-domain audio separation network for real-time, single-channel speech separation," in *Proc. IEEE ICASSP*, 2018.
- [14] K. Kinoshita, T. Ochiai, M. Delcroix, and T. Nakatani, "Improving noise robust automatic speech recognition with single-channel time-domain enhancement network," in *Proc. IEEE ICASSP*, 2020.
- [15] J. Heymann, L. Drude, and R. Haeb-Umbach, "Neural network based spectral mask estimation for acoustic beamforming," in *Proc. IEEE ICASSP*, 2016.
- [16] J. Heymann, M. Bacchiani, and T. N. Sainath, "Performance of mask based statistical beamforming in a smart home scenario," in *Proc. IEEE ICASSP*, 2018.
- [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. IEEE ICASSP*, 2019.
- [18] A. Pandey, C. Liu, Y. Wang, and Y. Saraf, "Dual application of speech enhancement for automatic speech recognition," in *Proc. IEEE SLT Workshop*, 2021.

- [19] H. Sato, T. Ochiai, M. Delcroix, K. Kinoshita, T. Moriya, and N. Kamo, "Should we always separate?: Switching between enhanced and observed signals for overlapping speech recognition," arXiv preprint arXiv:2106.00949, 2021.
- [20] 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," arXiv preprint arXiv:2201.03881, 2022.
- [21] 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 preprint arXiv:2201.06685, 2022.
- [22] Y. Koizumi, S. Karita, A. Narayanan, S. Panchapagesan, and M. Bacchiani, "SNRi target training for joint speech enhancement and recognition," arXiv preprint arXiv:2111.00764, 2021.
- [23] N. Howard, A. Park, T. Z. Shabestary, A. Gruenstein, and R. Prabhavalkar, "A neural acoustic echo canceller optimized using an automatic speech recognizer and large scale synthetic data," in *Proc. IEEE ICASSP*, 2021.
- [24] A. Narayanan and D. L. Wang, "Joint noise adaptive training for robust automatic speech recognition," in *Proc. IEEE ICASSP*, 2014.
- [25] Q. Wang, K. A. Lee, T. Koshinaka, K. Okabe, and H. Ya-mamoto, "Task-aware warping factors in mask-based speech enhancement," in *Proc. EUSIPCO*, 2021.
- [26] G. Pundak and T. N. Sainath, "Lower frame rate neural network acoustic models," in *Proc. Interspeech*, 2016.
- [27] J. Caroselli, A. Naranayan, and T. O'Malley, "Cleanformer: A microphone array configuration-invariant, streaming, multichannel neural enhancement frontend for ASR," under review, Interspeech 2022.
- [28] Y. A. Huang, T. Z. Shabestary, and A. Gruenstein, "Hotword cleaner: dual-microphone adaptive noise cancellation with deferred filter coefficients for robust keyword spotting," in *Proc. IEEE ICASSP*, 2019.
- [29] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu et al., "Conformer: Convolutionaugmented transformer for speech recognition," arXiv preprint arXiv:2005.08100, 2020.
- [30] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. IEEE ICASSP*, 2015.
- [31] J. Kearns, "Librivox: Free public domain audiobooks," Reference Reviews, 2014.
- [32] T. N. Sainath, Y. He, B. Li, A. Narayanan, R. Pang, A. Bruguier, S.-y. Chang, W. Li, R. Alvarez, Z. Chen et al., "A streaming on-device end-to-end model surpassing server-side conventional model quality and latency," in *Proc. IEEE ICASSP*, 2020.
- [33] J. Shen, P. Nguyen, Y. Wu, Z. Chen, M. X. Chen, Y. Jia, A. Kannan, T. Sainath, Y. Cao, C.-C. Chiu et al., "Lingvo: a modular and scalable framework for sequence-to-sequence modeling," arXiv preprint arXiv:1902.08295, 2019.
- [34] S. Panchapagesan, A. Narayanan, T. Z. Shabestary, S. Shao, N. Howard, A. Park, J. Walker, and A. Gruenstein, "A Conformerbased waveform-domain neural acoustic echo canceller optimized for ASR accuracy," under review, Interspeech 2022.
- [35] H. Zhang and D. Wang, "Deep learning for acoustic echo cancellation in noisy and double-talk scenarios," in *Proc. Interspeech*, 2018.