

Low-Rank and Sparse Model Merging for Multi-Lingual Speech Recognition and Translation

Anonymous ACL submission

Abstract

Language diversity presents a significant challenge in speech-to-text (S2T) tasks, such as automatic speech recognition and translation. Traditional multitask training approaches aim to address this by jointly optimizing multiple speech recognition and translation tasks across various languages. While models like Whisper, built on these strategies, demonstrate strong performance, they still face issues of high computational cost, language interference, suboptimal training configurations, and limited extensibility. To overcome these challenges, we introduce LoRS-Merging (low-rank and sparse model merging), a novel technique designed to efficiently integrate models trained on different languages or tasks while preserving performance and reducing computational overhead. LoRS-Merging combines low-rank and sparse pruning to retain essential structures while eliminating redundant parameters, mitigating language and task interference, and enhancing extensibility. Experimental results across a range of languages demonstrate that LoRS-Merging reduces the word error rate by 10% and improves BLEU scores by 4% compared to conventional multilingual multitask training baselines. Our findings suggest that model merging, particularly LoRS-Merging, is a scalable and effective complement to traditional multilingual training strategies for S2T applications.

1 Introduction

Language diversity poses a significant challenge in speech-to-text (S2T) tasks, such as automatic speech recognition (ASR) (Prabhavalkar et al., 2023) and speech translation (ST) (Xu et al., 2023). With over 7,000 languages spoken worldwide, developing robust S2T systems that generalise across varied linguistic structures remains a fundamental research goal (Liu and Niehues, 2024; Cheng et al., 2023; Sun et al., 2023; Saif et al., 2024; Wang et al., 2021; Le et al., 2021). The advent of end-to-end

(E2E) models (Chan et al., 2016; Gulati et al., 2020; Barrault et al., 2023) has marked a paradigm shift in S2T tasks, enabling direct mapping from speech to text across multiple languages within a unified framework. A prominent example is Whisper (Radford et al., 2023), an advanced multi-lingual speech model trained on a large-scale, diverse dataset covering multiple languages and tasks. Despite these advances, existing multi-lingual models still encounter significant challenges in scalability, efficiency, and performance trade-offs.

To address these challenges, multi-lingual training strategies (Saif et al., 2024; Xiao et al., 2021; Bai et al., 2018) have been adopted, aiming to enhance model generalisation across languages. These approaches typically rely on joint optimisation of diverse S2T tasks across multiple languages, leveraging shared representations to improve performance. Nevertheless, multi-lingual training is subject to inherent limitations, including substantial training costs, complex model configurations, and limited access to training data across multiple languages and tasks. Moreover, when handling new languages, the training methods typically require training from scratch.

To mitigate these issues, this paper proposes to use model merging (Ilharco et al., 2023; Yang et al., 2024a; Khan et al., 2024) to integrate models trained on different languages or tasks while maintaining performance and reducing computational overhead. Model merging merges the parameters of multiple separate models with different capabilities to build a universal model. With its high flexibility, model merging enables the seamless incorporation of new languages or tasks without the need for retraining the entire model. Additionally, since model merging allows models for different languages or tasks to be trained independently, it can effectively alleviate negative transfer issues (Wang et al., 2019; Zhang et al., 2022b; Wang et al., 2020b) commonly observed in multi-lingual

training. This training independence also enables the use of optimal training configurations for each language or task instead of the unified settings required in multi-lingual training.

Moreover, we propose **Low-Rank and Sparse model Merging** (LoRS-Merging), which uses a low-rank component to capture the compact structure and a sparse component to capture the scattered details in the weights. LoRS-Merging retains effective parts of structure and details while reducing redundant parts to reduce task interference. Specifically, coarse-grained singular value pruning is used to retain the low-rank structure, while fine-grained magnitude pruning is used to remove redundant details. The main contribution of this paper can be summarised as follows.

- We propose LoRS-Merging, a low-rank and sparse model merging method for multi-lingual ASR and speech translation. To the best of our knowledge, LoRS-Merging is the first work that explores model merging for speech models.
- LoRS-Merging exploits the combination of low-rank structure and sparsity of language-specific and task-specific weights in model merging, minimising the parameter redundancy and conflicts as well as providing an efficient way to incorporate new knowledge from a task or language-specialised model.
- Experiments are performed across 10 different languages where LoRS-Merging achieves **10%** relative WER reduction and **4%** relative BLEU increase compared to the multi-lingual multi-task training baseline. Moreover, we show that negative interference largely exists in multi-lingual training and LoRS-Merging alleviates this issue.

2 Related Work

2.1 Multi-Lingual ASR and ST

Multi-lingual speech models inherently face a trade-off between knowledge sharing and negative interference. Early studies adopted hand-picked sub-network sharing strategies, such as language-specific decoders (Dong et al., 2015), attention heads (Zhu et al., 2020), and layer norm/linear transformation (Zhang et al., 2020). Recent research has shifted toward approaches such as mixture-of-experts (Kwon and Chung, 2023; Wang et al., 2023), adapters (Le et al., 2021; Kannan et al., 2019), and pruning (Lu et al., 2022; Lai et al., 2021). To enhance multi-lingual representation

learning, language tokens (Johnson et al., 2017), embeddings (Di Gangi et al., 2019) or output factorizations (Zhang et al., 2023) are introduced to encode language identity, helping the model distinguish between languages.

The more effective approach is to adopt multi-lingual training strategies, such as multi-objective optimisation (Saif et al., 2024; Zhang et al., 2022a), adversarial learning (Xiao et al., 2021), meta learning (Hsu et al., 2020), and reinforcement learning (Bai et al., 2018). Moreover, large-scale pre-training by leveraging massive amounts of multi-lingual and multi-task data enables models to learn robust and transferable representations across languages, e.g. Whisper (Radford et al., 2023), SeamlessM4T (Barrault et al., 2023), and AudioPaLM (Rubenstein et al., 2023). LoRS-Merging, as an efficient post-training method proposed in this paper, further advances multi-lingual ASR and ST based on pre-trained speech models.

2.2 Model Merging

Model merging (Yang et al., 2024a; Khan et al., 2024) is an efficient post-training technique that integrates knowledge from models trained on different domains. One stream of research focuses on the loss landscape geometry (Khan et al., 2024) and studies the linear mode connectivity (LMC) (Frankle et al., 2020; Draxler et al., 2018) property that demonstrates the existence of a linearly connected path between local minima within the same loss basin. Many studies (Nagarajan and Kolter, 2019; Izmailov et al., 2018; Frankle et al., 2020) indicate that if two neural networks share part of their optimisation trajectory, such as different finetuned models from the same pretrained model, they typically satisfy LMC, allowing interpolation without sacrificing accuracy and forming the basis of our model merging method. For local minima in different loss basins, inspired by the permutation invariance (Entezari et al., 2021) of neural networks, neuron alignment techniques (Ainsworth et al., 2022; Singh and Jaggi, 2020; Tatro et al., 2020) can be used to place them into the same basin, thereby reducing merging loss.

Another stream considers the model spaces, including activation spaces and weight spaces. Research on activation spaces seeks to align the output representations or loss of the merged model with those of each single model as closely as possible (Yang et al., 2024b; Wei et al., 2025; Xiong et al., 2024). Studies based on weight spaces aim

to remove redundant parameters or localise effective parameters to resolve task interference. TIES-Merging (Yadav et al., 2024) and DARE (Yu et al., 2024) perform magnitude or random pruning on each single model to significantly remove redundant parameters. TALL-masks (Wang et al., 2024) and Localise-and-Stitch (He et al., 2024) optimise binary masks to localise sparse and effective task-specific parameters. In contrast, LoRS-Merging explores weight space merging by considering not only the detailed parameter redundancy as well as maintaining the effective structure of the weight space via low-rank pruning.

3 Methodology

3.1 Preliminaries

3.1.1 Task Arithmetic

Among diverse model merging methods, Task Arithmetic (TA) (Ilharco et al., 2023) has become a fundamental technique in this field due to its simplicity and effectiveness. TA introduces the concept of "task vector", defined as the delta parameter derived by subtracting pretrained weights from finetuned weights. By performing simple arithmetic operations on task vectors, TA enables task learning, forgetting, and analogising.

Assume that $\theta = \{W_l\}_{l=1}^L$ represents the parameters of the model, where W_l is the weight of l -th layer, and L is the total number of layers. Given a pretrained model θ_0 and a model θ_i finetuned on task t_i , the task vector is computed as $\tau_i = \theta_i - \theta_0$. Multiple task vectors can be summed to form a multi-task model, expressed as $\theta_{\text{merged}} = \theta_0 + \lambda \sum_{i=1}^n \tau_i$, where λ is a scaling coefficient for the task vectors.

3.1.2 Pruning

Given that neural networks are typically over-parameterised and exhibit high redundancy, a considerable number of neurons or connections can be pruned without affecting accuracy (LeCun et al., 1989). In model merging, pruning methods can reduce redundant parameters to mitigate task interference, thereby improving the merging performance.

Magnitude Pruning (MP) is an unstructured pruning method that prunes connections based on the magnitude of parameters as a measure of importance. Specifically, MP prunes the parameters according to a specific ratio p , as follows.

$$M_{ij} = \begin{cases} 1 & \text{if } |w_{ij}| \in \text{top } p\% \\ 0 & \text{o.w.} \end{cases} \quad (1)$$

$$W_{\text{pruned}} = M \odot W \quad (2)$$

where $W, M \in \mathbb{R}^{d \times k}$, and \odot denotes the element-wise multiplication. However, MP only focuses on the redundancy at the parameter level, overlooking the crucial structural information, which may lead to the disruption of the weight structure.

Singular Value Pruning (SVP) is a structured pruning method that removes smaller singular values and their corresponding singular vectors. In particular, SVP retains only the top r singular values while discarding the others.

$$W = U \Sigma V^T \quad (3)$$

$$W_{\text{pruned}} = U_r \Sigma_r V_r^T \quad (4)$$

where $U \in \mathbb{R}^{d \times d}$ and $V \in \mathbb{R}^{k \times k}$ are the left and right singular vector matrices of W , and U_r, V_r denote their first r columns. Although SVP preserves a compact weight structure, its coarse pruning granularity makes it challenging to reduce redundancy at a fine-grained parameter level.

3.2 Model Merging for Speech Models

The model merging process for speech model on S2T tasks with LoRS-Merging as an example is shown in Fig. 1, which comprises four steps. In step 1, a suitable pre-trained speech model is selected. In step 2, for each target language and target task combination, e.g. Catalan ASR, the pre-trained model is finetuned with the task-language-specific data and the delta weight is obtained. In step 3, weight pruning is applied to remove redundant and conflicting delta parameters. In step 4, task arithmetic is applied to combine pruned delta weights into each single merged matrix and hence obtain the merged model.

Model merging allows new language or task knowledge to be integrated into the model in a flexible post-training manner. When a new set of data for a specific language is obtained, model merging incorporates such knowledge by fine-tuning with the new data alone with data-specific configuration, which also releases the burden of requiring other data to avoid catastrophic forgetting. This benefit is thoroughly demonstrated in our experiments.

3.3 Low-Rank and Sparse Model Merging

The weights of neural networks contain information on both structure and details. Structural information is coherent, compact, and coarse-grained, whereas detail information is incoherent, scattered,

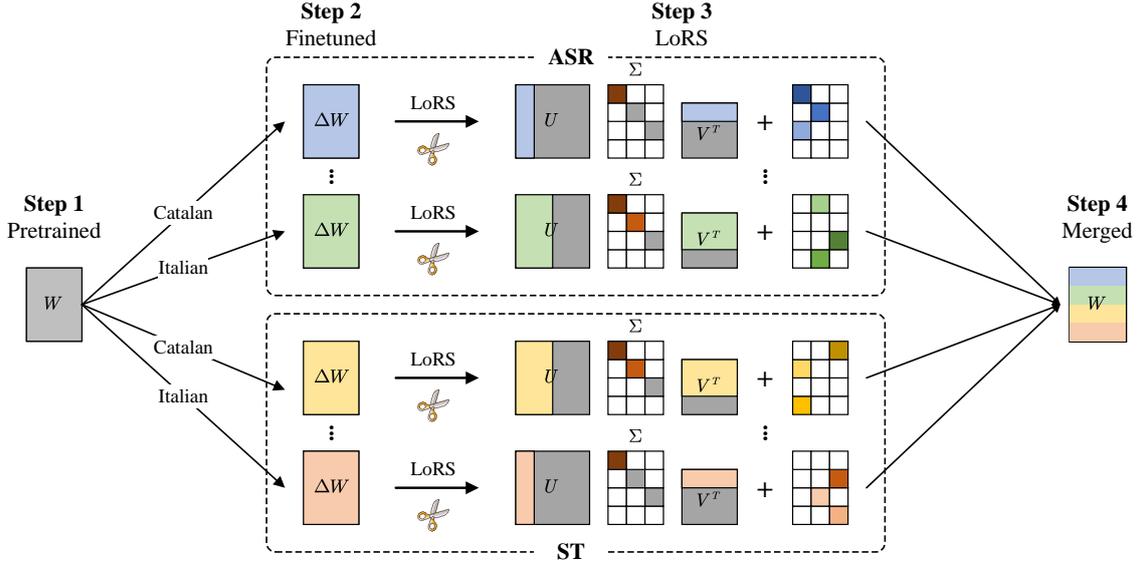


Figure 1: Model merging process with the proposed LoRS-Merging for speech models on multi-lingual ASR and ST tasks. In step 1, a suitable pre-trained speech model is selected. In step 2, the pre-trained model is finetuned with the task-language-specific data. In step 3, apply LoRS to the delta parameters to reduce model redundancy. In step 4, merge the delta parameters to get a multi-lingual and multi-task merged model.

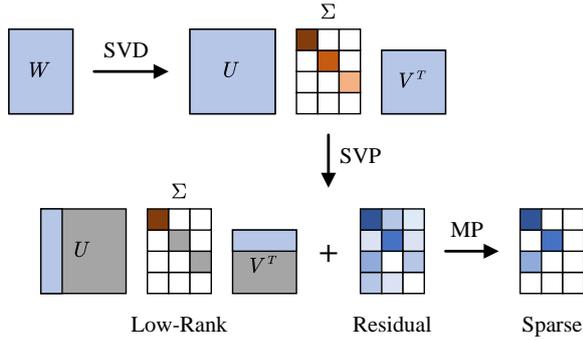


Figure 2: Illustration of LoRS-Merging method in detail. SVD stands for singular value decomposition and SVP for singular value pruning. MP is magnitude pruning operating on residual of the original weight matrix and the low-rank matrix.

and fine-grained. Both structural and detail information include effective and redundant parts. To reduce redundant parts in both the structure and detail aspects of the weights while retaining effective parts, the LoRS-Merging method is introduced as shown in detail in Fig. 2, which exploits the combination of low-rank structure by SVP and sparsity by MP. SVP performs coarse-grained pruning at the structure level, while MP enables fine-grained pruning at the detail level.

In the implementation, we approximate the original weights as the sum of a low-rank component and a sparse component, where the low-rank component captures the compact structure, and the sparse component captures the scattered details,

as shown in Eqn. (5).

$$W \approx L + S \quad (5)$$

where L represents the low-rank component, and S represents the sparse component. Specifically, L is the low-rank matrix obtained by retaining the top r singular values and their corresponding singular vectors from W :

$$L = U_r \Sigma_r V_r^T \quad (6)$$

and S is the sparse matrix obtained by performing MP on the residual of W and L :

$$S = M \odot (W - L) \quad (7)$$

To simplify the description, we refer to this entire process as $\text{LoRS}(\cdot)$. In this manner, SVP decouples the structure and details of the weight, preserving a compact structure while allowing fine-grained MP to remove redundant parts in the details.

For each model finetuned on single specific language or task data, we apply $\text{LoRS}(\cdot)$ to its task vector as a preprocessing step to reduce language or task interference in model merging. A multi-lingual or multi-task model can be achieved through simple merging, expressed as:

$$\theta_{\text{merged}} = \theta_0 + \lambda \sum_{i=1}^n \text{LoRS}(\tau_i) \quad (8)$$

Compared to multi-lingual or multi-task training methods, model merging is a simpler and more efficient approach, enabling the seamless incorporation of new languages or tasks without the need for

retraining. Additionally, due to its training independence, it mitigates negative transfer and provides optimal training configurations for each language or task to improve performance.

4 Experimental Setup

4.1 Data

CoVoST-2 (Wang et al., 2020a) is a large-scale multi-lingual ST corpus based on Common Voice. It covers translations from English into 15 languages and from 21 languages into English, with a total of 2,880 hours of speech from 78k speakers. We selected 5 high-resource languages and 5 low-resource languages as two language sets to investigate their ASR tasks and the from X to English ST tasks. The high-resource language set includes Catalan (ca), German (de), Spanish (es), French (fr), and Italian (it), while the low-resource language set includes Indonesian (id), Dutch (nl), Portuguese (pt), Russian (ru), and Swedish (sv). Due to the more abundant data in the high-resource language set, our main experimental results are obtained on the high-resource language set, while the low-resource language set serves as an auxiliary validation set. To balance the amount of data across different languages, we fixed the duration of training data for each language, with 5 hours for the high-resource language set and 1 hour for the low-resource language set. The dev and test sets of both language sets are 1 hour in duration.

4.2 Model and Training Specifications

Whisper (Radford et al., 2023) is a general-purpose multi-lingual ASR and ST model, a Transformer-based model trained on 680k hours of diverse audio. We chose the small version as the foundation model for the experiments because it achieves a good balance between performance and cost. It has 244 million parameters, with the encoder and decoder each consisting of 12 Transformer blocks. The weight matrices of the attention layers are all 768 by 768, and the MLP layers are 768 by 3072.

For each language-specific or task-specific fine-tuned model, we use a different, optimal learning rate for each during training, and these models are subsequently used for model merging. Finetuning involves updating all parameters. We choose Adam as the optimiser, set the batch size to 8, the accumulation iterations to 4, and train for 10 epochs. We also select the proportions of low-rank parameters retained by SVP from {1%, 2%, 3%, 5%}

and sparse parameters retained by MP from {10%, 20%, 40%, 60%}. The beam size for decoding is set to 20 across all languages and tasks. We use Sclite and SacreBLEU tools to score the ASR and ST results, respectively. See Appendix A for more details on hyper-parameter settings. Our experiments are performed on a single RTX 4090 GPU where training on one language and one task with 5 hours of speech data requires 1 hour.

4.3 Baseline and Merging Methods

We use **Multi-lingual and multi-task training** as the baseline for comparison with model merging methods, where training is conducted on data mixed from both multi-lingual and multi-task sets. To ensure a fair comparison, the same amount of training data is used from each language and each task. Note that for 5 different languages with both ASR and ST tasks, multi-lingual and multi-task training is performed on 10 times more data and hence 10 times more computational resources.

In addition to LoRS-Merging, we investigate the following model merging methods:

Weight Averaging (WA) merges multiple single models by unweighted averaging of their weights, $\theta_{\text{merged}} = \frac{1}{n} \sum_{i=1}^n \theta_i$.

Task Arithmetic (TA) uses a scaling factor to weight multiple task vectors estimated on a small development set, $\theta_{\text{merged}} = \theta_0 + \lambda \sum_{i=1}^n \tau_i$.

MP-Merging performs fine-grained magnitude pruning on task vectors to reduce redundancy at the detail level, $\theta_{\text{merged}} = \theta_0 + \lambda \sum_{i=1}^n \text{MP}(\tau_i)$.

SVP-Merging performs coarse-grained singular value pruning on task vectors to reduce redundancy at the structure level, $\theta_{\text{merged}} = \theta_0 + \lambda \sum_{i=1}^n \text{SVP}(\tau_i)$. (see Section 3.1.2).

Moreover, we compare methods against the performance of fine-tuning on each language-task combination. This is the topline of all merging methods since the model is completely adapted to a specific language for a specific task with optimised configurations and without any language conflicts.

5 Evaluation Results and Analysis

5.1 Multi-Lingual Model Merging

First, we investigate the merging of finetuned models for different languages on the same task, which corresponds to *multi-lingual single-task* learning.

Language knowledge interference yields imbalanced improvements: Table 1 shows the multi-lingual results of the ASR task with the high-

Table 1: Multi-lingual ASR model merging. Finetuned is the topline where the model is finetuned on each language independently, and Avg. averages WER directly.

System	WER↓					Avg.
	ca	de	es	fr	it	
Pretrained	20.6	19.6	14.7	24.5	19.4	19.88
Finetuned	19.5	19.7	14.4	22.1	19.2	19.05
Multi-lingual training	17.1	21.8	15.1	22.6	21.9	19.69
Weight Averaging	19.1	19.1	14.2	24.5	20.3	19.55
Task Arithmetic	19.1	18.8	13.9	24.0	19.8	19.23
MP-Merging	19.4	19.3	14.0	23.8	18.1	19.03
SVP-Merging	19.5	19.3	14.2	23.6	18.4	19.11
LoRS-Merging	19.0	18.8	13.9	23.5	18.5	18.85

Table 2: Multi-lingual ST model merging. Finetuned is the topline where the model is finetuned on each language independently, and Avg. averages BLEU directly.

System	BLEU↑					Avg.
	ca	de	es	fr	it	
Pretrained	21.1	24.1	28.6	26.8	26.8	25.48
Finetuned	22.6	24.6	29.2	27.2	27.3	26.18
Multi-lingual training	21.4	24.4	28.8	26.8	27.2	25.72
Weight Averaging	22.3	24.1	28.6	27.2	26.9	25.82
Task Arithmetic	22.1	24.3	28.9	27.3	26.8	25.88
MP-Merging	22.1	24.7	28.9	27.3	26.9	25.98
SVP-Merging	22.1	24.7	29.0	27.4	26.8	26.00
LoRS-Merging	22.2	24.8	29.0	27.5	26.9	26.08

resource language set. On average, multi-lingual training slightly improves the pretrained model but significantly underperforms the finetuned models and merging methods. This may be due to negative interference between the knowledge of different languages, leading to gradient conflicts during training (Wang et al., 2020b). From a per-language perspective, it is observed that ca and fr achieve the largest improvements during fine-tuning while still showing significant improvements in multi-lingual training, whereas languages with smaller improvements during finetuning exhibit a substantial performance drop in multi-lingual training, even worse than the pretrained model. This indicates a strong language conflict in multi-lingual training, with ca and fr dominating. Additionally, we observe that the optimal learning rates for finetuned models vary significantly across languages (see Appendix A), while the unified learning rate configuration required by multi-lingual training prevents each language from reaching its optimal performance.

Model merging mitigates language conflicts: In contrast, model merging methods show obvious improvements across almost all languages, demonstrating reduced conflict and better stability.

Table 3: Multi-task model merging performed on each language independently and WER/BLEU scores are averaged across languages. Finetuned is the topline where the model is finetuned on each language and task combination independently. Per-language results are shown in Appendix C.

System	Avg. WER↓	Avg. BLEU↑
Pretrained	19.88	25.48
Finetuned	19.05	26.18
Multi-task training	19.00	25.90
Weight Averaging	18.84	26.18
Task Arithmetic	18.76	26.30
MP-Merging	18.62	26.40
SVP-Merging	18.72	26.38
LoRS-Merging	18.45	26.48

Among model merging methods, TA outperforms WA due to its flexible scaling factor. Both MP-Merging and SVP-Merging further improve the performance of TA by reducing redundancy, and MP-Merging slightly outperforms SVP-Merging due to its finer-grained pruning. Combining the advantages of SVP and MP, LoRS-Merging achieves the best performance.

Table 2 provides the multi-lingual results on ST task with the high-resource language set. The main conclusion is consistent with the ASR task: model merging methods still significantly outperform multi-lingual training, with LoRS-Merging achieving the best performance.

5.2 Multi-Task Model Merging

Next, we merge finetuned models for different tasks (ASR and ST) with the same language which corresponds to *multi-task single-language* learning.

ASR and ST tasks for the same language can mutually benefit from each other: Table 3 presents the multi-task results with the high-resource language set. In general, multi-task training performs similarly to finetuned models on ASR but is a lot worse on ST. This is likely due to the substantial differences in optimal hyper-parameter configurations between the two tasks. Model merging methods clearly outperform finetuned models, which not only demonstrates their effectiveness but also shows the mutual benefits between ASR and ST. In terms of performance gains, the improvement in ASR is greater than in ST. We attribute this to the fact that ASR is inherently simpler than ST and can be viewed as a step in the ST task. Furthermore, as before, model merging methods combined with pruning further improve performance, and the

Table 4: Multi-lingual multi-task model merging. Finetuned is the topline where the model is finetuned on each language and task combination independently, and Avg. averages WER or BLEU scores directly.

System	WER↓						BLEU↑					
	ca	de	es	fr	it	Avg.	ca	de	es	fr	it	Avg.
Pretrained	20.6	19.6	14.7	24.5	19.4	19.88	21.1	24.1	28.6	26.8	26.8	25.48
Finetuned	19.5	19.7	14.4	22.1	19.2	19.05	22.6	24.6	29.2	27.2	27.3	26.18
ML and MT training	20.5	19.7	14.6	24.5	19.4	19.86	21.3	24.3	28.3	27.1	26.9	25.58
ML and MT Task Arithmetic	18.9	19.2	14.1	23.7	18.4	18.96	22.2	24.4	29.0	27.3	26.9	25.96
ML and MT LoRS-Merging	18.7	19.1	14.0	23.8	18.0	18.82	22.2	24.8	29.0	27.5	27.0	26.10
MT training	17.0	19.7	14.4	24.2	19.4	19.00	22.3	24.6	28.7	27.0	26.9	25.90
↔ + ML Task Arithmetic	18.1	19.0	14.2	24.5	20.6	19.37	22.7	24.7	28.6	27.3	26.5	25.96
↔ + ML LoRS-Merging	18.1	19.0	14.1	24.2	20.3	19.23	22.4	24.5	29.1	27.6	26.7	26.06
ML training	17.1	21.8	15.1	22.6	21.9	19.69	21.4	24.4	28.8	26.8	27.2	25.72
↔ + MT Task Arithmetic	17.1	18.5	13.3	22.7	18.0	18.00	22.6	25.0	29.2	27.5	26.9	26.24
↔ + MT LoRS-Merging	16.9	18.3	13.3	22.4	17.8	17.82	22.8	25.2	29.3	27.6	27.0	26.38

proposed LoRS-Merging achieves the best performance across the table.

5.3 Multi-Lingual Multi-Task Model Merging

Then, we investigate the merging of finetuned models for both different languages and tasks, which correspond to *multi-lingual (ML) and multi-task (MT)* learning. Specifically, we explore 4 different training and merging settings:

ML and MT training: Fine-tuning on all languages and both tasks jointly.

ML and MT merging: Fine-tuning on each language for each task separately and merging all.

MT training and ML merging: Fine-tuning both tasks jointly for each language, and merging models from different languages.

ML training and MT merging: Fine-tuning on all languages jointly for each task, and merging models from different tasks.

Table 4 displays the multi-lingual and multi-task results with the high-resource language set. Multi-lingual and multi-task training shows little improvement over the pretrained model, due to negative interference during training and the use of a unified training configuration for all languages and tasks. Nevertheless, the performance of multi-lingual and multi-task merging is on par with that of finetuned models, further underscoring the superiority of model merging. As a result, LoRS-Merging achieved the best performance when performing ML training followed by MT merging, which consistently outperforms Task Arithmetic. Overall, **10%** relative WER reduction and **4%** relative BLEU increase are achieved using LoRS-Merging compared to ML and MT training base-

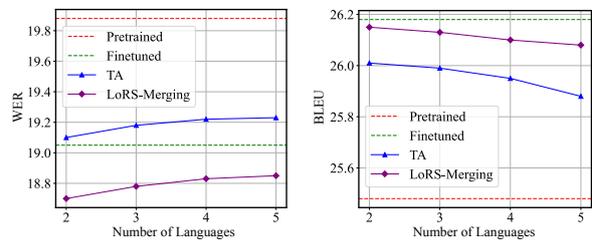


Figure 3: WER and BLEU against the number of languages. Performance is averaged across all languages and all training runs of language combinations.

line. We provide additional experiments on a set of low-resource languages in Appendix B to demonstrate the robustness and generalizability of model merging and LoRS-Merging.

5.4 Effect of Numbers of Languages

To further demonstrate the robustness of LoRS-Merging to language selection, experiments are performed using different numbers of languages. Figure 3 shows the average performance across all languages and all training runs with possible combinations of 2, 3, 4 or 5 languages.

LoRS-Merging improvements are consistent across different numbers of languages: As the number of languages increases, the performance of both TA and LoRS-Merging degrades due to negative interference between languages. LoRS-Merging consistently outperforms TA in both ASR and ST tasks, and even surpasses the finetuned models in the ASR task. This is likely due to LoRS-Merging further reducing model redundancy, therefore alleviating negative interference. Additionally, we observe that the optimal learning rate

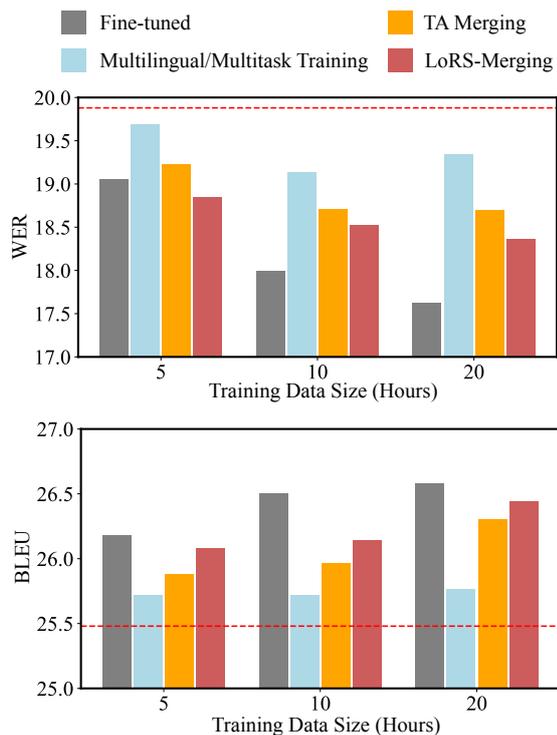


Figure 4: Performance variation against different training data sizes (number of hours for each language) on ASR (top) and ST (bottom) tasks.

for the finetuned ASR model is significantly larger compared to the ST task. This may lead to overfitting, whereas LoRS-Merging improves model generalization through model merging while reducing language interference, thus outperforming the finetuned models for the ASR task.

5.5 Effect of Language Data Scale

We then demonstrate the robustness of merging methods to different training data sizes for both tasks. Fig. 4 shows the WER (top) and BLEU (bottom) scores for ASR and ST at different data scales, respectively. As the data scale increases, the performance of multi-lingual training does not always improve. This may be because the pretrained model already performs well, and the significant language interference and conflict in multi-lingual training hinder the effective improvement of multi-language performance. Furthermore, the performance loss of model merging increases with data scale, compared to finetuned models. It can be explained by the fact that larger training data tends to increase the divergence in the optimisation trajectories of different finetuned models, resulting in the breakdown of linear mode connectivity, which leads to a greater performance loss. Moreover, LoRS-Merging still achieves obvious and stable

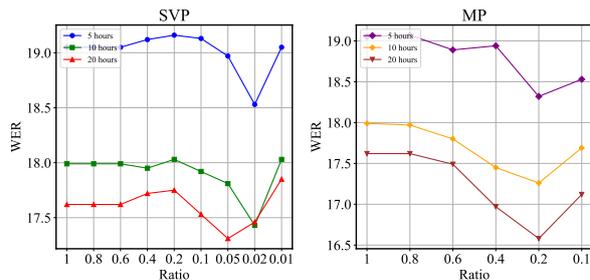


Figure 5: Model performance against the retain ratio (1 means to retain all weights and 0 means to prune all weights) in SVP (left) and MP (right) for ASR finetuned models. Three different training data sizes are used.

improvement compared to TA.

5.6 Analysis of Model Redundancy

Furthermore, we justify the necessity of SVP and MP to remove model redundancy by showing the model performance against the pruning ratio of finetuned models for ASR as shown in Fig. 5. As shown, both SVP and MP significantly improve the performance of finetuned models, indicating the presence of substantial redundancy in the structure and details of the finetuned models, respectively. The model performance reaches the best at a high pruning level, indicating that the redundancy is particularly large for ASR. We observed a much smaller redundancy in ST, which also explains the observation that LoRS-Merging achieves more salient improvement on ASR than ST. Moreover, redundancy increases with training data, possibly due to the accumulation of gradient noise during training. MP achieves greater performance gains than SVP, indicating more redundancy at the detail level, which is better addressed by fine-grained MP.

6 Conclusion

This paper explores model merging for multi-lingual ASR and ST on pre-trained speech models and proposes the LoRS-Merging approach. LoRS-Merging combines low-rank and sparse pruning to retain essential structures, eliminate redundant parameters and mitigate language and task interference. Experiments across five languages show that LoRS-Merging effectively alleviates language interference, and achieves a 10% relative WER reduction and a 4% relative BLEU score improvement compared to multi-lingual multi-task training baselines.

7 Limitations

There are three main limitations of this work. First, as a common limitation of all model merging methods, the same model structure is required across all tasks and languages. This is less of a concern under the current trend of using the same Transformer structure, but methods need to be developed in the future to accommodate subtle structural differences. Second, reasonably-sized training sets are required for each language, and low-resource languages may suffer from reduced improvements. Third, this work mainly explores the two most popular S2T tasks. Other possible tasks can be explored in future work, including spoken language understanding and speaker adaptation.

References

- Samuel Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. 2022. Git re-basin: Merging models modulo permutation symmetries. In *The Eleventh International Conference on Learning Representations*.
- He Bai, Yu Zhou, Jiajun Zhang, Liang Zhao, Mei-Yuh Hwang, and Chengqing Zong. 2018. Source-critical reinforcement learning for transferring spoken language understanding to a new language. *arXiv preprint arXiv:1808.06167*.
- Loïc Barrault, Yu-An Chung, Mariano Cora Meglioli, David Dale, Ning Dong, Paul-Ambroise Duquenne, Hady Elsahar, Hongyu Gong, Kevin Heffernan, John Hoffman, et al. 2023. Seamless4t-massively multilingual & multimodal machine translation. *arXiv preprint arXiv:2308.11596*.
- William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals. 2016. Listen, attend and spell: A neural network for large vocabulary conversational speech recognition. In *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 4960–4964. IEEE.
- Yong Cheng, Yu Zhang, Melvin Johnson, Wolfgang Macherey, and Ankur Bapna. 2023. Mu2slam: multitask, multilingual speech and language models. In *Proceedings of the 40th International Conference on Machine Learning*, pages 5504–5520.
- Mattia A Di Gangi, Matteo Negri, and Marco Turchi. 2019. One-to-many multilingual end-to-end speech translation. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 585–592. IEEE.
- Daxiang Dong, Hua Wu, Wei He, Dianhai Yu, and Haifeng Wang. 2015. Multi-task learning for multiple language translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational*

Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1723–1732.

- Felix Draxler, Kambis Veschgini, Manfred Salmhofer, and Fred Hamprecht. 2018. Essentially no barriers in neural network energy landscape. In *International conference on machine learning*, pages 1309–1318. PMLR.
- Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. 2021. The role of permutation invariance in linear mode connectivity of neural networks. In *International Conference on Learning Representations*.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. 2020. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pages 3259–3269. PMLR.
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. 2020. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*.
- Yifei He, Yuzheng Hu, Yong Lin, Tong Zhang, and Han Zhao. 2024. Localize-and-stitch: Efficient model merging via sparse task arithmetic. *arXiv preprint arXiv:2408.13656*.
- Jui-Yang Hsu, Yuan-Jui Chen, and Hung-yi Lee. 2020. Meta learning for end-to-end low-resource speech recognition. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7844–7848. IEEE.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. 2018. Averaging weights leads to wider optima and better generalization. In *34th Conference on Uncertainty in Artificial Intelligence 2018, UAI 2018*, pages 876–885. Association For Uncertainty in Artificial Intelligence (AUAI).
- Melvin Johnson, Mike Schuster, Quoc V Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, et al. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351.
- Anjuli Kannan, Arindrima Datta, Tara N Sainath, Eugene Weinstein, Bhuvana Ramabhadran, Yonghui Wu, Ankur Bapna, Zhifeng Chen, and Seungji Lee. 2019. Large-scale multilingual speech recognition

702	with a streaming end-to-end model. <i>Interspeech</i>	Paul K Rubenstein, Chulayuth Asawaroengchai,	757
703	2019.	Duc Dung Nguyen, Ankur Bapna, Zalán Borsos,	758
704	Arham Khan, Todd Nief, Nathaniel Hudson, Mansi	Félix de Chaumont Quitry, Peter Chen, Dalia El	759
705	Sakarvadia, Daniel Grzenda, Aswathy Ajith, Jordan	Badawy, Wei Han, Eugene Kharitonov, et al. 2023.	760
706	Pettyjohn, Kyle Chard, and Ian Foster. 2024. Sok:	Audiopalm: A large language model that can speak	761
707	On finding common ground in loss landscapes using	and listen. <i>arXiv preprint arXiv:2306.12925</i> .	762
708	deep model merging techniques. <i>arXiv preprint</i>		
709	<i>arXiv:2410.12927</i> .	AFM Saif, Lisha Chen, Xiaodong Cui, Songtao Lu,	763
710	Yoohwan Kwon and Soo-Whan Chung. 2023. Mole:	Brian Kingsbury, and Tianyi Chen. 2024. M2asr:	764
711	Mixture of language experts for multi-lingual auto-	Multilingual multi-task automatic speech recognition	765
712	matic speech recognition. In <i>ICASSP 2023-2023</i>	via multi-objective optimization. In <i>Interspeech</i> , vol-	766
713	<i>IEEE International Conference on Acoustics, Speech</i>	ume 2024, pages 1240–1244.	767
714	<i>and Signal Processing (ICASSP)</i> , pages 1–5. IEEE.	Sidak Pal Singh and Martin Jaggi. 2020. Model fusion	768
715	Cheng-I Jeff Lai, Yang Zhang, Alexander H Liu, Shiyu	via optimal transport. <i>Advances in Neural Informa-</i>	769
716	Chang, Yi-Lun Liao, Yung-Sung Chuang, Kaizhi	<i>tion Processing Systems</i> , 33:22045–22055.	770
717	Qian, Sameer Khurana, David Cox, and Jim Glass.	Haoran Sun, Xiaohu Zhao, Yikun Lei, Deyi Xiong, et al.	771
718	2021. Parp: Prune, adjust and re-prune for self-	2023. Towards a deep understanding of multilingual	772
719	supervised speech recognition. <i>Advances in Neural</i>	end-to-end speech translation. In <i>The 2023 Confer-</i>	773
720	<i>Information Processing Systems</i> , 34:21256–21272.	<i>ence on Empirical Methods in Natural Language</i>	774
721	Hang Le, Juan Pino, Changhan Wang, Jiatao Gu, Didier	<i>Processing</i> .	775
722	Schwab, and Laurent Besacier. 2021. Lightweight	Norman Tatro, Pin-Yu Chen, Payel Das, Igor Melnyk,	776
723	adapter tuning for multilingual speech translation.	Prasanna Sattigeri, and Rongjie Lai. 2020. Optimiz-	777
724	In <i>Proceedings of the 59th Annual Meeting of the</i>	ing mode connectivity via neuron alignment. <i>Ad-</i>	778
725	<i>Association for Computational Linguistics and the</i>	<i>vances in Neural Information Processing Systems</i> ,	779
726	<i>11th International Joint Conference on Natural Lan-</i>	33:15300–15311.	780
727	<i>guage Processing (Volume 2: Short Papers)</i> , pages	Changhan Wang, Morgane Riviere, Ann Lee, Anne Wu,	781
728	817–824.	Chaitanya Talnikar, Daniel Haziza, Mary Williamson,	782
729	Yann LeCun, John Denker, and Sara Solla. 1989. Opti-	Juan Pino, and Emmanuel Dupoux. 2021. Voxpop-	783
730	mal brain damage. <i>Advances in neural information</i>	uli: A large-scale multilingual speech corpus for rep-	784
731	<i>processing systems</i> , 2.	resentation learning, semi-supervised learning and	785
732	Danni Liu and Jan Niehues. 2024. Recent highlights in	interpretation. <i>arXiv preprint arXiv:2101.00390</i> .	786
733	multilingual and multimodal speech translation. In	Changhan Wang, Anne Wu, and Juan Pino. 2020a. Cov-	787
734	<i>Proceedings of the 21st International Conference on</i>	ost 2 and massively multilingual speech-to-text trans-	788
735	<i>Spoken Language Translation (IWSLT 2024)</i> , pages	lation. <i>arXiv preprint arXiv:2007.10310</i> .	789
736	235–253.	Ke Wang, Nikolaos Dimitriadis, Guillermo Ortiz-	790
737	Yizhou Lu, Mingkun Huang, Xinghua Qu, Pengfei	Jimenez, François Fleuret, and Pascal Frossard. 2024.	791
738	Wei, and Zejun Ma. 2022. Language adaptive	Localizing task information for improved model	792
739	cross-lingual speech representation learning with	merging and compression. In <i>Forty-first Interna-</i>	793
740	sparse sharing sub-networks. In <i>ICASSP 2022-2022</i>	<i>tional Conference on Machine Learning</i> .	794
741	<i>IEEE International Conference on Acoustics, Speech</i>	Wenxuan Wang, Guodong Ma, Yuke Li, and Binbin	795
742	<i>and Signal Processing (ICASSP)</i> , pages 6882–6886.	Du. 2023. Language-routing mixture of experts for	796
743	IEEE.	multilingual and code-switching speech recognition.	797
744	Vaishnavh Nagarajan and J Zico Kolter. 2019. Uniform	<i>arXiv preprint arXiv:2307.05956</i> .	798
745	convergence may be unable to explain generalization	Zirui Wang, Zihang Dai, Barnabás Póczos, and Jaime	799
746	in deep learning. <i>Advances in Neural Information</i>	Carbonell. 2019. Characterizing and avoiding nega-	800
747	<i>Processing Systems</i> , 32.	tive transfer. In <i>Proceedings of the IEEE/CVF con-</i>	801
748	Rohit Prabhavalkar, Takaaki Hori, Tara N Sainath, Ralf	<i>ference on computer vision and pattern recognition</i> ,	802
749	Schlüter, and Shinji Watanabe. 2023. End-to-end	pages 11293–11302.	803
750	speech recognition: A survey. <i>IEEE/ACM Transac-</i>	Zirui Wang, Zachary C Lipton, and Yulia Tsvetkov.	804
751	<i>tions on Audio, Speech, and Language Processing</i> .	2020b. On negative interference in multilingual mod-	805
752	Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-	els: Findings and a meta-learning treatment. <i>arXiv</i>	806
753	man, Christine McLeavey, and Ilya Sutskever. 2023.	<i>preprint arXiv:2010.03017</i> .	807
754	Robust speech recognition via large-scale weak su-	Yongxian Wei, Anke Tang, Li Shen, Feng Xiong, Chun	808
755	pervision. In <i>International conference on machine</i>	Yuan, and Xiaochun Cao. 2025. Modeling multi-	809
756	<i>learning</i> , pages 28492–28518. PMLR.	task model merging as adaptive projective gradient	810
		descent. <i>arXiv preprint arXiv:2501.01230</i> .	811

812 Yubei Xiao, Ke Gong, Pan Zhou, Guolin Zheng, Xiao-
813 dan Liang, and Liang Lin. 2021. Adversarial meta
814 sampling for multilingual low-resource speech recog-
815 nition. In *Proceedings of the AAAI Conference on*
816 *Artificial Intelligence*, pages 14112–14120.

817 Feng Xiong, Runxi Cheng, Wang Chen, Zhanqiu Zhang,
818 Yiwen Guo, Chun Yuan, and Ruifeng Xu. 2024.
819 Multi-task model merging via adaptive weight disen-
820 tanglement. *arXiv preprint arXiv:2411.18729*.

821 Chen Xu, Rong Ye, Qianqian Dong, Chengqi Zhao, Tom
822 Ko, Mingxuan Wang, Tong Xiao, and Jingbo Zhu.
823 2023. Recent advances in direct speech-to-text trans-
824 lation. In *Proceedings of the Thirty-Second Inter-
825 national Joint Conference on Artificial Intelligence*,
826 pages 6796–6804.

827 Prateek Yadav, Derek Tam, Leshem Choshen, Colin A
828 Raffel, and Mohit Bansal. 2024. Ties-merging: Re-
829 solving interference when merging models. *Ad-
830 vances in Neural Information Processing Systems*,
831 36.

832 Enneng Yang, Li Shen, Guibing Guo, Xingwei Wang,
833 Xiaochun Cao, Jie Zhang, and Dacheng Tao. 2024a.
834 Model merging in llms, mllms, and beyond: Meth-
835 ods, theories, applications and opportunities. *arXiv*
836 *preprint arXiv:2408.07666*.

837 Enneng Yang, Li Shen, Zhenyi Wang, Guibing Guo,
838 Xiaojun Chen, Xingwei Wang, and Dacheng Tao.
839 2024b. Representation surgery for multi-task model
840 merging. In *Forty-first International Conference on*
841 *Machine Learning*.

842 Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin
843 Li. 2024. Language models are super mario: Absorb-
844 ing abilities from homologous models as a free lunch.
845 In *Forty-first International Conference on Machine*
846 *Learning*.

847 Biao Zhang, Philip Williams, Ivan Titov, and Rico
848 Sennrich. 2020. Improving massively multilingual
849 neural machine translation and zero-shot translation.
850 *arXiv preprint arXiv:2004.11867*.

851 Chao Zhang, Bo Li, Tara Sainath, Trevor Strohman,
852 Sepand Mavandadi, Shuo yiin Chang, and Parisa
853 Haghani. 2022a. Streaming end-to-end multilingual
854 speech recognition with joint language identification.
855 In *Interspeech*.

856 Chao Zhang, Bo Li, Tara N. Sainath, Trevor Strohman,
857 and Shuo yiin Chang. 2023. Uml: A universal mono-
858 lingual output layer for multilingual asr. In *ICASSP*.

859 Wen Zhang, Lingfei Deng, Lei Zhang, and Dongrui Wu.
860 2022b. A survey on negative transfer. *IEEE/CAA*
861 *Journal of Automatica Sinica*, 10(2):305–329.

862 Yun Zhu, Parisa Haghani, Anshuman Tripathi, Bhu-
863 vana Ramabhadran, Brian Farris, Hainan Xu, Han
864 Lu, Hasim Sak, Isabel Leal, Neeraj Gaur, et al. 2020.
865 Multilingual speech recognition with self-attention
866 structured parameterization. In *INTERSPEECH*,
867 pages 4741–4745.

868
869
870
871

A Hyper-parameter Details

The detailed hyper-parameter settings for each language are shown in Table 5 for ASR and Table 6 for ST, respectively.

Table 5: ASR hyper-parameters for high-resource languages.

System	ASR				
	ca	de	es	fr	it
Finetuned					
learning rate	1×10^{-6}	5×10^{-8}	1×10^{-7}	1×10^{-6}	5×10^{-6}
Multi-lingual training					
learning rate	1×10^{-5}				
Task Arithmetic					
scaling factor λ	0.15				
LoRS-Merging					
scaling factor λ	0.15				
SVP ratio r	5%	3%	2%	1%	1%
MP ratio p	40%	60%	40%	10%	10%

Table 6: ST hyper-parameters for high-resource languages.

System	ST				
	ca	de	es	fr	it
Finetuned					
learning rate	1×10^{-6}	2×10^{-8}	2×10^{-8}	5×10^{-8}	5×10^{-8}
Multi-lingual training					
learning rate	5×10^{-9}				
Task Arithmetic					
scaling factor λ	0.15				
LoRS-Merging					
scaling factor λ	0.15				
SVP ratio r	5%	3%	5%	2%	1%
MP ratio p	60%	40%	20%	20%	20%

872
873
874
875
876

B Results of Low-Resource Language Set

The results of the low-resource language set are shown in this section. Specifically, Table 7 and 8 show the multi-lingual single task training and merging for ASR and ST respectively.

Table 7: Multi-lingual ASR model merging. Finetuned is the topline where the model is finetuned on each language independently, and Avg. averages WER directly.

System	WER↓					Avg.
	id	nl	pt	ru	sv	
Pretrained	16.9	16.0	10.1	17.1	17.1	15.43
Finetuned	15.0	14.8	9.7	16.8	14.7	14.20
Multi-lingual training	16.7	15.5	10.0	17.0	16.6	15.14
Weight Averaging	15.7	15.2	10.1	17.1	15.8	14.77
Task Arithmetic	15.7	15.1	9.9	17.0	15.8	14.69
MP-Merging	15.7	15.1	10.0	16.7	15.7	14.63
SVP-Merging	15.7	15.1	9.9	16.9	15.7	14.65
LoRS-Merging	15.7	15.1	9.7	16.8	15.6	14.57

Table 8: Multi-lingual ST model merging. Finetuned is the topline where the model is finetuned on each language independently, and Avg. averages BLEU directly.

System	BLEU↑					Avg.
	id	nl	pt	ru	sv	
Pretrained	32.5	31.6	43.3	35.5	32.1	35.00
Finetuned	35.2	34.0	43.8	36.7	37.6	37.46
Multi-lingual training	32.3	33.2	43.5	35.4	34.3	35.74
Weight Averaging	33.6	32.2	43.2	35.3	34.2	35.70
Task Arithmetic	33.9	32.8	43.1	35.5	34.3	35.92
MP-Merging	33.8	32.8	43.5	35.8	34.0	35.98
SVP-Merging	33.6	32.6	43.4	35.6	34.3	35.90
LoRS-Merging	33.9	32.8	43.2	35.9	34.5	36.06

Then, Table 9 shows the uni-lingual multi-task training and merging performance (c.f. compare to 3 for high-resource languages).

Last, Table 10 shows the results of multi-lingual and multi-task training and merging results for low-resource languages (compare to Table 4 for high-resource languages.). LoRS-Merging achieved the best performance across all merging and training methods in all tables.

C Detailed Results on Multi-task merging

Detailed per-language results of Table 3 are shown in Table 11.

877
878
879
880
881
882
883
884
885
886
887
888

Table 9: Multi-task model merging. Finetuned is the topline where the model is finetuned on each language and task combination independently, and Avg. averages WER or BLEU scores directly.

System	WER↓						BLEU↑					
	id	nl	pt	ru	sv	Avg.	id	nl	pt	ru	sv	Avg.
Pretrained	16.9	16.0	10.1	17.1	17.1	15.43	32.5	31.6	43.3	35.5	32.1	35.00
Finetuned	15.0	14.8	9.7	16.8	14.7	14.20	35.2	34.0	43.8	36.7	37.6	37.46
Multi-task training	15.4	15.0	9.3	16.6	14.3	14.12	35.3	33.7	43.6	36.2	35.8	36.92
Weight Averaging	14.7	14.9	9.3	16.6	13.8	13.88	35.4	33.9	44.1	36.3	35.9	37.12
Task Arithmetic	14.6	14.9	9.3	16.5	14.0	13.88	35.3	33.8	44.3	36.1	36.4	37.18
MP-Merging	14.4	14.7	9.4	16.5	13.8	13.78	35.7	33.9	44.3	36.1	36.1	37.22
SVP-Merging	14.6	14.8	9.2	16.4	13.9	13.80	35.3	33.9	44.3	36.2	36.3	37.20
LoRS-Merging	14.4	14.7	9.2	16.4	13.8	13.72	35.6	33.9	44.3	36.3	36.4	37.30

Table 10: Multi-lingual multi-task model merging. Finetuned is the topline where the model is finetuned on each language and task combination independently, and Avg. averages WER or BLEU scores directly.

System	WER↓						BLEU↑					
	id	nl	pt	ru	sv	Avg.	id	nl	pt	ru	sv	Avg.
Pretrained	16.9	16.0	10.1	17.1	17.1	15.43	32.5	31.6	43.3	35.5	32.1	35.00
Finetuned	15.0	14.8	9.7	16.8	14.7	14.20	35.2	34.0	43.8	36.7	37.6	37.46
ML and MT training	16.9	15.7	9.6	17.0	16.3	15.08	32.8	32.9	43.3	35.4	32.6	35.40
ML and MT Task Arithmetic	16.4	15.5	9.6	16.8	15.7	14.79	33.7	33.1	43.2	35.7	34.9	36.12
ML and MT LoRS-Merging	16.1	15.5	9.5	16.8	15.7	14.72	33.7	33.2	43.5	35.8	34.9	36.22
MT training	15.4	15.0	9.3	16.6	14.3	14.12	35.3	33.7	43.6	36.2	35.8	36.92
↔ + ML Task Arithmetic	16.0	15.5	9.5	16.9	15.4	14.66	34.1	32.8	43.7	35.6	33.3	35.90
↔ + ML LoRS-Merging	16.1	15.3	9.4	16.8	15.3	14.57	34.2	32.7	43.8	35.8	33.5	36.00
ML training	16.7	15.5	10.0	17.0	16.6	15.14	32.3	33.2	43.5	35.4	34.3	35.74
↔ + MT Task Arithmetic	17.1	15.5	9.5	17.0	15.5	14.89	32.1	33.1	43.6	35.7	33.6	35.62
↔ + MT LoRS-Merging	16.9	15.5	9.4	16.8	15.5	14.80	32.6	33.2	43.6	35.9	33.6	35.78

Table 11: Multi-task model merging. Finetuned is the topline where the model is finetuned on each language and task combination independently, and Avg. averages WER or BLEU scores directly.

System	WER↓						BLEU↑					
	ca	de	es	fr	it	Avg.	ca	de	es	fr	it	Avg.
Pretrained	20.6	19.6	14.7	24.5	19.4	19.88	21.1	24.1	28.6	26.8	26.8	25.48
Finetuned	19.5	19.7	14.4	22.1	19.2	19.05	22.6	24.6	29.2	27.2	27.3	26.18
Multi-task training	17.0	19.7	14.4	24.2	19.4	19.00	22.3	24.6	28.7	27.0	26.9	25.90
Weight Averaging	17.1	19.6	13.9	23.7	19.6	18.84	22.9	24.4	29.0	27.7	26.9	26.18
Task Arithmetic	17.2	19.3	14.0	23.3	19.7	18.76	23.4	24.5	28.9	27.7	27.0	26.30
MP-Merging	17.8	19.5	14.4	23.8	17.2	18.62	23.1	24.5	29.1	27.9	27.4	26.40
SVP-Merging	18.0	19.4	14.4	23.7	17.7	18.72	22.9	24.7	29.1	27.8	27.4	26.38
LoRS-Merging	17.5	19.4	14.2	23.1	17.7	18.45	23.1	24.5	29.3	27.9	27.6	26.48