

ON ARCHITECTURES AND TRAINING FOR RAW WAVEFORM FEATURE EXTRACTION IN ASR

Peter Vieting¹, Christoph Lüscher^{1,2}, Wilfried Michel^{1,2}, Ralf Schlüter^{1,2}, Hermann Ney^{1,2}

¹Human Language Technology and Pattern Recognition Group,
Computer Science Department, RWTH Aachen University, 52074 Aachen, Germany

²AppTek GmbH, 52062 Aachen, Germany

ABSTRACT

With the success of neural network based modeling in automatic speech recognition (ASR), many studies investigated acoustic modeling and learning of feature extractors directly based on the raw waveform. Recently, one line of research has focused on unsupervised pre-training of feature extractors on audio-only data to improve downstream ASR performance. In this work, we investigate the usefulness of one of these front-end frameworks, namely wav2vec, in a setting without additional untranscribed data for hybrid ASR systems. We compare this framework both to the manually defined standard Gammatone feature set, as well as to features extracted as part of the acoustic model of an ASR system trained supervised. We study the benefits of using the pre-trained feature extractor and explore how to additionally exploit an existing acoustic model trained with different features. Finally, we systematically examine combinations of the described features in order to further advance the performance.

Index Terms— Speech recognition, learnable feature extraction, unsupervised training, LibriSpeech

1. INTRODUCTION

Traditionally, automatic speech recognition (ASR) systems rely on hand-crafted feature extraction methods such as log Mel filterbank, Mel-frequency cepstral coefficients (MFCC) or Gammatone features [1, 2, 3]. The advance of neural modeling in ASR raised the question whether feature extraction should rather be a part of the neural acoustic model (AM). This could avoid potential information loss which might result in suboptimal performance. Besides some early methods using feed-forward neural networks (FFNNs) [4], mainly convolutional approaches have been proposed. They seem especially suitable for this task as they can operate similar to previous feature extraction methods but take advantage of learnable filters [5, 6, 7, 8]. In [9], the authors propose to use parameterized sinc functions to implement band-pass filters with few parameters in a way that is easy to interpret.

Recently, one line of research on learnable features for ASR focused on architectures which can be pre-trained on

audio-only data in an unsupervised fashion [10, 11, 12, 13, 14]. Untranscribed audio data is easier to obtain and the pre-training enables exploiting data which could otherwise not be incorporated. This is shown to be helpful in multiple cases, especially for low resource languages [13, 14]. Yet, it remains unclear whether these approaches offer additional benefits for research tasks with a closed set of training data. Additionally, none of these features has been tested with a hybrid neural network (NN)/hidden Markov model (HMM) system beyond basic connectionist temporal classification (CTC) to the best of our knowledge.

While lots of different architectures for ASR systems have been emerging in recent years [1, 15, 16, 17], hybrid NN/HMM systems constitute the state-of-the-art on many tasks like TED-LIUM release 2 [2], CHiME 6 [3] or AMI [18] and achieve competitive results on LibriSpeech [19]. It is therefore worthwhile to investigate how they can be further improved and how they compare to other high performance systems. This work aims at investigating and comparing the usefulness of different features – including learnable ones – for hybrid NN/HMM systems. Besides using a pre-trained feature extractor (FE), we explore how a strong existing AM trained with different features on the same task can be exploited to reduce the training effort needed in order to deploy the new front-end. In the same way, i-vectors are integrated to retrofit the system with speaker adaptation.

The feature extraction methods used in this work as well as some considerations regarding the hybrid ASR systems are presented in Section 2. The experiments are described in Section 3 and discussed in Section 4. Section 5 concludes the paper.

2. METHODS

2.1. Feature Extraction

In the following, the feature extraction methods used in this work are explained. Note that for traditional hand-crafted feature extraction methods, there is a clear border between the FE and the subsequent AM, as depicted in Figure 1. If the neural network operates directly on the audio samples, this

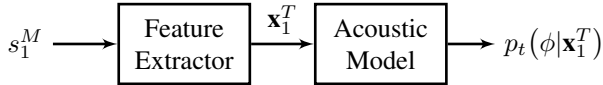


Fig. 1: Overall NN architecture with waveform samples s_1^M as input, the FE output \mathbf{x}_1^T and the AM which models the context dependent phoneme (CDp) label posterior probability $p_t(\phi|\mathbf{x}_1^T)$.

border blurs. Nevertheless, for the sake of clarity, we refer to the layers that replace the hand-crafted feature extraction as FE and to the remaining layers which are unchanged as AM.

2.1.1. Gammatone Features

Gammatone features are based on the Gammatone filter [20] which is designed to mimic the human auditory filter. They were first introduced for large vocabulary ASR in [21]. After pre-emphasizing the speech signal, a filterbank of Gammatone filters with center frequencies sampled from the Greenwood function [22] is applied. Temporal integration of the filter outputs' absolute values is typically performed using a Hanning window of 25 ms width with 10 ms shifts. After 10^{th} root compression, a cepstral decorrelation using the discrete cosine transform (DCT) and normalization techniques are applied.

2.1.2. Supervised Convolutional Features

The described Gammatone feature extraction pipeline motivated the approach in [8]. First, a convolutional layer for the purpose of time-frequency (tf) decomposition is applied on the raw waveform. In contrast to other works which train acoustic models directly on the waveform, the following envelope extraction is not performed by non-parameterized function such as max pooling, but rather with a rectification followed by low-pass filters which are shared between tf filters. By introducing multiple (N) low-pass filters, a multi-resolutional processing can be achieved. Finally, an additional non-linearity, e.g. logarithmic or root compression, may be applied, and the resulting features can be interpreted as critical band energies. Since the convolutional layers of this method are trained supervised, we refer to the features as SC features. The architecture is shown in Figure 2 (a).

2.1.3. wav2vec Features

The wav2vec features were first introduced in [10] and aim at learning audio representations in an unsupervised fashion from unlabeled data. These representations are supposed to improve downstream ASR performance. The network architecture consists of an encoder and a context network. A regular model with 5 layers in the encoder and 9 layers in the context network as well as a large model with 7 and 12 layers

respectively were presented. All convolutional layers have 512 channels and use the ReLU activation function. In the encoder, skip connections are used to allow for better convergence of the large model. The output of the context network can be used as input to the ASR system. Figure 2 (b) depicts the model.

A contrastive loss is used as the objective during unsupervised pre-training. As no pre-training is conducted in this work, we refrain from giving the details. There have been follow up works using quantization and self-attention [11, 12], however, we focus on the fully convolutional version.

2.1.4. i-vectors

A popular speaker adaptation method in ASR is the addition of i-vectors [23, 24, 25] which aim at conveying speaker characteristics. We follow the procedure described in [23] to train and extract the i-vectors. Specifically, the silence frames in the input waveform are filtered out first. Then, Gammatone features with a context of 9 frames and a subsequent linear discriminant analysis (LDA) to reduce the feature dimension to 60 are used to estimate the i-vectors. The dimensionality of the i-vectors themselves is set to 200 and they are normalized to have unit Euclidean norm. We extract one i-vector per utterance.

2.2. Hybrid ASR Systems

The ASR experiments in this work are carried out using hybrid NN/HMM systems. We take the model described in [26] as the baseline for our work and use the same CDps and the same alignments. The AM architecture consists of 6 bi-directional long-short term memory (BLSTM) layers with 1000 units for each direction followed by a linear layer with softmax activation which maps to the 12001 CDp labels clustered by a classification and regression tree (CART) [26].

In [26], all AMs trained with the frame-wise cross entropy (CE) loss were trained from scratch, i.e., the network parameters were initialized randomly at the beginning of the training. However, if a pre-trained model is already available, its parameters may be used in the same way a pre-trained FE is used. This can reduce the training effort needed for deploying a new set of features.

A possible example of how to exploit this is presented in [27], where the goal is to obtain noise invariant features. A model trained on clean data is virtually divided into FE and classifier by grouping lower and upper layers. Then, the model is further trained on noisy data while keeping the classifier constant or updating its parameters with a small learning rate. The learning rate for the FE is not reduced and its parameters are either initialized randomly or from the clean model.

In this work, we use the best model obtained by frame-wise CE training in [26] as a baseline to use the parameters

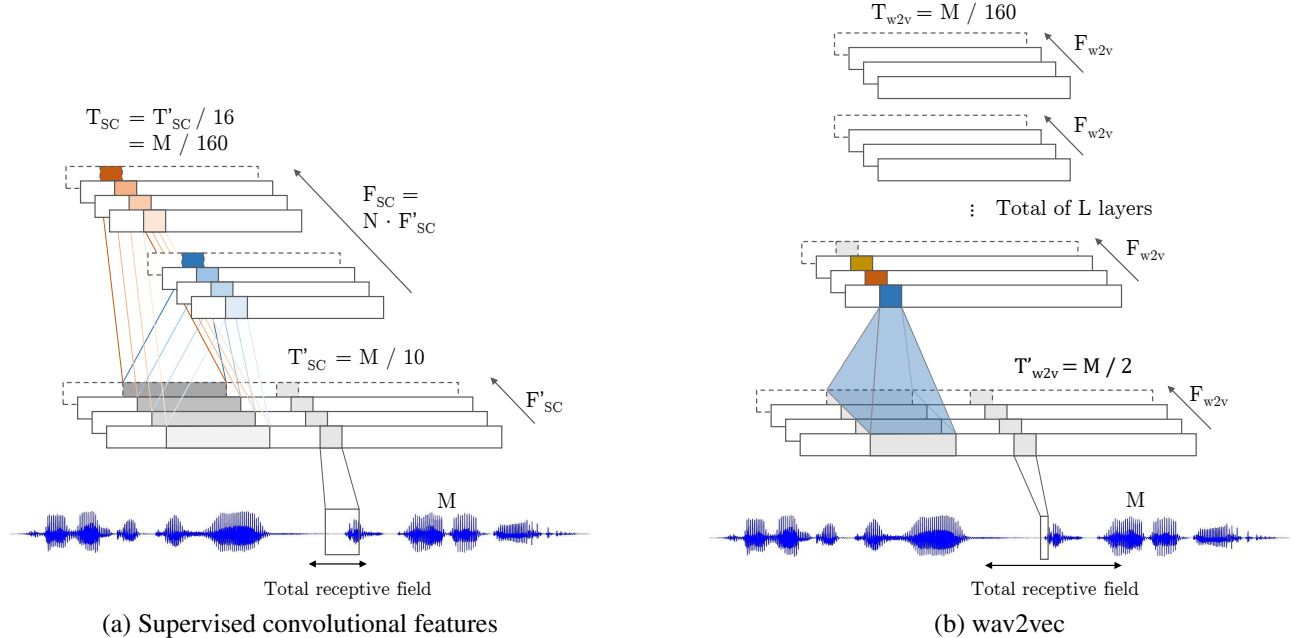


Fig. 2: Illustration of the supervised convolutional and wav2vec architectures. The horizontal rectangles represent feature channels over time. The last dashed rectangle in each layer indicates that there are many channels per layer. The connecting lines between layers show how an entry is calculated based on the previous layer’s activations. Typically, an entry is calculated based on all feature channels of the previous layer. E.g., in wav2vec, the first filter is represented in blue and the entries are computed using all previous channels. The other filters represented by other colors work analogously. In the supervised convolutional features extraction, the second layer pools each channel over time with multiple low pass filters. So the entries in the first channel are computed only based on values from the previous first channel. M is the number of waveform samples and the subsampling in time can be seen from T_{SC} and T_{w2v} respectively (ignoring boundary effects). The final sequence length is the same, $T_{SC} = T_{w2v}$. The number of features is F_{w2v} and F_{SC} respectively where N is the number of low-pass filters for pooling in the second supervised convolutional (SC) layer. While the SC architecture contains 2 layers, wav2vec large comprises $L = 19$ layers in total. Activation functions and normalizations are not depicted here. Note that the sizes are not drawn to scale.

of the BLSTM layers as well as the softmax output layer and continue training the model after resetting the learning rate. The baseline model was trained with 50-dimensional Gammatone features. If the dimension of the features differs from this, the first layer’s parameters can not be used because the weight matrix has a different shape. We circumvent this by initializing the first layer randomly and only adopt the parameters of the subsequent layers. Optionally, a random initialization of further subsequent layers may be applied. Unlike [27], we update all parameters with the same learning rate.

While most experiments are done using the frame-wise CE loss, we also deploy a lattice-based version of the state-level minimum Bayes risk (sMBR) criterion to conduct sequence discriminative training. The lattices are created using the frame-wise CE-trained system. Then, we reset the learning rate and continue training now using the sMBR criterion and an additional CE smoothing. Note that this learning rate reset and continued training resembles the procedure described above. However, none of the BLSTM layers is ini-

tialized randomly.

2.3. Language Models

During recognition, the official 4-gram language model (LM) as well as a custom LSTM LM are used in the first pass decoding [28]. Additionally, the lattices are rescored utilizing a Transformer LM [29]. The LSTM LM consists of a linear layer followed by two LSTM layers while the Transformer LM is built by 96 layers with 8 attention heads and both have an output softmax layer over the full 200k vocabulary. Both are identical with the ones used in [26]. The perplexities can be found in Table 1.

3. EXPERIMENTS

The experiments carried out for this work are presented in the following section. As in [26], the RWTH open-source

Table 1: LM perplexities on LibriSpeech test sets [29].

LM	dev		test	
	clean	other	clean	other
4-gram	151.7	140.6	158.1	145.7
LSTM	60.2	60.2	64.8	61.7
Transformer	53.2	54.2	57.6	55.0

toolkits RASR [30, 31] and RETURNN¹ [32, 33] were used for training and recognition.

3.1. LibriSpeech Dataset

The experiments are carried out on LibriSpeech which contains English read speech sampled at 16 kHz. Both the supervised ASR training as well as the unsupervised pre-training in [10] are done on the 960 hours training data, so no external data was used for the wav2vec features. The LMs are trained on the transcriptions of the 960 hours as well as the additional 800M-word text-only data. Testing is done on the common *dev-clean*, *dev-other*, *test-clean* and *test-other* subsets.

3.2. Features

All features are extracted from the raw waveform inside RETURNN. This allows for an easy integrability but also makes the feature extraction a fully differentiable part of the neural network which can be helpful in case the ASR system is integrated with neural pre-processing steps such as speech separation or enhancement.

For the Gammatone features, we therefore apply a convolutional layer with 50 Gammatone filters and a size of 640 samples which corresponds to 40 ms. This is in contrast to the Gammatone feature extraction in [26], where RASR is used to extract Gammatone features and an implementation using infinite impulse response (IIR) filters is deployed. The remaining steps are carried out as described in Section 2.1.1. Finally, we deploy batch normalization.

For the supervised convolutional features, we apply 150 filters for the tf decomposition each with a size of 256 samples and a stride of 10 and take the absolute values of the outputs. The envelope extraction is performed by 5 filters with a size of 40 samples and a stride of 16. After taking the root of the absolute values, a layer normalization completes the feature extraction.

The wav2vec features are extracted as described in [10]. We adapt the architecture denoted as *Large* in the original work, i.e., a convolutional encoder network consisting of 7 layers and a convolutional context network of 12 layers are used. The resulting features are 512-dimensional. Parameters of pre-trained models were made available by the authors

¹Training configuration files for the AM will be made available at <https://github.com/rwth-i6/returnn-experiments/>.

Table 2: Overview of the possibility to update the parameters during unsupervised and supervised training as well as the number of parameters for Gammatone (GT), supervised convolutional (SC) and wav2vec (w2v) features. The number of parameters in the AM varies only because of the different input dimensions to the first BLSTM layer.

Features	Trainable		#params	
	unsup.	sup.	FE	AM
GT	✗	✗	35k	152M
SC		✓	40k	158M
w2v	✓		29M	156M

and we use them in our experiments. As stated before, the pre-training was conducted on the LibriSpeech training data and no additional untranscribed audio data was used. The unsupervised pre-training is not subject of this work, however, we continue training the parameters using the supervised loss during ASR training.

Table 2 summarizes these features and outlines that while the Gammatone features are always extracted with a fixed set of parameters, the parameters of the supervised convolutional features can be updated during the supervised training and the wav2vec model can additionally even be pre-trained in an unsupervised fashion. Additionally, we can observe that the number of parameters for Gammatone and supervised convolutional features is comparable, while the wav2vec architecture exceeds this by a large factor. Note that the number of parameters for the Gammatone feature extraction is based on the finite impulse response (FIR) implementation described above.

The subsampling in the wav2vec encoder is chosen such that a feature frame is obtained every 10 ms which corresponds to the frame rate typically used for Gammatone features, e.g. in [26]. This is convenient for training the hybrid ASR system with the frame-wise CE criterion as the alignments obtained using Gammatone features can still be used. Additionally, it facilitates frame-wise feature combination. One frame of the wav2vec features in the regular model has a total receptive field of about 210 ms whereas it increases to about 810 ms for the large model. This is significantly larger than for the convolutional implementation of the Gammatone features or the features of the supervised convolutional approach, where the total receptive field corresponds to about 65 ms or 40 ms of audio, respectively.

When comparing different features, it seems natural to apply feature combination for possible further improvements. This is done using a frame-wise concatenation of all features here. The i-vectors are computed per utterance, therefore, this single vector is repeated and concatenated in each frame.

Table 3: Word error rates (WERs) [%] with and without unsupervised pre-training of wav2vec features. Pre-training and supervised training are conducted on the 960h of LibriSpeech and 16 epochs of supervised training were performed here. All results are obtained on *dev-other* using the 4-gram LM.

Features	Epochs unsup.	WER
w2v	0	10.1
	129	9.7

4. EXPERIMENTAL RESULTS

4.1. Effect of Training Strategies

The effect of performing unsupervised pre-training of the wav2vec features is shown in Table 3. The experiment with 0 epochs of unsupervised pre-training corresponds to the case where the wav2vec model is initialized randomly and only updated in supervised training. We can observe that it helps to use pre-training even in a setting where no additional untranscribed data is used. However, the relative gains are smaller than the ones on Wall Street Journal (WSJ) in [10] when using additional data which is intuitive.

In the baseline model of [26], the pre-extracted input sequence of Gammatone features is split into chunks using a windowing process during mini-batch construction to speed up the training. The window comprises 50 frames and is shifted by 25 frames, which corresponds to 0.5 s and 0.25 s of the input signal, respectively. A collection of chunks obtained by this windowing process then makes up one mini-batch. Now, since the network input is the raw waveform, we perform chunking on the audio directly. As described in Section 2.1.3, the total receptive field of the wav2vec large model is about 810 ms, which exceeds the chunk size of 0.5 s. Therefore, chunking was performed with a size of 3 s and a shift of 1.5 s for our experiments in Table 3. We observed, however, that reducing it to a size of 1 s and a shift of 0.5 s reduced the training time by about 25% while sacrificing only 0.1% absolute in WER which is why we apply this setting for the subsequent experiments.

Table 4 depicts the effect of using a pre-trained AM. We use the baseline from [26] which was trained for 12.5 epochs and train for 8 additional epochs here. Consistent improvements across all features can be seen. The parameters of the pre-trained AM are exploited as described in Section 2.2 and only the first BLSTM layer is initialized randomly. Random initialization of further BLSTM layers resulted in inferior performance.

4.2. Comparison of Architectures

As observed in Table 4, the performance of the ASR system using wav2vec features is on par with the baseline in [26] when training the AM from scratch as well as with a

Table 4: WERs [%] with and without supervised pre-training of the AM. In case of pre-training, the AM was trained for 12.5 epochs with GT features and we train for 8 further epochs with the stated features (see Section 2.2). Pre-training and supervised training are conducted on the 960h of LibriSpeech. All results are obtained on *dev-other* using the 4-gram LM.

Features	Pretrained GT-AM (12.5 epochs)	Epochs		WER
		unsup.	sup.	
GT [26]	no	-	12.5	9.6
GT	yes		8	8.7
SC	no		16	9.6
	yes	8	9.0	
w2v	no	129	16	9.7
	yes		8	8.8

pre-trained AM. This proves the features’ suitability for hybrid systems. Similarly, the supervised convolutional features’ performance is in the same ballpark, although slightly worse when using the pre-trained AM. The gap compared to the Gammatones matches the degradation of 3-4% relative found in [8] for FFNN. However, the results using an LSTM back-end were significantly worse in [8] which we do not observe here. A possible explanation is that the front-end parameters were trained with a FFNN and kept fixed during the LSTM training in [8] while we train them together with the BLSTM layers. To be fair, it should be noted that the wav2vec features might benefit from a system combination effect as parameters trained by the authors of [10] are used in our experiments.

4.3. Feature Combination

Next, Table 5 shows an overview of the results obtained by different features and their combinations. Combining Gammatone and wav2vec features improves the performance while adding supervised convolutional features to the other ones shows only marginal differences. A significant gain of 5-8% relative is attained when adding i-vectors to each of the previous features. This demonstrates the feasibility of retrofitting an AM with speaker adaptation by integrating i-vectors in the presented way. It also shows that even for the learnable and pre-trained wav2vec features, additional speaker information is helpful. The best performance can be achieved with a combination of Gammatone, wav2vec and i-vectors. A further iteration of resetting the learning rate and initializing the first layer randomly did not yield an improvement for this model.

We perform sequence discriminative training for 1.5 epochs and recognition with the LSTM LM as well as lattice rescoring using the Transformer LM for the best model. Additionally, an optimization of decoding parameters is con-

Table 5: WERs [%] for different feature combinations. The pre-trained AM was used in all cases and the wav2vec model where applicable. Pre-training and supervised training are conducted on the 960h of LibriSpeech. All results are obtained on *dev-other* using the 4-gram LM.

Features				WER
GT	w2v	SC	i-vec	
•				8.7
	•			8.8
		•		9.0
•	•			8.4
•		•		8.8
	•	•		8.7
•			•	8.2
	•		•	8.1
		•	•	8.4
•	•		•	8.0
•	•	•	•	8.1

ducted for the sequence discriminatively trained models. The results are depicted in Table 6. The comparison shows that the improvements over [26] that we achieve for the frame-wise CE trained model are stronger on the *other* set than on the *clean* one. Furthermore, the benefit of sequence discriminative training is smaller than it was in [26]. We suspect that while the positive influence of the new training criterion persists, the effect of a learning rate reset and some further epochs of training is already included in the frame-wise CE trained model here, unlike in [26], which could explain the observation.

In contrast, the relative improvements obtained by using the LSTM LM remain almost constant. As observed for the frame-wise CE-trained model, the gains over our previous best system are stronger on the *other* set, where we achieve a relative improvement of 6%, while the result is 4% relative better on *test-clean*.

The results presented in [19] are the best ones currently published on LibriSpeech with a hybrid system. The main reasons for the remaining gap in performance of our model are likely the different layer architecture as well as the time spent for supervised training of the AM. While 200 epochs were carried out in [19], we achieve our final results after training the AM for only 22 epochs. A further gap can be observed when comparing to the state-of-the-art performance reported in [34] where the number of training epochs conducted is unclear.

5. CONCLUSION

We apply learnable feature extractors, namely wav2vec and supervised convolutional features, in a hybrid BLSTM/HMM ASR system and show their suitability for this type of model.

Table 6: WERs [%] with a combination of Gammatone, wav2vec and i-vectors (comb.) on *test-clean* and *test-other* compared to other models from the literature. The pre-trained AM (12.5 epochs) is used and we train further 8 epochs. Sequence discriminative training contributes another 1.5 epochs resulting in 22 supervised epochs in total. Trf denotes Transformer and Cnf-Trnsd stands for the Conformer Transducer from [34]. Pre-training and supervised training are conducted on the 960h of LibriSpeech.

Features	AM	Epochs		Seq. dscr.	LM	WER	
		Unsup.	Sup.			clean	other
GT [26]	BLSTM-Hybrid	-	12.5	no	4-gram	4.4	10.0
Comb.		129	20.5	yes		3.9	8.6
GT [26]		-	14			3.8	8.8
Comb.		129	22		3.5	8.1	
GT [26]		-	14		LSTM	2.6	5.5
Comb.		129	22		2.5	5.2	
GT [26]	-	14	+Trf		2.3	5.0	
Comb.	129	22	2.2	4.7			
LgMel [19]	Trf-Hybrid	-	200	no	4gr+Trf	2.1	4.2
LgMel [34]	Cnf-Trnsd				LSTM	1.9	3.9

Using a pre-trained wav2vec module helps, however, it does not outperform the non-trainable features and features trained only with a supervised loss in a setting without additional untranscribed data. It is demonstrated that an existing AM trained on Gammatone features can be exploited in a similar way for all features. A combination of Gammatone and wav2vec features as well as i-vectors leads to the best results outperforming our previous best system by 4% and 6% relative on *test-clean* and *test-other*, respectively.

6. ACKNOWLEDGEMENTS

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement n° 694537, project ”SEQCLAS”). The work reflects only the authors’ views and the European Research Council Executive Agency (ERCEA) is not responsible for any use that may be made of the information it contains.

This work was partially supported by the project HYKIST funded by the German Federal Ministry of Health on the basis of a decision of the German Federal Parliament (Bundestag).

We thank Markus Kitzka for extracting the i-vectors as well as Wei Zhou for his help on the LSTM LM decoding. We also thank the authors of wav2vec [10] for providing their model parameters.

7. REFERENCES

- [1] William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals, “Listen, attend and spell: a neural network for large vocabulary conversational speech recognition,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, Shanghai, China, Mar. 2016.
- [2] Wei Zhou, Wilfried Michel, Kazuki Irie, Markus Kitza, Ralf Schlüter, and Hermann Ney, “The RWTH ASR system for TED-LIUM release 2: Improving hybrid HMM with specaugment,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, Barcelona, Spain, May 2020, pp. 7839–7843.
- [3] Jun Du, Yan-Hui Tu, Lei Sun, Li Chai, Xin Tang, Mao-Kui He, Feng Ma, Jia Pan, Jian-Qing Gao, Dan Liu, et al., “The USTC-NELSLIP systems for CHiME-6 challenge,” in *Proc. The 6th International Workshop on Speech Processing in Everyday Environments (CHiME 2020)*, 2020, pp. 19–23.
- [4] Zoltán Tüske, Pavel Golik, Ralf Schlüter, and Hermann Ney, “Acoustic modeling with deep neural networks using raw time signal for LVCSR,” in *Proc. Interspeech*, Singapore, Sept. 2014, pp. 890–894, ISCA Best Student Paper Award.
- [5] Dimitri Palaz, Mathew Magimai Doss, and Ronan Collobert, “Convolutional neural networks-based continuous speech recognition using raw speech signal,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 4295–4299.
- [6] Pavel Golik, Zoltán Tüske, Ralf Schlüter, and Hermann Ney, “Convolutional Neural Networks for Acoustic Modeling of Raw Time Signal in LVCSR,” in *Proc. Interspeech*, Dresden, Germany, Sept. 2015, pp. 26–30.
- [7] Tara N Sainath, Ron J Weiss, Andrew Senior, Kevin W Wilson, and Oriol Vinyals, “Learning the speech front-end with raw waveform CLDNNs,” in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [8] Zoltán Tüske, Ralf Schlüter, and Hermann Ney, “Acoustic modeling of speech waveform based on multi-resolution, neural network signal processing,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, Calgary, Canada, Apr. 2018, pp. 4859–4863.
- [9] Mirco Ravanelli and Yoshua Bengio, “Speaker recognition from raw waveform with SincNet,” in *2018 IEEE Spoken Language Technology Workshop (SLT)*, 2018, pp. 1021–1028.
- [10] Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli, “wav2vec: Unsupervised pre-training for speech recognition,” in *Proc. Interspeech*, Graz, Austria, Sept. 2019, pp. 3465–3469.
- [11] Alexei Baevski, Steffen Schneider, and Michael Auli, “vq-wav2vec: Self-supervised learning of discrete speech representations,” in *International Conference on Learning Representations*, Sept. 2019.
- [12] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *arXiv preprint arXiv:2006.11477*, 2020.
- [13] Alexis Conneau, Alexei Baevski, Ronan Collobert, Abdelrahman Mohamed, and Michael Auli, “Unsupervised cross-lingual representation learning for speech recognition,” *arXiv preprint arXiv:2006.13979*, 2020.
- [14] Kazuya Kawakami, Luyu Wang, Chris Dyer, Phil Blunsom, and Aaron van den Oord, “Learning robust and multilingual speech representations,” *arXiv preprint arXiv:2001.11128*, 2020.
- [15] Alex Graves, “Sequence transduction with recurrent neural networks,” in *Representation Learning Workshop, Int. Conf. on Machine Learning (ICML)*, Edinburgh, Scotland, June 2012.
- [16] Wei Zhou, Simon Berger, Ralf Schlüter, and Hermann Ney, “Phoneme based neural transducer for large vocabulary speech recognition,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, May 2021, To appear.
- [17] Ehsan Variani, David Rybach, Cyril Allauzen, and Michael Riley, “Hybrid autoregressive transducer (HAT),” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6139–6143.
- [18] Naoyuki Kanda, Yusuke Fujita, Shota Horiguchi, Rintaro Ikeshita, Kenji Nagamatsu, and Shinji Watanabe, “Acoustic modeling for distant multi-talker speech recognition with single-and multi-channel branches,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 6630–6634.
- [19] Frank Zhang, Yongqiang Wang, Xiaohui Zhang, Chunxi Liu, Yatharth Saraf, and Geoffrey Zweig, “Fast, simpler and more accurate hybrid ASR systems using word-pieces,” *arXiv preprint arXiv:2005.09150*, 2020.
- [20] AMHJ Aertsen, Peter IM Johannesma, and DJ Hermes, “Spectro-temporal receptive fields of auditory neurons in the grassfrog,” *Biological Cybernetics*, vol. 38, no. 4, pp. 235–248, 1980.

- [21] R. Schlüter, I. Bezrukov, H. Wagner, and H. Ney, “Gammatone features and feature combination for large vocabulary speech recognition,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, Honolulu, HI, USA, Apr. 2007, pp. 649–652.
- [22] Donald D Greenwood, “A cochlear frequency-position function for several species – 29 years later,” *The Journal of the Acoustical Society of America*, vol. 87, no. 6, pp. 2592–2605, 1990.
- [23] Markus Kitzka, Pavel Golik, Ralf Schlüter, and Hermann Ney, “Cumulative adaptation for BLSTM acoustic models,” in *Proc. Interspeech*, Graz, Austria, Sept. 2019, pp. 754–758.
- [24] Wayne Xiong, Lingfeng Wu, Fil Alleva, Jasha Droppo, Xuedong Huang, and Andreas Stolcke, “The Microsoft 2017 conversational speech recognition system,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5934–5938.
- [25] Naoyuki Kanda, Rintaro Ikeshita, Shota Horiguchi, Yusuke Fujita, Kenji Nagamatsu, Xiaofei Wang, Vimal Manohar, Nelson Enrique Yalta Soplin, Matthew Maciejewski, Szu-Jui Chen, et al., “The Hitachi/JHU CHiME-5 system: Advances in speech recognition for everyday home environments using multiple microphone arrays,” in *Proc. CHiME-5*, 2018, pp. 6–10.
- [26] Christoph Lüscher, Eugen Beck, Kazuki Irie, Markus Kitzka, Wilfried Michel, Albert Zeyer, Ralf Schlüter, and Hermann Ney, “RWTH ASR systems for LibriSpeech: Hybrid vs Attention,” in *Proc. Interspeech*, Graz, Austria, Sept. 2019.
- [27] Shucong Zhang, Cong-Thanh Do, Rama Doddipatla, and Steve Renals, “Learning noise invariant features through transfer learning for robust end-to-end speech recognition,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 7024–7028.
- [28] Eugen Beck, Wei Zhou, Ralf Schlüter, and Hermann Ney, “LSTM language models for LVCSR in first-pass decoding and lattice-rescoring,” *arXiv preprint arXiv:1907:NN*, July 2019.
- [29] Kazuki Irie, Albert Zeyer, Ralf Schlüter, and Hermann Ney, “Language modeling with deep transformers,” in *Proc. Interspeech*, Graz, Austria, Sept. 2019, pp. 3905–3909, ISCA Best Student Paper Award.
- [30] David Rybach, Stefan Hahn, Patrick Lehnen, David Nolden, Martin Sundermeyer, Zoltán Tüske, Simon Wiesler, Ralf Schlüter, and Hermann Ney, “RASR - the RWTH Aachen University open source speech recognition toolkit,” in *Proc. IEEE Automatic Speech Recog. and Understanding Workshop (ASRU)*, Waikoloa, HI, USA, Dec. 2011.
- [31] Simon Wiesler, Alexander Richard, Pavel Golik, Ralf Schlüter, and Hermann Ney, “RASR/NN: The RWTH neural network toolkit for speech recognition,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, 2014, pp. 3281–3285.
- [32] Patrick Doetsch, Albert Zeyer, Paul Voigtlaender, Ilya Kulikov, Ralf Schlüter, and Hermann Ney, “RE-TURNN: the RWTH extensible training framework for universal recurrent neural networks,” in *Proc. IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, New Orleans, LA, USA, 2017.
- [33] Albert Zeyer, Tamer Alkhouli, and Hermann Ney, “RE-TURNN as a generic flexible neural toolkit with application to translation and speech recognition,” in *Annual Meeting of the Assoc. for Computational Linguistics*, Melbourne, Australia, July 2018.
- [34] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al., “Conformer: Convolution-augmented transformer for speech recognition,” *arXiv preprint arXiv:2005.08100*, 2020.