

Revisiting Weight Averaging for Model Merging

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

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

Abstract

Model merging aims to build a multi-task learner by combining the parameters of individually fine-tuned models without additional training. While a straightforward approach is to average model parameters across tasks, this often results in suboptimal performance due to interference among parameters across tasks. In this paper, we present intriguing results that weight averaging implicitly induces task vectors centered around the weight averaging itself and that applying a low-rank approximation to these centered task vectors significantly improves merging performance. Our analysis shows that centering the task vectors effectively separates core task-specific knowledge and nuisance noise within the fine-tuned parameters into the top and lower singular vectors, respectively, allowing us to reduce inter-task interference through its low-rank approximation. We evaluate our method on eight image classification tasks, demonstrating that it outperforms prior methods by a significant margin, narrowing the performance gap with traditional multi-task learning to within 1-3%.

1. Introduction

Model merging has emerged as an efficient way for constructing multi-task learners [18]. Unlike traditional multi-task learning approaches that directly train a single model on multiple tasks [1], model merging leverages individually fine-tuned models and fuses their parameters to create a model that preserves their original capabilities. This approach eliminates the need to prepare training data or store separate sets of parameters for each task, thereby reducing the costs associated with modern deep neural networks, which often require large amounts of data and numerous parameters. Consequently, model merging has been favored in various applications, including federated learning [25], model compression [32], and continual learning [10].

Alongside fundamental observations about mode connectivity in parameter space [6–8], model merging through direct interpolation between fine-tuned parameters has become prevalent in the literature. For example, the most

straightforward approach is simply averaging the parameters of different models, a method known as weight averaging [3, 10, 13, 20, 33]. However, when the knowledge encoded in each model differs significantly, weight averaging can lead to substantial interference between parameters. This interference often results in considerable performance degradation compared to traditional multi-task learning approaches, especially as the number of merged models increases.

To mitigate the limitations of the weight averaging, recent approaches have leveraged the task arithmetic framework [12], which allows for the extrapolation of the parameters. In this framework, models are assumed to be fine-tuned from a common initialization, enabling the definition of *task vectors* as the directions pointing from the initialization to the fine-tuned parameters in the parameter space. Among several arithmetic operations on task vectors, addition with a scaling coefficient has shown to be effective in merging models trained on various tasks. Subsequent approaches have improved this arithmetic by either resolving interference between task vectors [35] or applying test-time adaptation techniques for better scaling coefficient [36].

In this paper, we revisit the weight averaging strategy through the lens of task arithmetic. We begin by formulating weight averaging as a task arithmetic that induces centered task vectors around the weight average itself. We then observe that applying low-rank approximations on the centered task vectors dramatically improves the performance of the merged model produced by the corresponding task arithmetic. When an appropriate rank is chosen, this approach largely outperforms the original task arithmetic and even advanced task vector-based variants. We also note that the optimal rank can be consistently found throughout different experimental settings. To understand this surprisingly high performance, we provide both theoretical insights and empirical evidence based on the spectral analysis of task vectors. We conjecture that centering the task vectors leads to the isolation of core task-specific knowledge from the interfering noise present in their top and lower singular vectors.

Based on this observation, we propose a novel training-free model merging approach called **Centered Arithmetic**

with **Rank-reduced Task vectors (CART)**, which consistently outperforms existing baselines across various tasks and model architectures. This method can be seamlessly integrated into existing merging approaches based on task vectors, such as those using test-time adaptation. In our experiments, we evaluate CART using eight image classification tasks and demonstrate that it significantly outperforms existing model merging approaches, both in the settings with and without test-time adaptation. Our contributions are summarized as follows:

- We propose revisiting weight averaging from the perspective of task arithmetic and present intriguing properties that emerge when combined with low-rank approximations.
- We conduct a spectral analysis of the task vectors induced by weight averaging, revealing that the singular vectors effectively isolate core knowledge from interfering noise in each task.
- Our experiments on eight image classification tasks demonstrate that this simple yet effective approach significantly outperforms prior model merging methods, reducing the performance gap with traditional multi-task learning to within 1–3%.

The remaining sections are organized as follows. We first review model merging approaches including weight averaging and task arithmetic in Section 3. Section 4 introduces intriguing observations of weight averaging along with our analyses, leading to the derivation of our method. Section 5 presents and discusses experimental results on image classification tasks.

2. Related Work

Multi-task Learning Multi-task learning (MTL) aims to construct a single model that efficiently handles multiple tasks in terms of parameter and inference efficiency, typically by training the model using all available training data for each task [1]. However, the increasing size of deep learning models [21, 30] makes retraining computationally infeasible when new tasks are added. Moreover, accessing the data distribution of each task is often challenging due to privacy concerns, rendering it practically difficult to utilize all training data. While MTL opts to improve generalization and enhance the performance of each task by learning multiple tasks simultaneously, it often results in decreased performance compared to models trained on individual tasks due to task conflicts [19, 27].

Parameter Manipulation and Task Arithmetic To overcome the limitations of traditional multi-task learning, various methods have been proposed to construct a

single model capable of handling multiple tasks by manipulating model parameters alone [11, 28, 33]. However, due to the non-convexity of loss functions and the complexity of the loss landscape in deep learning tasks, naively summing the parameters of models trained on different tasks often leads to significant performance degradation on each task [17, 23]. Intriguingly, Ilharco et al. [12] demonstrated that arithmetic operations in parameter space using *task vectors*—differences between the parameters of task-specific fine-tuned models and those of a common pretrained model—consistently influence model behavior. Specifically, adding task vectors enhances performance on multiple tasks (task addition), while subtracting a task vector diminishes performance on the associated task (task negation). This suggests that task vectors capture task-specific knowledge in a way that can be algebraically manipulated. Ortiz-Jimenez et al. [24] analyzed this phenomenon, attributing it to *weight disentanglement*—where different directions in the weight space of a model independently govern distinct regions in the input space, allowing the model to manipulate these regions separately for specific tasks—and described it as an emergent property inherent to pretrained parameters.

Model Merging and Task Interference Building upon the findings of Ilharco et al. [12], multi-task models can be constructed without retraining or access to training data by performing a weighted sum of task vectors. However, task interference still leads to reduced performance on individual tasks. To mitigate this issue, recent studies have proposed methods such as Ties-Merging [35], which trims parameters with small magnitudes in the task vectors and resolves sign conflicts among parameters with large magnitudes through majority voting across tasks. Similarly, Yang et al. [36] introduced a method that employs test-time adaptation [31] to dynamically determine the coefficients for the weighted sum of task vectors based on the input data distribution, thereby reducing task interference. Nonetheless, these methods are limited in further reducing task interference because they rely on the fixed definition of task vectors from the pretrained parameters.

3. Preliminary

Problem Setup We address the problem of model merging, which seeks to combine multiple task-specific models into a single model capable of performing all tasks. Let $\theta_1, \dots, \theta_T$ denote the parameters of the models trained for T different tasks and $\mathcal{L}_t(\theta)$ represent the corresponding loss function for task t . Our objective is to merge the parameters to produce θ_* that can perform all tasks, *i.e.*, $\mathcal{L}_t(\theta_*) \approx \mathcal{L}_t(\theta_t)$ for each task $t = 1, \dots, T$. By merging models in the parameter space, we need neither an access to

training data nor maintaining all individual model parameters. We adopt the standard setting [12, 13, 20] where the models share the same architecture and are fine-tuned from a common initialization θ_0 , *e.g.*, a pre-trained vision transformer [5, 26].

Weight Averaging A simple way to merge the fine-tuned models is to average their parameters across tasks, which is called weight averaging [3, 10, 33].

$$\theta_{\text{avg}} = \frac{1}{T} \sum_{t=1}^T \theta_t. \quad (1)$$

While weight averaging has been shown to increase robustness in merging models trained for the same task with different hyper-parameters [33], it often underperforms in merging models for different tasks [20]. Several approaches replace the simple averaging with weighted sum to consider importance factors of each task [13, 20].

Task Arithmetic Task arithmetic [12] proposes an alternative approach for model merging using the concept of *task vectors*. By viewing the parameters θ as a vector in the Euclidean space, a task vector τ_t is defined as a difference between the fine-tuned parameters from the pre-trained parameters: $\tau_t = \theta_t - \theta_0$. Each task vector τ_t represents a task-specific knowledge for task t , thus merging the knowledge from different tasks is achieved by performing a simple arithmetic on their task vectors:

$$\mathcal{A}(\lambda) = \theta_0 + \lambda \sum_{t=1}^T (\theta_t - \theta_0), \quad (2)$$

where $\lambda \in \mathbb{R}$ is a hyper-parameter controlling the contribution of the task vectors. Note that Eq. (2) reduces to the weight averaging (Eq. (1)) when $\lambda = \frac{1}{T}$. With the proper choice of λ , task arithmetic deviates from the weight averaging and has been shown to improve merging performance by amplifying the task-specific knowledge encoded in the task vectors. Several variants has been proposed to improve the task arithmetic by resolving interference between task vectors [35] or adaptively choosing the coefficient λ via test-time adaptation [36].

4. Revisiting Weight Averaging

In this section, we revisit the weight averaging from the perspective of task vectors. First, we empirically show that the weight average θ_{avg} serves as a good initialization for defining task vectors with low-rank approximations (Section 4.1). Then we provide an in-depth analysis of the intriguing properties (Section 4.2). Based on these, we propose a novel model merging approach by inducing task vectors from the weight average.

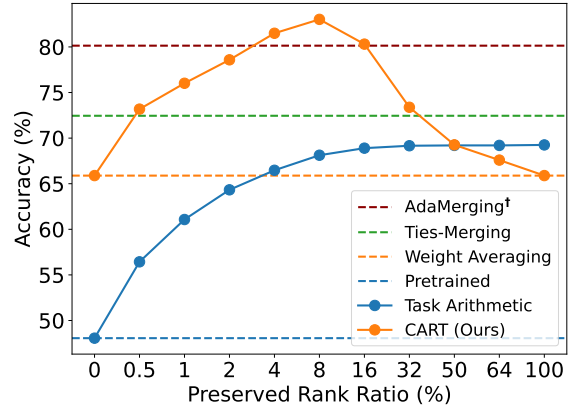


Figure 1. An average performance of model merging with low-rank approximations of task vectors. † denotes the method with test-time adaptation.

4.1. Intriguing Properties of Weight Averaging

We begin by showing that the weight averaging in Eq. (1) can be rewritten in the form of task arithmetic in Eq. (2) as follows:

$$\bar{\mathcal{A}}(\lambda) = \theta_{\text{avg}} + \lambda \sum_{t=1}^T (\theta_t - \theta_{\text{avg}}), \quad (3)$$

Note that $\bar{\mathcal{A}}(\lambda) = \theta_{\text{avg}}$ always holds regardless of λ . Eq. (3) provides a useful insight into understand weight averaging: the weight average θ_{avg} is *itself* the result of task arithmetic by posing θ_{avg} as an initial point for the vectors rather than θ_0 . The induced task vectors $\bar{\tau}_t = \theta_t - \theta_{\text{avg}}$ can be viewed as *centered*, *i.e.*, summing to zero, while the original task vectors from θ_0 are generally uncentered. Considering the inferior performance of weight averaging over task arithmetic, it may also seem to suggest that centering the task vectors with θ_{avg} is disadvantageous.

However, we observe an intriguing trend when we apply the rank reduction on the centered task vectors $\bar{\tau}_t$. Suppose the model consists of L layers and let $\theta^l \in \mathbb{R}^{m \times n}$ be the weight matrix at l -th layer of rank r . Then we apply the centered task arithmetic defined in Eq. (3) layer-wise, with low-rank approximation on the task vectors as follows:

$$\bar{\mathcal{A}}_k(\lambda) = \theta_{\text{avg}}^l + \lambda \sum_{t=1}^T \text{SVD}_k(\theta_t^l - \theta_{\text{avg}}^l), \quad \forall l \leq L, \quad (4)$$

where $\text{SVD}_k(\theta^l)$ denotes low-rank approximation of θ^l with top- k singular vectors, *i.e.*, $\text{SVD}_k(\theta^l) = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ where \mathbf{u}_i and \mathbf{v}_i denote the i -th left and right singular vectors obtained by Singular Value Decomposition (SVD), respectively.

Figure 1 illustrates the average accuracy of the merged model obtained by Eq. (4) with varying rank k (orange solid

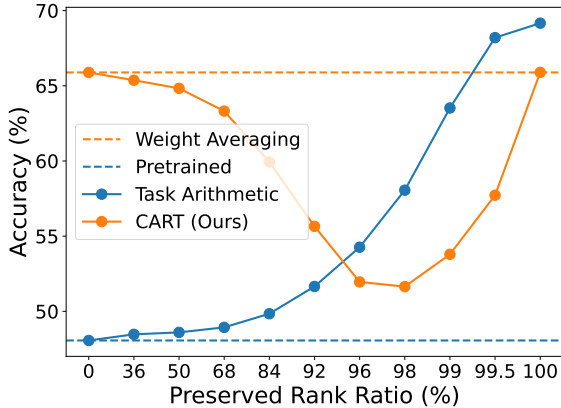


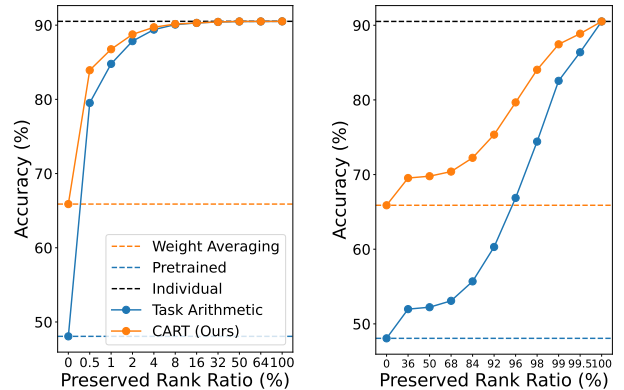
Figure 2. An average performance of model merging with reverse-ordered low-rank approximations of task vectors.

line). As the rank increases from zero to full, we observe a rapid performance improvement beyond standard weight averaging, followed by a decline back to the weight averaging level at full rank. The accuracies at the endpoints coincide with Eq. (4) since it reduces to weight averaging of Eq. (3) when $k = r$ and the summation of task vectors becomes zero when $k = 0$. However, the “bell-curve” pattern at $0 < k < r$ is intriguing since the peak is formed at surprisingly high accuracy, significantly outperforming the task arithmetic [12], its advanced variants such as TiesMerging [35] and even the one with test-time adaptation such as AdaMerging [36]. As we show in the experiment (Figure 5), this trend is consistently observed across different tasks and model sizes, with the optimal performance achieved around the same rank ($\approx 8\%$).

Notably, such properties are only observed when we perform the task arithmetic with centered task vectors $\bar{\tau}_t$. When we apply rank reduction on the original task vectors τ_t , the task arithmetic performance consistently decreases as rank reduces, collapsing to the pre-trained parameters at rank zero (blue solid line). This motivates us to further investigate the task vectors through the spectral components.

4.2. Analysis on Singular Vectors

To understand the intriguing properties of weight averaging discussed in Section 4.1, we propose a conjecture on the intrinsic structures of the task vectors with further analysis. We hypothesize that the interesting behavior of centered task vectors comes from a separation of knowledge contained in their singular vectors. The type of information represented by task vectors can be regarded as either *core knowledge* that is necessary to solve each task, or *noise* that comes from the randomness of fine-tuning. Our conjecture is that in the centered task vectors $\bar{\tau}_t$, core knowledge is concentrated in top singular vectors, while noise is dis-



(a) Proper-ordered singular vectors (b) Reverse-ordered singular vectors

Figure 3. Average performances of direct routing with proper- and reverse-ordered low-rank approximation of task vectors.

persed (ambient) in lower singular vectors. Conversely, in the original task vectors τ_t , we suspect that both core knowledge and noise are spread across all singular vectors.

Such isolation of core knowledge and noise in centered task vectors would allow their low-rank approximations to effectively retain the essential task-specific information with less interference, which explains the “bell-curve” behavior in Figure 1. When the rank is small, the performance of the centered task arithmetic rapidly increases by continuously absorbing the core knowledge. After absorbing all core knowledge, the remaining noise components cause severe interference in the parameter space, thus the performance drops with larger rank. At the full rank, the core knowledge and noise exactly counteract each other and the merged model reduces to the weight average. On the other hand, this explains why the performance of the original task arithmetic monotonically increases with rank; the core knowledge is distributed across all singular vectors.

To consolidate the conjecture, we perform the low-rank approximation in *reverse* order, where we preserve lowest k singular vectors for computing SVD_k in Eq. (4). The result of the arithmetic with reverse-ordered rank reduction is shown in Figure 2. We observe an inverted bell-curve of the centered task arithmetic (orange solid line), where the accuracy gradually decreases as the rank increases, followed by a rapid rise back to the weight averaging at high ranks near the full rank. This supports our conjecture, where the gradual performance drop is attributed to the noise from lower singular vectors, and the rapid rise near the full rank is attributed to the core knowledge from top singular vectors. The behavior of the original task arithmetic (blue solid line) is also consistent with our conjecture, where the performance still increases monotonically as rank increases, due to the distributed knowledge.

To further validate the conjecture, we apply direct rout-

ing of task vectors with low-rank approximations. In this setting, we maintain each task vector separately rather than merging the task vectors into a single model. Then to perform each task t , we route the corresponding task vector with low-rank approximation:

$$\mathcal{R}_k(t; \theta) = \theta^l + \text{SVD}_k(\theta_t^l - \theta^l), \quad (5)$$

with either $\theta^l = \theta_{\text{avg}}^l$ or $\theta^l = \theta_0^l$. Figure 3a and 3b illustrate the average accuracy of the routed models with proper and reverse-ordered low-rank approximations, respectively. From Figure 3a, we observe that the performance using centered task vectors (orange solid line) saturates in lower rank than those using original task vectors (blue solid line). This shows again that the core knowledge is concentrated to top singular vectors of centered task vectors. Similarly, Figure 3b supports that their lower singular vectors contain noise, which contributes to the slow increases in performance for most of the ranks.

4.3. Model Merging with CART

Based on the previous discussions, we propose a training-free model merging approach called **C**entered **A**rithmetic with **R**ank-reduced **T**ask vectors (**CART**). The method is simple; given T parameters $\theta_1, \dots, \theta_T$ individually trained for each task t , we apply Eq. (4) to obtain merged weight matrices of the model. Compared to the original task arithmetic, CART introduces an additional hyperparameter k . However, as we discuss in the experiment (Figure 5a), retaining 8% of the rank yields stable performance across different settings.

Thanks to the generality, CART can be plugged in any merging method that leverages the task arithmetic framework. For example, AdaMerging [36] exploits test-time adaptation to improve the performance. This introduces task-wise and layer-wise merging coefficients λ_t^l and adapts them during test time by minimizing the Shannon entropy on the test data. Combined with this variant, Eq. (4) can be transformed into

$$\tilde{\mathcal{A}}_k = \theta_{\text{avg}}^l + \sum_{t=1}^T \lambda_t^l \cdot \text{SVD}_k(\theta_t^l - \theta_{\text{avg}}^l), \quad \forall l \leq L, \quad (6)$$

which we refer as CART++ in our experiments.

5. Experiments

In this section, we present our experimental results. First, we describe the experimental setup in Section 5.1, followed by a comparison of our method with state-of-the-art model merging techniques in Section 5.2. Finally, we provide an in-depth analysis of our method, exploring its unique behavior, scalability, and sensitivity in Section 5.3.

5.1. Experimental Setup

Tasks and Models Following the prior arts in model merging [12, 35, 36], we conducted experiments on eight vision classification tasks spanning diverse domains and varying numbers of labels: Cars [14], DTD [4], EuroSAT [9], GTSRB [29], MNIST [15], RESISC45 [2], SUN397 [34], and SVHN [22]. Additionally, to examine the effects across different model scales, we employed the ViT-B-32 and ViT-L-14 architectures of CLIP [26]. Following the previous settings [12, 35, 36], we exploit task-specific classification heads obtained by the CLIP text encoder, to evaluate the merged model in each task. We report the classification accuracy at each task and their average.

Baseline Methods Following Yang et al. [36], we categorize our baselines into several groups. The first group consists of non-merging approaches, including zero-shot application of the pre-trained model (Pretrained), individually fine-tuned models for each task (Individual), and multi-task learning with joint training (Traditional MTL), which respectively represent the lower- and upper-performance bounds for model merging. The second group includes weight averaging methods, such as naive Weight Averaging and its advanced variations, including Fisher Merging [20] and RegMean [13]. The third group includes approaches based on task vectors such as Task Arithmetic [12] and Ties-Merging [35] that aims to reduce task interference. For task vector-based approaches, we additionally consider the method with test-time adaptation that adapts the merging coefficient layer-wise such as AdaMerging [36].

Implementation Details When merging the model, we applied the low-rank approximation of Eq. (4) to only the matrix component of the parameters such as weight matrices in MLP and project matrices in attention layers, while non-matrix components such as biases or ones in normalization layers are set to standard weight averaging. For hyperparameters, we observe that choosing $k = 0.08 \cdot r$, (*i.e.*, 8% of full rank) and $\lambda = 1.0$ works consistently the best across datasets and models.

5.2. Main Results

Table 1 presents the comparison with various model merging baselines with ViT-B-32 backbone. First, it shows that the performance of naive weight averaging is far below the upper bound (Individual and Traditional MTL) and generally performs worse compared to the approaches based on task vector (Task Arithmetic and Ties-Merging). However, we observe that applying the proper low-rank approximation to weight averaging by our method significantly boosts the performance, surpassing the prior art based on task vector (Ties-Merging) up to 10.6%. Considering that

Table 1. Multi-Task Performance on Eight Vision Tasks with Merged ViT-B-32 with CART.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	Average
Pretrained	62.3	59.7	60.7	45.5	31.4	32.6	48.5	43.8	48.0
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	88.9
Without Test-Time-Adaptation									
Weight Averaging	65.2	63.4	71.5	71.9	64.2	52.8	87.5	50.7	65.9
Fisher Merging [20]	68.6	69.2	70.7	66.4	72.9	51.1	87.9	59.9	68.3
RegMean [13]	65.3	63.5	75.6	78.6	78.1	67.4	93.7	52.0	71.8
Task Arithmetic [12]	55.2	54.9	66.7	78.9	80.2	69.7	97.3	50.4	69.1
Ties-Merging [35]	59.8	58.6	70.7	79.7	86.2	72.1	98.3	54.2	72.4
CART (Ours)	68.1	71.1	84.7	96.0	87.5	89.4	98.9	68.4	83.0
With Test-Time-Adaptation									
AdaMerging [36]	64.5	68.1	79.2	93.8	87.0	91.9	97.5	59.1	80.1
AdaMerging++ [36]	66.6	68.3	82.2	94.2	89.6	89.0	98.3	60.6	81.1
CART++ (Ours)	69.5	75.1	89.3	95.7	93.0	96.8	98.9	74.5	86.6

Table 2. Multi-Task Performance on Eight Vision Tasks with Merged ViT-L-14 with CART.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	Average
Pretrained	66.8	77.7	71.0	59.0	58.4	50.5	76.3	55.3	64.5
Individual	82.3	92.4	97.4	100.0	98.1	99.2	99.7	84.1	94.2
Traditional MTL	80.8	90.6	96.3	96.3	97.6	99.1	99.6	84.4	93.5
Without Test-Time-Adaptation									
Weight Averaging	72.1	81.6	82.6	91.9	78.2	70.7	97.1	62.8	79.6
Fisher Merging [20]	69.2	88.6	87.5	93.5	80.6	74.8	93.3	70.0	82.2
RegMean [13]	73.3	81.8	86.1	97.0	88.0	84.2	98.5	60.8	83.7
Task Arithmetic [12]	73.9	82.1	86.6	94.1	87.9	86.7	98.9	65.6	84.5
Ties-Merging [35]	76.5	85.0	89.3	95.7	90.3	83.3	99.0	68.6	86.0
CART (Ours)	77.3	89.9	93.5	98.8	94.3	95.9	99.4	77.4	90.8
With Test-Time-Adaptation									
AdaMerging [36]	79.0	90.3	90.8	96.2	93.4	98.0	99.0	79.9	90.8
AdaMerging++ [36]	79.4	90.3	91.6	97.4	93.4	97.5	99.0	79.2	91.0
CART++ (Ours)	80.0	92.0	94.0	98.8	96.1	98.5	99.2	82.0	92.6

Ties-Merging employed carefully hand-designed strategies to reduce conflicts among tasks, the improvement suggests that our method can be more effective in resolving inter-task conflicts by filtering out noisy spectral components of task vectors centered at weight averaging.

Our method outperforms even stronger task arithmetic approaches that involve test-time adaptation. For example, AdaMerging optimizes layer-wise coefficients of task vectors during test time, and AdaMerging++ introduces additional mechanisms to reduce inter-task interference. Without such techniques, our method achieves superior performance. When combined with the test-time adaptation (Eq. (6)), our method exhibits further improvement, closing the gap to traditional multi-task learning within 3%. This

result indicates that our method can be effectively combined with various merging techniques within the task arithmetic framework, suggesting that it may further close the performance gap with multi-task learning when integrated with more advanced techniques.

Table 2 presents the results using a larger backbone (ViT-L-14). While the performance of multi-task learning improves with a larger model, we observe even more pronounced gains across all model merging approaches, suggesting that model merging becomes increasingly effective with larger backbones. Consistent with previous observations, we also notice that our method significantly outperforms all baselines, both with and without test-time adaptation. Notably, applying test-time adaptation to our method

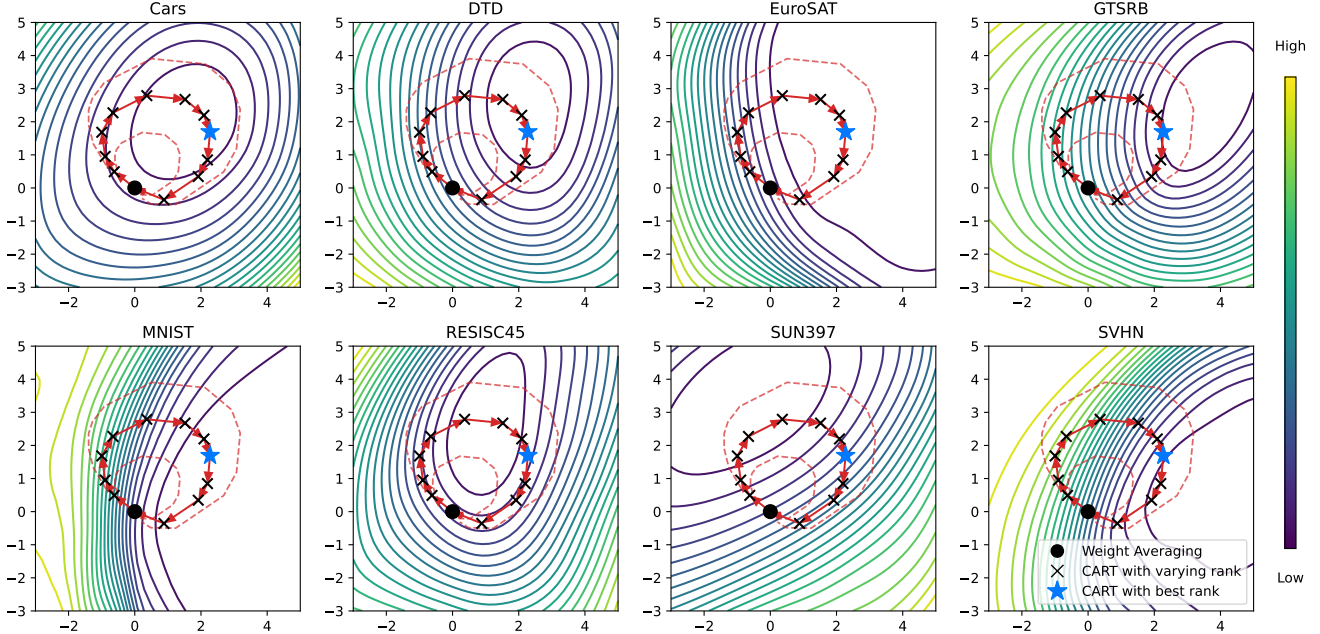


Figure 4. Loss landscape visualization of CART with varying rank k and the scaling coefficient λ for ViT-B-32 model. We plot the trajectory of the parameters as we increase k from zero to full rank, whose direction is represented as red arrows. The black circle indicates the weight average θ_{avg} and the blue star indicates CART with 8% of the full rank, which is used to report the accuracies in Table 1. CART with other ranks are indicated as black crosses. Trajectories with different λ are plotted as red dashed lines.

(CART++) further narrows the gap to the traditional multi-task learning performance within 1%. These results demonstrate that the effectiveness of our method remains consistent across different backbone sizes.

5.3. Analysis

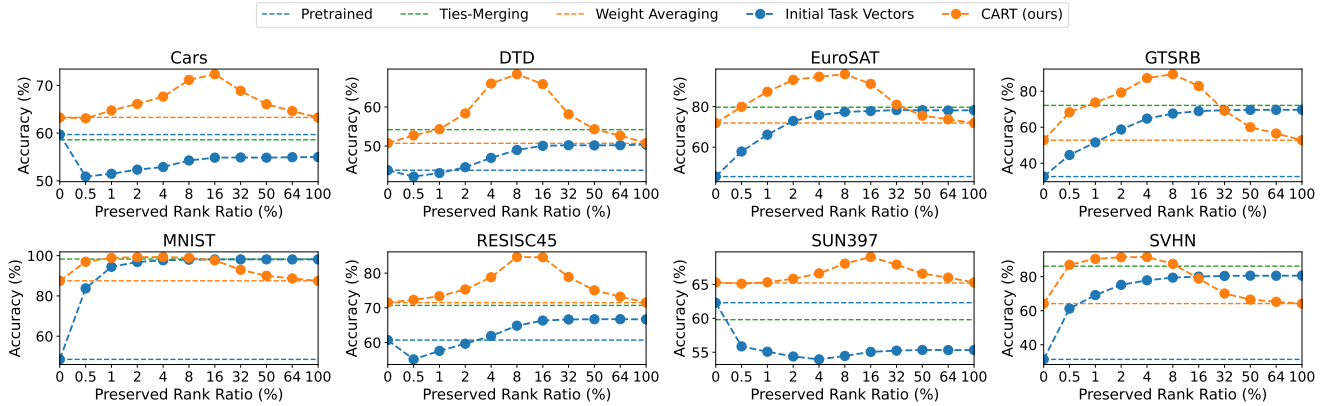
Loss Landscape Visualization To investigate the effect of rank reduction in CART, we visualize the loss landscapes around the weight average θ_{avg} . Following Li et al. [16], we vectorize parameters from all layers into a single flat vector, then project them onto a 2D space spanned by the first two principal components of the parameters. Figure 4 shows contour plots of the loss landscapes for each task, illustrating the trajectory of CART as the rank k from zero to full. Consistent to Figure 5, the trajectory forms a circular path revolves around the weight average, visiting the loss basins of each task. As discussed in Section 4.2, the top and lower singular vectors of the centered task vectors play distinct roles in the parameter space. We observe a common basin of all tasks near the weight average (located in the top-right area of each plot), where the top singular vectors of the centered task vectors guide θ_{avg} toward this common basin. Conversely, the lower singular vectors introduce noise that causes the parameters to deviate from the common basin. Combined with the scaling coefficient λ , which controls the radius of the trajectory, CART enables more dynamic exploration of the loss landscapes than the

linear movements of original task arithmetic, leading to a significant improvement in performance.

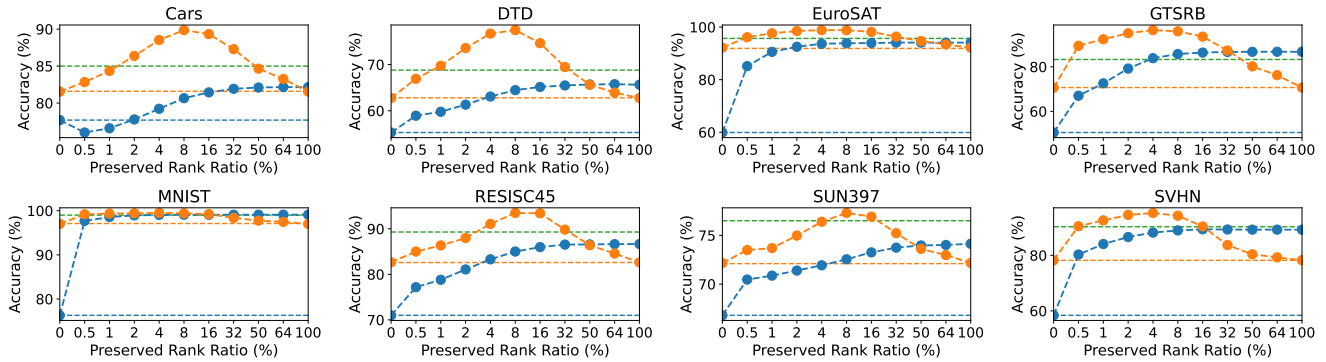
Scalability on Number of Tasks In this experiment, we assess the impact of the number of tasks on model merging performance. To this end, we construct a subset consisting of N tasks sampled from eight tasks, and report the average performance obtained by all combinations of subsets *i.e.*, $\binom{8}{N}$, while increasing N . Figure 6 illustrates the results.

We observe that methods based on simple averaging and task arithmetic exhibit significant performance degradation as the number of tasks increases. This indicates that the inter-task interference intensifies with the increasing number of tasks. On the other hand, the performance of our method remains relatively consistent across varying numbers of tasks. As discussed in Section 4.2, it suggests that eliminating less significant singular vectors of centered task vectors can effectively reduce task interference, leading to a much more scalable approach for model merging.

Optimal Rank for Low-Rank Approximation Since the performance of our method depends on the low-rank approximation of the task vector, we investigate the impact of rank (k in Eq. (4)) in merging performance. Figure 5 presents the per-task performance with varying ranks on two models, ViT-B-32 and ViT-L-14. Although the over-



(a) Performance plots for each task on ViT-B-32.



(b) Performance plots for each task on ViT-L-14.

Figure 5. Performance evaluation of CART and Task Arithmetic across eight distinct tasks with low-rank approximation on their task vectors. The plots depict the accuracy (%) of each task as a function of the preserved rank ratio (%) with two different model sizes. CART consistently outperforms Task Arithmetic at a specific rank ratio (e.g., 8%), demonstrating its effectiveness in the model merging process.

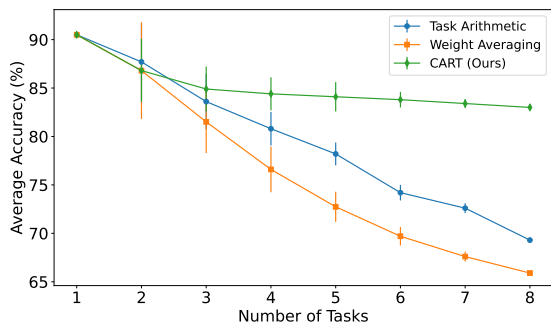


Figure 6. An average performance of merged model with increasing number of tasks. We evaluate the model with all possible combinations at each number and report the mean and standard error.

all performance depends largely on the rank of task vectors, we observe that the optimal performance lies around 8% of total dimensionality consistently over different tasks and models. It suggests that our method is relatively robust to hyperparameter tuning, reducing the need for extensive

searches to find the optimal rank.

6. Conclusion

In this study, we have revisited the effectiveness of weight averaging from the perspective of model merging. By conducting a detailed analysis to understand why it performs well, and observing these phenomena in the loss landscape, we confirmed that the centered task vectors induced from weight averaging effectively isolates the core knowledge from interfering noise in their singular vectors. Our proposed method, CART, significantly outperformed previous methodologies by a substantial margin across various classification tasks in 8 vision classification datasets. Particularly, through exhaustive experiments on various combinations of tasks, we demonstrated that our method is scalable with respect to the number of merging tasks. Moreover, since our method can be expressed within the framework of task arithmetic, it allows for incremental application of methodologies in this field. We have shown that CART++, an extension of our method, can achieve performance comparable to traditional multi-task learning.

Acknowledgement This work was in part supported by the National Research Foundation of Korea (RS-2024-00351212, RS-2024-00436165), and Institute of Information and communications Technology Planning and Evaluation (IITP) grant (RS-2024-00509279, RS-2022-II220926, RS-2022-II220959) funded by the Korean government (MSIT).

References

- [1] Rich Caruana. Multitask learning. *Machine learning*, 28: 41–75, 1997. 1, 2
- [2] 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, 1
- [3] Leshem Choshen, Elad Venezian, Noam Slonim, and Yoav Katz. Fusing finetuned models for better pretraining. *arXiv preprint arXiv:2204.03044*, 2022. 1, 3
- [4] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3606–3613, 2014. 5, 1
- [5] 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. 3, 1
- [6] Felix Draxler, Kambis Veschgini, Manfred Salmhofer, and Fred Hamprecht. Essentially no barriers in neural network energy landscape. In *International conference on machine learning*, pages 1309–1318. PMLR, 2018. 1
- [7] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pages 3259–3269. PMLR, 2020.
- [8] Timur Garipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry P Vetrov, and Andrew G Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in neural information processing systems*, 31, 2018. 1
- [9] 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, 1
- [10] Gabriel Ilharco, Mitchell Wortsman, Samir Yitzhak Gadre, Shuran Song, Hannaneh Hajishirzi, Simon Kornblith, Ali Farhadi, and Ludwig Schmidt. Patching open-vocabulary models by interpolating weights. *Advances in Neural Information Processing Systems*, 35:29262–29277, 2022. 1, 3
- [11] Gabriel Ilharco, Mitchell Wortsman, Samir Yitzhak Gadre, Shuran Song, Hannaneh Hajishirzi, Simon Kornblith, Ali Farhadi, and Ludwig Schmidt. Patching open-vocabulary models by interpolating weights. *Advances in Neural Information Processing Systems*, 35:29262–29277, 2022. 2
- [12] 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, 4, 5, 6
- [13] Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. Dataless knowledge fusion by merging weights of language models. In *International Conference on Learning Representations*, 2023. 1, 3, 5, 6
- [14] 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, 1
- [15] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998. 5, 1
- [16] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018. 7
- [17] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018. 2
- [18] Weishi Li, Yong Peng, Miao Zhang, Liang Ding, Han Hu, and Li Shen. Deep model fusion: A survey. *arXiv preprint arXiv:2309.15698*, 2023. 1
- [19] Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. *Advances in Neural Information Processing Systems*, 34:18878–18890, 2021. 2
- [20] Michael S Matena and Colin A Raffel. Merging models with fisher-weighted averaging. *Advances in Neural Information Processing Systems*, 2022. 1, 3, 5, 6
- [21] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124003, 2021. 2
- [22] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bis-sacco, 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. 5, 1
- [23] Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learning? *Advances in neural information processing systems*, 33:512–523, 2020. 2
- [24] Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. Task arithmetic in the tangent space: Improved editing of pre-trained models. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [25] Pian Qi, Diletta Chiaro, Antonella Guzzo, Michele Ianni, Giancarlo Fortino, and Francesco Piccialli. Model aggregation techniques in federated learning: A comprehensive survey. *Future Generation Computer Systems*, 150:272–293, 2024. 1
- [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,

- Amanda Aspell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 2021. [3](#), [5](#), [1](#)
- [27] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018. [2](#)
- [28] Sidak Pal Singh and Martin Jaggi. Model fusion via optimal transport. *Advances in Neural Information Processing Systems*, 33:22045–22055, 2020. [2](#)
- [29] Johannes Stalldkamp, 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](#), [1](#)
- [30] Pablo Villalobos, Jaime Sevilla, Tamay Besiroglu, Lennart Heim, Anson Ho, and Marius Hobbhahn. Machine learning model sizes and the parameter gap. *arXiv preprint arXiv:2207.02852*, 2022. [2](#)
- [31] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020. [2](#)
- [32] Ke Wang, Nikolaos Dimitriadis, Guillermo Ortiz-Jimenez, François Fleuret, and Pascal Frossard. Localizing task information for improved model merging and compression. In *Forty-first International Conference on Machine Learning*, 2024. [1](#)
- [33] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International conference on machine learning*, pages 23965–23998. PMLR, 2022. [1](#), [2](#), [3](#)
- [34] 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. [5](#), [1](#)
- [35] 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, 2024. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#)
- [36] Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guibing Guo, Xingwei Wang, and Dacheng Tao. Adamerging: Adaptive model merging for multi-task learning. *International Conference on Learning Representations*, 2024. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#)

A. Implementation Details

We provide the implementation details of our experiments.

Detailed Experimental Settings The baseline codes and pre-trained weights used in our experiments are available in the repository provided by task arithmetic [12]. This includes fine-tuned checkpoints for the eight vision classification tasks (Cars [14], DTD [4], EuroSAT [9], GT-SRB [29], MNIST [15], RESISC45 [2], SUN397 [34], and SVHN [22]) across both ViT-B-32 and ViT-L-14 architectures [5]. We used the CLIP [26] implementations from the OpenCLIP¹ library, which provides the pretrained weights.

Searching the Scaling Coefficient In CART, we performed a grid search for the scaling coefficient $\lambda \in (0.05, 0.1, 0.2, \dots, 1.0)$ and selected λ that yielded the best performance, following the task arithmetic settings [12]. In CART++, which employs test-time adaptation to optimize the scaling coefficients, we adopted the settings from AdaMerging [36]. Using the test dataset for each task, we iteratively optimized the scaling coefficients separately until the evaluation accuracy converged to an upper bound. We provide the pseudocode for CART and CART++ in Algorithm 1, where the test-time adaptation of CART++ is described in Algorithm 2.

B. Further Analysis on Loss Landscapes

We provide further analysis on the loss landscapes near the model parameters obtained by CART and task arithmetic.

Multi-Task Loss Landscape To directly compare the behaviors of CART and task arithmetic in the parameter space, we visualize the loss landscapes of them with respect to the multi-task loss (e.g., $\mathcal{L} = \sum_{t=1}^T \mathcal{L}_t$). As observed in the left panel of Figure 7, CART can closely approach the common basin induced by the multi-task loss, with its trajectory rotating from θ_{avg} as the rank increases. In contrast, for task arithmetic, we observed that even as the rank increases (up to full rank), it converges to a point farther from the common basin than CART, which is consistent with Figure 1. In the right panel of Figure 7, we examine the trajectory of task arithmetic with varying λ . We observe that when $\lambda = 0.125$, the trajectory moves toward the weight average, and as λ increases (e.g., to 1), there is a tendency to move away from the common basin. This trend aligns with the results presented by Ilharco et al. [12].

Task-wise Loss Landscapes In Figure 8, we investigate the loss landscape around the fine-tuned weights θ_t for each task. We confirm that the θ_t are located within the basin

¹https://github.com/mlfoundations/open_clip

Algorithm 1 Model Merging with CART

Require: Fine-tuned parameters $\{\theta_t\}_{t=1}^T$, rank pruning ratio γ , scaling coefficient λ , test data $\mathcal{D}_{\text{test}}$ (optional)

Ensure: Merged model parameters θ_*

- 1: **for** each layer l **do**
- 2: Compute the average of parameters:

$$\theta_{\text{avg}}^l = \frac{1}{T} \sum_{t=1}^T \theta_t^l$$

- 3: **for** each task t **do**
- 4: Compute the centered task vectors:

$$\bar{\tau}_t^l = \theta_t^l - \theta_{\text{avg}}^l$$

- 5: Perform SVD: $\bar{\tau}_t^l = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$
- 6: Compute the pruning rank $k = \lceil r \times \gamma \rceil$
- 7: Compute the low-rank approximation:

$$\bar{\tau}_t^l \leftarrow \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$$

- 8: **end for**
- 9: **end for**
- 10: **if** Test-Time Adaptation is enabled **then**
- 11: **Call** Algorithm 2 to optimize $\{\lambda_t^l\}$ using $\mathcal{D}_{\text{test}}$
 ▷ CART++
- 12: **else**
- 13: Set $\lambda_t^l = \lambda, \quad \forall t, l$
 ▷ CART
- 14: **end if**
- 15: **for** each layer l **do**
- 16: **Merge** models using optimized coefficients:

$$\theta_*^l = \theta_{\text{avg}}^l + \sum_{t=1}^T \lambda_t^l \bar{\tau}_t^l$$

- 17: **end for**
 - 18:
 - 19: **return** Merged model parameters θ_*
-

for most tasks. Additionally, by examining the positions of CART and task arithmetic, we observe that CART is closer to the loss basin of the single tasks, which is consistent with the results shown in Figure 5a.

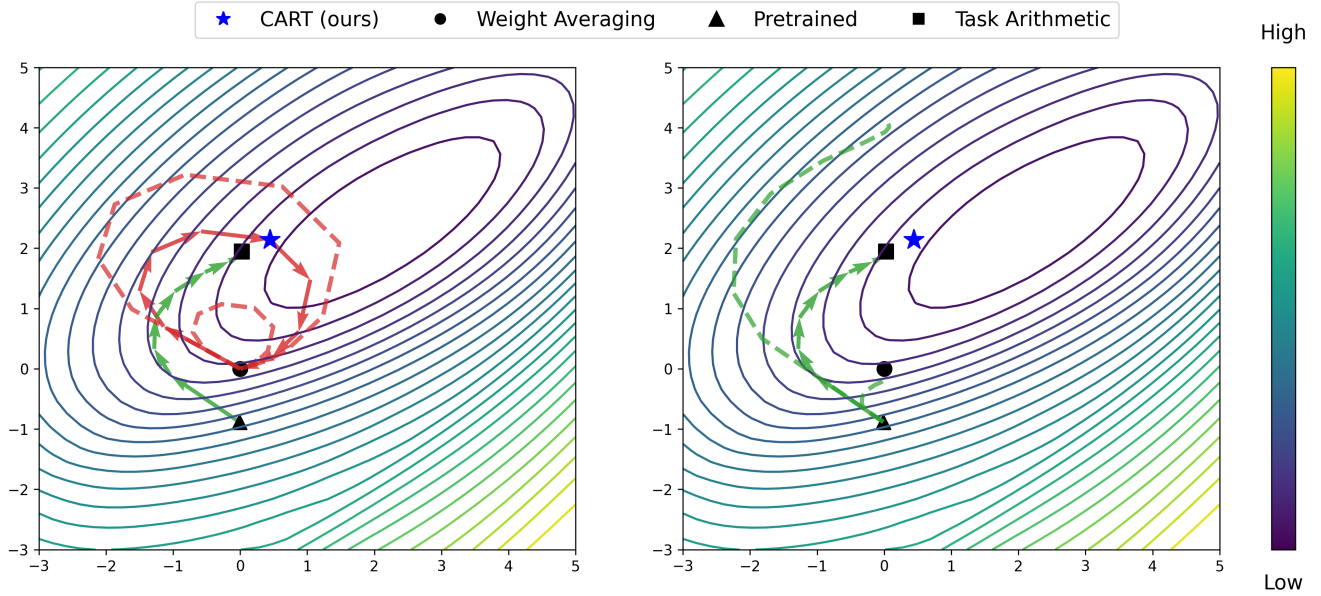


Figure 7. We present a visualization of the multi-task loss landscape (i.e., $\sum_t \mathcal{L}_t(\theta)$) for CART and task arithmetic using the ViT-B-32 model with varying values of k . For both model merging methodologies, we plot the parameter trajectories by incrementally increasing k from 0 to the full rank. Specifically, in our approach, the trajectory begins and returns to the weight average θ_{avg} , represented by a red line. In contrast, task arithmetic’s trajectory progresses from the pretrained model to the merged weights, depicted by a green line. Additionally, the black circle indicates the weight average θ_{avg} , the blue star marks CART with 8% of the full rank, the black triangle represents the pretrained model, and the black square denotes the merged model obtained through task arithmetic. Trajectories corresponding to different λ values are illustrated as dashed lines.

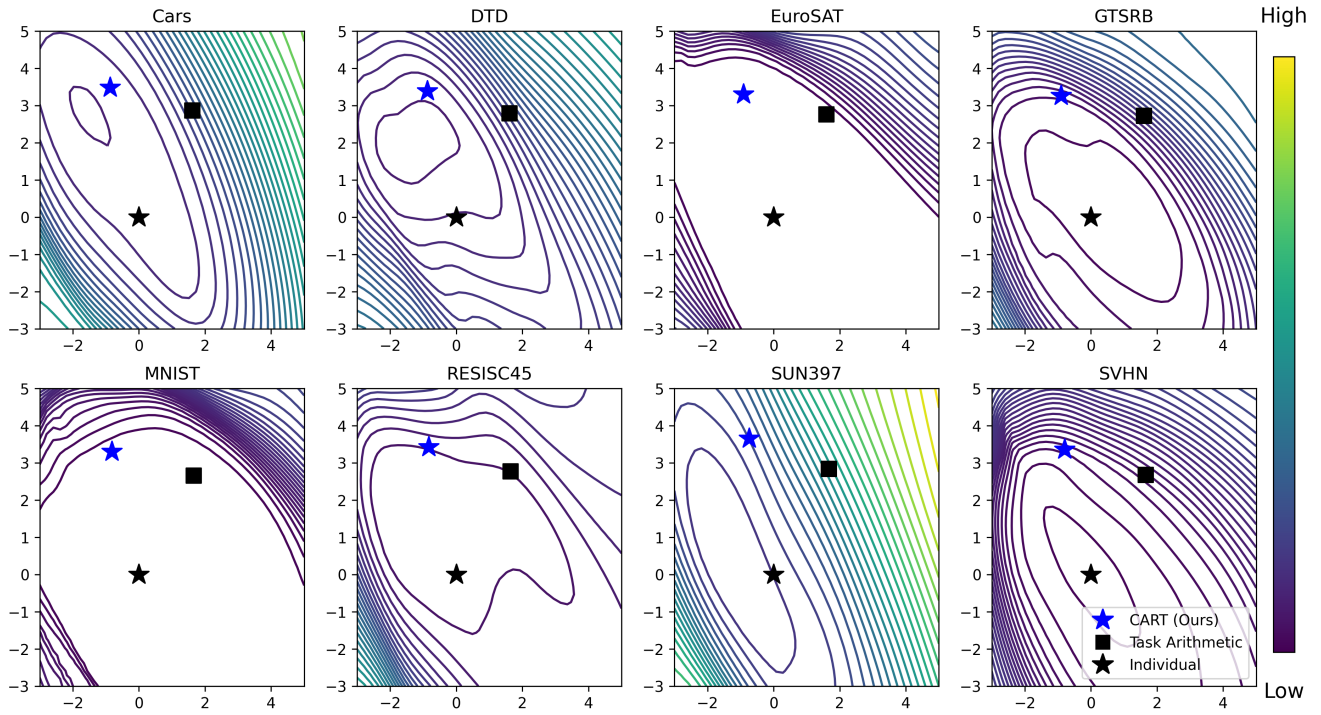


Figure 8. Visualizations of the loss landscape surrounding the individually fine-tuned models for each task using the ViT-B-32 model. The black star represents the individual parameters θ_t , the blue star denotes CART with 8% of the full rank, and the black square indicates the merged model obtained through task arithmetic.

Algorithm 2 Entropy-Based Test-Time Adaptation

Require: Test data $\mathcal{D}_{\text{test}}$, task vectors $\{\bar{\tau}_t^l\}$, average parameters θ_{avg}^l , learning rate η , iterations N

- 1: Initialize scaling coefficients $\{\lambda_t^l\}$
- 2: **for** $n = 1$ **to** N **do**
- 3: **for** each layer l **do**
- 4: Compute merged model parameters:

$$\theta_*^l = \theta_{\text{avg}}^l + \sum_{t=1}^T \lambda_t^l \bar{\tau}_t^l$$

- 5: **end for**
- 6: Compute entropy loss on test data:

$$\mathcal{L}_{\text{entropy}} = -\frac{1}{|\mathcal{D}_{\text{test}}|} \sum_{\mathbf{x} \in \mathcal{D}_{\text{test}}} \sum_c p(c|\mathbf{x}; \theta_*) \log p(c|\mathbf{x}; \theta_*)$$

- 7: **for** each task l **do**
- 8: **for** each task t **do**
- 9: Compute gradient $\nabla_{\lambda_t^l} \mathcal{L}_{\text{entropy}}$
- 10: Update scaling coefficients:

$$\lambda_t^l \leftarrow \lambda_t^l - \eta \nabla_{\lambda_t^l} \mathcal{L}_{\text{entropy}}$$

- 11: **end for**
 - 12: **end for**
 - 13: **end for**
 - 14:
 - 15: **return** Optimized scaling coefficients $\{\lambda_t^l\}$
-