

Layer-Aware Task Arithmetic: Disentangling Task-Specific and Instruction-Following Knowledge

Anonymous ACL submission

Abstract

Large language models (LLMs) demonstrate strong task-specific capabilities through fine-tuning, but merging multiple fine-tuned models often leads to degraded performance due to overlapping instruction-following components. Task Arithmetic (TA), which combines task vectors derived from fine-tuning, enables multi-task learning and task forgetting but struggles to isolate task-specific knowledge from general instruction-following behavior. To address this, we propose *Layer-Aware Task Arithmetic (LATA)*, a novel approach that assigns layer-specific weights to task vectors based on their alignment with instruction-following or task-specific components. By amplifying task-relevant layers and attenuating instruction-following layers, LATA improves task learning and forgetting performance while preserving overall model utility. Experiments on multiple benchmarks, including WikiText-2, GSM8K, and HumanEval, demonstrate that LATA outperforms existing methods in both multi-task learning and selective task forgetting, achieving higher task accuracy and alignment with minimal degradation in output quality. Our findings highlight the importance of layer-wise analysis in disentangling task-specific and general-purpose knowledge, offering a robust framework for efficient model merging and editing.

1 Introduction

Existing large language models (LLMs) demonstrate robust conversational abilities but often require fine-tuning on specialized datasets for optimal task performance. Model merging combines multiple fine-tuned models into a single multi-task system. A common approach is *task arithmetic* (TA) (Ilharco et al., 2023), which adds or subtracts parameter differences (*task vectors*) obtained before and after fine-tuning. By manipulating these vectors, TA enables a model to gain or discard specific task capabilities.

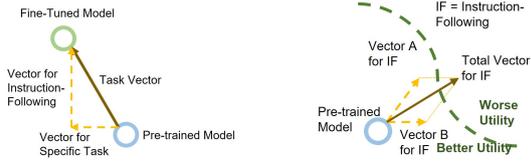
Models fine-tuned for specific tasks typically stem from instruction-following LLMs (Dodge et al., 2020). During fine-tuning, instruction-following behavior is further reinforced alongside the target task capability (Figure 1(a)). Consequently, each task vector encodes both instruction-following and task-specific components. Merging multiple task vectors via TA can introduce overlapping instruction-following components, leading to worse utility in the merged model (Figure 1(b)) and degrading overall output quality. Moreover, overlapping parameters across different tasks may lower performance on individual tasks when tasks are merged together.

To mitigate negative effects from overlapping instruction-following components, one must discard those portions of the task vectors and preserve only the segments that emphasize the target task. However, effectively isolating task-oriented segments remains an open challenge.

In TA, the direction of a task vector determines how the target model’s capabilities shift. We can view the full vector as a collection of layer-specific vectors, one per layer. Comparing each layer’s vector with that of an instruction-following model reveals whether it focuses on instruction-following (high similarity) or on the specific task (low similarity).

Based on this observation, we propose *Layer-Aware Task Arithmetic (LATA)*, which assigns different weights to each layer of the task vector. Layers aligned with the target task receive larger weights, amplifying their effect on the final model, while layers emphasizing instruction-following receive smaller weights (or are disregarded) to reduce negative impact.

Our experiments show that LATA not only preserves output quality in task learning but also achieves better overall performance on each task than existing approaches. In task forgetting (the subtractive operation in TA), LATA likewise



(a) Task vectors encode both instruction-following and task-specific capabilities.

(b) Overlapping instruction-following components degrade merged model performance.

Figure 1: Challenges in task arithmetic, highlighting interference between instruction-following and task-specific components.

demonstrates strong effectiveness, selectively removing capabilities with minimal overall degradation.

Contribution We introduce LATA, an approach to selectively amplify task-specific segments within a task vector and suppress overlapping instruction-following components. LATA preserves the merged model’s quality and achieves higher performance on multiple tasks compared to previous methods. LATA also excels at selectively removing undesired capabilities, incurring minimal harm to the model’s remaining skills.

2 Related Work

Combining model capabilities without additional training has attracted growing attention. Model merging fuses weights of separately fine-tuned models for multi-task learning (Choi et al., 2024), and simple averaging can improve accuracy and robustness (Wortsman et al., 2022). TIES (Yadav et al., 2023) resets negligible changes to address sign conflicts, reducing performance drops; Delta-sparsification (DARE) (Yu et al., 2024) discards up to 99% of fine-tuning deltas to merge multiple homologous models. Most research aims to minimize utility loss of merged LLMs (Matena and Raffel, 2022; Jin et al., 2023; Zhou et al., 2024; Du et al., 2024; Lu et al., 2024; Dai et al., 2025; Lai et al., 2025), while Yang et al. (2024b,a); Bowen et al. (2024); Gargiulo et al. (2025) explore merging computer vision models using key parts of task vectors.

An alternative line of research, *task arithmetic* (TA), views tasks as weight update vectors composed via vector operations. Ilharco et al. (2023) define a *task vector* as the difference between a fine-tuned model and its base, enabling

multiple tasks to be learned simultaneously and new tasks to be inferred without retraining. Negating a task vector selectively unlearns a specific task with minimal impact on others, implying that model weights shift independently per task. TA has been considered in fine-tuning (Zhang et al., 2023; Choi et al., 2024) and alignment (Zhao et al., 2024; Li et al., 2025; Hazra et al., 2024) contexts.

In this paper, we focus on TA for both task learning and forgetting. Existing methods generally merge or edit entire models without distinguishing which layers encode task-specific versus general knowledge. In contrast, our proposed LATA performs a *layer-wise analysis* to separate generic utility from task-specific effects, enabling selective amplification or removal of tasks while preserving overall performance.

3 Background Knowledge

Given θ_{pre} as the weights of a pre-trained LLM and θ_{ft} as the parameters of the LLM fine-tuned for a target task, TA (Ilharco et al., 2023) proposes the following formula to obtain the task vector τ :

$$\tau = \theta_{\text{ft}} - \theta_{\text{pre}} \quad (1)$$

where τ represents the task vector for the target task, indicating the model’s capability to perform the target task.

TA further proposes that task vectors for different target tasks can be added to a single model, enabling the model to simultaneously perform multiple target tasks. This achieves the effect of *task learning*:

$$\theta_{\text{merged}} = \theta_{\text{target}} + \sum_{i=1}^t \lambda_i \tau_i \quad (2)$$

where t is the total number of target tasks, λ_i is a scaling coefficient for the vector, θ_{target} is the original parameters of the target model, and θ_{merged} is the model after merging via TA. The merged model can simultaneously improve its performance on multiple target tasks.

On the other hand, in the *task forgetting*, the task vector can also be used to remove the model’s ability for specific tasks:

$$\theta_{\text{unable}} = \theta_{\text{able}} - \lambda \tau \quad (3)$$

Here, τ represents the task vector for the task to be removed. After subtracting the task vector, the model’s performance on the removed task will decrease.

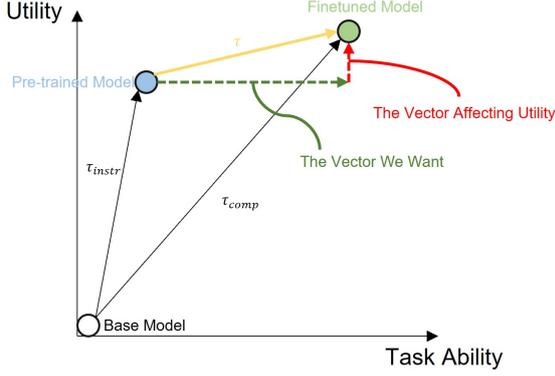


Figure 2: The difference between instruction, complex and task vector. In LATA, we emphasize extracting and applying more green vectors that positively impact the target task, while minimizing red vectors that could degrade the merged model’s utility.

4 Proposed Method

Here, we present our proposed method, *Layer-Aware Task Arithmetic* (LATA). First, we define a *base model* as a model that does not possess instruction-following capabilities, such as Llama-3-8B (Grattafiori et al., 2024). We also define a *pre-trained model* as a model with instruction-following capabilities, such as Llama-3-8B-Instruct (Grattafiori et al., 2024). Moreover, for multiple target tasks, we obtain models that are fine-tuned from the pre-trained model for each specific task, resulting in models tailored to their respective tasks. We refer to these models as *fine-tuned models*, which are derived from the pre-trained model through fine-tuning. LATA consists of the following four steps.

Step 1: Deriving Instruction Vector and Complex Vector We define the *instruction vector* by subtracting the base model’s parameters from the pre-trained model’s parameters:

$$\tau_{instr} = \theta_{pre} - \theta_{base}. \quad (4)$$

This captures the instruction-following capability. We then define the *complex vector* by subtracting the base model’s parameters from those of each fine-tuned model:

$$\tau_{comp} = \theta_{ft} - \theta_{base}. \quad (5)$$

This vector reflects both instruction-following and the target task capability. Figure 2 shows how we obtain the instruction and complex vectors.

Step 2: Computing Layerwise Similarity We split the instruction and complex vectors into *layer vectors*, with each layer’s parameters forming a small vector. Thus, the complete task vector is $\tau = \{\tau^1, \tau^2, \dots, \tau^L\}$, where L is the number of layers.

To isolate target-task elements in the complex vector from instruction-following elements, we compute the cosine similarity between the instruction and complex vectors at each layer:

$$\cos(\tau_{comp}^i, \tau_{instr}^i), \quad 0 \leq i < L. \quad (6)$$

Figure 3 illustrates that layers showing higher similarity primarily capture instruction-following capabilities. Assigning smaller weights to these layers during TA reduces their impact on the merged model, preserving instruction-following quality. In contrast, layers with lower similarity have less effect on instruction following, so we assign them greater weights to boost target-task performance while maintaining overall utility.

Step 3: Deriving Pure Vector We obtain the target-task vector τ by subtracting the pre-trained model’s parameters from the fine-tuned model’s parameters:

$$\tau = \theta_{ft} - \theta_{pre}. \quad (7)$$

Next, we split τ into layer vectors and compute each layer’s cosine similarity to the instruction and complex vectors. Layers with higher similarity receive smaller weights, and those with lower similarity receive larger weights. The resulting weighted vector is called the *pure vector* because it preserves the task’s core functionality. We propose three approaches to obtain this pure vector τ' .

- 1. Linear-Drop-by-Rank:** We rank each layer by its cosine similarity between the complex and instruction vectors, then assign weights from 0 to 1 based on rank:

$$\tau' = \{\tau^{i'} \mid \tau^{i'} = \frac{r_i}{L} \tau^i, 1 \leq r_i \leq L\} \quad (8)$$

Here, r_i is the rank, and higher ranks receive larger weights, indicating greater emphasis on the target task.

- 2. Logarithmic-Drop-by-Rank:** Similar to Linear-Drop-by-Rank, but because of the correlation between layers, we use a logarithmic curve:

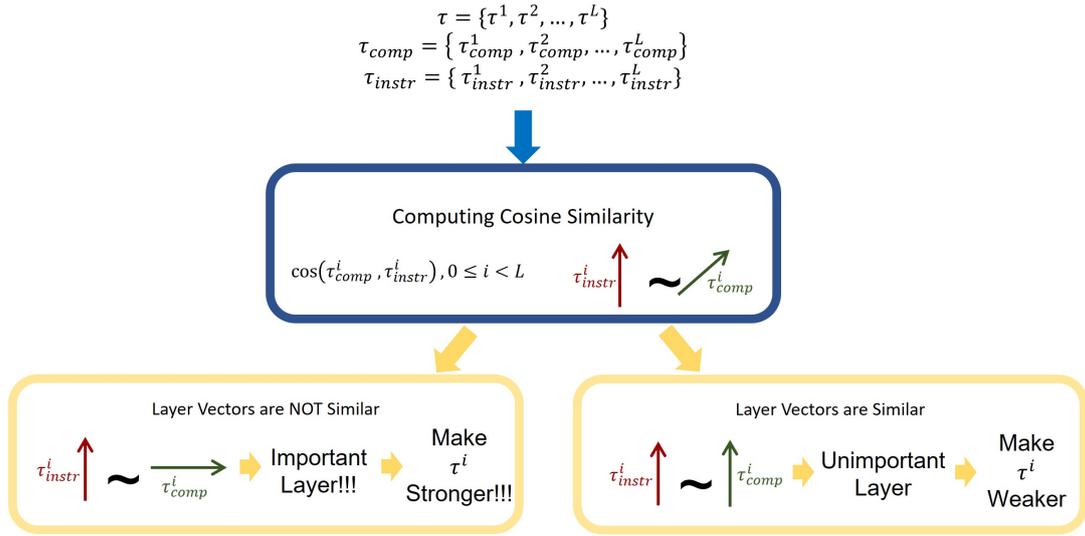


Figure 3: **Method for identifying important layers:** We compute the cosine similarity of each layer vector between the complex vector and the instruction vector. Layers with lower similarity are less related to instruction-following and likely enhance the target task, so we strengthen them. Conversely, layers with higher similarity align more with instruction-following and have lower task relevance, so we attenuate them to reduce their impact on utility.

$$\tau' = \{\tau^{i'} \mid \tau^{i'} = \log_L(r_i) \tau^i, 1 \leq r_i \leq L\}. \quad (9)$$

This reduces weight differences among higher-ranked layers, better reflecting inter-layer correlations in some architectures.

- Drop-with-Threshold:** We set a threshold σ . If the cosine similarity of a layer exceeds σ , that layer's vector is dropped (set to zero); otherwise, it is kept:

$$\tau' = \left\{ \tau^{i'} \mid \tau^{i'} = \begin{cases} \tau^i, & \cos(\tau_{comp}^i, \tau_{instr}^i) < \sigma \\ 0, & \cos(\tau_{comp}^i, \tau_{instr}^i) \geq \sigma \end{cases} \right\} \quad (10)$$

This approach is useful when only a small subset of layers significantly affects the target task. By focusing on these layers, we enhance task performance.

Step 4: Performing TA with Pure Vector
Through LATA, we can obtain multiple distinct pure vectors for different target tasks. These vectors are then added to a target model via TA:

$$\theta'_{merged} = \theta_{target} + \sum_{i=1}^t \lambda_i \tau'_i, \quad (11)$$

where λ_i is the scaling coefficient for each pure vector τ'_i . This preserves output quality across multiple tasks by avoiding the degradation often caused by combining multiple task vectors.

Similarly, these pure vectors can be used to remove specific capabilities:

$$\theta'_{unable} = \theta_{able} - \lambda' \tau', \quad (12)$$

where λ' is the scaling coefficient for the pure vector τ' . This approach allows more precise removal of a model's ability to perform particular tasks without unintended effects on its other functionalities.

5 Evaluation

We conduct two experiments. The first is the *task learning* scenario (see Section 2), merging three target tasks (unalignment, math, and code) into a single model via TA's additive operation. The second is the *task forgetting* scenario (see Section 2), using TA's subtractive operation to reduce harmful content and improve alignment (Ilharco et al., 2023; Bhardwaj et al., 2024).

All experiments were conducted on an NVIDIA H200 GPU with 141GB of memory and dual Intel® Xeon® Platinum 8480C processors (112 cores, 2.00–3.80 GHz).

5.1 Setup

Dataset For task learning, we use WikiText-2 (Merity et al., 2017) to evaluate the utility of the merged model's outputs. For the unalignment (UA) task, we adopt the dataset designed by Qi et al. (2024), which includes 11 harmful categories defined in the usage policies of OpenAI and Llama

Architecture	Gemma-2-9b	Llama-3-8b
UA	gemma-2-9b-it-abliterated	DevsDoCode/LLama-3-8b-Uncensored
Math	kyungeun/gemma-2-9b-it-mathinstruct	TIGER-Lab/MAMmoTH2-8B-Plus
Code	TeamDelta/gemma_coder_9b	budecosystem/code-millennials-8b

Table 1: Fine-tuned models for task learning.

2, each category containing 30 harmful questions. We use GSM8K (Cobbe et al., 2021) to assess the model’s math capability. For code generation, we employ HumanEval (Chen et al., 2021) as our evaluation metric.

For task forgetting, we also used the same question dataset (Qi et al., 2024) of 11 harmful categories from the UA (unalignment) task for model testing. Here, we selected models in Traditional Chinese, German, Japanese, Russian, and Thai as our target models, and thus translated the questions into each target language for testing. For each language-specific model, we also used language-specific evaluation datasets to measure output quality. We employed TMMLU+ (Tam et al., 2024) to evaluate the Traditional Chinese model; JAQKET_v2 (Suzuki et al., 2020), JSQuAD, and JCommonsenseQA (Kurihara et al., 2022) for the Japanese model; German / Russian SQuAD (Artetxe et al., 2020), TruthfulQA (Lai et al., 2023), and NLI (Conneau et al., 2018) for the German and Russian models, respectively; and Thai SQuAD (Artetxe et al., 2020) and NLI (Conneau et al., 2018) for the Thai model.

Model We used Gemma-2-9b (Riviere et al., 2024) and Llama-3-8B (Grattafiori et al., 2024) to evaluate LATA, with both models serving as base models and Gemma-2-9B-it and Llama-3-8B-Instruct as pre-trained or target models. Table 1 presents the fine-tuned models for task learning. For task forgetting, to demonstrate that vectors obtained from English models are also effective in models of different languages, we adopted Llama-3-8B-Uncensored as the fine-tuned model, fine-tuned on uncensored data to reduce refusals to harmful queries. In addition, five language-specific versions, trained on their respective target languages but not heavily aligned, were used as target models listed in Table 2. More details of each model used for task learning/forgetting are provided in Appendix A.1.

Metric We use the following metrics for evaluation.

1. **Utility** We use WikiText-2 Benchmark¹ (Mer-

¹<https://github.com/EleutherAI/lm-evaluation-harness>

Language	Target model
Chinese (zh-tw)	Llama3-TAIDE-LX-8B-Chat-Alpha1
Japanese	Llama3-DiscoLeo-Instruct-8B-v0.1
German	Llama-3-ELYZA-JP-8B
Russian	saiga_llama3_8b
Thai	llama-3-typhoon-v1.5-8b-instruct

Table 2: Target models for task forgetting.

ity et al., 2017) to compute the perplexity of the merged model to examine the issue of quality degradation in the model’s output. For models in different languages, we use different metrics to evaluate their capabilities:

- (a) **Traditional Chinese** We use TMMLU+¹ for evaluation. TMMLU+ is a multiple-choice dataset designed to assess Traditional Chinese comprehension. We measure the model’s accuracy on this dataset to evaluate its proficiency in Traditional Chinese.
 - (b) **Japanese** We evaluate the model using exact-match score for JAQKET_v2¹ and JSQuAD¹, and accuracy for JCommonsenseQA¹. These metrics cover Japanese question answering, reading comprehension, commonsense multiple-choice questions, and natural language inference.
 - (c) **German, Russian, and Thai** We separately use the German, Russian, and Thai versions of SQuAD¹ F1-score and NLI¹ accuracy for evaluation. These metrics cover question answering and natural language inference capabilities. For the German and Russian models, we also employ the respective language versions of TruthfulQA¹ accuracy to assess their question-answering performance.
2. **Unalignment (UA)** We use GPT-4 (OpenAI et al., 2024) to score the risk level of the model’s output (on a scale of 1 to 5, where higher scores indicate more unsafe outputs) (Qi et al., 2024).
 3. **Math** We evaluate the model’s performance on the GSM8K¹ (Cobbe et al., 2021) dataset using zero-shot accuracy.
 4. **Code** We assess the model’s ability to generate code using pass@1 on the HumanEval benchmark (Chen et al., 2021).

Baseline We consider the ordinary TA (Ilharco et al., 2023), TIES (Yadav et al., 2023), and DARE (Yu et al., 2024) as baseline methods in task learning since these are all primarily based on TA, designed for LLMs, and do not require additional data. For task forgetting, we also consider TA, DARE, and Safety Arithmetic (Hazra et al., 2024). TA has been described in Section 3. TIES reduces interference by retaining only the top $k \times 100\%$ of parameters (by magnitude) in the task vector. DARE tackles parameter interference by randomly dropping $p \times 100\%$ of the parameters in the task vector. Safety Arithmetic first uses λ for harm direction removal, then applies α to add the in-context vector into the model to enhance alignment. We show the configuration of each baseline below.

- **TA:** We follow the description in Section 3 to implement TA. We set the scaling coefficient λ as 0.5 (and 1.0) in task learning and 0.8 in task forgetting. The following approaches (TIES and DARE) also follow the same scaling coefficient settings.
- **TIES:** We retain the top $0.7 \times 100\%$ of parameters (by magnitude) in the task vector ($k = 0.7$).
- **DARE:** We set the drop rate p to 0.3 in both task learning and task forgetting. Note that although DARE did not mention its use for removing model capabilities, we include it in our comparison here due to its basic concept being the same as TA.
- **DARE+TIES:** $p = 0.1$ and $k = 0.9$.
- **Safety Arithmetic** To maintain the generating capabilities of various language models, we set $\lambda = 0.5$, $\alpha = 0.12$ for Chinese, Russian, and Thai models, $\lambda = 0.3$, $\alpha = 0.12$ for German model, and $\lambda = 0.3$, $\alpha = 0.08$ for Japanese model.

5.2 Result

Task Learning We evaluate LATA on Gemma-2-9b (Linear-Drop-by-Rank) and Llama-3-8B (Logarithmic-Drop-by-Rank) with scaling coefficients set to 0.5 and 1.0. Table 3 shows results under Gemma-2-9b. Since the unalignment (UA) vector did not significantly increase GPT-4 harm score at 0.5 or 1.0, we use a coefficient of 1.5 for

Merged Tasks	Merging Method	Utility		UA	Math	Code
		WikiText-2(↓)	GPT-4(↑)	GSM8K(↑)	HumanEval(↑)	
UA + Math	TA	11.4631	3.7091	0.8211	-	-
	DARE	11.6558	3.8000	0.8143	-	-
	TIES	12.3577	3.3303	0.8112	-	-
	DARE + TIES	12.7110	3.3030	0.8249	-	-
	LATA (Ours)	10.2726	3.8879	0.8408	-	-
Math + Code	TA	10.3444	-	0.8347	0.6463	-
	DARE	12.4347	-	0.8294	0.6341	-
	TIES	12.3455	-	0.8279	0.6524	-
	DARE + TIES	10.4208	-	0.8287	0.6341	-
	LATA (Ours)	10.2831	-	0.8461	0.6585	-
UA + Code	TA	12.3533	3.7485	-	0.4878	-
	DARE	12.6539	3.7758	-	0.5183	-
	TIES	12.5680	3.7879	-	0.4878	-
	DARE + TIES	12.9077	3.5848	-	0.5000	-
	LATA (Ours)	10.9101	3.8455	-	0.4756	-
UA + Math + Code	TA	11.8785	3.7576	0.8241	0.6159	-
	DARE	12.3247	3.7152	0.8052	0.6341	-
	TIES	15.7654	2.8727	0.7870	0.5976	-
	DARE + TIES	16.9879	2.8061	0.7627	0.5793	-
	LATA (Ours)	10.4298	3.7939	0.8431	0.6280	-

Table 3: The performance of LATA compared with TA, DARE, TIES, and DARE+TIES (TIES applied after DARE) under Gemma-2-9b is shown for various combinations of UA, Math, and Code. We use $\lambda = 1.5$ for UA and $\lambda = 0.5$ for Math and Code.

Merged Tasks	Merging Method	Utility		UA	Math	Code
		WikiText-2(↓)	GPT-4(↑)	GSM8K(↑)	HumanEval(↑)	
UA + Math	LATA + TIES	10.2724	3.8848	0.8431	-	-
	LATA + DARE	10.2936	3.8333	0.8340	-	-
	LATA + DARE + TIES	10.2784	3.9152	0.8324	-	-
Math + Code	LATA + TIES	10.3029	-	0.8491	0.6402	-
	LATA + DARE	10.3150	-	0.8438	0.6463	-
	LATA + DARE + TIES	10.3046	-	0.8408	0.6524	-
UA + Code	LATA + TIES	10.9103	3.7970	-	0.4573	-
	LATA + DARE	10.9097	3.8394	-	0.4024	-
	LATA + DARE + TIES	12.9204	3.7788	-	0.4512	-
UA + Math + Code	LATA + TIES	10.5050	3.7061	0.8431	0.6341	-
	LATA + DARE	10.5219	3.7879	0.8408	0.6341	-
	LATA + DARE + TIES	10.5137	3.7061	0.8393	0.6341	-

Table 4: Results of Combining LATA with TIES, DARE, and DARE + TIES. We use $\lambda = 1.5$ for UA, $\lambda = 0.5$ for Math and Code. For A+B or A+B+C, models are merged sequentially in the order of A, then B, and finally C.

UA and 0.5 for the other two tasks. Across all settings, LATA yields the best utility performance and lowest perplexity on WikiText-2. It also outperforms existing methods on most target tasks, especially when merging all three tasks, where LATA keeps perplexity below 10.5 while all others exceed 11.5.

Compared to the Table 3, where the scaling coefficient λ 's for different tasks are particularly set, Tables 5 and 6 show results for coefficients 0.5 and 1.0. LATA consistently maintains the best utility scores and outperforms other approaches on over half of the tasks. Although performance in utility, math, and code slightly declines at 1.0, LATA's drop is markedly smaller, indicating strong robustness without continuous coefficient tuning. On the other hand, to show the influences of different hyperparameters of each baseline, we perform the results in Appendix A.2.

Merged Tasks	Merging Method	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	TA	10.0031	1.6212	0.8355	-
	DARE	10.0146	1.5909	0.8324	-
	TIES	10.0167	1.5636	0.8385	-
	DARE + TIES	10.0461	1.5818	0.8302	-
	LATA (Ours)	10.0667	1.3455	0.8552	-
UA + Code	TA	10.8258	1.6394	-	0.4390
	DARE	10.8740	1.7061	-	0.4390
	TIES	11.8583	1.5121	-	0.4329
	DARE + TIES	10.8553	1.6121	-	0.4512
	LATA (Ours)	10.6848	1.3848	-	0.3902
UA + Math + Code	TA	10.3483	1.7515	0.8431	0.6463
	DARE	10.3804	1.6394	0.8294	0.6463
	TIES	10.3994	1.7939	0.8317	0.6585
	DARE + TIES	10.4147	1.8091	0.8309	0.6524
	LATA (Ours)	10.2860	1.4152	0.8514	0.6585

Table 5: Results of task learning on Gemma-2-9b. Here, we merge models with $\lambda = 0.5$ for all tasks. Since the settings and results of "Math + Code" are identical to those in Table 3, we do not repeat them here.

Merged Tasks	Merging Method	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	TA	10.7850	3.9758	0.7377	-
	DARE	10.8753	3.9424	0.7437	-
	TIES	10.7638	3.3606	0.7475	-
	DARE + TIES	10.8560	3.6424	0.7445	-
	LATA (Ours)	10.2638	2.2909	0.8158	-
Math + Code	TA	12.1416	-	0.7248	0.5793
	DARE	12.4347	-	0.7111	0.5305
	TIES	12.3455	-	0.7165	0.5366
	DARE + TIES	12.5877	-	0.6914	0.5061
	LATA (Ours)	11.0208	-	0.8317	0.6280
UA + Code	TA	11.9674	3.8545	-	0.3171
	DARE	12.1404	3.8394	-	0.3537
	TIES	11.9742	3.3545	-	0.2866
	DARE + TIES	12.0392	3.5061	-	0.2866
	LATA (Ours)	11.2949	2.5667	-	0.3293
UA + Math + Code	TA	12.2611	3.6364	0.7172	0.5671
	DARE	12.5596	3.6939	0.7005	0.5061
	TIES	12.5602	3.5061	0.7104	0.5366
	DARE + TIES	12.8658	3.4788	0.6914	0.5122
	LATA (Ours)	11.0486	2.6394	0.8271	0.6280

Table 6: Results of task learning on Gemma-2-9b. Here, we merge models with $\lambda = 1.0$ for all tasks.

We also investigate whether LATA can enhance DARE and TIES. Table 4 shows that combining LATA with these methods often yields superior utility. However, LATA+DARE+TIES typically underperforms LATA+DARE or LATA+TIES alone, mirroring the observation that DATA+TIES is weaker than DARE or TIES. Moreover, in most cases, these three-method combinations in Table 4 are worse than LATA alone (Table 3), as TIES and DARE may zero out crucial layer vectors selected by LATA. Hence, using LATA by itself remains the best choice.

Table 7 presents results under Llama-3-8B and additional results with different hyperparameters for each baseline can be found in Appendix A.3. Owing to its smaller size, we adopt Logarithmic-Drop-by-Rank to account for higher interdependence among layers. LATA still sustains superior overall utility while achieving competitive or best scores in several tasks. This demonstrates LATA’s effectiveness across different architectures. In sum-

Merged Tasks	Merging Method	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	TA	9.0025	3.6303	0.8089	-
	DARE	9.0559	3.5606	0.8074	-
	TIES	9.1528	3.3788	0.8036	-
	DARE + TIES	9.0055	3.5515	0.7923	-
	LATA (Ours)	9.3160	3.7667	0.7847	-
Math + Code	TA	10.0648	-	0.6664	0.3415
	DARE	10.1674	-	0.6558	0.3415
	TIES	10.0103	-	0.6914	0.3293
	DARE + TIES	10.2170	-	0.6778	0.2927
	LATA (Ours)	9.9947	-	0.7491	0.2439
UA + Code	TA	10.4806	3.7879	-	0.2317
	DARE	10.4840	3.7273	-	0.2256
	TIES	10.2491	3.5667	-	0.2256
	DARE + TIES	10.5172	3.6606	-	0.1951
	LATA (Ours)	10.4579	3.5030	-	0.2500
UA + Math + Code	TA	9.9398	3.6333	0.6626	0.2987
	DARE	10.0415	3.8000	0.6732	0.3171
	TIES	9.9066	3.5030	0.6793	0.3171
	DARE + TIES	10.1180	3.6727	0.6634	0.3537
	LATA (Ours)	9.9057	3.7939	0.7316	0.2378

Table 7: Results of task learning on Llama-3-8b. Here, we merge models with $\lambda = 0.5$ for all tasks.

mary, LATA consistently shows clear advantages in merging multiple models.

Task Forgetting We set the scaling coefficient λ to 1.0 and use Drop-with-Threshold at the threshold σ of 0.95 (see more discussion in Section 6). Figure 4 shows that applying TA’s subtractive operation to reduce harmful content substantially improves alignment. LATA consistently outperforms existing methods, reducing GPT-4 harm scores below 2 for all tested languages, notably from 3.60 to 2.57 in German. Meanwhile, utility remains on par with the original model. These results suggest LATA precisely targets task vectors for removal and, in some cases (see Section 6), adjusting a minimal subset of parameters is sufficient to eliminate specific capabilities.

6 Discussion

Distribution of Important Layers for Target Tasks Figure 5 shows layer-wise similarity rankings between the three target tasks’ complex vectors for Gemma-2-9b and the instruction vector. The vertical axis is the similarity ranking, and the horizontal axis is the layer index. Layers with lower similarity (thus more impact on the target task) generally appear after layer 20, especially between layers 26 and 30. We hypothesize that earlier layers focus more on processing the input instructions, making them closer to the instruction vector and less crucial to the target task. Conversely, later layers generate outputs based on the earlier layers’ interpretations, causing parameter changes there to have greater impact on the target task and thus lower similarity with the instruction vector.

Another notable observation is the significant

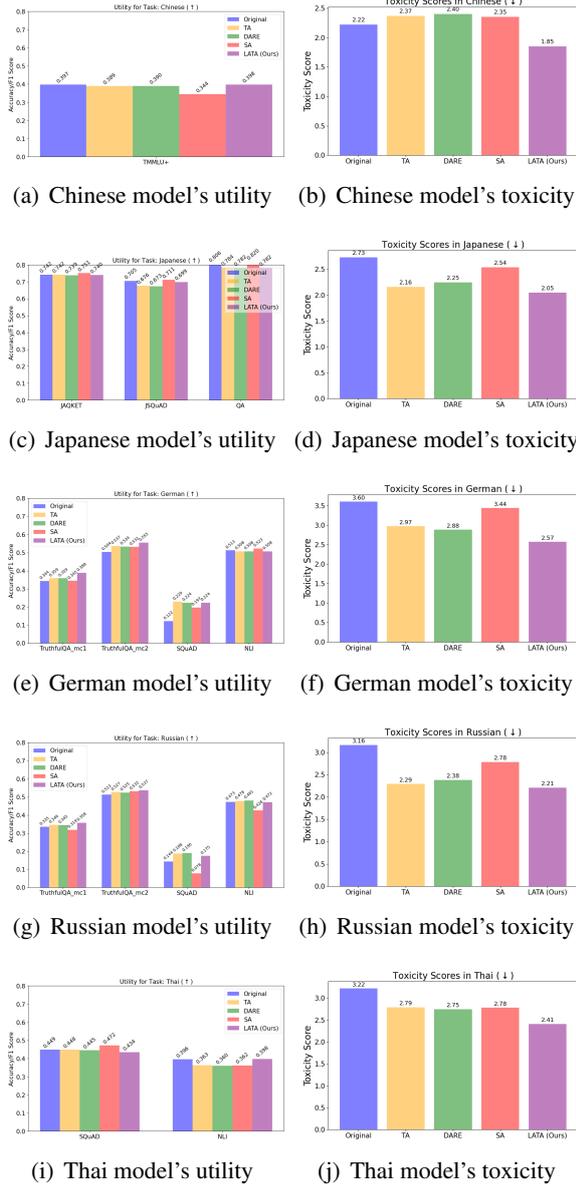


Figure 4: Result of task forgetting

494 overlap in similarity rankings for math and code
 495 tasks. We suspect a strong intrinsic similarity be-
 496 tween these two tasks, reflected in our experiments:
 497 when merging them simultaneously (*math + code*,
 498 *UA + math + code*), both tasks outperform their
 499 single-task scenarios (*UA + math*, *UA + code*), par-
 500 ticularly for code. This suggests that when task vec-
 501 tors share substantial similarity, merging them con-
 502 currently can further enhance the resulting model’s
 503 performance on each individual task.

504 **Significant Impact from a Small Fraction of**
 505 **Layer Vectors** In our task forgetting experiment,
 506 we set the threshold σ to 0.95, meaning that layer
 507 vectors with similarities above 0.95 were discarded.

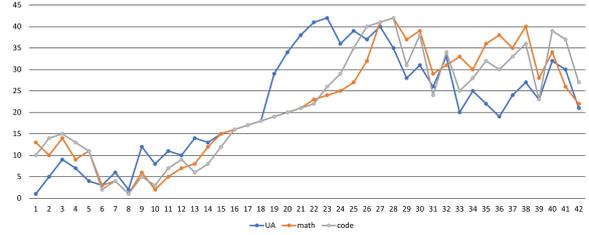


Figure 5: The graph illustrates the similarity rankings among layer vectors, with the x-axis representing the layer number and the y-axis indicating the similarity rank.

We arrived at this threshold because we observed extremely high similarity between the complex and instruction vectors for each layer, with only a handful of layer vectors showing similarity below 0.9. Even with the threshold fixed at 0.95, only about 10% of the layer vectors were retained as *pure vectors*, while the remaining 90% had similarities greater than 0.95. Under the DARE concept, discarding 90% of the vectors would ordinarily require rescaling the remaining 10% by a factor of $\frac{1}{1-0.9}$ (i.e., $10\times$). However, we merely applied $\lambda = 1.0$ to slightly increase these vectors, already achieving performance surpassing that of the original TA method. This finding indicates that a complete task vector indeed contains a subset of parameters that are highly critical to the target task, while a substantial portion is less significant. LATA successfully isolates these crucial and non-crucial segments from the task vector.

7 Conclusion

In this work, we introduced a novel approach (LATA) to TA, demonstrating its effectiveness in merging and fine-tuning LLMs across diverse tasks. LATA leverages dynamic task representations to achieve improved alignment and utility without compromising model performance. Through extensive experiments on benchmark datasets such as WikiText-2, GSM8K, and HumanEval, we showed that our approach consistently outperforms existing methods like DARE and TIES in balancing task-specific performance and generalization. Notably, our framework enables efficient model merging while mitigating interference between tasks, as evidenced by superior results in multi-task scenarios. Our findings highlight the potential of TA as a scalable and adaptable solution for optimizing LLMs in multi-task and cross-lingual settings.

545
546
547
548
549
550
551
552

553

554
555
556
557
558
559

560
561
562
563
564
565
566
567

568
569
570
571

572
573
574
575
576
577

578
579
580
581

582
583
584
585
586
587

588
589
590
591
592
593
594
595

596
597

Limitations

LATA relies on task arithmetic, so all models must share the same architecture (identical hidden dimensions and layer structures), which limits cross-family applications. Moreover, improper scaling coefficients of task vectors (λ) can lead to instability, potentially degrading model performance or causing catastrophic forgetting.

References

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. [On the cross-lingual transferability of monolingual representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics.

Rishabh Bhardwaj, Duc Anh Do, and Soujanya Poria. 2024. [Language models are Homer simpson! safety re-alignment of fine-tuned language models through task arithmetic](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14138–14149, Bangkok, Thailand. Association for Computational Linguistics.

Tian Bowen, Lai Songning, Wu Jiemin, Shuai Zhihao, Ge Shiming, and Yue Yutao. 2024. [Beyond task vectors: Selective task arithmetic based on importance metrics](#).

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. [Evaluating large language models trained on code](#). *arXiv preprint arXiv:2107.03374*.

Jiho Choi, Donggyun Kim, Chanhyuk Lee, and Seunghoon Hong. 2024. [Revisiting weight averaging for model merging](#). *arXiv preprint arXiv:2412.12153*.

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 problems](#). *ArXiv*, abs/2110.14168.

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. [XNLI: Evaluating cross-lingual sentence representations](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2475–2485, Brussels, Belgium. Association for Computational Linguistics.

Rui Dai, Sile Hu, Xu Shen, Yonggang Zhang, Xinmei Tian, and Jieping Ye. 2025. [Leveraging submodule](#)

linearity enhances task arithmetic performance in LLMs. In *The Thirteenth International Conference on Learning Representations (ICLR)*. 598
599
600

Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. 2020. [Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping](#). *arXiv preprint arXiv:2002.06305*. 601
602
603
604
605

Guodong Du, Junlin Lee, Jing Li, Runhua Jiang, Yifei Guo, Shuyang Yu, Hanting Liu, Sim Kuan Goh, Ho-Kin Tang, Daojing He, and Min Zhang. 2024. [Parameter competition balancing for model merging](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS)*. 606
607
608
609
610
611

Antonio Andrea Gargiulo, Donato Crisostomi, Maria Sofia Bucarelli, Simone Scardapane, Fabrizio Silvestri, and Emanuele Rodolà. 2025. [Task singular vectors: Reducing task interference in model merging](#). 612
613
614
615
616

Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vladimir Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. 2024. [Arcee’s MergeKit: A toolkit for merging large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 477–485, Miami, Florida, US. Association for Computational Linguistics. 617
618
619
620
621
622
623
624
625

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, 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 Al-lonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, 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, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-teng Jia, Kalyan Vasuden Alwala, Karthik Prasad, 626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657

658	Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth	Xu, Davide Testuggine, Delia David, Devi Parikh,	721
659	Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,	Diana Liskovich, Didem Foss, Dingkang Wang, Duc	722
660	Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	723
661	Lakhotia, Lauren Rantala-Yearly, Laurens van der	Elaine Montgomery, Eleonora Presani, Emily Hahn,	724
662	Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,	Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban	725
663	Louis Martin, Lovish Madaan, Lubo Malo, Lukas	Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	726
664	Blecher, Lukas Landzaat, Luke de Oliveira, Madeline	Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat	727
665	Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	728
666	Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew	Seide, Gabriela Medina Florez, Gabriella Schwarz,	729
667	Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-	Gada Badeer, Georgia Swee, Gil Halpern, Grant	730
668	badur, Mike Lewis, Min Si, Mitesh Kumar Singh,	Herman, Grigory Sizov, Guangyi, Zhang, Guna	731
669	Mona Hassan, Naman Goyal, Narjes Torabi, Niko-	Lakshminarayanan, Hakan Inan, Hamid Shojanazeri,	732
670	lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,	Han Zou, Hannah Wang, Hanwen Zha, Haroun	733
671	Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	734
672	Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic,	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	735
673	Peter Weng, Prajjwal Bhargava, Pratik Dubal,	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,	736
674	Praveen Krishnan, Punit Singh Koura, Puxin Xu,	Irina-Elena Veliche, Itai Gat, Jake Weissman, James	737
675	Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj	Geboski, James Kohli, Janice Lam, Japhet Asher,	738
676	Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer	739
677	Robert Stojnic, Roberta Raileanu, Rohan Maheswari,	Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy	740
678	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-	Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe	741
679	nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-	742
680	Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	743
681	hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal,	744
682	hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-	Katayoun Zand, Kathy Matosich, Kaushik	745
683	ran Narang, Sharath Rapparth, Sheng Shen, Shengye	Veeraraghavan, Kelly Michelena, Keqian Li, Kiran	746
684	Wan, Shruti Bhosale, Shun Zhang, Simon Vanden-	Jagadeesh, Kun Huang, Kunal Chawla, Kyle	747
685	hende, Soumya Batra, Spencer Whitman, Sten	Huang, Lailin Chen, Lakshya Garg, Lavender A,	748
686	Sootla, Stephane Collot, Suchin Gururangan, Sydney	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	749
687	Borodinsky, Tamar Herman, Tara Fowler, Tarek	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt,	750
688	Sheasha, Thomas Georgiou, Thomas Scialom, Tobias	Madian Khabsa, Manav Avalani, Manish Bhatt,	751
689	Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal	Martynas Mankus, Matan Hasson, Matthew Lennie,	752
690	Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh	Matthias Reso, Maxim Groshev, Maxim Naumov,	753
691	Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	754
692	Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic,	Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel,	755
693	Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney	Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	756
694	Meers, Xavier Martinet, Xiaodong Wang, Xiaofang	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	757
695	Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng	Mo Metanat, Mohammad Rastegari, Munish Bansal,	758
696	Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag,	Nandhini Santhanam, Natascha Parks, Natasha	759
697	Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	760
698	Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	761
699	Delpierre Coudert, Zheng Yan, Zhengxing	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	762
700	Chen, Zoe Papakipos, Aaditya Singh, Aayushi	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	763
701	Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro	764
702	Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,	Rittner, Philip Bontrager, Pierre Roux, Piotr	765
703	Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei	Dollar, Polina Zvyagina, Prashant Ratanchandani,	766
704	Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani,	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	767
705	Amos Teo, Anam Yunus, Andrei Lupu, Andres	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	768
706	Alvarado, Andrew Caples, Andrew Gu, Andrew	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	769
707	Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani,	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	770
708	Annie Dong, Annie Franco, Anuj Goyal, Aparajita	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,	771
709	Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	772
710	Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan,	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	773
711	Beau James, Ben Maurer, Benjamin Leonhardi,	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	774
712	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay,	775
713	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock,	Shaun Lindsay, Sheng Feng, Shenghao Lin,	776
714	Bram Wasti, Brandon Spence, Brani Stojkovic,	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	777
715	Brian Gamido, Britt Montalvo, Carl Parker, Carly	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	778
716	Burton, Catalina Mejia, Ce Liu, Changan Wang,	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	779
717	Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	780
718	Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer,	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	781
719	Cynthia Gao, Damon Civin, Dana Beaty, Daniel	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	782
720	Kreymer, Daniel Li, David Adkins, David	Subramanian, Sy Choudhury, Sydney Goldman, Tal	783

784	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	<i>The Eleventh International Conference on Learning</i>	843
785	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	<i>Representations</i> .	844
786	Matthews, Timothy Chou, Tzook Shaked, Varun		
787	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai		
788	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	Kentaro Kurihara, Daisuke Kawahara, and Tomohide	845
789	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	Shibata. 2022. JGLUE: Japanese general language	846
790	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	understanding evaluation . In <i>Proceedings of the Thir-</i>	847
791	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	<i>teenth Language Resources and Evaluation Confer-</i>	848
792	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo	<i>ence</i> , pages 2957–2966, Marseille, France. European	849
793	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,	Language Resources Association.	850
794	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,		
795	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	Kunfeng Lai, Zhenheng Tang, Xinglin Pan, Peijie Dong,	851
796	Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary	Xiang Liu, Haolan Chen, Li Shen, Bo Li, and Xi-	852
797	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	aowen Chu. 2025. Mediator: Memory-efficient llm	853
798	Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd	merging with less parameter conflicts and uncertainty	854
799	of models .	based routing .	855
800	Hasan Abed Al Kader Hammoud, Umberto Michieli,		
801	Fabio Pizzati, Philip Torr, Adel Bibi, Bernard	Viet Lai, Chien Nguyen, Nghia Ngo, Thuat Nguyen,	856
802	Ghanem, and Mete Ozay. 2024. Model merging and	Franck Dernoncourt, Ryan Rossi, and Thien Nguyen.	857
803	safety alignment: One bad model spoils the bunch .	2023. Okapi: Instruction-tuned large language mod-	858
804	In <i>Findings of the Association for Computational</i>	els in multiple languages with reinforcement learning	859
805	<i>Linguistics: EMNLP 2024</i> , pages 13033–13046, Mi-	from human feedback . In <i>Proceedings of the 2023</i>	860
806	ami, Florida, USA. Association for Computational	<i>Conference on Empirical Methods in Natural Lan-</i>	861
807	Linguistics.	<i>guage Processing: System Demonstrations</i> , pages	862
808	Rima Hazra, Sayan Layek, Somnath Banerjee, and Sou-	318–327, Singapore. Association for Computational	863
809	janya Poria. 2024. Safety arithmetic: A framework	Linguistics.	864
810	for test-time safety alignment of language models		
811	by steering parameters and activations . In <i>Proceed-</i>	Shen Li, Liuyi Yao, Lan Zhang, and Yaliang Li. 2025.	865
812	<i>ings of the 2024 Conference on Empirical Methods in</i>	Safety layers in aligned large language models: The	866
813	<i>Natural Language Processing</i> , pages 21759–21776,	key to LLM security . In <i>The Thirteenth International</i>	867
814	Miami, Florida, USA. Association for Computational	<i>Conference on Learning Representations</i> .	868
815	Linguistics.		
816	Masato Hirakawa, Shintaro Horie, Tomoaki Nakamura,	Zhenyi Lu, Chenghao Fan, Wei Wei, Xiaoye Qu, Dan-	869
817	Daisuke Oba, Sam Passaglia, and Akira Sasaki. 2024.	gyang Chen, and Yu Cheng. 2024. Twin-merging:	870
818	elyza/llama-3-elyza-jp-8b .	Dynamic integration of modular expertise in model	871
819	Chia-Yi Hsu, Yu-Lin Tsai, Chih-Hsun Lin, Pin-Yu Chen,	merging . In <i>Advances in Neural Information Pro-</i>	872
820	Chia-Mu Yu, and Chun-Ying Huang. 2024. Safe	<i>cessing Systems (NeurIPS)</i> .	873
821	loRA: The silver lining of reducing safety risks when		
822	finetuning large language models . In <i>The Thirty-</i>	Michael S. Matena and Colin A. Raffel. 2022. Merging	874
823	<i>eight Annual Conference on Neural Information</i>	models with fisher-weighted averaging . In <i>Advances</i>	875
824	<i>Processing Systems</i> .	in Neural Information Processing Systems (NeurIPS) .	876
825	Shih-Cheng Huang, Pin-Zu Li, Yu-chi Hsu, Kuang-		
826	Ming Chen, Yu Tung Lin, Shih-Kai Hsiao, Richard	Stephen Merity, Caiming Xiong, James Bradbury, and	877
827	Tsai, and Hung-yi Lee. 2024. Chat vector: A simple	Richard Socher. 2017. Pointer sentinel mixture mod-	878
828	approach to equip LLMs with instruction following	els . In <i>International Conference on Learning Repre-</i>	879
829	and model alignment in new languages . In <i>Proceed-</i>	<i>sations</i> .	880
830	<i>ings of the 62nd Annual Meeting of the Association</i>		
831	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	881
832	<i>pers)</i> , pages 10943–10959, Bangkok, Thailand. As-	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	882
833	sociation for Computational Linguistics.	man, Diogo Almeida, Janko Altschmidt, Sam Alt-	883
834	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Worts-	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	884
835	man, Suchin Gururangan, Ludwig Schmidt, Han-	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	885
836	naneh Hajishirzi, and Ali Farhadi. 2023. Editing	ing Bao, Mohammad Bavarian, Jeff Belgum, Ir-	886
837	models with task arithmetic . In <i>Proceedings of the</i>	wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,	887
838	<i>11th International Conference on Learning Representa-</i>	Christopher Berner, Lenny Bogdonoff, Oleg Boiko,	888
839	<i>tions (ICLR)</i> .	Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-	889
840	Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and	man, Tim Brooks, Miles Brundage, Kevin Button,	890
841	Pengxiang Cheng. 2023. Dataless knowledge fu-	Trevor Cai, Rosie Campbell, Andrew Cann, Brittany	891
842	sion by merging weights of language models . In	Carey, Chelsea Carlson, Rory Carmichael, Brooke	892
		Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully	893
		Chen, Ruby Chen, Jason Chen, Mark Chen, Ben	894
		Chess, Chester Cho, Casey Chu, Hyung Won Chung,	895
		Dave Cummings, Jeremiah Currier, Yunxing Dai,	896
		Cory Decareaux, Thomas Degry, Noah Deutsch,	897
		Damien Deville, Arka Dhar, David Dohan, Steve	898
		Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	899

900	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	963
901	Simón Posada Fishman, Juston Forte, Isabella Ful-	964
902	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	965
903	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	966
904	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	967
905	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	
906	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	
907	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	
908	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	
909	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	
910	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	
911	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	
912	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	
913	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	
914	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	
915	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	
916	Christina Kim, Yongjik Kim, Jan Hendrik Kirch-	
917	ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo,	
918	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	
919	stantinidis, Kyle Kopic, Gretchen Krueger, Vishal	
920	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	
921	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	
922	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	
923	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	
924	Anna Makanju, Kim Malfacini, Sam Manning, Todor	
925	Markov, Yaniv Markovski, Bianca Martin, Katie	
926	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	
927	McKinney, Christine McLeavey, Paul McMillan,	
928	Jake McNeil, David Medina, Aalok Mehta, Jacob	
929	Menick, Luke Metz, Andrey Mishchenko, Pamela	
930	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	
931	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	
932	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	
933	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	
934	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	
935	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	
936	tista Parascandolo, Joel Parish, Emy Parparita, Alex	
937	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	
938	man, Filipe de Avila Belbute Peres, Michael Petrov,	
939	Henrique Ponde de Oliveira Pinto, Michael, Poko-	
940	rny, Michelle Pocrass, Vitchyr H. Pong, Tolly Pow-	
941	ell, Alethea Power, Boris Power, Elizabeth Proehl,	
942	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	
943	Cameron Raymond, Francis Real, Kendra Rimbach,	
944	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	
945	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	
946	Girish Sastry, Heather Schmidt, David Schnurr, John	
947	Schulman, Daniel Selsam, Kyla Sheppard, Toki	
948	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	
949	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	
950	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	
951	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	
952	lipe Petroski Such, Natalie Summers, Ilya Sutskever,	
953	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	
954	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	
955	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	
956	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	
957	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	
958	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	
959	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	
960	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	
961	Clemens Winter, Samuel Wolrich, Hannah Wong,	
962	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	
	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	963
	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	964
	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	965
	Zheng, Juntang Zhuang, William Zhuk, and Barret	966
	Zoph. 2024. Gpt-4 technical report .	967
	Kunat Pipatanakul, Phatrasek Jirabovonvisut, Potsawee	968
	Manakul, Sittipong Sripaisarnmongkol, Ruangsak	969
	Patomwong, Pathomporn Chokchainant, and Kasima	970
	Tharnpipitchai. 2023. Typhoon: Thai large language	971
	models . <i>arXiv preprint arXiv:2312.13951</i> .	972
	Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi	973
	Jia, Prateek Mittal, and Peter Henderson. 2024. Fine-	974
	tuning aligned language models compromises safety,	975
	even when users do not intend to! In <i>The Twelfth In-</i>	976
	<i>ternational Conference on Learning Representations</i> .	977
	Morgane Riviere, Shreya Pathak, Pier Giuseppe	978
	Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard	979
	Hussenot, Thomas Mesnard, Bobak Shahriari,	980
	Alexandre Ramé, Johan Ferret, Peter Liu, Pouya	981
	Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos,	982
	Ravin Kumar, Charline Le Lan, Sammy Jerome, An-	983
	ton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan	984
	Girgin, Nikola Momchev, Matt Hoffman, Shantanu	985
	Thakoor, Jean-Bastien Grill, Behnam Neyshabur,	986
	Olivier Bachem, Alanna Walton, Aliaksei Severyn,	987
	Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin	988
	Abdagic, Amanda Carl, Amy Shen, Andy Brock,	989
	Andy Coenen, Anthony Laforge, Antonia Pater-	990
	son, Ben Bastian, Bilal Piot, Bo Wu, Brandon	991
	Royal, Charlie Chen, Chintu Kumar, Chris Perry,	992
	Chris Welty, Christopher A. Choquette-Choo, Danila	993
	Sinopalnikov, David Weinberger, Dimple Vijayku-	994
	mar, Dominika Rogozińska, Dustin Herbison, Elisa	995
	Bandy, Emma Wang, Eric Noland, Erica Moreira,	996
	Evan Senter, Evgenii Eltyshev, Francesco Visin,	997
	Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus	998
	Martins, Hadi Hashemi, Hanna Klimczak-Plucińska,	999
	Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda	1000
	Mein, Jack Zhou, James Svensson, Jeff Stanway,	1001
	Jetha Chan, Jin Peng Zhou, Joana Carrasqueira,	1002
	Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost	1003
	van Amersfoort, Josh Gordon, Josh Lipschultz, Josh	1004
	Newlan, Ju yeong Ji, Kareem Mohamed, Kartikeya	1005
	Badola, Kat Black, Katie Millican, Keelin McDonell,	1006
	Kelvin Nguyen, Kiranbir Sodhia, Kish Greene,	1007
	Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre,	1008
	Lena Heuermann, Leticia Lago, Lilly McNealus,	1009
	Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon,	1010
	Luciano Martins, Machel Reid, Manvinder Singh,	1011
	Matt Iverson, Martin Görner, Mat Velloso, Mateo	1012
	Wirth, Matt Davidow, Matt Miller, Matthew Rahtz,	1013
	Matthew Watson, Meg Risdal, Mehran Kazemi,	1014
	Michael Moynihan, Ming Zhang, Minsuk Kahng,	1015
	Minwoo Park, Mofi Rahman, Mohit Khatwani, Na-	1016
	talie Dao, Nenshad Bardoliwalla, Nesh Devanathan,	1017
	Neta Dumai, Nilay Chauhan, Oscar Wahltinez,	1018
	Pankil Botarda, Parker Barnes, Paul Barham, Paul	1019
	Michel, Pengchong Jin, Petko Georgiev, Phil Cull-	1020
	iton, Pradeep Kuppala, Ramona Comanescu, Ramona	1021
	Merhej, Reena Jana, Reza Ardeshir Rokni, Rishabh	1022
	Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc	1023

1024	Carthy, Sarah Cogan, Sarah Perrin, Sébastien M. R. Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. 2024. Gemma 2: Improving open language models at a practical size.		
1025			learning. In <i>The Twelfth International Conference on Learning Representations.</i>
1026			1082
1027			1083
1028		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>Proceedings of the 41st International Conference on Machine Learning (ICML).</i>	1084
1029			1085
1030			1086
1031			1087
1032			1088
1033		Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhui Chen. 2024. Mammoth2: Scaling instructions from the web. <i>arXiv preprint arXiv:2405.03548.</i>	1089
1034			1090
1035			1091
1036		Jinghan Zhang, Shiqi Chen, Junteng Liu, and Junxian He. 2023. Composing parameter-efficient modules with arithmetic operation. In <i>Thirty-seventh Conference on Neural Information Processing Systems.</i>	1092
1037			1093
1038			1094
1039			1095
1040		Wei Zhao, Zhe Li, Yige Li, Ye Zhang, and Jun Sun. 2024. Defending large language models against jailbreak attacks via layer-specific editing. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 5094–5109, Miami, Florida, USA. Association for Computational Linguistics.	1096
1041			1097
1042	Masatoshi Suzuki, Jun Suzuki, Koji Matsuda, Kyosuke Nishida, and Naoya Inoue. 2020. JAQKET: Construction of a japanese QA dataset on the subject of quizzes. In <i>Proceedings of the Annual Meeting of the Association for Natural Language Processing</i> , volume 26, pages 237–240.		1098
1043			1099
1044			1100
1045			1101
1046		Yuyan Zhou, Liang Song, Bingning Wang, and Weipeng Chen. 2024. MetaGPT: Merging large language models using model exclusive task arithmetic. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1711–1724, Miami, Florida, USA. Association for Computational Linguistics.	1102
1047			1103
1048	Zhi Rui Tam, Ya Ting Pai, Yen-Wei Lee, Hong-Han Shuai, Jun-Da Chen, Wei Min Chu, and Sega Cheng. 2024. TMMLU+: An improved traditional chinese evaluation suite for foundation models. In <i>First Conference on Language Modeling.</i>		1104
1049			1105
1050			1106
1051			1107
1052			1108
1053	Jiongxiao Wang, Jiazhao Li, Yiquan Li, Xiangyu Qi, Junjie Hu, Yixuan Li, Patrick McDaniel, Muhao Chen, Bo Li, and Chaowei Xiao. 2024. Backdooralign: Mitigating fine-tuning based jailbreak attack with backdoor enhanced safety alignment. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems.</i>		1109
1054			1110
1055			1111
1056			1112
1057			1113
1058			1114
1059			1115
1060	Mitchell Wortsman, Gabriel Ilharco, Samir Ya 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>Proceedings of the 39th International Conference on Machine Learning</i> , pages 23965–23998. PMLR.		1116
1061			1117
1062			1118
1063			1119
1064			1120
1065			1121
1066			1122
1067			1123
1068			1124
1069	Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. 2023. Ties-merging: Resolving interference when merging models. In <i>Advances in Neural Information Processing Systems (NeurIPS).</i>		1125
1070			1126
1071			1127
1072			1128
1073			1129
1074	Enneng Yang, Li Shen, Zhenyi Wang, Guibing Guo, Xiaojun Chen, Xingwei Wang, and Dacheng Tao. 2024a. Representation surgery for multi-task model merging. In <i>Forty-first International Conference on Machine Learning.</i>		1130
1075			1131
1076			1132
1077			1133
1078			1134
1079	Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guibing Guo, Xingwei Wang, and Dacheng Tao. 2024b. Adamerging: Adaptive model merging for multi-task		1135
1080			1136
1081			1137
			1138
			1139
			1140
			1141
			1142
			1143
			1144
			1145
			1146
			1147
			1148
			1149
			1150
			1151
			1152
			1153
			1154
			1155
			1156
			1157
			1158
			1159
			1160
			1161
			1162
			1163
			1164
			1165
			1166
			1167
			1168
			1169
			1170
			1171
			1172
			1173
			1174
			1175
			1176
			1177
			1178
			1179
			1180
			1181
			1182
			1183
			1184
			1185
			1186
			1187
			1188
			1189
			1190
			1191
			1192
			1193
			1194
			1195
			1196
			1197
			1198
			1199
			1200
			1201
			1202
			1203
			1204
			1205
			1206
			1207
			1208
			1209
			1210
			1211
			1212
			1213
			1214
			1215
			1216
			1217
			1218
			1219
			1220
			1221
			1222
			1223
			1224
			1225
			1226
			1227
			1228
			1229
			1230
			1231
			1232
			1233
			1234
			1235
			1236
			1237
			1238
			1239
			1240
			1241
			1242
			1243
			1244
			1245
			1246
			1247
			1248
			1249
			1250
			1251
			1252
			1253
			1254
			1255
			1256
			1257
			1258
			1259
			1260
			1261
			1262
			1263
			1264
			1265
			1266
			1267
			1268
			1269
			1270
			1271
			1272
			1273
			1274
			1275
			1276
			1277
			1278
			1279
			1280
			1281
			1282
			1283
			1284
			1285
			1286
			1287
			1288
			1289
			1290
			1291
			1292
			1293
			1294
			1295
			1296
			1297
			1298
			1299
			1300

Pre-Trained / Target Model: Meta-Llama-3-8B-Instruct⁸

Fine-Tuned Models:

UA: LLama-3-8b-Uncensored⁹

Math: MAMmoTH2-8B-Plus¹⁰

Code: code-millennials-8b¹¹

A.1.2 Task Forgetting

We show more details of models used for task forgetting.

Base Model: Meta-Llama-3-8B⁷

Pre-Trained Model: Meta-Llama-3-8B-Instruct⁸

Fine-Tuned Models: LLama-3-8b-Uncensored⁹

Target Models:

Traditional Chinese: Llama3-TAIDE-LX-8B-Chat-Alpha¹²

German: Llama3-DiscoLeo-Instruct-8B-v0.1¹³

Japanese: Llama-3-ELYZA-JP-8B¹⁴

Russian: saiga_llama3_8b¹⁵

Thai: llama-3-typhoon-v1.5-8b-instruct_8b¹⁶

A.2 Results with Different Hyperparameters on Gemma-2-9b

In this section, we show different hyperparameters of DARE, TIES, and DARE+TIES across different scaling coefficients on Gemma-2-9b. The results explain why the hyperparameters we used in the main text are the most effective for all baselines.

DARE. Table 8 follows the same settings as Table 5 while demonstrating the performance with varying drop rates. DARE achieves better results when the drop rate is 0.3.

On the other hand, we also consider different values of scaling coefficients. Following the settings of Table 6, in Table 9, we show the performance of DARE with different drop rates and the coefficient fixed at 1.0. Overall, compared with Table 6, we obtain the best result for DARE when the drop rate is

⁸<https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>

⁹<https://huggingface.co/DevsDoCode/LLama-3-8b-Uncensored>

¹⁰<https://huggingface.co/TIGER-Lab/MAMmoTH2-8B-Plus>

¹¹<https://huggingface.co/budecosystem/code-millennials-8b>

¹²<https://huggingface.co/taide/Llama3-TAIDE-LX-8B-Chat-Alpha>

¹³<https://huggingface.co/DiscoResearch/Llama3-DiscoLeo-Instruct-8B-v0.1>

¹⁴<https://huggingface.co/elyza/Llama-3-ELYZA-JP-8B>

¹⁵https://huggingface.co/IlyaGusev/saiga_llama3_8b

¹⁶<https://huggingface.co/scb10x/llama-3-typhoon-v1.5-8b-instruct>

Merged Tasks	Drop Rate	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	0.3	10.0146	1.5909	0.8324	-
	0.6	10.0979	1.6333	0.8241	-
	0.9	10.5492	1.6242	0.8112	-
Math + Code	0.3	10.4055	-	0.8294	0.6341
	0.6	10.4918	-	0.8393	0.6220
	0.9	11.4782	-	0.7703	0.5671
UA + Code	0.3	10.8740	1.7061	-	0.4390
	0.6	10.9795	1.6818	-	0.4634
	0.9	11.3437	1.8303	-	0.5366
UA + Math + Code	0.3	10.3804	1.6394	0.8294	0.6463
	0.6	10.4883	1.7091	0.8249	0.6280
	0.9	11.6337	1.8636	0.7453	0.5305

Table 8: Results of task learning with DARE under Gemma-2-9b. All scaling coefficients here are set as 0.5.

set to 0.3. This is why we choose these parameters in Table 5 and 6.

Merged Tasks	Drop Rate	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	0.3	10.8753	3.9424	0.7437	-
	0.6	11.2912	3.8515	0.7036	-
	0.9	15.3662	3.7636	0.4594	-
Math + Code	0.3	12.4347	-	0.7111	0.5305
	0.6	13.2541	-	0.6626	0.4329
	0.9	49.7530	-	0.0205	0.0183
UA + Code	0.3	12.1404	3.8394	-	0.3537
	0.6	12.3582	3.7697	-	0.3354
	0.9	16.0632	3.3121	-	0.3171
UA + Math + Code	0.3	12.5596	3.6939	0.7005	0.5061
	0.6	13.5538	3.6061	0.6262	0.3902
	0.9	61.5858	×	0.0091	0.0183

Table 9: Results of task learning with DARE under Gemma-2-9b. All scaling coefficients here are set as 1.0. The cross sign indicates that the model can only generate gibberish.

TIES. In Table 10, we follow the same settings with Table 5, but show more results for different k of TIES. TIES obtain better utilities across different combinations of task merging when $k = 0.7$.

Merged Tasks	Top k	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	0.5	10.0067	1.5394	0.8347	-
	0.7	10.0167	1.5636	0.8385	-
	0.9	10.0267	1.5788	0.8340	-
Math + Code	0.5	10.3486	-	0.8317	0.6707
	0.7	10.3763	-	0.8279	0.6524
	0.9	11.3946	-	0.8309	0.6524
UA + Code	0.5	10.8511	1.5455	-	0.4451
	0.7	10.8583	1.5121	-	0.4329
	0.9	10.8495	1.5455	-	0.4390
UA + Math + Code	0.5	10.3727	1.6515	0.8264	0.6585
	0.7	10.3994	1.7939	0.8317	0.6585
	0.9	10.4196	1.8636	0.8309	0.6463

Table 10: Results of task learning with TIES under Gemma-2-9b. All scaling coefficients here are set as 0.5.

Apart from the hyperparameter of TIES, we also take the scaling coefficient into account. Therefore, Table 11 uses the same settings as Table 6, with the only difference being the top k . In comparison with Table 6, the results are better when k is 0.7.

Therefore, the proper hyperparameters of TIES are setting k as 0.7.

Merged Tasks	Top k	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	0.5	10.6077	3.3455	0.7635	-
	0.7	10.7638	3.3606	0.7475	-
	0.9	10.8087	3.5394	0.7362	-
Math + Code	0.5	11.9588	-	0.7460	0.5732
	0.7	12.3455	-	0.7165	0.5366
	0.9	12.5847	-	0.6892	0.5427
UA + Code	0.5	11.8724	3.4182	-	0.3049
	0.7	11.9742	3.3545	-	0.2866
	0.9	11.9892	3.4455	-	0.2927
UA + Math + Code	0.5	12.1112	3.3848	0.7263	0.5732
	0.7	12.5602	3.5061	0.7104	0.5366
	0.9	12.8162	3.5727	0.6907	0.5244

Table 11: Results of task learning with TIES under Gemma-2-9b. All scaling coefficients here are set as 1.0.

DARE + TIES. Here, we show more different hyperparameter combinations of DARE+TIES with the scaling coefficient set to 0.5 in Table 12. Across different tasks, the results with $(p, k) = (0.1, 0.9)$ outperform other settings. These are the parameters we use in the main text as well.

A.3 Results with Different Hyperparameters on Llama-3-8B

In this section, we present various hyperparameters for DARE, TIES, and DARE+TIES on Llama-3-8B. The results demonstrate why the hyperparameters chosen in the main text are the most optimal across all baselines.

DARE. In the main text, we show the results of DARE when the drop rate is 0.3 and the scaling coefficient is 0.5 on the Llama-3-8B model. Table 13 presents additional results of DARE using different drop rate settings. However, Table 13 demonstrates that DARE can get the best result when the drop rate is set as 0.3.

TIES. Fixing the scaling coefficient at 0.5, we conduct more experiments of TIES on Llama-3-8B for different values of k , and results are shown in Table 14. Most of results with $k = 7$ surpass the other values of k .

DARE+TIES. We run more experiments of DARE+TIES on Llama-3-8B to show the impacts of different combinations of the drop rate and top k . Table 15 shows the results, with $(p, k) = (0.1, 0.9)$ achieving the best performance in most cases. This indicates that the parameters we use in the main text are the most favorable for this method.

Merged Tasks	Drop Rate p / Top k	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	0.7 / 0.3	10.1749	1.6000	0.8127	-
	0.4 / 0.6	10.0813	1.5545	0.8332	-
	0.1 / 0.9	10.0461	1.5818	0.8302	-
Math + Code	0.7 / 0.3	10.7264	-	0.8173	0.6341
	0.4 / 0.6	10.4415	-	0.8294	0.6524
	0.1 / 0.9	10.4208	-	0.8287	0.6341
UA + Code	0.7 / 0.3	10.9549	1.6091	-	0.4329
	0.4 / 0.6	10.8592	1.6030	-	0.4512
	0.1 / 0.9	10.8553	1.6121	-	0.4512
UA + Math + Code	0.7 / 0.3	10.7478	1.7333	0.8089	0.6280
	0.4 / 0.6	10.4725	1.7727	0.8279	0.6646
	0.1 / 0.9	10.4147	1.8091	0.8309	0.6524

Table 12: Results of task learning with DARE + TIES under Gemma-2-9b. All scaling coefficients here are set as 0.5.

Merged Tasks	Drop Rate	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	0.3	9.0559	3.5606	0.8074	-
	0.6	9.2150	3.6576	0.7900	-
	0.9	10.3136	3.3394	0.7195	-
Math + Code	0.3	10.1674	-	0.6558	0.3415
	0.6	10.5052	-	0.6626	0.2195
	0.9	14.0010	-	0.4814	0.1707
UA + Code	0.3	10.4840	3.7273	-	0.2256
	0.6	10.6566	3.6303	-	0.1463
	0.9	11.6871	3.7939	-	0.1646
UA + Math + Code	0.3	10.0415	3.8000	0.6732	0.3171
	0.6	10.2933	3.5061	0.6467	0.2866
	0.9	13.3988	3.9152	0.4723	0.1159

Table 13: Results of task learning with DARE under Llama-3-8B. All scaling coefficients here are set as 0.5.

Merged Tasks	Top k	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	0.5	9.1528	3.3788	0.8036	-
	0.7	9.0490	3.4909	0.8089	-
	0.9	9.0009	3.6636	0.7983	-
Math + Code	0.5	10.0103	-	0.6914	0.3293
	0.7	10.1618	-	0.6831	0.3354
	0.9	10.2093	-	0.6732	0.3232
UA + Code	0.5	10.2491	3.5667	-	0.2256
	0.7	10.4020	3.6818	-	0.1646
	0.9	10.5076	3.6727	-	0.1707
UA + Math + Code	0.5	9.9066	3.5030	0.6793	0.3171
	0.7	10.0566	3.5697	0.6694	0.2622
	0.9	10.1106	3.5818	0.6535	0.3171

Table 14: Results of task learning with TIES under Llama-3-8B. All scaling coefficients here are set as 0.5.

Merged Tasks	Drop Rate p / Top k	Utility WikiText-2(↓)	UA GPT-4(↑)	Math GSM8K(↑)	Code HumanEval(↑)
UA + Math	0.7 / 0.3	9.3173	3.7091	0.7983	-
	0.4 / 0.6	9.0607	3.5471	0.7945	-
	0.1 / 0.9	9.0055	3.5515	0.7923	-
Math + Code	0.7 / 0.3	10.8194	-	0.6391	0.2683
	0.4 / 0.6	10.3354	-	0.6520	0.3293
	0.1 / 0.9	10.2170	-	0.6778	0.2927
UA + Code	0.7 / 0.3	10.8128	3.6030	-	0.0976
	0.4 / 0.6	10.5554	3.7242	-	0.2256
	0.1 / 0.9	10.5172	3.6606	-	0.1951
UA + Math + Code	0.7 / 0.3	10.7319	3.8182	0.6224	0.3232
	0.4 / 0.6	10.2063	3.6303	0.6535	0.3293
	0.1 / 0.9	10.1180	3.6727	0.6634	0.3537

Table 15: Results of task learning with DARE + TIES under Llama-3-8B. All scaling coefficients here are set as 0.5.