# LUPET: Incorporating Hierarchical Information Path into Multilingual ASR

Wei Liu\*,1, Jingyong Hou2, Dong Yang2, Muyong Cao2, Tan Lee1

 $^{\rm 1}$  Department of Electronic Engineering, The Chinese University of Hong Kong  $^{\rm 2}$  GVoice, Tencent

louislau\_1129@link.cuhk.edu.hk, {jingyonghou,daviddyang,locwellcao}@tencent.com, tanlee@cuhk.edu.hk

#### **Abstract**

Toward high-performance multilingual automatic speech recognition (ASR), various types of linguistic information and model design have demonstrated their effectiveness independently. They include language identity (LID), phoneme information, language-specific processing modules, and crosslingual self-supervised speech representation. It is expected that leveraging their benefits synergistically in a unified solution would further improve the overall system performance. This paper presents a novel design of a hierarchical information path, named LUPET, which sequentially encodes, from the shallow layers to deep layers, multiple aspects of linguistic and acoustic information at diverse granularity scales. The path starts from LID prediction, followed by acoustic unit discovery, phoneme sharing, and finally token recognition routed by a mixture-ofexpert. ASR experiments are carried out on 10 languages in the Common Voice corpus. The results demonstrate the superior performance of LUPET as compared to the baseline systems. Most importantly, LUPET effectively mitigates the issue of performance compromise of high-resource languages with low-resource ones in the multilingual setting.

**Index Terms**: Multilingual ASR, language identity, self-supervised speech representation learning, mixture-of-expert

# 1. Introduction

Conventionally an automatic speech recognition (ASR) system is developed to transcribe speech into text for a specific language. Toward multilingual ASR, recent research is focused on building a unified model that covers multiple languages [1-4]. One practical advantage is the reduction of training and deployment costs, as compared to building a separate monolingual model for each language. It also facilitates sharing of linguistic knowledge among the languages and may help elevate the recognition performance on those with limited data resources [5]. It was shown that training a fully shared end-toend (E2E) model on a multilingual speech corpus is a simple yet effective solution (vanilla) [4]. The multilingual corpus is made by mixing the corpora of different languages, and a shared vocabulary is used. Due to the heterogeneous nature of different languages [4, 6], the vanilla scheme exhibits the issue that the recognition performance on the high-resource languages inevitably compromises in the multilingual training, in order to attain reasonable performance on the low-resource languages.

To mitigate the issue of performance compromise, there have been attempts that incorporate language identity (LID) information [7–10], phoneme information [11–13], and language-specific architecture, e.g., mixture-of-expert (MoE) [14–18].

These approaches typically require a supervised training process, which requires labeled training data. On the other hand, self-supervised learning (SSL) is believed to be an effective way of cross-lingual data sharing. The representative works include wav2vec2.0 [19], HuBERT [20], XLSR [21], and many others. The two-stage training scheme, i.e., pre-training and finetuning, has been widely applied in SSL. In [22], joint unsupervised and supervised training (JUST) was shown to outperform two-stage training on multilingual ASR. JUST uses a contrastive loss and a masked language model (MLM) loss to learn discrete units for better contextualized representations.

In the present study, a hierarchical information path is developed to combine multiple useful factors synergistically to boost the overall performance of a multilingual ASR system. The path comprises a sequence of prediction modules that incorporate linguistic and acoustic information at diverse granularity levels into the recognition process. These modules are namely, LID, Acoustic Unit discovery, Phoneme sharing, and mixture of Experts for Token recognition. We use the acronym LUPET to denote the proposed design, in which each alphabet represents one of the information components in the path. The path LUPET can be easily integrated into a vanilla ASR architecture by unfolding with the encoder layers. Within this path, information from a shallow layer can benefit those that occur in deeper ones, and hence it is considered a hierarchical flow. Here the shallow layer refers to an encoder layer close to the input. Importantly, the required information labels are either straightforward to obtain or can be derived via an SSL process.

The effectiveness of LUPET is evaluated by experiments on 10 languages in the Common Voice [23]. The results show that, compared to the vanilla system, LUPET can achieve 19.7% and 12.3% relative reduction of average word error rate with CTC and attention decoding, respectively. LUPET also outperforms previous baseline systems. In particular, it demonstrates superior performance on the high-resource languages as its performance compromise to low-resource languages is alleviated.

# 2. LUPET

#### 2.1. Vanilla E2E Multilingual ASR

The E2E multilingual ASR architecture we adopt is the hybrid CTC-Attention conformer [24,25]. It consists of three components, namely encoder, decoder, and CTC [26] layer. The encoder takes an acoustic feature sequence  $\mathbf{X} = \{\mathbf{x}_t\}_{t=1}^T$  as input and converts it to hidden representation  $\mathbf{H} = \{\mathbf{h}_t\}_{t=1}^{T'}$ , where T and T' denote the number of original frames and the subsampled frames. The  $\mathbf{H}$  is then forwarded to two classification branches for predicting the token sequence  $\mathbf{Y} = \{y_u \in \mathcal{V}\}_{u=1}^U$ , where  $\mathcal{V}$  is a shared multilingual vocabulary built by BPE [27]

<sup>\*</sup> This work was done during an internship at Tencent.

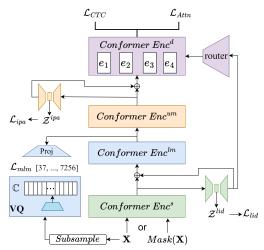


Figure 1: The overall architecture of our proposed LUPET multilingual ASR. LUPET information path unfolds with the encoder layers.  $\{Enc^s, Enc^{lm}, Enc^{um}, Enc^d\}$  represent shallow, lower-middle, upper-middle, deep layers, respectively.  $Enc^s$  and  $Enc^{um}$  are used for LID and IPA phoneme prediction.  $Enc^{lm}$  performs acoustic unit discovery with a random-projection quantizer, where  $\mathbb C$  denotes the codebook for vector quantization (VQ).  $Enc^d$  denotes conformer layers modified with MoE which consists of four experts and a router. All trapezoid modules refer to linear projection.

and U denotes the number of tokens. One classification branch, i.e., decoder, conditions  $\mathbf H$  to autoregressively compute tokenlevel posterior  $p(y_u|\mathbf H,y_{1:u-1})$  via a cross attention mechanism. The attention loss is given as

$$\mathcal{L}_{Attn} = -\sum_{u=1}^{U} log p(y_u | \mathbf{H}, y_{1:u-1}). \tag{1}$$

Another classification branch, i.e., the CTC layer, simultaneously derives the frame-level posteriors  $p(\mathbf{z}_t|\mathbf{h}_t)$ . The CTC loss is formulated as follows:

$$\mathcal{L}_{CTC} = CTC(\mathbf{Z}, \mathbf{Y}) = -\sum_{\mathbf{Z} \in B^{-1}(\mathbf{Y})} \sum_{t=1}^{T'} logp(\mathbf{z}_t | \mathbf{h}_t), \quad (2)$$

where  $\mathbf{z}_t$  is the logits over  $\mathcal{V} \cup \emptyset$  and  $B^{-1}$  is the inverse function that gives all valid alignment paths between input sequence  $\mathbf{H}$  and output sequence  $\mathbf{Y}$ . The blank token  $\emptyset$  is specially designed by CTC for aligning  $\mathbf{H}$  and  $\mathbf{Y}$ .  $\mathbf{Z} = \{\mathbf{z}_t\}_{t=1}^{T'}$  represents one possible alignment path. The training objective of hybrid CTC-Attention is a linear combination of Eq. 1 and Eq. 2:

$$\mathcal{L}_{CTC-Attn} = (1 - \lambda)\mathcal{L}_{Attn} + \lambda\mathcal{L}_{CTC}, \tag{3}$$

where  $\lambda$  is a coefficient to control the weight of CTC loss.

# 2.2. Incorporating LUPET

Recent studies have pointed out (1) Both LID and phoneme information are beneficial for multilingual training [10, 13]; (2) Design language-specific modules to process language-specific information is useful to reduce language interference [16, 18]; (3) The success of self-supervised cross-lingual representation learning applied in ASR [21, 28]. Although these factors' effectiveness has been verified separately, how to synergistically combine them to contribute a better solution from a unified perspective remains an open question. The proposed LUPET provides a novel view from the multilingual hierarchical information path. From LID to acoustic unit followed by phoneme then

go through MoE routing to the final token, the information that occurred in the early position of the path is assumed to contribute to the prediction of later information of the path.

As shown in Fig. 1, the full encoder is composed of  $\{Enc^s, Enc^{lm}, Enc^{um}, Enc^d\}$ , from shallow layers to deep layers. LUPET information path unfolds with the encoder layers. The shallow layers of encoder  $Enc^s$  are used to identify the spoken language. Denote the output of  $Enc^s$  as shallow representations  $\mathbf{H}^s$ .  $\mathbf{H}^s$  is then projected to the LID logits  $\mathbf{Z}^{lid}$  via a linear transformation. The logits dimension  $dim(\mathbf{Z}^{lid}) = \#LID + 1$ , where 1 represents a special blank token for CTC as mentioned in Sec. 2.1. The LID prediction loss is formulated as:

$$\mathcal{L}_{lid} = CTC(\mathbf{Z}^{lid}, LID_{seq}), \tag{4}$$

where the sequential LID labels  $LID_{seq}$  are constructed by repeating the single LID label to the number of output tokens.

The predicted LID information then is propagated to subsequent layers of the encoder via self-conditioning

$$\mathbf{H}^{s'} = \mathbf{H}^s + LIN(\mathbf{Z}^{lid}),\tag{5}$$

where LIN denotes a linear layer to keep the hidden dimension and  $\mathbf{H}^{s'}$  is the input representations of  $Enc^{lm}$ .

The lower-middle layers of encoder  $Enc^{lm}$  are utilized to perform acoustic unit discovery. Similar to BEST-RQ [28], a random-projection quantizer including a projection matrix  $Proj^c$  and codebook  $\mathbb C$  is applied and none of the parameters are trainable. Vector quantization (VQ) is carried out on acoustic features  $\mathbf X$  to produce discrete labels

$$\mathbf{Lab}_{u} = \mathbb{C}(Proj^{c}(Sub(\mathbf{X})))), \tag{6}$$

where Sub represents the subsample operation and  $Proj^c$  performs projection from the speech feature dimension to the code vector dimension. The output of  $\mathbb C$  are the indices of the nearest code vectors of the codebook to the input vectors.

With probability p,  $\mathbf{X}$  is randomly masked to feed the encoder. Masked language modeling (MLM) is then performed to predict  $\mathbf{Lab}_u$ . Denote  $\mathbf{H}_M^{lm}$  as the output representation of  $Enc^{lm}$ , where the under-script M means having Mask( $\mathbf{X}$ ) as input. Let  $\mathbf{mi}$  be the masked indices on  $\mathbf{H}_M^{lm}$ , the MLM loss can be written as:

$$\mathcal{L}_{mlm} = CE(Proj^{u}(\mathbf{H}_{M}^{lm}[\mathbf{mi}]), \mathbf{Lab}_{u}[\mathbf{mi}]),$$
 (7)

where  $Proj^u$  projects the hidden dimension to the size of the codebook and MLM loss is the cross-entropy (CE) between logits over the codebook and labels at the masked positions.

Discrete acoustic units discovered by  $Enc^{lm}$  are expected to facilitate pronunciation learning to incorporate phonetic information for subsequent layers. Similar to LID prediction, the output representation  $\mathbf{H}^{um}$  of upper-middle encoder  $Enc^{um}$  is used to predict phoneme sequence IPA.  $\mathbf{Z}^{ipa}$ , projected by  $\mathbf{H}^{um}$ , is the logits over IPA phonemes and an additional blank token. Eq. 8 gives the loss of IPA prediction:

$$\mathcal{L}_{ipa} = CTC(\mathbf{Z}^{ipa}, IPA). \tag{8}$$

Following Eq. 5, self-conditioning is similarly applied to obtain  $\mathbf{H}^{um'}$  to propagate the predicted phonetic information.

$$\mathbf{H}^{um'} = \mathbf{H}^{um} + LIN(\mathbf{Z}^{ipa}) \tag{9}$$

Lastly, the deep layers of encoder  $Enc^d$  are modified with MoE. Multiple FFN experts and a routing network are included in

Table 1: Training and testing hours of 10 languages from the Common Voice 13.0 corpus in our experiments.

LID	en	fr	es	zh	it
Train	2279.98	872.19	448.45	359.91	286.61
Test	26.90	26.08	26.65	30.44	26.27
LID	ru	pt	tr	nl	tt
Train	178.78	125.35	69.08	73.82	19.95
Test	15.73	11.91	12.05	14.58	5.70

the MoE structure. The language self-condition representation  $(LIN(\mathbf{Z}^{lid}))$  in Eq. 5) is regarded as the LID embedding to feed the routing network. The output of the routing network is a softmax distribution over the number of experts. Followed [16], the top-2 experts with the highest probabilities are dynamically routed to process each frame based on the frame-level LID information from shallow encoder layers.

To incorporate the hierarchical information, from LID, acoustic unit, phoneme, and token, the objective function of LUPET is given as a linear combination of Eq. (3, 4, 7, 8):

$$\mathcal{L}_{LUPET} = \mathcal{L}_{CTC-Attn} + w_1 \mathcal{L}_{lid} + w_2 \mathcal{L}_{mlm} + w_3 \mathcal{L}_{ipa},$$
(10)

where  $w_1, w_2, w_3$  are the weights of the corresponding losses.

### 3. Experimental Setup

#### 3.1. Dataset

The 10 languages, namely English (en), French (fr), Spanish (es), Chinese (zh), Italian (it), Russian (ru), Portuguese (pt), Turkish (tr), Dutch (nl) and Tatar (tt) from the public available Common Voice 13.0 [23] are selected for our multilingual ASR experiments. The language coverage includes high-resource languages, e.g., English with around 2,280 hours of training data, and low-resource languages, e.g., Tatar with only about 20 hours. The detailed training and testing statistics are listed in Tab. 1. Note that zh includes Mandarin, Taiwanese, and Cantonese. A standard text normalization (the same as in Whisper [2]) is applied to all transcriptions of the dataset.

### 3.2. Multilingual ASR Configurations

# 3.2.1. Vanilla

The vanilla model adopts a hybrid CTC-Attention architecture. The encoder has 12 conformer layers with 8 attention heads and 512 hidden dimensions, while the decoder has 6 transformer layers [29]. The CTC weight  $\lambda$  in Eq. 3 is set to 0.3. The input acoustic feature to the network is the typical 80-dimensional log-Mel filterbank. The output vocabulary used is derived from Whisper's tokenizer. This tokenizer was obtained by BPE using UTF-8 bytes of the entire training dataset of Whisper.

# 3.2.2. LUPET

Compared to vanilla, the encoder architecture has several modifications by incorporating LUPET. The output positions of  $\{Enc^s, Enc^{lm}, Enc^{um}, Enc^d\}$  are at the  $\{3\text{-th}, 6\text{-th}, 9\text{-th}, 12\text{-th}\}$  layer of the original encoder, respectively. When MLM loss takes effect, acoustic feature  $\mathbf{X}$  is randomly masked consecutive 20 frames with probability p=0.01. Codebook  $\mathbb C$  of the random-projection quantizer has a size of 8192 and a dimension of 16. IPA sequence per-utterance is obtained using an open-sourced toolkit phonemizer [30]. In each layer of  $Enc^d$ , MoE including 8 FFN experts is used to replace the end-FFN of

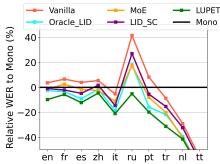


Figure 2: Relative WER changes of different systems to monolingual systems on 10 languages by CTC greedy decoding.

the original conformer layer. In Eq. 10, the weight coefficients  $w_1$ ,  $w_2$ , and  $w_3$  are set to 0.3, 0.07, and 0.3, respectively.

#### 3.2.3. Baselines

Several baselines are used for comparison: (1) Mono, monolingual ASR with vanilla architecture and 256 hidden dimensions is trained per language. (2)  $Oracle\_LID$ , append the pre-known LID embedding to input acoustic feature for multilingual ASR training. (3) MoE [16], keep the  $Enc^d$  of LUPET and remove other auxiliary losses, hidden representation  $\mathbf{H}^{um}$  is used as the input to the routing network. (4)  $LID\_SC$  [10], LID prediction by CTC and LID information self-conditioning (SC) are performed over the vanilla model. (5) Whisper, whisper-large-v2 with oracle LID is used for decoding.

#### 3.3. Training Scheme and Evaluation Metric

We implement the vanilla and our proposed LUPET methods on the Wenet toolkit [31]. The model is trained with Adam [32] optimizer with a learning rate (LR) of 1e-3. LR schedule has a warmup step of 15000. Batch size is set to 12 with  $accum\_grad=16$ . 8 V100 GPUs are used for DDP training. Each multilingual model is trained for 50 epochs and each monolingual model is trained for 100 epochs. If not specified otherwise, MLM takes effect from epoch 5 to 30. The final model for decoding is obtained by averaging the 10 best models with the lowest validation losses. Character error rate (CER) for Chinese and word error rate (WER) for other languages are used to measure the system's performance.

### 4. Results and Analysis

# 4.1. Performance Comparison to Monolingual System

A desirable multilingual ASR is expected to present better performance than its monolingual counterparts (Mono). As shown in Fig. 2, LUPET and other baselines are used to compare performance to monolingual systems on the test sets of 10 languages. The x-axis follows an order from high-resource to low-resource languages. The y-axis denotes the relative WER to Mono, where the more negative value represents the lower WER. It is clear to see all curves are basically decreasing except the peak at ru. The decreasing trend is straightforward as the more low-resourced languages can achieve more significant performance gains. Multilingual training brings degradation to the Russian (ru) language in most cases. We speculate this may be due to the compromised phenomenon from Russian (ru) to Tarta (tt). The language tt achieves above 60% relative WER reduction via multilingual training at the cost of side effects on language ru, which shares the same language family as tt.

It is worth noted that the recognition performance of *Vanilla* system on four high-resource languages (en, fr, es, zh) can-

Table 2: WER (%) results of different systems on 10 languages of Common Voice by CTC greedy decoding. avg 5high denotes the averaged WER results of top-5 high-resourced languages and avg 5low is similar for the low-resourced case. In the LUPET block, the backslash / represents the ablation study by removing the following component, where U = acoustic unit discovery, P = IPA sharing, and L = LID prediction.  $w_2 = 1$  means the weight coefficient of  $\mathcal{L}_{mlm}$  is set as 1. Uto50ep means U takes effect until 50 epochs.

Model	en	fr	es	zh	it	ru	pt	tr	nl	tt	avg	avg w/o tt	avg 5high	avg 5low
Mono	13.03	12.51	9.37	13.01	11.15	11.55	11.16	25.73	19.34	83.62	21.05	14.09	11.81	30.28
Vanilla	13.50	13.33	9.74	13.71	10.56	16.34	12.07	23.49	13.72	36.78	16.32	14.05	12.17	20.48
Oracle_LID	12.69	12.08	8.51	12.8	9.07	13.64	9.39	20.12	11.72	30.29	14.03	12.22	11.03	17.03
LID_SC	12.94	12.24	8.84	13.46	9.36	14.72	10.89	22.53	12.74	34.19	15.19	13.08	11.37	19.01
MoE	12.86	12.81	9.23	12.67	9.91	13.56	10.26	20.55	11.46	30.47	14.38	12.59	11.50	17.26
LUPET	11.75	11.79	8.22	12.41	8.81	10.95	8.95	17.71	11.32	29.12	13.10	11.32	10.60	15.61
LUPET / U	12.33	12.32	8.73	12.42	9.45	10.58	9.62	18.00	10.86	27.76	13.21	11.59	11.05	15.36
LUPET / P	12.35	12.22	8.67	12.31	9.39	11.92	10.51	21.92	12.08	27.23	13.86	12.37	10.99	16.73
LUPET / UP	12.71	12.38	8.82	12.16	9.54	11.9	10.34	20.89	11.72	26.47	13.69	12.27	11.12	16.26
LUPET / LU	11.96	12.02	8.46	12.21	9.08	10.99	9.44	19.49	11.03	31.90	13.66	11.63	10.75	16.57
LUPET $w_2 = 1$	11.80	11.86	8.54	12.33	9.10	12.38	10.25	17.85	10.30	33.11	13.75	11.60	10.73	16.78
LUPET Uto50ep	11.72	12.09	8.40	12.73	9.02	11.90	9.98	19.06	11.96	31.35	13.82	11.87	10.79	16.85

not surpass the corresponding monolingual system, demonstrating the general compromised phenomenon towards low-resource languages in multilingual training. *Oracle\_LID* system brings consistent improvements over *Vanilla* across all languages, which proves the benefits of LID information. Both *MoE* and *LID\_SC* also show obviously better performance than *Vanilla*, while being inferior to *Oracle\_LID*. *LUPET* outperforms all other baselines, serving as the only system that gives WER reduction on all languages compared to *Mono*. The advantage of *LUPET* is highlighted by the superior performance on high-resource languages, which largely mitigates the compromised phenomenon during multilingual training.

#### 4.2. LUPET's Effectiveness Verification

Tab. 2 presents the WER results of different systems on 10 languages. Averaged WER are calculated for overall comparison. The top-5 languages, with more training data, are roughly referred to as high-resources (5high), while the remaining 5 languages are low-resources (5low). Having similar observation from Fig. 2, LUPET gives significantly better performance on all languages compared to other baselines. The benefits of LID prediction and MoE routing structure have been well verified by system LID\_SC and MoE.

To further illustrate the effectiveness of LUPET, several ablation studies are carried out to investigate the remained components. By removing U (acoustic unit discovery) from LUPET, it can be clearly observed that WERs on high-resource languages consistently increase. Contrary to high-resource, some low-resource languages especially for tt, achieve somewhat improvements. It demonstrates the quality of discrete units that discovered by MLM is related to the amount of data. Hence, MLM can usually bring positive gains to high-resource languages. With a longer MLM effective period (Uto50ep), low-resource languages have obvious performance degradation. The gains towards high-resource gradually converge to the language en. When increasing the weight coefficient of  $\mathcal{L}_{mlm}$  ( $w_2 = 1$ ), inferior results compared to the original LUPET setting are presented for most of languages.

Disabling both U and P (IPA sharing prediction) exhibits worse results. Tab. 2 provides two comparison views to understand the independent effect of the component P. (1) LU-PET/LU can be seen as "MoE + P". Comparing it with MoE, IPA sharing is found to be beneficial for all languages except tt. (2) When only removing P from LUPET, WERs on all languages basically degrade. Low-resource languages clearly give the worse results, where tr increase the absolute WER above 4%. One possible reason is that the independent U would lead

to information loss due to the masking mechanism in MLM, especially for low-resource languages. In LUPET, with the help of intermediate phoneme prediction (*P*), the loss of information is largely mitigated, thus not presenting much worse results on low-resource languages.

Table 3: Averaged WER (%) of different systems by attention decoding. Note that Whisper-large-v2 decodes in a greedy manner, while other systems utilize beam search with beam\_size=20.

Model	avg	avg w/o tt	avg 5high	avg 5low
Whisper	20.86	11.52	13.21	28.52
Mono	17.89	10.02	9.52	26.26
Vanilla	10.43	8.78	8.88	11.99
Oracle_LID	9.20	7.89	8.31	10.08
LID_SC	10.22	8.62	8.62	11.83
MoE	9.61	8.25	8.46	10.76
LUPET	9.15	7.77	7.95	10.35

#### 4.3. Results on Attention Decoding

Tab. 3 presents the averaged WER metrics of different systems by attention decoding. Not surprisingly, LUPET exhibits the overall best performance, especially on high-resource languages, and significantly outperforms its CTC decoding counterpart by 3.94% absolute WER. Whisper is introduced as an external reference. As can be seen, the zero-shot performance of Whisper is easily surpassed even by Mono, illustrating the importance of in-domain training. Furthermore, it is noted that the overall performance gaps between systems in attention decoding are far less than that of CTC decoding, e.g., comparing Oracle\_LID and LUPET. We hypothesize that the attention decoder served as a language model may help overfit the pattern in the specific domain. This also explains why attention decoding clearly outperforms CTC decoding in our experiments.

### 5. Conclusions

This paper presents a novel view to seamlessly incorporate hierarchical information into multilingual ASR. Multiple information in different granularity, i.e., LID, acoustic unit, phoneme, and token, form a path LUPET that unfolds with encoder layers. Experiments carried out on 10 languages of Common Voice corpus illustrate the effectiveness of LUPET, even outperforming the system with oracle LID information. Different components in LUPET are proved to be useful in ablation studies. It is found that the acoustic unit discovery and phoneme prediction significantly help the recognition on high-resource languages, largely mitigating the compromised phenomenon.

#### 6. References

- [1] B. Li, R. Pang, T. N. Sainath, A. Gulati, Y. Zhang, J. Qin, P. Haghani, W. R. Huang, M. Ma, and J. Bai, "Scaling end-to-end models for large-scale multilingual asr," in *Proc. ASRU*. IEEE, 2021, pp. 1011–1018.
- [2] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *Proc. ICML*. PMLR, 2023, pp. 28 492–28 518.
- [3] Y. Zhang, W. Han, J. Qin, Y. Wang, A. Bapna, Z. Chen, N. Chen, B. Li, V. Axelrod, G. Wang et al., "Google USM: Scaling automatic speech recognition beyond 100 languages," arXiv preprint arXiv:2303.01037, 2023.
- [4] V. Pratap, A. Sriram, P. Tomasello, A. Y. Hannun, V. Liptchinsky, G. Synnaeve, and R. Collobert, "Massively multilingual ASR: 50 languages, 1 model, 1 billion parameters," in *Proc. Interspeech*. ISCA, 2020, pp. 4751–4755.
- [5] H. Yadav and S. Sitaram, "A survey of multilingual models for automatic speech recognition," in *Proc. LREC*. European Language Resources Association, 2022, pp. 5071–5079.
- [6] B. Li, R. Pang, Y. Zhang, T. N. Sainath, T. Strohman, P. Haghani, Y. Zhu, B. Farris, N. Gaur, and M. Prasad, "Massively multilingual asr: A lifelong learning solution," in *Proc. ICASSP*. IEEE, 2022, pp. 6397–6401.
- [7] S. Watanabe, T. Hori, and J. R. Hershey, "Language independent end-to-end architecture for joint language identification and speech recognition," in *Proc. ASRU*. IEEE, 2017, pp. 265–271.
- [8] C. Zhang, B. Li, T. N. Sainath, T. Strohman, S. Mavandadi, S. Chang, and P. Haghani, "Streaming end-to-end multilingual speech recognition with joint language identification," in *Proc. Interspeech*. ISCA, 2022, pp. 3223–3227.
- [9] L. Zhou, J. Li, E. Sun, and S. Liu, "A configurable multilingual model is all you need to recognize all languages," in *Proc. ICASSP*. IEEE, 2022, pp. 6422–6426.
- [10] W. Chen, B. Yan, J. Shi, Y. Peng, S. Maiti, and S. Watanabe, "Improving massively multilingual asr with auxiliary CTC objectives," in *Proc. ICASSP*. IEEE, 2023, pp. 1–5.
- [11] H. B. Sailor and T. Hain, "Multilingual speech recognition using language-specific phoneme recognition as auxiliary task for indian languages." in *Proc. Interspeech*, 2020, pp. 4756–4760.
- [12] C. Zhu, K. An, H. Zheng, and Z. Ou, "Multilingual and crosslingual speech recognition using phonological-vector based phone embeddings," in *Proc. ASRU*. IEEE, 2021, pp. 1034–1041.
- [13] C. Wang, Y. Wu, Y. Qian, K. Kumatani, S. Liu, F. Wei, M. Zeng, and X. Huang, "Unispeech: Unified speech representation learning with labeled and unlabeled data," in *Proc. ICML*. PMLR, 2021, pp. 10 937–10 947.
- [14] N. Gaur, B. Farris, P. Haghani, I. Leal, P. J. Moreno, M. Prasad, B. Ramabhadran, and Y. Zhu, "Mixture of informed experts for multilingual speech recognition," in *Proc. ICASSP*. IEEE, 2021, pp. 6234–6238.
- [15] Z. You, S. Feng, D. Su, and D. Yu, "Speechmoe2: Mixture-ofexperts model with improved routing," in *Proc. ICASSP*. IEEE, 2022, pp. 7217–7221.
- [16] K. Hu, B. Li, T. N. Sainath, Y. Zhang, and F. Beaufays, "Mixture-of-expert conformer for streaming multilingual asr," arXiv preprint arXiv:2305.15663, 2023.
- [17] W. Wang, G. Ma, Y. Li, and B. Du, "Language-routing mixture of experts for multilingual and code-switching speech recognition," arXiv preprint arXiv:2307.05956, 2023.
- [18] E. Sun, J. Li, Y. Hu, Y. Zhu, L. Zhou, J. Xue, P. Wang, L. Liu, S. Liu, E. Lin, and Y. Gong, "Building high-accuracy multilingual ASR with gated language experts and curriculum training," in *Proc. ASRU*. IEEE, 2023, pp. 1–7.
- [19] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proc. NeurIPS*, 2020, pp. 12449–12460.

- [20] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhut-dinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [21] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, "Unsupervised cross-lingual representation learning for speech recognition," in *Proc. Interspeech*. ISCA, 2021, pp. 2426–2430.
- [22] J. Bai, B. Li, Y. Zhang, A. Bapna, N. Siddhartha, K. C. Sim, and T. N. Sainath, "Joint unsupervised and supervised training for multilingual asr," in *Proc. ICASSP*. IEEE, 2022, pp. 6402–6406.
- [23] R. Ardila, M. Branson, K. Davis, M. Kohler, J. Meyer, M. Henretty, R. Morais, L. Saunders, F. M. Tyers, and G. Weber, "Common voice: A massively-multilingual speech corpus," in *Proc. LREC*. European Language Resources Association, 2020, pp. 4218–4222.
- [24] S. Watanabe, T. Hori, S. Kim, J. R. Hershey, and T. Hayashi, "Hybrid CTC/attention architecture for end-to-end speech recognition," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1240–1253, 2017.
- [25] A. Gulati, J. Qin, C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, and R. Pang, "Conformer: Convolution-augmented transformer for speech recognition," in *Proc. Interspeech*. ISCA, 2020, pp. 5036–5040.
- [26] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proc. ICML*, 2006, pp. 369–376.
- [27] V. Zouhar, C. Meister, J. L. Gastaldi, L. Du, T. Vieira, M. Sachan, and R. Cotterell, "A formal perspective on byte-pair encoding," in *Proc. ACL*. Association for Computational Linguistics, 2023, pp. 598–614.
- [28] C.-C. Chiu, J. Qin, Y. Zhang, J. Yu, and Y. Wu, "Self-supervised learning with random-projection quantizer for speech recognition," in *Proc. ICML*. PMLR, 2022, pp. 3915–3924.
- [29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. NeurIPS*, 2017, pp. 5998–6008.
- [30] M. Bernard and H. Titeux, "Phonemizer: Text to phones transcription for multiple languages in python," *Journal of Open Source Software*, vol. 6, no. 68, p. 3958, 2021. [Online]. Available: https://doi.org/10.21105/joss.03958
- [31] Z. Yao, D. Wu, X. Wang, B. Zhang, F. Yu, C. Yang, Z. Peng, X. Chen, L. Xie, and X. Lei, "Wenet: Production oriented streaming and non-streaming end-to-end speech recognition toolkit," in *Proc. Interspeech.* ISCA, 2021, pp. 4054–4058.
- [32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015.