

# Extracting and Transferring Abilities For Building Multi-lingual Ability-enhanced Large Language Models

Anonymous ACL submission

## Abstract

Multi-lingual ability transfer has become increasingly important for the broad application of large language models (LLMs). Existing work highly relies on training with the multi-lingual ability-related data, which may be not available for low-resource languages. To solve it, we propose a **Multi-lingual Ability Extraction and Transfer** approach, named as **MAET**. Our key idea is to decompose and extract language-agnostic ability-related weights from LLMs, and transfer them across different languages by simple addition and subtraction operations without training. Specially, our MAET consists of the extraction and transfer stages. In the extraction stage, we firstly locate key neurons that are highly related to specific abilities, and then employ them to extract the transferable ability-specific weights. In the transfer stage, we further select the ability-related parameter tensors, and design the merging strategy based on the linguistic and ability specific weights, to build the multi-lingual ability-enhanced LLM. To demonstrate the effectiveness of our proposed approach, we conduct extensive experiments on mathematical and scientific tasks in both high-resource lingual and low-resource lingual scenarios. Experiment results have shown that MAET can effectively and efficiently extract and transfer the advanced abilities, and outperform training-based baselines methods. Our code and data will be publicly released.

## 1 Introduction

Large language models (LLMs) have recently shown remarkable performance on various general tasks, *e.g.*, text generation and question answering (Zhao et al., 2023; OpenAI, 2023; Dubey et al., 2024). Despite the success, LLMs are still struggling to solve complex tasks (*e.g.*, mathematical reasoning), which require LLMs to possess specific advanced abilities (*e.g.*, deductive reasoning) and knowledge (*e.g.*, mathematical theory) (Yue et al.,

2024; Lu et al., 2022). To address it and further improve LLMs, existing work either collects the related data to train LLMs (Du et al., 2024; Chen et al., 2024a), or merges the parameters of existing well-performed LLMs to transfer their advanced abilities into one single model (Ilharco et al., 2023; Yadav et al., 2023; Yu et al., 2024a).

Despite the success, it is not easy to collect sufficient training corpus or well-trained LLMs related to specific abilities, especially in multi-lingual scenarios. Especially, some popular languages (*e.g.*, English) have dominated the linguistic expressions of the open web data, and the amount of available domain-specific data for low-resource languages (*e.g.*, Bengali or Telugu) is highly limited (Magueresse et al., 2020; Patzelt, 2024; Mirashi et al., 2024). Fortunately, existing work (Zhao et al., 2024; Schäfer et al., 2024) has revealed that the learned knowledge from one language by LLMs could be inherited and leveraged by other languages. For example, Llama-series LLMs are trained mainly on English texts, while they can also solve the tasks based on other languages. Such a finding has been widely explored in either improving the overall performance of multi-lingual LLMs (Schäfer et al., 2024) or enhancing fine-grained knowledge (Chen et al., 2024a). However, the related work mostly relies on training with ability-related corpus in the target language, which is not always available for low-resource languages.

To conduct a more effective ability transfer, our idea is to learn and extract the “*ability-specific weights*” that preserves the knowledge about specific abilities for the LLM. If such ability-specific and language-specific weights could be decomposed, it is achievable to transfer the required abilities into target languages by just combining the corresponding weights, even building a multi-lingual ability-enhanced LLM like building blocks. Based on this idea, in this paper, we propose a **Multi-lingual Ability Extraction and Transfer** approach,

named as **MAET**. Specifically, our approach consists of two major stages, *i.e.*, ability extraction and transfer stage. In the extraction stage, we first locate the abilities-related neurons and leverage related corpus in a reference language to continually pre-train the LLM on these identified neurons. Then, based on the LLM trained on the general corpus, we devise the formula to extract the ability-specific weights. In the transfer stage, we utilize the ability-related weights to select related parameter tensors, and design a specific model merging strategy by interpolating linguistic and ability-specific weights. In our approach, we only need ability-specific corpus from any rich-resource language and general multi-lingual corpus, which can effectively mitigate the data scarcity issues in low-resource languages.

To assess the effectiveness of our approach, we conduct the evaluation based on two comprehensive reasoning benchmarks, namely Multi-lingual Grade School Math (MGSM) (Shi et al., 2023) and science tasks from multi-lingual MMLU (Lai et al., 2023) as the evaluation benchmarks. According to the evaluation results, the proposed approach MAET outperforms other competitive baseline methods (*e.g.*, continual pre-training (Gururangan et al., 2020) and model merging methods with task vectors (Ilharco et al., 2023), achieving the 9.1% relative improvement compared to the base LLM. In addition, our approach can work well with relatively fewer training data, demonstrating an improved efficiency in practice.

## 2 Related Work

We introduce the related work from the following three perspectives:

**Continual Pre-training.** Although LLMs have shown remarkable performance on various downstream work, they still struggle in several specific tasks, *e.g.*, complex reasoning tasks (Paster et al., 2024; Shao et al., 2024) or low-resource lingual scenarios (Hedderich et al., 2021; Panchbhai and Pankanti, 2021). To adapt LLMs pre-trained on the general corpus to multi-lingual scenarios or specific tasks, existing work (Luo et al., 2022; Taylor et al., 2022; Zhao et al., 2022; Zhang et al., 2024a) has collected the corresponding corpus to continually pre-train (CPT) LLMs. During the continual pre-training process, the mixture strategy between the general corpus and the ability-related corpus should be carefully considered to prevent hurting

the general abilities of LLMs (Ye et al., 2024; Xie et al., 2023; Chen et al., 2024a; Siriwardhana et al., 2024). However, previous study (Chang et al., 2024; Lu et al., 2023) has pointed out that it is difficult to collect the required corpus, especially for low-resource language corpus. Therefore, synthesizing data from powerful LLMs is widely utilized to expand the task-specific training corpus (Chen et al., 2021b; Yu et al., 2024b; Zhou et al., 2024a). Besides, because of the limitation of computation resources, a series of approaches (Hu et al., 2022; Li and Liang, 2021; Dettmers et al., 2023) only train several parameters to reduce the expenses. In this work, we focus on adapting LLMs to multi-lingual complex reasoning scenarios through continually pre-training LLMs on the single-lingual task-specific corpus.

**Knowledge Editing.** According to lottery ticket hypothesis (Frankle and Carbin, 2019), training a small number of model parameters will achieve comparable or even better performance on downstream tasks. Existing study (Du et al., 2024; Wang et al., 2024b; Gong et al., 2024) has leveraged the inner information of LLMs, *e.g.*, gradient or cosine similarity between different hidden states, to select and train the related sub-network. Besides, the probe (*i.e.*, a newly initialized parameter) can be implemented to detect the knowledge of LLMs and process targeted repair (Wang et al., 2024a; Jiang et al., 2024). However, Since these approaches need to calculate and select a sub-network of each training instance, which might cause the instability of the training process, several study (Chen et al., 2024b; Zhang et al., 2024b) pointed out that the task-related sub-network can be determined before the training process, and only updating the value of the corresponding neurons can achieve better performance. In this work, we focus on editing the task-specific neurons of LLMs to improve the corresponding capacities in multi-lingual tasks.

**Model Merging.** Given that the CPT process will bring huge computational expenses, previous work leveraged model merging techniques to integrate different abilities (*e.g.*, mathematical reasoning and code synthesizing) into one model (Yang et al., 2024; Xu et al., 2024b; Stoica et al., 2024). During the merging process, the interference between different LLMs might be conflict with each other and affect the final performance. Therefore, researchers proposed the clip (Yadav et al., 2023) or

randomly dropout (Yu et al., 2024a) mechanism to mitigate the performance decrease. Moreover, the selection of the hyper-parameters (*e.g.*, weight of each model) is the challenge of the model merging process, and previous work (Zhou et al., 2024b; Matena and Raffel, 2022) utilized the inner parameters of LLMs or external matrixes to determine the hyper-parameters. Furthermore, a series of work has studied improving the reasoning ability of LLMs in non-English scenarios by merging the reasoning-specialized model and multi-lingual model (Huang et al., 2024; Yoon et al., 2024). Inspired by the above work, we try to locate the task-related sub-networks of LLMs and transfer the advanced abilities.

### 3 Preliminary

Despite that LLMs exhibit remarkable performance on general tasks, they still have limited advanced abilities, *e.g.*, mathematical and scientific reasoning abilities. A typical approach to enhance these abilities is to continually pre-train (CPT) LLMs with ability-related corpus. However, such training data might not always be available or sufficient, especially for minor domains (*e.g.*, Bengali). In this work, we focus on the task of *ability extraction and transfer* by continual pre-training and merging LLMs. Concretely, LLMs are trained on the collected corpus from a certain domain, and we aim to only transfer its learned advanced capabilities to target domains (Zhuang et al., 2021; Farahani et al., 2021) without further training. In this work, we study the cross-lingual scene where the linguistic-agnostic advanced ability and linguistic abilities should be extracted and transferred, to build a unified multi-lingual ability-enhanced LLM.

Formally, for a certain ability  $A_i$  and a set of languages  $L = \{L_0, L_1, \dots, L_n\}$ , we assume that the general corpus of all languages can be collected, denoted as  $C_{\text{general}} = \{C_{L_0}, C_{L_1}, \dots, C_{L_n}\}$ , while the ability-related corpus is only available in language  $L_0$  (*i.e.*, English), denoted as  $C_{L_0, A_i}$ . Based on the above corpora, our goal is to extract and transfer the advanced ability  $A_i$  from language  $L_0$  and linguistic abilities from other languages  $L_1, \dots, L_n$ , into a unified LLM.

### 4 Approach

In this section, we propose the **Multi-lingual Ability Extraction and Transfer** approach, named as **MAET**, which can effectively transfer the ad-

vanced abilities from single-lingual LLMs, to build a multi-lingual ability-enhanced LLM. The key motivation of our approach is to identify and extract ability-related neurons or weights, and transfer the target abilities into a LLM in an efficient way. The framework of MAET is presented in Figure 1.

#### 4.1 Ability-related Weights Extraction

In this part, we aim to locate and learn ability-related parameter weights within an LLM, to enable efficient transferring of the ability into other LLMs. Concretely, it consists of two major steps, *i.e.*, key neurons locating and ability-related parameter weights learning, which are detailed in the following.

**Locating the Key Neurons.** The gradient of each neuron in LLMs can be utilized to estimate its correlation degree with specific task ability (Pruthi et al., 2020; Chen et al., 2024b; Xia et al., 2024), we select those with high gradient values as key neurons. To this end, we first use the ability-related corpus  $C'_{L_0, A_i}$  to continually pre-train the LLM, while sampling a small amount to train the model can be also applied to reduce the computation consumption. During training, the LLM learns the language modeling task and each neuron is updated by the gradients associated by the training instances. Due to the high cost of calculating the accumulation of gradient at each training step, we calculate the value changes of the LLM neurons before and after the training process to approximate the importance. Formally, the importance function  $I(A_i, \theta_j)$  of neurons can be computed as:

$$I(A_i, \theta_j) = \sum_{d_k \in C'_{L_0, A_i}} \text{Grad}(\theta_j, d_k) \approx \frac{\|\tilde{\theta}_j - \theta_j\|}{\text{LearningRate}}, \quad (1)$$

where  $d_k$  denotes the  $k$ -th instance of training corpus  $C'_{L_0, A_i}$  and  $\tilde{\theta}_j$  denote the value of the  $j$ -th neuron of LLM after training, respectively. Based on it and inspired by previous work (Yadav et al., 2023), we rank all neurons according to their importance scores, and then select the top  $k_1\%$  ones into the set  $\mathcal{N}_{A_i}$  as the key neurons.

**Learning Ability-related Weights.** Based on the identified key neurons in  $\mathcal{N}_{A_i}$ , we further learn the ability-related parameter weights. Our motivation is to decompose the parameter weights according to their changes *before* and *after* the LLM has mastered a specific ability, which is achievable owing to the modularity and composition nature

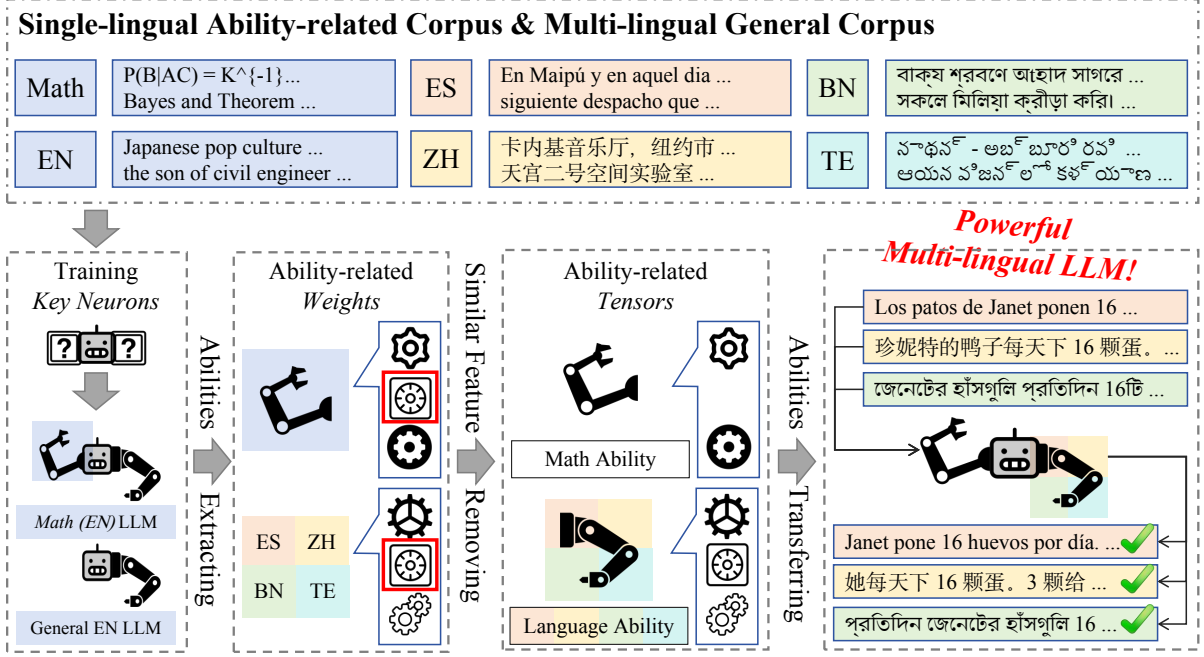


Figure 1: The framework of our approach MAET, including extraction and transfer stages. In the extraction stage, we first locate the key neurons, and utilize the single-lingual ability-related corpus and general corpus to train the LLM on these neurons to obtain the ability-related weight. Then, we remove the parameter tensors related to language knowledge in the ability weight and transfer the remaining to the base LLM. After these stages, we can obtain a powerful LLM with advanced abilities that can solve the corresponding tasks in multi-lingual scenarios.

of the LLM parameter matrices (Yu et al., 2024a; Shazeer et al., 2017). First, we utilize the key neurons locating method mentioned above to extract the ability-related neuron set  $\mathcal{N}_{A_i}$ , and also obtain the language-related neuron set  $\mathcal{N}_{L_0}$  via the same way. Then, we train the LLM with the mixture of ability-related corpus and general corpus on the key neuron set  $\mathcal{N}_{A_i} \cup \mathcal{N}_{L_0}$  and  $\mathcal{N}_{L_0}$  respectively, to obtain two specific models, denoted as  $\text{LLM}_{A_i, L_0}$  with parameters  $\Theta_{A_i, L_0}$  and  $\text{LLM}_{L_0}$  with parameters  $\Theta_{L_0}$ . Next, we measure the parameter changes between the backbone and the trained models, and obtain the ability-related weights via the parameter decomposition operation as:

$$R(A_i) = \alpha \cdot \underbrace{(\Theta_{A_i, L_0} - \Theta_o)}_{\text{Ability \& language difference}} - \beta \cdot \underbrace{(\Theta_{L_0} - \Theta_o)}_{\text{Language difference}}, \quad (2)$$

where  $\alpha$  and  $\beta$  are tunable coefficients to balance the two parts of weight differences, and  $\Theta_o$  denote the original parameters of the LLM, which serves as the reference for parameter decomposition. As we only train the parameters within the neuron set, its weight difference should preserve the knowledge about the corresponding ability. Thus, it can be regarded as the *ability-related parameter representations*, and is promising to transfer the ability

into other LLMs by the addition operation.

## 4.2 Multi-lingual Ability Transfer

After obtaining the ability-related weights, we utilize them to transfer and integrate the abilities, to build a multi-lingual ability-enhanced LLM.

**Ability-related Parameter Tensor Selection.** Although we can locate the ability-related key neurons, it is still hard to avoid the involvement of irrelevant ones. Our empirical studies have found that neuron-level features are easy to be affected by the noisy data. Therefore, we consider identifying ability-related parameter tensors, which correspond to the parameter matrices within the LLM. Specifically, we firstly leverage the ability-related weights of languages  $R(L_1), \dots, R(L_n)$  to obtain the multi-lingual weight  $R_{Lang}$ . Given that large models have varying levels of proficiency in different languages, we use the hyper-parameters  $\mu_1, \dots, \mu_n$  to tune this process as:

$$R_{Lang} = \sum_{i=1}^n \mu_i \cdot R(L_i), \quad (3)$$

where  $R(L_i)$  preserves the linguistic ability of language  $L_i$  learned based on Equation 2. Therefore,  $R_{Lang}$  can be considered as the general language



ability of LLMs that spans multiple languages. As we aim to find the parameter tensors that have low linguistic effects but focus on the desired abilities (*e.g.*, mathematical reasoning), we rank all the tensors according to their similarities with  $R_{Lang}$ , and pick up the last  $k_2\%$  ones. Formally, for tensor  $\tau_i$ , we calculate the cosine similarity of this parameter between  $R(A_i)$  and  $R_{Lang}$ , as follows,

$$S(\tau_i) = \text{sim}(R(A_i)[\tau_i], R_{Lang}[\tau_i]), \quad (4)$$

where we use the cosine similarity to implement the similarity function  $\text{sim}(\cdot)$ . After obtaining the similarity of all tensors, we rank them in a descending order based on the similarity values, and then select the last  $k_2\%$  parameters into the set  $\mathcal{T}$  as the ability-related parameters.

### Building Multi-lingual Ability-enhanced LLM.

Based on the selected ability-related tensors  $\mathcal{T}$ , we design the model merging process by interpolating ability weights and multi-lingual weights, to build the multi-lingual ability-enhanced LLM. Formally, the final parameter tensors of the target LLM are computed as:

$$\tilde{\tau}_i = \tau_i^{(o)} + \begin{cases} \gamma \cdot R(A_i)[\tau_i] + \eta \cdot R_{Lang}[\tau_i], & \tau_i \in \mathcal{T} \\ R_{Lang}[\tau_i], & \tau_i \notin \mathcal{T} \end{cases}, \quad (5)$$

where  $\tau_i^{(o)}$  denotes the original value of parameter tensor  $\tau_i$ , and  $\gamma$  and  $\eta$  are tunable hyper-parameters. This formula can be explained in two different cases. When a parameter tensor serves as the major role for specific abilities, we update it by adding both ability- and linguistic-related weights; otherwise, we simply enhance it with multi-lingual weights. In this way, we can derive a more powerful LLM that is equipped with the multi-lingual abilities and specific advanced abilities.

### 4.3 The Overall Procedure

To better demonstrate our approach, we present key concepts in Table 4 for further clarifying and provide the complete procedure in Algorithm 1 in the pseudo-code form. The procedure of MAET consists of two main stages, *i.e.*, ability-related weights extraction and multi-lingual ability transferring. For the extraction stage, we first utilize the accumulated gradient to estimate the importance of each neuron by Equation 1. Then, we leverage the model trained on the general corpus to remove the influence of language and obtain the ability-related weight through Equation 2. In the

Approaches	MLAR	TPara	AC	AT
CPT	Yes	Full	No	No
MoE	Yes	Full	No	No
LoRA	Yes	Low-Rank	No	No
MoL	Yes	Low-Rank	No	No
TV	Yes	Full	Yes	No
MAET	No	Ability-related	Yes	Yes

Table 1: The difference between our MAET and the methods in previous work (*i.e.*, CPT (Hu et al., 2022), Mixture-of-Expert (MoE) (Shazeer et al., 2017), LoRA (Hu et al., 2022), Mixture-of-LoRA (MoL) (Feng et al., 2024), and Task Vector (TV) (Ilharco et al., 2023). MLAR, TPara, AC, and AT denote the abbreviation of multi-lingual ability-related corpus, parameters for training, ability composition, and ability transfer.

transfer stage, we utilize Equation 3 and Equation 4 to obtain the multi-lingual weight and identify the ability-related parameter tensors in LLM. After it, we leverage Equation 5 to fulfill the multi-lingual ability transfer, to build the multi-lingual ability-enhanced LLM.

To highlight the difference between our approach and previous work, we present the comparison of these methods in Table 1. To adapt LLMs to multi-lingual scenarios, most of the existing methods (*e.g.*, CPT and TV) require the multi-lingual ability-related corpus (*i.e.*, ability-related corpus is required for each language) for training the LLM parameters. In comparison, our proposed approach only trains and modifies the ability-related parameters, which can efficiently focus on enhancing the specific ability. A major novelty of our work is that we identify the key units and implement the sparse update in the model training and merging procedure, which can effectively decompose, extract, and transfer the abilities of LLMs. In addition, compared with the LoRA-based methods (*i.e.*, LoRA and MoL) that also sparsely update the LLM parameters, our approach selectively updates the ability-related neurons, while LoRA-based methods utilize the low-rank matrices to approximate the original parameters.

## 5 Experiment

### 5.1 Experimental Settings

In this part, we introduce the details of our experimental settings, including the datasets utilized in the training and evaluation process, and the baseline methods. Moreover, we present the implementation details of our approach in Appendix B.

**Datasets.** In this work, we focus on transferring the advanced abilities (*i.e.*, mathematical and scientific reasoning abilities) of LLMs from English scenarios to multi-lingual scenarios, including high-resource languages (*i.e.*, Spanish and Chinese) and low-resource languages (*i.e.*, Bengali and Telugu). For the training corpus, we extract the corpus of the corresponding languages from the dataset proposed by previous work (Yang et al., 2023; Scao et al., 2022; Laurençon et al., 2022) as the general training corpus, and utilize OpenWebMath (Paster et al., 2024) and the arXiv papers (Soldaini et al., 2024) as the ability-related corpus for mathematical tasks and scientific tasks respectively. For the evaluation benchmark, we follow the evaluation settings in previous work (OpenAI, 2023), utilizing *Multi-lingual Grade School Math (MGSM)* (Shi et al., 2023) and science tasks from *multi-lingual MMLU* (Lai et al., 2023) (*i.e.*, college biology, college chemistry, college physics, high school biology, high school chemistry, and high school physics) as the downstream tasks for multi-lingual scenarios. The statistical information of the datasets is presented in Table 6.

**Baselines.** In our evaluation, a baseline can be represented as three parts, *i.e.*, training parameters, training approach, and training data. First, we conduct the full parameters training and the LoRA training (Hu et al., 2022) in our evaluation, denoted as the “*F*” and “*L*” at the prefix of the training approaches, respectively. For the training approach, we employ *continual pre-training (CPT)* (Gururangan et al., 2020), *domain adaption (DA)* (Taylor et al., 2022), and *model merging with task vector (TV)* (Ilharco et al., 2023). Besides, for the training data, “*L*”, “*A*”, and “*T*” refer to the multi-lingual general corpus, the single-lingual ability-related corpus, and the translated multi-lingual ability-related corpus from GPT-4o (Hurst et al., 2024), respectively. Moreover, to conduct a more comprehensive evaluation, we also present the performance of different LLMs, *i.e.*, Baichuan-2 7B (Yang et al., 2023), Mistral 7B (Jiang et al., 2023), LLaMA-2 7B (Touvron et al., 2023), and LLaMA-3 8B (Dubey et al., 2024).

## 5.2 Main Results

To comprehensively evaluate our proposed MAET, we employ MAET on mathematical and scientific tasks in multi-lingual scenarios and present the results in Table 2.

First, MAET outperforms other baselines in the average performance of all downstream tasks. In our experiment, continual pre-training LLMs on the mixture of multi-lingual general corpus and single-lingual ability-related corpus (*i.e.*, F-CPT<sub>L&A</sub>) can enhance the specific ability of LLMs, achieving the second best performance. However, when adapting LLMs to a new domain or enhancing a new ability of LLM, CPT-based methods should retrain the LLMs on the ability-related and multi-lingual corpus, showing that CPT is leaked of transferability and requires more computational resources. For the issue of new domain adapting, MAET only utilizes a small amount of single-lingual ability-related corpus (*i.e.*, English corpus in practice) to obtain the ability weight, which can be employed to transfer the corresponding advanced ability, achieving both effectiveness and efficiency.

Second, LoRA-based methods (Hu et al., 2022) (*e.g.*, L-CPT<sub>L&A</sub>, L-CPT<sub>L</sub>, L-TV) initialize the low-rank matrices and only update these matrices, performing sparsely optimize on LLM. Since the trainable parameters in LoRA represent the whole parameters of LLM rather than ability-related sub-network, it cannot perform well on the multi-lingual scenarios, indicating the failure of the advanced abilities transferring. In contrast, MAET first identifies the ability-related sub-networks and utilizes the corresponding sub-networks to perform the following operations. Because of the decomposing of the inner abilities of LLMs, MAET can help LLMs improve their specific ability.

Third, translation-based methods are the strong baselines to enhance the LLM performance in low-resource languages. In the experiment, we utilize GPT-4o to translate the ability-related corpus from English to other languages, and present the prompt in Appendix D. According to the experimental results in the above table, we can observe that our MAET outperforms the translation-based method. The translation-based method consumes more computational resources and cannot achieve better performance. The reason might be that the translated corpus shares similar knowledge of the specific domain, which makes LLM overfit the corresponding knowledge and cannot really understand the specific knowledge. In contrast, our approach decomposes the scientific ability and language ability, and transfers the scientific ability from one language to another, preventing overfitting, decreasing the expense, and improving performance.

Methods	#Tokens ( $\downarrow$ )	Multilingual Mathematical Tasks ( $\uparrow$ )					Multilingual Scientific Tasks ( $\uparrow$ )				
		ES	ZH	BN	TE	Avg.	ES	ZH	BN	TE	Avg.
Baichuan-2 7B	-	17.20	28.00	4.80	2.40	13.10	42.27	46.43	30.17	26.21	36.27
Mistral 7B	-	38.80	34.40	9.60	2.80	21.40	52.08	45.33	32.91	27.96	39.57
LLaMA-2 7B	-	7.60	12.00	1.60	0.00	5.30	34.16	31.68	24.56	22.15	28.14
LLaMA-3 8B	-	<u>48.40</u>	38.80	28.80	20.40	34.10	55.06	47.24	36.63	29.26	42.05
+ F-CPT <sub>L&amp;A</sub>	20B	46.80	<b>42.00</b>	28.40	<b>27.60</b>	<u>36.20</u>	55.92	<u>48.57</u>	<u>36.84</u>	30.10	<u>42.86</u>
+ L-CPT <sub>L&amp;A</sub>	20B	44.80	37.60	28.80	23.60	33.70	54.77	46.81	36.41	29.88	41.97
+ F-CPT <sub>A&amp;T</sub>	20B	48.00	40.00	28.40	35.50	25.60	53.73	46.30	35.06	31.73	41.71
+ F-CPT <sub>A</sub>	4B	47.20	40.80	20.00	13.20	30.30	51.90	45.71	33.35	29.41	40.09
+ F-CPT <sub>T</sub>	20B	48.00	41.20	27.20	24.40	35.20	50.35	45.36	34.54	<b>34.46</b>	41.18
+ F-CPT <sub>L</sub>	8B	38.80	35.60	28.00	23.60	31.50	53.56	47.14	35.89	30.64	41.81
+ F-CPT <sub>L</sub> & DA	12B	41.60	39.60	<b>34.40</b>	<b>27.60</b>	35.80	52.71	48.05	35.49	28.62	41.11
+ L-CPT <sub>L</sub>	8B	46.40	39.20	28.40	22.80	34.20	55.04	48.09	36.66	30.43	42.56
+ L-CPT <sub>L</sub> & DA	12B	46.80	37.60	28.00	<u>27.20</u>	34.90	55.65	<b>49.10</b>	36.48	29.65	42.72
+ F-TV	12B	42.00	32.40	16.00	10.40	25.20	53.36	46.57	36.70	<u>30.73</u>	41.84
+ L-TV	12B	45.60	39.20	30.80	25.60	35.30	55.46	48.27	36.65	<u>30.44</u>	42.71
+ MAET (Ours)	12B	<b>49.60</b>	<u>41.60</u>	<u>32.40</u>	25.20	<b>37.20</b>	<b>56.20</b>	48.00	<b>37.64</b>	30.38	<b>43.06</b>

Table 2: The performance comparison of different approaches on multilingual mathematical tasks and multilingual scientific tasks. ES, ZH, BN, and TE denote Spanish, Chinese, Bengali, and Telugu, respectively. #Tokens denotes the number of training tokens. The best is in bold and the second best is underlined.

Last, compared with the model merging based approaches (*i.e.*, F-TV and L-TV), experimental results have shown that MAET performs better than these baseline methods, since we decompose the relation between ability and the language of the training corpus. In the previous model merging approaches, they mainly added the parameters of different models to obtain the final model, without considering the the relation between language and abilities. Due to the extraction mechanism of MAET, we mitigate the effect of languages and make the weight more related to ability, which can be transferred in multi-lingual scenarios.

### 5.3 Detailed Analysis

To comprehensively evaluate our proposed approach MAET and analyze its features, we conduct several experiments and detailed analysis in this part, including the ablation study, analysis of the transfer ratio of LLM parameters, and the generalization of MAET.

**Ablation Study.** To assess the effectiveness of each component of our proposed MAET, we conduct the ablation study and present the results in Figure 2. We implement MAET on multi-lingual mathematical and scientific tasks without each module of MAET, *i.e.*, key neurons locating (*i.e.*, Eq. 1), ability weight obtaining (*i.e.*, Eq. 2), ability-related parameter tensor identifying (*i.e.*, Eq. 4), and advanced abilities transferring (Eq. 5). First, in most

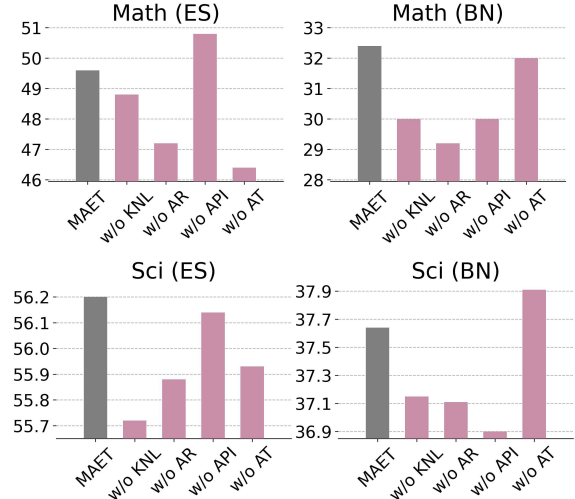


Figure 2: The results of ablation study. KNL, AW, API, and AT denote key neurons locating (Eq. 1), ability weights obtaining (Eq. 2), ability-related tensors identifying (Eq. 4), and advanced abilities transferring (Eq. 5).

downstream scenarios, removing any module of MAET will affect the final performance, which has verified the effectiveness of the MAET process. Second, without ability weight obtaining, *i.e.*, directly utilizing the difference between LLM trained on the ability-related corpus and the backbone LLM as the ability weight, we can observe that the performance is seriously hurt in both scenarios, indicating this process can significantly extract the advanced abilities from the single-lingual

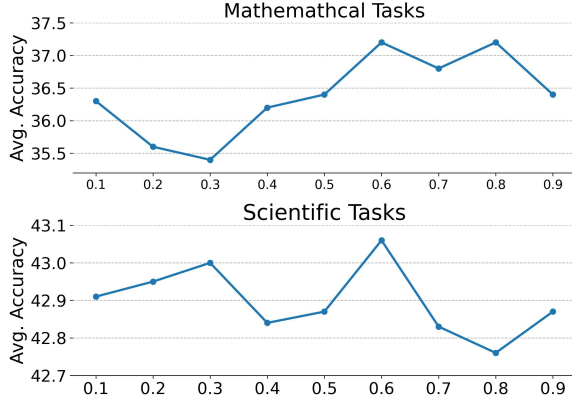


Figure 3: The performance of different proportions for the ability-related parameters identification.

corpus and decrease the influence of the language of the training corpus. Third, comparing the results of the models whether adopting the ability transferring process, experimental results show that LLM with the multi-lingual abilities enhanced cannot well solve multi-lingual mathematical and scientific tasks, and leveraging the ability weight provided by MAET can improve the LLM performance on advanced tasks.

**Ratio of Key Parameters During Transferring Stage.** Identifying and updating the ability-related sub-network of LLMs is the key point of our MAET. We conduct experiments to analyze the influence of the transferring ratio  $k_2\%$  and show the results in Figure 3. Observing the results, the performance of LLM has decreased in a lower and higher ratio of the ability-related parameters identifying process. The main reason is that the lower proportion transfers incomplete knowledge to the model and makes LLM unable to possess the corresponding ability, affecting the performance on the downstream tasks. In contrast, the higher proportion cannot extract the ability weight precisely and will transfer too much language-related knowledge to the model, making the conflict with the LLM inner knowledge and hurting the multi-lingual scenarios performance.

**Out-of-Domain Performance of MAET.** We conduct experiments about adapting mathematical ability on the general LLM through MAET, and assess the performance on out-of-domain (OOD) tasks (*i.e.*, MMLU (Hendrycks et al., 2021), HumanEval (Chen et al., 2021a), MBPP (Austin et al., 2021), and OpenbookQA (Mihaylov et al., 2018)), which can reflect and assess different abilities of

Methods	MMLU	MBPP	OpenbookQA
LLaMA-3 8B	60.85	46.60	65.00
+ CPT	-2.39	-7.00	-3.60
+ MAET	+0.22	+0.80	+0.00

Table 3: The out-of-domain performance of different training methods to train LLaMA-3 8B on OpenWebMath. During the ability-enhancing process, previous methods will hurt the OOD abilities of LLM, while our MAET can maintain the corresponding abilities.

LLMs. Results are presented in Table 3. We can observe that the performance of LLM on all evaluation tasks has decreased through the CPT training process, and the maximum decrease has been achieved 7.32% on the HumanEval task. One of the possible reasons is that LLaMA-3 has been trained on OpenWebMath during pre-training and the CPT process makes it overfit and forget the knowledge of other domains, hurting the performance on OOD tasks. In contrast, our proposed MAET achieves comparable and even better performance with backbone LLM in all downstream scenarios. Since we identify and update the key neurons related to the specific ability, the ability of LLM can be precisely enhanced, and this strategy also helps the OOD tasks needed for mathematical ability, *e.g.*, mathematical tasks in MMLU and MBPP.

## 6 Conclusion

In this paper, we presented MAET, which extracted the advanced ability-related weights from the LLM and supported simple addition and subtraction operations to transfer the ability across different languages. Concretely, MAET included two main stages, *i.e.*, extraction and transfer. For the extraction stage, we located the key neurons and extracted the ability-related weights. Then, in the transfer stage, we identified the key parameter tensors and leveraged them to transfer the advanced ability into other LLMs. In this process, the multi-lingual ability-related training corpus is not required, and the experimental results have shown that our approach outperformed competitive baselines.

As future work, we will consider better methods to identify the ability-related sub-network to decompose the abilities of LLMs and utilize an automated approach to determine the hyper-parameter. Besides, we will implement MAET on larger-scale models, and scenarios with more languages and requiring more abilities to evaluate its effectiveness.



## Limitations

In this section, we discuss the limitations of our work. First, we only implement our approach MAET on 8B LLMs (*i.e.*, LLaMA-3 8B), and do not adopt the LLMs with larger scales (*e.g.*, 13B or 70B LLMs) in the experiment, due to the limitation of computational resources. We will test the effectiveness of our approach on these LLMs in the future. Second, we only evaluate our approach on two downstream tasks (*i.e.*, mathematical and scientific reasoning tasks) in multi-lingual scenarios. Although they are challenging and widely-used testbeds, it is still meaningful to verify our methods on other tasks. Whereas, as we test the performance on diverse high-resource and low-resource languages, it can also provide comprehensive performance estimation for our approach in multi-lingual scenarios. Finally, we do not consider the potential risk and ethics issues that might hurt the alignment of LLMs when using our approach. Actually, our approach is also applicable to transfer the alignment ability across languages. We will investigate to it in the future.

## References

- Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. Program synthesis with large language models. *CoRR*, abs/2108.07732.
- Hsin-Yu Chang, Pei-Yu Chen, Tun-Hsiang Chou, Chang-Sheng Kao, Hsuan-Yun Yu, Yen-Ting Lin, and Yun-Nung Chen. 2024. A survey of data synthesis approaches. *CoRR*, abs/2407.03672.
- Jie Chen, Zhipeng Chen, Jiapeng Wang, Kun Zhou, Yutao Zhu, Jinhao Jiang, Yingqian Min, Wayne Xin Zhao, Zhicheng Dou, Jiaxin Mao, Yankai Lin, Ruihua Song, Jun Xu, Xu Chen, Rui Yan, Zhewei Wei, Di Hu, Wenbing Huang, and Ji-Rong Wen. 2024a. Towards effective and efficient continual pre-training of large language models. *CoRR*, abs/2407.18743.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021a. Evaluating large language models trained on code. *CoRR*, abs/2107.03374.
- Richard J Chen, Ming Y Lu, Tiffany Y Chen, Drew FK Williamson, and Faisal Mahmood. 2021b. Synthetic data in machine learning for medicine and healthcare. *Nature Biomedical Engineering*, 5(6):493–497.
- Zhipeng Chen, Kun Zhou, Wayne Xin Zhao, Jingyuan Wang, and Ji-Rong Wen. 2024b. Low-redundant optimization for large language model alignment. *CoRR*, abs/2406.12606.
- Pei Cheng, Xiayang Shi, and Yinlin Li. 2024. Enhancing translation ability of large language models by leveraging task-related layers. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 6110–6121.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. 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*.
- Wenyu Du, Shuang Cheng, Tongxu Luo, Zihan Qiu, Zeyu Huang, Ka Chun Cheung, Reynold Cheng, and Jie Fu. 2024. Unlocking continual learning abilities in language models. *CoRR*, abs/2406.17245.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Alonsoius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,

732	Jeet Shah, Jelmer van der Linde, Jennifer Billock,	Zixian Huang, Wenhao Zhu, Gong Cheng, Lei Li, and	788
733	Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,	Fei Yuan. 2024. Mindmerger: Efficient boosting	789
734	Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,	LLM reasoning in non-english languages. <i>CoRR</i> ,	790
735	Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph	abs/2405.17386.	791
736	Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,		
737	Kalyan Vasuden Alwala, Kartikeya Upasani, Kate	Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Day-	792
738	Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and	iheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang,	793
739	et al. 2024. The llama 3 herd of models. <i>CoRR</i> ,	Bowen Yu, Kai Dang, An Yang, Rui Men, Fei Huang,	794
740	abs/2407.21783.	Xingzhang Ren, Xuancheng Ren, Jingren Zhou, and	795
		Junyang Lin. 2024. Qwen2.5-coder technical report.	796
741	Abolfazl Farahani, Behrouz Pourshojae, Khaled	<i>CoRR</i> , abs/2409.12186.	797
742	Rasheed, and Hamid R. Arabnia. 2021. A concise		
743	review of transfer learning. <i>CoRR</i> , abs/2104.02144.	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam	798
		Perelman, Aditya Ramesh, Aidan Clark, AJ Os-	799
744	Wenfeng Feng, Chuzhan Hao, Yuewei Zhang, Yu Han,	trow, Akila Welihinda, Alan Hayes, Alec Radford,	800
745	and Hao Wang. 2024. Mixture-of-loras: An efficient	et al. 2024. Gpt-4o system card. <i>arXiv preprint</i>	801
746	multitask tuning method for large language models.	<i>arXiv:2410.21276</i> .	802
747	In <i>Proceedings of the 2024 Joint International Con-</i>		
748	<i>ference on Computational Linguistics, Language Re-</i>	Gabriel Ilharco, Marco Túlio Ribeiro, Mitchell Worts-	803
749	<i>sources and Evaluation, LREC/COLING 2024, 20-25</i>	man, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali	804
750	<i>May, 2024, Torino, Italy</i> , pages 11371–11380.	Farhadi. 2023. Editing models with task arithmetic.	805
		In <i>The Eleventh International Conference on Learn-</i>	806
751	Jonathan Frankle and Michael Carbin. 2019. The lottery	<i>ing Representations, ICLR 2023, Kigali, Rwanda,</i>	807
752	ticket hypothesis: Finding sparse, trainable neural	<i>May 1-5, 2023</i> .	808
753	networks. In <i>7th International Conference on Learn-</i>		
754	<i>ing Representations, ICLR 2019, New Orleans, LA,</i>	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	809
755	<i>USA, May 6-9, 2019</i> .	sch, Chris Bamford, Devendra Singh Chaplot, Diego	810
		de Las Casas, Florian Bressand, Gianna Lengyel,	811
756	Zhuocheng Gong, Ang Lv, Jian Guan, Junxi Yan, Wei	Guillaume Lample, Lucile Saulnier, L��lio Re-	812
757	Wu, Huishuai Zhang, Minlie Huang, Dongyan Zhao,	nard Lavaud, Marie-Anne Lachaux, Pierre Stock,	813
758	and Rui Yan. 2024. Mixture-of-modules: Reinvent-	Teven Le Scao, Thibaut Lavril, Thomas Wang, Timo-	814
759	ing transformers as dynamic assemblies of modules.	th��e Lacroix, and William El Sayed. 2023. Mistral	815
760	<i>CoRR</i> , abs/2407.06677.	7b. <i>CoRR</i> , abs/2310.06825.	816
761	Suchin Gururangan, Ana Marasovic, Swabha	Yuxin Jiang, Yufei Wang, Chuhan Wu, Wanjun Zhong,	817
762	Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,	Xingshan Zeng, Jiahui Gao, Liangyou Li, Xin Jiang,	818
763	and Noah A. Smith. 2020. Don’t stop pretraining:	Lifeng Shang, Ruiming Tang, Qun Liu, and Wei	819
764	Adapt language models to domains and tasks. In	Wang. 2024. Learning to edit: Aligning llms with	820
765	<i>Proceedings of the 58th Annual Meeting of the</i>	knowledge editing. In <i>Proceedings of the 62nd An-</i>	821
766	<i>Association for Computational Linguistics, ACL</i>	<i>nual Meeting of the Association for Computational</i>	822
767	<i>2020, Online, July 5-10, 2020</i> , pages 8342–8360.	<i>Linguistics (Volume 1: Long Papers), ACL 2024,</i>	823
		<i>Bangkok, Thailand, August 11-16, 2024</i> , pages 4689–	824
768	Michael A. Hedderich, Lukas Lange, Heike Adel, Jan-	4705.	825
769	nik Str��tgen, and Dietrich Klakow. 2021. A survey		
770	on recent approaches for natural language process-	Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo,	826
771	ing in low-resource scenarios. In <i>Proceedings of</i>	Thuat Nguyen, Franck Dernoncourt, Ryan A. Rossi,	827
772	<i>the 2021 Conference of the North American Chap-</i>	and Thien Huu Nguyen. 2023. Okapi: Instruction-	828
773	<i>ter of the Association for Computational Linguistics:</i>	tuned large language models in multiple languages	829
774	<i>Human Language Technologies, NAACL-HLT 2021,</i>	with reinforcement learning from human feedback.	830
775	<i>Online, June 6-11, 2021</i> , pages 2545–2568.	In <i>Proceedings of the 2023 Conference on Empirical</i>	831
		<i>Methods in Natural Language Processing, EMNLP</i>	832
776	Dan Hendrycks, Collin Burns, Steven Basart, Andy	<i>2023 - System Demonstrations, Singapore, December</i>	833
777	Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-	<i>6-10, 2023</i> , pages 318–327.	834
778	hardt. 2021. Measuring massive multitask language		
779	understanding. In <i>9th International Conference on</i>	Hugo Lauren��on, Lucile Saulnier, Thomas Wang,	835
780	<i>Learning Representations, ICLR 2021, Virtual Event,</i>	Christopher Akiki, Albert Villanova del Moral,	836
781	<i>Austria, May 3-7, 2021</i> .	Teven Le Scao, Leandro von Werra, Chenghao Mou,	837
		Eduardo Gonz��lez Ponferrada, Huu Nguyen, J��rg	838
782	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	Frohberg, Mario Sasko, Quentin Lhoest, Angelina	839
783	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	McMillan-Major, G��rard Dupont, Stella Biderman,	840
784	Weizhu Chen. 2022. Lora: Low-rank adaptation of	Anna Rogers, Loubna Ben Allal, Francesco De Toni,	841
785	large language models. In <i>The Tenth International</i>	Giada Pistilli, Olivier Nguyen, Somaieh Nikpoor,	842
786	<i>Conference on Learning Representations, ICLR 2022,</i>	Maram Masoud, Pierre Colombo, Javier de la Rosa,	843
787	<i>Virtual Event, April 25-29, 2022</i> . OpenReview.net.	Paulo Villegas, Tristan Thrush, Shayne Longpre, Se-	844
		bastian Nagel, Leon Weber, Manuel Mu��oz, Jian	845

846	Zhu, Daniel van Strien, Zaid Alyafeai, Khalid Al-	Kale. 2024. On importance of pruning and distilla-	903
847	mubarak, Minh Chien Vu, Itziar Gonzalez-Dios,	tion for efficient low resource nlp. <i>arXiv preprint</i>	904
848	Aitor Soroa, Kyle Lo, Manan Dey, Pedro Ortiz	<i>arXiv:2409.14162</i> .	905
849	Suarez, Aaron Gokaslan, Shamik Bose, David Ife-		
850	oluwa Adelani, Long Phan, Hieu Tran, Ian Yu, Suhas	OpenAI. 2023. GPT-4 technical report. <i>CoRR</i> ,	906
851	Pai, Jenny Chim, Violette Lepercq, Suzana Ilic,	abs/2303.08774.	907
852	Margaret Mitchell, Alexandra Sasha Luccioni, and		
853	Yacine Jernite. 2022. The bigscience ROOTS corpus:	Anand Panchbhai and Smarana Pankanti. 2021. Explor-	908
854	A 1.6tb composite multilingual dataset. In <i>Advances</i>	ing large language models in a limited resource sce-	909
855	<i>in Neural Information Processing Systems 35: Annual</i>	nario. <i>2021 11th International Conference on Cloud</i>	910
856	<i>Conference on Neural Information Processing</i>	<i>Computing, Data Science &amp; Engineering (Conflu-</i>	911
857	<i>Systems 2022, NeurIPS 2022, New Orleans, LA, USA,</i>	<i>ence)</i> , pages 147–152.	912
858	<i>November 28 - December 9, 2022.</i>		
859	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	Keiran Paster, Marco Dos Santos, Zhangir Azerbayev,	913
860	Optimizing continuous prompts for generation. In	and Jimmy Ba. 2024. Openwebmath: An open	914
861	<i>Proceedings of the 59th Annual Meeting of the Asso-</i>	dataset of high-quality mathematical web text. In	915
862	<i>ciation for Computational Linguistics and the 11th</i>	<i>The Twelfth International Conference on Learning</i>	916
863	<i>International Joint Conference on Natural Language</i>	<i>Representations, ICLR 2024, Vienna, Austria, May</i>	917
864	<i>Processing, ACL/IJCNLP 2021, (Volume 1: Long</i>	<i>7-11, 2024.</i>	918
865	<i>Papers)</i> , Virtual Event, August 1-6, 2021, pages 4582–		
866	4597. Association for Computational Linguistics.	Tim Patzelt. 2024. Medical concept normaliza-	919
		tion in a low-resource setting. <i>arXiv preprint</i>	920
		<i>arXiv:2409.14579</i> .	921
867	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-	Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund	922
868	Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Pe-	Sundararajan. 2020. Estimating training data influ-	923
869	ter Clark, and Ashwin Kalyan. 2022. Learn to ex-	ence by tracing gradient descent. In <i>Advances in</i>	924
870	plain: Multimodal reasoning via thought chains for	<i>Neural Information Processing Systems 33: Annual</i>	925
871	science question answering. In <i>Advances in Neural</i>	<i>Conference on Neural Information Processing Sys-</i>	926
872	<i>Information Processing Systems 35: Annual Confer-</i>	<i>tems 2020, NeurIPS 2020, December 6-12, 2020,</i>	927
873	<i>ence on Neural Information Processing Systems 2022,</i>	<i>virtual.</i>	928
874	<i>NeurIPS 2022, New Orleans, LA, USA, November 28</i>		
875	<i>- December 9, 2022.</i>	Morgane Rivière, Shreya Pathak, Pier Giuseppe	929
		Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard	930
876	Yingzhou Lu, Huazheng Wang, and Wenqi Wei. 2023.	Hussenot, Thomas Mesnard, Bobak Shahriari,	931
877	Machine learning for synthetic data generation: a	Alexandre Ramé, Johan Ferret, Peter Liu, Pouya	932
878	review. <i>CoRR</i> , abs/2302.04062.	Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos,	933
		Ravin Kumar, Charline Le Lan, Sammy Jerome, An-	934
879	Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng	ton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan	935
880	Zhang, Hoifung Poon, and Tie-Yan Liu. 2022.	Girgin, Nikola Momchev, Matt Hoffman, Shantanu	936
881	Biogpt: generative pre-trained transformer for	Thakoor, Jean-Bastien Grill, Behnam Neyshabur,	937
882	biomedical text generation and mining. <i>Briefings</i>	Olivier Bachem, Alanna Walton, Aliaksei Severyn,	938
883	<i>Bioinform.</i> , 23(6).	Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin	939
		Abdagic, Amanda Carl, Amy Shen, Andy Brock,	940
884	Alexandre Magueresse, Vincent Carles, and Evan	Andy Coenen, Anthony Laforge, Antonia Pater-	941
885	Heetderks. 2020. Low-resource languages: A re-	son, Ben Bastian, Bilal Piot, Bo Wu, Brandon	942
886	view of past work and future challenges. <i>CoRR</i> ,	Royal, Charlie Chen, Chintu Kumar, Chris Perry,	943
887	abs/2006.07264.	Chris Welty, Christopher A. Choquette-Choo, Danila	944
		Sinopalnikov, David Weinberger, Dimple Vijayku-	945
888	Michael Matena and Colin Raffel. 2022. Merging mod-	mar, Dominika Rogozinska, Dustin Herbison, Elisa	946
889	els with fisher-weighted averaging. In <i>Advances in</i>	Bandy, Emma Wang, Eric Noland, Erica Moreira,	947
890	<i>Neural Information Processing Systems 35: Annual</i>	Evan Senter, Evgenii Eltyshev, Francesco Visin,	948
891	<i>Conference on Neural Information Processing Sys-</i>	Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus	949
892	<i>tems 2022, NeurIPS 2022, New Orleans, LA, USA,</i>	Martins, Hadi Hashemi, Hanna Klimczak-Plucinska,	950
893	<i>November 28 - December 9, 2022.</i>	Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda	951
		Mein, Jack Zhou, James Svensson, Jeff Stanway,	952
894	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish	Jetha Chan, Jin Peng Zhou, Joana Carrasqueira,	953
895	Sabharwal. 2018. Can a suit of armor conduct elec-	Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost	954
896	tricity? A new dataset for open book question an-	van Amersfoort, Josh Gordon, Josh Lipschultz,	955
897	swering. In <i>Proceedings of the 2018 Conference on</i>	Josh Newlan, Ju-yeong Ji, Kareem Mohamed, Kar-	956
898	<i>Empirical Methods in Natural Language Processing,</i>	tikeya Badola, Kat Black, Katie Millican, Keelin	957
899	<i>Brussels, Belgium, October 31 - November 4, 2018,</i>	McDonnell, Kelvin Nguyen, Kiranbir Sodhia, Kish	958
900	pages 2381–2391.	Greene, Lars Lowe Sjöstrand, Lauren Usui, Laurent	959
		Sifre, Lena Heuermann, Leticia Lago, and Lilly Mc-	960
901	Aishwarya Mirashi, Purva Lingayat, Srushti Sona-	Nealus. 2024. Gemma 2: Improving open language	961
902	vane, Tejas Padhiyar, Raviraj Joshi, and Geetanjali	models at a practical size. <i>CoRR</i> , abs/2408.00118.	962

963	Teven Le Scao, Angela Fan, Christopher Akiki, El-	an open corpus of three trillion tokens for language	1022
964	lie Pavlick, Suzana Ilic, Daniel Hesslow, Roman	model pretraining research. In <i>Proceedings of the</i>	1023
965	Castagné, Alexandra Sasha Luccioni, François Yvon,	<i>62nd Annual Meeting of the Association for Compu-</i>	1024
966	Matthias Gallé, Jonathan Tow, Alexander M. Rush,	<i>tational Linguistics (Volume 1: Long Papers), ACL</i>	1025
967	Stella Biderman, Albert Webson, Pawan Sasanka Am-	<i>2024, Bangkok, Thailand, August 11-16, 2024</i> , pages	1026
968	manamanchi, Thomas Wang, Benoît Sagot, Niklas	15725–15788. Association for Computational Lin-	1027
969	Muennighoff, Albert Villanova del Moral, Olatunji	guistics.	1028
970	Ruwase, Rachel Bawden, Stas Bekman, Angelina		
971	McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile	George Stoica, Daniel Bolya, Jakob Bjorner, Pratik	1029
972	Saulnier, Samson Tan, Pedro Ortiz Suarez, Vic-	Ramesh, Taylor Hearn, and Judy Hoffman. 2024.	1030
973	tor Sanh, Hugo Laurençon, Yacine Jernite, Julien	Zipit! merging models from different tasks without	1031
974	Launay, Margaret Mitchell, Colin Raffel, Aaron	training. In <i>The Twelfth International Conference</i>	1032
975	Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri	<i>on Learning Representations, ICLR 2024, Vienna,</i>	1033
976	Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg	<i>Austria, May 7-11, 2024</i> .	1034
977	Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue,		
978	Christopher Klammer, Colin Leong, Daniel van Strien,	Tianyi Tang, Wenyang Luo, Haoyang Huang, Dong-	1035
979	David Ifeoluwa Adelani, and et al. 2022. BLOOM:	dong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei,	1036
980	A 176b-parameter open-access multilingual language	and Ji-Rong Wen. 2024. Language-specific neurons:	1037
981	model. <i>CoRR</i> , abs/2211.05100.	The key to multilingual capabilities in large language	1038
		models. In <i>Proceedings of the 62nd Annual Meeting</i>	1039
982	Anton Schäfer, Shauli Ravfogel, Thomas Hofmann,	<i>of the Association for Computational Linguistics (Vol-</i>	1040
983	Tiago Pimentel, and Imanol Schlag. 2024. Lan-	<i>ume 1: Long Papers), ACL 2024, Bangkok, Thailand,</i>	1041
984	guage imbalance can boost cross-lingual generali-	<i>August 11-16, 2024</i> , pages 5701–5715.	1042
985	sation. <i>CoRR</i> , abs/2404.07982.		
986	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas	1043
987	Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu,	Scialom, Anthony Hartshorn, Elvis Saravia, An-	1044
988	and Daya Guo. 2024. Deepseekmath: Pushing the	drew Poulton, Viktor Kerkez, and Robert Stojnic.	1045
989	limits of mathematical reasoning in open language	2022. Galactica: A large language model for science.	1046
990	models. <i>CoRR</i> , abs/2402.03300.	<i>CoRR</i> , abs/2211.09085.	1047
991	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziars,		
992	Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	1048
993	Dean. 2017. Outrageously large neural networks:	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	1049
994	The sparsely-gated mixture-of-experts layer. <i>arXiv</i>	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	1050
995	preprint <i>arXiv:1701.06538</i> .	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-	1051
996	Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang,	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	1052
997	Suraj Srivats, Soroush Vosoughi, Hyung Won Chung,	Jude Fernandes, Jeremy Fu, Wenyan Fu, Brian Fuller,	1053
998	Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das,	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	1054
999	and Jason Wei. 2023. Language models are multi-	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	1055
1000	lingual chain-of-thought reasoners. In <i>The Eleventh</i>	Inan, Marcin Kardas, Viktor Kerkez, Madihan Khabza,	1056
1001	<i>International Conference on Learning Representa-</i>	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	1057
1002	<i>tions, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> .	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	1058
1003	Shamane Siriwardhana, Mark McQuade, Thomas Gau-	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	1059
1004	thier, Lucas Atkins, Fernando Fernandes Neto, Luke	tiniet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	1060
1005	Meyers, Anneketh Vij, Tyler Odenthal, Charles God-	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	1061
1006	dard, Mary MacCarthy, and Jacob Solawetz. 2024.	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	1062
1007	Domain adaptation of llama3-70b-instruct through	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	1063
1008	continual pre-training and model merging: A com-	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	1064
1009	prehensive evaluation. <i>CoRR</i> , abs/2406.14971.	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	1065
1010	Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	1066
1011	Schwenk, David Atkinson, Russell Authur, Ben Bo-	Melanie Kambadur, Sharan Narang, Aurélien Ro-	1067
1012	gin, Khyathi Raghavi Chandu, Jennifer Dumas, Yanai	driguez, Robert Stojnic, Sergey Edunov, and Thomas	1068
1013	Elazar, Valentin Hofmann, Ananya Harsh Jha, Sachin	Scialom. 2023. Llama 2: Open foundation and fine-	1069
1014	Kumar, Li Lucy, Xinxin Lyu, Nathan Lambert, Ian	tuned chat models. <i>CoRR</i> , abs/2307.09288.	1070
1015	Magnusson, Jacob Morrison, Niklas Muennighoff,		
1016	Aakanksha Naik, Crystal Nam, Matthew E. Peters,	Huanqian Wang, Yang Yue, Rui Lu, Jingxin Shi, An-	1071
1017	Abhilasha Ravichander, Kyle Richardson, Zejiang	drew Zhao, Shenzhi Wang, Shiji Song, and Gao	1072
1018	Shen, Emma Strubell, Nishant Subramani, Oyvind	Huang. 2024a. Model surgery: Modulating llm’s	1073
1019	Tafjord, Evan Pete Walsh, Luke Zettlemoyer, Noah A.	behavior via simple parameter editing. <i>CoRR</i> ,	1074
1020	Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groen-	abs/2407.08770.	1075
1021	eveland, Jesse Dodge, and Kyle Lo. 2024. Dolma:		
		Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi,	1076
		Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi	1077
		Yang, Jindong Wang, and Huajun Chen. 2024b.	1078
		Detoxifying large language models via knowledge	1079
		editing. In <i>Proceedings of the 62nd Annual Meeting</i>	1080



1081	<i>of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 3093–3118.	
1082		
1083		
1084	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
1085		
1086		
1087		
1088		
1089		
1090		
1091		
1092		
1093		
1094		
1095		
1096	Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. LESS: selecting influential data for targeted instruction tuning. In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> .	
1097		
1098		
1099		
1100		
1101		
1102	Chaojun Xiao, Zhengyan Zhang, Chenyang Song, Dazhi Jiang, Feng Yao, Xu Han, Xiaozhi Wang, Shuo Wang, Yufei Huang, Guanyu Lin, et al. 2024. Configurable foundation models: Building llms from a modular perspective. <i>arXiv preprint arXiv:2409.02877</i> .	
1103		
1104		
1105		
1106		
1107	Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V. Le, Tengyu Ma, and Adams Wei Yu. 2023. Doremi: Optimizing data mixtures speeds up language model pretraining. In <i>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</i> .	
1108		
1109		
1110		
1111		
1112		
1113		
1114		
1115	Haoyun Xu, Runzhe Zhan, Derek F. Wong, and Lidia S. Chao. 2024a. Let’s focus on neuron: Neuron-level supervised fine-tuning for large language model. <i>CoRR</i> , abs/2403.11621.	
1116		
1117		
1118		
1119	Zhengqi Xu, Ke Yuan, Huiqiong Wang, Yong Wang, Mingli Song, and Jie Song. 2024b. Training-free pretrained model merging. <i>CoRR</i> , abs/2403.01753.	
1120		
1121		
1122	Prateek Yadav, Derek Tam, Leshem Choshen, Colin A. Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models. In <i>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</i> .	
1123		
1124		
1125		
1126		
1127		
1128		
1129	Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, Juntao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng	
1130		
1131		
1132		
1133		
1134		
1135		
1136		
1137		
	Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. 2023. Baichuan 2: Open large-scale language models. <i>CoRR</i> , abs/2309.10305.	1138
		1139
		1140
		1141
		1142
		1143
		1144
	Enneng Yang, Li Shen, Guibing Guo, Xingwei Wang, Xiaochun Cao, Jie Zhang, and Dacheng Tao. 2024. Model merging in llms, mllms, and beyond: Methods, theories, applications and opportunities. <i>CoRR</i> , abs/2408.07666.	1145
		1146
		1147
		1148
		1149
	Jiasheng Ye, Peiju Liu, Tianxiang Sun, Yunhua Zhou, Jun Zhan, and Xipeng Qiu. 2024. Data mixing laws: Optimizing data mixtures by predicting language modeling performance. <i>CoRR</i> , abs/2403.16952.	1150
		1151
		1152
		1153
	Dongkeun Yoon, Joel Jang, Sungdong Kim, Seung-gone Kim, Sheikh Shafayat, and Minjoon Seo. 2024. Langbridge: Multilingual reasoning without multilingual supervision. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 7502–7522. Association for Computational Linguistics.	1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
	Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024a. Language models are super mario: Absorbing abilities from homologous models as a free lunch. In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> .	1162
		1163
		1164
		1165
		1166
		1167
	Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2024b. Metamath: Bootstrap your own mathematical questions for large language models. In <i>The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024</i> .	1168
		1169
		1170
		1171
		1172
		1173
		1174
	Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhui Chen. 2024. Mammoth2: Scaling instructions from the web. <i>CoRR</i> , abs/2405.03548.	1175
		1176
		1177
	Wenxuan Zhang, Hou Pong Chan, Yiran Zhao, Mahani Aljunied, Jianyu Wang, Chaoqun Liu, Yue Deng, Zhiqiang Hu, Weiwen Xu, Yew Ken Chia, Xin Li, and Lidong Bing. 2024a. <a href="#">Seallms 3: Open foundation and chat multilingual large language models for southeast asian languages</a> . <i>CoRR</i> , abs/2407.19672.	1178
		1179
		1180
		1181
		1182
		1183
	Zhihao Zhang, Jun Zhao, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024b. Unveiling linguistic regions in large language models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 6228–6247.	1184
		1185
		1186
		1187
		1188
		1189
		1190
	Jun Zhao, Zhihao Zhang, Luhui Gao, Qi Zhang, Tao Gui, and Xuanjing Huang. 2024. Llama beyond english: An empirical study on language capability transfer. <i>CoRR</i> , abs/2401.01055.	1191
		1192
		1193
		1194

Wayne Xin Zhao, Kun Zhou, Zheng Gong, Beichen Zhang, Yuanhang Zhou, Jing Sha, Zhigang Chen, Shijin Wang, Cong Liu, and Ji-Rong Wen. 2022. Jiuzhang: A chinese pre-trained language model for mathematical problem understanding. In *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 - 18, 2022*, pages 4571–4581.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen 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. 2023. A survey of large language models. *CoRR*, abs/2303.18223.

Kun Zhou, Beichen Zhang, Jiapeng Wang, Zhipeng Chen, Wayne Xin Zhao, Jing Sha, Zhichao Sheng, Shijin Wang, and Ji-Rong Wen. 2024a. Jiuzhang3.0: Efficiently improving mathematical reasoning by training small data synthesis models. *CoRR*, abs/2405.14365.

Yuyan Zhou, Liang Song, Bingning Wang, and Weipeng Chen. 2024b. Metagpt: Merging large language models using model exclusive task arithmetic. *CoRR*, abs/2406.11385.

Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. 2021. A comprehensive survey on transfer learning. *Proc. IEEE*, 109(1):43–76.

## A Empirical Study

A surge of work (Zhang et al., 2024b; Xiao et al., 2024; Tang et al., 2024) has pointed out that LLMs sparsely activate the specific sub-modules to perform corresponding tasks. Based on these findings, we conduct empirical experiments to explore whether the specific sub-module, which is related to advanced abilities, can be extracted and transferred. We utilize the forum corpus (*i.e.*, Zhihu for Chinese forum corpus and Reddit for English forum corpus) to continually pre-train LLMs, and then assess the training performance (*i.e.*, the value of loss function) and similarity of LLM neurons.

The forum corpus can be considered as containing the question-answering (QA) ability, which is necessary and important for LLMs. The results from Figure 4a have shown that only training the top 5% relevant neurons of LLMs can achieve the lower training loss and fit into the training set more quickly, indicating that LLMs contain the sub-module corresponding to the QA ability. Moreover, from Figure 4b and Figure 4c, we can observe that the LLM trained on Zhihu has shown higher similarity with the LLM trained on Reddit than the LLM trained on Github (*i.e.*, lower L1 Norm and higher cosine similarity), and the cosine similarity of different layers in LLM are largely different.

According to the above results, we have found that the different sub-networks of LLMs control the different abilities, and precisely selecting the correct sub-module of LLMs will help the extraction of advanced abilities from the single-lingual corpus and the transfer of these abilities to multi-lingual scenarios. Concretely, although Zhihu and Reddit are in different languages, they will influence the similar sub-modules of LLM and make these sub-networks show high similarity with each other. These sub-networks can be referred to the ability-related sub-networks, which are slightly influenced by languages.

## B Implementation Details

In the experiment, we adapt LLaMA-3 8B as the backbone LLM, and employ Transformers (Wolf et al., 2020) and DeepSpeed () framework to perform the training process. And we also present the evaluation results of different backbone LLM (*i.e.*, Qwen2.5 0.5B (Hui et al., 2024) and Gemma2 2B (Rivière et al., 2024)) in Appendix E. For the training process, the learning rate, batch size, and training step are set as  $5 \times 10^{-5}$ , 1M tokens, and 2B

Concepts	Meaning
Key Neurons	Neuron refers to one of the trainable values of the tensors in LLMs. As previous work pointed out (Xu et al., 2024a), different neurons might control the different abilities of LLMs. Following this finding, in our work, we define the neurons that control the specific ability as the "Key Neurons". Key neurons can be regarded as a set without duplication, and a neuron belonging to the set means that this neuron can control the specific ability (Chen et al., 2024b). During the following training process, only the neurons belonging to the key neurons will be trained and optimized.
Ability-related Weights	Ability-related weights refer to the value of the whole neuron in LLM, which can represent the corresponding ability of LLM (Yu et al., 2024a; Ilharco et al., 2023). In MAET, we obtain the ability-related weights through equation 2. The ability-related weights contain the value of all neurons. Since only the key neurons will be trained during the training process, the value of the neurons not belonging to key neurons is zero in the ability-related weights.
Ability-related Tensors	Ability-related tensors can be regarded as a set of LLM tensors, which is related to the corresponding ability. Previous work has studied how the LLM layers influence the ability (Cheng et al., 2024). Different from key neurons, ability-related tensors focus on higher-level information, integrating the sparse neurons into a coarser-grained element (Xiao et al., 2024). A tensor belonging to the ability-related tensors denotes that this tensor is highly related to the corresponding ability and can control this ability.
Language-specific Weights	Similar to the ability-related weights, language-specific weights also refer to the value of the whole neurons in LLMs (Zhang et al., 2024b). However, language-specific weights represent the language abilities of LLM that include multiple abilities (i.e., one language can be regarded as one ability) (Tang et al., 2024), and the method of obtaining them is also different from ability-specific weights. In MAET, we first calculate the ability-related weights of each language and then Integrating these weights together to obtain the language-specific.

Table 4: The key concepts of our approach.

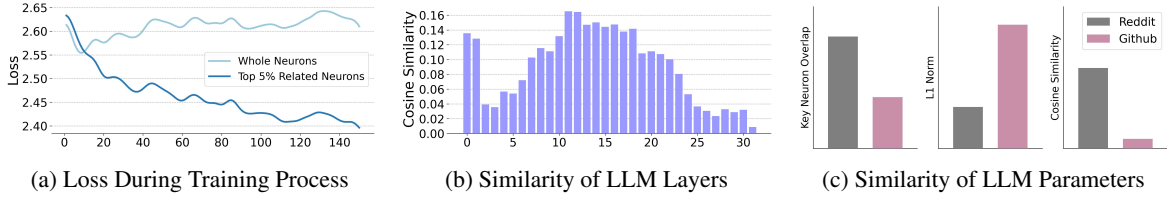


Figure 4: The results of empirical experiments. We present the loss of different training methods during the training process, the cosine similarity of LLM layers after being trained on Zhihu and Reddit, and the similarity of LLMs being trained on different training corpus.

tokens, respectively. Besides, for the key neurons locating, we select the top 5% relevant neurons as the key neuron set  $\mathcal{N}$  for both stages and identify the last 80% and 60% similar tensor as the key sub-network  $\mathcal{T}$  for mathematical reasoning tasks and scientific reasoning tasks respectively.

**Hyper-parameters Selection.** we released all of the hyper-parameters during our experiment in Table 5, to reproduce our proposed approach better. The hyperparameters discussed in the paper can be categorized into two types: training-related parameters (e.g., learning rate, batch size) and training-independent parameters (i.e.,  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\eta$ , and  $\mu$ ). Training-related parameters do not require extensive hyperparameter tuning, as existing studies (Dubey et al., 2024; Hui et al., 2024) provide

clear guidelines for setting them. On the other hand, training-independent parameters are used to construct ability-related weights, tensors, and language-specific weights. These techniques are similar to those employed in model merging (Ilharco et al., 2023; Yadav et al., 2023), and the hyperparameter setting approach outlined in the paper can be applied. A limited number of hyperparameter sets can be defined and validated, as the process primarily involves simple additions and subtractions of model parameters, making it computationally inexpensive.

## C Details of Dataset

We present the statistical information of the datasets in Table 6. We mainly consider English,

---

**Algorithm 1:** The complete procedure of our proposed approach.

---

**Input :** Single-lingual ability-related corpus  $C_{L_0, A_i}$ , multi-lingual general corpus  $C_{L_0}, C_{L_1}, \dots, C_{L_n}$ , and the parameters of the backbone model  $\Theta_o$ .  
**Output :** A well-trained multi-lingual ability-enhanced LLM.

```

// Ability-related Weights Extraction
1  $\theta' \leftarrow \text{CPT}(C_{L_0, A_i}, \Theta_o)$ ;
2 for  $j$ -th neuron in  $\Theta_o$  do
3   | Calculate the importance score of the corresponding neuron using Eq. 1;
4   Identify the key neuron set  $\mathcal{N}_{A_i}$ ;
5    $\text{LLM}_{A_i, L_0} \leftarrow \text{CPT}(C_{L_0, A}, \Theta_o, \mathcal{N}_{A_i} \cup \mathcal{N}_{L_0})$ ;
6    $\text{LLM}_{L_0} \leftarrow \text{CPT}(C_{L_0}, \Theta_o, \mathcal{N}_{L_0})$ ;
7   Learning the ability-related weight  $R(A_i)$  using Eq. 2;

// Multi-lingual Ability Transfer
8   Obtaining the multi-lingual weight  $R_{Lang}$  using Eq. 3;
9   for  $j$ -th parameter tensor in LLM do
10    | Calculate the correlation using Eq. 4;
11   Identify the ability-related parameters  $\mathcal{T}$ ;
12   Transfer the ability to multi-lingual scenarios using Eq. 5;
13   Obtain the well-trained multi-lingual ability-enhanced LLM.

```

---

Spanish, Chinese, Bengali, and Telugu in our experiment, and utilized English as the in-domain language while others as the out-of-domain languages. For the evaluation datasets, we select MGSM and multi-lingual MMLU as the evaluation benchmarks, which contain the parallel data in different languages and are useful for multi-lingual complex tasks evaluation.

## D Prompt for Translation

You should translate the following text from English to {TARGET LANGUAGE} and should not modify the latex code or website code. You should not add any details that are not mentioned in the original text.

```

## English
{ENGLISH TEXT}

## {TARGET LANGUAGE}

```

## E Performance of Small Scale LLMs

We conduct the different LLMs with different sizes (*i.e.*, Qwen2.5-0.5B and Gemma2-2B) in our experiment to valid the practicality of our approach. We assess MAET and baselines on multi-lingual scientific reasoning tasks and present the evaluation

results in Table 7. Comparing the performance of MAET and the baseline methods, we can observe that MAET can also enhance the performance of small scale models and outperform competitive baselines. Therefore, the evaluation results have shown the effectiveness of MAET and verified that MAET is a general LLM enhancement technology.

## F Ability-related Sub-networks of LLM

To assess and probe the ability-related sub-networks of LLMs, we only transfer the specific tensors (*i.e.*, tensors in self-attention and MLP mechanism) from the ability weight to the final models through Eq. 5, to analyze the LLM inner abilities. The experimental results are presented in Table 8. From the experiment, we can observe that although the proportion of MLP layers (41.38%) is lower than the attention layers (45.26%), only transferring the MLP layers outperforms transferring the attention layers, indicating that the MLP layers are more related to the advanced abilities and stores the corresponding knowledge. In the MLP layers of LLM, the gate mechanism (*i.e.*, MLP Gate) will control the transmission of information and the down project mechanism (*i.e.*, MLP Down) will integrate the knowledge from previous layers, so that transferring the MLP layers can achieve better performance on the downstream tasks.



Stage	Hyper-Parameter	Mathematical Tasks	Scientific Tasks
Extracting	Learning Rate	$5 \times 10^{-5}$	$5 \times 10^{-5}$
	Batch Size	1M Tokens	1M Tokens
	Training Steps	2B Tokens	2B Tokens
	$\alpha$ in Extraction	0.8	0.8
	$\beta$ in Extraction	0.2	0.2
	Ratio of Key Neurons	5%	5%
Transferring	Learning Rate	$5 \times 10^{-5}$	$5 \times 10^{-5}$
	Batch Size	1M Tokens	1M Tokens
	Training Steps	2B Tokens	2B Tokens
	$\gamma$ in Transferring	0.2	0.2
	$\eta$ in Transferring	1.0	1.0
	Ratio of Key Neurons	80%	60%
	$\mu$ for Spanish	1.5	1.5
	$\mu$ for Chinese	2.0	2.0
	$\mu$ for Bengali	1.2	1.2
	$\mu$ for Telugu	1.2	1.2

Table 5: The details of hyper-parameters in the training and evaluation process.

Language	Training Dataset (Tokens)		Evaluation Dataset (Instances)	
	General Corpus	Ability-related Corpus	Mathematical Tasks	Scientific Tasks
English	1.81B	1.30B (Math) / 1.82B (Sci)	250	1,245
Spanish	1.81B	-	250	1,232
Chinese	1.80B	-	250	1,229
Bengali	1.81B	-	250	1,137
Telugu	1.81B	-	250	1,036

Table 6: The statistical information of the training and evaluation datasets.

Methods	Qwen2.5 0.5B			Gemma2 2B		
	ES	TE	Avg.	ES	TE	Avg.
Backbone LLM	36.64	25.69	31.17	43.41	30.01	36.71
+ F-CPT <sub>L&amp;A</sub>	32.90	22.43	27.67	38.48	<b>30.39</b>	34.62
+ F-CPT <sub>A</sub>	32.62	25.26	28.94	37.83	25.39	31.61
+ MAET w/o API	36.72	28.91	32.82	43.23	29.59	36.41
+ MAET (Ours)	<b>36.91</b>	<b>29.62</b>	<b>33.27</b>	<b>43.62</b>	30.37	<b>37.00</b>

Table 7: The performance comparison of different LLMs on multilingual scientific tasks.

LLM Tensors	Proportion of $\mathcal{T}$	ES	ZH	BN	TE	Avg.
All Tensors	100.00%	49.60	41.60	32.40	25.20	37.20
Attention All	45.26%	48.80	41.60	28.80	26.40	36.40
Attention Q	12.07%	47.60	40.80	30.80	26.40	36.40
Attention K	10.34%	47.20	42.40	29.60	24.40	35.90
Attention V	9.48%	47.60	42.40	28.80	25.20	36.00
Attention O	13.36%	48.00	40.40	30.80	27.20	36.60
MLP All	41.38%	48.80	39.60	31.60	27.60	36.90
MLP Up	13.79%	50.00	40.00	28.80	25.20	36.00
MLP Gate	13.79%	46.00	41.20	30.00	24.00	35.30
MLP Down	13.79%	49.60	41.60	30.40	26.00	36.90

Table 8: The effect of only merging the specific LLM tensors during the transferring process (*i.e.*, Eq.5) on multilingual mathematical tasks.