Training-free LLM Merging for Multilingual Multi-task Learning

Anonymous ACL submission

Abstract

Large Language Models (LLMs) have demonstrated exceptional capabilities across diverse natural language processing (NLP) tasks. The release of open-source LLMs like LLaMA and Qwen has triggered the development of numerous fine-tuned models tailored for various tasks and languages. In this paper, we explore an important question: is it possible to combine these specialized models to create a unified model with multi-task and multi-lingual capabilities. We introduces Hierarchical Iterative Merging 011 (Hi-Merging), a training-free method for uni-012 fying multiple specialized LLMs into a single model. Specifically, Hi-Merging employs 015 model-wise and layer-wise pruning and scaling, guided by contribution analysis, to mitigate parameter conflicts. Extensive experiments on En-017 glish and Chinese datasets, covering multiplechoice and question-answering tasks, validate Hi-Merging across three paradigms: multilingual merging, multi-task merging, and multilingual multi-task merging. The results demonstrate that Hi-Merging consistently outperforms existing merging techniques and surpasses the 025 performance of models fine-tuned on combined datasets in most scenarios. Code is available at this anonymous link¹.

1 Introduction

Large Language Models (LLMs) have revolutionized Natural Language Processing (NLP) by demonstrating unprecedented capabilities in capturing and utilizing world knowledge (Zhao et al., 2024). Recent advances in architecture design and training methodologies have enabled models like GPT-4 (OpenAI, 2023) to engage in human-like dialogue and solve real-world problems, enabling breakthroughs in healthcare, education, and scientific research.

With the advent of open-source large language models (LLMs) like LLaMA-3 (Dubey et al., 2024)

and Qwen (Yang et al., 2024a), significant research efforts have been dedicated to fine-tuning these models for specific tasks, domains, and languages. As a result, Hugging Face² now hosts over one million specialized LLMs across various languages and tasks, and this number continues to grow rapidly. These models represent a vast repository of task-specific and language-specific expertise, ranging from Chinese medical applications (Chen et al., 2023) to English financial question and answering (Cheng et al., 2024). A natural question arises: is it possible to combine these language-specific and task-specific fine-tuned LLMs into a single unified model with broad capabilities, including multi-lingual and multi-task functionalities? If achievable, the deployment of such a unified model could perform multiple tasks that currently require multiple LLMs, thereby significantly enhancing the application of LLMs. One potential solution is to gather all fine-tuning data and retrain the LLMs from scratch. However, this approach has three significant disadvantages: 1) the availability of fine-tuning data, as the models are often public but the data is not; 2) retraining large LLMs requires substantial computational resources; and 3) balancing the training data from different tasks and languages to achieve overall optimal performance is a non-trivial challenge.

041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

Based on the above considerations, model merging (Yang et al., 2024b) emerges as a promising solution for unifying multiple specialized models while preserving their individual capabilities. However, current model merging methods face two fundamental challenges. First, interference between merged models can arise from noise introduced by data bias (Tsuchiya, 2018) or the training process, such as overfitting, impairs the merged model's generalization. Second, models trained independently follow distinct optimization trajectories, leading

¹https://anonymous.4open.science/r/hi-merging

²https://huggingface.co/



Figure 1: Illustration of three paradigms for our LLM merging: merging models that specialize in different languages (left), merging models that excel at different tasks (middle), and merging models that exhibit expertise in both different languages and different tasks (left). Through such merging, a single model can inherit the combined capabilities of both original models, enabling broader applicability and enhanced performance.

to different knowledge alignments in their parameter spaces (Ilharco et al., 2023). These misaligned parameters become incompatible for direct combination without additional training.

087

100

101

102

103

104

105

106

108

110

111

112

To address these challenges, we propose Hi-Merging, a Hierarchical Iterative Merging method. It first applies model-wise pruning and scaling to the delta vectors (parameter differences between fine-tuned models and the foundation model) to eliminate noisy parameters introduced during finetuning. Then, we apply layer-wise pruning and scaling iteratively for the knowledge misalignment, starting from the most conflicted layers. To identify the severity of layer-wise conflicts, we develop contribution analysis - a method that quantifies each layer's contribution by measuring how adding or removing specific layers affects model capabilities. By analyzing how our contribution metrics change before and after a pre-merging process, we can identify potential conflicts, thereby guiding our iterative optimization process to resolve parameter incompatibilities without additional training systematically.

Our contributions can be summarized as follows:

- We pioneer the use of model merging to enhance LLMs' multilingual and multi-task capabilities without additional training.
- We propose Hi-Merging, a hierarchical iterative approach that effectively reduces the interference of noise and knowledge alignment conflicts during model merging.
- Extensive experiments on four datasets demonstrate the effectiveness of Hi-Merging under three

merging paradigms: multilingual merging, multitask merging, and multilingual multi-task merging, consistently achieving superior performance across all settings. 113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

2 Preliminary for LLM Merging

In this section, we detail notations and introduce existing model (LLM) merging solutions as the preliminary.

Model merging aims to combine multiple models with distinct capabilities as a single model, which has all the strengths of these models. In this paper, we user two-model merging for illustration: Given models \mathcal{M}_A and \mathcal{M}_B with parameters θ_A and θ_B , both fine-tuned from a foundation model \mathcal{M}_F with parameters θ_F for tasks t_A and t_B respectively, model merging aims to combine them into a single model \mathcal{M}_{merge} with parameters θ_{merge} that preserves capabilities for both tasks.

Typical model merging strategies include weighted averaging and delta vector-based merging. The former combines model parameters through a weighted sum (Wortsman et al., 2022):

$$\boldsymbol{\theta}_{\text{merge}} = \sum_{m \in \{A,B\}} \omega_m \boldsymbol{\theta}_m$$
 (1)

where ω_m is the weight to balance different capabilities constrained to $\sum_{m \in \{A,B\}} \omega_m = 1, \omega_m > 0$. And $m \in \{A, B\}$ is the model identifier.

The second strategy merges models based on delta vectors, the parameter differences between fine-tuned models and their foundation model,

143

148

149

150

151

152

153

154

155

156

158

159

162

163

164

165

166

169

170

171

172

173

174

175

176

177

178

179

182

183

186

which can be mathematically defined as:

$$\boldsymbol{\delta}_m = \boldsymbol{\theta}_m - \boldsymbol{\theta}_F \tag{2}$$

144 Delta vectors δ_m defined in Equation (2) reveal 145 model-specific updates from the foundation model, 146 enabling a delta-weighted merging strategy (II-147 harco et al., 2023):

$$\boldsymbol{\theta}_{\text{merge}} = \boldsymbol{\theta}_F + \sum_{m \in \{A,B\}} \omega_m \boldsymbol{\delta}_m$$
 (3)

where $\omega_m > 0$. Note that both strategies, illustrated in Equation (1) and Equation (3), can be easily extended to multiple model merging scenarios by expanding the model list $\{A, B\}$.

3 Method

In this section, we introduce the proposed method, which consists of two major components: (1) model-wise pruning and scaling that removes noisy and redundant parameters and moderate excessive ones and (2) layer-wise pruning and scaling iterating on conflicted layers to address knowledge misalignment issues.

3.1 Model-wise Pruning and Scaling

This section introduces two operations to process delta vectors: pruning and scaling.

During the fine-tuning, models can accumulate noisy parameters and learn sharp parameters for the specific fine-tuning task. We introduce the pruning and scaling operations to tackle these two problems, respectively, which are controlled by the following hyperparameters:

• **Pruning Threshold** (*p*): This parameter specifies the proportion of the delta vector that should be preserved. By retaining the largest *p* percentage of the vector's components and rendering the remaining (1 - p) to zero, the pruning operation can eliminates trivial parameter updates (data-specific noise) while preserving meaningful task-specific knowledge.

• Scaling Factor (s): This factor controls the magnitude of the delta vector. With this parameter, the scaling operation contributes to addressing over-aggressive parameters by scaling down sharp updates, which may result from the overfitting during fine-tuning. The pruning does not apply to large parameter changes as they likely encode essential knowledge. The scaling provides a way to moderate their excessive influence.



Figure 2: The accuracy of the fine-tuned Qwen2-7B-Instruct on the MedQA dataset after the model-wise pruning and scaling process with different combinations of the pruning threshold p and the scaling factor s.

With these hyperparameters, the pruning and scaling cooperatively process the delta vectors in a complementary manner: pruning eliminates negligible parameter changes while scaling moderates the significant ones. Note that both p and s constrained to [0, 1].

We empirically validate the effectiveness of the pruning and scaling operations by iterating p and s from [0.1, 1]. The result is visualized in Figure 2. We can find that the individual model can maintain or even improve performance with appropriate pruning and scaling. For example, p = 0.1, s = 0.9 (preserving 10% of parameters and scaling all delta values with 0.9) can defeat the original model (p = 1, s = 1). This finding supports our idea of conducting model-wise pruning and scaling to overcome noisy and radical parameter updates.

Next, we introduce the model-wise pruning and scaling details. Specifically, the delta vector (defined in Equation (2)) for a given LLM \mathcal{M}_m can be defined as $\delta_m = [\delta_{m,1}, \delta_{m,2}, \ldots, \delta_{m,N}]$, where $m \in \{A, B\}$ is the model identifier and N indicates the size of trainable parameters.

The **pruning** operation Top_p retains the $\lceil p \cdot N \rceil$ elements of δ_m with the largest absolute value and zeros out the rest, resulting in $\tilde{\delta}_m$:

$$\boldsymbol{\delta}_m = \operatorname{Top}_p(\boldsymbol{\delta}_m) \tag{4}$$

In detail, the *n*-th component of δ_m is:

$$\tilde{\delta}_{m,n} = \begin{cases} \delta_{m,n}, & \text{if } n \in \{\pi(1), \pi(2), \dots, \pi(\lceil p \cdot N \rceil)\} \\ 0, & \text{otherwise} \end{cases}$$
(5) 216

3

187

212 213 214

215

217

218

240

241 242

244

246

247

249

251

253

256

259

261

$$\alpha_{m1,m2}^{l} = \mathbf{P}_{t_{m1}}(\boldsymbol{\theta}_{m2} - \boldsymbol{\hat{\delta}}_{m2}^{l}) - \mathbf{P}_{t_{m1}}(\boldsymbol{\theta}_{m2})$$

$$\beta_{m1,m2}^{l} = \mathbf{P}_{t_{m1}}(\hat{\boldsymbol{\theta}}_{F} + \boldsymbol{\delta}_{m2}^{l}) - \mathbf{P}_{t_{m1}}(\boldsymbol{\theta}_{F})$$
(9)

where $\pi(n)$ represents the index of the *n*-th largest component of δ_m in absolute value, such that:

$$\left|\delta_{m,\pi(1)}\right| \ge \left|\delta_{m,\pi(2)}\right| \ge \dots \ge \left|\delta_{m,\pi(N)}\right| \quad (6)$$

The scaling operation adjusts the magnitude of the pruned delta vector $\boldsymbol{\delta}_m$ by multiplying it with the scaling factor $s \in [0, 1]$ as $s\delta_m$.

Regarding the different setting of p and s for each model, the model-wise pruning and scaling can be compactly expressed as:

$$\hat{\boldsymbol{\delta}}_m = s_m \cdot \operatorname{Top}_{p_m}(\boldsymbol{\delta}_m) = s_m \tilde{\boldsymbol{\delta}}_m$$
 (7)

where $\hat{\delta}_m$ represents the delta vector after the model-wise pruning and scaling.

Through model-wise processing with pruning and scaling, we effectively identify noisy and excessive parameter updates from the fine-tuning, maintaining and moderating the key knowledge about the fine-tuning task for the subsequent merging.

3.2 Layer-wise Pruning and Scaling

In this section, we conduct the layer-wise model merging with pruning and scaling operations with a novel contribution analysis method to measure the parameter conflict.

3.2.1 Contribution Analysis

Directly merging the model-wise processed delta vectors $\{\hat{\boldsymbol{\delta}}_m\}_{m \in \{A,B\}}$ as in Equation 1 or Equation 3 will encounter the weight misalignment problem, which is overlooked by existing methods.

To investigate potential conflicts when merging a specific layer, we measure its contribution by calculating the performance difference before and after the merge. Precisely, we assess the merging contribution from two directions:

- **Deletion Impact** (α) : To estimate this impact, we first construct a merged model \mathcal{M}_G that merges all layers using the merging process mentioned in Equation (1) or Equation (3). Then, we calculate the performance degradation caused by removing the delta vector for a specific layer.
- Addition Impact (β): This impact is measured by the performance improvement of adding the delta vector for a specific layer to the pre-trained foundation model \mathcal{M}_F .

These impacts can be mathematically represented as:

where
$$m1 \in \{A, B\}$$
 is the task capability identifier
and $m_2 \in \{A, B, G\}$ is the model identifier. We
investigate the layer-wise contribution so that l is
the layer index. $\hat{\delta}_{m2}^{l}$ is the delta vector for \mathcal{M}_{m2} at
layer l . $P_{t_{m1}}(\cdot)$ represents the performance metric
on the task t_{m1} . For example, BLEU-4 (Papineni

We sum up two impacts as the overall contribution:

$$c_{m1,m2}^{l} = \alpha_{m1,m2}^{l} + \beta_{m1,m2}^{l}$$
(10)

263

264

265

268

269

270

271

272

273

274

276

277

278

279

280

281

284

287

288

290

291

292

294

295

We

l is

3.2.2 **Iterative Conflict Elimination**

et al., 2002) score for the OA task.

The contribution analysis method defined in Equation (8)-(10) provides a solution to measure the importance of merging specific layers. We can then define the conflict resulted by model $\mathcal{M}_m(m \in$ $\{A, B\}$) within the layer l of the merged model as:

$$\gamma_m^l = c_{m,m}^l - c_{m,G}^l \tag{11}$$

In this formula, we set the capability identifier m1 = m as we expect the merged model can maintain the performance of \mathcal{M}_m on t_m . We can then identify the most severe conflicting layers that impair the fine-tuned performance by sorting $\Gamma^l = \sum_{m \in \{A,B\}} \gamma^l_m.$

To mitigate the parameter misalignment, we iteratively merge the most conflicting layers (with the largest Γ^l). Specifically, to process a specific layer, there are three types of conflict as illustrated in Figure 3:

1. Severe Conflict: $\gamma_A^l > 0$ and $\gamma_B^l > 0$, indicating both capabilities are impaired by the merging. In such cases, only the delta vector with a larger contribution is retained, e.g., dropping $\hat{\boldsymbol{\delta}}_{B}^{l}$ in the figure. Namely, $\hat{\boldsymbol{\delta}}_{B}^{l}$ is set to zero.





Figure 3: The demonstration of different conflict elimi-

nation strategies for three pre-merging conditions.

4

(8)

2. **Partial Conflict:** $\gamma_A^l * \gamma_B^l < 0$, i.e., one of the delta vectors leads to the parameter misalignment. The solution for this case is to prune and scale the conflict delta vector again, as we defined in Section 3.1. For example, in Figure 3, the overfitting on t_A ($\gamma_A^l < 0$ and $\gamma_B^l > 0$) leads to the degradation of the ability for t_B . As a result, we prune and scale $\hat{\delta}_A^l$ again as³:

$$\hat{\boldsymbol{\delta}}_{A}^{l} = s_{A} \cdot \operatorname{Top}_{p_{A}}(\hat{\boldsymbol{\delta}}_{A}^{l}),$$
 (12)

3. Mutual Enhancement: If $\gamma_A^l \leq 0$ and $\gamma_B^l \leq 0$, the merging process improves for both capabilities. In this case, no further adjustment is necessary for this layer.

After resolving the conflicts of all layers, the parameters of the final merged model \mathcal{M}_{merge} is:

$$\boldsymbol{\theta}_{\text{merge}} = \boldsymbol{\theta}_{\text{F}} + \hat{\boldsymbol{\delta}}_A + \hat{\boldsymbol{\delta}}_B$$
 (13)

4 Experiments

296

297

298

301

304

311

312

313

314

315

316

317

318

320

321

322

324

325

326

327

330

332

335

336

In this section, we conduct comprehensive experiments to evaluate the effectiveness of Hi-Merging on multilingual multi-task merging problems by answering following research questions (RQ):

- **RQ1:** How does Hi-Merging compare with baselines on merging LLMs for different languages?
- **RQ2:** How do these methods perform for merging LLMs for different tasks?
- **RQ3:** Are these methods applicable to combine LLMs on different taks and languages?
- **RQ4:** Is Hi-Merging able to merge open-source models from the Hugging Face?

4.1 Experimental Settings

4.1.1 Datasets

We select four datasets listed in Table 1 that cover multilingual multi-task capabilities, including English and Chinese languages, with multiple-choice question answering (MCQA) and open-domain question answering (QA) tasks.

4.1.2 Baselines

In our experiments, we compare Hi-Merging with weighted averaging (Model Soups (Wortsman et al., 2022)) and delta vector-based methods (Arithmetic (Ilharco et al., 2023), TIES-Merging (Yadav et al., 2023), DARE (Yu et al., 2024), DELLA (Deep et al., 2024), and Model Breadcrumbs (Davari and Belilovsky, 2024)), as

Table 1: The brief description and statistics of the four datasets (MedQA (Jin et al., 2020), CMExam (Liu et al., 2023), HealthCareMagic (Li et al., 2023), and cMedQA2) (Zhang et al., 2018) used for fine-tuning.

Name	Task	Language	Train	Validation	Test
MedQA	MCQA	English	10,000	400	400
CMExam	MCQA	Chinese	50,000	4,000	4,000
HealthCareMagic	QA	English	30,000	1,000	1,000
cMedQA2	QA	Chinese	30,000	1,000	1,000

well as the multilingual and multi-task models finetuned on combined datasets. Details of these baselines are in Appendix A.1.1.

341

342

343

344

345

346

347

348

349

350

351

352

353

354

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

374

375

376

4.1.3 Implementation Details

We use Qwen2-7B-Instruct as foundation modelswith results for other foundation models presented in Appendix A.2. For fine-tuning, we employ LLaMA-Factory ⁴ with LoRA (rank=8, alpha=16, dropout=0.01) and a batch size of 64. The learning rate is 1.0^{-4} with cosine decay and warm-up. LLM merging is performed using mergekit ⁵. For hyperparameter tuning, both *p* and *s* in model-wise processing range from 0.1 to 1.0 with a step of 0.1. In layer-wise processing, the pruning threshold *p* and scaling factor *s* are successively set to half of their model-wise values.

4.1.4 Evaluation Metrics

For the MCQA task, accuracy is employed to measure the proportion of correct answers (Devlin et al., 2019).For the QA task, we use BLEU-4 (Papineni et al., 2002) to evaluate the precision of the generation, and ROUGE-1,2,L (Lin, 2004) to assess the overlap and coherence with the ground truth. The averaged performance across all metrics is reported as "Avg.". Additionally, we report the relative improvement over the foundation model.

4.2 Multilingual Merging (RQ1)

We first verify the effectiveness of Hi-Merging on multilingual LLM merging. Here, we merge models trained on datasets in different languages (English and Chinese) but for the same task type (QA in Table 2 and MCQA results in Table 3).

With the foundation model Qwen2-7B-Instruct (Yang et al., 2024a), we observe more moderate performance degradation with our method, while Llama3-8B-Instruct (Dubey et al., 2024) exhibits persistent conflicts (results in

³We use the same notation $\hat{\delta}_{A}^{l}$ for clarity.

⁴https://github.com/hiyouga/LLaMA-Factory

⁵https://github.com/arcee-ai/mergekit

Types	Methods		L1 (Health	CareMagic)			L2 (cN	ledQA2)		Δνσ	Impr
Types		BLEU-4	ROUGE-1	ROUGE-2	ROUGE-1	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-1		
Pre-trained	Qwen2-7B-Instruct	30.1209	26.3524	5.3280	15.7451	1.7090	14.1527	1.7822	9.0934	13.0355	-
Fine-tuned	Model A (L1)	35.5717	30.2512	8.9044	20.3625	3.7609	19.1370	3.1364	15.1441	17.0335	+30.66%
	Model B (L2)	24.8587	24.9841	4.1492	15.1967	4.4159	21.2210	4.0680	17.4600	14.5442	+11.57%
	Multilingual	35.7637	29.9781	8.6687	20.1184	3.7660	20.9869	3.7784	16.8850	17.4932	+34.19%
	Model Soups	33.2627	28.8258	7.5487	18.9459	4.6801	21.5564	4.0502	17.5380	17.0510	+30.80%
	Task Arithmetic	33.0398	28.7169	7.5726	18.9600	4.7181	21.4108	4.0503	17.6772	17.0182	+30.55%
	TIES	33.6571	29.0496	7.7769	19.1503	4.3751	20.8551	3.7518	17.1978	16.9767	+30.23%
	DARE	33.3031	28.9575	7.8222	19.1702	4.7578	21.0865	3.8996	17.2488	17.0307	+30.64%
Merged	DARE+TIES	26.8091	26.0330	5.2307	16.5201	4.2456	20.6276	3.7531	17.1445	15.0455	+15.41%
	Model Breadcrumbs	34.3247	29.4403	8.1518	19.6443	4.4092	20.9365	3.8138	17.1378	17.2323	+32.19%
	DELLA	33.4207	28.9234	7.6728	18.9674	4.6827	21.1596	4.0709	17.4775	17.0469	+30.77%
	DELLA+TIES	27.2331	26.1723	5.4339	16.6009	4.7130	21.2275	4.2944	17.7694	15.4306	+18.37%
	Hi-Merging (Ours)	35.9500	29.9826	8.8738	20.3844	4.7009	21.1752	3.9704	17.2361	17.78	+36.42%

Table 2: Performance comparison of merging methods for multilingual QA using Qwen2-7B-Instruct.

Note: (1) Model A is fine-tuned on HealthCareMagic (L1: English), Model B is fine-tuned on cMedQA2 (L2: Chinese), Multilingual model is fine-tuned on both datasets; (2) Merged models are obtained by merging Model A and B; (3) The overall best result is marked in bold and the best merging result is underlined.

Table 3: Performance	comparison	of merging	methods
for multilingual MCQ	A using Qw	en2-7B-Inst	ruct.

Types	Methods	L1 (MedQA)	L2 (CMExam)	Avg.	Impr.
	Qwen2-7B-Instruct	51.4062	74.6217	63.0140	-
Pre-trained	Yi-1.5-9B	46.8185	58.6499	52.7342	-16.31%
	Baichuan2-7B	6.4415	7.1439	6.7927	-89.22%
Fine-tuned	Model A (L1)	59.1406	83.7771	71.4589	+13.40%
	Model B (L2)	54.4531	88.6171	71.5351	+13.52%
	Multilingual	60.0781	88.2246	74.1514	+17.67%
	Model Soups	59.6094	88.6926	74.1510	+17.67%
	Task Arithmetic	59.5312	88.7681	74.1497	+17.67%
	TIES	59.0625	88.7832	73.9229	+17.31%
	DARE	58.6719	88.6926	73.6823	+16.93%
Merged	DARE + TIES	58.9063	88.6021	73.7542	+17.04%
	Model Breadcrumbs	58.8281	88.6322	73.7302	+17.00%
	DELLA	58.9844	88.7681	73.8763	+17.24%
	DELLA + TIES	58.2812	88.7530	73.5171	+16.67%
	Hi-Merging (Ours)	<u>61.0156</u>	88.6501	74.8329	+18.76%

Note: (1) Model A is fine-tuned on MedQA (L1: English). Model B is fine-tuned on CMExam (L2: Chinese), Multilingual model is fine-tuned on both datasets. (2) Merged models are obtained by merging Model A and B. (3) The overall best result is in bold and the best merging result is underlined.

Appendix A.2). This difference stems from the capability gap of foundation models. Qwen2's stronger language ability enables more structured parameter updates, where new knowledge aligns better. In contrast, weaker foundation models lead to dispersed parameter updates and knowledge misalignment. In addition, we investigate the impact of different training sample sizes on LLM merging in Appendix A.3.

Baselines like Model Soups and Task Arithmetic that combine models without considering noises and conflicts achieve relatively stable average performance. Methods that reduce conflicts, such as TIES and DARE, occasionally achieve the best results on individual metrics by discarding delta vectors. However, their average performance needs more systematic processing strategies. In contrast, our Hi-Merging method, with hierarchical pruning and scaling approach, not only achieves the best average performance but attains optimal results in about half of the individual metrics. 393

394

395

396

397

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

494

425

4.3 Multi-task Merging (RQ2)

For multi-task merging, we combine models trained on the same task but in different languages (e.g., English MCQA with Chinese MCQA), as shown in Table 4. The results show that merged models consistently outperform their individual fine-tuned counterparts, with many merging methods even surpassing multi-task fine-tuned models. Notably, our Hi-Merging approach achieves a 1.84% relative improvement over the multi-task fine-tuned model. We attribute this success to three factors. 1) During multi-task fine-tuning with limited data (compared to pre-training), tasks can interfere with each other due to the "seesaw effect". In contrast, model merging allows parameters to be optimized independently before integration, avoiding such interference. 2) Since both models are fine-tuned from the same foundation model, their parameter updates tend to follow similar optimization trajectories, making successful merging more likely. 3) The inherent sparsity of large language models provides sufficient parameter space to accommodate multi-task knowledge from both models without significant conflicts.

4.4 Multilingual Multi-task Merging (RQ3)

For multilingual multi-task merging, we combine models trained on completely different tasks and languages. Specifically, we merge a model

377

378

			1	L1 (English)				L	2 (Chinese)				
Types	Methods	T1 (MedQA)		T2 (Health	CareMagic)		T1 (CMExam)		T2 (cMedQA2)			Avg.	Impr.
		Accuracy	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-1	Accuracy	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-1		
Pre-trained	Qwen2-7B-Instruct	51.4062	30.1209	26.3524	5.3280	15.7451	74.6217	1.7090	14.1527	1.7822	9.0934	23.0312	-
Fine-tuned	Model A (T1)	59.1406	34.6533	28.7482	6.9168	17.9525	88.6171	2.8064	16.8617	2.5603	12.0561	27.0313	+17.36%
	Model B (T2)	53.0469	35.5717	30.2512	8.9044	20.3625	81.5670	4.4159	21.2210	4.0680	17.4600	27.6869	+20.21%
	Multi-task	59.2188	35.6009	30.2101	9.1375	20.4645	88.6926	3.7790	20.5919	3.8096	16.9265	28.8431	+25.23%
Merged	Model Soups	58.5156	36.4411	30.5654	9.1754	20.4259	88.8285	4.3912	21.0216	4.0040	17.2843	29.0653	+26.19%
	Task Arithmetic	58.5938	36.3290	30.6624	9.1945	20.5406	88.7983	4.3018	20.6467	3.7496	16.9995	29.0493	+26.13%
	TIES	60.4688	35.7851	30.3243	9.0310	20.3723	88.6171	4.5434	21.5629	4.1910	17.4909	29.1996	+26.78%
	DARE	58.4375	36.5802	30.5488	9.0818	20.3945	88.7681	4.5487	21.3255	3.8403	17.4471	29.0865	+26.29%
	DARE+TIES	59.3750	35.7062	30.1950	8.7840	20.0878	88.8285	4.1587	21.1291	3.8124	17.2868	28.9361	+25.63%
	Model Breadcrumbs	57.8906	36.4620	30.2173	8.7845	20.0169	88.8889	4.4472	21.2492	3.8931	17.2846	28.9128	+25.53%
	DELLA	58.5938	36.3494	30.1715	8.8125	20.1879	88.8134	4.3718	21.0226	3.9300	17.3403	28.9818	+25.83%
	DELLA+TIES	59.5312	36.0774	30.4743	9.1151	20.4599	88.6021	4.3202	21.2269	4.0164	17.3779	29.0115	+25.96%
	Hi-Merging (Ours)	60.5469	36.4926	30.5467	9.1231	20.3523	88.9795	4.6781	21.5367	4.2165	17.5038	29.2673	+27.07%

Table 4: Performance comparison of merging methods for multi-task learning.

Note: (1) Model A is fine-tuned on English datasets (T1: MedQA or HealthCareMagic), Model B is fine-tuned on Chinese datasets (T2: CMExam or cMedQA2), Multi-task models are fine-tuned on datasets with the same language; (2) Merged models are obtained by merging Model A and B; (3) The overall best result is marked in bold and the best merging result is underlined.

Table 5: Performance comparison of merging methods for multilingual multi-task learning.

Types	Methods	T1, L1 T2, L2 (MedQA) (cMedQA2)				T1, L2 (CMExam)	n) T2, L1 (HealthCareMagic)				Avg.	Impr.	
		Accuracy	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-1	Accuracy	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-1		
Pre-trained	Qwen2-7B-Instruct	51.4062	1.709	14.1527	1.7822	9.0934	74.6217	30.1209	26.3524	5.328	15.7451	23.0312	-
Fine-tuned	Model A (T1) Model B (T2) Multilingual Multi-task	59.1406 54.4922 60.7812	2.8064 4.4159 3.8473	16.8617 21.221 20.8741	2.5603 4.068 4.0434	12.0561 17.46 16.9525	88.6171 79.6875 88.9795	34.6713 35.5717 35.7429	28.4279 30.2512 30.1735	6.6122 8.9044 8.9153	18.1117 20.3625 20.3902	26.9865 27.6434 29.0700	+17.17% +20.03% +26.22%
Merged	Model Soups Task Arithmetic TTES DARE DARE DARE+TIES Model Breadcrumbs DELLA DELLA+TIES Hi-Merging (Ours)	58.3584 58.0469 59.6094 57.8906 58.75 57.1094 58.0469 59.0625 60.2344	4.6592 4.6682 4.3764 4.5671 4.4929 4.7217 <u>4.8065</u> <u>4.4854</u> 4.7743	21.2316 21.2618 21.0083 21.1856 21.3194 21.4192 21.5135 20.9954 21.1954	4.0559 4.0984 3.9002 3.9549 4.0824 4.1477 4.1356 4.0491 4.1749	17.3805 17.4231 17.4194 17.2328 17.4826 17.4182 17.4962 17.563 17.3991	88.6322 88.7379 88.7228 88.6322 88.5568 88.6021 88.6167 88.6624 88.6624	36.1765 36.1222 35.7708 35.8639 34.8223 36.4961 36.0159 35.0176 36.5223	30.7169 30.2256 30.5143 30.1489 29.7597 30.3911 30.3747 29.9666 30.3932	9.2702 8.757 8.8994 8.815 8.3004 9.0696 9.0414 8.658 8.7882	20.5227 20.1357 20.3487 20.1025 19.7624 20.4108 20.3929 20.1406 20.1619	29.1004 28.9477 29.0570 28.8394 28.7329 28.9786 29.0440 28.8601 29.2442	+26.35% +25.69% +26.16% +25.22% +24.76% +25.82% +26.11% +25.31% +27.02%

Note: (1) Model A is fine-tuned on MCQA datasets (T1: MedQA or CMExam), Model B is fine-tuned on QA datasets (T2: cMedQA2 or HealthCareMagic); (2) Merged models are obtained by merging Model A and B; (3) The overall best result is marked in bold and the best merging result is underlined.

trained for MCQA in one language (Model A: MedQA in English or CMExam in Chinese) with another model trained for QA in the opposite language (Model B: cMedQA2 in Chinese or Health-CareMagic in English), as illustrated in Table 5.

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

An intriguing phenomenon observed in our experiments is that multilingual multi-task fine-tuning tends to predominantly affect QA task performance, while model merging methods typically significantly impact MCQA task performance. We attribute this to two key factors. First, the QA task requires free-form text generation, demanding a more complex representation space than the relatively constrained choice selection in the MCQA task. Therefore, this complexity makes QA performance more vulnerable to degradation during joint fine-tuning. Second, model merging directly combines model parameters. Since MCQA tasks need exact boundaries for classification, these boundaries are more easily disrupted during the merging process, making MCQA performance more vulnerable to merging operations.

4.5 Open-source LLM Merging (RQ4)

To validate the generality of our merging approach, we conduct experiments using two opensource medical models from Hugging Face: Echelon-AI/Med-Qwen2-7B⁶, fine-tuned on English datasets for tasks such as medical QA and information retrieval (IR), and shtdbb/qwen2-7bmed⁷, fine-tuned on Chinese datasets for dialogue generation. Both models are derived from Qwen2-7B-Instruct. Figure 4 illustrates the performance comparison across 12 medical datasets, with metrics normalized for better visualization. 448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

Our approach demonstrates robust performance across the task spectrum. In 7 out of 12 datasets, Hi-Merging achieves the best performance among all models, with only two datasets showing apparent degradation compared to the better-performing individual model. These results demonstrate Hi-

⁶https://huggingface.co/Echelon-AI/ Med-Qwen2-7B

⁷https://huggingface.co/shtdbb/qwen2-7b-med



502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544



Figure 4: Performance of Hi-Merging on two opensource medical models, Echelon-AI/Med-Qwen2-7B and shtdbb/qwen2-7b-med, which are fine-tuned from the foundation model Qwen/Qwen2-7B-Instruct.

Merging's ability to effectively fuse medical knowledge while maintaining or enhancing performance across diverse languages and tasks. Detailed implementation setup and unprocessed numerical results can be found in Appendix A.1.2 and A.4.

5 Related Works

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491 492

493

494

495

496

5.1 Multilingual Task-Oriented LLMs

Multilingual tasks in NLP encompass a broad range of challenges, including machine translation (Wang et al., 2022), multilingual text summarization (Gambhir and Gupta, 2017), and sentiment analysis (Dashtipour et al., 2016) etc. Recently, LLMs have greatly contributed to advancing multilingual tasks by leveraging massive amounts of multilingual data (Brown et al., 2020; Devlin et al., 2019; Xue et al., 2021). Despite their success, LLMs struggle with multilingual limitations due to imbalanced pre-training data, resulting in better performance for high-resource languages over low-resource ones.

To enhance multilingual capabilities, LLMs employ continual training on specific languages, as seen in models like Chinese-LLaMA (Cui et al., 2023) and EuroLLM (Martins et al., 2024). Additionally, supervised fine-tuning techniques, such as LoRA in Chinese-Alpaca (Cui et al., 2023) and incorporating translation tasks in XGLM-7B (Li et al., 2024), further improve multilingual understanding. However, LLMs are usually enhanced for one language at a time, resulting in multiple isolated models.

5.2 Model Merging

Model merging aims to integrate knowledge from multiple fine-tuned models into a single one. These methods are categorized into two types: weightedbased merging and interference mitigation.

Weighted-based merging focuses on combining model parameters effectively. This includes simple techniques like parameter averaging, such as Model Soups (Wortsman et al., 2022), Fisherweighted merging (Matena and Raffel, 2022) and RegMean (Jin et al., 2023). While computationally efficient, these methods often miss conflicting parameter updates, leading to performance degradation. Therefore, Task Arithmetic (Ilharco et al., 2023) proposes manipulating delta vectors. AdaMerging (Yang et al., 2024c) and evolutionary algorithms (Akiba et al., 2024) optimize merging coefficients and blend diverse models, respectively.

Interference mitigation techniques aim to reduce parameter conflicts based on the overparameterization and sparsity of LLM. DARE (Yu et al., 2024) and SparseGPT (Frantar and Alistarh, 2023) show high LLM performance despite significant parameter pruning. DELLA (Deep et al., 2024) introduces MAGPRUNE for selective pruning and parameter rescaling. TALL-masks (Wang et al., 2024) isolate task-specific parameters to minimize interference. However, these techniques focus mainly on individual parameter-level operations without considering the structural relationships and knowledge dependencies across model layers.

6 Conclusion

In this paper, we proposed Hi-Merging, a novel approach for merging LLMs for multilingual multitask learning. Hi-Merging leverages model-wise and layer-wise pruning and scaling strategy to minimize the conflict between fine-tuned models' delta vectors. The model-wise process eliminates the fine-tuning noise and overfitting parameters of the original models. Then, the layer-wise process analyzes the contribution of each layer's delta vector to the fine-tuning performance, reducing the interference of conflicts in several key layers. Extensive experiments on the MCQA and QA datasets demonstrated that Hi-Merging outperforms traditional merging techniques and even surpasses models trained on multiple datasets. Future work will explore finer-grained conflict analysis strategies.

547

551

553

555

557

563

564

565

566

567

568

569

570

571

572

573

574

575

577

580

581

582

583

584

586

587

588

589

590

591

592

596

7 Limitations

While our proposed Hi-Merging method demonstrates promising results, several limitations should be acknowledged. First, our evaluation is currently limited to two task types (MCQA and QA) and two languages (English and Chinese). The effectiveness of Hi-Merging on a broader range of NLP tasks and language families remains to be investigated.

Second, our method focuses on merging models fine-tuned from the same foundation model. The applicability and performance of Hi-Merging when merging models from different architectural families or pre-training approaches is yet to be explored. This limitation becomes particularly relevant as the field continues to see diverse model architectures and training paradigms.

Finally, our current implementation assumes relatively balanced task importance. The method might need adaptation for scenarios where certain tasks or languages should be prioritized over others, potentially requiring a more flexible weighting mechanism in the merging process.

References

- Takuya Akiba, Makoto Shing, Yujin Tang, Qi Sun, and David Ha. 2024. Evolutionary optimization of model merging recipes. *CoRR*, abs/2403.13187.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Junying Chen, Xidong Wang, Anningzhe Gao, Feng Jiang, Shunian Chen, Hongbo Zhang, Dingjie Song, Wenya Xie, Chuyi Kong, Jianquan Li, Xiang Wan, Haizhou Li, and Benyou Wang. 2023. Huatuogptii, one-stage training for medical adaption of Ilms. *Preprint*, arXiv:2311.09774.
- Daixuan Cheng, Shaohan Huang, and Furu Wei. 2024. Adapting large language models via reading comprehension. In *The Twelfth International Conference on Learning Representations*.

- Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.
- Kia Dashtipour, Soujanya Poria, Amir Hussain, Erik Cambria, Ahmad Y. A. Hawalah, Alexander F. Gelbukh, and Qiang Zhou. 2016. Multilingual sentiment analysis: State of the art and independent comparison of techniques. *Cogn. Comput.*, 8(4):757–771.
- MohammadReza Davari and Eugene Belilovsky. 2024. Model breadcrumbs: Scaling multi-task model merging with sparse masks. In *Computer Vision - ECCV* 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXXV, volume 15133 of Lecture Notes in Computer Science, pages 270–287. Springer.
- Pala Tej Deep, Rishabh Bhardwaj, and Soujanya Poria. 2024. Della-merging: Reducing interference in model merging through magnitude-based sampling. *CoRR*, abs/2406.11617.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 4171–4186. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Elias Frantar and Dan Alistarh. 2023. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 10323–10337. PMLR.
- Mahak Gambhir and Vishal Gupta. 2017. Recent automatic text summarization techniques: a survey. *Artif. Intell. Rev.*, 47(1):1–66.
- Gabriel Ilharco, Marco Túlio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? A large-scale open domain question answering dataset from medical exams. *CoRR*, abs/2009.13081.

598 599 600

601

602

603

604

597

612

613

614

615

616

617

630

631

632

633

634

625

635 636

637

638

639

640

641

642

643

644

645

646

647

648

649

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

707

651

- 65 65
- 66 66
- 664 665 666 667
- 668 669 670 671 672 673
- 674 675 676 677
- 678 679 680 681
- 68 68
- 6
- 68

6

6

- 69
- 6

6

- 6
- 700

7

70

- Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. 2023. Dataless knowledge fusion by merging weights of language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5,* 2023. OpenReview.net.
- Jiahuan Li, Hao Zhou, Shujian Huang, Shanbo Cheng, and Jiajun Chen. 2024. Eliciting the translation ability of large language models via multilingual finetuning with translation instructions. *Trans. Assoc. Comput. Linguistics*, 12:576–592.
- Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve Jiang, and You Zhang. 2023. Chatdoctor: A medical chat model fine-tuned on a large language model meta-ai (llama) using medical domain knowledge. *Cureus*, 15(6).
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Junling Liu, Peilin Zhou, Yining Hua, Dading Chong, Zhongyu Tian, Andrew Liu, Helin Wang, Chenyu You, Zhenhua Guo, Lei Zhu, et al. 2023. Benchmarking large language models on cmexam–a comprehensive chinese medical exam dataset. *arXiv preprint arXiv:2306.03030*.
- Pedro Henrique Martins, Patrick Fernandes, João Alves, Nuno M. Guerreiro, Ricardo Rei, Duarte M. Alves, José Pombal, Amin Farajian, Manuel Faysse, Mateusz Klimaszewski, Pierre Colombo, Barry Haddow, José G. C. de Souza, Alexandra Birch, and André F. T. Martins. 2024. Eurollm: Multilingual language models for europe. *Preprint*, arXiv:2409.16235.
- Michael Matena and Colin Raffel. 2022. Merging models with fisher-weighted averaging. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- OpenAI. 2023. GPT-4 technical report. CoRR, abs/2303.08774.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA, pages 311–318. ACL.
- Masatoshi Tsuchiya. 2018. Performance impact caused by hidden bias of training data for recognizing textual entailment. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Haifeng Wang, Hua Wu, Zhongjun He, Liang Huang, and Kenneth Ward Church. 2022. Progress in machine translation. *Engineering*, 18:143–153.

- Ke Wang, Nikolaos Dimitriadis, Guillermo Ortiz-Jiménez, François Fleuret, and Pascal Frossard. 2024. Localizing task information for improved model merging and compression. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.
- Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. 2022. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 23965–23998. PMLR.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings* of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021, pages 483–498. Association for Computational Linguistics.
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin A. Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models. In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024a. Qwen2 technical report. CoRR, abs/2407.10671.
- Enneng Yang, Li Shen, Guibing Guo, Xingwei Wang, Xiaochun Cao, Jie Zhang, and Dacheng Tao. 2024b. Model merging in llms, mllms, and beyond: Methods, theories, applications and opportunities. *CoRR*, abs/2408.07666.
- Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guibing Guo, Xingwei Wang, and Dacheng Tao. 2024c. Adamerging: Adaptive model merging for multi-task learning. In *The Twelfth International Conference*

866

867

818

819

820

767

766

- 774 775 776
- 778
- 780
- 781
- 784

790

796

797

801

809

810

813

814

815

817

А

A.1 Experimental Settings

arXiv:2303.18223.

Appendix

2024. OpenReview.net.

A.1.1 Baselines

In our experiments, we compare it against a comprehensive set of baseline methods, including traditional weighted averaging techniques and stateof-the-art approaches specifically developed for fine-tuned models.

on Learning Representations, ICLR 2024, Vienna,

Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin

Li. 2024. Language models are super mario: Absorb-

ing abilities from homologous models as a free lunch.

In Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27,

Sheng Zhang, Xin Zhang, Hui Wang, Lixiang Guo, and

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,

Xiaolei Wang, Yupeng Hou, Yingqian Min, Be-

ichen Zhang, Junjie Zhang, Zican Dong, Yifan Du,

Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao

Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang

Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen.

2024. A survey of large language models. Preprint,

selection. IEEE Access, 6:74061-74071.

Shanshan Liu. 2018. Multi-scale attentive interac-

tion networks for chinese medical question answer

Austria, May 7-11, 2024. OpenReview.net.

- Multilingual Multi-task Training This approach trains a single model on the combined datasets of multiple languages simultaneously, without distinguishing between tasks.
- Model Soups (Wortsman et al., 2022) Uniform Soup is a simple merging method where the parameters of the fine-tuned models are averaged based on their importance.
- Task Arithmetic (Ilharco et al., 2023) This method performs arithmetic operations on the parameter differences between the pre-trained and fine-tuned models.
- TIES (Yadav et al., 2023) The Task Interference Elimination Strategy (TIES) minimize negative transfer and task interference by pruning redundant parameters and using a chosen sign to determine parameter update directions.
- DARE (Yu et al., 2024) Delta Alignment for Robust Ensemble (DARE) reduces the interference across tasks by randomly drop the delta vectors.
- · Model Breadcrumbs (Davari and Belilovsky, 2024) This approach tracks and prunes maxima and minima in delta vectors to retain critical taskspecific features.

• DELLA (Deep et al., 2024) DELLA follows DARE and assign drop rates to delta vectors according to their absolute values, improving performance stability.

A.1.2 Implementation Details

Our experimental environment consisted of a CentOS Linux 7 operation system with Python 3.12.4 and CUDA 12.2. All model training and inference operations were implemented using PyTorch 2.4.0. The hardware setup included 8 Tesla V100 GPUs, each equipped with 32GB of memory, enabling efficient parallel processing of large-scale models.

For model adaptation, we applied LoRA to all linear networks in the model. The learning rate schedule was carefully designed with a 100-step warm-up phase followed by cosine decay, which helped achieve stable convergence while maintaining optimal model performance. This configuration proved effective in balancing training efficiency and model quality across both multilingual and multi-task scenarios.

In addition to Qwen2-7B-Instruct, we also experimented with other foundation models including Llama-3-8B-Instruct (results shown in A.2). However, Qwen2-7B-Instruct demonstrated more consistent performance, particularly in handling both English and Chinese tasks, making it the preferred choice for our main experiments.

For visualization in Figure 4, we normalized the performance metrics to facilitate clear comparisons. The performance values of the models on each dataset represent the average of the QA task metrics (BLEU-4, ROUGE-1, ROUGE-2, and ROUGE-L). We scaled the pre-trained Qwen2-7B-Instruct's performance to 20 and the better-performing finetuned model's performance to 80 for each task. The performance values of the other fine-tuned model and our merged model were then proportionally adjusted within this range to maintain their relative differences.

A.2 Multilingual Merging

We further conduct experiments on Llama-3-8B-Instruct, as presented in Table 6 and 7. The results in Table 7 show that the performance of merged models based on Llama-3-8B-Instruct is generally inferior to that of the pre-merged fine-tuned models. This indicates that the effectiveness of the merging process is strongly influenced by the quality of the foundation models.

The observed degradation in performance can

896

Table 6: Performance comparison of merging methods for multilingual MCQA using Llama-3-8B-Instruct.

Types	Methods	L1 (MedQA)	L2 (CMExam)	Avg.
	Llama-3-8B-Instruct	57.9733	17.2821	37.6277
Pre-trained	GLM-4-9B	54.7656	69.5194	62.1425
	Gemma-2-9B	14.2583	2.7698	8.5141
Fine-tuned	Model A (L1)	60.4688	52.2706	56.3697
	Model B (L2)	60.3906	60.5525	61.0575
	Multilingual	62.8906	61.0356	61.9631
	Model Soups	61.2500	61.0507	61.1504
	Task Arithmetic	61.2500	<u>61.8750</u>	61.5625
	TIES	61.7188	61.3225	61.5207
	DARE	61.5625	61.3678	61.4652
Merged	DARE + TIES	60.9375	59.4656	60.2016
	Model Breadcrumbs	61.0156	60.4318	60.7237
	DELLA	60.8594	60.7186	60.7890
	DELLA + TIES	61.9531	61.3527	61.6529
	Hi-Merging (Ours)	62.2656	61.0757	<u>61.6707</u>

Note: (1) Model A is fine-tuned on MedQA (L1: English). Model B is fine-tuned on CMExam (L2: Chinese), Multilingual model is fine-tuned on both datasets. (2) Merged models are obtained by merging Model A and B. (3) The overall best result is in bold and the best merging result is underlined.

be attributed to several factors. First, weaker foundation models, such as Llama-3-8B-Instruct, tend to produce delta vectors with more dispersed and less coherent parameter distributions during finetuning. These delta vectors often carry noisy or conflicting information, which makes the merging process prone to parameter conflicts. Second, the weaker representational capacity of these models limits their ability to encode robust and semantically aligned knowledge, further exacerbating the challenges of merging.

In contrast, stronger foundation models, such as Qwen2-7B-Instruct, exhibit fewer conflicts during merging and demonstrate consistent performance improvements across tasks and languages. This is because their fine-tuned delta vectors are more compact and carry knowledge that is better aligned with the foundation model's semantic space, making the integration process more effective.

A.3 Number of training samples

We examine the impact of varying the number of training samples on the conflict during model merging, as shown in Figure 5. In the experiment, we use two QA datasets, HealthCareMagic (English) and cMedQA2 (Chinese), sampling 10k, 20k, 30k, 40k, and 50k training examples from each to produce a series of fine-tuned models, five per dataset. This setup evaluates how the number of training samples influences both individual model performance and compatibility during merging. The xaxis of Figure 5 represents the number of training samples, while the y-axis denotes the average performance metrics, including BLEU-4, ROUGE-1, ROUGE-2, and ROUGE-L.

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

However, Figures 5b and 5c show that merged models through either Model Soups or Task Arithmetic suffer from performance drops driven by the increasing size of training sample as further training leads to conflicting highly specialized models. Figure 5d shows the opposite: our method retains performance trends in line with fine-tuned models and addresses conflicts to retain improving performance through larger training sets.

These results highlight the robustness of our method in resolving merging conflicts, ensuring that the merged models retain the strengths of individual models while achieving stable and superior performance across training sample sizes.

A.4 Open Source LLM Merging

Table 8 presents the detailed numerical results for all models across the 12 medical datasets. The datasets cover a wide range of medical tasks and languages, allowing us to comprehensively evaluate the models' capabilities and the effectiveness of our merging approach.

Types	Methods		L1 (Healt	hCareMagic)	1		L2 (cM	(ledQA2)		Avg.	Impr
Types		BLEU-4	ROUGE-1	ROUGE-2	ROUGE-L	BLEU-4	ROUGE-1	ROUGE-2	ROUGE-L		
Pre-trained	Llama-3-8B-Instruct	16.3118	21.6011	3.1389	10.8666	0.0225	0.4710	0.0211	0.2343	6.5834	-
	Model A (L1)	36.0325	30.4111	9.2743	20.7236	0.0185	0.1288	0.0025	0.0841	12.0844	+83.5%
Fine-tuned	Model B (L2)	3.9950	7.7673	0.9267	4.6358	3.0638	20.4016	3.4178	16.3067	7.5643	+14.9%
	Multilingual	35.6154	30.5447	9.2156	20.4271	3.0250	20.3136	3.4964	16.0911	17.3411	+163.5%
	Model Soups	32.1199	28.2278	6.3715	18.2456	3.3256	19.5499	2.8670	15.4688	15.7720	+139.5%
	Task Arithmetic	31.6679	27.7646	6.0354	18.0448	3.3805	19.6475	2.9507	15.4806	15.6215	+137.2%
	TIES	32.1494	28.0527	6.7440	18.2913	3.2238	19.5369	2.8854	15.2112	15.7618	+139.4%
	DARE	25.9679	25.6716	4.5173	16.3803	3.5337	20.8586	3.1736	16.6716	14.5968	+121.7%
Merged	DARE+TIES	26.6707	25.9106	5.2031	16.5525	3.2236	19.8564	2.9963	15.5967	14.5012	+120.2%
	Model Breadcrumbs	26.9844	26.1004	4.7037	16.3247	3.3307	20.7442	3.3069	16.2874	14.7228	+123.6%
	DELLA	25.6313	25.6792	4.5522	16.1313	3.6612	20.9176	3.3286	16.7355	14.5796	+121.4%
	DELLA+TIES	27.1246	26.0186	5.3163	16.6170	3.3433	19.9122	3.0848	15.9942	14.6764	+122.9%
	Hi-Merging (Ours)	33.5960	28.4141	7.2167	18.8804	3.1967	19.8207	2.9509	15.7833	16.2324	+146.5%

Table 7: Performance comparison of merging methods for multilingual QA using Llama-3-8B-Instruct.

Note: (1) Model A is fine-tuned on HealthCareMagic (L1: English), Model B is fine-tuned on cMedQA2 (L2: Chinese), Multilingual model is fine-tuned on both datasets; (2) Merged models are obtained by merging Model A and B; (3) The overall best result is marked in bold and the best merging result is underlined.



Figure 5: Impact of training sample size on model merging conflicts. Blue and orange lines represent the average performance metrics for HealthCareMagic and cMedQA2, respectively.

Table 8: Numerical performance of Hi-Merging on two open-source models, Echelon-AI/Med-Qwen2-7B and shtdbb/qwen2-7b-med.

Models	MedQA	MediQA	Medical Flashcards	Health Advice	Pubmed	WikiDoc	WikiDoc Patient	CORD 19	iCliniq	HealthCareMagic	ChatMed	MedChatZH
Qwen2-7B-Instruct	37.3868	17.3595	22.7668	2.8205	5.7994	17.6217	18.785	39.1748	19.3292	28.7051	9.9138	8.0654
Echelon-AI/Med-Qwen2-7B	64.2862	32.052	41.1081	97.7523	92.9898	20.7237	26.9203	40.7167	26.5593	30.3212	15.1218	9.2714
shtdbb/qwen2-7b-med	40.1598	27.1442	29.85	4.096	11.5017	20.5808	21.1528	41.2026	27.332	33.3678	19.4513	11.2665
Hi-Merging (Ours)	64.9011	31.9421	45.1714	97.755	92.1692	21.0211	26.3293	40.9803	28.7816	31.6779	19.8074	11.2958