

# Two Front-Ends, One Model : Fusing Heterogeneous Speech Features for Low Resource ASR with Multilingual Pre-Training

Anonymous ACL submission

## Abstract

Transfer learning is widely applied in various deep learning-based speech tasks, especially for tasks with a limited amount of data. Recent studies in transfer learning mainly focused on either supervised or self-supervised perspectives. This work, however, seeks to incorporate the two schemes together towards low-resource automatic speech recognition (ASR) for minor and endangered language (EL) communities. We propose a general framework to use learned transformations to resolve time resolution differences between any speech features, allowing for fusion of any self-supervised representations or spectral features used in multilingual pre-training. Our experiments over two low-resource languages and three ELs demonstrate that the proposed framework can significantly improve the absolute average word error rate from 45.4% to 35.5%.

## 1 Introduction

End-to-end (E2E) approaches to ASR have shown promising results compared to hybrid approaches for not only high-resourced scenarios (Chiu et al., 2018; Karita et al., 2019; Pham et al., 2019; Guo et al., 2021), but also certain low-resource scenarios in which linguistic documentations are insufficient for building lexicon-dependent models (Grenoble et al., 2011; Zahrer et al., 2020; Shi et al., 2021a). On the other hand, end-to-end approaches to low-resource ASR are distinctly disadvantaged by a lower data efficiency (Lüscher et al., 2019) and language-mismatch with powerful self-supervised representations (Hsu et al., 2021).

One direction towards mitigating these low-resource issues is to incorporate knowledge from several languages into multilingual end-to-end models (Watanabe et al., 2017; Toshniwal et al., 2018; Kannan et al., 2019). When there is no training data available for the target languages, these systems can be applied in a zero-shot manner (Li

et al., 2020; Yan et al., 2021; Xu et al., 2021). Fortunately, many languages have small amounts of data which can be used to fine-tune large-scale multilingual models towards target languages, resulting in further improvements (Hou et al., 2020; Pratap et al., 2020; Adams et al., 2019; Li et al., 2021).

Another direction is to use self-supervised learning representations (SSLR) trained on large untranscribed corpora as a front-end feature for ASR, replacing conventional spectral features like log Mel filterbank coefficients (FBank) (Yi et al., 2020; Wu et al., 2020; Baevski et al., 2020; N et al., 2021; Chang et al., 2021; Liu et al., 2021). Although these approaches have shown improvements across many languages, performance depends on the relatedness between the SSLR training languages and the target language (Conneau et al., 2019).

In this work, we are interested in leveraging both multilingual pre-trained (MPT) models with conventional speech feature front-ends and various SSLRs as resources for our low-resource ASR systems. In particular, we seek to efficiently incorporate multiple speech features, which may or may not have the same time resolution, as fused inputs to our end-to-end models. We propose a general framework to fuse such heterogeneous speech features and investigate several different learnable transformations for the fusion (Sec 3). Then we describe one instance following this front-end fusion framework which combines HuBERT features (Hsu et al., 2021) with an MPT model trained on FBank features (Sec 4.2). We demonstrate experimentally that our method improves absolute average WER by 9.9% on three endangered languages, and two of other low-resource languages in Sec 4.3.

Further, our data, pre-trained models, and reproducible methods are released open-source<sup>1</sup> to promote future developments on several endangered (Totonac, Yoloxóchitl Mixtec, and Highland Puebla Nahuatl) and low-resourced (Arabic and Tamil) lan-

<sup>1</sup>Available after the double-blind review period.

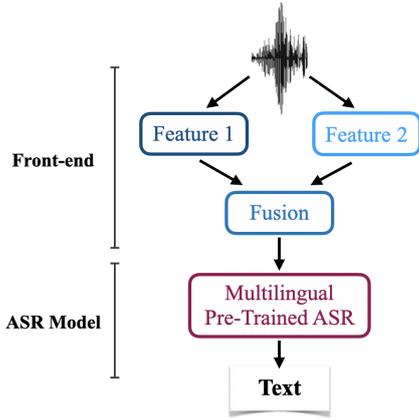


Figure 1: The architecture of our proposed model. The general formulation can be found in Sec 3.1.

081 guages. Notably, our released Totonac ASR data is  
 082 the first publicly available annotated speech corpus.

## 083 2 Motivation

084 Multilingual models have not only reached large-  
 085 scales, but they have also demonstrated high effi-  
 086 ciency in modeling multiple languages. For in-  
 087 stance, (Li et al., 2021) found multilingual training  
 088 boosted low-resourced languages while also avoid-  
 089 ing degradation of high-resource languages. Sep-  
 090 arately, SSLRs encode general-purpose informa-  
 091 tion about speech that can apply to various down-  
 092 stream tasks, including ASR (Yang et al., 2021).  
 093 For low-resource scenarios, Yi et al. (2020) found  
 094 that wav2vec2 was useful for 6 low-resource lan-  
 095 guages, suggesting that SSLR can replace spectral  
 096 features like FBanks in these cases.

097 Rather than viewing MPT and SSLR as two dis-  
 098 tinct techniques, we argue that for low-resourced  
 099 ASR these ideas are critically intertwined for sev-  
 100 eral reasons. (1) It is difficult to maintain broad  
 101 compatibility of supervised MPT models if differ-  
 102 ent front-ends are preferred for different scenarios.  
 103 (2) For low-resource ASR, there are often domain  
 104 mismatches between SSLR pre-training data and  
 105 ASR fine-tuning for a target language, leading to  
 106 unstable performances over when relying only on  
 107 only an SSLR front-end. (3) SSLRs are often pre-  
 108 trained exclusively over major languages like En-  
 109 glish (e.g., Hubert), again leading to unstable per-  
 110 formances depending on cross-lingual similarities.

## 111 3 Methodology

### 112 3.1 Heterogeneous Speech Feature Fusion

113 Historically, multi-layer perceptron tandem fea-  
 114 tures are concatenated with spectral features to

115 reach better performances (Hermansky et al., 2000;  
 116 Zhu et al., 2004; Lal and King, 2013). Following  
 117 their insights, Chen et al. (2021) performed the con-  
 118 catenation between STFT and SSLR features for  
 119 speech separation by repeating the SSLR feature  
 120 across the time-domain. While this repeat-based  
 121 method ameliorates the dimension mismatch of  
 122 features, it does not necessarily produce optimally  
 123 fused features for a particular task at hand.

124 We propose a general framework to fuse any  
 125 SSLRs and spectral features through learnable  
 126 transformations, allowing for the joint use of  
 127 different supervised pre-training models with  
 128 self-supervised representations. Such learnable  
 129 fusions have been employed in various multi-  
 130 source/multimodal applications previously (Li-  
 131 bovický and Helcl, 2017; Hori et al., 2017). The  
 132 framework is illustrated in Figure 1 and formulat-  
 133 ed as follows:

134 For an utterance, denote  $\mathbf{X}^{\mathcal{F}_1} = (\mathbf{x}_t^{(1)} \in$   
 135  $\mathbb{R}^{D_1} | t = 1, \dots, T_1)$  as a feature type  $\mathcal{F}_1$ , where  
 136  $D_1$  is the dimension of the feature at each frame  
 137 and  $T_1$  stands for the number of frames. Similarity,  
 138 we could define  $\mathbf{X}^{\mathcal{F}_2}$  along with  $D_2$  and  $T_2$ . As  $D_1$   
 139 and  $D_2$ ,  $T_1$  and  $T_2$  may not necessarily be the same,  
 140 combining two front-ends cannot be achieved by  
 141 simply concatenation. In Figure 1, we introduce  
 142 a fusion block to fuse features from two hetero-  
 143 geneous features (noted  $\mathcal{F}_1$  and  $\mathcal{F}_2$ ). Specifically,  
 144 we perform learnable transformations over speech  
 145 features  $\mathbf{X}^{\mathcal{F}_1}$  and  $\mathbf{X}^{\mathcal{F}_2}$  with linear, recurrent, con-  
 146 volution, or attention-based neural architectures.  
 147 We then use RESHAPE so that the transforms of  
 148  $\mathbf{X}^{\mathcal{F}_1}$  and  $\mathbf{X}^{\mathcal{F}_2}$  have the same dimensions, as shown  
 149 in Equation (1).

$$\begin{aligned} \tilde{\mathbf{X}}^{\mathcal{F}_1} &= \text{RESHAPE}(\text{TRANSFORM}(\mathbf{X}^{\mathcal{F}_1})) \\ \tilde{\mathbf{X}}^{\mathcal{F}_2} &= \text{RESHAPE}(\text{TRANSFORM}(\mathbf{X}^{\mathcal{F}_2})) \end{aligned} \quad (1)$$

151 After that,  $\tilde{\mathbf{X}}^{\mathcal{F}_1}$  and  $\tilde{\mathbf{X}}^{\mathcal{F}_2}$  are concatenated at the  
 152 feature-dimension. In the next sub-section, we in-  
 153 troduce one particular instance of this framework.

### 154 3.2 Fusion between SSLR and FBank for 155 Multilingual Pre-trained (MPT) ASR

156 As discussed in Sec 2, several issues may raise  
 157 when performing a simple combination between  
 158 MPT and SSLR, though they can improve low-  
 159 resource ASR, respectively. Therefore, we seek  
 160 to fuse the requisite front-ends for our end-to-  
 161 end model, following the formulation in Sec 3.1.  
 162 Firstly, we pre-train a multilingual encoder-decoder

model with language identification and hybrid CTC/Attention objectives (Watanabe et al., 2017; Hou et al., 2020). This architecture is built upon the FBank, so we define our first front-end feature  $\mathbf{X}^{\text{FB}} = (\mathbf{x}_t^{\text{FB}} \in \mathbb{R}^{D_{\text{FB}}} | t = 1, \dots, T_{\text{FB}})$ . Secondly, we use HuBERT as our SSLR front-end feature (Hsu et al., 2021),  $\mathbf{X}^{\text{HUB}} = (\mathbf{x}_t^{\text{HUB}} \in \mathbb{R}^{D_{\text{HUB}}} | t = 1, \dots, T_{\text{HUB}})$ . Following our framework in Eq. (2), we compute  $\tilde{\mathbf{X}}^{\text{FB}}$  and  $\tilde{\mathbf{X}}^{\text{HUB}}$  before ultimately obtaining fused features  $\mathbf{X}^{\text{FUSE}}$  as follows:

$$\mathbf{X}^{\text{FUSE}} = \tilde{\mathbf{X}}^{\text{FB}} \oplus \tilde{\mathbf{X}}^{\text{HUB}}, \quad (2)$$

where  $\oplus$  denotes feature-dimension concatenation.

## 4 Experiments

### 4.1 Datasets

We use a combination of Commonvoice 5.1 (Ardila et al., 2020) and Voxforge (www.voxforge.org) for MPT. The corpus results in 5,029 hours of training data, including 52 languages from different language families.

We enroll three endangered languages which are not included in the multilingual corpus for testing, including YoloXóchitl Mixtec (YM), Highland Puebla Nahuatl (HPN), and Totonac. Though endangered, YM and HPN have around 100 hours of transcribed speech in their released version (Shi et al., 2021b,a). However, to simulate a low-resource scenario, we randomly select 5,000 utterances (around 10 hours) from the official training sets, but used the same validation and test sets, as introduced in (Shi et al., 2021b,a). Totonac is another EL, spoken in the northern sierras of Puebla and adjacent areas of Veracruz. In this work, we release a public available version of Totonac speech resources. The corpus includes 10 hours of speech (86 long recordings) with fine-grained transcription. We randomly select 70 recordings as the training set, 8 for validation, and 8 for testing.

In addition to the three endangered languages mentioned above, we perform experiments on Arabic (AR) and Tamil (TA) corpora from Commonvoice to assert the robustness of our proposed methods in an in-domain low resource scenario. Both Arabic and Tamil have 20 hours of speech.

### 4.2 Experimental Setups

**Baseline (A):** For all the languages, our baseline (namely Exp A in later sections) adopts the same transformer-based encoder-decoder architecture with CTC/Attention hybrid training (Kim et al.,

2017). The front-end in Exp A extracts FBank features at a frame length of 20ms and a frameshift of 8ms. The extracted FBank features are firstly subsampled with a convolutional block and then fed into the encoder-decoder. The encoder contains 12-layer self-attention blocks with 4-head attention and 512-dimensional hidden sizes. While the decoder has 6 cross-attention transformer blocks. Specaugmentation (Park et al., 2019) and speed perturbation are employed for data augmentation. For training, we use Adam optimizer and Noam scheduler with a 1.0 learning rate at peak. The warm-up step is set to 4,000, considering the low-resource scenario. All the parameters are initialized with Xavier uniform distribution (Glorot and Bengio, 2010). The ASR model is trained on byte-pair-encoding (BPE) units of 250. The same architecture and training configuration are aligned for the following experiments.

**Multilingual Pre-training (B):** MPT is performed on a large-scale corpus introduced in Sec 4.1. We follow the same pre-training strategy as (Hou et al., 2020). For each utterance, the model needs to generate a language ID token prior to ASR transcription. To keep a necessary coverage for transcribing all 52 languages, we set a BPE size of 7,000. Large batch size is applied here in order to stabilize the multilingual training. After the pre-training, Exp B is conducted with parameter initialization from the pre-trained model.

**Self-supervised Representation (C):** In our experiments, we employ HuBERT as the front-end.<sup>2</sup> To fully explore the potential of HuBERT, we select the HuBERT-large model pre-trained over 60k hours of LibriLight (Kahn et al., 2020; Ott et al., 2019). The SSLR wrapper provided in (Yang et al., 2021) is applied to extract high-dimensional features with 20ms frameshift. In Exp C, the model directly applies HuBERT representation for training, which is the same approach as in (Chang et al., 2021). As for ablation purposes, we also conduct experiments on only SSLRs front-end with the initialization from the MPT model. We name these experiments as Exp C' in the next section.

**Joint-system (D&E):** The joint-system incorporates both FBank and SSLR in model front-end. According to our settings, the resolution ratio be-

<sup>2</sup>We also conduct experiments on Wav2vec2 and Wav2vec2-XLSR (Conneau et al., 2019; Baevski et al., 2020). However, the performances are not stable for ASR training.

Exp	Front-end			ASR	CER/WER					
	FBank	SSLR	Align	MPT	Totonac	YM	HPN	AR	TA	Avg.
<b>A</b>	✓	✗	-	✗	17.3/50.6	26.2/50.8	51.5/77.6	15.4/29.2	6.1/ <b>19.0</b>	23.3/45.4
<b>B</b>	✓	✗	-	✓	17.9/50.3	24.8/47.5	34.4/64.6	12.9/26.7	8.2/24.0	19.6/42.6
<b>C</b>	✗	✓	-	✗	17.1/48.3	38.8/61.2	29.4/58.3	15.1/29.2	6.2/19.7	21.3/43.3
<b>D</b>	✓	✓	✓	✗	14.6/46.7	<b>19.1/42.6</b>	<b>23.1/52.4</b>	<b>8.4/22.4</b>	6.1/24.5	<b>14.3/37.7</b>
<b>E</b>	✓	✓	✓	✓	<b>14.4/45.6</b>	<b>20.0/40.0</b>	<b>25.1/52.1</b>	<b>9.2/20.2</b>	<b>5.9/19.4</b>	<b>14.9/35.5</b>

Table 1: Results comparing our proposed fused front-end models (**D**, **E**) with various single front-end baselines (**A**, **B**, **C**) for 5 low-resourced or endangered languages, as measured by Character (CER) and Word (WER) Error Rates.

Exp	Front-End	MPT	CER	WER
<b>A</b>	FBANK	✗	17.3	50.6
<b>B</b>	FBANK	✓	17.9	50.3
$\Delta(\mathbf{A} \rightarrow \mathbf{B})$	-	-	-0.6	+0.3
<b>C</b>	HUBERT	✗	17.1	48.3
<b>C'</b>	HUBERT	✓	20.9	58.4
$\Delta(\mathbf{C} \rightarrow \mathbf{C}')$	-	-	+3.8	+10.1

Table 2: Ablation study comparing improvement/degradation on Totonac when incorporating FBank-based MPT with FBank-based fine-tuning ( $\Delta(\mathbf{A} \rightarrow \mathbf{B})$ ) vs. incorporating FBank-based MPT with HuBERT-based fine-tuning ( $\Delta(\mathbf{C} \rightarrow \mathbf{C}')$ ).

Exp	Fusion type	CER	WER
<b>D</b>	Linear	<b>14.6</b>	<b>46.7</b>
<b>D1</b>	Repeat	15.8	48.2
<b>D2</b>	RNN	65.1	86.9
<b>D3</b>	Convolution	16.1	50.8
<b>D4</b>	Attention	17.8	52.4

Table 3: Ablation study comparing performance on Totonac of several fusion types (**D1-4**) with our proposed linear fusion model without ASR pre-training (**D**).

tween FBank and HuBERT feature is 5:2. As discussed in Sec 3.1, we apply transformations to both features and then reshape them into the same time resolution. The linear layer contains 400 units in our experiments. The model then consumes the concatenation of both HuBERT and FBank features as inputs. We name the experiments with the fusion block as Exp **D**. We refer to the experiments as Exp **E** if it is initialized with the MPT model. We default to using the linear fusion block. But, to investigate other potential methods for feature fusion, we conduct ablation studies (i.e., Exp **D1-D4**) over four other approaches, including simple repeating, convolution, recurrent, and attention-based fusion.

### 4.3 Results and Discussion

Table 1 provides results of our main experiments over the five low resource languages introduced in Sec 4.1. Our proposed model (Exp **E**) reaches the best performances, which improves 8.4% absolute

average CER and 9.9% absolute average WER than the baseline in Exp **A**. Besides, Exp **B** with MPT model leads to notable improvements over the baseline for some languages such as HPN, even though HPN was not in the set of languages used by the multilingual pre-training model. SSLR could also benefit some languages (e.g., Totonac) as indicated from Exp **C**. According to Exp **D**, the proposed fusion module demonstrates better performances for most languages, and also reaches the best average CER across the five languages.

Table 2 shows ablation study of Exp **C'** where we only use SSLR features with initialization from the MPT model originally built upon FBanks. Exp **C'** is degraded compared to Exp **A**, **B**, and **C**, suggesting incompatibility between SSLR and MPT model as the latter is trained on FBank features.

Table 3 provides results for Exp **D**, which consider the various fusion strategies discussed in Sec 3.1. It shows that linear fusion outperforms simple repeat method (i.e., Exp **D1**). Recurrent, convolution and attention-based networks strategies are also less effective than the linear approach in our context.

## 5 Conclusion

In this work, we suggest that self-supervised learning and supervised pre-training can jointly improve the ASR performances in low-resource scenarios. We propose a framework to align features with different time-domain resolutions and demonstrate the effectiveness of fusing various front-ends features. We also release a Totonac ASR corpus, serving for the purpose of endangered language documentation, and we show that our reproducible methods enable to get very good results in very low resource scenarios. In future works, we will investigate (1) multilingual pre-training with fused SSLR features; (2) zero-shot learning, especially for EL documentation purposes.

316  
317  
318  
319  
320  
321  
322  
323  
  
324  
325  
326  
327  
328  
  
329  
330  
331  
332  
333  
  
334  
335  
336  
337  
338  
339  
  
340  
341  
342  
343  
344  
345  
  
346  
347  
348  
349  
350  
351  
352  
  
353  
354  
355  
356  
357  
358  
  
359  
360  
361  
362  
363  
364  
  
365  
366  
  
367  
368  
369  
370  
371

## References

Oliver Adams, Matthew Wiesner, Shinji Watanabe, and David Yarowsky. 2019. Massively multilingual adversarial speech recognition. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 96–108.

Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M. Tyers, and Gregor Weber. 2020. Common voice: A massively-multilingual speech corpus. In *LREC*.

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems*, 33.

Xuankai Chang, Takashi Maekaku, Pengcheng Guo, Jing Shi, Yen-Ju Lu, Aswin Shanmugam Subramanian, Tianzi Wang, Shu-wen Yang, Yu Tsao, Hung-yi Lee, et al. 2021. An exploration of self-supervised pretrained representations for end-to-end speech recognition. *arXiv preprint arXiv:2110.04590*.

Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2021. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *arXiv preprint arXiv:2110.13900*.

Chung-Cheng Chiu, Tara N Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjuli Kannan, Ron J Weiss, Kanishka Rao, Ekaterina Gonnina, et al. 2018. State-of-the-art speech recognition with sequence-to-sequence models. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4774–4778.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *Proc. Interspeech*.

Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings.

Lenore A Grenoble, Peter K Austin, and Julia Sallabank. 2011. Handbook of endangered languages.

Pengcheng Guo, Florian Boyer, Xuankai Chang, Tomoki Hayashi, Yosuke Higuchi, Hirofumi Inaguma, Naoyuki Kamo, Chenda Li, Daniel Garcia-Romero, Jiatong Shi, et al. 2021. Recent developments on espnet toolkit boosted by conformer.

In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5874–5878. IEEE. 372  
373  
374

Hynek Hermansky, Daniel PW Ellis, and Sangita Sharma. 2000. Tandem connectionist feature extraction for conventional hmm systems. In *2000 IEEE international conference on acoustics, speech, and signal processing. Proceedings (Cat. No. 00CH37100)*, volume 3, pages 1635–1638. IEEE. 375  
376  
377  
378  
379  
380

Chiori Hori, Takaaki Hori, Teng-Yok Lee, Ziming Zhang, Bret Harsham, John R Hershey, Tim K Marks, and Kazuhiko Sumi. 2017. Attention-based multimodal fusion for video description. In *Proceedings of the IEEE international conference on computer vision*, pages 4193–4202. 381  
382  
383  
384  
385  
386

Wenxin Hou, Yue Dong, Bairong Zhuang, Longfei Yang, Jiatong Shi, and Takahiro Shinozaki. 2020. Large-Scale End-to-End Multilingual Speech Recognition and Language Identification with Multi-Task Learning. In *Proc. Interspeech 2020*, pages 1037–1041. 387  
388  
389  
390  
391  
392

Wei-Ning Hsu, Yao-Hung Hubert Tsai, Benjamin Bolte, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: How much can a bad teacher benefit asr pre-training? In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6533–6537. IEEE. 393  
394  
395  
396  
397  
398

J. Kahn, M. Rivière, W. Zheng, E. Kharitonov, Q. Xu, P. E. Mazaré, J. Karadayi, V. Liptchinsky, R. Collobert, C. Fuegen, T. Likhomanenko, G. Synnaeve, A. Joulin, A. Mohamed, and E. Dupoux. 2020. Libri-light: A benchmark for asr with limited or no supervision. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7669–7673. <https://github.com/facebookresearch/libri-light>. 407  
408

Anjuli Kannan, Arindrima Datta, Tara N. Sainath, Eugene Weinstein, Bhuvana Ramabhadran, Yonghui Wu, Ankur Bapna, Zhifeng Chen, and Seungji Lee. 2019. Large-Scale Multilingual Speech Recognition with a Streaming End-to-End Model. In *Proc. Interspeech 2019*, pages 2130–2134. 409  
410  
411  
412  
413  
414

Shigeki Karita, Nanxin Chen, Tomoki Hayashi, Takaaki Hori, Hirofumi Inaguma, Ziyang Jiang, Masao Someki, Nelson Enrique Yalta Soplín, Ryuichi Yamamoto, Xiaofei Wang, et al. 2019. A comparative study on transformer vs rnn in speech applications. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 449–456. IEEE. 415  
416  
417  
418  
419  
420  
421  
422

Suyoun Kim, Takaaki Hori, and Shinji Watanabe. 2017. Joint ctc-attention based end-to-end speech recognition using multi-task learning. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4835–4839. IEEE. 423  
424  
425  
426  
427

428	Partha Lal and Simon King. 2013. Cross-lingual automatic speech recognition using tandem features. <i>IEEE Transactions on Audio, Speech, and Language Processing</i> , 21(12):2506–2515.	483
429		484
430		485
431		486
432	Bo Li, Ruoming Pang, Tara N Sainath, Anmol Gulati, Yu Zhang, James Qin, Parisa Haghani, W Ronny Huang, Min Ma, and Junwen Bai. 2021. Scaling end-to-end models for large-scale multilingual asr. <i>arXiv preprint arXiv:2104.14830</i> .	487
433		488
434		489
435		490
436		
437	Xinjian Li, Siddharth Dalmia, Juncheng Li, Matthew Lee, Patrick Littell, Jiali Yao, Antonios Anastasopoulos, David R. Mortensen, Graham Neubig, Alan W Black, and Florian Metze. 2020. <a href="#">Universal phone recognition with a multilingual allophone system</a> .	491
438		492
439		493
440		494
441		495
442	Jindřich Libovický and Jindřich Helcl. 2017. <a href="#">Attention strategies for multi-source sequence-to-sequence learning</a> . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 196–202, Vancouver, Canada. Association for Computational Linguistics.	496
443		497
444		498
445		499
446		500
447		501
448	Andy T Liu, Shang-Wen Li, and Hung-yi Lee. 2021. Tera: Self-supervised learning of transformer encoder representation for speech. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 29:2351–2366.	502
449		503
450		
451		498
452		499
453	Christoph Lüscher, Eugen Beck, Kazuki Irie, Markus Kitza, Wilfried Michel, Albert Zeyer, Ralf Schlüter, and Hermann Ney. 2019. <a href="#">Rwth asr systems for librispeech: Hybrid vs attention</a> . <i>Interspeech 2019</i> .	500
454		501
455		502
456		503
457	Krishna D. N, Pinyi Wang, and Bruno Bozza. 2021. <a href="#">Using Large Self-Supervised Models for Low-Resource Speech Recognition</a> . In <i>Proc. Interspeech 2021</i> , pages 2436–2440.	504
458		505
459		506
460		507
461	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. <a href="#">fairseq: A fast, extensible toolkit for sequence modeling</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)</i> , pages 48–53.	508
462		509
463		
464		510
465		511
466		512
467		513
468	Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le. 2019. <a href="#">SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition</a> . In <i>Proc. Interspeech 2019</i> , pages 2613–2617.	514
469		515
470		516
471		
472		517
473	Ngoc-Quan Pham, Thai-Son Nguyen, Jan Niehues, Markus Müller, and Alex Waibel. 2019. Very deep self-attention networks for end-to-end speech recognition. <i>Proceedings of Interspeech 2019</i> , pages 66–70.	518
474		519
475		520
476		
477		521
478	Vineel Pratap, Anuroop Sriram, Paden Tomasello, Awni Hannun, Vitaliy Liptchinsky, Gabriel Synnaeve, and Ronan Collobert. 2020. <a href="#">Massively Multilingual ASR: 50 Languages, 1 Model, 1 Billion Parameters</a> . In <i>Proc. Interspeech 2020</i> , pages 4751–4755.	522
479		523
480		524
481		525
482		526
		527
		528
		529
		530
		531
		532
		533
		534
		535
		536
		537
		538
		539
		483
		484
		485
		486
		487
		488
		489
		490
		491
		492
		493
		494
		495
		496
		497
		498
		499
		500
		501
		502
		503
		504
		505
		506
		507
		508
		509
		510
		511
		512
		513
		514
		515
		516
		517
		518
		519
		520
		521
		522
		523
		524
		525
		526
		527
		528
		529
		530
		531
		532
		533
		534
		535
		536
		537
		538
		539

540 Qifeng Zhu, Barry Chen, Nelson Morgan, and An-  
541 dreas Stolcke. 2004. On using mlp features in lvcsr.  
542 In *Eighth International Conference on Spoken Lan-  
543 guage Processing*.