

LANGUAGE ADAPTIVE CROSS-LINGUAL SPEECH REPRESENTATION LEARNING WITH SPARSE SHARING SUB-NETWORKS

Yizhou Lu, Mingkun Huang, Xinghua Qu, Pengfei Wei, Zejun Ma

Speech & Audio Team, ByteDance AI Lab

{luyizhou.ritter, huangmingkun, xinghua.qu, pengfei.wei, mazejun}@bytedance.com

ABSTRACT

Unsupervised cross-lingual speech representation learning (XLSR) has recently shown promising results in speech recognition by leveraging vast amounts of unlabeled data across multiple languages. However, standard XLSR model suffers from language interference problem due to the lack of language specific modeling ability. In this work, we investigate language adaptive training on XLSR models. More importantly, we propose a novel language adaptive pre-training approach based on sparse sharing sub-networks. It makes room for language specific modeling by pruning out unimportant parameters for each language, without requiring any manually designed language specific component. After pruning, each language only maintains a sparse sub-network, while the sub-networks are partially shared with each other. Experimental results on a downstream multilingual speech recognition task show that our proposed method significantly outperforms baseline XLSR models on both high resource and low resource languages. Besides, our proposed method consistently outperforms other adaptation methods and requires fewer parameters.

Index Terms— representation learning, multilingual, language adaptation, speech recognition

1. INTRODUCTION

Automatic speech recognition (ASR) techniques have documented many success stories and been widely deployed in the wild [1, 2]. However, building a reasonable ASR system for practical application still requires tens of thousands of hours annotated data [2]. It remains a challenging task to extend ASR system to around 7000 languages in the world as most languages are lack of training data. Self-supervised learning, where training targets are derived from the input itself, has recently shown promising results in speech recognition [3–8]. Recent work on self-supervised learning can be categorized into two main approaches: contrastive learning approaches that distinguish a true sample from negative ones in a latent space [3, 4, 8], and reconstruction-based approaches that directly predict future frames [5, 7] or reconstruct masked inputs [6, 9]. Self-supervised learning provides an efficient way to learn from unlabeled data, and it allows the network to learn with fewer labeled data and generalize better. Its feasibility for monolingual speech recognition with limited amounts of labeled data has been shown in [8].

Self-supervised learning is further applied to multilingual setting in [10], namely cross-lingual speech representation learning (XLSR). Multilingual pre-training simplifies the procedure of training individual seed models for each language by supporting multiple languages with a single model. Recent studies [10, 11] also show that multilingual pre-training outperforms monolingual pre-training in low resource languages. From a multitask learning

perspective, the common knowledge learned from each task shall help the other related tasks learn and generalize better [12]. However, due to the vast diversity of pronunciation styles across different languages, the shared network often struggles in optimizing various languages simultaneously [13–15]. While low resource languages benefit from joint training with similar languages, high resource languages often suffer from the negative transfer problem, resulting in inferior performance [10]. Such performance degeneration becomes more significant as the model expands to more languages [14] or more training data [16], which obstacles the application of pre-trained models on downstream multilingual tasks [17].

In this work, we study language adaptive cross-lingual speech representation learning to alleviate the aforementioned interference problem. More importantly, we propose a novel language adaptive pre-training approach based on sparse sharing sub-networks (S3Net). Sparse sharing architecture is initially proposed in [18] to jointly learn sub-networks for multiple tasks and further applied in the field of neural machine translation [19–21]. Inspired by that, we extract a sub-network for each language and all the sparsely shared sub-networks are jointly trained to learn language adaptive speech representations, with each language only updating its corresponding sub-network. The key idea of this work is that redundant parameters can be pruned out separately for each language with minor or no performance degradation, and these parameters discarded by one language can be further utilized by other languages to learn better representations. Besides, it automatically distributes both shared and language specific parameters at each layer, without requiring any additional language specific component [22]. Given these designs, the learnt representations are expected to benefit more from positive transfer, while the negative transfer effects from dissimilar languages are mitigated.

There are several ways of extracting sub-networks, and in this paper we mainly focus on two ideas: one follows the procedure of lottery ticket hypothesis (LTH) [23] and the other is based on taylor expansion (TE) [24]. We compare our S3Net with baseline XLSR model and several other adaptation methods [22, 25]. Experimental results on a downstream multilingual task show the effectiveness of our proposed method. Specifically, S3Net yields an average 9.8% and 7.4% relative error reduction for XLSR base and large models respectively. Notably, it significantly improves the performance of high resource languages, achieving an average 17.8% and 16.7% relative error reduction. Moreover, S3Net also consistently outperforms other adaptation methods and requires fewer parameters.

2. RELATED WORK

Language interference problem has been investigated in many prior work on multilingual ASR. From a capacity perspective, [10] finds

that enlarging model size alleviates the interference problem, and [16] scales up their model to 10 billion parameters to accommodate multiple languages. Another line of this research tends to retain the language specific modeling ability. It is observed that simply adding a one-hot language identity (LID) vector to condition the multilingual model can boost the performance [26]. To better capture the language specific knowledge, previous studies usually augment networks with additional manually designed components, such as language specific weight matrices [15], light weight adapters [22], decoupled multilingual encoders [13] or decoders [14]. However, the inserted module size, structure and injection position are all important factors to consider, which requires additional efforts to fuse those additional modules into the original network [25].

We follow the line of language adaptation in this study. While previous work are mainly focused on supervised setting, to the best of our knowledge, this is the first work to apply language adaptive training to unsupervised pre-training models.

3. LANGUAGE ADAPTIVE PRE-TRAINING WITH SPARSE SHARING SUB-NETWORKS

In this section, we describe the proposed language adaptive pre-training approach based on sparse sharing sub-networks, namely S3Net. The proposed S3Net mainly includes three parts as shown in Figure 1: XLSR pre-training, extracting sparse sub-networks and language adaptive pre-training.

3.1. Pre-training of XLSR model

XLSR model [10] extends wav2vec 2.0 framework [8] that basically consists of three components: feature encoder, context network and quantization module. The multi-layer convolutional feature encoder takes raw waveform as input and maps the input into latent speech representations $\mathbf{z} = \mathbf{z}_1, \dots, \mathbf{z}_T$. Each \mathbf{z}_t represents an approximately 25ms wide audio with the frame stride of 20ms. The transformer [27] based context network utilizes the latent speech representations, and outputs contextual representations $\mathbf{c} = \mathbf{c}_1, \dots, \mathbf{c}_T$ that capture contextual information from full sequence. To provide the training targets, the latent representation \mathbf{z}_t is discretized to \mathbf{q}_t with a quantization module. The quantizer has $G = 2$ codebooks with $V = 320$ entries each, resulting in a set of over 100K codewords in total. Gumbel softmax enables the model to choose discrete codebook entries in a fully differentiable way [28].

The model is trained by solving a contrast learning task. Following [8], we randomly sample from all time steps with probability $s = 0.065$ as the starting indices, and then continuously mask the subsequent $M = 10$ steps. The goal of the contrast task is to distinguish a true encoded sample \mathbf{q}_t among distractors \mathbf{Q}_t that are sampled uniformly from other masked steps of the same sequence:

$$\mathcal{L} = -\log \frac{\exp(\text{sim}(\mathbf{c}_t, \mathbf{q}_t)/\kappa)}{\sum_{\tilde{\mathbf{q}} \in \mathbf{Q}_t} \exp(\text{sim}(\mathbf{c}_t, \tilde{\mathbf{q}})/\kappa)} + \lambda \mathcal{L}_d \quad (1)$$

where $\text{sim}(\cdot, \cdot)$ denotes cosine similarity, λ is the weight factor and diversity loss $\mathcal{L}_d = \frac{1}{GV} \sum_{g=1}^G \sum_{v=1}^V \bar{p}_{g,v} \log \bar{p}_{g,v}$ is employed to increase the use of codebook representations by maximizing the entropy over the selection probability of entry v in codebook g .

We form multilingual batches [29] by sampling with a multinomial distribution $p_l \sim \left(\frac{n_l}{N}\right)^\alpha$, where n_l is the number of hours for language l , N is the total number of hours and α is the sampling factor. We upsample the data from low resource languages with $\alpha = 0.5$ and for high resource languages we use natural sampling probability with $\alpha = 1.0$.

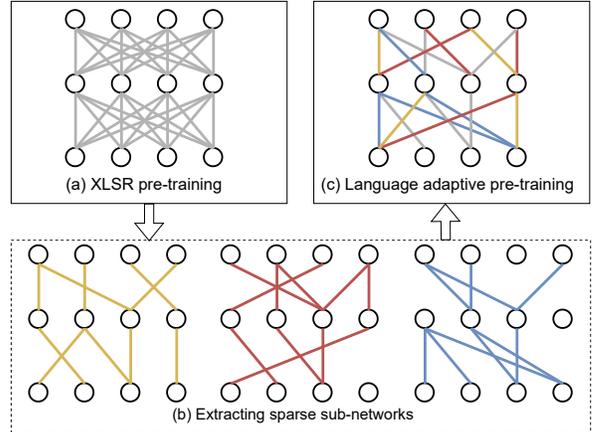


Fig. 1. Training procedure of the proposed S3Net. Connections between two nodes can be shared by a set of languages (gray lines), or be occupied by one specific language (colored lines). There also exists empty connections in the final S3Net model.

3.2. Extracting sparse sub-networks

Lottery ticket hypothesis [23] suggests that, a randomly initialized dense network contains small sub-networks (winning tickets) such that, when trained in isolation, can match the accuracy of full network. This motivates us to prune out redundant parameters for each language to make room for language specific modeling, and we hypothesize that these parameters discarded by one language can be further utilized by other languages to learn better representations. However, [18] shows that finding winning tickets from randomly initialized network is not stable, and a multitask warmup stage helps stabilize the process. Thus we extract sub-networks on top of a pre-trained XLSR model. Two different approaches of extracting sub-networks are explored in this work. The first one mainly follows the procedure of LTH [23] and the second one is based on importance scores that are estimated with first order Taylor expansion [24].

Extracting sub-networks with LTH. For the first approach we adopt a simple one-shot magnitude pruning (OMP) strategy instead of the iterative pruning strategy used in [23], as iterative pruning requires several rounds of training, pruning and resetting for each language, which is computationally expensive and time consuming. For each language l , we train XLSR model θ with only the specific language data \mathcal{D}^l for additional few steps to obtain $\hat{\theta}^l$. Sub-network is extracted from $\hat{\theta}^l$ and a proportion of p parameters with lowest magnitude are considered as unimportant for this language, thus being pruned out. For the sake of convenience, we use a fix pruning rate p for all languages in this work. The structure of sub-network for language l is controlled by a binary mask matrix \mathbf{m}^l which is initialized with $\mathbf{1}$. If a parameter $\hat{\theta}_i^l$ is pruned out, the corresponding item m_i^l in the mask is updated to zero. The sub-network employed by language l thus can be denoted as $\theta^l = \mathbf{m}^l \odot \theta$, where \odot denotes the element-wise product operation.

Extracting sub-networks with TE. Following [24], the importance of a parameter can be quantified by the error induced by removing it. A squared difference of prediction errors with and without the parameter θ_i is used to measure the importance score \mathcal{I}_i^l of language l :

$$\mathcal{I}_i^l = \left[\mathcal{L}(\mathcal{D}^l, \theta) - \mathcal{L}(\mathcal{D}^l, \theta | \theta_i = 0) \right]^2 \quad (2)$$

Model	Pre-trained data	es	fr	it	ky	nl	ru	sv	tt	zh	Avg
Number of unlabeled audio data		168h	353h	90h	17h	29h	55h	3h	17h	50h	
<i>Baselines from XLSR [10]</i>											
XLSR-Monolingual	CV-Mono*	6.8	10.4	10.9	29.6	37.4	11.6	63.6	21.4	31.4	24.8
XLSR-10	CV-Multi*	9.4	13.4	13.8	8.6	16.3	11.2	21.0	8.3	24.5	14.1
XLSR-10 (Large)	CV-Multi*	7.7	12.2	11.6	7.0	13.8	9.3	20.8	7.3	22.3	12.4
<i>Re-run baselines and our models</i>											
XLSR-10	CV-Multi	10.8	12.8	15.1	8.5	15.4	11.8	22.1	8.1	24.2	14.3
S3Net-TE		9.9	12.0	14.4	7.8	14.7	11.3	22.1	7.7	23.9	13.8
S3Net-LTH		8.7	10.8	12.4	7.5	14.1	10.1	22.0	7.2	22.9	12.9
XLSR-10 (Large)	CV-Multi	9.0	10.6	12.7	6.8	12.8	10.1	19.9	6.6	21.5	12.2
S3Net-TE (Large)		8.4	10.5	12.4	6.7	12.5	10.1	19.6	6.3	21.6	12.0
S3Net-LTH (Large)		7.3	9.2	10.4	6.3	12.1	9.4	19.5	6.1	21.5	11.3

Table 1. Evaluation results on CommonVoice dataset. The last column is the averaged PER on nine languages. Re-run baselines and our models are all pre-trained on ten languages, and evaluated on nine languages with shared vocabulary using CTC criterion. *: They use different version of the CommonVoice dataset, but the data size is the same as ours.

where \mathcal{D}^l denotes the training data of language l , and $\mathcal{L}(\mathcal{D}^l, \theta)$ is the averaged loss on \mathcal{D}^l calculated by Equation 1. However, directly calculating the importance score in Equation 2 requires $|\theta|$ times calculation of the loss function, which is infeasible. Fortunately, we can approximate the score with first order Taylor expansion [24]:

$$\mathcal{I}_i^l \approx (g_i^l \theta_i)^2 \quad (3)$$

where $g_i^l = \frac{\partial \mathcal{L}(\mathcal{D}^l, \theta)}{\partial \theta_i}$ is the gradient for θ_i that can be efficiently calculated with backward propagation. For each language, those parameters with top p lowest importance scores are pruned out, which is similar to the LTH based approach.

3.3. Language adaptive pre-training with S3Net

Once we obtain the masks $\mathbf{m}_1, \dots, \mathbf{m}_L$ for all languages, we continue to train the model θ with all multilingual data to learn language adaptive speech representations. We form batches that only contain utterances from one language and those batches are randomly sampled with the same sampling strategy as XLSR pre-training. It is noted that each language only maintains a sub-network, and for each batch only the sub-network from the corresponding language will participate the forward computation and be updated. In contrast to other adaptation methods that typically require manually designed language specific components, S3Net automatically distributes both shared and language specific parameters at each layer, and we expect that the model capacity is better allocated for each language.

4. EXPERIMENTS

4.1. Experimental setup

We use CommonVoice dataset ¹ [30] for pre-training. For a fair comparison with XLSR [10], we consider the following nine lan-

¹<https://commonvoice.mozilla.org/datasets>. We use the December 2020 release version.

guages for evaluation ²: Spanish (*es*), French (*fr*), Italian (*it*), Kyrgyz (*ky*), Dutch (*nl*), Russian (*ru*), Swedish (*sv*), Tatar (*tt*) and Chinese (*zh*). We pre-train base and large models on 1350 hours unlabeled multilingual dataset (CV-Multi) as in [10], which is composed of 782 hours data from above nine language plus additional 568 hours English (*en*) data. For fine-tuning, we use the evaluation splits from [31], which contains 1 hour labeled training data, 20 minutes validation data and 1 hour evaluation data for each language.

For our baseline models, we use the same model structure and training hyperparameters as XLSR [10]. The models are all trained on 64 GPUs, and we train 250k steps for base model and 400k steps for large model. For fine-tuning, we adopt Connectionist Temporal Classification (CTC) [32] criterion and evaluate the multilingual performance of the pre-trained model. A randomly initialized output layer with shared vocabulary is added on top of the pre-training model. We use Adam optimizer and the learning rate is warmed up for 2k updates to $5e-5$, keeps constant for 8k updates and then linearly decay for 10k updates. Phone error rate (PER) is reported following previous work.

4.2. Evaluation of proposed method

We denote S3Net using LTH for sub-network extraction as S3Net-LTH, and similarly S3Net-TE. For S3Net-LTH, we train the XLSR model with additional 50k steps separately for each language to extract the sub-networks. For S3Net-TE, we freeze the parameters and directly calculate importance scores with the gradients and weights. The prune rate is set to 0.4 for all base and large models, and by default the unstructured pruning is done layer by layer for each linear layer in context network. After obtaining the sub-network masks, we restart from the XLSR model and jointly train the sub-networks with all multilingual data for additional 50k steps. All hyper-parameters are tuned based on the performance of validation set.

The results of baseline XLSR models and the proposed S3Net are shown in Table 1. Specifically, we reproduce similar results

²There are initially ten languages for evaluation, but files of Turkish (*tr*) in the test set are missing in CommonVoice December 2020 release, so we exclude this language.

as [10], with an average PER of 14.3% and 12.2% for base and large models. However, the results on each individual language are different from [10] due to the different versions of pre-training data. We can see that both S3Net-TE and S3Net-LTH consistently outperform baseline XLSR models on each language, and S3Net-LTH models perform the best, with an average 9.8% and 7.4% relative PER reduction respectively for base and large models. Since high resource languages suffer more from language interference problem, it can be seen from the table that the proposed method achieves more performance improvements on high resource languages (es, fr, it), with 17.8% and 16.7% average relative PER reduction. We also observe that S3Net-TE models perform worse than S3Net-LTH models. This is because S3Net-TE employs a simpler sub-network extraction strategy, which may not yield satisfactory sub-networks. Besides, we also conduct monolingual fine-tuning experiments individually on each languages, and we find that it performs slightly better than multilingual fine-tuning, with average PER of 12.8% and 11.2% for S3Net-LTH base and large models.

4.3. Comparison of S3Net with other adaptation methods

We further compare the proposed method with other adaptation methods such as gating network [25] and adapter [22]. For gating network, one-hot LID embedding is used to learn the language specific scaling and biasing vectors [25]. We add gating network module directly after the output of feature encoder to modulate the latent representations. As for the adapter based language adaptation method, we add adapters before the input of each transformer layer. The structure of adapter module follows [22], and the projection dimension is set to 256 for base and large model. The gating network and adapter modules are inserted into the XLSR model and further trained with 50k steps as S3Net models. We show in Table 2 that all the adaptation methods improve the performance of baseline XLSR model, and our proposed S3Net-LTH outperforms all other adaptation methods while requiring fewer parameters.

Table 2. Comparison of different adaptation methods. Multilingual evaluation results are averaged on high resource languages (High), low resource languages (Low) and all nine languages (Avg).

Model	#Params	CV-Eval		
		High	Low	Avg
XLSR-10	95M	12.9	15.0	14.3
+ Gating Network	95M	12.2	14.7	13.9
+ Adapter	143M	11.5	14.1	13.2
S3Net-LTH	95M	10.6	14.0	12.9
XLSR-10 (Large)	317M	10.8	13.0	12.2
+ Gating Network	317M	10.4	12.8	12.0
+ Adapter	444M	10.4	12.9	12.1
S3Net-LTH (Large)	317M	9.0	12.5	11.3

4.4. Ablation study of sparse sub-networks

To analyze the influence of different sub-networks on the performance of S3Net, we perform extensive ablation studies in Table 3. In the standard setup, we individually extract sub-networks for each language (#Mask=10). To verify that the improvements of S3Net come from the language specific modeling, we jointly extract one sub-network for all languages (#Mask=1). We also experiment

with individually extracting sub-networks for high resource languages (en, es, fr, it) but jointly extracting one sub-network for all low resource languages (#Mask=5). Besides, we find that layerwise pruning performs slightly better compared with global pruning. We also conduct randomly pruning experiments, which demonstrates the effectiveness of extracting sub-networks with proposed strategies.

Table 3. Analysis of different sub-networks. Models are trained with base structure and prune rate is set to 0.4 throughout the experiments.

Model	#Mask	Type	Strategy	CV-Eval		
				High	Low	Avg
XLSR-10	N/A	N/A	N/A	12.9	15.0	14.3
S3Net	1	Global	LTH	13.0	15.3	14.5
	5	Global	LTH	10.8	15.0	13.6
	10	Global	LTH	10.8	14.0	13.0
	10	Global	Random	14.2	16.9	16.0
	10	Layerwise	TE	12.1	14.6	13.8
	10	Layerwise	LTH	10.6	14.0	12.9

4.5. Analysis of different prune rate

We further investigate the influence of different prune rates. We gradually increase the prune rate p from 0.0 to 0.8 for S3Net-LTH base model, and evaluate the multilingual performance on nine languages. Training XLSR model with additional 50k steps ($p = 0.0$) slightly improved the PER from 14.3% to 14.1%. It is shown that the best choice of p is 0.4 and the performance starts to degrade as the prune rate continues to increase. The figure also shows that we can prune out as many as 70% parameters for each language, while the performance of S3Net-LTH still outperforms baseline XLSR model. Similar phenomenon is also observed in large models.

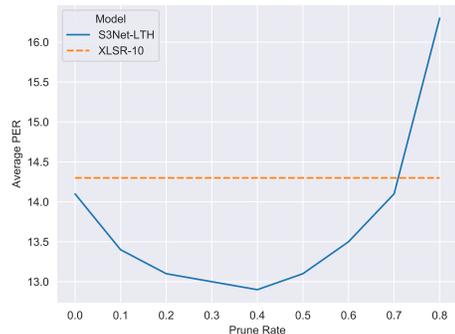


Fig. 2. Evaluation results of different prune rate for S3Net-LTH.

5. CONCLUSION AND FUTURE WORK

In this work, we study language adaptive cross-lingual speech representation learning. We investigate different approaches of extracting sub-networks and show that the proposed S3Net helps alleviating the language interference problem, especially for high resource languages. In the future, we plan to experiment with larger scale multilingual data and the application of multilingual pre-trained models on downstream tasks. We also plan to study structured sparsity and N:M sparsity for network acceleration.

6. REFERENCES

- [1] George Saon, Gakuto Kurata, Tom Sercu, Kartik Audhkhasi, Samuel Thomas, Dimitrios Dimitriadis, Xiaodong Cui, Bhuvana Ramabhadran, Michael Picheny, Lynn-Li Lim, et al., “English conversational telephone speech recognition by humans and machines,” in *Interspeech*, 2017.
- [2] Chung-Cheng Chiu, Tara N Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjali Kannan, Ron J Weiss, Kanishka Rao, Ekaterina Gonina, et al., “State-of-the-art speech recognition with sequence-to-sequence models,” in *ICASSP*, 2018.
- [3] Aaron van den Oord, Yazhe Li, and Oriol Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748*, 2018.
- [4] Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli, “wav2vec: Unsupervised pre-training for speech recognition,” in *Interspeech*, 2019.
- [5] Yu-An Chung, Wei-Ning Hsu, Hao Tang, and James R Glass, “An unsupervised autoregressive model for speech representation learning,” in *Interspeech*, 2019.
- [6] Andy T Liu, Shu-wen Yang, Po-Han Chi, Po-chun Hsu, and Hung-yi Lee, “Mockingjay: Unsupervised speech representation learning with deep bidirectional transformer encoders,” in *ICASSP*, 2020.
- [7] Shaoshi Ling, Yuzong Liu, Julian Salazar, and Katrin Kirchhoff, “Deep contextualized acoustic representations for semi-supervised speech recognition,” in *ICASSP*, 2020.
- [8] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” in *NeurIPS*, 2020.
- [9] Weiran Wang, Qingming Tang, and Karen Livescu, “Unsupervised pre-training of bidirectional speech encoders via masked reconstruction,” in *ICASSP*, 2020.
- [10] 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.
- [11] Chengyi Wang, Yu Wu, Yao Qian, Kenichi Kumatani, Shujie Liu, Furu Wei, Michael Zeng, and Xuedong Huang, “Unispeech: Unified speech representation learning with labeled and unlabeled data,” *arXiv preprint arXiv:2101.07597*, 2021.
- [12] Rich Caruana, “Multitask learning,” *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [13] Yizhou Lu, Mingkun Huang, Hao Li, Jiaqi Guo, and Yanmin Qian, “Bi-encoder transformer network for mandarin-english code-switching speech recognition using mixture of experts,” in *Interspeech*, 2020.
- [14] Vineel Pratap, Anuroop Sriram, Paden Tomasello, Awni Hannun, Vitaliy Liptchinsky, Gabriel Synnaeve, and Ronan Collobert, “Massively multilingual asr: 50 languages, 1 model, 1 billion parameters,” in *Interspeech*, 2020.
- [15] Ngoc-Quan Pham, Tuan-Nam Nguyen, Sebastian Stueker, and Alexander Waibel, “Efficient weight factorization for multilingual speech recognition,” in *Interspeech*, 2021.
- [16] Bo Li, Ruoming Pang, Tara N Sainath, Anmol Gulati, Yu Zhang, James Qin, Parisa Haghani, W Ronny Huang, and Min Ma, “Scaling end-to-end models for large-scale multilingual asr,” *arXiv preprint arXiv:2104.14830*, 2021.
- [17] Xian Li, Changhan Wang, Yun Tang, Chau Tran, Yuqing Tang, Juan Pino, Alexei Baevski, Alexis Conneau, and Michael Auli, “Multilingual speech translation with efficient finetuning of pretrained models,” *arXiv preprint arXiv:2010.12829*, 2020.
- [18] Tianxiang Sun, Yunfan Shao, Xiaonan Li, Pengfei Liu, Hang Yan, Xipeng Qiu, and Xuanjing Huang, “Learning sparse sharing architectures for multiple tasks,” in *AAAI*, 2020.
- [19] Jianze Liang, Chengqi Zhao, Mingxuan Wang, Xipeng Qiu, and Lei Li, “Finding sparse structures for domain specific neural machine translation,” in *AAAI*, 2021.
- [20] Zehui Lin, Liwei Wu, Mingxuan Wang, and Lei Li, “Learning language specific sub-network for multilingual machine translation,” *arXiv preprint arXiv:2105.09259*, 2021.
- [21] Wanying Xie, Yang Feng, Shuhao Gu, and Dong Yu, “Importance-based neuron allocation for multilingual neural machine translation,” in *ACL*, 2021.
- [22] Anjali Kannan, Arindrima Datta, Tara N Sainath, Eugene Weinstein, Bhuvana Ramabhadran, Yonghui Wu, Ankur Bapna, Zhifeng Chen, and Seungji Lee, “Large-scale multilingual speech recognition with a streaming end-to-end model,” in *Interspeech*, 2019.
- [23] Jonathan Frankle and Michael Carbin, “The lottery ticket hypothesis: Finding sparse, trainable neural networks,” in *ICLR*, 2018.
- [24] Pavlo Molchanov, Arun Mallya, Stephen Tyree, Iuri Frosio, and Jan Kautz, “Importance estimation for neural network pruning,” in *CVPR*, 2019.
- [25] Xun Gong, Yizhou Lu, Zhikai Zhou, and Yanmin Qian, “Layer-wise fast adaptation for end-to-end multi-accent speech recognition,” in *Interspeech*, 2021.
- [26] Shubham Toshniwal, Tara N Sainath, Ron J Weiss, Bo Li, Pedro Moreno, Eugene Weinstein, and Kanishka Rao, “Multilingual speech recognition with a single end-to-end model,” in *ICASSP*, 2018.
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017.
- [28] Eric Jang, Shixiang Gu, and Ben Poole, “Categorical reparameterization with gumbel-softmax,” *arXiv preprint arXiv:1611.01144*, 2016.
- [29] Alexis Conneau and Guillaume Lample, “Cross-lingual language model pretraining,” in *NeurIPS*, 2019.
- [30] Rosana Ardila, Megan Branson, Kelly Davis, Michael Henretty, Michael Kohler, Josh Meyer, Reuben Morais, Lindsay Saunders, Francis M Tyers, and Gregor Weber, “Common voice: A massively-multilingual speech corpus,” *arXiv preprint arXiv:1912.06670*, 2019.
- [31] Morgane Riviere, Armand Joulin, Pierre-Emmanuel Mazaré, and Emmanuel Dupoux, “Unsupervised pretraining transfers well across languages,” in *ICASSP*, 2020.
- [32] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *ICML*, 2006.