

Exploiting Adapters for Cross-lingual Low-resource Speech Recognition

Wenxin Hou, Han Zhu, Yidong Wang, Jindong Wang, Tao Qin, Renjun Xu, and Takahiro Shinozaki

Abstract—Cross-lingual speech adaptation aims to solve the problem of leveraging multiple rich-resource languages to build models for a low-resource target language. Since the low-resource language has limited training data, speech recognition models can easily overfit. Adapter is a versatile module that can be plugged into Transformer for parameter-efficient learning. In this paper, we propose to use adapters for parameter-efficient cross-lingual speech adaptation. Based on our previous MetaAdapter that implicitly leverages adapters, we propose a novel algorithm called SimAdapter for explicitly learning knowledge from adapters. Our algorithms can be easily integrated into the Transformer structure. MetaAdapter leverages meta-learning to transfer the general knowledge from training data to the test language. SimAdapter aims to learn the similarities between the source and target languages during fine-tuning using the adapters. We conduct extensive experiments on five-low-resource languages in the Common Voice dataset. Results demonstrate that MetaAdapter and SimAdapter can reduce WER by 2.98% and 2.55% with only 2.5% and 15.5% of trainable parameters compared to the strong full-model fine-tuning baseline. Moreover, we show that these two novel algorithms can be integrated for better performance with up to 3.55% relative WER reduction.

Index Terms—speech recognition, cross-lingual adaptation, meta-learning, parameter-efficiency

I. INTRODUCTION

AUTOMATIC speech recognition (ASR) based on end-to-end (E2E) models has made remarkable progress by training on large-scale data [1], [2]. We can use a single E2E ASR system for a large number of languages [3], [4] without complicated language-specific processing. Nevertheless, a well-known limitation of E2E ASR methods is that they require a considerable amount of training data to achieve superior performances among various domains [5] and languages [6]. Therefore, it remains a challenge for E2E ASR models to achieve good performance for most of the low-resource languages among about 7,000 languages in the world.

Some research has indicated that the performances of low-resource languages benefit from transferring the common knowledge from rich-resource languages in ASR [7]. For instance, as shown in Fig. 1, given Romanian as a low-resource target language, cross-lingual ASR aims to build

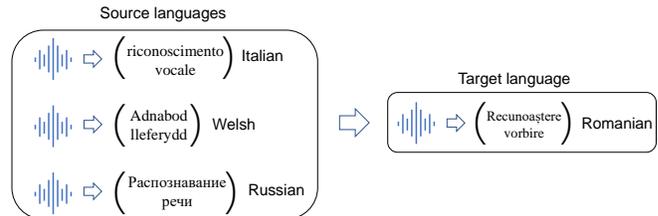


Fig. 1. Illustration of the cross-lingual speech recognition task. Given three rich-resource languages as source languages (Italian, Welsh, and Russian), how to learn the transferable knowledge from them to build cross-lingual ASR models for the target language Romanian?

models by leveraging the available rich-resource languages such as Italian, Welsh, and Russian as source languages. Knowledge transfer can be achieved in three avenues: (1) pre-training on the rich-resource languages and then fine-tuning on the low-resource tasks [3], [8]; (2) performing multi-task training using all languages [9]; and (3) learning the general common knowledge and then rapidly adapting to the low-resource languages using meta-learning [10]. A possible explanation is that different languages intrinsically share some beneficial information like speaker, environment and linguistic information. In this paper, we mainly focus on the fine-tuning methods.

Due to the limited training data in low-resource languages, direct re-training makes the model easily overfit. These problems make the transfer-based methods inefficient [11], [12]. Recently, the adapter module was proposed for parameter-efficient fine-tuning in multilingual or cross-lingual settings [11]–[13], which can mitigate overfitting. Adapter is an add-on module to the encoder and decoder layers in Transformer mainly composed of layer normalization and fully-connected layers. During fine-tuning, we can freeze the backbones of the pre-trained models and only train the adapters which have a small number of task-specific parameters. Pfeiffer et al. [14] studied the fusion of adapters in natural language processing where they linearly combine the outputs of multiple adapters for target classification task adaptation. However, it remains unexplored to investigate the performance of multiple adapters on cross-lingual ASR tasks.

In our previous work [11], we proposed *MetaAdapter* to learn general and transferable speech representations using model-agnostic meta-learning (MAML) [15] and achieved promising results on extremely low-resource languages. However, it is unclear whether MetaAdapter can handle the non-extreme cases where there are moderate training data. Moreover, MetaAdapter relies on meta-learning to implicitly learn

W. Hou, Y. Wang, and T. Shinozaki are with Tokyo Institute of Technology, Tokyo, Japan. Email: {hou.w.aa, wang.y.ca}@m.titech.ac.jp, shino@ict.e.titech.ac.jp.

W. Hou is also with Microsoft, work done at Tokyo Institute of Technology and Microsoft Research Asia. Email: wenxinhou@microsoft.com

H. Zhu is with Institute of Acoustics, Chinese Academy of Sciences, China. Email: zhuhan@hcl.ioa.ac.cn.

J. Wang and T. Qin are with Microsoft Research Asia. Email: {jindong.wang, taoqin}@microsoft.com.

R. Xu is with Zhejiang University. Email: rux@zju.edu.cn.

Correspondence to Jindong Wang and Takahiro Shinozaki.

from source languages, which makes no assumptions on the relationship between source and target languages that may weaken its interpretability. Therefore, in this paper, we comprehensively investigate the potential of leveraging multiple source adapters in cross-lingual speech recognition. Based on our analysis on MetaAdapter, we propose a novel algorithm: *SimAdapter*, to learn the similarity between the source and target languages using the attention mechanism. Our key motivation is that different languages in the world are sharing similarities based on their similar geological characteristics or evolution [16]–[18]. Therefore, it is feasible to explicitly model such similarities in the ASR models.

Both of the two algorithms are parameter-efficient and thus can prevent the overfitting problem. To our best knowledge, there is no existing research that tries to model the cross-lingual ASR tasks by studying their relationship using meta-learning and transfer learning-based adapters. In addition, the MetaAdapter and SimAdapter are compatible, thus can be integrated for better performance.

Our contributions can be summarized as follows:

- We comprehensively analyze our previously proposed MetaAdapter and propose a novel algorithm for cross-lingual low-resource ASR: SimAdapter.
- Experiments on five low-resource languages demonstrated a relative WER improvement of 2.98% with MetaAdapter and 2.55% with SimAdapter using only 2.5% and 15.5% trainable parameters compared with the strong full-model fine-tuning baseline, respectively.
- These two algorithms can be integrated to achieve better performance with up to 3.55% relative improvement.

This paper is substantially an extended version of our previously published paper [11] at ICASSP 2021. Compared to the previous version, we make heavy extensions as follows: (1) We propose a parallel new algorithm called SimAdapter for cross-lingual ASR. (2) We investigate the difference and integration between the MetaAdapter and SimAdapter algorithms. (3) We conduct extensive experiments on cross-lingual ASR datasets to validate the effectiveness of these algorithms.

The structure of this paper is as follows. In Section II, we review the related work to multilingual, cross-lingual ASR and adapters. Section III introduces our main ideas. Section IV and Section V introduce the details of MetaAdapter and SimAdapter algorithms, and their integration. Section VI presents experimental design details and Section VII reports our experimental results and analysis. Finally, in Section VIII, we conclude this paper and present some future work.

II. RELATED WORKS

A. Multilingual and Cross-lingual Speech Recognition

Multilingual E2E ASR is getting increasing attention over the years to handle multiple languages with a single model. Watanabe et al. [19] proposed a language-independent architecture based on hybrid CTC-attention structure [20] with augmented vocabulary for character-based E2E ASR and joint language identification. Toshniwal et al. [21] found that multilingual training leads to a significant relative improvement of recognition performance and the results can be further

boosted by conditioning the model on a language identifier. Some attempts take a step towards realizing language-universal ASR. Li et al. [22] proposed to replace the characters with the Unicode bytes as the output. Datta et al. [23] unified different writing systems through a many-to-one transliteration transducer. Recently, large-scale multilingual ASR systems have been investigated [3], [4], [8], [12], [24]. [3] proposed jointly training on 16,000 hours of speech data of 51 languages with up to 1 billion parameters. Inspired by [19], Hou et al. presented LID-42 [4], a large-scale multilingual acoustic Transformer model trained on 11 mixed corpora of 42 languages.

Cho et al. [25] validated the effectiveness of cross-lingual transfer learning for improving ASR performance. And this advantage can be further revealed by large-scale pre-training [3], [8]. For example, LID-42 can achieve a relative WER reduction of up to 28.1% on low-resource languages by cross-lingual transfer [4]. Yi et al. [26] introduced an adversarial learning objective to learn language-agnostic features. They appended a language discriminator after the shared encoder to distinguish which language the bottleneck features belong to. The objective of the discriminator is to correctly identify the language while the adversarial objective of the encoder is to fool the discriminator. The adversarial training process is realized with the use of the gradient reversal layer (GRL) [27]. Adams et al. [8] performed experiments to analyze the impacts of language similarity, context-independent phoneme CTC objective and the aforementioned language-adversarial classification objective during multilingual pre-training to encourage language-agnostic features for better cross-lingual adaptation.

Other than learning the language-agnostic features, the optimization-based meta-learning approaches [15], [28] that aim to find a proper initialization for rapid adaptation have also been explored for cross-lingual ASR [11]. Hsu et al. [10] proposed to apply the model-agnostic meta-learning (MAML) [15] as the pre-training method and achieved significant improvement over the conventional multilingual pre-training baseline. Xiao et al. [29] proposed the Adversarial Meta Sampling framework by introducing a policy network (task sampler) to dynamically sample languages based on their task difficulty. The ASR model is trained to minimize the loss while the task sampler learns to choose the most difficult languages that can maximize the loss. As a consequence, the learned initialization has a more balanced distance to all languages and shows a good generalization capacity in low-resource speech tasks.

B. Adapters

Due to the large quantity of parameters contained in the Transformer-based models [4], [30]–[32], recent literature proposed the *Adapter* structure [33], [34] for parameter-efficient adaptation of pre-trained Transformers [30], [35] on various downstream tasks including language understanding [36] and neural machine translation (NMT) [35], etc. Adapter is a versatile module that can be plugged into the Transformer blocks. The general philosophy for adapter-based fine-tuning is to freeze the parameters θ_b of the Transformer backbone and

only tune the parameters θ_a of the adapter. Compared to fine-tuning the whole Transformer model, fine-tuning the adapters is significantly efficient with acceptable performance loss [33]. Therefore, adapters have been adopted as a fine-tuning technique in few-shot domain adaptation for NMT [37] and unsupervised cross-lingual transfer [38] or domain adaptation [39] of large-scale pre-trained language models like BERT [30] and XLM [40]. Li et al. [41] proposed a hypernetwork that could generate parameters of task-specific adapters from task descriptions to enable zero-shot learning [42]. More recently, [14] introduced AdapterFusion to fuse adapters trained on different tasks to share the knowledge. The difference between our work and theirs is that we focus on the cross-lingual sequence-to-sequence ASR task while they experiment on text classification based on BERT [30].

Some researchers have proposed to apply the Adapters to the E2E ASR tasks. In [12], Kannan et al. proposes to use the adapters to handle the data imbalance problem for large-scale multilingual ASR. After obtaining the model trained on the union of data from all languages, they trained the language-dependent adapters on each of the languages again so that the multilingual backbone shares information across languages while the adapters could allow for per-language specialization. Winata et al. [13] extends this idea by further introducing a common adapter for all languages to learn language-agnostic information in the multilingual data. On the other hand, Hou et al. [11] investigates the possibility of applying adapters to cross-lingual ASR under the assumption that a large-scale pre-trained multilingual model [4] should have contained adequate general acoustic and linguistic knowledge and could be adapted to any unseen target language with moderate feature adaptation. Furthermore, they proposed to pre-train the adapters with meta-learning to obtain the MetaAdapter that provides a proper initialization for fast adaptation.

SimAdapter is similar to Mixture of Expert (MoE) [43]. MoE is often used to scale up the model size while retaining the computing efficiency. Therefore, there are often many experts and the expert outputs are often “sparsely activated” by using routing layers. In practice, only top- k experts are selected where $k = 1$ or 2. Also, MoE is trained on large-scale data along with the whole model where the experts acquire specific knowledge by themselves while each adapter component in our SimAdapter is “taught”, and the SimAdapter is applied to parameter-efficient adaptation to low-resource languages. However, we believe that some idea behind MoE is helpful to us. For example, as we observe that SimAdapter could distract its attention when the number of languages increases. We could also apply a similar routing mechanism to our SimAdapter. We will leave this for our future exploration.

III. EXPLOITING ADAPTERS FOR CROSS-LINGUAL ASR

A. Problem Definition

The goal of cross-lingual speech recognition is to transfer the knowledge from the existing languages to the new language. Formally speaking, given N rich-resource languages $\{S_1, S_2, \dots, S_N\}$, cross-lingual ASR aims at adapting the pre-trained model to an unseen target low-resource language L_T .

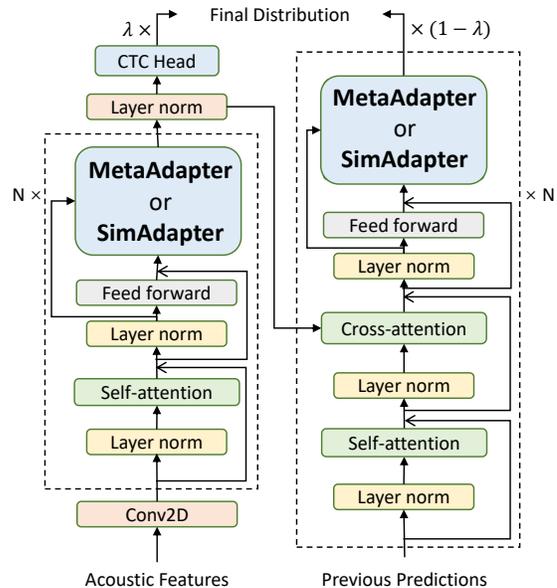


Fig. 2. Illustration of the MetaAdapter and SimAdapter module injected in the Speech-Transformer. Note that the residual connection between the feed-forward layer and layer normalization only applies to the SimAdapter.

Each language S_i is composed of the speech-text pairs and we typically use X and y to denote them, respectively, i.e., $S_i = \{X_j, y_j\}_{j=1}^{N_i}$, where N_i is the total number of training data. Also note that the target language is low-resource compared to the training languages, i.e., $N_T \ll N_i, \forall 1 \leq i \leq N$.

B. Overview

In this paper, we comprehensively investigate the potential of adapters to achieve parameter-efficient cross-lingual speech recognition. On the one hand, the parameters of the adapters are the only trainable parameters in the model with the rest parameters frozen, which remains parameter-efficient; on the other hand, the adapters module can also help reduce overfitting on the low-resource cross-lingual data.

To exploit adapters for cross-lingual ASR, it is important to study the relationship between different languages. In this paper, we comprehensively analyze the *MetaAdapter* as well as the newly proposed *SimAdapter* algorithms that learn and exploit the inter-language relationships to improve cross-lingual ASR. Generally speaking, the *MetaAdapter* is based on the meta-learning algorithm to extract general latent knowledge from existing training tasks and then adapt the knowledge to the target task. On the other hand, the *SimAdapter* algorithm is to directly explore the similarity between the source and target languages and then exploit such similarity to dynamically fuse the useful knowledge to the target language. Finally, we show that these two algorithms are not independent, but can be integrated for better performance. As shown in Fig. 2, our *MetaAdapter* and *SimAdapter* can be easily plugged into the Transformer backbone for implementation.

C. Backbone: Super Multilingual Transformer ASR Model

The super language-independent 42-lingual ASR model (LID-42) is proposed by Hou et al. in [4]. LID-42 is based

on the *big* Speech-Transformer [44] and joint CTC-attention structure [20]. We elaborate the model details below.

As model inputs, LID-42 accepts the 83-dimensional acoustic features (filter banks with pitch) computed with 10 ms frame shift and 25 ms frame length. The acoustic features are firstly subsampled by a factor of 4 by 2 convolution layers with kernel size 3 and stride 2. The resulted features have a receptive field of 100 milliseconds for each frame. Then the following encoder layers process the subsampled features by self-attention and feed-forward as illustrated in [35]. Apart from self-attention and feed-forward, the decoder layers further accept the encoder outputs and perform cross-attention.

For the CTC-attention hybrid structure, an auxiliary CTC task [45] is introduced for encoder outputs in order to encourage the monotonic alignment and accelerate the convergence speed [20]. In training, a weighted sum of the sequence-to-sequence attention loss \mathcal{L}_{ATT} and the CTC loss \mathcal{L}_{CTC} is employed:

$$\mathcal{L}_{ASR} = (1 - \lambda)\mathcal{L}_{ATT} + \lambda\mathcal{L}_{CTC}, \quad (1)$$

where λ denotes the weight of the CTC module.

Similarly, during decoding, the CTC module outputs are used to re-score the beam search results of the Transformer decoder on-the-fly:

$$\hat{Y} = \arg \max_{Y \in \mathcal{Y}} (1 - \lambda) \log P_{ATT}(Y|X) + \lambda \log P_{CTC}(Y|X), \quad (2)$$

where X are the 83-dimensional acoustic features (filter banks with pitch), \mathcal{Y} denotes the set of the decoding hypotheses.

As model outputs, a shared vocabulary including characters/subwords and language tokens (e.g., <en>, <fr>) of 42 languages is adopted to realize language-independent training and recognizing. Furthermore, a language token is inserted to the beginning of each training label as an auxiliary language identification target. The model is trained to firstly identify the language before recognizing the speech contents. It is worth noting that we focus on monolingual transfer in this work. Therefore, language-specific heads are used and the language identification objective is dropped during fine-tuning.

LID-42 is trained on around 5000-hour labeled speech data mixing 11 corpora covering 42 languages and has revealed a strong performance on cross-lingual transfer learning tasks as shown in previous works [4], [11].

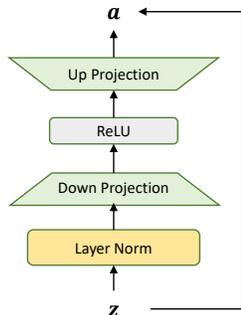


Fig. 3. Architecture of the adapter module.

D. Adapters

As shown in Fig. 3, a commonly-used adapter structure includes layer normalization, a down-projection layer, a non-linear activation function, and an up-projection layer. There is also a residual connection that allows the adapter to keep the original representation unchanged. Thus, the adapter function is formulated as:

$$\mathbf{a}^l = \text{Adapter}(\mathbf{z}^l) = \mathbf{z}^l + \mathbf{W}_u^l \text{ReLU}(\mathbf{W}_d^l (\text{LN}(\mathbf{z}^l))), \quad (3)$$

where \mathbf{z}^l represents the outputs of layer l , LN denotes layer normalization. \mathbf{W}_u , \mathbf{W}_d are weight parameters for up projection and down projection.

We will introduce these two algorithms and their integration in next sections.

IV. METAADAPTER

In this section, we introduce MetaAdapter in detail. MetaAdapter is inspired by the idea of meta-learning [46] for fast adaptation to the new target tasks. In our previous work [11], we investigated two meta-learning algorithms: Model-Agnostic Meta-Learning (MAML) [15] and Reptile [28]. We found that MAML is more robust to the overfitting problem brought by the variance of adaptation data size and pre-training epochs. Therefore, we adopt the MAML as our meta-training algorithm in this work.

However, it is expensive to perform meta-learning directly on the full Speech-Transformer model since the model has heavy parameters that could easily overfit the low-resource target data. To resolve this issue, MetaAdapter utilizes the adapters to significantly reduce the adaptation parameters by learning aiming a proper initialization for faster adaptation.

A. Architecture

The process of MetaAdapter is illustrated in Fig. 4. Given a pre-trained backbone speech-Transformer ASR model, MetaAdapter is composed of two phases: (i) meta-train the MetaAdapter on a bunch of source tasks; (ii) fine-tune the pre-trained adapter on unseen target tasks.

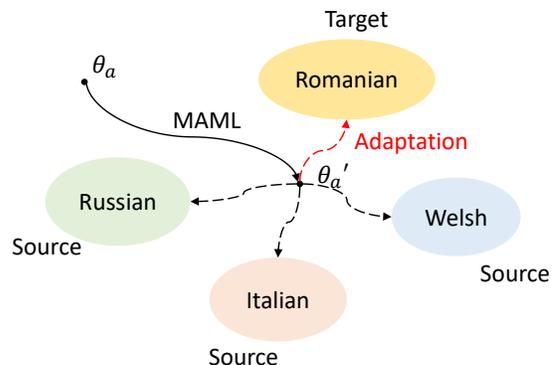


Fig. 4. Illustration of MetaAdapter.

To use meta-learning, we view different languages as different tasks. We split the parameters of MetaAdapter into two types: the backbone parameters θ_b (i.e., vanilla Transformer) and the parameters of all adapters θ_a . Thus, given

N different source languages $\{S_1, S_2, \dots, S_N\}$, we pre-train the MetaAdapter module f_{θ_a} to obtain good initialization parameters θ_a that could generalize for fast adaptation given any unseen target language. Meanwhile, parameters of the pre-trained backbone θ_b are frozen during both the pre-training and the fine-tuning.

B. Training MetaAdapter

In each pre-training episode, two subsets are randomly sampled from each source training language S_i , namely meta-training set S_i^{tr} and meta-validation set S_i^{val} , i.e., $S_i^{tr} \cap S_i^{val} = \emptyset$. One episode is composed of two iterations: an inner iteration and an outer iteration. In the inner iteration, MAML updates the adapter parameters θ_a by performing one or more gradient descent on S_i^{tr} . For notation simplicity, the updated parameter for language S_i using the inner gradient descent iteration is:

$$\theta'_{a,i} = \theta_a - \epsilon \nabla \mathcal{L}_{S_i^{tr}}(f_{\theta_a}), \quad (4)$$

where \mathcal{L} is the ASR loss function as introduced in section III-C and ϵ is the fast adaptation learning rate. In the outer iteration, the adapter parameters are then optimized to improve the performance of $f_{\theta'_{a,i}}$ with respect to θ_a across all the meta-validation sets S_i^{val} . The meta-optimization objective of the outer iteration is:

$$\mathcal{L}_{S_i^{val}}(f_{\theta'_{a,i}}) = \mathcal{L}_{S_i^{val}}(f_{\theta_a - \epsilon \nabla_{\theta_a} \mathcal{L}_{S_i^{tr}}(f_{\theta_a})}). \quad (5)$$

We optimize the meta-optimization objective through gradient descent as:

$$\theta_a = \theta_a - \mu \sum_{i=1}^N \nabla_{\theta_a} \mathcal{L}_{S_i^{val}}(f_{\theta'_{a,i}}), \quad (6)$$

where μ is the meta step size.

After pre-training, the MetaAdapter should obtain a proper initialization for any unseen target language(s). The complete training procedure of the MetaAdapter is presented in Algo. 1.

Algorithm 1 Learning algorithm of the MetaAdapter

Input: Rich-resource languages $\{S_1, \dots, S_N\}$, low-resource task L_T .

- 1: Train language-specific heads on source languages S_i .
 - 2: Initialize the MetaAdapter.
 - 3: **while** meta-learning not done **do**
 - 4: Optimizing the MetaAdapter using Eq. (6).
 - 5: **end while**
 - 6: Train the target head on target language L_T .
 - 7: Fine-tune the MetaAdapter using ASR loss Eq. (1).
 - 8: **return** Cross-lingual ASR model.
-

V. SIMADAPTER

We propose SimAdapter to improve the adapter-based cross-lingual adaptation as well as the model interpretability by explicitly leveraging the knowledge of the source languages from the adapter modules. Here, ‘Sim’ refers to similarity.

SimAdapter is inspired by existing research on language and speech origins [16]–[18], which implies that different

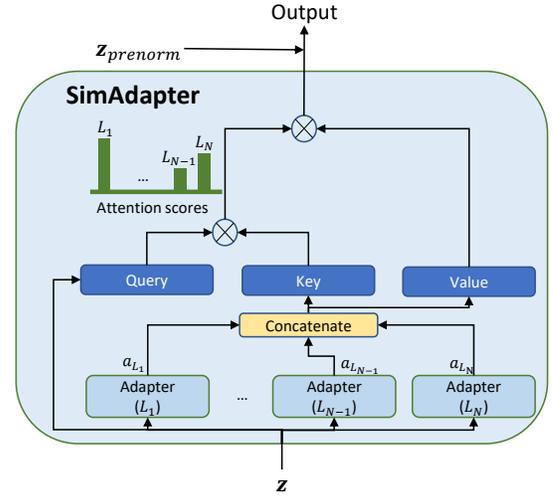


Fig. 5. Illustration of the SimAdapter module. The language-specific features a_{L_k} of different languages $L_k \in \{L_1, L_2, \dots, L_N\}$ are attended by the language-agnostic features \mathbf{z} to extract better features for the target language.

languages in the world are sharing similarities based on their similar geological characteristics or cultural developments. Thus, it is feasible to leverage the knowledge of multilingual adapters for new target languages.

A. Architecture

SimAdapter is a parameter-efficient algorithm that learns the similarities between existing language-specific adapters and the target low-resource language based on the attention mechanism [35]. Similar to the adapters, SimAdapter can also be easily integrated with existing pre-trained models and adapters.

The detailed composition of the SimAdapter is shown in Fig. 5. By taking the language-agnostic representations from the backbone model as the query, and the language-specific outputs from multiple adapters as the keys and values, the final output for SimAdapter over attention are computed as (For notation simplicity, we omit the layer index l below):

$$\text{SimAdapter}(\mathbf{z}, \mathbf{a}_{\{S_1, S_2, \dots, S_N\}}) = \sum_{i=1}^N \text{Attn}(\mathbf{z}, \mathbf{a}_{S_i}) \cdot (\mathbf{a}_{S_i} \mathbf{W}_V), \quad (7)$$

where $\text{SimAdapter}(\cdot)$ and $\text{Attn}(\cdot)$ denotes the SimAdapter and attention operations, respectively. Specifically, the attention operation is computed as:

$$\text{Attn}(\mathbf{z}, \mathbf{a}) = \text{Softmax}\left(\frac{(\mathbf{z} \mathbf{W}_Q)(\mathbf{a} \mathbf{W}_K)^\top}{\tau}\right), \quad (8)$$

where τ is the temperature coefficient, $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ are attention matrices. Note that while $\mathbf{W}_Q, \mathbf{W}_K$ are initialized randomly, \mathbf{W}_V is initialized with a diagonal of ones and the rest of the matrix with small weights ($1e-6$) to retain the adapter representations. Furthermore, a regularization term is introduced to avoid drastic feature changes:

$$\mathcal{L}_{\text{reg}} = \sum_{i,j} ((\mathbf{I}_V)_{i,j} - (\mathbf{W}_V)_{i,j})^2, \quad (9)$$

where \mathbf{I}_V is the identity matrix with the same size as \mathbf{W}_V .

In our cross-lingual setting, the SimAdapter module is expected to utilize language-specific knowledge from existing language adapters.

B. Fusion Guide Loss

Although SimAdapter aims to benefit from the similar knowledge of other languages, we believe that the most crucial information is stored in the adapter of the target language. However, since the weights of source and target adapters are initialized equally, SimAdapter often distracts its attention significantly from the target language during adaptation and generally does not perform well in our experiments. To alleviate this problem, we propose a fusion guide loss to encourage the model to focus on the corresponding adapters for the similarity learning. Specifically, for each language fusion layer f , we average the cross entropy of adapter attention scores among all K time steps and S samples. The layer-wise guide losses are added up as:

$$\mathcal{L}_{\text{guide}}^f = -\frac{1}{K \times S} \sum_{s=1}^S \sum_{k=1}^K \log \alpha_{f,k}^s, \quad (10)$$

$$\mathcal{L}_{\text{guide}} = \sum_f \mathcal{L}_{\text{guide}}^f. \quad (11)$$

Note that K represents the number of frames in the encoder and the number of tokens in the decoder side, $\alpha_{f,k}^s$ denotes the attention score of target language's Adapter. In this way, the attention scores are optimized via conventional Empirical Risk Minimization (ERM) [47].

C. Training SimAdapter

A difference between the previous application of Adapter-Fusion [14] and our SimAdapter for cross-lingual ASR is that a language-specific language head is required to be trained for the unseen target language. However, training the Adapters together with the language heads may result in the insufficient learning of semantic information in the adapters. Therefore, in this work, we introduced a three-stage training strategy for SimAdapter-based ASR cross-lingual adaptation.

In the first stage, different from the previous work [11], SimAdapter trains the language-specific heads for each source language S_i as well as the target language separately. This step aligns the language heads to the same latent semantic space of the backbone model. In the second stage, adapters are trained based on the pre-trained heads to learn the information. In the third stage, SimAdapter leverages the fusion of source languages for better adaptation to the target language. Only the parameters of the SimAdapter are trained in this stage.

By adding the \mathbf{W}_V regularization term weighted by η and the fusion guided loss weighted by γ , the final adaptation objective is given by:

$$\mathcal{L} = \mathcal{L}_{\text{ASR}} + \eta \mathcal{L}_{\text{reg}} + \gamma \mathcal{L}_{\text{guide}}. \quad (12)$$

The complete training procedure of SimAdapter is presented in Algorithm 2.

Algorithm 2 Learning algorithm of SimAdapter

Input: Rich-resource languages $\{S_1, \dots, S_N\}$, low-resource task L_T .

- 1: Train language-specific heads on the source languages S_i and the target language.
 - 2: Train Adapters A_t on top of language-specific heads.
 - 3: Initialize SimAdapter layers.
 - 4: **while** not done **do**
 - 5: Optimizing SimAdapter layers using Eq. (12).
 - 6: **end while**
 - 7: **return** Target ASR model.
-

D. Integration of MetaAdapter and SimAdapter

Although MetaAdapter and SimAdapter can both benefit cross-lingual adaptation by leveraging the knowledge of source languages, they are designed from different perspectives. MetaAdapter aims to obtain a proper initialization for fast adaptation by learning from the source languages, which can be regarded as a type of latent transfer. On the other hand, SimAdapter explicitly models the similarities between source and target languages with the attention mechanism. Therefore, MetaAdapter is good at handling extremely low-resource languages, while with more training data SimAdapter can capture the language similarities more precisely.

Moreover, note that MetaAdapter and SimAdapter are not independent, but can be integrated into one algorithm, which we denote as *SimAdapter+*. We can simply fuse the source adapters with the target adapter learned by the MetaAdapter using SimAdapter. This can be seen as a two-stage knowledge transfer process where we aim to learn general and transferable knowledge from meta-learning in the first stage; then, we perform adaptation using the SimAdapter algorithm for fine-grained knowledge transfer to achieve better performance.

VI. EXPERIMENTAL SETUP

A. Data Set

We adopt the Common Voice 5.1 [48] corpus for our experiments. We follow the official data splits for training, validation and testing. For the SimAdapter, we select five rich-resource languages as source languages and five low-resource languages as targets. Note that the source and target languages are all from European and some of them are spoken in geographically close regions to empirically analyze if the language similarities can be revealed by SimAdapter. The detailed data statistics are shown in TABLE I.

B. Compared Approaches

We consider the following fine-tuning-based approaches as well as both end-to-end and conventional hybrid models and trained from random initialization for comparison. To evaluate the efficiency of different methods, we also list numbers of trainable parameters in Table II. It is shown that our MetaAdapter and SimAdapter (and SimAdapter+) only use **2.5%** and **15%** of the training parameters from the full model, respectively, which are significantly parameter-efficient.

TABLE I
TRAINING / VALIDATION / TESTING HOURS

	Language	Train	Valid	Test
Source	Russian (ru)	80.61	11.78	12.61
	Welsh (cy)	74.84	4.35	4.25
	Italian (it)	88.74	19.74	20.85
	Basque (eu)	73.26	7.46	7.85
	Portuguese (pt)	37.40	5.40	5.85
Target	Romanian (ro)	3.04	0.42	1.66
	Czech (cs)	20.66	2.84	3.13
	Breton (br)	2.84	1.51	1.75
	Arabic (ar)	7.87	2.01	2.09
	Ukrainian (uk)	17.35	2.30	2.36

TABLE II
COMPARISON OF NUMBER OF TRAINABLE PARAMETERS.

Method	# Trainable Parameters
Hybrid DNN/HMM	14,247K
Full Model	27,235K
Head	77K
Head+(Meta-)Adapter	676K
Head+(Meta-)Adapter+SimAdapter	4,224K

1) Baselines without applying transfer learning:

- DNN/HMM: Standard hybrid DNN/HMM models are trained with lattice-free MMI [49] criterion using Kaldi [50]. Specifically, we use 9 layers TDNN [51] the acoustic model. The acoustic features are 100-dimensional i-vector [52] and 40-dimensional MFCC. 3-gram language model is applied for decoding.¹
- Trans.(B): We train a randomly-initialized big Transformer model of the same size and architecture as LID-42 from scratch.
- Trans.(S): To mitigate overfitting, we decrease the feed forward network from 2048 to 1024 so that the number of model parameters is reduced from 27,235K to 18,664K.
- Head: We keep the backbone model (LID-42) frozen as feature extractor and train the language-specific heads on top of it.

2) Fine-tuning based transfer:

- Full-FT: We fine-tune the full model on each target language individually, leading to 5 separate models.
- Full-FT+L2: We apply L2 regularizations to Full-FT to avoid overfitting.
- Part-FT: We make only the last 3 decoder layers trainable and freeze the rest parameters to fine-tune on the target languages to mitigate overfitting.

3) Adapter based transfer:

- Adapter: We inject and train the vanilla adapters while keeping the backbone model frozen.
- MetaAdapter: We inject the pre-trained MetaAdapter and train it as the vanilla adapters do.
- SimAdapter: We fuse the adapters of the source languages with the target language to improve the

performance.

- SimAdapter+: Specifically, we combine the MetaAdapter and the SimAdapter (namely SimAdapter+) to evaluate its performance and verify whether MetaAdapter and SimAdapter are compatible.

C. Implementation Details

We implement the E2E methods based on the ESPnet [53] codebase. The subword-based LID-42 model proposed in [4] is used as the backbone model for adaptation. The acoustic features are extracted following ESPnet. Numbers of SentencePiece [54] subwords are set to 150 and 100 for high- and low-resource languages, respectively.

We use Adam [55] as the optimizer. For the full-model fine-tuning, we follow the same learning rate scheduling strategy as [35] and warmup the learning rate to 0.2 in the first 10 epochs. The total number of training epochs is 200 for full-model fine-tuning and SimAdapter, and 100 for the other methods. Early stopping with patience 10 is adopted except for the training of source heads and adapters. The source languages heads and adapters are trained using a batch size of 1024 and a learning rate of 0.028. The target heads and adapters are trained using a batch size of 512 and a learning rate of 0.02. For the SimAdapter, we use a batch size of 128 and a learning rate of $2e - 5$. We adopt $\eta = 0.01$ for the regularization loss weight and 1.0 as the guide loss weight γ . The temperature coefficient τ is simply set to 1.0. We train the MetaAdapter for 30 epochs using Adam [55] with $\beta_1 = 0$ in the inner training loop and vanilla stochastic gradient descent (SGD) in the outer loop. The inner adaptation learning rate and initial meta step size μ are 0.028 and 1.0, respectively. The meta step size linearly annealed to 0 over the course of training. The weight of the CTC module λ is set to 0.3 throughout the experiments following ESPnet [53]. Beam size 10 is employed for joint decoding. Our source code is available at: <https://github.com/jindongwang/transferlearning/tree/master/code/ASR/Adapter>.

D. Evaluation Metrics

In this work, we use word error rate (WER) as our evaluation metric. We average the results on 5 languages to evaluate the overall performance of different methods by default. To reflect the performance on target languages according to their imbalanced test data duration (more test data often represents more training data), we also compute the weighted average WERs, which is more friendly to the methods that require relatively more training data.

VII. EXPERIMENTAL RESULTS

A. Cross-lingual speech recognition

Table III shows the main results on cross-lingual ASR. The first three columns show the non-fine-tuning-based baselines. First, it can be found that the hybrid DNN/HMM model outperforms Transformer (big) on 2 out of 4 languages (Romanian (ro), Arabic (ar)), and these 2 languages are with least training

¹We did not find proper pronunciation dictionary for Breton. Therefore, only results of the other 4 languages are presented.

TABLE III
WORD ERROR RATES (WER) ON THE CROSS-LINGUAL ASR TASKS

Target	DNN/HMM	Trans.(B)	Trans.(S)	Head	Full-FT	Full-FT+L2	Part-FT	Adapter	SimAdapter	MetaAdapter	SimAdapter+
Romanian (ro)	70.14	97.25	94.72	63.98	53.90	52.74	52.92	48.34	47.37	44.59	47.29
Czech (cs)	63.15	48.87	51.68	75.12	34.75	35.80	54.66	37.93	35.86	37.13	34.72
Breton (br)	-	97.88	92.05	82.80	61.71	61.75	66.24	58.77	58.19	58.47	59.14
Arabic (ar)	69.31	75.32	74.88	81.70	47.63	50.09	58.49	47.31	47.23	46.82	46.39
Ukrainian (uk)	77.76	64.09	67.89	82.71	45.62	46.45	66.12	50.84	48.73	49.36	47.41
AVG	-	76.68	76.24	77.26	48.72	49.37	59.69	48.64	47.48	47.27	46.99
Weighted AVG	-	72.28	72.50	77.54	46.72	47.50	59.43	47.38	46.08	46.12	45.45

data. The results indicate that the overfitting issue occurs in the Transformer model. Transformer (S) mitigates the problem to some extent but it is still far from satisfaction. It could further be inferred that even hybrid DNN/HMM has the overfitting problem on the extremely low-resource Romanian language, since lower WER is obtained with the linear head simply trained on top of the frozen but powerful LID-42 backbone.

On the other hand, from the fine-tuning- and adapter- based approaches presented on the middle- and right-hand sides, we can observe that the adapters successfully avoid the overfitting problem and outperform the Full-FT method on 3 very low-resource languages (Romanian, Breton, Arabic). Applying L2 regularization and partial fine-tuning both improve the performance on Romanian but degrades on the other 4 languages. It can be also found that the MetaAdapter and SimAdapter approaches can achieve similar and competitive results on the 5 target languages. Furthermore, we notice that both the MetaAdapter and SimAdapter consistently improve the performance over the adapters and narrow the gap with Full-FT on the languages with relatively more training data (Czech and Ukrainian). Meanwhile, the MetaAdapter method performs better on the extremely low-resource languages (ar, ro) and has lower average WER, while SimAdapter shows better results on moderate low-resource languages (br, cs) and obtains lower weighted average WER. Finally, by combining the MetaAdapter with SimAdapter, the SimAdapter+ method surpasses all the other approaches and obtains the best average performance on the 5 languages, indicating that the two methods are compatible since they leverage the source information in different ways. Combining the results from TABLE II where SimAdapter+ only uses 15.5% trainable parameters of the full model, we see that SimAdapter+ is both effective and parameter-efficient.

B. Ablation Study

1) *Impact of different training strategies:* We compare the impact brought by different adapter-training strategies, i.e., jointly training the adapter with head and the first two stages of the training strategy proposed in Section V-C. The results are presented in Table IV. It is clear that the proposed two-stage training strategy can consistently reduce the WERs of both the adapters and the SimAdapter.

2) *Impact of pre-training epochs for MetaAdapter:* To validate the meta-training effects for the MetaAdapter, we select checkpoints of 5 pre-trained epochs {10, 15, 20, 25, 30} and fine-tune them following the same setting as explained in Section VI. We present the results in Fig. 8(a). It could

TABLE IV
COMPARISON OF DIFFERENT ADAPTER TRAINING STRATEGIES.

Target	Joint	+SimAdapter	Two-stage	+SimAdapter
ro	52.92	53.88	48.34	47.37
cs	39.16	36.79	37.93	35.86
br	65.10	63.37	58.77	58.19
ar	50.53	49.31	47.31	47.23
uk	52.27	48.84	50.84	48.73
Average	52.00	50.44	48.64	47.48
+Weighted	50.35	48.57	47.38	46.08

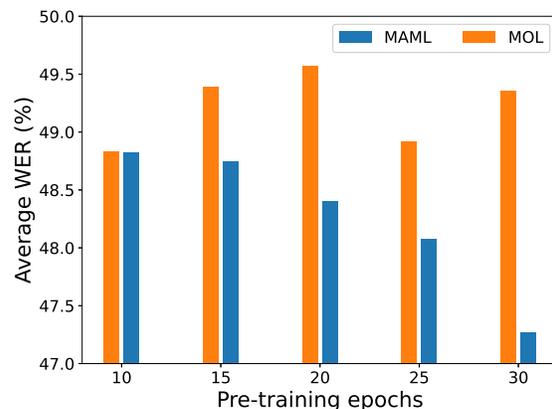


Fig. 6. Comparison between MAML and conventional multi-objective learning (MOL) approach for Adapter pre-training.

be found that the WERs are reduced with more pre-training epochs, indicating the effectiveness of meta-learning.

For comparison, we also conduct the same experiment on another adapter pre-trained on source languages using the conventional multi-objective learning (MOL) method and visualize the average WERs in Fig. 6. It is clear that with the more pre-training epochs, the MOL-trained adapter tends to overfit the source data and performs worse on the target languages.

3) *Analyzing the weight of guide loss for SimAdapter:* We then analyze the impacts of the weight γ of the proposed guide loss within $\{0, 0.001, 0.01, 0.1, 0.5, 0.75, 1.0\}$ for the SimAdapter. As shown in Fig. 8(b), the model performances on the 5 languages generally get improved with the increasing of γ when $\gamma < 0.5$. When $\gamma \geq 0.5$, the WER may vary among languages. The best overall performance is obtained when $\gamma = 1$. In real applications, the value of γ needs to tune on the target dataset.

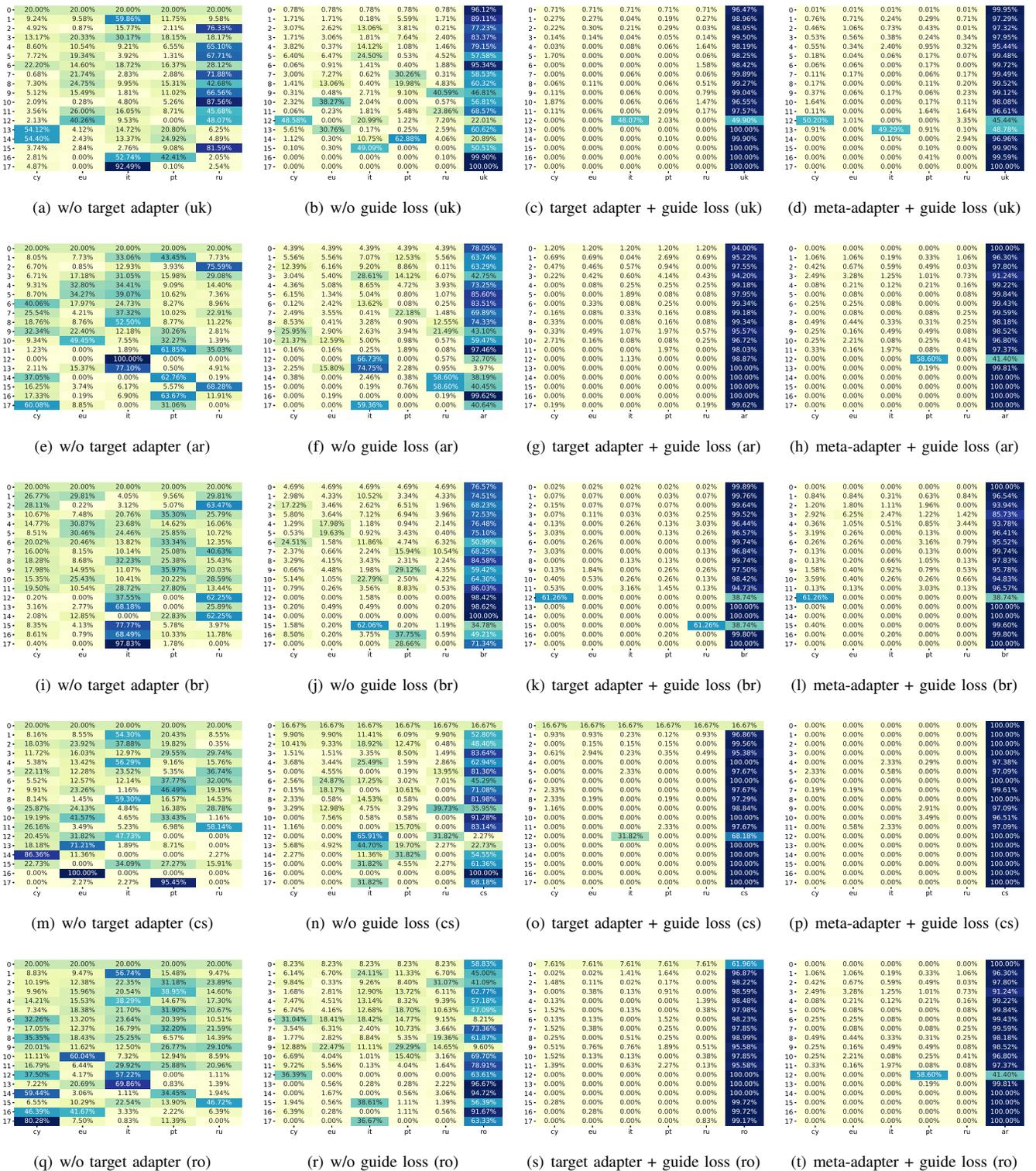


Fig. 7. Attention matrices of five low-resource target languages. A row in the figure denotes a language, whose four settings are: (1) without target adapter, (2) with target adapter but no guide loss ($\gamma = 0$), (3) with target adapter and guide loss, and (4) SimAdapter+. Column index indicates the Transformer layer number, where 0th to 11th layers are encoders, 12th to 17th are decoders. *Best viewed in color and zoomed in.*

4) *How much information can be shared across languages:* Although SimAdapter improves the WER results, we do not know whether and how much it could benefit from other languages. Therefore, we conduct two experiments to validate this. Firstly, we examine how much the other languages can contribute without using the adapters from languages to

see whether additional gains can be obtained with only source adapters. TABLE V shows the results. It can be found that even without the target adapter, SimAdapter can still improve the performance for most of the languages except for Romanian, indicating the effectiveness of learning language information from source adapters.

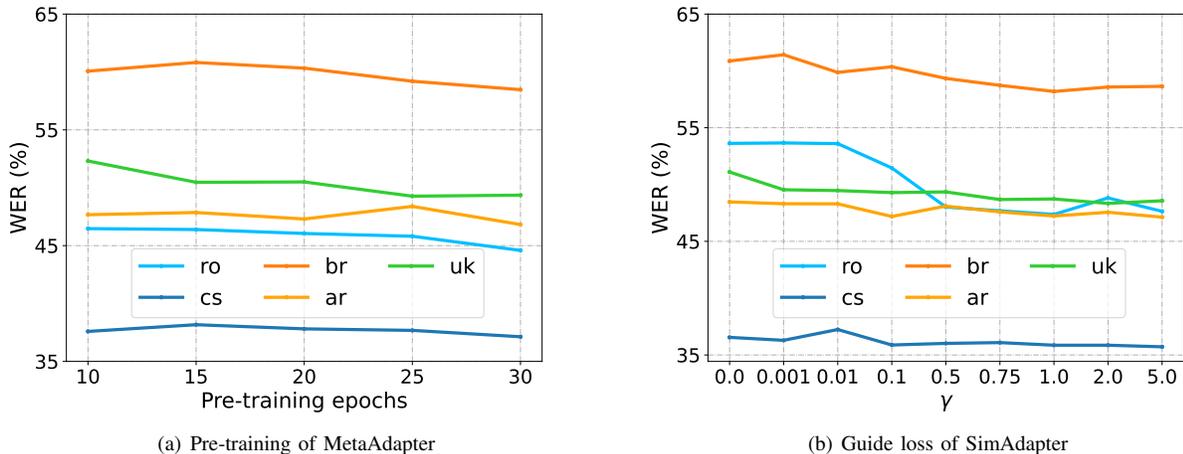


Fig. 8. Analysis of (a) pre-training epochs of MetaAdapter and (b) importance of guide loss in SimAdapter.

TABLE V
WER RESULTS OF SIMADAPTER WITH OR WITHOUT ADAPTER L_T .
FUSION GUIDE LOSS IS SET TO 0 FOR SIMADAPTER WITH ADAPTER L_T .

Target	Head	w/o Adapter L_T	w/ Adapter L_T
ro	63.98	67.83	53.62
cs	75.12	55.06	36.55
br	82.80	77.04	60.87
ar	81.70	64.68	48.47
uk	82.71	69.09	51.10
Average	77.26	66.74	50.12
+Weighted	77.54	65.33	48.39

In the second experiment, we train two different SimAdapter models on Ukrainian by fusing the target-language adapter and one source-language adapter to analyze the contributions of different source languages. Specifically, we choose Italian and Russian as the source languages. Since Russian is more similar to Ukrainian than Italian, we expect more gains of SimAdapter trained with the Russian adapter. The results align with our expectation. We observe that the SimAdapter with Italian adapter obtains a WER of 48.70, while with the Russian adapter, the WER is 47.73, indicating that SimAdapter could transfer more useful knowledge from Russian than Italian to model the Ukrainian language.

C. Attention Visualization

To further show the relationship between source and target languages, we visualize the attention maps for each target language. The attention value reflects their similarities. Fig. 7 shows the results of three different types of languages: (1) without target adapter, (2) with target adapter but no guide loss ($\gamma = 0$), (3) with target adapter and guide loss, and (4) with target MetaAdapter and guide loss.

We take the Ukrainian (uk) as an example. Firstly, from the figure on the left, we can observe a trend that SimAdapter layers tend to pay more attention to the Russian (ru)’s adapter, which could be because of the linguistic similarity between Ukrainian and Russian. However, after introducing the target

adapter, SimAdapter layers obviously turn to focus more on the target adapter, but there are still diverse attentions across other languages. By introducing the guide loss, the SimAdapter layers are forced to pay more attention to the target adapter and fusing less information from other languages.

We also notice that in the first encoder layer, the attention distribution seems to be uniform across the source languages. By analyzing the outputs, we found that the adapters in the first layer tend to keep the backbone representation unchanged via the residual connection. The same phenomenon can also be observed in the Czech (cs) target language. A possible reason could be that the first layer is to extract general acoustic features which are language-independent. Since we observe a similar trend in the first decoder layer (layer 12) that the attention distributions tend to be more distracted, we thus assume that adapters in the bottom layers in both the encoder and decoder are less important for cross-lingual adaptation, which we conduct experiments in next subsection to analyze the performance of fusing different adapters.

D. Do all Adapter layers need to be fused?

By observing the attention maps, we notice that for some layers, the attention seems to focus solely on the target adapter with a 100% attention score. This phenomenon occurs more frequently in the higher decoder layers, i.e., 12th to 17th layers in Fig. 7. In such cases, the fusion seems not to be necessary. We doubt whether we can achieve comparable performance while fusing adapters in part of the layers only. Therefore, we conduct the ablation experiments by only fusing part of the layers. The results are presented in TABLE VI. Although some languages (e.g., Breton) can retain the performance by only fusing 2 bottom layers, fusing more layers generally lead to better performance.

E. Training and inference time

Finally, we compare the average training time of full-model fine-tuning, MetaAdapter and SimAdapter methods per

TABLE VI
ABLATION STUDY OF THE ENCODER AND DECODERS

Target	Enc1-Dec1	Enc12-Dec1	Enc12-Dec6
ro	48.39	48.25	47.37
cs	37.31	36.30	35.86
br	57.85	59.08	58.19
ar	47.48	47.34	47.23
uk	50.58	48.98	48.73
Average	48.32	47.99	47.48
+Weighted	47.04	46.55	46.08

TABLE VII
AVERAGE TRAINING / INFERENCE TIME.

	Training Time (sec.)	RTF
Full-FT	0.253 (-)	0.045 (-)
MetaAdapter	0.143 (43.48%↓)	0.043 (4.06%↓)
SimAdapter	0.263 (3.95%↑)	0.055 (22.12%↑)

iteration as well as their inference real-time factor (RTF) on the 5 target languages. The RTF metric is used to evaluate the decoding time cost by computing the ratio of the model decoding time to the total utterance duration on the test data. The training and decoding are conducted on 1 GeForce RTX 2080 Ti GPU with batch size 64. The results are shown in TABLE VII.

It could be found that the MetaAdapter module significantly accelerates the training process while the SimAdapter introduces minor additional time cost compared with full-model fine-tuning. The RTFs of Full-FT and MetaAdapter are at the same level. The reason that MetaAdapter has slightly lower RTF could be due to its shorter average prediction lengths. On the other hand, the relative RTF increases of 22.12% brought by SimAdapter is also acceptable.

VIII. CONCLUSIONS AND FUTURE WORK

In this paper, we propose to exploit MetaAdapter and SimAdapter for adapter-based cross-lingual speech recognition. The proposed SimAdapter leverages the attention mechanism to learn the similarities between the source and target languages during fine-tuning using the adapters. We also show that the two algorithms can be integrated for better performance. Experiments on five low-resource languages from Common Voice dataset demonstrate the superiority of the two algorithms. In the future, we plan to extend these algorithms to other language families and further improve the training and inference speed of our methods.

REFERENCES

- [1] D. Wang, X. Wang, and S. Lv, "An overview of end-to-end automatic speech recognition," *Symmetry*, vol. 11, no. 8, p. 1018, 2019.
- [2] S. Karita, N. Chen, T. Hayashi, T. Hori, H. Inaguma, Z. Jiang, M. Someki, N. E. Y. Soplin, R. Yamamoto, X. Wang *et al.*, "A comparative study on transformer vs rnn in speech applications," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 449–456.
- [3] V. Pratap, A. Sriram, P. Tomasello, A. Hannun, V. Liptchinsky, G. Synnaeve, and R. Collobert, "Massively multilingual asr: 50 languages, 1 model, 1 billion parameters," *Proc. Interspeech 2020*, pp. 4751–4755, 2020.
- [4] W. Hou, Y. Dong, B. Zhuang, L. Yang, J. Shi, and T. Shinozaki, "Large-Scale End-to-End Multilingual Speech Recognition and Language Identification with Multi-Task Learning," in *Proc. Interspeech 2020*, 2020, pp. 1037–1041. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2020-2164>
- [5] W. Hou, J. Wang, X. Tan, T. Qin, and T. Shinozaki, "Cross-Domain Speech Recognition with Unsupervised Character-Level Distribution Matching," in *Proc. Interspeech 2021*, 2021, pp. 3425–3429.
- [6] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, "Unsupervised Cross-Lingual Representation Learning for Speech Recognition," in *Proc. Interspeech 2021*, 2021, pp. 2426–2430.
- [7] E. Chuangsuwanich, "Multilingual techniques for low resource automatic speech recognition," Massachusetts Institute of Technology Cambridge United States, Tech. Rep., 2016.
- [8] O. Adams, M. Wiesner, S. Watanabe, and D. Yarowsky, "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)*, 2019, pp. 96–108.
- [9] S. Zhou, S. Xu, and B. Xu, "Multilingual end-to-end speech recognition with a single transformer on low-resource languages," *arXiv preprint arXiv:1806.05059*, 2018.
- [10] J.-Y. Hsu, Y.-J. Chen, and H.-y. Lee, "Meta learning for end-to-end low-resource speech recognition," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7844–7848.
- [11] W. Hou, Y. Wang, S. Gao, and T. Shinozaki, "Meta-adapter: Efficient cross-lingual adaptation with meta-learning," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021.
- [12] A. Kannan, A. Datta, T. N. Sainath, E. Weinstein, B. Ramabhadran, Y. Wu, A. Bapna, Z. Chen, and S. Lee, "Large-Scale Multilingual Speech Recognition with a Streaming End-to-End Model," in *Proc. Interspeech 2019*, 2019, pp. 2130–2134.
- [13] G. I. Winata, G. Wang, C. Xiong, and S. Hoi, "Adapt-and-Adjust: Overcoming the Long-Tail Problem of Multilingual Speech Recognition," in *Proc. Interspeech 2021*, 2021, pp. 2451–2455.
- [14] J. Pfeiffer, A. Kamath, A. Rücklé, K. Cho, and I. Gurevych, "AdapterFusion: Non-destructive task composition for transfer learning," in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Online: Association for Computational Linguistics, Apr. 2021, pp. 487–503. [Online]. Available: <https://aclanthology.org/2021.eacl-main.39>
- [15] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *International Conference on Machine Learning*. PMLR, 2017, pp. 1126–1135.
- [16] J. Ke and J. H. Holland, "Language origin from an emergentist perspective," *Applied Linguistics*, vol. 27, no. 4, pp. 691–716, 2006.
- [17] P. F. MacNeilage, *The origin of speech*. Oxford University Press, 2010, no. 10.
- [18] D. W. Frayer and C. Nicolay, "14 fossil evidence for the origin of speech sounds," 2000.
- [19] S. Watanabe, T. Hori, and J. R. Hershey, "Language independent end-to-end architecture for joint language identification and speech recognition," in *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2017, pp. 265–271.
- [20] S. Kim, T. Hori, and S. Watanabe, "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)*. IEEE, 2017, pp. 4835–4839.
- [21] S. Toshniwal, T. N. Sainath, R. J. Weiss, B. Li, P. Moreno, E. Weinstein, and K. Rao, "Multilingual speech recognition with a single end-to-end model," in *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2018, pp. 4904–4908.
- [22] B. Li, Y. Zhang, T. Sainath, Y. Wu, and W. Chan, "Bytes are all you need: End-to-end multilingual speech recognition and synthesis with bytes," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 5621–5625.
- [23] A. Datta, B. Ramabhadran, J. Emond, A. Kannan, and B. Roark, "Language-agnostic multilingual modeling," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 8239–8243.
- [24] B. Li, R. Pang, T. N. Sainath, A. Gulati, Y. Zhang, J. Qin, P. Haghani, W. R. Huang, and M. Ma, "Scaling end-to-end models for large-scale multilingual asr," in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2021.

- [25] J. Cho, M. K. Baskar, R. Li, M. Wiesner, S. H. Mallidi, N. Yalta, M. Karafiat, S. Watanabe, and T. Hori, "Multilingual sequence-to-sequence speech recognition: architecture, transfer learning, and language modeling," in *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2018, pp. 521–527.
- [26] J. Yi, J. Tao, Z. Wen, and Y. Bai, "Adversarial multilingual training for low-resource speech recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4899–4903.
- [27] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *The journal of machine learning research*, vol. 17, no. 1, pp. 2096–2030, 2016.
- [28] A. Nichol, J. Achiam, and J. Schulman, "On first-order meta-learning algorithms," *arXiv preprint arXiv:1803.02999*, 2018.
- [29] Y. Xiao, K. Gong, P. Zhou, G. Zheng, X. Liang, and L. Lin, "Adversarial meta sampling for multilingual low-resource speech recognition," in *Association for the Advancement of Artificial Intelligence*, 2021.
- [30] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," 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)*, 2019, pp. 4171–4186.
- [31] A. Baevski, S. Schneider, and M. Auli, "vq-wav2vec: Self-supervised learning of discrete speech representations," in *International Conference on Learning Representations*, 2019.
- [32] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [33] N. Hounsby, A. Giurghi, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly, "Parameter-efficient transfer learning for nlp," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2790–2799.
- [34] A. Bapna and O. Firat, "Simple, scalable adaptation for neural machine translation," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 1538–1548.
- [35] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems*, vol. 30, pp. 5998–6008, 2017.
- [36] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, "Glue: A multi-task benchmark and analysis platform for natural language understanding," *EMNLP 2018*, p. 353, 2018.
- [37] A. Sharaf, H. Hassan, and H. Daumé III, "Meta-learning for few-shot nmt adaptation," in *Proceedings of the Fourth Workshop on Neural Generation and Translation*, 2020, pp. 43–53.
- [38] M. Artetxe, S. Ruder, and D. Yogatama, "On the cross-lingual transferability of monolingual representations," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 4623–4637.
- [39] J. Li, R. He, H. Ye, H. T. Ng, L. Bing, and R. Yan, "Unsupervised domain adaptation of a pretrained cross-lingual language model," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, C. Bessiere, Ed. International Joint Conferences on Artificial Intelligence Organization, 7 2020, pp. 3672–3678, main track. [Online]. Available: <https://doi.org/10.24963/ijcai.2020/508>
- [40] A. CONNEAU and G. Lample, "Cross-lingual language model pretraining," in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019. [Online]. Available: <https://proceedings.neurips.cc/paper/2019/file/c04c19c2c2474dbf5f7ac4372c5b9af1-Paper.pdf>
- [41] Q. Ye and X. Ren, "Learning to generate task-specific adapters from task description," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. Online: Association for Computational Linguistics, Aug. 2021, pp. 646–653. [Online]. Available: <https://aclanthology.org/2021.acl-short.82>
- [42] O. Weller, N. Lourie, M. Gardner, and M. Peters, "Learning from task descriptions," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 1361–1375.
- [43] Y. J. Kim, A. A. Awan, A. Muzio, A. F. C. Salinas, L. Lu, A. HENDY, S. Rajbhandari, Y. He, and H. H. Awadalla, "Scalable and efficient moe training for multitask multilingual models," 2021.
- [44] L. Dong, S. Xu, and B. Xu, "Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5884–5888.
- [45] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [46] J. Vanschoren, "Meta-learning: A survey," *arXiv preprint arXiv:1810.03548*, 2018.
- [47] V. Vapnik, "Statistical learning theory," *NY: Wiley*, vol. 1, p. 2, 1998.
- [48] R. Ardila, M. Branson, K. Davis, M. Kohler, J. Meyer, M. Henretty, R. Morais, L. Saunders, F. Tyers, and G. Weber, "Common voice: A massively-multilingual speech corpus," in *Proceedings of The 12th Language Resources and Evaluation Conference*, 2020, pp. 4218–4222.
- [49] D. Povey, V. Peddinti, D. Galvez, P. Ghahremani, V. Manohar, X. Na, Y. Wang, and S. Khudanpur, "Purely sequence-trained neural networks for asr based on lattice-free mmi," in *Interspeech*, 2016, pp. 2751–2755.
- [50] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, K. N. Goel, M. Hannemann, P. Motlíček, Y. Qian, P. Schwarz, J. Silovský, G. Stemmer, and K. Veselý, "The kaldi speech recognition toolkit," in *Proceedings of ASRU 2011*. IEEE Signal Processing Society, 2011, pp. 1–4. [Online]. Available: <https://www.fit.vut.cz/research/publication/11196>
- [51] V. Peddinti, D. Povey, and S. Khudanpur, "A time delay neural network architecture for efficient modeling of long temporal contexts," in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [52] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2010.
- [53] S. Watanabe, T. Hori, S. Karita, T. Hayashi, J. Nishitoba, Y. Unno, N.-E. Y. Soplin, J. Heymann, M. Wiesner, N. Chen *et al.*, "Espnet: End-to-end speech processing toolkit," *Proc. Interspeech 2018*, pp. 2207–2211, 2018.
- [54] T. Kudo and J. Richardson, "Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing," *EMNLP 2018*, p. 66, 2018.
- [55] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.