

Extend Model Merging from Fine-Tuned to Pre-Trained Large Language Models via Weight Disentanglement

Anonymous ACL submission

Abstract

Merging Large Language Models (LLMs) aims to amalgamate multiple homologous LLMs into one with all the capabilities. Ideally, LLMs sharing the same backbone should be mergeable, irrespective of whether they are Fine-Tuned (FT) with minor parameter changes or Pre-Trained (PT) with substantial parameter shifts. However, existing methods often manually assign the model importance, rendering them feasible only for LLMs with similar parameter alterations (e.g., multiple FT LLMs). The diverse parameter changed ranges between FT and PT LLMs pose challenges for current solutions in empirically determining the optimal combination. In this paper, we make a pioneering effort to broaden the applicability of merging techniques from FT to PT LLMs. We initially examine the efficacy of current methods in merging FT and PT LLMs, showing that they struggle to handle PT LLMs. Subsequently, we extend the merging scope based on **WeIght DisENtanglement** (WIDEN), which first disentangles model weights into magnitude and direction components, and then performs adaptive fusion by considering their respective contributions. In the experiments, we merge Qwen1.5-Chat (an FT LLM with instruction-following skills) with Sailor (a PT LLM with multilingual abilities) across 1.8B, 4B, 7B, and 14B model sizes. Results reveal that: (1) existing solutions usually fail when merging Sailor, either losing both abilities or only retaining instruction-following skills; (2) WIDEN successfully injects the multilingual abilities of Sailor into Qwen1.5-Chat and make it proficient in Southeast Asian languages, achieving enhancements in the fundamental capabilities.

1 Introduction

In recent years, model merging has sparked significant interest as a prominent topic, which intends to integrate multiple homologous models (sharing the same backbone) into a singular one that encapsulates all the abilities (Wortsman et al., 2022; Matena

and Raffel, 2022; Ilharco et al., 2023). Model merging is lauded for its computational frugality, especially when applied to Large Language Models (LLMs). Notably, it achieves integration without using additional training data or even GPUs, establishing a new paradigm for efficiently combining capabilities of LLMs (Yu et al., 2024).

Technically, there are predominantly two strategies to equip LLMs with desired capabilities (Zhao et al., 2023): fine-tuning to elicit existing skills (Wang et al., 2023; Zhang et al., 2023a) and pre-training to inject new abilities (Wu et al., 2024). Existing merging methods mainly focus on integrating the skills of Fine-Tuned (FT) LLMs with minor parameter changes relative to the backbone, typically within 0.002 (Yu et al., 2024). However, it is crucial to acknowledge that pre-training is the cornerstone for fundamentally enhancing the capabilities of LLMs. The practicality of merging techniques in scenarios where Pre-Trained (PT) LLMs undergo substantial parameter shifts remains unexplored, as depicted in Figure 1. Consequently, if the application of merging is restricted to FT LLMs, its potential for broader improvement would be significantly constrained.

In this work, we make two technical contributions. **We examine the feasibility of existing approaches in absorbing the abilities from PT LLMs.** We investigate the performance of widely used arithmetic-based (Wortsman et al., 2022; Ilharco et al., 2023), geometric-based (Shoemaker, 1985; Jang et al., 2024), and pruning-based (Yadav et al., 2023; Davari and Belilovsky, 2023; Yu et al., 2024) methods when merging FT and PT LLMs. As illustrated in Table 1, we find current methods either lose efficacy in retaining the abilities of PT LLMs (leading to a decrease of approximately 10 to 20 points on average) or fail to preserve both capabilities (resulting in an average degradation of about 15 and 30 points, respectively). One possible reason is that existing methods depend on man-

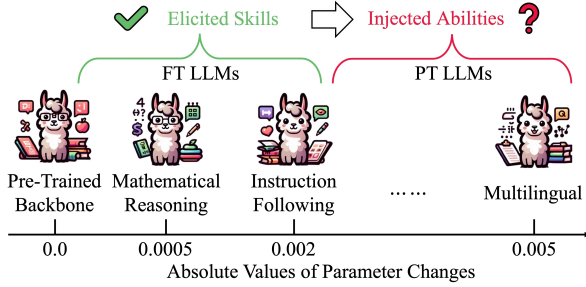


Figure 1: Issues of existing merging techniques.

	Instruction Following	Multilingual
Qwen1.5-14B-Chat	68.08	53.74
Sailor-14B	64.02	59.90
Arithmetic-based	66.30 (-1.78)	40.72 (-19.18)
Geometric-based	67.59 (-0.49)	49.52 (-10.38)
Pruning-based	51.72 (-16.36)	28.69 (-31.21)
WIDEN	66.75 (-1.33)	59.67 (-0.23)

Table 1: Average results of merging Qwen1.5-14B-Chat and Sailor-14B. Metrics of the best methods in Arithmetic, Geometric, and Pruning categories are reported.

ually assigned scaling terms to gauge the model contribution, which is only applicable when multiple LLMs depict comparable parameter alterations. Nonetheless, when confronted with diverse parameter changed ranges between FT and PT LLMs, deriving the optimal scaling factors according to human expertise becomes exceedingly arduous.

We propose a new solution grounded in WeIght DisEntanglement (WIDEN) to expand the scope of merging techniques from FT to PT LLMs. WIDEN tackles the drawbacks of existing works by automatically computing the model importance in the merging process without requiring manual specification, mitigating the influence induced by diverse parameter changed ranges between FT and PT LLMs. To be specific, WIDEN first disentangles each weight of a given LLM into two components: *magnitude* and *direction*. Then, the divergence of each component relative to the backbone is quantified to provide a numerical measure of how much each LLM has been altered. Next, WIDEN employs a ranking mechanism within each LLM to obtain the weight importance, tackling the diversity in parameter changed ranges between FT and PT LLMs. Finally, WIDEN performs adaptive merging on multiple LLMs by Softmax with the score calibration design.

We experiment with Qwen1.5-Chat (Bai et al., 2023) (an FT LLM with instruction-following skills) and Sailor (Dou et al., 2024) (a PT LLM with multilingual abilities for South-East Asia) across 1.8B, 4B, 7B, and 14B model scales¹. Experimental results indicate that WIDEN outperforms existing methods by not only absorbing the

¹As far as we know, Sailor is one of the few publicly accessible PT LLM that has undergone sufficient continued pre-training upon the open-source Qwen1.5 model (see Section A.8 and Section A.10 for more details), ideally suitable to our experimental scenarios. Therefore, Sailor and its homologous counterpart, Qwen1.5-Chat, are selected for our study.

multilingual abilities of Sailor but also preserving the instruction-following skills of Qwen1.5-Chat. For example, in Table 1, WIDEN slightly causes an average reduction of 0.23 and 1.33 points for Sailor-14B and Qwen1.5-14B-Chat, significantly better than existing methods. These observations demonstrate that WIDEN effectively extends the applicability of merging techniques from FT to PT LLMs. We further merge three FT LLMs including WizardLM-13B (Xu et al., 2024) for instruction following, WizardMath-13B (Luo et al., 2023) for mathematical reasoning, and llama-2-13b-code-alpaca (Chaudhary, 2023) for code generation. Results show that WIDEN is also feasible under the conventional setting and can strike a favorable balance among these capabilities. Resources are available at <https://anonymous.4open.science/r/MergeLLM-5E0D>.

2 Related Work

Fine-Tuning and Pre-Training of LLMs. LLMs can be adapted to various tasks via two strategies: fine-tuning and pre-training (Zhao et al., 2023). Fine-tuning is designed to elicit backbones with specific skills by optimizing models on a limited set of task-specific data, obtaining FT LLMs with skills such as instruction following (Rafailov et al., 2023; Song et al., 2024) and mathematical reasoning (Yuan et al., 2023; Luo et al., 2023). The fine-tuning process typically brings minor modifications to the model parameters (Yu et al., 2024). In contrast to fine-tuning, pre-training trains LLMs on large-scale raw corpora to enhance models with domain knowledge (Ke et al., 2022, 2023; Cheng et al., 2024), deriving PT LLMs with fundamental abilities like finance analysis (Xie et al., 2023) and law assistance (Colombo et al., 2024b). Pre-training often causes more obvious parameter shifts than fine-tuning due to the use of extensive data.

Merging of LLMs. Existing merging techniques for LLMs can be categorized into three groups, which are based on arithmetic, geometric, and pruning. Average Merging (Wortsman et al., 2022) and Task Arithmetic (Ilharco et al., 2023) belong to arithmetic-based approaches. As geometric-based methods, both SLERP (Shoemake, 1985) and Model Stock (Jang et al., 2024) consider the geometric properties in weight space. TIES-Merging (Yadav et al., 2023), Breadcrumbs (Davari and Belilovsky, 2023), and DARE (Yu et al., 2024) are methods based on pruning. Please see Section A.6 for more descriptions of the existing approaches. However, most current methods manually determine the importance of each model, suitable only for LLMs with similar parameter changes. When the parameter changed ranges are diverse between FT and PT LLMs, determining the optimal combination becomes overwhelmingly challenging. This paper initially verifies the limitations of existing methods in combining the abilities of PT LLMs. Then, an approach based on weight disentanglement is introduced to effectively expand the scope of merging techniques from FT to PT LLMs.

3 Methodology

3.1 Preliminaries

Merging Beyond FT LLMs. Given a collection of N homologous LLMs characterized by parameters $\{\Theta^1, \Theta^2, \dots, \Theta^N\}$, all of which share the same backbone with parameters Θ_{PRE} , model merging aims to amalgamate the parameters of N LLMs into a singular model with all the capabilities, denoted as Θ_{M} . Previous studies only focus on combining the skills of FT LLMs parameterized by $\{\Theta_{\text{FT}}^1, \Theta_{\text{FT}}^2, \dots, \Theta_{\text{FT}}^N\}$, where each model exhibits slight parameter changes, usually within 0.002 (Yu et al., 2024). In this paper, we extend the scope of merging techniques from FT to PT LLMs, intending to absorb the abilities of PT LLMs. Therefore, the parameters targeted for merging become $\{\Theta_{\text{TYPE}_1}^1, \Theta_{\text{TYPE}_2}^2, \dots, \Theta_{\text{TYPE}_N}^N\}$, where TYPE_n ($1 \leq n \leq N$) can be either FT or PT.

Weight Disentanglement. As outlined in Salimans and Kingma (2016); Liu et al. (2024), a weight $\mathbf{W} \in \mathbb{R}^{d \times k}$ can be disentangled into two components: a row vector $\mathbf{m} \in \mathbb{R}^{1 \times k}$ that captures the magnitudes and a matrix $\mathbf{D} \in \mathbb{R}^{d \times k}$ that stores the direction vectors. Here, d and k represent the output and input dimensions. Mathematically, the

disentanglement of weight \mathbf{W} is achieved by

$$\mathbf{W} = \mathbf{m}\mathbf{D} = \|\mathbf{W}\|_c \frac{\mathbf{W}}{\|\mathbf{W}\|_c} \in \mathbb{R}^{d \times k},$$

where $\|\cdot\|_c$ denotes the vector-wise l_c -norm of a matrix across each column. Such a decoupling operation guarantees that each column $\mathbf{D}_{:,j}$ ($1 \leq j \leq k$) is a unit vector, and scalar $m_j \in \mathbf{m}$ signifies the magnitude of direction vector $\mathbf{D}_{:,j}$. Since the primary challenge of extending merging scope to PT LLMs lies in the manual assignment of model importance, we employ weight disentanglement to initially decouple weights into magnitudes and directions, and then automatically compute the weight importance without human expertise based on these two components.

3.2 Exploring Efficacy of Current Methods When Merging PT LLMs

We investigate the efficacy of seven commonly used merging techniques when integrating the abilities of PT LLMs. To be specific, Average Merging (Wortsman et al., 2022) and Task Arithmetic (Ilharco et al., 2023) are arithmetic-based methods. SLERP (Shoemake, 1985) and Model Stock (Jang et al., 2024) belong to geometric-based approaches. TIES-Merging (Yadav et al., 2023), Breadcrumbs (Davari and Belilovsky, 2023) and DARE (Yu et al., 2024) are pruning-based solutions. Please see Section A.6 for detailed descriptions of these methods. To evaluate the performance, we attempt to combine the instruction-following skills of an FT LLM, Qwen1.5-Chat (Bai et al., 2023), and the multilingual abilities of a PT LLM, Sailor (Dou et al., 2024). Experimental setup, results, and analysis can be found in Section 4.

Since this part mainly concentrates on the feasibility of merging techniques when applied to PT LLMs, we highlight the key conclusion pertinent to PT LLMs: *existing merging approaches face difficulties in preserving the abilities of PT LLMs*. As evidenced in Table 2, the performance of all merging methods on the multilingual abilities significantly declines. This phenomenon is largely attributed to the reliance of most methods on manually assigned scaling factors to determine the contribution of each model at various levels throughout the merging process, encompassing model level (Ilharco et al., 2023; Yadav et al., 2023; Davari and Belilovsky, 2023), layer/module level (Goddard et al., 2024), and parameter level (Shoemake, 1985). The diverse parameter changed ranges between FT

and PT LLMs complicate the manual assignment of model importance, making it intractable to define optimal scaling factors case by case.

3.3 Extending Merging Scope to PT LLMs via Weight Disentanglement

We present a new approach based on **WeIght DisENtanglement (WIDEN)** to innovatively broaden the applicability of model merging techniques from FT to PT LLMs, whose key concept is to adaptively assess the importance of weights during the merging process for neutralizing the effects of diverse parameter changed ranges between FT and PT LLMs. As shown in Figure 4 in Section A.1, WIDEN mainly comprises four steps. Given the weights of LLMs (including the backbone as well as models to be merged), WIDEN 1) disentangles each weight into a row vector of magnitudes and a matrix of direction vectors; 2) estimates weight divergence relative to the backbone founded on absolute values of magnitude alterations and cosine similarities between direction vectors; 3) ranks the weights inside each LLM grounded in their divergence to derive the weight importance, thereby mitigating the impact of diverse parameter changed ranges; 4) merges multiple LLMs into a single one according to the obtained weight importance via Softmax with score calibration.

Disentangling Weights of LLMs. Given multiple homologous LLMs (each LLM can be obtained by either FT or PT) with parameters $\{\Theta^1, \Theta^2, \dots, \Theta^N\}$ as well as the backbone with parameters Θ_{PRE} , we first perform weight disentanglement for the parameters. Take $\mathbf{W}^n \in \Theta^n$ with shape $\mathbb{R}^{d \times k}$ as an example². \mathbf{W}^n can be decoupled into $\mathbf{m}^n = \|\mathbf{W}^n\|_c \in \mathbb{R}^{1 \times k}$ and $\mathbf{D}^n = \frac{\mathbf{W}^n}{\|\mathbf{W}^n\|_c} \in \mathbb{R}^{d \times k}$. After applying this disentanglement across all the LLMs, we can obtain the sets of row vectors of magnitudes $\{\mathbf{m}^n\}_{n=1}^N \cup \{\mathbf{m}_{\text{PRE}}\}$ and matrices of direction vectors $\{\mathbf{D}^n\}_{n=1}^N \cup \{\mathbf{D}_{\text{PRE}}\}$.

Estimating Weight Divergence Relative to Backbone. We estimate the weight divergence of each LLM relative to the backbone from the perspective of magnitudes and directions with two measurements. To be specific, we compute the absolute values of magnitude alterations and determine the changes between direction vectors based

²Note that Θ^n represents the collection of parameters of the n -th LLM, consisting of a multitude of weights.

on cosine similarities as follows,

$$\begin{aligned} \Delta \mathbf{m}^n &= |\mathbf{m}^n - \mathbf{m}_{\text{PRE}}| \in \mathbb{R}^{1 \times k}, \\ \Delta \mathbf{D}_j^n &= 1 - \text{CosineSimilarity}(\mathbf{D}_{:,j}^n, \mathbf{D}_{\text{PRE},:,j}) \in \mathbb{R}, \\ &\text{for } 1 \leq j \leq k, 1 \leq n \leq N, \end{aligned}$$

where $\text{CosineSimilarity}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|_2 \cdot \|\mathbf{y}\|_2}$. Thus, we obtain the divergences of the LLMs relative to the backbone in both magnitudes $\{\Delta \mathbf{m}^n \in \mathbb{R}^{1 \times k}\}_{n=1}^N$ and directions $\{\Delta \mathbf{D}^n \in \mathbb{R}^{1 \times k}\}_{n=1}^N$.

Ranking Weights Inside Each LLM. We design a ranking mechanism to alleviate the potential impact of diverse parameter changed ranges among various LLMs, which assigns importance to the weights within each LLM according to their divergence relative to the backbone (greater divergence indicates higher essentiality). The ranking mechanism is applied to both the magnitudes and the directions of weights. To illustrate, consider the magnitudes as an instance. Given $\Delta \mathbf{m}^n \in \mathbb{R}^{1 \times k}$ of the n -th LLM, we initially sort $\Delta \mathbf{m}^n$ in ascending order, yielding an index row vector $\mathbf{m}_{\text{IND}}^n \in \mathbb{R}^{1 \times k}$ that contains values ranging from 1 to k . Subsequently, we derive a row vector $\tilde{\mathbf{m}}^n \in \mathbb{R}^{1 \times k}$ that encapsulates normalized ranking scores based on $\mathbf{m}_{\text{IND}}^n$, which is computed by

$$\tilde{m}_{\text{IND}_j}^n = j/k, \text{ for } 1 \leq j \leq k.$$

$\tilde{\mathbf{m}}^n \in \mathbb{R}^{1 \times k}$ represents the normalized importance of each position within the range $[1, \dots, k]$ for the n -th LLM. Following the same procedure, the directions of weights can also be assigned with normalized importance, which can be denoted by $\tilde{\mathbf{D}}^n \in \mathbb{R}^{1 \times k}$. Such a ranking mechanism ensures that, within each LLM, the importance of magnitudes and directions is uniformly distributed between 0 and 1, thereby eliminating the potential influences arising from diverse parameter changed ranges between FT and PT LLMs. After applying the ranking operation for all the LLMs, we can ultimately obtain $\{\tilde{\mathbf{m}}^n \in \mathbb{R}^{1 \times k}\}_{n=1}^N$ and $\{\tilde{\mathbf{D}}^n \in \mathbb{R}^{1 \times k}\}_{n=1}^N$.

Merging LLMs via Softmax with Score Calibration. We employ an adaptive merging strategy for multiple LLMs through a Softmax function, complemented by score calibration. Initially, we calculate the importance scores for magnitudes and directions by applying the Softmax function to $\{\tilde{\mathbf{m}}^n \in \mathbb{R}^{1 \times k}\}_{n=1}^N$ and $\{\tilde{\mathbf{D}}^n \in \mathbb{R}^{1 \times k}\}_{n=1}^N$, yielding $\tilde{\mathcal{M}}, \tilde{\mathcal{D}} \in \mathbb{R}^{N \times k}$. Since the computation process for magnitudes and directions are similar, we

only illustrate the calculations relevant to magnitudes due to space limits. Please refer to Section A.2 for the calculations related to directions.

$$\widetilde{\mathcal{M}}_{n,j} = \frac{\exp(\widetilde{m}_j^n)}{\sum_{n'=1}^N \exp(\widetilde{m}_{j'}^{n'})} \in \mathbb{R}.$$

However, Softmax restricts the sum of parameter importance across multiple LLMs to 1, potentially diminishing the significance of crucial parameters in certain cases. Thus, we incorporate the score calibration to relax the constraint of Softmax for essential parameters. We identify crucial parameters as those whose importance exceeds the average level by a factor of t as follows,

$$\mathbb{P}_m^n = \{j | \widetilde{m}_j^n > \frac{t}{k} \cdot \sum_{j'=1}^k \widetilde{m}_{j'}^n\}. \quad (1)$$

Subsequently, we calibrate the scores using \mathbb{P}_m^n by

$$\mathcal{M}_{n,j} = \begin{cases} s, & \text{if } j \in \mathbb{P}_m^n \\ \widetilde{\mathcal{M}}_{n,j}, & \text{if } j \notin \mathbb{P}_m^n \end{cases},$$

where s regulates the numerical value of score calibration. Finally, we integrate the weights of multiple LLMs into $\mathbf{W}_M \in \mathbb{R}^{d \times k}$ by considering the adjusted contributions of magnitudes and directions,

$$\mathbf{W}_M = \mathbf{W}_{\text{PRE}} + \sum_{n=1}^N \frac{\mathcal{M}_{n,:} + \mathcal{D}_{n,:}}{2} \odot (\mathbf{W}^n - \mathbf{W}_{\text{PRE}}). \quad (2)$$

Note that t and s are designed to control the merging importance of parameters after applying the Softmax function. If more parameters are desired to be assigned with higher importance, t should be reduced and s should be increased. Conversely, t should be increased and s should be reduced.

Remark 1. The above procedure is designed to deal with two-dimensional weights within LLMs, accounting for both magnitudes and directions. For one-dimensional parameters, such as weights in normalization layers and biases in linear transformations, we handle them as vectors of magnitudes and estimate their changes relative to the backbone by absolute values of the differences.

Remark 2. Existing arithmetic-based merging methods including Average Merging (Wortsman et al., 2022) and Task Arithmetic (Ilharco et al., 2023), can be viewed as special instances of the proposed WIDEN. Specifically, when $t < 0.0$ and $s = 1/N$, WIDEN transforms into Average Merging; when $t < 0.0$ and $s = \lambda$, WIDEN represents Task Arithmetic. Detailed explanations are shown in Section A.3.

4 Experiments

We conduct experiments in two scenarios: 1) integrating FT and PT LLMs, a new setting not explored before; 2) combining FT LLMs as in previous research. As our main focus lies in merging FT and PT LLMs, we provide the details and results of the second setting in Section A.9.

4.1 Experimental Setup

Merging FT and PT LLMs. We choose Qwen1.5-Chat (Bai et al., 2023) with instruction-following skills as the FT LLM and select Sailor (Dou et al., 2024) with multilingual abilities for South-East Asia as the PT LLM. Both models adopt Qwen1.5 (Bai et al., 2023) as the backbone. Open LLM Leaderboard (Beeching et al., 2023) and benchmark for South-East Asian languages (Dou et al., 2024) are used for evaluating the performance of models across 1.8B, 4B, 7B, and 14B sizes. Please see Section A.5 for the overview and evaluation metrics of the benchmarks. Also, refer to Table 6 in Section A.4 for the details of FT and PT LLMs. We compare WIDEN with seven popular baselines for model merging, including Average Merging (Wortsman et al., 2022), Task Arithmetic (Ilharco et al., 2023), SLERP (Shoemaker, 1985), Model Stock (Jang et al., 2024), TIES-Merging (Yadav et al., 2023), Breadcrumbs (Davari and Belilovsky, 2023), and DARE (Yu et al., 2024). See Section 3.2 and Section A.6 for more descriptions.

Configurations of Merging Methods. We apply grid search to identify the optimal settings for various merging techniques. The proposed WIDEN utilizes l_2 normalization and involves two hyperparameters: s and t . For ease of implementation, the score calibration factor s is consistently fixed to 1.0 across all the cases. The factor t is determined by grid search. Please refer to Table 7 in Section A.7 for detailed information about the searched ranges.

4.2 Overall Performance of Merging LLMs

Table 2 and Table 11 show the results of merging Qwen1.5-Chat and Sailor on South-East Asian language benchmark. Since Average Merging is a special case of Task Arithmetic when the scaling term is 0.5, we thereby only report the results of Task Arithmetic, which inherently include the performance of Average Merging. Note that th, id, vi, and jv are abbreviations of Thai, Indonesian, Vietnamese, and Javanese. The best and second-best results are marked in **bold** and underlined fonts.

Size	Models	Merging Methods	XQuAD	TydiQA	XQuAD	XCOPA			Belebele			M3Exam	Average	Average
			th	id	vi	th	id	vi	th	id	vi	qv	Rank	Rank
7B	Qwen1.5	/	53.79/69.30	57.17/77.28	56.63/76.99	54.20	62.20	66.20	38.33	42.00	42.89	26.15	55.63	/
	Qwen1.5-Chat	/	24.28/46.77	42.30/67.57	45.51/69.91	56.20	66.80	70.40	38.67	43.11	47.11	28.30	49.76	/
	Sailor	/	57.88/71.06	60.53/75.42	53.81/74.62	59.00	72.20	72.20	41.56	44.33	45.33	32.88	58.52	/
	Qwen1.5-Chat & Sailor	Task Arithmetic	28.20/49.62	45.84/65.78	37.38/61.53	63.20	77.60	73.40	38.89	46.89	45.11	30.46	51.07	2.15
		SLERP	16.62/43.62	20.53/54.02	33.70/61.49	55.80	73.40	73.00	38.44	47.89	47.56	28.30	45.72	3.23
		Model Stock	26.72/52.69	24.78/58.88	43.80/69.50	54.60	66.00	69.40	37.33	42.78	43.67	27.76	47.53	3.31
		TIES-Merging	0.61/8.84	5.66/17.23	7.70/20.78	50.20	62.20	59.80	30.22	35.33	35.11	25.07	27.60	5.54
		Breadcrumbs	6.79/11.38	7.61/15.23	12.32/27.90	51.40	66.40	57.20	31.33	34.00	32.56	24.53	29.13	5.23
	WIDEN	42.65/64.21	45.84/73.37	48.42/73.17	60.20	77.40	73.60	40.11	51.11	48.56	32.88	56.27	1.15	
Qwen1.5	/	55.53/74.39	60.35/81.07	57.66/77.62	58.40	70.40	72.60	41.22	48.67	44.44	26.15	59.12	/	
Qwen1.5-Chat	/	33.59/59.98	37.17/65.46	44.14/71.91	61.80	75.20	71.80	44.00	51.00	52.67	29.92	53.74	/	
Sailor	/	49.43/70.01	58.94/77.85	57.74/77.34	62.60	77.60	78.60	40.89	47.67	47.11	32.88	59.90	/	
14B	Qwen1.5	/	8.53/24.39	13.45/33.54	13.52/25.75	59.80	82.40	78.20	46.00	56.33	53.78	33.69	40.72	2.54
	Qwen1.5-Chat	/	14.53/44.70	<u>22.48/61.67</u>	42.69/69.48	61.80	75.60	74.60	<u>43.22</u>	52.56	<u>50.56</u>	29.92	<u>49.52</u>	<u>2.46</u>
	Sailor	/	25.59/53.10	14.87/51.19	44.74/70.20	58.60	70.40	71.80	42.67	49.89	45.11	27.22	48.11	3.08
	Qwen1.5-Chat & Sailor	Task Arithmetic	0.44/8.78	1.42/12.87	0.00/6.95	55.20	69.20	67.20	32.78	39.00	37.11	27.22	27.55	5.46
		SLERP	1.22/6.48	2.30/20.88	3.17/14.46	52.20	64.60	63.40	34.78	42.11	40.67	26.68	28.69	5.23
		WIDEN	49.61/73.16	50.62/75.09	54.75/78.23	60.80	77.40	74.60	42.22	56.22	50.44	32.61	59.67	1.77

Table 2: Results of merging Qwen1.5-Chat and Sailor on South-East Asian language benchmark.

From Table 2, two conclusions can be summarized.

Firstly, *existing model merging approaches encounter significant challenges when incorporating the multilingual abilities of Sailor, leading to a marked decline in performance.* The downturn is probably attributed to the difficulty in determining the optimal combination due to diverse parameter changed ranges between Qwen1.5-Chat and Sailor. We also notice that the reduction is particularly pronounced in pruning-based methods, prompting us to conduct additional verifications. As demonstrated in Table 3, we find that the feasibility of pruning strategies such as DARE and Magnitude-based Pruning (MP) in TIES-Merging and Breadcrumbs is severely compromised with minor parameter drop rates on Sailor-7B, far below the levels reported results in the original studies (i.e., 0.9 in DARE, 0.8 in TIES-Merging, and 0.85 in Breadcrumbs), diminishing the effectiveness of pruning in alleviating parameter interference. As a result, DARE fails to serve as a plug-in for existing merging techniques when considering PT LLMs, and its inferior results are excluded.

Secondly, *WIDEN effectively assimilates the multilingual capabilities of Sailor, emerging as the top performer among all the merging techniques.* The key advantage of WIDEN lies in the adaptive computation of weight importance by considering both magnitudes and directions during the merging process, mitigating the effects of diverse parameter changed ranges between FT and PT LLMs.

Table 4 and Table 12 depict the merging performance on Open LLM Leaderboard. We find that geometric-based approaches (SLERP and Model

	Drop Rate	XQuAD	XCOPA	Belebele
Sailor-7B	/	53.81/74.62	72.20	45.33
DARE	0.1	47.56/66.95	64.20	41.00
	0.3	5.90/16.05	55.60	30.56
MP	0.1	54.23/75.16	72.80	45.44
	0.3	52.44/73.53	72.20	44.78
	0.5	49.19/70.11	70.00	43.67
	0.8	13.77/30.13	59.00	34.56

Table 3: Performance of pruning strategies on Sailor-7B for Vietnamese-related tasks.

Stock) excel in retraining the instruction-following skills of Qwen1.5-Chat, indicating that parameters of FT LLMs may potentially exhibit more evident properties in the geometric space. WIDEN shows competitive results alongside SLERP and Model Stock, underscoring its applicability in merging FT LLMs. Moreover, WIDEN outperforms arithmetic-based methods since it is a generalized format of these methods and offers greater flexibility through the adaptive computation of weight importance. The performance of WIDEN consistently improves with increasing model sizes, indicating its potential scalability. Although WIDEN achieves competitive but not state-of-the-art performance on the Open LLM Leaderboard, it consistently delivers satisfactory results across both benchmarks, while most baselines fail to do so, demonstrating the robustness and generalizability of WIDEN.

From Table 9 in Section A.9, we observe that the efficacy of certain baselines drastically fluctuates when integrating FT LLMs. For example, Model Stock appears to lose potency, whereas pruning-based methods including TIES-Merging and Breadcrumbs show competitive performance. WIDEN

Size	Models	Merging Methods	ARC	Hella-Swag	MMLU	Truthful-QA	Wino-grande	GSM8K	Average	Average Rank	
7B	Qwen1.5	/	54.86	78.45	60.60	51.09	71.03	56.79	62.14	/	
	Qwen1.5-Chat	/	56.14	78.71	60.18	53.61	67.48	54.21	61.72	/	
	Sailor	/	49.57	76.13	52.91	40.07	71.35	34.65	54.11	/	
	Qwen1.5-Chat & Sailor	Task Arithmetic		52.05	75.15	59.38	50.84	69.77	25.55	55.46	3.50
		SLERP		54.78	76.20	60.76	50.78	71.51	55.50	61.59	2.33
		Model Stock		55.12	76.29	61.18	49.33	71.43	55.80	<u>61.53</u>	2.00
		TIES-Merging		43.86	56.88	52.39	46.59	67.56	0.00	44.55	5.67
Breadcrumbs		47.18	49.99	52.66	52.05	64.88	0.45	44.53	4.67		
WIDEN		53.84	<u>76.25</u>	57.65	49.34	71.90	44.81	58.97	2.83		
14B	Qwen1.5	/	56.40	81.22	67.79	52.04	74.43	68.01	66.65	/	
	Qwen1.5-Chat	/	57.25	82.56	67.48	60.42	72.69	68.08	68.08	/	
	Sailor	/	55.46	80.31	62.95	46.64	76.80	61.94	64.02	/	
	Qwen1.5-Chat & Sailor	Task Arithmetic		56.57	81.59	67.52	62.93	75.22	53.98	66.30	<u>2.50</u>
		SLERP		55.72	79.94	<u>67.94</u>	57.51	75.14	69.29	67.59	3.00
		Model Stock		<u>57.00</u>	<u>80.50</u>	68.44	51.98	<u>76.01</u>	<u>66.72</u>	<u>66.77</u>	2.33
		TIES-Merging		49.74	67.23	60.54	47.43	72.14	0.30	49.56	5.67
Breadcrumbs		51.88	62.22	63.47	<u>57.90</u>	70.32	4.55	51.72	4.83		
WIDEN		57.17	80.05	66.00	54.85	76.09	66.34	66.75	2.67		

Table 4: Performance of merging Qwen1.5-Chat and Sailor on Open LLM Leaderboard.

consistently depicts results that are on par with established merging techniques in most situations, affirming its suitability in the standard setting of merging multiple FT LLMs. It is worth noting that WIDEN performs competitively but less prominently than baselines when merging multiple FT models. This is because WIDEN excels at merging LLMs with obvious differences in parameter changed ranges by disentangling parameters into magnitudes and directions. In the case of FT models with minor and similar parameter changes, treating weights holistically or disentangling them leads to minimal disparity, which makes the disentanglement operation less pronounced.

4.3 Investigations of Designs in WIDEN

The foundational designs in WIDEN consist of three components: weight disentanglement, ranking weights inside each model, and score calibration for Softmax. To assess the contribution of each module, we respectively remove the above components and measure the performance of the remaining parts. Specifically, we eliminate the disentanglement of weights by calculating the discrepancy between the weights of LLM and the corresponding backbone using cosine similarities, denoted as WIDEN w/o WD. We substitute the ranking mechanism with min-max normalization within each model, represented by WIDEN w/o RANK. We discard the score calibration and directly employ Softmax to compute importance scores, identified as WIDEN w/o SC. Figure 2 shows the impact of these three modifications, where OLL and SEA

are the abbreviations for Open LLM Leaderboard and South-East Asian language benchmark, respectively. Please notice that the reported results are the average of metrics across all the datasets within each benchmark.

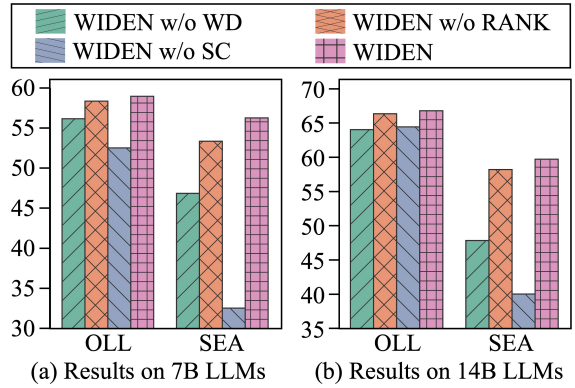


Figure 2: Effects of various designs in WIDEN.

From Figure 2, we find that each design in WIDEN contributes to enhancing the merging performance, particularly in absorbing the multilingual abilities on the South-East Asian language benchmark. Precisely, the weight disentanglement refines the estimation of weight importance at a granular level, considering both magnitude and direction. The ranking mechanism offers a smoother distribution of weight importance based on continuous indices, effectively mitigating the influence of diverse parameter changed ranges. The calibration of scores computed by Softmax reallocates importance to critical parameters, which maintains

Adjustments	Models	L→L	L→M	L→H	M→L	M→M	M→H	H→L	H→M	H→H
WIDEN w/o WD to WIDEN	Qwen1.5-7B-Chat	18.82%	11.09%	3.42%	13.97%	10.18%	9.18%	0.54%	12.06%	20.75%
	Sailor-7B	15.34%	10.50%	7.48%	17.80%	7.72%	7.80%	0.18%	15.10%	18.07%
WIDEN w/o SC to WIDEN	Qwen1.5-7B-Chat	24.78%	7.69%	0.85%	7.93%	17.51%	7.88%	0.62%	8.12%	24.61%
	Sailor-7B	22.01%	9.52%	1.80%	9.63%	15.14%	8.56%	1.69%	8.67%	22.99%

Table 5: Adjustments of weight importance made by WIDEN.

the characteristics of essential parameters across multiple models. In summary, the components of WIDEN are indispensable and improve performance with varied benefits; the removal of any module leads to diminished outcomes.

4.4 Analysis of Computed Weight Importance

We further delve into the properties of weight importance calculated by WIDEN from both qualitative and quantitative perspectives. Since Figure 2 demonstrates that the improvements in weight disentanglement and score calibration are notably more pronounced, we qualitatively depict the distribution of weight importance computed by WIDEN, WIDEN w/o WD, and WIDEN w/o SC on 7B model size in Figure 3. Our observations reveal that: 1) WIDEN exhibits a more balanced and reasonable weight importance distribution than WIDEN w/o WD, attributed to the disentanglement of weights. The distribution of WIDEN ranges approximately from 0.3 to 0.8 and 0.9 to 1.0, versus 0.3 to 0.6 and 0.9 to 1.0 for WIDEN w/o WD. WIDEN considers the collective contributions of magnitude and direction, rather than the individual impacts of weights, leading to a more holistic assessment of weight importance with increased numbers of weights falling within the importance range from 0.6 to 0.8. As a result, compared with WIDEN w/o WD, WIDEN achieves 4.98% and 20.08% improvements on average on the Open LLM Leaderboard and the South-East Asian language benchmark, respectively; 2) In contrast to WIDEN w/o SC, WIDEN distinguishes essential weights and assigns high importance within the range of 0.6 to 0.8 as well as 0.9 to 1.0 for certain weights, thanks to the design of score calibration. Therefore, WIDEN ensures the retention of essential weights in both Qwen1.5-7B-Chat and Sailor-7B, resulting in 12.25% and 72.87% average enhancements on the two benchmarks.

Furthermore, we categorize weight importance into three levels: Low (L), Medium (M), and High (H). The Low tier comprises the first third of weights when sorted by ascending importance,

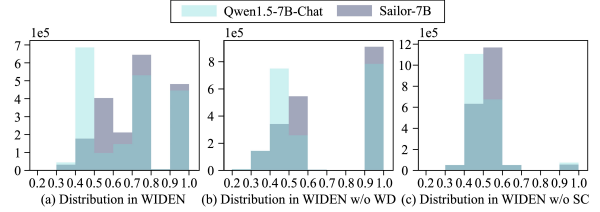


Figure 3: Distribution of weight importance computed by WIDEN and its variations.

indicating those with the least significance. The Medium tier includes weights from the 1/3 mark to the 2/3 mark, and the High tier contains weights from the 2/3 mark to the end. Table 5 quantitatively illustrates the adjustments of weight importance made by WIDEN when compared to WIDEN w/o WD and WIDEN w/o SC across three levels. We find that WIDEN effectively reallocates the weight importance via three aspects: 1) elevating weights of lower importance from Low to Medium; 2) either demoting or promoting weights of medium importance from Medium to Low or from Medium to High, respectively; 3) decreasing weights of high importance from High to Medium. These adjustments in weight importance explain how WIDEN brings improvements through the designs of weight disentanglement and score calibration.

5 Conclusion

In this study, we paved the way for extending the merging scope from FT to PT LLMs. Specifically, we first observed that existing methods struggled to integrate the abilities of PT LLMs and then introduced WIDEN, an innovative approach based on weight disentanglement, to effectively deploy merging strategies to PT LLMs. Experimental findings demonstrated that WIDEN not only exhibited an advantage in absorbing the abilities of PT LLMs but also preserved the skills of FT LLMs. We further offered a detailed analysis of the designs underlying WIDEN. This work makes the first attempt to broaden the sources of combinable abilities, which is expected to foster the broader application of model merging techniques.

625 Limitations

626 This paper investigates the merging task of LLMs,
627 no matter they are derived from fine-tuning or pre-
628 training. Although the proposed WIDEN effective-
629 ly extends the merging scope from FT to PT
630 LLMs, it remains applicable only to LLMs with
631 identical architectures (i.e., homologous LLMs).
632 For models derived from different backbones, our
633 method is no longer suitable. Furthermore, due
634 to the limited availability of LLMs that meet the
635 required criteria, the number of models attempted
636 to be merged in the experiments is restricted, with
637 a maximum of three. In the future, if more suitable
638 models become accessible, exploring the impact of
639 the number of LLMs to be merged on performance
640 would be a promising research direction.

641 Ethics Statement

642 Though this work has no direct ethical problems,
643 LLMs may still potentially generate harmful in-
644 formation including gender bias, fake news, and
645 private messages when equipped with our approach.
646 It is necessary and promising to design specialized
647 mechanisms to regulate these underlying issues.

648 References

649 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama.
650 2020. On the cross-lingual transferability of mono-
651 lingual representations. In *Proceedings of the 58th*
652 *Annual Meeting of the Association for Computational*
653 *Linguistics*, pages 4623–4637. Association for Com-
654 putational Linguistics.

655 Jacob Austin, Augustus Odena, Maxwell I. Nye,
656 Maarten Bosma, Henryk Michalewski, David Dohan,
657 Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le,
658 and Charles Sutton. 2021. Program synthesis with
659 large language models. *CoRR*, abs/2108.07732.

660 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
661 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
662 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,
663 Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,
664 Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,
665 Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong
666 Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang
667 Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian
668 Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi
669 Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang,
670 Yichang Zhang, Zhenru Zhang, Chang Zhou, Jin-
671 gren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023.
672 Qwen technical report. *CoRR*, abs/2309.16609.

673 Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel
674 Artetxe, Satya Narayan Shukla, Donald Husa, Naman
675 Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and

Madian Khabsa. 2023. The belebele benchmark: a
parallel reading comprehension dataset in 122 lan-
guage variants. *CoRR*, abs/2308.16884.

Edward Beeching, Clémentine Fourrier, Nathan Habib,
Sheon Han, Nathan Lambert, Nazneen Rajani, Omar
Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023.
Open llm leaderboard.

Sahil Chaudhary. 2023. Code alpaca: An instruction-
following llama model for code generation. <https://github.com/sahil280114/codealpaca>.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
Henrique Pondé de Oliveira Pinto, Jared Kaplan,
Harrison Edwards, Yuri Burda, Nicholas Joseph,
Greg Brockman, Alex Ray, Raul Puri, Gretchen
Krueger, Michael Petrov, Heidy Khlaaf, Girish Sas-
try, Pamela Mishkin, Brooke Chan, Scott Gray,
Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz
Kaiser, Mohammad Bavarian, Clemens Winter,
Philippe Tillet, Felipe Petroski Such, Dave Cum-
mings, Matthias Plappert, Fotios Chantzis, Eliza-
beth Barnes, Ariel Herbert-Voss, William Hebgem
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. 2021. Evaluat-
ing large language models trained on code. *CoRR*,
abs/2107.03374.

Daixuan Cheng, Shaohan Huang, and Furu Wei. 2024.
Adapting large language models via reading compre-
hension. In *The Twelfth International Conference on*
Learning Representations. OpenReview.net.

Jonathan H. Clark, Jennimaria Palomaki, Vitaly Niko-
laev, Eunsol Choi, Dan Garrette, Michael Collins,
and Tom Kwiatkowski. 2020. Tydi QA: A bench-
mark for information-seeking question answering in
typologically diverse languages. *Trans. Assoc. Com-
put. Linguistics*, 8:454–470.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,
Ashish Sabharwal, Carissa Schoenick, and Oyvind
Tafjord. 2018. Think you have solved question an-
swering? try arc, the AI2 reasoning challenge. *CoRR*,
abs/1803.05457.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, Christopher Hesse, and John Schulman.
2021. Training verifiers to solve math word prob-
lems. *CoRR*, abs/2110.14168.

Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf,
Dominic Culver, Rui Melo, Caio Corro, André F. T.
Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia
Morgado, and Michael Desa. 2024a. Saullm-7b: A
pioneering large language model for law. *CoRR*,
abs/2403.03883.

734	Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, Andre FT Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, and Michael Desa. 2024b. Saullm-7b: A pioneering large language model for law. <i>arXiv preprint arXiv:2403.03883</i> .	Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. 2023. Continual pre-training of language models. In <i>The Eleventh International Conference on Learning Representations</i> . OpenReview.net.	790 791 792 793 794
740	MohammadReza Davari and Eugene Belilovsky. 2023. Model breadcrumbs: Scaling multi-task model merging with sparse masks. <i>CoRR</i> , abs/2312.06795.	Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. Biomistral: A collection of open-source pretrained large language models for medical domains. In <i>Findings of the Association for Computational Linguistics</i> , pages 5848–5864. Association for Computational Linguistics.	795 796 797 798 799 800 801
743	Longxu Dou, Qian Liu, Guangtao Zeng, Jia Guo, Jiahui Zhou, Wei Lu, and Min Lin. 2024. Sailor: Open language models for south-east asia. <i>CoRR</i> , abs/2404.03608.	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics</i> , pages 3214–3252. Association for Computational Linguistics.	802 803 804 805 806 807
747	Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. 2024. Length-controlled alpaca-eval: A simple way to debias automatic evaluators. <i>CoRR</i> , abs/2404.04475.	Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024. Dora: Weight-decomposed low-rank adaptation. <i>CoRR</i> , abs/2402.09353.	808 809 810 811 812
751	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. A framework for few-shot language model evaluation.	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. <i>CoRR</i> , abs/2308.09583.	813 814 815 816 817 818
760	Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vlad Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. 2024. Arcee’s mergekit: A toolkit for merging large language models. <i>CoRR</i> , abs/2403.13257.	Michael Matena and Colin Raffel. 2022. Merging models with fisher-weighted averaging. In <i>NeurIPS</i> .	819 820
765	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding. In <i>9th International Conference on Learning Representations</i> . OpenReview.net.	Edoardo Maria Ponti, Goran Glavas, Olga Majewska, Qianchu Liu, Ivan Vulic, and Anna Korhonen. 2020. XCOPA: A multilingual dataset for causal common-sense reasoning. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> , pages 2362–2376. Association for Computational Linguistics.	821 822 823 824 825 826 827
770	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the MATH dataset. In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1</i> .	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In <i>Advances in Neural Information Processing Systems 36</i> .	828 829 830 831 832
776	Gabriel Ilharco, Marco Túlio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In <i>The Eleventh International Conference on Learning Representations</i> . OpenReview.net.	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale. In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence</i> , pages 8732–8740. AAAI Press.	833 834 835 836 837
781	Dong-Hwan Jang, Sangdoo Yun, and Dongyoon Han. 2024. Model stock: All we need is just a few fine-tuned models. <i>CoRR</i> , abs/2403.19522.	Tim Salimans and Diederik P. Kingma. 2016. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In <i>Advances in Neural Information Processing Systems 29</i> , page 901.	838 839 840 841 842
784	Zixuan Ke, Haowei Lin, Yijia Shao, Hu Xu, Lei Shu, and Bing Liu. 2022. Continual training of language models for few-shot learning. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 10205–10216. Association for Computational Linguistics.	Ken Shoemake. 1985. Animating rotation with quaternion curves. In <i>Proceedings of the 12th Annual Conference on Computer Graphics and Interactive Techniques</i> , pages 245–254. ACM.	843 844 845 846

847	Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. 2024. Preference ranking optimization for human alignment. In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence</i> , pages 18990–18998. AAAI Press.	904
848		905
849		906
850		907
851		908
852	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>CoRR</i> , abs/2307.09288.	909
853		910
854		911
855		912
856		913
857		
858		914
859		915
860		916
861		917
862		
863		918
864		919
865		920
866		921
867		922
868		923
869		
870		924
871		925
872		926
873		927
874		928
875	Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023. Aligning large language models with human: A survey. <i>CoRR</i> , abs/2307.12966.	929
876		930
877		931
878		932
879		933
880	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 <i>International Conference on Machine Learning</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 23965–23998. PMLR.	934
881		935
882		936
883		937
884		938
885		939
886		940
887		941
888		
889		
890	Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari. 2024. Continual learning for large language models: A survey. <i>CoRR</i> , abs/2402.01364.	
891		
892		
893		
894	Yong Xie, Karan Aggarwal, and Aitzaz Ahmad. 2023. Efficient continual pre-training for building domain specific large language models. <i>CoRR</i> , abs/2311.08545.	
895		
896		
897		
898	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. WizardLM: Empowering large pre-trained language models to follow complex instructions. In <i>The Twelfth International Conference on Learning Representations</i> .	
899		
900		
901		
902		
903		
	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</i> 36.	
	Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. 2024. Language models are super mario: Absorbing abilities from homologous models as a free lunch. In <i>International Conference on Machine Learning</i> . PMLR.	
	Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Chuanqi Tan, and Chang Zhou. 2023. Scaling relationship on learning mathematical reasoning with large language models. <i>CoRR</i> , abs/2308.01825.	
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics</i> , pages 4791–4800. Association for Computational Linguistics.	
	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2023a. Instruction tuning for large language models: A survey. <i>CoRR</i> , abs/2308.10792.	
	Wenxuan Zhang, Mahani Aljunied, Chang Gao, Yew Ken Chia, and Lidong Bing. 2023b. M3exam: A multilingual, multimodal, multilevel benchmark for examining large language models. In <i>Advances in Neural Information Processing Systems</i> 36.	
	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. <i>CoRR</i> , abs/2303.18223.	

A Appendix

A.1 Computation Process of WIDEN

Figure 4 illustrates the framework of WIDEN.

A.2 Computation Process Related to Directions

Calculations relevant to directions are executed by

$$\begin{aligned} \tilde{D}_{n,j} &= \frac{\exp(\tilde{D}_j^n)}{\sum_{n'=1}^N \exp(\tilde{D}_{j'}^{n'})} \in \mathbb{R}, \\ \mathbb{P}_D^n &= \{j \mid \tilde{D}_j^n > \frac{t}{k} \cdot \sum_{j'=1}^k \tilde{D}_{j'}^n\}, \\ D_{n,j} &= \begin{cases} s, & \text{if } j \in \mathbb{P}_D^n \\ \tilde{D}_{n,j}, & \text{if } j \notin \mathbb{P}_D^n \end{cases}. \end{aligned} \quad (3)$$

A.3 Connections Between WIDEN and Arithmetic-based Merging Methods

We discuss the relationships between WIDEN and existing arithmetic-based merging methods including Average Merging and Task Arithmetic. The computation procedure of Average Merging (Wortsman et al., 2022) for N LLMs is denoted by

$$W_M = \frac{1}{N} \sum_{n=1}^N W^n = W_{\text{PRE}} + \frac{1}{N} \sum_{n=1}^N (W^n - W_{\text{PRE}}).$$

Task Arithmetic (Ilharco et al., 2023) is implemented as follows,

$$W_M = W_{\text{PRE}} + \lambda \sum_{n=1}^N (W^n - W_{\text{PRE}}),$$

where λ denotes the scaling term. It is straightforward that in Equation (1) and Equation (3), if t is set to be minus, all the parameters can be considered crucial, with their importance scores calibrated to s . Thus, Equation (2) can be rewritten as

$$\begin{aligned} W_M &= W_{\text{PRE}} + \sum_{n=1}^N \frac{s+s}{2} (W^n - W_{\text{PRE}}) \\ &= W_{\text{PRE}} + s \sum_{n=1}^N (W^n - W_{\text{PRE}}). \end{aligned} \quad (4)$$

To this end, when $t < 0.0$ and $s = 1/N$, WIDEN aligns with Average Merging; when $t < 0.0$ and $s = \lambda$, WIDEN becomes Task Arithmetic.

A.4 Details of FT and PT LLMs

Table 6 depicts the versions and correspondences with backbones of FT and PT LLMs.

Types	Models	Backbones
FT LLM	Qwen1.5-1.8B-Chat ³	Qwen1.5-1.8B ⁴
PT LLM	Sailor-1.8B ⁵	Qwen1.5-1.8B ⁴
FT LLM	Qwen1.5-4B-Chat ⁶	Qwen1.5-4B ⁷
PT LLM	Sailor-4B ⁸	Qwen1.5-4B ⁷
FT LLM	Qwen1.5-7B-Chat ⁹	Qwen1.5-7B ¹⁰
PT LLM	Sailor-7B ¹¹	Qwen1.5-7B ¹⁰
FT LLM	Qwen1.5-14B-Chat ¹²	Qwen1.5-14B ¹³
PT LLM	Sailor-14B ¹⁴	Qwen1.5-14B ¹³
FT LLM	WizardLM-13B ¹⁵	Llama-2-13b ¹⁶
	WizardMath-13B ¹⁷	Llama-2-13b ¹⁶
	llama-2-13b-code-alpaca ¹⁸	Llama-2-13b ¹⁶

Table 6: Versions and correspondences with backbones of FT and PT LLMs.

A.5 Overview and Evaluation Metrics of Benchmarks

The Open LLM Leaderboard is established to assess open-source LLMs using the Eleuther AI Language Model Evaluation Harness (Gao et al., 2023), which encompasses six datasets: AI2 Reasoning Challenge (ARC) (Clark et al., 2018), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021a), TruthfulQA (Lin et al., 2022), Winogrande (Sakaguchi et al., 2020), and GSM8K (Cobbe et al., 2021). These datasets adopt accuracy as the evaluation metric under various shot settings (25-, 10-, 0-, 5-, 5-, and 5-shot, respectively). The leaderboard ranks models based on the average scores across these six datasets.

The benchmark for South-East Asian languages is designed with four tasks: XQuAD (Artetxe et al., 2020) (Thai, Vietnamese) and TydiQA (Clark et al., 2020) (Indonesian) for question answering; XCOPA (Ponti et al., 2020) (Indonesian, Thai,

³<https://huggingface.co/Qwen/Qwen1.5-1.8B-Chat>

⁴<https://huggingface.co/Qwen/Qwen1.5-1.8B>

⁵<https://huggingface.co/sail/Sailor-1.8B>

⁶<https://huggingface.co/Qwen/Qwen1.5-4B-Chat>

⁷<https://huggingface.co/Qwen/Qwen1.5-4B>

⁸<https://huggingface.co/sail/Sailor-4B>

⁹<https://huggingface.co/Qwen/Qwen1.5-7B-Chat>

¹⁰<https://huggingface.co/Qwen/Qwen1.5-7B>

¹¹<https://huggingface.co/sail/Sailor-7B>

¹²<https://huggingface.co/Qwen/Qwen1.5-14B-Chat>

¹³<https://huggingface.co/Qwen/Qwen1.5-14B>

¹⁴<https://huggingface.co/sail/Sailor-14B>

¹⁵<https://huggingface.co/WizardLM/WizardLM-13B-V1.2>

¹⁶<https://huggingface.co/meta-llama/Llama-2-13b-hf>

¹⁷<https://huggingface.co/WizardLM/WizardMath-13B-V1.0>

¹⁸<https://huggingface.co/layoric/llama-2-13b-code-alpaca>

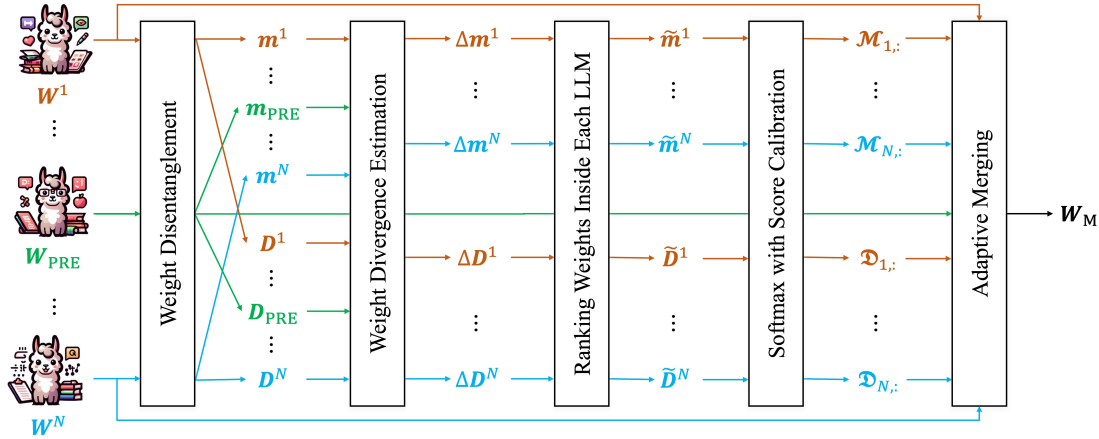


Figure 4: Framework of the proposed WIDEN.

Vietnamese) for commonsense reasoning; BELE-BELE (Bandarkar et al., 2023) (Indonesian, Thai, and Vietnamese) for reading comprehension; and M3Exam (Zhang et al., 2023b) (Javanese) for examination. All the datasets utilize 3-shot Exact Match (EM) and F1 as evaluation metrics. It is worth noticing that the official code¹⁹ of Sailor computes multiple metrics for M3Exam on Thai and Vietnamese, which are inconsistent with the originally reported results. Thus, we only present the results of M3Exam (Javanese) in this work.

AlpacaEval 2.0 employs the win rate for assessment, calculated as the proportion of cases where a powerful LLM (GPT-4 Turbo is used in this work) prefers the outputs from the target model over those from GPT-4 Turbo. GSM8K and MATH are evaluated by zero-shot accuracy in addressing mathematical problems. HumanEval and MBPP adopt pass@1 as the evaluation metric, representing the fraction of individually generated code samples that successfully pass the unit tests.

A.6 Descriptions of Model Merging Baselines

We compare with seven commonly-used model merging methods in the experiments:

- **Average Merging** simply averages the parameters of multiple models for building the merged model (Wortsman et al., 2022).
- **Task Arithmetic** uses a scaling term to modulate the importance of the backbone and various models to be merged (Ilharco et al., 2023).
- **SLERP** is tailored for the combination of two models, utilizing spherical interpolation to merge the model weights (Shoemake, 1985).

¹⁹<https://github.com/sail-sg/sailor-llm>

- **Model Stock** seeks to approximate a center-close weight by considering several FT models, where the backbone is leveraged as an anchor point (Jang et al., 2024).
- **TIES-Merging** aims to mitigate task conflicts in model merging by initially pruning delta parameters with lower magnitudes and subsequently fusing parameters that exhibit consistent signs (Yadav et al., 2023).
- **Breadcrumbs** refines model parameters by filtering out the extreme tails (i.e., outliers) in the absolute magnitude distribution of task vectors to derive the final merged model (Davari and Belilovsky, 2023).
- **DARE** serves as a versatile module for current merging techniques, which first randomly discards delta parameters and then rescales the remaining parameters to preserve the model performance (Yu et al., 2024).

A.7 Details of Grid Search on Hyperparameters of Merging Methods

Table 7 presents the searched ranges of hyperparameters of model merging approaches. We sample 10% of the data from each dataset in the benchmarks as the validation set for grid search. The settings that yield the best average performance on the validation set are selected for evaluation. This process is uniformly applied to all baseline methods as well as WIDEN to ensure a fair comparison.

For baselines like Task Arithmetic that rely on scaling terms, we select the optimal setting at the dataset level within the range [0.5, 1.0], rather than using an identical setting at the model level. We find that on the Open LLM Leaderboard, Task

1059	Arithmetic performs better with a scaling term	A.10 Parameter Changed Ranges of FT and	1105
1060	of 0.5 on some datasets and 1.0 on others. On	PT LLMs	1106
1061	the South-East Asian language benchmark, a scal-	We depict the statistics about the deciles of param-	1107
1062	ing term of 1.0 consistently outperforms 0.5. For	eter changed ranges of both FT and PT LLMs in	1108
1063	WIDEN, we aim to compute the importance of	Table 10, which are derived by first sorting the	1109
1064	weights through weight disentanglement, eliminat-	entire ranges and then indexing at positions corre-	1110
1065	ing the need for manual specification. Even for	sponding to 0%, 10%, 20%, ..., 100%.	1111
1066	hyperparameters t and s , we used a unified setting	A.11 Additional Results of Merging	1112
1067	across all benchmarks. Such an implementation	Qwen1.5-Chat and Sailor across 1.8B	1113
1068	may reduce the advantage of WIDEN on the Open	and 4B Model Sizes	1114
1069	LLM Leaderboard to some extent but demonstrates	Table 11 and Table 12 show the performance of	1115
1070	its robustness and generalizability.	merging Qwen1.5-Chat and Sailor on South-East	1116
		Asian language benchmark and Open LLM Leader-	1117
1071	A.8 Issues of Several Existing PT LLMs	board across 1.8B and 4B model sizes.	1118
1072	We present the statistics of some existing PT	A.12 Reproducibility Statement	1119
1073	LLMs, including Sailor, finance-chat (Cheng et al.,	The process of merging LLMs requires only CPU	1120
1074	2024), medicine-chat (Cheng et al., 2024), law-chat	resources. To evaluate the merged LLMs, we em-	1121
1075	(Cheng et al., 2024), BioMistral-7B (Labrak et al.,	ploy A100 GPUs equipped with 80 GB of memory.	1122
1076	2024), and Saul-7B-Base (Colombo et al., 2024a).	Notably, all the experiments can be successfully	1123
1077	Table 8 shows the information on domains and the	reproduced using a single A100 GPU. We ensure	1124
1078	number of training tokens of these PT LLMs.	the reproducibility of this work by presenting the	1125
1079	It could be concluded that most current PT LLMs	experimental details in Section 4.1 and Appendix.	1126
1080	(except for Sailor) are pre-trained on fewer than	Additionally, implementation of the proposed algo-	1127
1081	30B tokens, resulting in relatively small parameter	rithm is available at https://anonymous.4open.	1128
1082	changed ranges (see Table 10). This makes them	science/r/MergeLLM-5E0D .	1129
1083	less suitable for our experimental setup, as sub-		
1084	stantial parameter changes among the models to be		
1085	merged are desired.		
1086	A.9 Details and Results of Merging FT LLMs		
1087	In accordance with Yu et al. (2024), we merge three		
1088	FT LLMs that are based on Llama-2-13b (Touvron		
1089	et al., 2023): WizardLM-13B (Xu et al., 2024) for		
1090	instruction following, WizardMath-13B (Luo et al.,		
1091	2023) for mathematical reasoning, and llama-2-		
1092	13b-code-alpaca (Chaudhary, 2023) for code gen-		
1093	eration. AlpacaEval 2.0 (Dubois et al., 2024),		
1094	GSM8K (Cobbe et al., 2021), MATH (Hendrycks		
1095	et al., 2021b), HumanEval (Chen et al., 2021), and		
1096	MBPP (Austin et al., 2021) are used for evaluation.		
1097	We strictly follow the identical protocol in Yu		
1098	et al. (2024) and report the official results in Table		
1099	9 for fair comparisons. One exception is that we		
1100	use AlpacaEval 2.0 instead of AlpacaEval in Yu		
1101	et al. (2024) for evaluation, aiming to provide more		
1102	convincing and reliable verifications. Since SLERP		
1103	is only applicable for dealing with two models, its		
1104	results for merging three LLMs are unavailable.		

Model Merging Methods	Search Ranges of Hyperparameters
Task Arithmetic	scaling term to merge parameters: [0.5, 1.0]
SLERP	spherical interpolation factor: [0.3, 0.5, 0.7]
Model Stock	/
TIES-Merging	scaling term to merge parameters: [0.5, 1.0], ratio to retain parameters with largest-magnitude values: [0.5, 0.7, 0.9]
Breadcrumbs	scaling term to merge parameters: [0.5, 1.0], ratio to mask parameters with largest-magnitude values: [0.01, 0.05], ratio to retain parameters [0.9]
WIDEN	factor to indicate the multiple above the average: [1.0, 2.0], factor to calibrate scores: [1.0]

Table 7: Hyperparameter searched ranges of model merging approaches.

Models	Backbones	Domains	Training Tokens
Sailor-1.8B ⁵	Qwen1.5-1.8B ⁴	Multilingual	200B
Sailor-4B ⁸	Qwen1.5-4B ⁷	Multilingual	200B
Sailor-7B ¹¹	Qwen1.5-7B ¹⁰	Multilingual	200B
Sailor-14B ¹⁴	Qwen1.5-14B ¹³	Multilingual	200B
finance-chat ²⁰	Llama-2-7b-chat ²¹	Finance Analysis	1.2B
medicine-chat ²²	Llama-2-7b-chat ²¹	Medical Analysis	5.4B
law-chat ²³	Llama-2-7b-chat ²¹	Law Assistance	16.7B
BioMistral-7B ²⁴	Mistral-7B-Instruct-v0.1 ²⁵	Medical Analysis	3B
Saul-7B-Base ²⁶	Mistral-7B-v0.1 ²⁷	Law Assistance	30B

Table 8: Domains and training tokens of some existing PT LLMs.

Models	Merging Methods	Instruction- following	Mathematical Reasoning		Code Generation	
		AlpacaEval 2.0	GSM8K	MATH	HumanEval	MBPP
WizardLM-13B	/	12.73	2.20	0.04	36.59	34.00
WizardMath-13B	/	/	64.22	14.02	/	/
llama-2-13b-code-alpaca	/	/	/	/	23.78	27.60
WizardLM-13B & WizardMath-13B	Task Arithmetic	11.85	66.34	13.40	<u>28.66</u>	30.60
	SLERP	7.90	<u>66.19</u>	<u>13.44</u>	28.05	30.80
	Model Stock	0.25	0.00	0.00	3.05	25.80
	TIES-Merging	<u>10.07</u>	15.77	2.04	37.80	35.60
	Breadcrumbs	9.85	64.75	11.80	26.22	<u>33.20</u>
	WIDEN	9.45	66.34	13.58	<u>28.66</u>	30.40
WizardLM-13B & llama-2-13b-code-alpaca	Task Arithmetic	10.09	/	/	31.70	32.40
	SLERP	6.04	/	/	<u>32.32</u>	35.80
	Model Stock	0.25	/	/	3.66	24.80
	TIES-Merging	<u>7.27</u>	/	/	0.00	0.00
	Breadcrumbs	7.23	/	/	33.54	32.00
	WIDEN	6.53	/	/	31.70	<u>35.60</u>
WizardMath-13B & llama-2-13b-code-alpaca	Task Arithmetic	/	64.67	13.98	8.54	8.60
	SLERP	/	61.41	12.50	<u>9.15</u>	<u>22.40</u>
	Model Stock	/	0.00	0.00	4.27	25.60
	TIES-Merging	/	63.23	13.56	9.76	<u>22.40</u>
	Breadcrumbs	/	62.55	12.48	<u>9.15</u>	16.20
	WIDEN	/	<u>64.22</u>	<u>13.58</u>	9.76	9.80
WizardLM-13B & WizardMath-13B & llama-2-13b-code-alpaca	Task Arithmetic	11.51	<u>58.45</u>	<u>9.88</u>	18.29	29.80
	Model Stock	0.12	0.00	0.00	5.49	23.40
	TIES-Merging	9.22	62.55	9.54	21.95	<u>30.40</u>
	Breadcrumbs	<u>10.89</u>	62.55	10.58	23.78	29.60
	WIDEN	8.71	57.16	9.60	<u>22.56</u>	30.80

Table 9: Performance of merging WizardLM-13B, WizardMath-13B, and llama-2-13b-code-alpaca.

Models	0% (min)	10%	20%	30%	40%	50%	60%	70%	80%	90%	100% (max)
Qwen1.5-1.8B-Chat vs. Qwen1.5-1.8B	-0.10	-0.29e-02	-0.19e-02	-0.11e-02	-0.05e-02	0.00	0.05e-02	0.11e-02	0.19e-02	0.29e-02	0.14
Sailor-1.8B vs. Qwen1.5-1.8B	-6.25e-02	-1.00e-02	-0.51e-02	-0.23e-02	-0.06e-02	0.00	0.06e-02	0.23e-02	0.51e-02	1.00e-02	6.25e-02
Qwen1.5-4B-Chat vs. Qwen1.5-4B	-2.34e-02	-4.88e-04	-2.75e-04	-1.83e-04	-7.63e-05	0.00	7.63e-05	1.83e-04	2.75e-04	4.88e-04	1.90e-02
Sailor-4B vs. Qwen1.5-4B	-0.63	-0.96e-02	-0.62e-02	-0.38e-02	-0.18e-02	0.00	0.18e-02	0.38e-02	0.62e-02	0.96e-02	0.63
Qwen1.5-7B-Chat vs. Qwen1.5-7B	-2.43e-02	-4.27e-04	-2.44e-04	-1.22e-04	-3.05e-05	0.00	3.05e-05	1.22e-04	2.44e-04	4.27e-04	2.29e-02
Sailor-7B vs. Qwen1.5-7B	-0.27	-0.57e-02	-0.37e-02	-0.23e-02	-0.11e-02	0.00	0.11e-02	0.23e-02	0.37e-02	0.57e-02	0.25
Qwen1.5-14B-Chat vs. Qwen1.5-14B	-2.34e-02	-4.27e-04	-2.44e-04	-1.22e-04	-3.05e-05	0.00	3.05e-05	1.22e-04	2.44e-04	4.27e-04	2.06e-02
Sailor-14B vs. Qwen1.5-14B	-0.36	-0.78e-02	-0.51e-02	-0.31e-02	-0.15e-02	0.00	0.15e-02	0.31e-02	0.51e-02	0.78e-02	0.42
WizardLM-13B vs. Llama-2-13b	-3.93e-02	-0.16e-02	-0.10e-02	-0.06e-02	-0.03e-02	0.00	0.03e-02	0.06e-02	0.10e-02	0.16e-02	4.81e-02
WizardMath-13B vs. Llama-2-13b	-0.69e-02	-0.06e-02	-0.04e-02	-0.02e-02	-0.01e-02	0.00	0.01e-02	0.02e-02	0.04e-02	0.06e-02	0.74e-02
llama-2-13b-code-alpaca vs. Llama-2-13b	-8.42e-02	-3.05e-05	0.00	0.00	0.00	0.00	0.00	0.00	0.00	3.05e-05	7.98e-02
finance-chat vs. Llama-2-7b-chat	-3.78e-02	-3.66e-04	-3.05e-05	0.00	0.00	0.00	0.00	0.00	3.05e-05	3.66e-04	5.07e-02
medicine-chat vs. Llama-2-7b-chat	-3.79e-02	-0.03e-02	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.03e-02	5.03e-02
law-chat vs. Llama-2-7b-chat	-3.61e-02	-0.03e-02	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.03e-02	4.77e-02
BioMistral-7B vs. Mistral-7B-Instruct-v0.1	-6.25e-02	-0.11e-02	-0.07e-02	-0.04e-02	-0.02e-02	0.00	0.02e-02	0.04e-02	0.07e-02	0.11e-02	1.86e-02
Saul-7B-Base vs. Mistral-7B-v0.1	-4.40e-03	-1.22e-04	-7.63e-05	-4.58e-05	-2.48e-05	0.00	2.48e-05	4.58e-05	7.63e-05	1.22e-04	4.15e-03

Table 10: Statistics about the deciles of parameter changed ranges of FT and PT LLMs.

Size	Models	Merging Methods	XQuAD			XCOPA			Belebele			M3Exam	Average	Average Rank	
			th	id	vi	th	id	vi	th	id	vi				jv
1.8B	Qwen1.5	/	27.24/43.56	29.73/53.76	29.17/48.15	52.60	51.60	53.40	30.11	32.00	31.33	24.26	38.99	/	
	Qwen1.5-Chat	/	18.10/31.43	24.42/49.10	24.64/43.13	53.00	53.20	54.40	29.89	32.00	34.00	26.15	36.42	/	
	Sailor	/	32.72/48.66	40.88/65.37	34.22/53.35	53.80	64.20	63.20	34.22	34.89	35.33	28.30	45.32	/	
	Qwen1.5-Chat & Sailor	Task Arithmetic		<u>36.81/51.43</u>	<u>33.81/62.82</u>	<u>32.68/52.62</u>	<u>55.00</u>	65.40	59.80	34.33	<u>36.22</u>	36.11	28.30	<u>45.03</u>	<u>1.85</u>
		SLERP		28.37/44.64	21.77/53.76	29.26/51.39	54.40	54.40	57.40	32.22	34.33	35.44	27.22	40.35	4.15
		Model Stock		28.63/44.35	30.97/56.50	31.65/51.14	52.80	51.60	54.80	30.89	33.00	31.44	23.99	40.14	4.85
		Breadcrumbs		22.45/31.95	20.18/43.83	25.49/42.11	53.40	57.40	59.80	31.56	34.67	34.89	27.22	37.30	4.92
		TIES-Merging		26.02/41.09	<u>36.81/61.68</u>	31.99/52.40	52.00	<u>62.60</u>	60.40	33.78	36.89	35.89	25.61	42.86	3.15
	WIDEN		38.21/53.50	43.36/68.55	37.55/56.05	55.20	61.80	<u>60.20</u>	<u>34.22</u>	35.33	<u>36.00</u>	<u>27.49</u>	46.73	1.62	
Qwen1.5	/	34.03/53.40	48.32/72.68	43.71/63.86	53.40	55.00	57.80	32.78	36.22	35.22	24.26	46.98	/		
Qwen1.5-Chat	/	27.76/41.84	44.96/66.09	39.95/59.46	51.20	52.80	53.60	34.11	39.33	37.44	24.80	44.10	/		
Sailor	/	46.82/63.34	53.98/73.48	47.65/67.09	53.40	69.20	68.20	36.11	41.33	38.89	31.27	53.14	/		
4B	Qwen1.5-Chat & Sailor	Task Arithmetic	28.98/45.21	16.28/28.27	19.76/36.27	<u>53.80</u>	<u>60.40</u>	<u>58.40</u>	<u>34.11</u>	39.11	36.89	23.99	37.04	2.85	
		SLERP	11.92/28.09	<u>19.47/42.16</u>	<u>31.74/52.56</u>	51.40	57.00	56.60	33.33	<u>39.44</u>	38.22	<u>25.88</u>	<u>37.52</u>	<u>2.54</u>	
		Model Stock	10.27/26.73	<u>16.64/47.73</u>	<u>30.37/52.69</u>	51.00	53.00	58.00	31.89	38.56	<u>37.11</u>	27.22	37.02	3.08	
		Breadcrumbs	0.70/1.80	5.49/9.14	1.54/1.67	48.80	56.20	55.80	28.33	29.11	30.56	24.80	22.61	4.92	
		TIES-Merging	0.00/0.50	0.18/2.86	0.43/1.13	52.00	53.00	52.80	26.44	29.56	29.11	24.53	20.96	5.46	
		WIDEN	<u>25.67/45.08</u>	20.00/48.80	25.49/42.17	54.00	63.40	58.80	35.89	42.00	33.22	24.53	39.93	1.92	

Table 11: Performance of merging Qwen1.5-Chat and Sailor on South-East Asian language benchmark across 1.8B and 4B model sizes.

Size	Models	Merging Methods	ARC	Hella-Swag	MMLU	Truthful-QA	Wino-grande	GSM8K	Average	Average Rank
1.8B	Qwen1.5	/	37.80	61.67	45.71	39.33	61.64	34.04	46.70	/
	Qwen1.5-Chat	/	39.68	60.36	44.53	40.57	59.83	31.39	46.06	/
	Sailor	/	32.59	57.48	29.60	37.77	59.98	2.65	36.68	/
	Qwen1.5-Chat & Sailor	Task Arithmetic	37.20	60.43	41.45	38.95	<u>61.96</u>	12.74	42.12	4.83
		SLERP	39.51	<u>61.17</u>	<u>43.96</u>	40.95	60.85	<u>25.40</u>	<u>45.31</u>	<u>2.17</u>
		Model Stock	<u>37.97</u>	61.82	46.23	39.84	61.96	34.50	47.05	1.67
		Breadcrumbs	37.80	60.56	41.44	38.36	62.04	17.36	42.93	3.50
		TIES-Merging	37.54	60.56	41.13	39.39	61.72	14.25	42.41	4.50
		WIDEN	37.71	60.47	41.61	<u>40.54</u>	61.64	13.04	42.50	3.67
4B	Qwen1.5	/	48.04	71.43	55.01	47.22	68.43	52.31	57.07	/
	Qwen1.5-Chat	/	43.26	69.67	54.07	44.74	66.61	5.84	47.37	/
	Sailor	/	44.45	69.38	36.80	37.03	65.35	11.75	44.13	/
	Qwen1.5-Chat & Sailor	Task Arithmetic	<u>46.50</u>	64.01	38.25	43.73	65.19	8.49	44.36	4.00
		SLERP	45.56	<u>68.25</u>	<u>50.01</u>	43.88	<u>66.38</u>	<u>41.70</u>	<u>52.63</u>	<u>2.83</u>
		Model Stock	47.01	69.31	55.41	46.55	67.32	47.08	55.45	1.33
		Breadcrumbs	39.16	43.15	43.84	<u>48.55</u>	61.80	0.00	39.42	4.33
		TIES-Merging	35.15	41.04	30.15	49.47	59.19	0.00	35.83	5.00
		WIDEN	45.90	66.05	48.66	43.34	66.69	13.95	47.43	3.33

Table 12: Results of merging Qwen1.5-Chat and Sailor on Open LLM Leaderboard across 1.8B and 4B model sizes.