Regularizing Learnable Feature Extraction for Automatic Speech Recognition

Peter Vieting¹, Maximilian Kannen¹, Benedikt Hilmes^{1,2}, Ralf Schlüter^{1,2}, Hermann Ney^{1,2}

¹Machine Learning and Human Language Technology Group, RWTH Aachen University, Germany ²AppTek GmbH, Germany

{vieting, hilmes, schlueter, ney}@ml.rwth-aachen.de

Abstract

Neural front-ends are an appealing alternative to traditional, fixed feature extraction pipelines for automatic speech recognition (ASR) systems since they can be directly trained to fit the acoustic model. However, their performance often falls short compared to classical methods, which we show is largely due to their increased susceptibility to overfitting. This work therefore investigates regularization methods for training ASR models with learnable feature extraction front-ends. First, we examine audio perturbation methods and show that larger relative improvements can be obtained for learnable features. Additionally, we identify two limitations in the standard use of SpecAugment for these front-ends and propose masking in the short time Fourier transform (STFT)-domain as a simple but effective modification to address these challenges. Finally, integrating both regularization approaches effectively closes the performance gap between traditional and learnable features. Index Terms: speech recognition, feature extraction, raw wave-

form modeling

1. Introduction

Feature extraction is a key component of any automatic speech recognition (ASR) system. Conventional methods rely on handcrafted feature engineering which is often inspired by human auditory perception [1, 2]. While this has shown to be effective, these approaches risk discarding valuable speech information, potentially limiting the ASR performance. Learnable feature extraction methods address this limitation by optimizing features directly for the acoustic model and thus ensure a more tailored representation [3, 4, 5]. Beyond the final performance, the use of learnable features aligns with the broader goal of designing end-to-end ASR systems, where the entire processing pipeline – from raw waveform to predicted labels – is optimized jointly within a single monolithic neural framework.

In practice, learnable feature extraction front-ends often struggle to compete with classical methods especially when the amount of training data is limited [5]. While prior works have not explicitly considered overfitting as a problem, we hypothesize that learnable front-ends are more susceptible to it. Figure 1 clearly shows this effect. This makes overfitting an important concern that causes a performance gap. There are well-known methods to mitigate it during neural network training in general [6, 7]. However, their application to learnable feature extractors in ASR has received little attention. In this work, we investigate perturbations of the input audio during training as an effective regularization strategy to improve generalization. Furthermore, we demonstrate that the default use of SpecAugment [8] is suboptimal for learnable feature extractors and propose a refined method to address these limitations. Naturally, overfitting is particularly challenging when training data is limited. In this study, we focus on scenarios with substantially less data than typically used in previous studies that achieve competitive performance with learnable front-ends [3]. The Switchboard corpus with 311 h of data is a suitable choice for this. Beyond the size, this corpus is particularly challenging for learnable features because of its telephony speech nature with a limited bandwidth of 8 kHz. Our approach is thus viable for demanding low resource scenarios and situations where large pretrained models are prohibitive due to efficiency constraints.

The key contributions of this work are as follows:

- We show that overfitting is a problem when training ASR models with learnable features.
- A systematic study on the effect of audio perturbation for learnable front-ends is presented which, to the best of our knowledge, has not been explored before.
- A novel variant of SpecAugment is proposed to specifically address the challenges of learnable feature extraction.
- We demonstrate that combining both regularization methods allows the learnable front-end to perform on par with traditional features on Switchboard with 311 h of training data.



Figure 1: Train and dev connectionist temporal classification (CTC) scores for the baseline training with log Mel and supervised convolutional features (SCF), respectively, demonstrating the overfitting issue for the latter.

2. Related Work

While there exists a line of research on learnable front-ends for ASR and related tasks [3, 4, 5, 9], the mitigation of overfitting in this context has not been addressed so far. Parametric methods [10] may inherently suffer less from overfitting due to their restricted modeling capacity. However, this limitation may potentially affect the overall performance. As an alternative approach, we therefore aim to use regularization methods during the neural network training to mitigate overfitting in this work.

The literature provides a rich set of techniques for regularizing neural network training such as weight decay [6], dropout [7] or label smoothing [11]. Perturbation of the input is another important regularization technique. Common methods for audio perturbation are speed, tempo and pitch perturbation [12, 13, 14]. Further ideas to augment the audio input include signal companding [15] and mixup [16].

Additionally, SpecAugment [8] has been shown to be particularly effective for ASR. It works by masking blocks of time frames or feature channels during training as well as performing time warping. Following its success, numerous variations have been proposed, e.g. dropping [17] or swapping [18] blocks instead of masking, masking based on phonemes or other semantic regions [19, 20] and more [21, 22, 23]. [24] proposes a variation of SpecAugment's time warping for a wav2vec 2.0 model with learnable front-end. However, this modification does not appear to be specifically designed for learnable feature extraction and the experimental results do not isolate its individual contribution to the performance improvement.

The SpecAugment variation proposed in this work is different from previous work in that it works in the short time Fourier transform (STFT)-domain and caters specifically for the needs of learnable front-ends.

3. Methods

3.1. Feature Extraction

Typically, acoustic models for ASR operate on features that represent the input speech. One of the most common representations are log Mel filterbank features. They are obtained by first computing the STFT, in our case with window size 25 ms and shift 10 ms, and then taking the square of the magnitude. Using Mel warping, an 80-dim vector is obtained per frame and finally a logarithmic compression and normalization are applied.

As the learnable feature extraction front-end, we use supervised convolutional features similar to [5, 25, 26]. The main components are two convolutional layers that operate as timefrequency-decomposition and temporal integration. These operations resemble the functional principle of the Gammatone filterbank and the Hanning window used for Gammatone features [2]. However, all filters are randomly initialized and trained with the rest of the neural network with the same supervised objective function. The first convolutional layer takes the waveform as input. It has 150 filters with a size of 16 ms and a stride of 0.625 ms. The second layer has 5 filters with a size of 40 frames and a stride of 16, resulting in feature frames with a 10 ms shift. Since each of the 5 filters is applied to all 150 output channels of the first layer, the final feature dimension is 750. In contrast to Gammatones features, a multi-resolutional temporal integration is achieved because 5 different filters can be learned. Additionally, we use the absolute value as the elementwise activation function after the first layer and the 2.5th root of the absolute value after the second layer. Finally, layer normalization is applied. In this case of a learnable front-end, the border between feature extraction and acoustic model blurs and the neural network operates directly on the raw waveform.

Both feature extraction blocks are followed by SpecAugment (see Section 3.3) and VGG convolutions, resulting in a subsampling by a factor of 4 and a frame shift of 40 ms.

3.2. Audio Perturbation

For our work, we investigate multiple ways of modifying the audio signal in order to regularize the training of our model. Each of the perturbations is applied at the sequence level with a specified probability p, where p < 1 ensures to also include original audio in training. While most earlier works generate copies of the corpus with fixed perturbation factors, we sample random factors for the modification in a given range on-the-fly to increase the variability of perturbed examples. This allows the same sequence to pass through different modifications across multiple training iterations.

Speed, Tempo and Pitch Perturbation: The most common way of modifying an audio signal is changing its speed, tempo or pitch. For speed perturbation, we resample the audio by a factor *a*. This results in both the duration and the pitch being modified. Tempo perturbation only changes the tempo and thus the duration of the audio while keeping the pitch level constant. This is achieved using the waveform similarity overlapadd (WSOLA) algorithm [27]. Lastly, we can also only modify the pitch of the audio while keeping the duration constant. Again we use WSOLA, by first changing the tempo of the audio signal but then resample the signal back to the original duration.

Nonlinear Amplitude Perturbation: The next perturbation involves a nonlinear distortion of the audio signal's amplitudes. It is defined as:

$$\tilde{x}(t) = \operatorname{sign}(x(t)) \cdot |x(t)|^{\beta} \tag{1}$$

where x(t) is the original audio normalized to be between -1 and 1, β is the factor that controls the strength of the distortion and $\tilde{x}(t)$ is the perturbed audio. The perturbation either results in a more peaky signal or de-emphasizes outliers.

 μ -Law Perturbation: The next perturbation technique we investigate is based on μ -law companding. This companding algorithm reduces the dynamic range of a given signal and is typically used in telecommunication systems. However, previous work has also applied it as a method for speech data augmentation [15]. We use the continuous encoding equation which offers the possibility to control the strength of the augmentation effect by a single parameter μ :

$$\tilde{x}(t) = \operatorname{sgn}(x(t)) \cdot \frac{\ln(1 + \mu \cdot |x(t)|)}{\ln(1 + \mu)}$$
(2)

While in telecommunication usually μ is set to 255, we use significantly smaller values in order to avoid perturbed examples that deviate too far from the distribution of the real training data. Note that since $\mu > 0$, this is the only perturbation where in general $E[\tilde{x}(t)] \neq x(t)$.

Preemphasis: Preemphasis is a filter technique often used for ASR to emphasize higher frequencies in the audio signal. The filter output y(t) is defined as:

$$y(t) = x(t) - \alpha x(t-1) \tag{3}$$

where α is usually chosen to be close to 1. To avoid any mismatch, preemphasis is applied in the same way during training and inference. In addition to the standard usage of preemphasis, we experiment with a second level of preemphasis as a perturbation technique. For this, we apply Eq. (3) a second time with a random value $\tilde{\alpha}$ sampled from a pre-defined range close to 0. For inference, we only apply the regular preemphasis with α .

3.3. SpecAugment

SpecAugment [8] is a regularization technique that involves masking random regions of the extracted features along both the time and feature dimensions during training. Similar to dropout [7], this approach prevents overfitting by ensuring the



(b) Proposed masking in STFT-domain before feature extraction.

Figure 2: Baseline vs. proposed masking strategies.

model does not overly depend on any particular part of the features. However, unlike dropout, SpecAugment introduces structure by masking contiguous regions of data. Specifically, when applying time masking, blocks of consecutive time frames are masked. Since neighboring frames typically have a high correlation, masking a single frame results in minimal information loss. In contrast, masking entire regions pushes the model to learn dependencies over longer contexts.

For log Mel features, this logic trivially transfers to masking in the feature dimension since neighboring feature channels correspond to adjacent frequency regions. However, this is fundamentally different in the case of learnable features. While the model also typically learns a set of bandpass filters [5, 26], it cannot be assumed that adjacent channels represent similar frequencies since the order of filters is arbitrary. We hypothesize that this limits the efficacy of SpecAugment for learnable frontends. Furthermore, SpecAugment is applied after the feature extraction as depicted in Figure 2a. This causes the gradient for the masked positions to be 0, resulting in fewer updates for each front-end filter. As another consequence, learnable layers before and after the masking allow the model to shift information and bypass the masks resulting in ineffective regularization.

In this work, we investigate two different ways to solve these issues. The first way is to sort the filters based on the peak frequency in their frequency response and to mask channels with neighboring peak frequencies. However, this is complicated by the fact that the learned filters are not perfectly sharp bandpass filters but sometimes include wide or even multiple passbands [5, 26]. Thus, a proper masking of frequency areas cannot be guaranteed. Also, the masking is still applied after the feature extractor such that this approach only addresses the first issue described above.

For the second proposed approach, we apply SpecAugment before the feature extraction process. We apply masking in both time and frequency dimension in the STFT-domain, transform the signal back using the inverse STFT and use the resulting signal as input to our model. This way, we can apply SpecAugment before feature extraction and guarantee that both time and frequency areas are masked properly, tackling both issues described above. The information flow for the baseline and proposed masking strategies is visualized in Figure 2.

4. Experimental Setup

4.1. Data

The Switchboard-1 Release 2 dataset [28] serves as training data for our experiments. The corpus contains 311 h of conversational English telephony speech sampled at 8 kHz. We use Hub5'00 as the development set and Hub5'01 for evaluation.

4.2. ASR Model

For our experiments, we use a connectionist temporal classification (CTC)-based ASR model [29]. We follow the setup from a previous work [30] and use a 4-gram language model (LM) during recognition. Feature extraction is done as explained in Section 3.1 including VGG-style subsampling, followed by dropout and a final linear layer producing representations with a shift of 40 ms. Following results from preliminary experiments, standard preemphasis with $\alpha = 0.97$ is used for the supervised convolutional features, but not for log Mel. Also, for the log Mel baseline, it was beneficial to normalize the input waveform over time and the resulting features over the batch. Our encoder consists of 12 conformer [31] layers with a hidden dimension of 512. Compared to the original paper, we swap the position of convolution and multi-head self-attention modules. A final linear layer maps the encoder output to the target vocabulary, resulting in 74M parameters for the log Mel model. When using learnable features, there are around 21k parameters for the feature extractor and the number of parameters for the linear laver after VGG increases from 1.3M to 12.3M because of the almost 10 times larger feature dimension¹ leading to 85M parameters in total. We train our model using the NAdam optimizer applying L2 weight decay [6] with a factor of 0.01 for the VGG frontend and 0.0001 for the encoder. Further, dropout [7] is applied with a probability of 0.1 for regularization. We split the data into 6 sub-epochs and train a total of 450 subepochs, evaluating our model on sub-epochs 400 to 450 in steps of 10 and select the best checkpoint based on word error rate (WER) on Hub5'00 for test evaluation. We apply SpecAugment from the first sub-epoch. To increase training stability, perturbations are only enabled after the first 4 full epochs by initializing the model with the corresponding checkpoint of the baseline. To keep the total number of training epochs constant, we stop training with perturbations 4 full epochs earlier. We use a onecycle learning rate (OCLR) schedule starting from $1.325 \cdot 10^{-5}$ with a peak of $4.0 \cdot 10^{-4}$ after 180 sub-epochs going down to $1.0 \cdot 10^{-5}$ over another 180 sub-epochs and staying at that level for the final sub-epochs of training. The experiments can be run on a single consumer graphics processing unit (GPU) with 24 GB VRAM (e.g. Nvidia RTX 3090) so that the barrier for reproduction is low. Also, the software is publicly accessible [32, 33] and the recipes are available online².

5. Experimental Results

5.1. Audio Perturbation

We study the effect of the audio perturbation techniques introduced in Section 3.2 in Table 1. Due to hardware availability constraints, we run the ablation study for Table 1 on less powerful GPUs with only 11 GB of memory. To facilitate this, the batch size is reduced from 100 s to 50 s and the gradients are accumulated over two batches. Note that this is not strictly equivalent when the model contains batch normalization. Additionally, preemphasis with $\alpha = 1.0$ yielded better results for some trainings and we always report the lower WER in Table 1. The full ablation study contains a large number of experiments, however, in the interest of clarity, we present only those experiments here that deliver the best results for log Mel and super-

²https://github.com/rwth-i6/

returnn-experiments/tree/master/

¹The VGG output dimension is 32 times the feature dimension resulting in $750 \cdot 32 \cdot 512 \approx 12.3$ M parameters for the linear layer.

²⁰²⁵⁻regularizing-learnable-features

vised convolutional features for each perturbation type.

For each experiment, we define the probability p of perturbing the audio input as well as the range to sample the hyperparameter controlling the perturbation strength from. We check $p \in \{0.3, 0.7, 1.0\}$ for most perturbations to cover a wide range of probabilities. For speed perturbation, we sample from the range that is typically used, i.e., $a \sim U(0.9, 1.1)$ [12], but also from the less common $a \sim U(0.88, 1.12)$ [34]. Most other works use a range of [0.9, 1.1] for tempo perturbations as well [12]. However, our experiments showed that larger ranges are more beneficial and therefore our best results with tempo perturbation are obtained with $a \sim U(0.7, 1.3)$. For pitch perturbation, a change of -2 to 2 semitones roughly corresponds to speed changes of 0.9 to 1.1 and is therefore our starting point. Yet, larger ranges degraded the performance.

For non-standard perturbations, selecting appropriate hyperparameters is more challenging due to the lack of prior references. To address this, we determined the minimum and maximum perturbation factors such that the audio still sounded natural based on the authors' subjective perception. Based on initial experimental results, we further tuned the parameters.

Table 1 shows the obtained results. It is clear that tempo perturbation achieves the best WERs across both feature types. While the log Mel model also clearly benefits from the nonlinear amplitude, μ -law and preemphasis perturbation, only the latter yields a small improvement for the learnable features. However, the overall best results exhibit a larger relative improvement for the learnable front-end.

Table 1: Ablation study of different audio perturbation methods. Results are reported for both log Mel and supervised convolutional features (SCF) on Hub5'00. All results use standard SpecAugment.

Pertu	WER [%]				
Туре	p	Min	Max	log Mel	SCF
None	-	-	-	13.0	13.2
Speed	0.7	0.88	1.12	13.1	13.3
Tempo	1.0	0.7	1.3	12.6	12.6
Pitch	0.7	-2	2	13.2	13.4
Nonlin. Amp.	0.7	0.8	1.2	12.6	13.2
μ -law	0.3	1	5	12.6	13.2
Preemphasis	0.7	-0.05	0.05	12.7	13.0

5.2. SpecAugment Variations

The next set of experiments deals with the analysis of the previously proposed variations of SpecAugment. A maximum time mask size of 15 frames was tuned for the log Mel baseline model and previously copied for experiments with the learnable front-end. The maximum feature mask sizes were tuned individually for each row in Table 2 and are 8 or 15 in the baseline cases. Now, we first evaluate creating larger time masks while reducing the number of masks at the same time to keep the ratio of masked time frames constant. Because of the temporal correlation of the input, this has a stronger regularization effect and forces the model to learn longer context dependencies. As visible in Table 2, this has little effect for the log Mel features but yields a 0.4% absolute improvement for the learnable front-end. Applying SpecAugment in the STFT-domain achieves the same absolute gain of 0.4%. In combination with larger time masks, the best result is obtained. Interestingly, even the log Mel features benefit from masking in the STFT-domain in combination with larger time masks. Finally, feature masking based on the sorted order of learned filters does not improve the result.

Table 2: Comparison of different SpecAugment variations. The given mask sizes for time (T) and frequency/feature (F) are upper limits for random sampling of the actual sizes. Results are reported on Hub5'00.

Features	SpecAugment	Mask Size T F		WER [%]
log Mel	Baseline	15 30	8	12.8 12.8
	STET domain	15	4	12.8
	STI I-domain	30	8	12.5
SCF	Baseline	15	15	13.2
	Dasenne	30	8	12.8
	Sorted	15	0	13.3
	STET domain	15	4	12.8
	STI I-domain	30	8	12.6

Table 3: Final results on Hub5'00 and Hub5'01 for the combination of the previously best performing audio perturbations (tempo perturbation) and SpecAugment settings.

Features	Derturbation	Spec Augment	WER [%]		
	renturbation	SpecAuginent	Hub5'00	Hub5'01	
log Mel	no	default	12.8	11.5	
	yes	ueraun	12.4	11.2	
		STFT-domain	12.4	11.3	
SCF	no	default	13.2	12.4	
	yes	STFT-domain	12.5	11.3	

5.3. Final Combination

To obtain the final best model, we combine the best SpecAugment variation with the best audio perturbations. We select tempo perturbation here as it performed best on both log Mel and learnable features. While we also tried combining different perturbations, this resulted in an overall worse performance, likely because the regularization effect is too strong in this case. Table 3 presents the results. As visible in the first three lines, tempo perturbation improves the log Mel baseline results on both Hub5'00 and Hub5'01. In combination with tempo perturbation. SpecAugment in the STFT-domain is not beneficial for log Mel features. In contrast, the combination yields the best results for the learnable front-end which are now within 0.1% absolute from the best log Mel results. We also observe that the gap between train and dev scores is significantly smaller than in Figure 1, showing that overfitting has been mitigated effectively. Therefore, we demonstrate that the presented regularization methods for training an ASR system with learnable features allow closing the performance gap to traditional features.

6. Conclusion

This work focuses on regularization techniques for training CTC-based ASR models using both traditional and learnable feature extractors on the Switchboard dataset. The examined audio perturbations can mitigate the overfitting during training and the overall improvements are relatively larger for learnable features. Furthermore, we propose masking in the STFT-domain as a simple but effective modification of SpecAugment to allow masking of contiguous frequency regions and ensuring variability in the input which was previously not the case for neural front-ends. The final results with a combination of both techniques effectively close the performance gap between traditional and learnable features.

7. Acknowledgements

This work was partially supported by NeuroSys, which as part of the initiative "Clusters4Future" is funded by the Federal Ministry of Education and Research BMBF (funding ID 03ZU2106DD), and by the project RESCALE within the program *AI Lighthouse Projects for the Environment, Climate, Nature and Resources* funded by the Federal Ministry for the Environment, Nature Conservation, Nuclear Safety and Consumer Protection (BMUV), funding ID: 6KI32006A.

8. References

- S. S. Stevens and J. Volkmann, "The relation of pitch to frequency: A revised scale," *The American Journal of Psychology*, vol. 53, no. 3, pp. 329–353, 1940.
- [2] R. Schlüter, I. Bezrukov, H. Wagner, and H. Ney, "Gammatone features and feature combination for large vocabulary speech recognition," in *Proc. ICASSP*, Honolulu, HI, USA, Apr. 2007, pp. 649–652.
- [3] T. N. Sainath, R. J. Weiss, A. Senior, K. W. Wilson, and O. Vinyals, "Learning the speech front-end with raw waveform CLDNNs," in *Proc. Interspeech*, Dresden, Germany, Sep. 2015.
- [4] N. Zeghidour, N. Usunier, I. Kokkinos, T. Schaiz, G. Synnaeve, and E. Dupoux, "Learning filterbanks from raw speech for phone recognition," in *Proc. ICASSP*, Calgary, Canada, Apr. 2018, pp. 5509–5513.
- [5] Z. Tüske, R. Schlüter, and H. Ney, "Acoustic modeling of speech waveform based on multi-resolution, neural network signal processing," in *ICASSP*, Calgary, Canada, 2018, pp. 4859–4863.
- [6] S. Hanson and L. Pratt, "Comparing biases for minimal network construction with back-propagation," in *Proc. NIPS*, Denver, CO, USA, 1988, pp. 177–185.
- [7] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014.
- [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," Graz, Austria, pp. 2613–2617, Sep. 2019.
- [9] N. Zeghidour, O. Teboul, F. de Chaumont Quitry, and M. Tagliasacchi, "LEAF: A learnable frontend for audio classification," in *Proc. ICLR*, Vienna, Austria, May 2021.
- [10] M. Ravanelli and Y. Bengio, "Speaker recognition from raw waveform with SincNet," in *Proc. SLT*, Athens, Greece, Dec. 2018, pp. 1021–1028.
- [11] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 2016, pp. 2818–2826.
- [12] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur, "Audio augmentation for speech recognition," in *Proc. Interspeech*, Dresden, Germany, Sep. 2015, pp. 3586–3589.
- [13] N. Jaitly and G. E. Hinton, "Vocal tract length perturbation (VTLP) improves speech recognition," in *Proc. ICML workshop* on deep learning for audio, speech and language, vol. 117, Atlanta, GA, USA, Jun. 2013, p. 21.
- [14] N. Kanda, R. Takeda, and Y. Obuchi, "Elastic spectral distortion for low resource speech recognition with deep neural networks," in *Proc. ASRU*, Olomouc, Czech Republic, 2013, pp. 309–314.
- [15] R. K. Das, J. Yang, and H. Li, "Data augmentation with signal companding for detection of logical access attacks," in *Proc. ICASSP*, Toronto, Canada, Jun. 2021, pp. 6349–6353.
- [16] I. Medennikov, Y. Y. Khokhlov, A. Romanenko, D. Popov, N. A. Tomashenko, I. Sorokin, and A. Zatvornitskiy, "An investigation of mixup training strategies for acoustic models in ASR," in *Proc. Interspeech*, Hyderabad, India, Sep. 2018, pp. 2903–2907.

- [17] A. Jain, P. R. Samala, D. Mittal, P. Jyothi, and M. Singh, "SPLICEOUT: A simple and efficient audio augmentation method," in *Proc. Interspeech*, Incheon, Korea, Sep. 2022, pp. 2678–2682.
- [18] X. Song, Z. Wu, Y. Huang, D. Su, and H. Meng, "SpecSwap: A simple data augmentation method for end-to-end speech recognition," in *Proc. Interspeech*, Shanghai, China, Oct. 2020, pp. 581– 585.
- [19] C. Wang, Y. Wu, Y. Du, J. Li, S. Liu, L. Lu, S. Ren, G. Ye, S. Zhao, and M. Zhou, "Semantic mask for transformer based end-to-end speech recognition," in *Proc. Interspeech*, Shanghai, China, Oct. 2020, pp. 971–975.
- [20] X. Zhang, J. Wang, N. Cheng, K. Zhu, and J. Xiao, "Improving speech representation learning via speech-level and phonemelevel masking approach," in *Proc. International Conference on Mobility, Sensing and Networking*, Guangzhou, China, Dec. 2022, pp. 485–489.
- [21] H. Wang, Y. Zou, and W. Wang, "Specaugment++: A hidden space data augmentation method for acoustic scene classification," in *Proc. Interspeech*, Brno, Czech Republic, Sep. 2021, pp. 551–555.
- [22] R. Damania, C. Homan, and E. Prud'hommeaux, "Combining simple but novel data augmentation methods for improving conformer ASR," in *Proc. Interspeech*, Incheon, Korea, Sep. 2022, pp. 4890–4894.
- [23] J. Han, M. Matuszewski, O. Sikorski, H. Sung, and H. Cho, "Randmasking augment: A simple and randomized data augmentation for acoustic scene classification," in *Proc. ICASSP*, Rhodes, Greece, Jun. 2023, pp. 1–5.
- [24] H. Luo, X. Xie, P. Li, and F. Xin, "Data augmentation based unsupervised pre-training for low-resource speech recognition," in *Proc. Chinese Control and Decision Conference*, Xi'an, China, May 2024, pp. 5007–5012.
- [25] P. Vieting, C. Lüscher, W. Michel, R. Schlüter, and H. Ney, "On architectures and training for raw waveform feature extraction in ASR," in *Proc. ASRU*, Cartagena, Colombia, Dec. 2021, pp. 267– 274.
- [26] P. Vieting, R. Schlüter, and H. Ney, "Comparative analysis of the wav2vec 2.0 feature extractor," in *ITG Conference on Speech Communication*, Aachen, Germany, Sep. 2023, pp. 131–135.
- [27] W. Verhelst and M. Roelands, "An overlap-add technique based on waveform similarity (WSOLA) for high quality time-scale modification of speech," in *Proc. ICASSP*, Minneapolis, MN, USA, Apr. 1993, pp. 554–557.
- [28] J. Godfrey, E. Holliman, and J. McDaniel, "Switchboard: telephone speech corpus for research and development," in *Proc. ICASSP*, San Francisco, CA, USA, Mar. 1992, pp. 517–520.
- [29] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proc. ICML*, New York, NY, USA, Jun. 2006, p. 369–376.
- [30] W. Zhou, W. Michel, R. Schlüter, and H. Ney, "Efficient training of neural transducer for speech recognition," in *Proc. Interspeech*, Incheon, Korea, Sep. 2022, pp. 2058–2062.
- [31] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, and R. Pang, "Conformer: Convolution-augmented transformer for speech recognition," in *Proc. Interspeech*, Shanghai, China, Oct. 2020, pp. 5036–5040.
- [32] A. Zeyer, T. Alkhouli, and H. Ney, "RETURNN as a generic flexible neural toolkit with application to translation and speech recognition," in *Proc. ACL*, Melbourne, Australia, Jul. 2018.
- [33] W. Zhou, E. Beck, S. Berger, R. Schlüter, and H. Ney, "RASR2: The RWTH ASR toolkit for generic sequence-to-sequence speech recognition," in *Proc. Interspeech*, Dublin, Ireland, Aug. 2023, pp. 4094–4098.
- [34] G. Chen, X. Na, Y. Wang, Z. Yan, J. Zhang, S. Ma, and Y. Wang, "Data augmentation for children's speech recognition – the "ethiopian" system for the SLT 2021 children speech recognition challenge," Preprint arXiv:2011.04547, 2020.