

# AdaRank: Adaptive Rank Pruning for Enhanced Model Merging

Chanhyuk Lee, Jiho Choi, Chanryeol Lee, Donggyun Kim, Seunghoon Hong  
KAIST

{chan3684, jiho.choi, lcy9442, kdgyun425, seunghoon.hong}@kaist.ac.kr

## Abstract

*Model merging has emerged as a promising approach for unifying independently fine-tuned models into an integrated framework, significantly enhancing computational efficiency in multi-task learning. Recently, several SVD-based techniques have been introduced to exploit low-rank structures for enhanced merging, but their reliance on such manually designed rank selection often leads to cross-task interference and suboptimal performance. In this paper, we propose **AdaRank**, a novel model merging framework that adaptively selects the most beneficial singular directions of task vectors to merge multiple models. We empirically show that the dominant singular components of task vectors can cause critical interference with other tasks, and that naive truncation across tasks and layers degrades performance. In contrast, AdaRank dynamically prunes the singular components that cause interference and offers an optimal amount of information to each task vector by learning to prune ranks during test-time via entropy minimization. Our analysis demonstrates that such method mitigates detrimental overlaps among tasks, while empirical results show that AdaRank consistently achieves state-of-the-art performance with various backbones and number of tasks, reducing the performance gap between fine-tuned models to nearly 1%.*

## 1. Introduction

Recent advancements in machine learning have significantly enhanced the performance and diversity of pre-trained models, enabling the widespread availability of fine-tuned models tailored to specific tasks across various domains [12, 14]. These developments have made high-quality, task-specialized models increasingly accessible, often distributed through public repositories for easy deployment. However, utilizing each of these fine-tuned models independently remains computationally expensive and impractical, particularly as the number of tasks grows in real-world applications. Consequently, *model merging* has emerged as a promising solution to integrate individually

fine-tuned models into unified, efficient framework, facilitating scalable multi-task performance without the need for extensive retraining or resource-heavy infrastructure.

Among the pioneering techniques in model merging, Task Arithmetic (TA) [21] combines task vectors—defined as the difference between fine-tuned and pre-trained model weights—to integrate multiple models into a single framework, enhancing multi-task performance without requiring access to train data. Building on this baseline, several studies have proposed methods to modify these task vectors at the parameter level [19, 55, 60] to address cross-task interference, a key factor contributing to the performance gap between merged and fine-tuned models. More recently, researchers have adopted Singular Value Decomposition (SVD) to adjust task vectors in the spectral domain rather than through element-wise modifications [4, 16, 28, 32], reporting improved performance as SVD effectively preserves the structural integrity of task vector.

Despite these efforts, SVD-based methods do not fully bridge the performance gap with fine-tuned models, largely due to their heuristic selection of the low-rank subspace. As we demonstrate in this paper, this approach introduces several limitations, including the suboptimal selection of top singular components and the inflexible application of fixed ranks across diverse tasks and layers. These shortcomings constrain their effectiveness in multi-task model merging, but it is not their fundamental limitation; they could be mitigated through careful selection of singular components, which is our idea underlying this work.

To address these challenges, we propose **AdaRank** (**Adaptive Rank Pruning**), a method that dynamically identifies optimal singular components for each task vector using test-time adaptation [48, 54]. We introduce binary masks for each components to overcome the rigidity of the top- $k$  approximation strategy, effectively filtering out directions that exhibit significant interference across tasks. Drawing inspiration from AdaMerging [61], which successfully utilized entropy minimization to calibrate merged models, we employ the same objective to adapt our masks during test time, while eliminating the need for labeled training data.

We evaluate the effectiveness of AdaRank in merging models with various backbones and numbers of tasks, such as vision and language transformers, demonstrating that our method successfully reduces the performance gap with individual models. Especially, AdaRank reduces the performance gap between fine-tuned models to nearly 1%. Our method also seamlessly integrates with various model merging baselines, such as conventional Task Arithmetic [21] or CART [4], as well as prior test-time adaptation approaches like AdaMerging [61], independently enhancing overall performance while preserving the benefits of each method. These results highlight the effectiveness of AdaRank and underscore its potential as a versatile solution for model merging. We release our full code for reproducing the experiment results displayed in the paper.<sup>1</sup>

## 2. Related Work

**Model Merging.** Model merging seeks to integrate multiple independently trained models into a unified multi-task framework while preserving the performance of each individual component. Early approaches, such as Fisher Merging [33], which weights each model using the Fisher information matrix, and RegMean [22], which minimizes the  $L_2$  norm discrepancy with respect to each model’s parameters, were proposed as alternatives to simple weight averaging. More recently, Task Arithmetic (TA) [21] has emerged as a straightforward solution, defining *task vectors* as the difference between fine-tuned and pre-trained weights and merging them through simple arithmetic operations. However, TA suffers from cross-task interference, as task vectors are simply added together without addressing their potential interference, leading to suboptimal performance in multi-task scenarios.

**Task Arithmetic with Weight Sparsification.** To address cross-task interference, several methods [19, 55, 60, 63] proposed to sparsify task vectors by removing redundant components and preserving critical parameters in an element-wise manner. Notably, TIES-Merging [60] selects dominant parameters from task vectors based on their magnitude and constructs a sign vector reflecting the prevailing sign across models. Similarly, DARE [55] employs a Bernoulli distribution to randomly drop parameters and rescales the remaining ones to approximate the original task vector. Additionally, Consensus Merging [55] builds on prior methods by applying an additional fixed mask to task vectors, extracting task-specific information through the  $L_1$  norm difference between the merged model and individual task vectors. Despite the efforts, these methods based on element-wise sparsification still exhibit noticeable performance gap between fine-tuned models.

<sup>1</sup><https://github.com/david3684/AdaRank>

**Task Arithmetic in a Low-Rank Subspace.** Instead, recent works [4, 16, 28, 32] leverage Singular Value Decomposition (SVD) to address interference between task vectors. For example, CART [4] redefines task vectors as deviations from the average of fine-tuned weights rather than pre-trained weights, and showed applying low-rank approximation to these task vectors make merged model outperforms previous merging techniques without further modifications. Task Singular Vectors (TSV) [16] propose enforcing a low-rank structure on each task vector, followed by a whitening transformation to minimize interference among the truncated components. STAR [28] provides theoretical evidence that low-rank approximation of task vectors reduces the upper bound of interference in the merged model. These SVD-based model merging approaches have succeeded in narrowing the performance gap with fine-tuned models compared to earlier methods, but they still rely on heuristic low-rank approximation strategies, which remain a limitation.

**Model Merging with Test-Time Adaptation** Due to restrictions on accessing training data during model merging, several studies [32, 51, 61, 62] have employed Test-Time Adaptation (TTA) to modify task vectors or apply post-merging processes. This approach is strongly inspired by studies such as Test-Time Training (TTT) [48] and TENT [54], both of which seek a surrogate for the supervised loss during test-time. For instance, AdaMerging [61] leverages Shannon Entropy [44] as a surrogate for the supervised loss to refine task-specific merging coefficients during test-time. Building on this success, subsequent works like WeMOE [51] and Twin-Merging [32] have introduced a router module into the merged model to dynamically combine task-specific parameters during inference. These methods employ TTA with entropy minimization to train this module. Additionally, Representation Surgery [62] proposes a *post-merging* method that introduces an additional layer after the merged model to address representation bias, also trained with test data.

## 3. Preliminaries

**Task Arithmetic.** Given  $T$  heterogeneous tasks and model parameters  $\theta_i$  for  $i = 1, \dots, T$  fine-tuned from the same pre-trained backbone  $\theta_0$ , model merging aims to build a merged parameter  $\theta_m$  capable of performing all tasks. We consider the layer-wise Task Arithmetic (TA) [21] as a base approach that obtains the merged parameter by:

$$\theta_m^l = \theta_0^l + \lambda^l \sum_{i=1}^T \tau_i^l, \quad (1)$$

where  $l$  denotes the layer index,  $\theta^l \in \mathbb{R}^{d \times d'}$  is the parameter of  $l$ th layer, and  $\tau_i^l = \theta_i^l - \theta_0^l$  is a *task vector*<sup>2</sup> representing task-specific knowledge encoded in the difference between the fine-tuned and pre-trained parameters.

While TA has been effective in various model merging scenarios, it has also been widely observed that its merging performance is largely limited by cross-task interference in task vectors *i.e.*, adding a task vector degrades the performance of the other tasks. To address the issue, various attempts have been made to truncate the conflicting components of the task vector.

**Task Vector Truncation.** Early approaches propose to truncate task vectors by inspecting their *element-wise* contribution to the merged model. Inspired by model pruning, this is implemented by multiplying an element-wise binary mask to the task vector by  $\hat{\tau}_i = B_i \odot \tau_i$ , where such masks are carefully designed to reduce conflicts across task vectors [19, 60, 63]. However, such hand-designed, per-element pruning often breaks the inherent row-column correlations in the weight matrix, potentially destroying a low-dimensional structure of the fine-tuned parameters critical for individual tasks [11, 49, 52].

Recent studies have discovered that exploiting the low-rank structure of task vectors can be an effective alternative [4, 16, 28], often surpassing element-wise pruning. Under this framework, the model merging in Eq. (1) is modified by:

$$\theta_m^l = \theta_0^l + \lambda^l \sum_{i=1}^T \text{SVD}_k(\tau_i^l), \quad (2)$$

where  $\text{SVD}_k$  denotes the Singular Value Decomposition (SVD) with low-rank approximation on top- $k$  singular components. In contrast to element-wise pruning, truncating the singular components of task vector preserves the row-column correlations of the original matrix since the modified task vector remains confined to the same parameter subspace [52]. Meanwhile, restricting the rank of task vectors  $k$  to be low helps reduce the chances for interference across tasks.

Despite their success, the cross-task interference is still inherent in SVD-based methods. This is primarily because each task’s singular components are obtained via independent SVD per task, allowing correlations to persist in singular components across tasks. While the prior works introduced additional mechanisms to further ensure orthogonality of task vectors via whitening [16] or reorientation [4], it does not fully address the issue as their merging performance highly depends on the choice of rank  $k$ .

<sup>2</sup>In layer-wise approaches, each task vector  $\theta^l \in \mathbb{R}^{d \times d'}$  is actually a matrix but we follow the convention of calling it a vector [21].

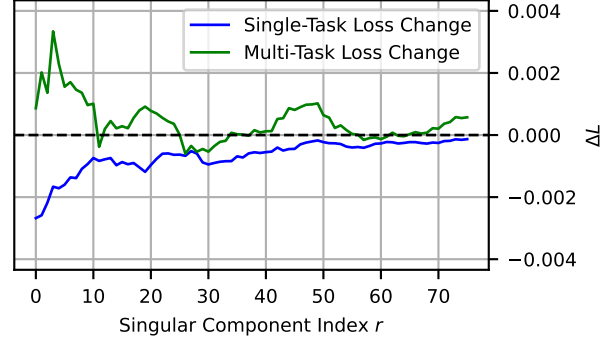


Figure 1. Impact of adding each singular component of a task vector to a model merged with full-rank task vectors (excluding the target task) on the net change in multi-task and single-task losses. Only the top 10% of indices are shown for clarity.

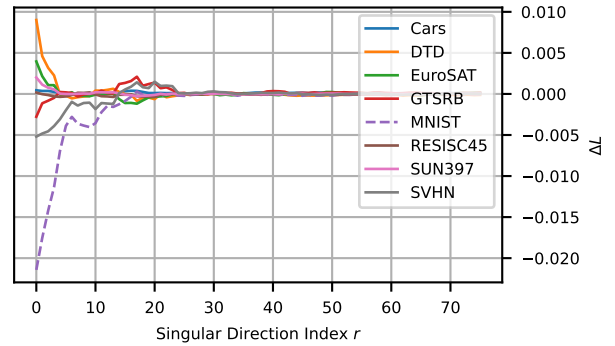


Figure 2. Impact on loss change of each task when adding singular components from the MNIST task vector. MNIST loss change is shown as a dotted line. Only the top 10% of indices are displayed for clarity.

## 4. Analysis on SVD of Task Vectors

Building on the low-rank paradigm, we seek an alternative direction for enhancing model merging. Specifically, we scrutinize the common practice of truncating each task vector to its top- $k$  singular components, posing two central questions: **(1)** *Is choosing the top singular components always beneficial for model merging?* **(2)** *Is enforcing a fixed rank across tasks and layers desirable for model merging?* In this section, we conduct an empirical analysis to investigate these questions, revealing the limitations of naive top- $k$  truncation.

**Limitations of Top Singular Components.** While the Eckart–Young theorem [15] guarantees that the top- $k$  singular components yield the best low-rank approximation in terms of reconstruction error for a single matrix, this optimality does not necessarily transfer to multi-task model merging. In a single-task scenario, these top components—characterized by large singular values—effectively

minimize the loss for that particular task. However, in a multi-task setting, the influence of top components extends beyond its own task and affects other tasks in unpredictable ways.

To measure the effect of adding singular components under other tasks’ interference, we set up the following experiment: for task  $i$ , we construct a merged model  $\theta_m$  using full-rank task vectors of the other tasks  $j \in \{1, \dots, T\} \setminus \{i\}$ . Thus,  $\theta_m$  is defined as

$$\theta_m = \theta_0 + \lambda \sum_{j \neq i} \tau_j. \quad (3)$$

Then, we measure how the multi-task loss changes by merging a  $r$ -th singular component of  $i$ -th task  $\mathbf{s}_{ir} = \sigma_{ir} \mathbf{u}_{ir} \mathbf{v}_{ir}^\top$  for each  $r$  one at a time by:

$$\Delta L(r) = \sum_{j=1}^T \Delta L_j(r) = \sum_{j=1}^T L_j(\theta_m + \lambda \mathbf{s}_{ir}) - L_j(\theta_m), \quad (4)$$

Figure 1 summarizes the results. It shows that top-singular components tend to contribute more in reducing the loss for its own task  $\Delta L_i$  (blue curve), while also introducing significant interference to the other tasks, often resulting in a net increase in multi-task loss  $\Delta L$  (green curve). It indicates that naively merging top singular components can often be suboptimal in multi-task model merging, since the interference with other tasks may outweigh the performance gains from adding the top singular component of a single task.

For deeper understanding of this phenomenon, we analyze the loss changes of individual tasks when the excluded task  $i$  is MNIST [27] (Figure 2). We observe that adding a singular component from the MNIST benefits semantically aligned tasks (e.g., SVHN [34], digit classification task), whereas dissimilar tasks (e.g., DTD [5], a texture classification task) experience increases in a net loss. These varying interactions lead to an unpredictable multi-task loss gain for each singular component; notably, top components with large singular values  $\sigma_{ir}$  can have a more pronounced negative impact on dissimilar tasks. Overall, these results suggest that top singular components are not always optimal for multi-task model merging.

**Limitations of Fixed Rank Truncation.** Another significant limitation lies in the use of a fixed top- $k$  truncation across diverse tasks and model layers. Figure 3 illustrates the effective rank [43] of the task vectors, defined as the number of singular components that preserve 95% of the total energy measured by the sum of the squared singular values. As shown, the effective rank varies considerably across task vectors from different models. Li et al. [29] established that the intrinsic dimension of neural network

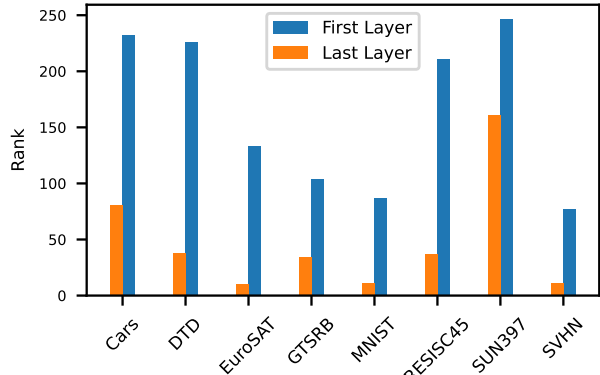


Figure 3. Effective rank [43] capturing 95% of total energy in the first and last MLP layer of ViT-B/32 task vectors obtained from 8 different fine-tuned weights.

representations increases with task complexity. Our experiment supports that this trend is also observed in task vectors, showing a strong correlation between effective rank and task demands. For instance, task vectors from models fine-tuned with SUN397 [59] (397 classes) require a higher effective rank compared to those from simpler tasks like MNIST [27] or SVHN [34]. Moreover, task vectors exhibit pronounced rank variation across layers (see Figure 3). Early layers, which capture task-agnostic features [26, 36, 41], show higher ranks with lower variance, reflecting shared information across tasks. In contrast, later layers, encoding task-specific representations, demonstrate greater rank variability and lower overall ranks. This divergence of rank among tasks and layers highlights the challenge of using a fixed top- $k$  truncation for model merging, since we may either discard important components that are critical for some tasks or keep unnecessary components that cause interference between tasks.

**Summary.** Summarizing these observations: **(1)** Some of top singular components benefit their own tasks but degrade performance in other tasks, making naive top- $k$  selection suboptimal on multi-task performance. **(2)** Task vectors exhibit diverse rank requirements across tasks and layers, so a fixed rank truncation fails to accommodate these differences. These findings highlight the necessity of an *adaptive* selection strategy that can *selectively preserve* critical singular components for each task and layer while mitigating negative interference.

## 5. Adaptive Rank Pruning

In this section, we introduce **AdaRank**, **Adaptive Rank Pruning**, a test-time adaptation [48, 54] method to find the optimal singular components that minimizes the cross-task interference in model merging, while preserving each task’s performance.

To selectively preserve or prune singular components, we introduce a binary mask that indicates the selection of each component. For each layer  $l$  and task  $i$ , we define a binary mask vector  $B_i^l \in \{0, 1\}^{1 \times d}$ , where each element indicates whether the corresponding singular component is preserved (1) or pruned (0). The set of binary masks for all tasks in layer  $l$  is denoted as  $B^l = \{B_1^l, B_2^l, \dots, B_T^l\}$ . For brevity, we denote the collection of all such masks across layers and tasks simply as  $B$ . Under a given state of  $B$ , the layer-wise merged model is expressed as:

$$\theta^l(B^l) = \theta_0^l + \lambda^l \sum_{i=1}^T U_i^l (\text{diag}(B_i^l) \odot \Sigma_i^l) V_i^{l\top}, \quad (5)$$

where  $U_i^l$  and  $V_i^l$  are the left and right singular vector matrices,  $\Sigma_i^l$  is the diagonal matrix of singular values, and  $\odot$  denotes the element-wise product.

Note that when  $B_{ir} = 1$  if  $r \leq k$  and  $B_{ir} = 0$  if  $r > k$  for all  $i$ , Eq. (5) reduces to the top- $k$  selection strategy (Eq. (2)). On the other hand, if all elements of  $B$  are set to one, it reduces to standard Task Arithmetic composed of full-rank task vectors (Eq. (1)). Allowing  $B$  to be an arbitrary binary matrix, we can effectively address issues of naive top- $k$  truncation discussed in Section 4: we can prune adversarial top singular components that increase the net multi-task loss, and allow variable ranks in task vectors across tasks and layers.

However, finding the optimal mask sets  $B$  is not straightforward since we cannot access training data or directly compute the loss gradient during model merging. Instead, we adopt test-time adaptation with Shannon Entropy minimization [44] as a surrogate objective for optimizing  $B$  without access to training data. Entropy minimization has been widely employed as a surrogate objective in model merging and compression [32, 51, 61], and it has generally proven effective due to its strong correlation with supervised multi-task loss [61].

Finally, our learning objective is minimizing the sum of output entropies  $H_i$  for each task:

$$\underset{B}{\operatorname{argmin}} \sum_{i=1}^T \sum_{x_i \in \mathcal{D}_i} H_i(f(\theta(B), x_i)), \quad (6)$$

where  $\mathcal{D}_i$  is a small *unlabeled test data* for task  $i$ , and  $f(\theta(B), x_i)$  is the model output parameterized by  $\theta(B)$  for input  $x_i \in \mathcal{D}_i$ .

We optimize the binary values of  $B$  using the Straight-Through Estimator (STE) [1] with a sigmoid function, treating  $B_i$  as a continuous parameter during the backward pass, consistent with prior works [9, 20, 31]. In the forward pass,  $B_i$  is rounded to  $\{0, 1\}$  to serve as a binary mask, while in the backward pass, it remains continuous to propagate gradients. Once the optimal set  $B$  is determined, the merged

model can be deployed by applying  $B^l$  to Eq. 5 without incurring additional storage overhead.

## 6. Experiments

### 6.1. Experimental Setup

**Baselines.** We compare AdaRank with prior approaches in model merging, categorized into two groups: *static* merging methods, which do not include an adaptation process (e.g., Task Arithmetic [21] and TIES-Merging [60]), and *adaptive* merging methods, which incorporate a test-time adaptation process (e.g., AdaMerging [61]). While approaches based on model compression are not directly comparable since they can access to individual task vectors and ground-truth task indices in test time, we include comparisons to these methods in the Appendix for completeness.

**Implementation Details.** To demonstrate the generality of our method, we integrate AdaRank to two SVD-based merging baselines: Task Arithmetic with SVD [16] and CART [4], which correspond to Eq. (2) with the difference in setting  $\theta_0$  as a pre-trained model or weight averaging, respectively. For optimization, we initialize the binary matrix  $B$  such that Eq. (5) is initialized to Eq. (1) for Task Arithmetic and to Eq. (2) for CART. Since the layer-wise coefficient  $\lambda^l$  can be learned jointly with the binary mask  $B$  via test-time adaptation as in AdaMerging [61], we learn both of them jointly in our default setting. The effect of  $\lambda$  and  $B$  is analyzed in Section 6.4.

### 6.2. Merging Vision Models

Following standard experiment protocol [55], we evaluate our method by merging two Vision Transformer [14] backbones of CLIP [40]: ViT-B/32 and ViT-L/14, which are fine-tuned on 8, 14, and 20 tasks. For the 8-task benchmark, we use the following datasets: Cars [24], DTD [5], EuroSAT [18], SVHN [34], GTSRB [47], MNIST [27], SUN397 [59], and RESISC45 [3]. The 14-task benchmark extends this set with six additional datasets: CIFAR100 [25], STL10 [7], Flowers102 [35], OxfordIIIT-Pet [37], PCAM [53], and FER2013 [17]. For the 20-task benchmark, we further include six more datasets: EMNIST [8], CIFAR10 [25], Food101 [2], FashionMNIST [58], RenderedSST2 [46], and KMNIST [6].

**Results** The average accuracy of merged model for 8, 14, 20 tasks and ViT-B/32, ViT-L/14 are presented in Tab. 1. For per-task breakdown and additional baselines, see Appendix B. Our method consistently achieves state-of-the-art results across all backbones and task counts, reducing the gap between individual models by nearly 1% on the 8-task benchmark. In particular, when integrated with Task Arithmetic, our method yields average gains of 18.6% for ViT-

Table 1. Average multi-task performance on 8, 14, 20 vision tasks with merged ViT-B/32 and ViT-L/14.

Method	ViT-B/32			ViT-L/14		
	8 tasks	14 tasks	20 tasks	8 tasks	14 tasks	20 tasks
Pretrained	48.0	59.6	56.0	65.0	68.4	65.4
Individual	90.5	90.3	90.5	94.2	93.3	94.0
Static Merging Methods						
Weight Averaging	65.9	64.3	60.9	79.6	76.8	71.7
Task Arithmetic (TA) [21]	69.2	65.4	61.0	84.5	79.6	74.2
Ties-Merging [60]	72.4	65.2	62.9	86.1	79.5	75.8
Consensus-Ties [55]	74.8	68.2	62.9	87.2	81.5	78.8
Consensus-TA [55]	75.2	70.0	65.0	86.6	81.9	77.6
TSV-M [16]	83.8	79.5	76.7	91.2	88.3	87.3
CART [4]	84.7	79.5	76.8	92.6	88.0	87.9
Adaptive Merging Methods						
AdaMerging [61]	80.1	76.7	69.2	90.8	88.0	86.8
TA+AdaRank	<b>87.9</b>	<b>82.1</b>	<b>81.4</b>	<b>92.9</b>	<b>89.4</b>	<b>89.1</b>
CART+AdaMerging [4]	85.9	82.3	82.7	93.1	90.4	91.3
CART+AdaRank	<b>89.2</b>	<b>86.2</b>	<b>86.4</b>	<b>93.4</b>	<b>91.4</b>	<b>91.8</b>

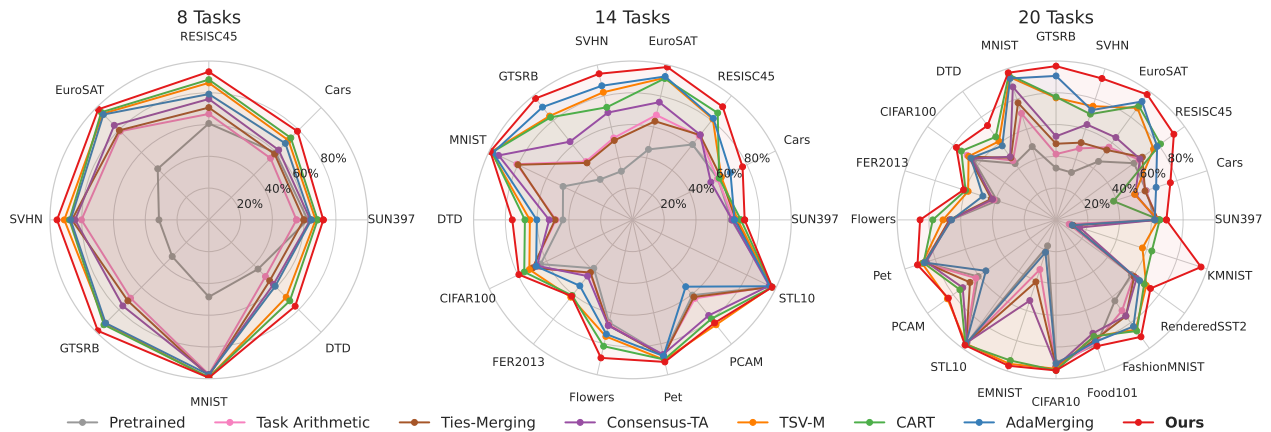


Figure 4. Individual task performance of merged ViT-B/32 model with 8, 14, 20 tasks, respectively. Detailed performances and results for ViT-L/14 is reported in the Appendix.

B/32 and 11.0% for ViT-L/14, surpassing AdaMerging [61] by 8.47% and 1.93%, respectively. Similarly, we observe consistent gain by applying our method to different SVD-based merging framework (CART vs. CART-AdaRank), showing the general advantage of choosing singular components over naive top- $k$  truncation.

Furthermore, compared to AdaMerging [61], which learns only the task-wise coefficient  $\lambda$  using full-rank task vectors, our method achieves significantly higher performance by merging models with task vectors that are optimally low-rank approximated through our proposed approach. This indicates that our masking strategy independently enhances performance, irrespective of whether  $\lambda$  is learned, a point we further explore in Section 6.4. More-

over, the observed improvement validates our approach discussed in section 4, concerning the degradation of multi-task performance due to interfering singular components.

### 6.3. Merging Language Models

We further evaluate AdaRank on RoBERTa-base [30] and GPT-2 [39] models, each fine-tuned on seven GLUE tasks. Following the setting from Fusionbench [50], we merge 7 weights finetuned for Text classification tasks respectively CoLA [56], SST-2 [45], MRPC [13], QQP [38], MNLI [57], QNLI [42], and RTE [10]. We report the Matthews correlation coefficient for CoLA task and accuracy for others, in both results.

Table 2. Multi-task performance on 7 NLP tasks with merged RoBERTa model.

Method	CoLA	SST2	MRPC	QQP	MNLI	QNLI	RTE	Average
Individual	0.6018	0.9404	0.8922	0.9141	0.872	0.9271	0.7906	0.8483
Weight Averaging	0.1808	0.8188	0.7794	0.7960	0.4383	0.7106	0.6173	0.6202
Task Arithmetic (TA) [21]	0.2330	0.8658	0.7868	0.8395	0.6371	0.7304	0.6101	0.6718
Ties-Merging [60]	0.2499	0.8349	0.7868	0.8515	0.6072	0.7580	0.4224	0.6444
CART [4]	0.3092	0.9197	0.8088	0.7953	0.5767	0.7772	0.7112	0.6997
AdaMerging [61]	-0.0359	0.9266	0.7721	0.8221	0.7880	0.7961	0.6643	0.6762
TA+AdaRank	0.1401	0.9151	0.777	0.7963	0.7814	0.8409	0.6715	<b>0.7032</b>
CART+AdaRank	0.3638	0.9278	0.7721	0.7956	0.7752	0.8753	0.6823	<b>0.7417</b>

Table 3. Multi-task performance on 7 NLP tasks with merged GPT-2 model.

Method	CoLA	SST2	MRPC	QQP	MNLI	QNLI	RTE	Average
Individual	0.4077	0.9118	0.8039	0.3964	0.8200	0.8827	0.6534	0.7680
Weight Averaging	0.1214	0.5252	0.5098	0.7670	0.5925	0.5761	0.4477	0.5057
Task Arithmetic (TA) [21]	-0.0019	0.8360	0.6961	0.8182	0.7188	0.7049	0.4729	0.6064
Ties-Merging [60]	0.0328	0.8177	0.6838	0.8284	0.7433	0.6957	0.4765	0.6112
CART [4]	0.1143	0.8624	0.5466	0.8177	0.7010	0.7620	0.5235	0.6182
AdaMerging [61]	0.0587	0.7982	0.7083	0.8104	0.6845	0.6758	0.4621	0.5997
TA+AdaRank	0.0617	0.8819	0.6275	0.7935	0.7539	0.8131	0.4982	<b>0.6328</b>
CART+AdaRank	0.1152	0.8853	0.6471	0.7982	0.7446	0.8473	0.5018	<b>0.6485</b>

**Results.** Results for RoBERTa-base and GPT-2 are presented in Table 2 and Table 3, respectively. As shown in table, applying AdaRank in language model merging achieves a consistent gain in both Task Arithmetic and CART baseline compared to the initial performance. These results indicate that our method demonstrates robust effectiveness not only on vision models of various sizes but also across different modalities and backbone architectures, thereby highlighting its scalability.

#### 6.4. Analysis

In this section, we further evaluate the effectiveness of our method. First, we visualize the learned binary masks from the best performing model to investigate the pruning patterns inherent in our approach. Next, we conduct an ablation study to validate our binary mask selection process is critical for enhance in multi-task performance.

**Mask Visualization.** Figure 5 presents the layer-wise masks obtained by our method, illustrating how it directly addresses the two primary limitations of naive rank truncation discussed in Section 4. First, masks from diverse layers and tasks reveal that a significant fraction of the top singular components are pruned, indicating that AdaRank actively discards those components when they are detrimental

to multi-task performance, which aligns with the cross-task interference we discussed earlier. Second, the masks clearly demonstrate that task vectors from different layers and tasks demand different ranks. Comparing masks from the first attention block (top row) with those from the last block (bottom row), we could observe a notable difference in the masking patterns. While early layers tend to preserve a broader range of indices and exhibit more uniform pruning patterns across tasks, deeper layers display greater variability in the preserved indices, reflecting task-specific needs. Notably, a similar pattern appears when our method is applied in both Task Arithmetic and CART (left and right column in Figure 5, respectively), with CART producing sparser results due to its low-rank initialization for the mask. Collectively, these findings confirm that AdaRank not only pinpoints and removes the interfering components left intact by naive top- $k$  approaches but also dynamically adapts the rank per layer and task to preserve essential information.

**Ablation Study.** Since we optimized task vector coefficient  $\lambda$  and binary mask  $B$  simultaneously, we conduct further ablation study to isolate the contributions of each component. Table 4 presents the impact on performance of applying test-time adaptation on  $\lambda$  and  $B$  independently for each starting method, Task Arithmetic (TA) and CART. As

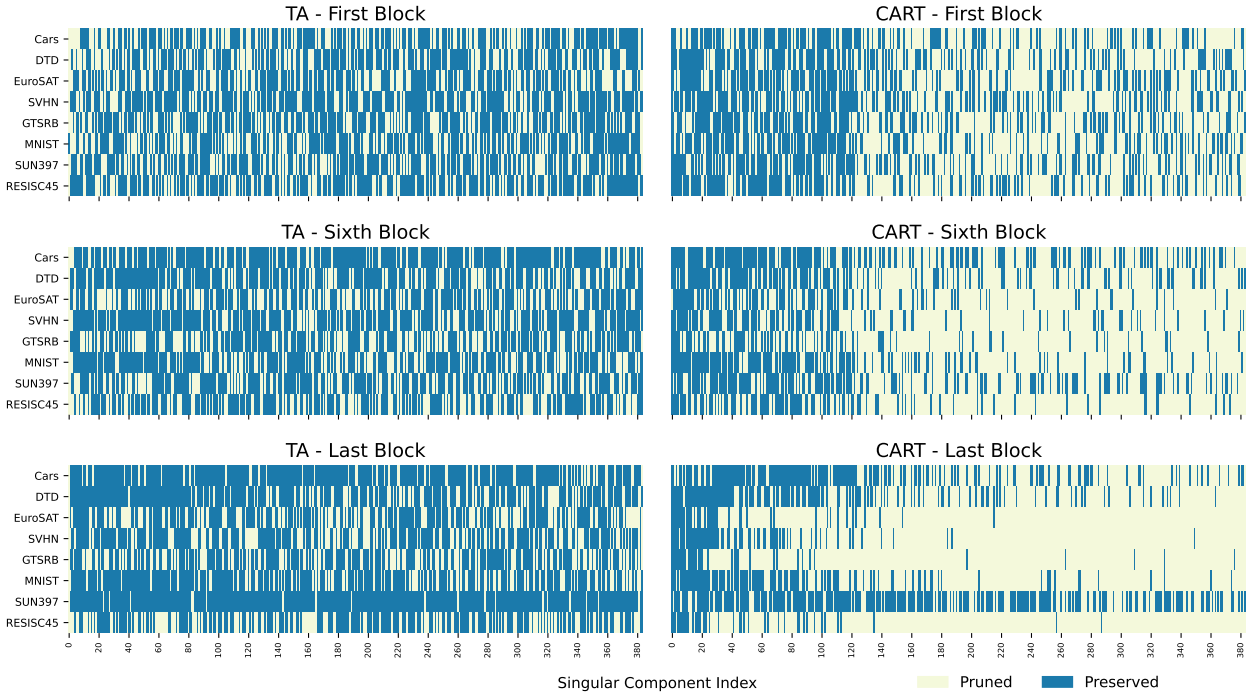


Figure 5. Binary masks derived from AdaRank for merging 8 fine-tuned ViT models. Left column shows AdaRank applied on Task Arithmetic (TA), and right column shows AdaRank applied on CART. Layers from the top, middle, and bottom blocks are plotted, showing only the top 50% of indices for clarity. Blue indicates preserved singular components, while yellow denotes pruned indices. Each rows of plotted masks correspond to individual tasks involved in merging.

Table 4. Ablation study results comparing the contributions of coefficient tuning and binary mask adaptation applied to TA (Task Arithmetic) and CART. Accuracy is reported as the average multi-task performance in the 8-task merging benchmark for ViT-B/32.

Learning Components		Baselines	
Coefficient $\lambda$	Binary Mask $B$	TA	CART
×	×	69.1	84.7
✓	×	80.1	85.9
×	✓	79.9	88.7
✓	✓	<b>87.9</b>	<b>89.2</b>

reported in AdaMerging [61], learning  $\lambda$  improves model merging performance in both both baselines. Regardless of that, optimizing  $B$  also significantly enhances the performance in both cases—yielding improvements comparable to learning  $\lambda$  in TA and even outperforming in CART. This results indicates that our method can independently contribute to substantial improvements in model merging performance, even without learning  $\lambda$ . Notably, in the CART setting, optimizing  $B$  alone achieves results within 2% of the best model. This observation aligns with CART’s finding that applying low-rank approximation to task vectors

derived from averaged weights is more effective than application in TA. Nonetheless, since we observed that combining both approaches orthogonally enhances multi-task performance, we employ both methods for our best-performing model.

## 7. Conclusion

In this paper, we present AdaRank, a dynamic extension of SVD-based rank truncation for multi-task model merging, designed to reduce cross-task interference from suboptimal low-rank approximations. Our analysis shows that dominant singular components often cause interference, while task vectors need varying ranks across layers and tasks. AdaRank tackles this by using test-time adaptation to learn a binary mask over singular components, guided by entropy minimization without labeled data. By applying AdaRank to various model merging baseline methods, we obtained significant performance gains on both vision and language benchmarks compared to prior methods. Furthermore, we showed that mask patterns obtained from our method well-align with our analysis on the limitation of prior works. Moving forward, we anticipate exploring more complex scenarios, such as online data streams, larger-scale multi-task scenarios, or multi-modal models. These future works would further validate robustness and scalability of our method in real-world applications.



## References

- [1] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. 5, 12
- [2] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 446–461, 2014. 5
- [3] Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017. 5
- [4] Jiho Choi, Donggyun Kim, Chanhyuk Lee, and Seunghoon Hong. Revisiting weight averaging for model merging. *arXiv preprint arXiv:2412.12153*, 2024. 1, 2, 3, 5, 6, 7, 12, 13, 14, 15, 16
- [5] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3606–3613, 2014. 4, 5
- [6] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha. Deep learning for classical japanese literature. *arXiv preprint arXiv:1812.01718*, 2018. 5
- [7] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 215–223, 2011. 5
- [8] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. EMNIST: An extension of MNIST to handwritten letters. In *Proceedings of the International Joint Conference on Neural Networks (IJCNN)*, pages 2921–2926, 2017. 5
- [9] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. *Advances in Neural Information Processing Systems*, 28, 2015. 5
- [10] Ido Dagan, Oren Glickman, and Bernardo Magnini. The PASCAL recognising textual entailment challenge. In *Machine Learning Challenges: Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Textual Entailment*, pages 177–190. Springer, 2006. 6
- [11] Emily L Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. *Advances in Neural Information Processing Systems*, 27, 2014. 3
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019. 1
- [13] William B. Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005. 6
- [14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. 1, 5
- [15] Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936. 3
- [16] Antonio Andrea Gargiulo, Donato Crisostomi, Maria Sofia Bucarelli, Simone Scardapane, Fabrizio Silvestri, and Emanuele Rodolà. Task singular vectors: Reducing task interference in model merging. *arXiv preprint arXiv:2412.00081*, 2024. 1, 2, 3, 5, 6, 12, 13, 14, 15, 16
- [17] Ian J. Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, Yingbo Zhou, Chetan Ramaiah, Fangxiang Feng, Ruifan Li, Fawaz Athar, Raghuraman Sabesan, Sivasankaran Rajaraman, Zhenhua Li, et al. Challenges in representation learning: A report on three machine learning contests. In *Proceedings of the International Conference on Neural Information Processing (ICANN)*, pages 117–124, 2013. 5
- [18] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019. 5
- [19] Chenyu Huang, Peng Ye, Tao Chen, Tong He, Xiangyu Yue, and Wanli Ouyang. EMR-merging: Tuning-free high-performance model merging. In *Advances in Neural Information Processing Systems*, 2024. 1, 2, 3, 12
- [20] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. *Advances in Neural Information Processing Systems*, 29, 2016. 5
- [21] Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In *International Conference on Learning Representations*, 2023. 1, 2, 3, 5, 6, 7, 12, 13, 14, 15, 16
- [22] Xisen Jin, Xiang Ren, Daniel Preotiu-Pietro, and Pengxiang Cheng. Dataless knowledge fusion by merging weights of language models. In *International Conference on Learning Representations*, 2023. 2, 13, 15
- [23] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. 12
- [24] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 554–561, 2013. 5
- [25] Alex Krizhevsky. Learning multiple layers of features from

- tiny images. Technical report, University of Toronto, 2009. 5
- [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems*, 25, 2012. 4
- [27] Yann LeCun. The MNIST database of handwritten digits. 1998. 4, 5
- [28] Yu-Ang Lee, Ching-Yun Ko, Tejaswini Pedapati, I Chung, Mi-Yen Yeh, Pin-Yu Chen, et al. STAR: Spectral truncation and rescale for model merging. *arXiv preprint arXiv:2502.10339*, 2025. 1, 2, 3
- [29] Chunyuan Li, Heerad Farkhor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic dimension of objective landscapes. In *International Conference on Learning Representations*, 2018. 4
- [30] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. 6
- [31] Christos Louizos, Max Welling, and Diederik P. Kingma. Learning sparse neural networks through  $L_0$  regularization. In *International Conference on Learning Representations*, 2018. 5
- [32] Zhenyi Lu, Chenghao Fan, Wei Wei, Xiaoye Qu, Danyang Chen, and Yu Cheng. Twin-merging: Dynamic integration of modular expertise in model merging. In *Advances in Neural Information Processing Systems*, 2024. 1, 2, 5
- [33] Michael S Matena and Colin A Raffel. Merging models with fisher-weighted averaging. *Advances in Neural Information Processing Systems*, 35:17703–17716, 2022. 2, 13
- [34] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, page 4. Granada, 2011. 4, 5
- [35] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Proceedings of the Sixth Indian Conference on Computer Vision, Graphics and Image Processing (ICVGIP)*, pages 722–729. IEEE, 2008. 5
- [36] Vardan Papyan, Xuemei Han, and David L. Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences of the United States of America*, 117:24652–24663, 2020. 4
- [37] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3498–3505, 2012. 5
- [38] Quora. Quora question pairs, 2017. 6
- [39] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI*, 2019. 6
- [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PmLR, 2021. 5
- [41] Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *International Conference on Machine Learning*, pages 2847–2854. PMLR, 2017. 4
- [42] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas, 2016. Association for Computational Linguistics. 6
- [43] Olivier Roy and Martin Vetterli. The effective rank: A measure of effective dimensionality. In *2007 15th European Signal Processing Conference*, pages 606–610. IEEE, 2007. 4
- [44] Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948. 2, 5
- [45] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA, 2013. Association for Computational Linguistics. 6
- [46] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1631–1642, 2013. 5
- [47] Johannes Stalkamp, Marc Schlipf, Jan Salmen, and Christian Igel. The German traffic sign recognition benchmark: A multi-class classification competition. In *The 2011 International Joint Conference on Neural Networks*, pages 1453–1460. IEEE, 2011. 5
- [48] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *International Conference on Machine Learning*, pages 9229–9248. PMLR, 2020. 1, 2, 4
- [49] Cheng Tai, Tong Xiao, Xiaogang Wang, and Weinan E. Convolutional neural networks with low-rank regularization. In *International Conference on Learning Representations*, 2016. 3
- [50] Anke Tang, Li Shen, Yong Luo, Han Hu, Bo Du, and Dacheng Tao. Fusionbench: A comprehensive benchmark of deep model fusion. *arXiv preprint arXiv:2406.03280*, 2024. 6
- [51] Anke Tang, Li Shen, Yong Luo, Nan Yin, Lefei Zhang, and Dacheng Tao. Merging Multi-Task Models via Weight-Ensembling Mixture of Experts. In *International Conference on Machine Learning*. OpenReview.net, 2024. 2, 5
- [52] Yuchuan Tian, Hanting Chen, Tianyu Guo, Chao Xu, and Yunhe Wang. Towards higher ranks via adversarial weight

- pruning. *Advances in Neural Information Processing Systems*, 36:1189–1207, 2023. [3](#)
- [53] Bastiaan S. Veeling, Jasper Linmans, Jim Winkels, Taco Cohen, and Max Welling. Rotation equivariant cnns for digital pathology. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 210–218, 2018. [5](#)
- [54] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *International Conference on Learning Representations*, 2021. [1](#), [2](#), [4](#)
- [55] Ke Wang, Nikolaos Dimitriadis, Guillermo Ortiz-Jimenez, François Fleuret, and Pascal Frossard. Localizing task information for improved model merging and compression. In *International Conference on Machine Learning*, 2024. [1](#), [2](#), [5](#), [6](#), [12](#), [13](#), [14](#), [15](#), [16](#)
- [56] Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. Neural network acceptability judgments. In *Transactions of the Association for Computational Linguistics*, pages 625–641, 2019. [6](#)
- [57] Adina Williams, Nikita Nangia, and Samuel R. Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana, 2018. Association for Computational Linguistics. [6](#)
- [58] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017. [5](#)
- [59] Jianxiong Xiao, Krista A Ehinger, James Hays, Antonio Torralba, and Aude Oliva. Sun database: Exploring a large collection of scene categories. *International Journal of Computer Vision*, 119:3–22, 2016. [4](#), [5](#)
- [60] Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. Ties-merging: Resolving interference when merging models. *Advances in Neural Information Processing Systems*, 36:7093–7115, 2023. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [13](#), [14](#), [15](#), [16](#)
- [61] Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guibing Guo, Xingwei Wang, and Dacheng Tao. AdaMerging: Adaptive Model Merging for Multi-Task Learning. In *International Conference on Learning Representations*, 2023. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#), [12](#), [13](#), [14](#), [15](#), [16](#)
- [62] Enneng Yang, Li Shen, Zhenyi Wang, Guibing Guo, Xiaojun Chen, Xingwei Wang, and Dacheng Tao. Representation surgery for multi-task model merging. In *International Conference on Machine Learning*, 2024. [2](#), [12](#)
- [63] Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Language Models are Super Mario: Absorbing Abilities from Homologous Models as a Free Lunch. In *International Conference on Machine Learning*, pages 5775–5775. PMLR, 2024. [2](#), [3](#), [12](#)

## Appendix

Table 5. Average multi-task performance on 8-task with merged ViT-B/32, ViT-L/14 compared to compression-based methods.

Method	ViT-B/32	ViT-L/14
TA+ <b>AdaRank</b>	87.9	92.94
CART+ <b>AdaRank</b>	89.2	93.4
EMR-Merging [19]	88.8	93.5
AdaMerging w/Surgery [62]	87.5	92.3
TSV-C [16]	89.0	93.8
TALL-Masks [55]	90.8	94.3

### A. Experiment Settings

**Hyperparameters.** As described in Section 6.1, we initialized the task-wise coefficient  $\lambda$  and the binary mask  $B$  based on the best-performing models for Task Arithmetic (TA) [21] and CART [4]. For TA, we used fixed values of  $\lambda = 0.3$  and  $B = 1$  across all indices, consistent with the reported best-performing configuration. For CART, we performed a grid search over  $\lambda \in \{0.1, 0.2, \dots, 3.0\}$  and selected the  $\lambda$  yielding the highest performance for each setting. For  $B$ , we initialized  $B = 1$  for the top 8%, 12%, 16%, and 32% of indices and chose the value that performed best.

For test-time adaptation (TTA), we followed the settings from AdaMerging [61], using the Adam [23] optimizer with a learning rate of 0.001 for vision model merging and  $5 \times 10^{-5}$  for language model merging. Both configurations used momentum values of (0.9, 0.999) and a batch size of 16.

**Details of Straight-Through Estimator.** We add a detailed explanation of how the Straight-Through Estimator [1] works with our binary masks. In the case of a single mask parameter  $B_i^l \in \{0, 1\}$ , we maintain a corresponding continuous parameter  $\tilde{B}_i^l \in \mathbb{R}$ . During the forward pass, we first apply the sigmoid function to constrain  $\tilde{B}_i^l$  between 0 and 1, and then round the result to obtain a binary mask using a threshold of 0.5. Specifically, we compute  $B_i^l = \mathbf{1}\{\sigma(\tilde{B}_i^l/T) \geq 0.5\}$ , where  $\sigma$  denotes sigmoid function, and  $T$  is the temperature. During the backward pass, we simply pass the gradient through continuous value of  $\tilde{B}_i^l$ . In our experiment, we constantly applied  $T = 10$  which performed the best among  $T \in \{1, 2, 5, 10\}$ .

**Datasets for Tasks.** In vision model merging experiments, we integrated codebases from Ilharco et al. [21] and Yang et al. [61]. We used the same fine-tuned checkpoints as Ilharco et al. [21] for the 8-task benchmark, while

for the 14- and 20-task benchmarks, we utilized checkpoints provided by Wang et al. [55]. For language model merging experiments, we adapted the codebase from Huang et al. [19] and employed checkpoints from Yu et al. [63].

**Computational Resources.** All merging, TTA, and evaluation experiments were conducted on a single NVIDIA RTX A6000 with 48GB of memory, except for the vision 8-task benchmark, which was performed on an NVIDIA RTX 3090 with 24GB of memory.

### B. Additional Experimental Results

#### Comparison with Other Variants of Model Merging

Table 5 compares our method with variants of model merging that require additional parameters or knowledge of the task index during inference. Specifically:

- **EMR-Merging** [19] involves three steps: Elect, Mask, and Rescale-Merging. It constructs and stores task-specific modulators during merging, applying them based on the task input index.
- **Representation Surgery** [62] introduces a task-specific module after the output layer to mitigate representation bias, training it with Test-Time Adaptation (TTA).
- **TSV-C** [16], a variant of TSV-M, employs low-rank task vectors to reduce storage costs.
- **TALL-Masks** [55], a variant of Consensus Merging, stores task-specific binary masks and applies them based on the task index rather than performing direct merging.

Compared to these methods, AdaRank requires no additional storage or memory after optimizing the binary mask  $B$ . Moreover, it does not depend on task index knowledge, a key distinction between model merging and multi-task *compression*. Notably, applying AdaRank to baseline methods often outperforms approaches such as EMR-Merging and Representation Surgery.

**Full Results** We present comprehensive experimental results, including individual task performances and additional baselines. For a radar plot of the ViT-L/14 experiment corresponding to Figure 4, see Figure 6. Full performance tables are provided from Table 6 to Table 11.

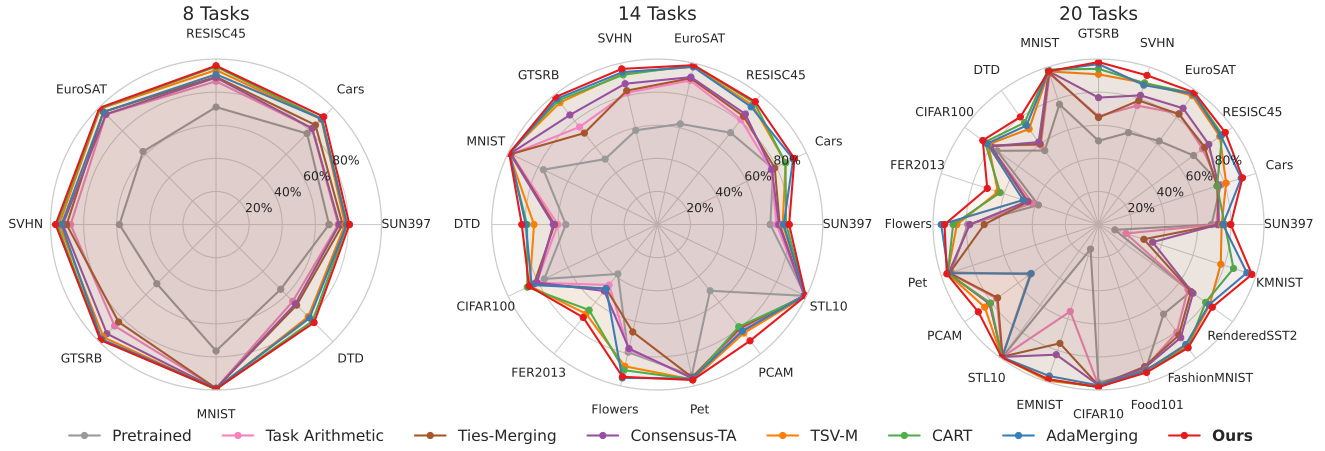


Figure 6. Individual task performance of merged ViT-L/14 model with 8, 14, 20 tasks, respectively.

Table 6. Multi-Task performance comparison on 8 Vision Tasks with Merged ViT-B/32.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	Avg.
Pretrained	62.3	59.7	60.7	45.5	31.4	32.6	48.5	43.8	48.1
Individual	75.3	77.7	96.1	99.7	97.5	98.7	99.7	79.4	90.5
Traditional MTL	73.9	74.4	93.9	98.2	95.8	98.9	99.5	77.9	89.1
Weight Averaging	65.2	63.4	71.5	71.9	64.2	52.8	87.5	50.7	65.9
Fisher Merging [33]	68.6	69.2	70.7	66.4	72.9	51.1	87.9	59.9	68.3
RegMean [22]	65.3	63.5	75.6	78.6	78.1	67.4	93.7	52.0	71.8
Task Arithmetic [21]	55.2	54.9	66.7	78.9	80.2	69.7	97.3	50.4	69.2
Ties-Merging [60]	59.8	58.6	70.7	79.7	86.2	72.1	98.3	54.2	72.5
Consensus-Ties [55]	62.5	61.8	76.3	81.6	82.0	80.5	97.3	56.0	74.8
Consensus-TA [55]	63.9	62.2	76.1	84.2	84.2	76.6	97.4	57.5	75.2
TSV-M [16]	67.2	70.8	86.3	94.6	91.0	92.3	99.3	68.9	83.8
CART [4]	68.5	73.0	88.3	95.8	87.8	93.4	99.1	72.1	84.7
Adamerging [61]	64.5	68.1	79.2	93.8	87.0	91.9	97.5	59.1	80.1
Adamerging++ [61]	66.6	68.3	82.2	94.2	89.6	89.0	98.3	60.6	81.1
<b>TA+AdaRank</b>	71.1	79.1	91.3	97.2	94.2	98.3	99.2	72.7	<b>87.9</b>
CART+AdaMerging [4]	69.5	75.1	89.3	95.7	93.0	96.8	98.9	68.4	85.8
<b>CART+AdaRank</b>	72.1	78.9	93.3	98.4	95.6	98.8	99.4	76.9	<b>89.2</b>

Table 7. Multi-Task performance comparison on 14 Vision Tasks with Merged ViT-B/32.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST
Pretrained	62.3	59.7	60.7	45.5	31.4	32.6	48.5
Individual	75.3	77.7	96.1	99.7	97.5	98.7	99.7
Weight Averaging	64.2	60.7	67.2	64.6	49.4	43.5	76.2
Task Arithmetic [21]	63.9	59.5	67.5	67.7	52.9	47.0	80.8
Ties-Merging [60]	65.1	61.8	68.3	63.7	51.3	45.9	80.0
Consensus-Ties [55]	63.6	58.8	69.7	71.9	56.2	61.2	88.3
Consensus-TA [55]	62.8	54.8	68.5	76.0	69.3	63.0	93.5
TSV-M [16]	66.3	62.1	81.2	91.7	82.4	83.6	98.8
CART [4]	68.3	60.6	86.1	91.3	72.7	82.6	98.1
AdaMerging [61]	64.3	68.5	81.7	92.6	86.6	90.8	97.5
TA+AdaRank	69.2	77.3	91.3	95.9	94.1	97.1	99.1
CART+AdaMerging [4]	67.4	72.5	87.8	96.0	90.9	95.6	98.6
CART+AdaRank	70.7	77.0	91.1	98.7	94.4	97.8	99.3

Method	DTD	CIFAR100	FER2013	Flowers	Pet	PCAM	STL10
Pretrained	43.8	64.2	39.0	66.3	87.4	60.6	97.1
Individual	79.4	89.3	73.0	90.5	91.1	87.9	98.0
Weight Averaging	47.2	69.8	41.6	68.2	88.1	61.9	97.2
Task Arithmetic [21]	48.2	69.6	42.9	67.6	87.5	63.2	96.7
Ties-Merging [60]	48.7	69.7	42.4	68.1	88.0	62.1	97.2
Consensus-Ties [55]	51.8	67.9	95.4	65.7	86.2	72.3	45.3
Consensus-TA [55]	52.4	66.6	45.3	68.3	86.9	77.0	95.6
TSV-M [16]	64.6	72.0	62.3	75.3	90.4	84.5	97.2
CART [4]	67.7	75.6	60.9	81.6	90.2	79.7	97.6
AdaMerging [61]	60.2	67.3	53.1	73.8	87.9	53.8	96.3
TA+AdaRank	72.5	75.6	45.4	82.1	92.2	59.5	97.5
CART+AdaMerging [4]	71.2	71.5	60.0	80.5	87.8	75.6	96.3
CART+AdaRank	75.6	79.4	61.4	89.1	91.7	83.0	97.7

Table 8. Multi-Task performance comparison on 20 Vision Tasks with Merged ViT-B/32.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	CIFAR100	FER2013
Pretrained	62.3	59.7	60.7	45.5	31.4	32.6	48.5	43.8	64.2	39.0
Individual	75.3	77.7	96.1	99.7	97.5	98.7	99.7	79.4	89.3	73.0
Weight Averaging	59.6	46.0	56.3	41.3	70.0	64.6	64.0	69.3	96.9	66.5
Task Arithmetic [21]	64.1	59.4	64.6	56.6	47.3	41.4	70.5	46.2	69.2	41.0
Ties-Merging [60]	64.5	57.0	68.8	59.4	48.7	48.0	78.3	49.5	70.6	43.3
Consensus-Ties [55]	64.4	58.9	67.2	54.4	51.1	47.9	77.5	48.4	67.6	96.3
Consensus-TA [55]	63.6	52.5	65.3	64.1	63.1	52.6	88.0	49.1	65.6	42.0
TSV-M [16]	64.3	52.0	75.9	87.1	75.2	76.8	94.6	61.1	68.1	58.2
CART [4]	65.3	38.1	81.3	88.7	70.0	77.4	96.2	64.6	73.7	59.9
AdaMerging [61]	62.1	66.3	78.7	92.1	72.7	90.6	93.6	57.6	66.3	48.4
TA+AdaRank	68.1	74.4	90.7	95.6	92.0	96.0	96.9	68.5	75.6	43.8
CART+AdaMerging [4]	67.3	71.2	86.3	96.6	88.3	95.0	96.4	71.2	72.4	54.4
CART+AdaRank	69.5	75.7	91.7	97.6	93.6	96.8	97.4	73.4	77.6	61.2

Method	Flowers	Pet	PCAM	STL10	EMNIST	CIFAR10	Food101	FashionMNIST	RenderedSST2	KMNIST
Pretrained	66.3	87.4	60.6	97.1	17.2	89.8	82.6	63.0	58.6	9.8
Individual	90.5	91.1	87.9	98.0	99.8	97.9	89.1	95.5	74.4	98.6
Weight Averaging	87.6	62.2	40.8	31.6	92.8	81.1	70.8	60.5	8.5	47.5
Task Arithmetic [21]	66.7	87.7	62.4	96.9	32.9	92.7	81.1	70.7	60.4	8.7
Ties-Merging [60]	71.6	85.3	64.4	96.0	39.9	93.5	75.9	72.7	64.7	12.4
Consensus-Ties [55]	67.1	86.9	67.0	42.8	41.0	92.4	79.4	74.9	60.8	11.4
Consensus-TA [55]	66.4	85.9	72.6	95.4	53.4	92.3	75.1	74.7	62.7	15.6
TSV-M [16]	71.2	88.6	84.5	96.4	95.3	93.8	77.3	85.4	70.2	57.2
CART [4]	77.6	87.7	74.8	97.0	93.1	94.8	77.1	86.5	68.6	63.4
AdaMerging [61]	65.7	87.0	54.6	96.6	21.5	90.5	80.7	82.9	65.0	10.8
TA+AdaRank	75.5	91.9	60.7	97.3	95.9	94.6	83.8	91.0	69.3	66.3
CART+AdaMerging [4]	79.6	86.4	80.0	96.9	95.1	91.5	79.8	82.7	70.0	92.5
CART+AdaRank	85.6	91.7	83.8	97.5	96.6	94.8	83.6	91.0	73.5	96.0

Table 9. Multi-Task performance comparison on 8 Vision Tasks with Merged ViT-L/14.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	Avg.
Pretrained	68.3	77.8	71.0	62.4	58.4	50.6	76.4	55.3	65.0
Individual	82.3	92.4	97.4	99.9	98.1	99.2	99.7	84.1	94.1
Traditional MTL	80.8	90.6	96.3	96.3	97.6	99.1	99.6	84.4	93.1
Weight Averaging	72.1	81.6	82.6	91.9	78.2	70.7	97.1	62.8	79.6
Fisher Merging	69.2	88.6	87.5	93.5	80.6	74.8	93.3	70.0	82.2
RegMean [22]	73.3	81.8	86.1	97.0	88.0	84.2	98.5	60.8	83.7
Task Arithmetic [21]	73.9	82.1	86.6	94.1	87.9	86.7	98.9	65.6	84.5
Ties-Merging [60]	76.5	85.0	89.3	96.3	90.3	83.3	99.0	68.9	86.1
Consensus-Ties [55]	74.9	83.6	88.7	96.6	90.5	93.2	99.1	71.1	87.2
Consensus-TA [55]	74.5	82.2	88.8	94.2	92.6	93.3	99.2	67.8	86.6
TSV-M [16]	78.0	90.0	93.4	99.0	94.8	96.3	99.5	78.8	91.2
CART [4]	79.3	90.4	95.4	99.3	96.1	98.3	99.6	82.5	92.6
AdaMerging [61]	79.0	90.3	90.8	96.2	93.4	98.0	99.0	79.9	90.8
AdaMerging++ [61]	79.4	90.3	91.6	97.4	93.4	97.5	99.0	79.2	91.0
TA+AdaRank	80.4	92.4	94.5	98.8	96.6	99.1	99.4	82.3	<b>92.9</b>
CART+AdaMerging [4]	80.1	91.5	94.7	99.3	96.8	98.9	99.5	83.6	93.1
CART+AdaRank	80.6	92.1	96.0	99.7	97.0	98.8	99.4	83.8	<b>93.4</b>

Table 10. Multi-Task performance comparison on 14 Vision Tasks with Merged ViT-L/14.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST
Pretrained	68.3	77.8	71.0	62.4	58.4	50.6	76.4
Individual	82.3	92.4	97.4	99.9	98.1	99.2	99.7
Weight Averaging	70.9	79.7	78.0	84.1	72.8	61.7	93.9
Task Arithmetic [21]	72.1	76.4	81.1	89.3	81.7	75.4	97.9
Ties-Merging [60]	74.0	78.8	83.7	90.7	83.0	70.7	98.1
Consensus-Ties [55]	72.1	75.6	84.6	95.4	87.8	83.4	97.9
Consensus-TA [55]	73.5	76.8	85.4	91.4	87.3	84.7	98.7
TSV-M [16]	75.8	86.1	92.3	98.0	93.6	94.3	99.5
CART [4]	77.9	86.0	94.1	98.8	92.8	95.9	99.5
AdaMerging [61]	76.4	91.2	91.0	97.7	94.5	97.2	98.9
TA+AdaRank	78.7	92.6	94.8	98.4	95.7	98.5	99.0
CART+AdaMerging [4]	79.8	91.8	94.5	98.2	95.2	98.1	99.1
CART+AdaRank	79.7	92.0	95.0	98.8	96.5	98.6	99.3
Method	DTD	CIFAR100	FER2013	Flowers	Pet	PCAM	STL10
Pretrained	55.3	75.8	38.2	79.1	93.6	51.2	99.4
Individual	84.1	93.3	77.0	97.9	95.5	90.3	99.5
Weight Averaging	59.7	82.7	42.5	80.5	94.7	74.2	99.4
Task Arithmetic [21]	60.1	81.1	46.7	77.5	95.1	81.1	98.8
Ties-Merging [60]	62.1	82.7	49.6	66.6	94.7	80.1	98.9
Consensus-Ties [55]	65.5	80.4	47.5	76.7	94.4	80.7	98.5
Consensus-TA [55]	62.9	80.7	51.4	76.9	95.3	82.6	98.6
TSV-M [16]	74.5	85.6	69.0	87.9	96.1	83.9	99.5
CART [4]	78.7	87.2	66.1	90.3	96.0	79.0	99.6
AdaMerging [61]	79.5	84.3	49.5	95.1	95.5	82.4	99.1
TA+AdaRank	79.9	86.9	52.1	93.3	96.1	86.8	99.4
CART+AdaMerging [4]	82.0	86.8	75.2	94.5	96.5	74.7	99.4
CART+AdaRank	81.9	86.0	71.8	94.4	96.4	89.9	99.5

Table 11. Multi-Task performance comparison on 20 Vision Tasks with Merged ViT-L/14.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	CIFAR100	FER2013
Pretrained	68.3	77.8	71.0	62.4	58.4	50.6	76.4	55.3	75.8	38.2
Individual	82.3	92.4	97.4	99.9	98.1	99.2	99.7	84.1	93.3	77.0
Weight Averaging	70.3	78.6	76.2	79.0	70.8	59.0	92.6	58.1	82.6	40.6
Task Arithmetic [21]	71.3	76.6	77.8	82.9	75.6	65.4	95.8	59.3	81.8	41.9
Ties-Merging [60]	72.5	75.4	79.3	82.6	78.8	64.7	96.6	60.2	80.7	44.8
Consensus-TA [55]	72.6	76.2	82.4	86.9	82.1	76.7	97.4	61.6	80.5	45.4
Consensus-Ties [55]	72.1	71.5	80.6	85.1	82.5	78.5	96.1	63.1	77.4	44.5
TSV-M [16]	74.4	81.1	90.6	96.3	90.0	90.8	97.3	71.4	82.4	63.9
CART [4]	76.3	75.3	92.4	97.9	89.9	94.1	98.5	76.1	84.8	62.5
AdaMerging [61]	75.2	90.7	91.4	97.6	88.6	97.0	97.7	74.0	83.2	47.9
TA+AdaRank	77.3	91.7	94.7	97.4	93.2	98.1	98.0	75.9	85.0	54.4
CART+AdaMerging [4]	79.3	91.1	93.8	98.4	93.9	97.5	97.7	81.4	85.9	74.1
CART+AdaRank	79.9	91.5	94.7	98.6	94.8	98.2	97.4	80.4	86.5	70.7
Method	Flowers	Pet	PCAM	STL10	EMNIST	CIFAR10	Food101	FashionMNIST	RenderedSST2	KMNIST
Pretrained	79.1	93.6	51.2	99.4	15.6	95.6	92.3	66.9	68.9	10.4
Individual	97.9	95.5	90.3	99.5	99.8	99.2	95.5	95.8	85.4	98.8
Weight Averaging	80.0	94.5	71.0	99.4	36.3	97.3	92.5	76.3	67.4	11.5
Task Arithmetic [21]	78.2	94.9	76.1	99.0	55.2	97.4	90.9	80.5	66.8	17.5
Ties-Merging [60]	69.1	94.7	75.4	98.7	75.5	97.3	90.3	82.6	69.1	28.8
Consensus-TA [55]	77.8	95.4	81.5	98.9	82.7	97.1	90.9	84.5	70.6	34.4
Consensus-Ties [55]	74.8	94.6	78.9	98.3	79.7	96.3	87.5	81.6	65.9	43.9
TSV-M [16]	85.6	95.9	85.0	99.3	99.3	97.9	92.3	91.0	82.9	77.7
CART [4]	87.9	95.8	80.7	99.3	98.5	98.3	92.6	91.8	80.0	85.8
AdaMerging [61]	95.1	95.4	50.3	99.1	96.3	97.2	92.7	89.7	82.5	94.3
TA+AdaRank	92.0	95.9	66.3	99.3	97.5	97.8	93.8	91.6	85.7	97.0
CART+AdaMerging [4]	95.7	96.5	79.2	99.4	98.1	97.8	93.4	91.5	85.4	96.7
CART+AdaRank	93.1	96.4	89.8	99.4	98.3	98.2	94.0	92.1	85.0	97.5