003

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

# TASK ARITHMETIC IN TRUST REGION: A TRAINING-FREE MODEL MERGING APPROACH TO NAVIGATE KNOWLEDGE CONFLICTS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Multi-task model merging offers an efficient solution for integrating knowledge from multiple fine-tuned models, mitigating the significant computational and storage demands associated with multi-task training. As a key technique in this field, Task Arithmetic (TA) defines task vectors by subtracting the pre-trained model ( $\theta_{pre}$ ) from the fine-tuned task models in parameter space, then adjusting the weight between these task vectors and  $\theta_{pre}$  to balance task-generalized and taskspecific knowledge. Despite the promising performance of TA, conflicts can arise among the task vectors, particularly when different tasks require distinct model adaptations. In this paper, we formally define this issue as knowledge conflicts, characterized by the performance degradation of one task after merging with a model fine-tuned for another task. Through in-depth analysis, we show that these conflicts stem primarily from the components of task vectors that align with the gradient of task-specific losses at  $\theta_{pre}$ . To address this, we propose Task Arithmetic in Trust Region (TATR), which defines the trust region as dimensions in the model parameter space that cause only small changes (corresponding to the task vector components with gradient orthogonal direction) in the task-specific losses. Restricting parameter merging within this trust region, TATR can effectively alleviate knowledge conflicts. Moreover, TATR serves as both an independent approach and a plug-and-play module compatible with a wide range of TAbased methods. Extensive empirical evaluations on eight distinct datasets robustly demonstrate that TATR improves the multi-task performance of several TA-based model merging methods by an observable margin.

## 1 INTRODUCTION

The growing adoption of large foundation models is ac-037 companied by significant practical challenges in terms of computational and storage demands (Kaplan et al., 2020). To address these challenges, multi-task model 040 merging (Matena & Raffel, 2022) has emerged as a 041 promising solution. For example, Task Arithmetic (Il-042 harco et al., 2023b) merges models by summing the task 043 vectors from multiple tasks and applying them to the 044 pre-trained model. Here task vectors are the difference in model parameters between the pre-trained foundation model and its fine-tuned version on a specific task. This 046 approach builds a high-performance multi-task model 047 by simple arithmetic operations in the model parameter 048 space, thereby reducing computational overheads associated with fine-tuning on multiple tasks. 050

051 Despite their successes, task arithmetic and its variants
(Yadav et al., 2023; Wang et al., 2024; Yang et al., 2024b;a) still suffer from conflicts between task vectors. As illustrated in Figure 1, adding task vectors pointing to



Figure 1: Illustration of knowledge conflicts between task vectors. In scenario (a), the two task vectors contain largemagnitude components in opposite directions. In scenario (b), the difference in vector magnitudes causes the merged model to be dominated by one task. Both lead to suboptimal performance in one or more tasks.

largely opposite directions may lead to catastrophic forgetting, and inconsistent task vector magni-055 tudes may cause unbalanced merging, allowing the resulting model to be disproportionately influ-056 enced by a small subset of tasks. We refer to this issue as **knowledge conflicts**, represented as the 057 expected performance variation of one task observed before and after merging another task vector. 058 Knowledge conflicts differ from the typical notion of negative transfer (Yang et al., 2022; Meng et al., 2021; Liu et al., 2021b; Wang et al., 2023), as the former specifically refers to conflicts between predetermined, static task vectors, whereas the latter typically describes dynamic interference 060 among tasks during training. Although current methods like sign alignment and test-time adaptation 061 partially address knowledge conflicts, a thorough analysis of the root causes and a dedicated solution 062 remain elusive. 063

In this paper, we propose a novel trust-region criterion for model merging, Task Arithmetic in the
 Trust Region (TATR), which addresses the knowledge conflict problem. The trust region contains
 dimensions in the model parameter space that cause only small changes in the task-specific losses.
 When merging models, only the components of task vectors in the trust region are added to the
 pre-trained weights; other dimensions are discarded. TATR can be used independently or jointly
 with other techniques like Ties-Merging (Yadav et al., 2023), AdaMerging (Yang et al., 2024b), and
 Surgery (Yang et al., 2024a).

071 TATR utilizes the first-order Taylor series to compute the changes in the task-specific losses. It contrasts with a simplistic approach that selects components of task vectors that align with the 072 negative gradient direction. While the simplistic approach is intuitive, empirical evidence reveals 073 that it usually does not alleviate knowledge conflicts. We contend that, as the task vectors have 074 large magnitudes, the first-order Taylor series fails to approximate the function well, leading to 075 performance degradation. In contrast, TATR identifies directions that are orthogonal to the gradient, 076 along which minimal cross-task interference happens. Due to overparameterization and parameter 077 redundancy in the models (Dalvi et al., 2020; Chen et al., 2022b), such directions are usually not 078 difficult to find.

- 079
- 081
- 082
- 083 084
- 085
- 086
- 087

80

090

092

In summary, the contributions of this paper are as follows:

- We conduct an analysis of knowledge conflicts that arise during model merging. Our investigation reveals that the components of task vectors aligned with the gradient of task-specific losses are the primary source of knowledge conflicts.
  - We propose an approach, Task Arithmetic in the Trust Region (TATR), which defines a trust region to address knowledge conflicts. TATR serves as both an independent approach and a plug-and-play module compatible with a wide range of TA-based methods.
- We evaluate TATR through experiments across eight datasets. The experimental results demonstrate that TATR effectively mitigates knowledge conflicts and improves the performance of several TA-based model merging methods by an observable margin.
- 2 RELATED WORK

## 093 094 2.1 Traditional Multi-Task Learning

Multi-task learning (MTL) aims to improve performance by sharing knowledge across related tasks 096 (Zhang & Yang, 2022). A significant challenge for MTL is negative transfer (Liu et al., 2017; Zhang 097 et al., 2023b), where joint training on conflicting tasks yields performance lower than training on the 098 tasks individually. Various solutions to negative transfer have been proposed, such as modularization (Tang et al., 2020; Ma et al., 2018), sparsification (Ding et al., 2021; Sun et al., 2020; Liu et al., 2019), and soft parameter sharing (Gao et al., 2020; Hazimeh et al., 2021). Other strategies focus 100 on optimizing task interactions, such as adjusting task-specific loss weights (Sener & Koltun, 2018; 101 Liu et al., 2019; 2022; Hu et al., 2023; Chen et al., 2022a), resolving gradient direction conflicts (Yu 102 et al., 2020; Chen et al., 2020; Liu et al., 2021a; Javaloy & Valera, 2022; Navon et al., 2022), or 103 preventing the dominance of certain tasks (Chen et al., 2018; He et al., 2022; Yang et al., 2023). 104

Traditional MTL are not well-suited for merging foundation models. First, retraining these models
 using vast amounts of data incurs significant computational costs. Large-scale foundation models
 are already resource-intensive, and training them with multi-task objectives further amplifies these
 demands, requiring immense computation and time. Additionally, retraining from scratch wastes

valuable knowledge optimized in each individual expert model. These considerations have driven
 the development of model merging as an alternative to multi-task learning.

110 111 112

113

## 2.2 MULTI-TASK LEARNING THROUGH MODEL MERGING

114 Model merging techniques, which aim to integrate knowledge across models, have attracted in-115 creasing attention in recent years. As a precursor, Stochastic Weight Averaging (SWA) (Izmailov 116 et al., 2018) averages model weights near the end of training. This concept was further advanced by approaches like SWAD (Cha et al., 2021) and Ensemble of Averages (EoA) (Arpit et al., 2022). 117 Empirical evidence from Ilharco et al. (2023a) demonstrates that parameter averaging effectively 118 integrates knowledge from models trained on diverse tasks. DLCPA (Sun et al., 2023) proposes 119 to apply cumulative parameter averaging (CPA) to continually assimilate knowledge across dis-120 tinct tasks. Fisher-Merging (Matena & Raffel, 2022) leverages the Fisher information matrix Fisher 121 (1925) to measure the importance of model parameters and merge models using weighted averaging. 122 Additionally, RegMean (Jin et al., 2023) formulates an optimal merging model by minimizing the 123 distance to each model in the parameter space. 124

Recently, Task Arithmetic (TA) (Ilharco et al., 2023b) innovatively proposes the concept of "task 125 vector", defined as the vector from a pre-trained model to its fine-tuned counterpart in the parameter 126 space. By weighting these task vectors and adding them back to the pre-trained model, TA strikes a 127 harmonious balance between generalized knowledge from the pre-train model and the task-specific 128 knowledge in the task vectors. Following this insight, Ties-Merging (Yadav et al., 2023) refines the 129 fusion process by discarding parameters deemed insignificant or of low magnitude. PEFT (Zhang 130 et al., 2023a) and MoLE (Wu et al., 2024) further extend TA by integrating it with LoRA (Hu et al., 131 2022) modules. Furthermore, Ortiz-Jimenez et al. (2023) suggests fine-tuning models in the tangent 132 space, which can effectively mitigate conflict between task vectors.

Furthermore, several approaches combine test-time adaptation techniques with TA, yielding superior MTL performance. These test-time adaptation-based methods typically allocate merging weights and fine-tune them during testing using unsupervised test data. For instance, AdaMerging (Yang et al., 2024b) trains a set of merging coefficients for layers, while other methods fit lightweight adapter modules, such as representation surgery (Yang et al., 2024a) and MoE router (Tang et al., 2024).

139 140

## 3 PRELIMINARIES

## 3.1 PROBLEM SETTING

Formally, let  $\theta_{\text{pre}} \in \mathbb{R}^N$  denote the set of N parameters of a pre-trained model, which is initially trained using a diverse, large-scale dataset to encapsulate generalized, task-agnostic knowledge. Subsequently,  $\theta_{\text{pre}}$  undergoes fine-tuning for K distinct downstream tasks, yielding a set of finetuned parameters  $\{\theta_k\}_{k=1}^K$ , where each  $\theta_k$  is tailored to a specific task k.

The objective of model merging is to integrate these fine-tuned parameters from the task-specific models  $\{\theta_k\}_{k=1}^K$  into a single model  $\theta_{MTL}$ . This merged model  $\theta_{MTL}$  aims to achieve effective generalization across all K tasks without resorting to trivial solutions such as retraining from scratch or requiring full access to the training datasets of all tasks.

154 155

156

## 3.2 TASK ARITHMETIC

Task arithmetic (TA) (Ilharco et al., 2023b) is known as a competitive baseline for model merging by leveraging **task vectors**, which are defined as the differential parameters between the pre-trained model  $\theta_{pre}$  and each fine-tuned model. Specifically, the task vector for task k is given by:

$$\Delta_k = \theta_k - \theta_{\rm pre}.\tag{1}$$

TA posits that these task vectors encapsulate essential task-specific knowledge. The merged model is then constructed by adding the cumulative task vector from all tasks back to  $\theta_{pre}$ :

$$\theta_{\rm TA} = \underbrace{\theta_{\rm pre}}_{\rm task-generalized} + \lambda \underbrace{\sum_{k} \Delta_k}_{\rm task-specific} , \qquad (2)$$

where  $\lambda > 0$  is a pre-defined hyper-parameter that governs the influence of task-specific adjustments.

This method is favored over direct averaging of fine-tuned parameters as it seeks a balance between
generalized and task-specific knowledge, contributing to its competitive advantage. Nevertheless,
task vectors can encode conflicting adaptations across different tasks, leading to potential knowledge
conflicts that may result in information loss and diminished performance. This issue, termed as
"knowledge conflict", will be dissected further in the subsequent section.

## 4 ANALYZING KNOWLEDGE CONFLICT

Knowledge conflict frequently arises when merging MTL models, as the expert models encapsulate diverse, sometimes conflicting, knowledge. We formally define knowledge conflict as follows:

**Definition 1** (Knowledge Conflict). *Given a pre-trained model*  $\theta_{\text{pre}}$  *and a set of fine-tuned, task-specific models*  $\{\theta_k\}_{k=1}^K$ , where  $\theta_k$  represents the parameters optimized for task k, the knowledge conflict on task j caused by task i can be quantified by the change in performance of task j when task i is included in the model merging process. Formally, the knowledge conflict is defined as

$$\mathcal{C}_{j|i} := L_j \left( \theta_{\text{MTL}}(\{\theta_k\}_{k=1}^K) \right) - L_j \left( \theta_{\text{MTL}}(\{\theta_k\}_{k\neq i}) \right),$$

where  $L_j(\theta)$  denotes the loss for task j with model parameters  $\theta$ , and  $\theta_{MTL}(\{\theta_k\}_{k \neq i})$  represents the merged model parameters excluding the model fine-tuned for task i. The overall knowledge conflict, C, is computed as the sum of  $C_{j|i}$  across all task pairs (i, j):

 A higher value of C indicates a greater degree of conflict, as it reflects a larger negative impact on task j's performance when task i is incorporated into the merging process.

 $\mathcal{C} := \sum_{i 
eq j} \mathcal{C}_{j|i}.$ 

Knowledge conflict can be regarded as a special case of negative transfer, although these concepts emphasize different aspects. In traditional MTL, negative transfer typically refers to the *dynamic* interference between tasks during joint training, where conflicting gradients impede the model from learning effective representations (Zhang et al., 2023b). In contrast, the knowledge conflict defined here highlights a static nature among the fine-tuned model parameters, where further training to resolve task interference is prohibited. Each fine-tuned model has already encoded task-specific knowledge, which may be inherently incompatible with that of other tasks. As a result, knowl-edge conflict in model merging presents a unique challenge, necessitating methods that can align or reconcile parameters without resorting to retraining.

In the context of TA, knowledge conflict can be further articulated through task vectors:

$$C_{\mathrm{TA}j|i} = L_j \left( \theta_{\mathrm{pre}} + \lambda \sum_k \Delta_k \right) - L_j \left( \theta_{\mathrm{pre}} + \lambda \sum_{k \neq i} \Delta_k \right).$$
(3)

An intuitive hypothesis is that task vector components aligned with the gradient ascent direction contribute to knowledge conflicts. More formally, we apply a Taylor expansion around  $\theta_{pre}$  on the

227

228

229

230

231

232

233

234 235



Figure 2: (a) Performance comparison across eight datasets (Cf. Section 6.1) when merging negative components (e.g., components aligned with the loss descent direction) and orthogonal components (e.g., components orthogonal to the gradient) of task vectors, corresponding to  $\theta_{TATR}^{neg}$  and  $\theta_{TATR}^{orth}$ , respectively. (b) Loss landscape of the EuroSAT dataset and the components of the cumulative task vector from the remaining seven datasets. (c) The total loss landscape over all eight datasets, along with the components of the cumulative across task vectors. Note that in both (b) and (c), the loss landscape is visualized in a hyperplane going through the three points:  $\theta_{pre} + \Delta^{\perp}, \theta_{pre} + \Delta^{+}$ , and  $\theta_{pre} + \Delta^{-}$ . Refer Section B in Appendix for more details.

right-hand side of Eq. (3):

239 240 241

$$\approx L_{j}(\theta_{\rm pre}) + \left\langle \nabla_{\theta}L_{j}(\theta_{\rm pre}), \lambda \sum_{k} \Delta_{k} \right\rangle - L_{j}(\theta_{\rm pre}) - \left\langle \nabla_{\theta}L_{j}(\theta_{\rm pre}), \lambda \sum_{k \neq i} \Delta_{k} \right\rangle$$
(4)

243 244 245

257

258 259

260

261

262

264

265

242

$$= \lambda \left\langle \nabla_{\theta} L_{j}\left(\theta_{\text{pre}}\right), \Delta_{i} \right\rangle = \lambda \sum_{n=1}^{N} \nabla_{\theta} L_{j}\left(\theta_{\text{pre}}\right) [n] \cdot \Delta_{i}[n]$$

 $L_j \left( \theta_{\text{pre}} + \lambda \sum \Delta_k \right) - L_j \left( \theta_{\text{pre}} + \lambda \sum \Delta_k \right)$ 

246 where v[n] selects the *n*-th component of the vector v. Equation (4) suggests that task vector com-247 ponents aligned with the gradient ascent direction are primarily responsible for knowledge conflicts, 248 while those in line with the gradient descent direction should be prioritized during model merging. 249 That is, we should avoid merging the *n*-th component of the task vector  $\Delta_i$  if  $\nabla_{\theta} L_i(\theta_{\text{pre}})[n] \cdot \Delta_i[n]$ 250 is large.

251 On the other hand, perhaps counter-intuitively, empirical evidence (Figure 2(a)) shows that merging 252 the gradient **descent** components (i.e.,  $\nabla_{\theta} L_{i}(\theta_{pre})[n] \cdot \Delta_{i}[n] < 0$ ) causes a significant performance 253 drop. A potential explanation for this phenomenon is that the task vectors have large magnitudes, 254 thereby the first-order Taylor expansion cannot offer a good approximation of the task loss  $L_i$ . As a 255 result, even if we merge a component in the gradient descent direction, we can overshoot the local 256 optimum and end up increasing the task loss (Ruder, 2017).

To facilitate analysis, we decompose the task vector  $\Delta_i$  into three components:

- Orthogonal component, which contains elements with near-zero inner product  $\Delta_i^{\perp} =$  $\Delta_i \odot \mathbb{1}_{\{\nabla_{\theta} L_i(\theta_{\rm pre}) \odot \Delta_i \approx 0\}};$
- **Positive component**, with elements having a positive inner product  $\Delta_i^+ = \Delta_i \odot$  $\mathbb{1}_{\{\nabla_{\theta} L_j(\theta_{\text{pre}}) \odot \Delta_i > 0\}};$
- Negative component, defined by elements with a negative inner product  $\Delta_i^- = \Delta_i \odot$  $\mathbb{1}_{\{\nabla_{\theta} L_j(\theta_{\text{pre}}) \odot \Delta_i < 0\}}.$

Here,  $\odot$  denotes the Hadamard (element-wise) product, and  $\mathbb{1}_{\{p\}} \in \mathbb{R}^N$  is an indicator vector that 267 takes the value 1 in the dimension that p is true and 0 otherwise. 268

We illustrate the impact of the three components within the loss landscape in Figure 2 (b). It is 269 evident that the positive component leads to an increase in loss, as the model moves in the gradient 270 ascent direction. The orthogonal component results in relatively smooth changes in the loss. In-271 terestingly, while the negative component initially follows the descent direction of the loss, it 272 overshoots the local optimum, ultimately leading to an increase in loss. As a result, the total loss 273 across all tasks shown in Figure 2 (c) highlights that the orthogonal component is more beneficial 274 for knowledge fusion than either the positive or negative components.

#### 5 TASK ARITHMETIC IN THE TRUST REGION

278 The above observations are reasonable since neural network parameters, particularly those in pre-279 trained foundation models, often exhibit high redundancy (Dalvi et al., 2020; Chen et al., 2022b). Additionally, task-specific knowledge is often low-rank Hu et al. (2022), i.e., only a few parameter 281 directions are critical for learning the task. In order to identify a small set of critical parameters that 282 should not be altered during model merging and alleviate knowledge conflicts, we propose defining 283 the following trust region:

**Definition 2** (Trust Region for Knowledge Conflict). Given a pre-trained model  $\theta_{pre}$ , the trust region specific in the dimension space is defined as follows:

$$\mathcal{TR} := \left\{ n \bigg| \sum_{i \neq j} \big| \nabla_{\theta} L_j(\theta_{\text{pre}})[n] \cdot \Delta_i[n] \big| < \epsilon \right\},\tag{5}$$

where  $n \leq N$  indexes the dimensions of the parameter space,  $\epsilon$  represents the sensitivity threshold, and any dimension exceeding this threshold will be excluded from the trust region and not permitted to merge.

293 Dimensions outside the trust region (corresponding components of task vector that are collinear 294 with the gradient direction, regardless of whether the directions are aligned or opposite) are likely to 295 cause knowledge conflicts. Conversely, when  $\nabla_{\theta} L_i(\theta_{\text{pre}})$  and  $\Delta_i$  are orthogonal, their projections 296 minimally interfere with each other, thereby reducing knowledge conflict. 297

We are now ready to present the TATR method. TATR mitigates knowledge conflict by restricting 298 merging within the trust region, involving the following three key steps. 299

**Calculating task-specific gradients.** The first step involves computing the gradient for each task. 300 Since accessing the full training data for each task is often impractical, we approximate the gradient 301 using an exemplar set for each task, denoted as  $\{S_1, \ldots, S_K\}$ . For each task, the absolute gradient 302 of the loss function  $L_k(.)$  (cross-entropy loss in our experiments) is computed as follows: 303

$$\nabla_{\theta} L_k\left(\theta_{\text{pre}}\right) \approx \mathbb{E}_{x_k \in S_k} \left| \nabla_{\theta} L_k\left(x_k; \theta_{\text{pre}}\right) \right|. \tag{6}$$

305 Notably, we place the expectation outside the absolute value operation, drawing inspiration from the 306 Fisher Information Matrix (Wasserman, 2013). This design captures absolute gradients that reflect 307 the average variation of parameters, facilitating the measurement of knowledge conflict across every 308 exemplar. Additionally, the exemplar size can be remarkably small. Our empirical results in Figure 3 (a) show that even in a one-shot setting, we achieve a competitive average accuracy of 72.3%, which 310 is close to the highest accuracy of 72.8% obtained with 16 samples. Similar results are observed 311 when TATR is integrated into AdaMerging Yang et al. (2024b).

312 We also propose a zero-shot version, where the task vector is used to estimate the gradient. Although 313 there may be estimation errors, this approach still offers performance improvements for TA-based 314 methods in most scenarios: 315

$$\left|\nabla_{\theta} L_k\left(\theta_{\rm pre}\right)\right| \approx \left|\Delta_k\right|.\tag{7}$$

Establishing the trust region. Next, we aim to identify the trust region with minimal knowledge 317 conflict, with a key requirement being the determination of the sensitivity threshold  $\epsilon$ . However, 318 manually specifying the exact value of  $\epsilon$  becomes complex and tedious. Therefore, we employ a 319 ranking method to infer  $\epsilon$ . To achieve this, we derive the sensitivity of each dimension that may 320 cause knowledge conflict, based on Definition 2: 321

322

316

275 276

277

284

285

286 287

289

290

291

292

304

323

 $\Omega^{\text{Trust}} = \sum_{i \neq j} \left| \nabla_{\theta} L_{j} \left( \theta_{\text{pre}} \right) \odot \Delta_{i} \right| = \sum_{i \neq j} \left| \nabla_{\theta} L_{j} \left( \theta_{\text{pre}} \right) \right| \odot \left| \Delta_{i} \right|.$ (8) Algorithm 1: The model merging process of TATR **Input:** Pre-trained model  $\theta_{\text{pre}}$ ; Task vectors  $\{\Delta_1, \dots, \Delta_K\}$ ; Exemplar-set  $\{S_1, \dots, S_K\}$ **Output:** Merged model  $\theta_{\text{TATR}}$ 1 // Deriving gradients for each task **2** for k = 1, ..., K do  $\mathbf{3} \mid G_k = \mathbb{E}_{x_k \in S_k} \left| \nabla_{\theta} L_k \left( x_k; \theta_{\text{pre}} \right) \right|$ 4 // Establishing the trust region  $\epsilon = \text{proportion\_selection}(\Omega^{\text{Trust}}, \tau)$  $\tau \ \mathcal{TR} = \{n \mid \Omega^{\mathrm{Trust}}[n] < \epsilon\}$ 8 // Merging  $\theta_{\text{TATR}} = \theta_{\text{pre}} + \lambda \sum_{k} \Delta_k \odot \mathbb{1}_{\{n \in \mathcal{TR}\}}$ 10 return  $\theta_{\text{TATR}}$ 

 Next, the sensitivity threshold  $\epsilon$  of the trust region is determined through a proportional selection operation:

$$\epsilon = \text{proportion\_selection}(\Omega^{\text{Trust}}, \tau).$$
(9)

In this process,  $\Omega^{\text{Trust}}$  is sorted in descending order, and the values corresponding to the predefined ratio  $\tau$  are selected as the sensitivity threshold  $\epsilon$ . Based on  $\epsilon$ , we are able to establish the trust region  $\mathcal{TR}$  according to Definition 2.

**Merging the task vectors.** The final step involves merging the task vectors using TA, where the merging occurs within the dimensions confined to the trust region:

$$\theta_{\text{TATR}} = \theta_{\text{pre}} + \lambda \sum_{k} \Delta_k \odot \mathbb{1}_{\{n \in \mathcal{TR}\}},\tag{10}$$

where  $\mathbb{1}_{\{n \in \mathcal{TR}\}} \in \mathbb{R}^N$  is an indicator vector whose value is 1 at index *n* if *n* belongs to the trust region and 0 otherwise. The detailed workings of TATR are outlined in Algorithm 1. The entire merging process does not rely on any additional training process.

Moreover, the techniques introduced in TATR selectively limit the merging process to a subset of
 model parameters, allowing it to function as a plug-and-play module that seamlessly integrates with
 a wide range of TA-based approaches, such as:

• **Ties-Merging & TATR:** Ties-Merging (Yadav et al., 2023) partially reduces knowledge conflicts by pruning low-magnitude parameters and aligning the signs of task vectors. However, this approach overlooks conflicts that may arise from high-magnitude parameters. This bias can lead to knowledge conflicts, where some tasks dominate the model's behavior. The combination of TATR with Ties-Merging refines the process, as shown in the following formula:

$$\theta_{\text{Ties+TATR}} = \theta_{\text{pre}} + \lambda \sum_{k} \Phi(\Delta_k) \odot \mathbb{1}_{\{n \in \mathcal{TR}\}}, \qquad (11)$$

where  $\Phi(.)$  indicates the TrIm, Elect Sign, and Merge operation of Ties-Merging.

 AdaMerging & TATR: AdaMerging (Yang et al., 2024b) adaptively learns merging coefficients but does not inherently resolve knowledge conflicts between task vectors. This can lead to interference during coefficient learning, especially when tasks require opposing parameter adaptations. TATR addresses this by pre-filtering task vectors to retain only those components within the trust region, ensuring that AdaMerging operates in a conflictreduced parameter space:

$$\theta_{\text{Ada+TATR}} = \theta_{\text{pre}} + \sum_{k} \lambda_k \Delta_k \odot \mathbb{1}_{\{n \in \mathcal{TR}\}},$$
(12)

where  $\lambda_1, \ldots, \lambda_K$  represent the learnable coefficients for AdaMerging.

• **Surgery & TATR:** Similarly, Surgery (Yang et al., 2024a) introduces additional modules to align task-specific features during merging. TATR complements Surgery by pre-selecting components of task vectors that reside in the trust region. The integrated approach is formalized as:

$$\theta_{\text{Surgery+TATR}} = \left\{ \theta_{\text{surgery}}, \theta_{\text{pre}} + \lambda \sum_{k} \Delta_{k} \odot \mathbb{1}_{\{n \in \mathcal{TR}\}} \right\},$$
(13)

where  $\theta_{surgery}$  denotes the additional parameters introduce by the Surgery module.

## 6 EXPERIMENTS

6.1 Settings

**Datasets.** Following prior works (Ilharco et al., 2023b; Yadav et al., 2023; Yang et al., 2024b;a), we perform model merging on the following eight datasets: SUN397 (Xiao et al., 2016), Cars (Krause et al., 2013), RESISC45 (Cheng et al., 2017), EuroSAT (Helber et al., 2019), SVHN (Netzer et al., 2011), GTSRB (Stallkamp et al., 2011), MNIST (LeCun & Cortes, 2010), DTD (Cimpoi et al., 2014).

395 **Baselines.** We compare our approach against a diverse set of methods, categorized into basic base-396 line methods, test-time training-based model merging methods, and training-free model merging 397 methods. Basic baseline methods include the Pre-trained model, Individual task model, and the Tra-398 ditional Multi-Task Learning model. For test-time training-based methods, we provide AdaMerg-399 ing, AdaMerging++ (Yang et al., 2024b), and Surgery (Yang et al., 2024a). Among the training-free 400 methods, we consider the simple Weight Average, Fisher Merging (Matena & Raffel, 2022), Reg-401 Mean (Jin et al., 2023), Task Arithmetic (Ilharco et al., 2023b), and Ties-Merging (Yadav et al., 402 2023).

Implementation details. Our implementation strictly follows task arithmetic (Ilharco et al., 2023b) and AdaMerging (Yang et al., 2024b). We apply the ViT-B/32 and ViT-L/14 in CLIP (Radford et al., 2021) as the pre-trained model. Task vectors are derived from task arithmetic (Ilharco et al., 2023b) which is fine-tuned on each specific dataset. We report the accuracy of each task after merging the models, along with the average accuracy (i.e., Avg ACC). The hyper-parameter  $\tau$  is tuned within the range [0.1%, 0.2%, 0.5%, 1.0%, 2.0%, 5.0%], while the size of the exemplar set is fixed at 128. Additional implementation details can be found in our supplementary code.

410 411

378

379

380

382

384

385 386

387

388

389 390

391

392

393

394

6.2 PERFORMANCE COMPARISON

The performance of all baselines using the ViT-B/32 and ViT-L/14 architectures is presented in
Table 1 and Table 2, respectively. We report the performance metrics for each task after merging, as
well as the overall average performance.

As illustrated in the tables, the pre-trained model exhibits the lowest performance across all methods,
 due to the absence of task-specific supervision. In contrast, the Individual models achieve the highest
 performance, as they are exclusively trained for each specific task, which thus represents the upper bound performance for model merging. Traditional MTL encounters knowledge conflict issues,
 resulting in slightly lower performance compared to the Individual models.

Among the model merging methods, the simplest Weight Averaging suffers significant knowledge
 conflicts, resulting in a worse performance. Fisher Merging and RegMean improve Weight Averag ing by incorporating parameter importance weight into the averaging process. TA and its enhanced
 version, Ties-merging, demonstrate substantial performance improvements by better balancing the
 pre-trained and task-specific knowledge. Additionally, owing to the additional training process,
 test-time training-based models (AdaMerging and Surgery) generally outperform the training-free
 methods.

Our proposed TATR method belongs to the training-free model merging approach. As the techniques
 of TATR are orthogonal to existing model merging methods, we also report performance when
 TATR is plugged into strong baselines. The experimental results demonstrate that TATR consistently
 enhances all TA-based methods. When incorporated into task arithmetic, both TATR and its zero shot version lead to significant performance improvements, increasing average accuracy by 3.7%

and 1.5% on ViT-B/32, and by 0.8% and 0.1% on ViT-L/14, respectively. The best results are obtained when TATR is combined with layer-wise AdaMerging++, achieving an average accuracy of 82.5% on ViT-B/32 and 91.5% on ViT-L/14. 

Table 1: Multi-task performance when merging ViT-B/32 models on eight tasks. The column of "# Best" indicates the number of datasets on which the proposed method achieved the best performance, and the best and second-best performance are highlighted with **bold** and <u>underline</u>.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	# Best	Avg Acc			
Basic baseline methods													
Pre-trained	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			
	Test-t	ime trai	ning based me	thods									
TW AdaMerging	58.0	53.2	68.8	85.7	81.1	84.4	92.4	44.8	0	71.1			
TW AdaMerging++	60.8	56.9	73.1	83.4	87.3	82.4	95.7	50.1	0	73.7			
LW AdaMerging	64.5	68.1	79.2	93.8	87.0	91.9	97.5	59.1	1	80.1			
LW AdaMerging++	66.6	68.3	82.2	94.2	89.6	89.0	98.3	60.6	0	81.1			
Surgery Merging	63.8	59.9	83.3	97.9	87.0	87.0	98.6	69.4	1	80.9			
LW AdaMerging++ & TATR zero-shot (Ours)	72.0	70.8	81.5	88.9	84.9	84.2	99.3	66.7	3	81.0			
LW AdaMerging++ & TATR (Ours)	<u>69.8</u>	<u>70.3</u>	83.7	93.7	90.0	<u>90.2</u>	98.3	63.7	1	82.5			
Surgery & TATR zero-shot (Ours)	64.2	60.4	82.7	96.9	86.4	86.5	98.5	68.7	0	80.5			
Surgery & TATR (Ours)	67.1	62.2	87.1	<u>97.4</u>	87.3	88.5	<u>98.7</u>	70.9	2	82.4			
		Training	g-free methods										
Weight Averaging	65.3	63.4	71.4	71.7	64.2	52.8	87.5	50.1	0	65.8			
Fisher Merging	68.6	69.2	70.7	66.4	72.9	51.1	87.9	59.9	3	68.3			
RegMean	65.3	63.5	75.6	78.6	78.1	67.4	93.7	52.0	0	71.8			
Task Arithmetic	55.2	54.9	66.7	78.9	80.2	69.7	97.3	50.4	0	69.1			
Ties-Merging	59.8	58.6	70.7	79.7	86.2	72.1	98.3	54.2	2	72.4			
TATR zero-shot (Ours)	59.0	56.6	69.2	80.2	79.0	70.5	97.0	53.5	0	70.6			
TATR (Ours)	62.7	59.3	72.3	82.3	80.5	72.6	97.0	55.4	1	72.8			
Ties-Merging & TATR zero-shot (Ours)	64.9	64.2	74.7	76.4	81.2	69.3	96.5	54.3	1	72.7			
Ties-Merging & TATR (Ours)	<u>66.3</u>	<u>65.9</u>	75.9	79.4	79.9	68.1	96.2	54.8	1	73.3			

## Table 2: Multi-task performance when merging ViT-L/14 models on eight tasks.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	# Best	Avg Acc			
Basic baseline methods													
Pre-trained	66.8	77.7	71.0	59.9	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			
Test-time training based methods													
AdaMerging	79.0	90.3	90.8	96.2	93.4	98.0	99.0	79.9	2	90.8			
AdaMerging++	79.4	90.3	91.6	97.4	93.4	97.5	99.0	79.2	1	91.0			
Surgery Merging	75.7	84.4	93.1	98.8	91.3	93.4	99.1	76.1	1	89.0			
AdaMerging & TATR zero-shot (Ours)	80.7	<u>95.3</u>	<u>95.0</u>	94.9	84.7	92.4	99.8	86.0	1	91.1			
AdaMerging++ & TATR (Ours)	81.6	95.9	95.8	95.5	83.2	92.6	<u>99.7</u>	87.5	4	91.5			
Surgery & TATR zero-shot (Ours)	75.6	85.1	93.8	98.5	91.0	93.1	99.2	76.3	0	89.1			
Surgery & TATR (Ours)	76.3	85.8	93.8	98.8	91.4	93.0	99.2	77.9	1	89.5			
		Trair	ning-free meth	ods									
Weight Averaging	72.1	81.6	82.6	91.9	78.2	70.7	97.1	62.8	0	79.6			
Fisher Merging	69.2	88.6	87.5	93.5	80.6	74.8	93.3	70.0	1	82.2			
RegMean	73.3	81.8	86.1	97.0	88.0	84.2	98.5	60.8	1	83.7			
Task Arithmetic	73.9	82.1	86.6	94.1	87.9	86.7	98.9	65.6	0	84.5			
Ties-Merging	76.5	85.0	89.3	95.7	90.3	83.3	99.0	68.8	3	86.0			
TATR zero-shot (Ours)	74.3	81.5	86.6	92.7	88.6	88.1	99.1	66.0	0	84.6			
TATR (Ours)	74.6	83.7	87.6	93.7	88.6	88.1	99.0	66.8	0	85.3			
Ties-Merging & TATR zero-shot (Ours)	75.8	85.3	89.2	94.7	89.1	87.1	99.0	68.6	0	86.1			
Ties-Merging & TATR (Ours)	76.3	85.3	88.8	94.4	90.8	88.7	99.2	68.8	4	86.5			

## 

#### SENSITIVITY ANALYSIS OF HYPERPARAMETERS 6.3

This section presents an analysis of the model's sensitivity to two hyperparameters: the number of exemplar samples and the proportion  $\tau$  in Eq. (9). As shown in Figure 3, the performance of TATR remains stable with respect to both hyperparameters. Furthermore, Figure 3 (a) demonstrates that even in a one-shot setting, TATR achieves a competitive average accuracy of 72.3%, evidently outperforming Task Arithmetic (0 exemplars in Figure 3 (a)) and comparable to the highest accuracy of 72.8%. Similarly, experiments plugged into AdaMerging also support this phenomenon, where the one-shot scenario achieves an average accuracy of 82.3%, nearly matching the peak accuracy of 82.5%. Additionally, Figure 3 (b) suggests that excluding a small proportion of parameters (less than 1%) is sufficient to alleviate the knowledge conflicts.



Figure 3: Average accuracy (%) of TATR on eight tasks versus the number of exemplars (a) and  $\tau$ (b).



Figure 4: The average sensitivity of each dataset to task vectors across layers.

### ANALYSIS OF SENSITIVITY $\Omega^{\text{Trust}}$ for Knowledge Conflict 6.4

Figure 4 illustrates the average sensitivity of each dataset to task vectors across different layers. Three key characteristics can be observed. Firstly, the shallow layers exhibit greater sensitivity than other layers. Shallow layers typically encode task-generalized knowledge, and the increased sensi-tivity highlights the importance of preserving this information in the TATR method. Secondly, the sensitivity exhibits periodic variations across layers, with bias layers generally exhibiting higher sensitivity than weight layers. This trend is reasonable, as bias layers have a more pronounced impact on network outputs, making them more susceptible to knowledge conflicts. Lastly, datasets com-prising digit data (e.g., SVHN and MNIST) show relatively lower sensitivity to knowledge conflicts, which can be attributed to their significant domain differences from other real-world datasets. 

### CONCLUSION

In this paper, we delve deep into the critical challenge of knowledge conflict in multi-task model merging with a focus on task arithmetic. We began by formalizing the concept of knowledge conflict as the degradation in model performance caused by the interference between task vectors. Our analysis and empirical findings suggest that components of task vectors orthogonal to the gradient direction exhibit minimal knowledge conflict. This insight motivates us to define a trust region based on orthogonality and propose Task Arithmetic in the Trust Region (TATR). Extensive experiments across eight diverse datasets demonstrate that TATR effectively mitigates the knowledge conflict, enhancing the overall multi-task performance of task arithmetic-based methods.

### **REPRODUCIBILITY STATEMENT**

We have included the complete source code in the supplementary materials, and will release the full codebase as open source upon publication. All datasets and settings are documented for clarity.

# 540 REFERENCES

552

565

566

567

568 569

570

571

572

584

585

586

588

589

- 542 Devansh Arpit, Huan Wang, Yingbo Zhou, and Caiming Xiong. Ensemble of averages: Improving
   543 model selection and boosting performance in domain generalization. In <u>Advances in Neural</u>
   544 Information Processing Systems, volume 35, pp. 8265–8277, 2022.
- Junbum Cha, Sanghyuk Chun, Kyungjae Lee, Han-Cheol Cho, Seunghyun Park, Yunsung Lee, and
  Sungrae Park. Swad: Domain generalization by seeking flat minima. In <u>Advances in Neural</u> Information Processing Systems, volume 34, pp. 22405–22418, 2021.
- Shuxiao Chen, Koby Crammer, Hangfeng He, Dan Roth, and Weijie J Su. Weighted training for cross-task learning. In <u>International Conference on Learning Representations</u>, 2022a. URL https://openreview.net/forum?id=ltM1RMZntpu.
- Tianlong Chen, Zhenyu Zhang, Yu Cheng, Ahmed Awadallah, and Zhangyang Wang. The principle of diversity: Training stronger vision transformers calls for reducing all levels of redundancy. In
   Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 12020–12030, June 2022b.
- Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, and Andrew Rabinovich. GradNorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In <u>International Conference</u> on Machine Learning, volume 80, pp. 794–803, 10–15 Jul 2018.
- Zhao Chen, Jiquan Ngiam, Yanping Huang, Thang Luong, Henrik Kretzschmar, Yuning Chai, and
   Dragomir Anguelov. Just pick a sign: Optimizing deep multitask models with gradient sign
   dropout. In <u>Advances in Neural Information Processing Systems</u>, volume 33, pp. 2039–2050,
   2020.
  - Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. <u>Proceedings of the IEEE</u>, 105(10):1865–1883, 2017. doi: 10.1109/JPROC.2017.2675998.
  - Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In <u>IEEE Conference on Computer Vision and Pattern Recognition</u>, June 2014.
- Fahim Dalvi, Hassan Sajjad, Nadir Durrani, and Yonatan Belinkov. Analyzing redundancy in pretrained transformer models. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 4908–4926, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.398. URL https://aclanthology.org/2020. emnlp-main.398.
- 579 Ke Ding, Xin Dong, Yong He, Lei Cheng, Chilin Fu, Zhaoxin Huan, Hai Li, Tan Yan, Liang Zhang, Xiaolu Zhang, and Linjian Mo. Mssm: A multiple-level sparse sharing model for efficient multi-task learning. In <u>International ACM SIGIR Conference on Research and Development</u>
  582 <u>in Information Retrieval</u>, pp. 2237–2241, 2021. URL https://doi.org/10.1145/ 3404835.3463022.
  - Ronald Aylmer Fisher. Theory of statistical estimation. In <u>Mathematical proceedings of the</u> Cambridge philosophical society, volume 22, pp. 700–725, 1925.
  - Yuan Gao, Haoping Bai, Zequn Jie, Jiayi Ma, Kui Jia, and Wei Liu. Mtl-nas: Task-agnostic neural architecture search towards general-purpose multi-task learning. In <u>IEEE Conference on Computer</u> Vision and Pattern Recognition, June 2020.
- Hussein Hazimeh, Zhe Zhao, Aakanksha Chowdhery, Maheswaran Sathiamoorthy, Yihua Chen,
   Rahul Mazumder, Lichan Hong, and Ed Chi. Dselect-k: Differentiable selection in the mixture of
   experts with applications to multi-task learning. In <u>Advances in Neural Information Processing</u>
   Systems, volume 34, pp. 29335–29347, 2021.

594 595 596 597	Yun He, Xue Feng, Cheng Cheng, Geng Ji, Yunsong Guo, and James Caverlee. Metabalance: Improving multi-task recommendations via adapting gradient magnitudes of auxiliary tasks. In <u>International ACM Web Conference</u> , pp. 2205–2215, 2022. doi: 10.1145/3485447.3512093. URL https://doi.org/10.1145/3485447.3512093.
598 599 600 601 602	Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. <u>IEEE Journal of Selected</u> <u>Topics in Applied Earth Observations and Remote Sensing</u> , 12(7):2217–2226, 2019. doi: 10. 1109/JSTARS.2019.2918242.
603 604 605 606	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In <u>International</u> <u>Conference on Learning Representations</u> , 2022. URL https://openreview.net/forum? id=nZeVKeeFYf9.
607 608 609	Yuzheng Hu, Ruicheng Xian, Qilong Wu, Qiuling Fan, Lang Yin, and Han Zhao. Revisiting scalar- ization in multi-task learning: A theoretical perspective. In <u>Advances in Neural Information</u> <u>Processing Systems</u> , volume 36, pp. 48510–48533, 2023.
610 611 612 613	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In <u>International Conference on Learning Representations</u> , 2023a.
614 615 616	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In <u>International Conference on Learning</u> <u>Representations</u> , 2023b. URL https://openreview.net/forum?id=6t0Kwf8-jrj.
617 618 619	Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. In <u>Conference on Uncertainty</u> <u>in Artificial Intelligence</u> , pp. 876–885, 2018.
620 621 622 623	Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and gener- alization in neural networks. In <u>Advances in Neural Information Processing Systems</u> , volume 31, 2018.
624 625 626	Adrián Javaloy and Isabel Valera. Rotograd: Gradient homogenization in multitask learning. In <u>International Conference on Learning Representations</u> , 2022. URL https://openreview.net/forum?id=T8wHz4rnuGL.
627 628 629	Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. Dataless knowledge fusion by merging weights of language models. In <u>International Conference on Learning Representations</u> , 2023. URL https://openreview.net/forum?id=FCnohuR6AnM.
630 631 632 633	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. <u>arXiv preprint arXiv:2001.08361</u> , 2020.
634 635	Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In <u>IEEE International Conference on Computer Vision Workshops</u> , June 2013.
636 637 638	Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL http://yann.lecun.com/exdb/mnist/.
639 640 641	Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. In <u>Advances in Neural Information Processing Systems</u> , volume 34, pp. 18878–18890, 2021a.
642 643 644	Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. In <u>Advances in Neural Information Processing Systems</u> , volume 34, pp. 18878–18890, 2021b.
646 647	Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. Adversarial multi-task learning for text classifi- cation. In <u>Association for Computational Linguistics</u> , pp. 1–10, July 2017. doi: 10.18653/v1/ P17-1001. URL https://aclanthology.org/P17-1001.

648 649	Shikun Liu, Edward Johns, and Andrew J. Davison. End-to-end multi-task learning with attention. In IEEE Conference on Computer Vision and Pattern Recognition, June 2019.
651 652 653	Shikun Liu, Stephen James, Andrew Davison, and Edward Johns. Auto-lambda: Disentangling dynamic task relationships. <u>Transactions on Machine Learning Research</u> , 2022. ISSN 2835- 8856. URL https://openreview.net/forum?id=KKeCMim5VN.
654 655 656 657	Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H. Chi. Modeling task relation- ships in multi-task learning with multi-gate mixture-of-experts. In <u>International ACM SIGKDD</u> <u>International Conference on Knowledge Discovery &amp; Data Mining</u> , pp. 1930–1939, 2018. doi: 10.1145/3219819.3220007. URL https://doi.org/10.1145/3219819.3220007.
658 659 660	Michael S Matena and Colin A Raffel. Merging models with fisher-weighted averaging. In <u>Advances</u> <u>in Neural Information Processing Systems</u> , volume 35, pp. 17703–17716, 2022.
661 662 663	Ze Meng, Xin Yao, and Lifeng Sun. Multi-task distillation: Towards mitigating the negative transfer in multi-task learning. In IEEE International Conference Image Processing, pp. 389–393, 2021. doi: 10.1109/ICIP42928.2021.9506618.
664 665 666 667	Aviv Navon, Aviv Shamsian, Idan Achituve, Haggai Maron, Kenji Kawaguchi, Gal Chechik, and Ethan Fetaya. Multi-task learning as a bargaining game. In <u>International Conference on Machine</u> <u>Learning</u> , volume 162, pp. 16428–16446, 17–23 Jul 2022. URL https://proceedings. mlr.press/v162/navon22a.html.
668 669 670 671	Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In <u>NIPS workshop on deep</u> <u>learning and unsupervised feature learning</u> , volume 2011, pp. 4, 2011.
672 673 674	Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. Task arithmetic in the tangent space: Improved editing of pre-trained models. In <u>Advances in Neural Information Processing</u> <u>Systems</u> , volume 36, pp. 66727–66754, 2023.
675 676 677 678 679	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In International Conference on Machine Learning, volume 139, pp. 8748–8763, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/radford21a.html.
680 681 682	Sebastian Ruder. An overview of gradient descent optimization algorithms. <u>arXiv:1609.04747</u> , 2017. URL https://arxiv.org/abs/1609.04747.
683 684	Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In <u>Advances</u> <u>in Neural Information Processing Systems</u> , volume 31. Curran Associates, Inc., 2018.
685 686 687 688	Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. The german traffic sign recog- nition benchmark: A multi-class classification competition. In <u>International Joint Conference on</u> <u>Neural Networks</u> , pp. 1453–1460, 2011.
689 690 691	Wenju Sun, Qingyong Li, Wen Wang, and Yangli-ao Geng. Towards plastic and stable exemplar- free incremental learning: A dual-learner framework with cumulative parameter averaging. <u>arXiv</u> preprint arXiv:2310.18639, 2023.
692 693 694	Ximeng Sun, Rameswar Panda, Rogerio Feris, and Kate Saenko. Adashare: Learning what to share for efficient deep multi-task learning. In <u>Advances in Neural Information Processing Systems</u> , volume 33, pp. 8728–8740, 2020.
695 696 697 698	Anke Tang, Li Shen, Yong Luo, Nan Yin, Lefei Zhang, and Dacheng Tao. Merging multi-task models via weight-ensembling mixture of experts. In <u>International Conference on Machine Learning</u> , 2024. URL https://openreview.net/forum?id=nLRKn074RB.
699 700 701	Hongyan Tang, Junning Liu, Ming Zhao, and Xudong Gong. Progressive layered extraction (ple): A novel multi-task learning (mtl) model for personalized recommendations. In <u>ACM Conference</u> on <u>Recommender Systems</u> , pp. 269–278, 2020. doi: 10.1145/3383313.3412236. URL https: //doi.org/10.1145/3383313.3412236.

- Jingyao Wang, Yi Ren, Zeen Song, Jianqi Zhang, Changwen Zheng, and Wenwen Qiang. Hacking task confounder in meta-learning. In <u>IJCAI</u>, 2023.
- Ke Wang, Nikolaos Dimitriadis, Guillermo Ortiz-Jimenez, François Fleuret, and Pascal Frossard. Localizing task information for improved model merging and compression. In <u>International</u> <u>Conference on Machine Learning</u>, 2024. URL https://openreview.net/forum?id= DWT9uiGjxT.
- <sup>709</sup> L Wasserman. All of statistics: A concise course in statistical inference, 2013.
- Xun Wu, Shaohan Huang, and Furu Wei. Mixture of loRA experts. In <u>International Conference</u> on <u>Learning Representations</u>, 2024. URL https://openreview.net/forum?id= uWvKBCYh4S.
- 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.
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. TIES-merging: Resolving interference when merging models. In <u>Advances in Neural Information Processing</u> Systems, 2023.
- Chenxiao Yang, Junwei Pan, Xiaofeng Gao, Tingyu Jiang, Dapeng Liu, and Guihai Chen. Cross-task knowledge distillation in multi-task recommendation. In <u>AAAI Conference on Artificial</u> Intelligence, volume 36, pp. 4318–4326, Jun. 2022. doi: 10.1609/aaai.v36i4.20352.
- Enneng Yang, Junwei Pan, Ximei Wang, Haibin Yu, Li Shen, Xihua Chen, Lei Xiao, Jie Jiang, and Guibing Guo. Adatask: A task-aware adaptive learning rate approach to multi-task learning.
   <u>AAAI Conference on Artificial Intelligence</u>, 37(9):10745–10753, Jun. 2023. URL https://ojs.aaai.org/index.php/AAAI/article/view/26275.
- Finneng Yang, Li Shen, Zhenyi Wang, Guibing Guo, Xiaojun Chen, Xingwei Wang, and Dacheng Tao. Representation surgery for multi-task model merging. In <u>International Conference on</u> <u>Machine Learning</u>, 2024a. URL https://openreview.net/forum?id=Sbl2keQEML.
- Finneng 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, 2024b. URL https://openreview.net/forum?id=
   nZP6NgD3QY.
- Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.
   Gradient surgery for multi-task learning. In <u>Advances in Neural Information Processing Systems</u>, volume 33, pp. 5824–5836, 2020.
  - Jinghan Zhang, shiqi chen, Junteng Liu, and Junxian He. Composing parameter-efficient modules with arithmetic operation. In Advances in Neural Information Processing Systems, volume 36, pp. 12589–12610, 2023a.
- Wen Zhang, Lingfei Deng, Lei Zhang, and Dongrui Wu. A survey on negative transfer. <u>IEEE/CAA</u>
   Journal of Automatica Sinica, 10(2):305–329, 2023b. doi: 10.1109/JAS.2022.106004.
  - Yu Zhang and Qiang Yang. A survey on multi-task learning. IEEE Transactions on Knowledge and Data Engineering, 34(12):5586–5609, 2022. doi: 10.1109/TKDE.2021.3070203.
- 752 753 A Experiment Details

738

742

743

744

745

748

749

750 751

754

755 This section provides details of experiments, including the description of the experimental environment, datasets, and baselines.

## A.1 ENVIRONMENT

All experiments detailed in our manuscript and appendix were conducted on a workstation running
Ubuntu 16.04, equipped with 18 Intel Xeon 2.60GHz CPUs, 256 GB of memory, and 6 NVIDIA
RTX3090 GPUs. Python 3.8 was used to implement all the methods.

A.2 DATASETS

761 762

763

769

770

771 772

773

774

775

776

777

779

780

781

782

783

784

785

786

787 788

797

798

799

800

801 802

804

- Our experiments strictly follow Task Arithmetic (Ilharco et al., 2023b) and AdaMerging (Yang et al., 2024b), utilizing eight widely-used image classification datasets. The information of these datasets is described as follows:
  - **SUN397** (Xiao et al., 2016): A scene classification dataset containing 108,754 images across 397 classes. Each class includes at least 100 images.
    - **Stanford Cars** (**Cars**) (Krause et al., 2013): A car classification dataset featuring 16,185 images of 196 car categories. The dataset is evenly split between training and test sets.
  - **RESISC45** (Cheng et al., 2017): A remote sensing image classification dataset comprising 31,500 images across 45 scene categories, with approximately 700 images per class.
  - **EuroSAT** (Helber et al., 2019): A satellite image classification dataset consisting of 27,000 labeled and geo-referenced images distributed among 10 categories.
  - SVHN (Netzer et al., 2011): A real-world digit classification dataset derived from house numbers in Google Street View images. It includes 10 classes, with a training set of 73,257 images, a test set of 26,032 images, and an additional 531,131 samples available for extended training.
  - **GTSRB** (Stallkamp et al., 2011): A traffic sign classification dataset comprising more than 50,000 images across 43 traffic sign categories.
    - **MNIST** (LeCun & Cortes, 2010): A well-known benchmark for handwritten digit classification, containing 60,000 training images and 10,000 test images, evenly distributed among 10 classes of digit numbers.
    - **DTD** (Cimpoi et al., 2014): A texture classification dataset consisting of 5,640 images distributed across 47 texture classes, with approximately 120 images per class.
- 789 A.3 BASELINES.

In our experiments, we compare our methods with several baseline approaches, which are grouped into four categories: basic baseline methods, test-time training-based methods, training-free methods, and our proposed methods. The details of these methods are as follows:

- i) Basic baseline methods:
  - **Pre-trained** directly employs a pre-trained model to predict across multiple tasks. Since it does not incorporate any downstream task-specific information during model training, its performance on downstream tasks is typically suboptimal.
    - **Individual**. In this approach, an independent fine-tuned model is used for each task. While it avoids interference between tasks, it cannot perform multiple tasks simultaneously. It serves as a reference *upper bound* for model merging approaches.
      - **Traditional MTL** aggregates the original training data from all tasks to train a single multitask model.
- ii) Test-time training-based methods:
  - AdaMerging (Yang et al., 2024b) leverages an unlabeled test set to adaptively learn the merging coefficients at either a layer-wise or task-wise level in Task Arithmetic.
- AdaMerging++ (Yang et al., 2024b) an enhanced version of AdaMerging, integrates the principles of Ties-Merging (Yadav et al., 2023).

810		• Surgery (Yang et al. 2024a) introduces a feature transformation module, trained to align
811		features during the merging process. In this work, we adopt the basic version of Surgery
812		combined with task arithmetic for evaluation
813	•••	
814	m) 1	raining-free methods:
815		• Weight Averaging directly everyges model perspecters from multiple tasks into a single
816		model enabling multi-task learning without additional training
817		$\mathbf{F} = \mathbf{M} + $
818		• Fisher Merging (Matena & Raffel, 2022) leverages the Fisher information matrix to assess
819		parameter importance, merging model parameters based on this importance.
820 821		• <b>RegMean</b> (Jin et al., 2023) refines weight matrices by adjusting and linearly combining rows utilizing statistical information derived from the training data
822		
823		• Task Arithmetic (Ilharco et al., 2023b) introduces the concept of a "task vector," defined
824		as the difference between fine-tuned model parameters and pre-trained model parameters.
825		multi-task learning
826		
827		• <b>Ties-Merging</b> (Yadav et al., 2023) eliminates unimportant parameters from the task vector
828		and resolves sign connicts among parameters, reducing interference during the linal task
829		vector merging process.
830	iv) O	Our methods:
831	., -	
832		• TATR. This is the core method introduced in our work, which applies task arithmetic within
833		the trust region.
834		• TATR zero-shot A zero-shot variant of TATR that utilizes task vectors to estimate the
835		gradient as described in Eq.(7). Other zero-shot variations follow a similar approach.
836		• Ties-Merging & TATR integrates TATR into the Ties-Merging framework by applying
837		TATR's mask on the task vectors after processing them with Ties-Merging.
838		• AdaMerging & TATR plugged TATR into AdaMerging, where TATR is applied prior to
839		training the AdaMerging coefficients.
841		• Surgery & TATR. Similarly, TATR is integrated into Surgery by applying it before training
842		the additional modules introduced by Surgery.
843		
844	В	VISUALIZATION OF LOSS LANDSCAPE
845		
846	<b>B</b> .1	METHODOLOGY FOR VISUALIZING THE LOSS LANDSCAPE
847	T. (1.	' and a set d'a dama data tao far 's all' 's data to ta tao a l'it 's data a
848	In the	is section, we outline the methodology for visualizing the loss landscape, which involves three
849	KCy S	steps.
850	Task	vector decomposition. To effectively visualize the loss landscape, we first decompose a task-
851	speci	fic vector into three essential components: a positive component (aligned with the gradient
852	direc	to the gradient). These component (opposed to the gradient), and an orthogonal component (orthog-
853	vecto	or and the gradient, as detailed in Section 5. This decomposition allows us to explore how dif-
854	feren	it aspects of the task vector interact with the gradient, each contributing uniquely to the overall
855	optin	nization behavior.
856	Corr	structing the 2D plane. We among these three services the 2D seculiant set of the
857	the p	surviving the 2D plane. We arrange these three components in a 2D coordinate system, with $a_{0}$ or $a_{1}$ and $a_{2}$ or $a_{1}$ and $a_{2}$ or $a_{2}$ and $a$
858	none	nt at $(1,0)$ Additionally we project the parameters of the pre-trained model onto this plane
859	using	g linear combinations of the three components. Although this projection is an approximation, as
860	the p	re-trained model's parameters may not lie perfectly within the plane defined by these vectors.
861	it pro	ovides sufficient insight into the interaction between task vectors and knowledge conflict.
862	-	

**Contour plot generation**. Finally, we sample points across the plane by selecting coordinates within the range of [-0.2, 1.2] for both axes at intervals of 0.1. For each sampled point, we adjust

the model's parameters accordingly and compute the corresponding loss value. These loss values are then used to generate a contour plot, providing a visual representation of the loss landscape.

## B.2 LOSS LANDSCAPE FOR EACH INDIVIDUAL TASK

In this section, we present the loss landscape for each individual task, along with the components of task vectors from the other seven tasks within the landscape. From Figure 5, we observe the same patterns as described in the manuscript. Specifically, the positive component tends to ascend along the gradient direction, while the negative component, despite aligning with the gradient descent direction, often overshoots local optima, leading to performance degradation. The orthogonal component, in general, shows little sensitivity to performance changes. These findings further support the generality of our conclusions and provide additional evidence for the effectiveness of the TATR method.

876 877 878

879

880

882

883 884 885

886 887

888

889

890

891

867

868

## B.3 LOSS LANDSCAPE FOR ALL TASKS

Furthermore, we visualize the overall loss landscape across all tasks, including the components of the task vectors. We present the loss landscape under different mask ratios. As shown in Figure 6, we observe similar patterns: both the positive and negative components negatively impact the model's overall multi-task performance, leading to knowledge conflicts. In contrast, the orthogonal component contributes to improving the model's multi-task capability.

## C ADDITIONAL EXPERIMENTS

## C.1 COMPARISON ON VIT-B/16

Table 3 presents the results of various model merging methods using the ViT-B/16 architecture. As we can see, TATR significantly improves the multi-task performance of Task Arithmetic, raising the average performance from 73.8% to 77.0%. Additionally, the zero-shot version also provides a certain degree of improvement, ultimately reaching 74.1%.

896 897

Table 3: Multi-task performance when merging ViT-B/16 models on eight tasks.

Method	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	Avg Acc
Pre-trained	63.8	64.6	65.7	54.5	52.0	43.3	51.7	45.1	55.0
Individual	81.8	86.8	96.9	99.7	97.8	99.1	99.7	82.0	92.9
Weight Averaging	67.7	70.0	75.3	79.5	74.9	60.1	94.4	43.8	70.7
Fisher Merging	68.5	69.9	75.2	80.4	73.2	61.2	94.5	50.7	71.7
RegMean	69.1	71.6	77.6	88.8	83.7	70.2	96.9	54.6	76.6
Task Arithmetic	61.1	65.9	74.0	76.2	88.0	73.9	98.4	53.0	73.8
Ties-Merging	69.1	72.5	80.5	84.0	85.0	71.5	98.1	54.9	77.0
TATR zero-shot (Ours)	60.5	63.4	73.0	78.8	88.4	75.8	98.4	54.6	74.1
TATR (Ours)	67.4	70.4	77.9	81.7	87.6	77.2	98.3	55.6	77.0

## C.2 GENERALIZATION COMPARISON

This section explores the generalization ability of TATR. Specifically, we merge models using task vectors from six tasks and evaluate their performance on two unseen tasks. We conduct two experiments: in the first, MNIST and EuroSAT are set as unseen tasks, while in the second, RESISC45 and SVHN are treated as unseen. The results in Table 4 show that TATR outperforms Task Arithmetic on the unseen datasets, with an average performance improvement of 0.8% and 1.3%, respectively. This improvement in generalization is attributed to TATR's ability to handle knowledge conflicts, ensuring that model updates move toward a more globally optimal direction.

913

## 914 C.3 ANALYSIS OF EXEMPLAR NUMBER 915

In this section, we further investigate the sensitivity of TATR to the number of exemplars. Table 5
 reports the merging performance with varying exemplar numbers. As shown, the zero-shot version consistently outperforms Task Arithmetic across all tasks, achieving an average performance



Figure 5: Loss landscape for each dataset and the components of the cumulative task vector from the remaining seven datasets.



Figure 6: Loss landscape for each datasets and the components of the cumulative task vector from the remaining seven datasets.

Table 4: Generalization results on two unseen tasks when merging ViT-B/16 models on six tasks.

Method	SUN397	Cars	RESISC45	DTD	SVHN	GTSRB	Avg Acc	MNIST	EuroSAT	Avg A
Task Arithmetic TATR (Ours)	63.3 66.0	62.4 64.1	75.1 77.9	57.8 60.1	84.6 83.9	80.4 81.8	70.6 72.3	77.2 77.2	46.2 47.7	61.7 62.5
Method	SUN397	Cars	GTSRB	EuroSAT	DTD	MNIST	Avg Acc	RESISC45	SVHN	Avg A
Task Arithmetic	64.0	64.0	75.2	87.7	57.0	95.7	73.9	52.3	44.9	51.
TATR (Ours)	66.5	65.2	76.8	87.9	59.5	95.6	75.3	54.7	50.0	52.

improvement of 1.7%. In the one-shot scenario, TATR significantly boosts performance, with an average increase of 3.4% per task, nearing optimal performance. As the number of exemplars in-creases, the performance improves across all tasks, obtaining the best performance at the number 16.

Table 5: Impact of the number of exemplars when merging ViT-B/32 models on eight tasks.

42	Method	Exemplar number	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	Avg Acc
10	Task Arithmetic	-	55.2	54.9	66.7	78.9	80.2	69.7	97.3	50.4	69.1
3	TATR zero-shot (Ours)	-	59.0	56.6	69.2	80.2	79.0	70.5	97.0	53.5	70.6
1	TATR (Ours)	1	62.0	59.0	71.6	81.8	80.3	72.4	96.9	54.7	72.3
	TATR (Ours)	2	62.3	59.2	71.6	81.5	80.5	72.4	97.0	55.4	72.5
	TATR (Ours)	4	62.3	59.3	71.8	82.4	80.5	72.7	97.0	55.1	72.6
	TATR (Ours)	8	62.6	59.3	72.2	82.3	80.1	72.6	97.0	55.3	72.7
	TATR (Ours)	16	62.7	59.5	72.3	82.4	80.4	72.6	97.0	55.3	72.8
	TATR (Ours)	32	62.7	59.5	72.4	82.4	80.4	72.5	97.0	55.4	72.8
	TATR (Ours)	64	62.7	59.3	72.3	82.5	80.4	72.7	97.0	55.4	72.8
3	TATR (Ours)	128	62.7	59.4	72.3	82.5	80.4	72.7	97.0	55.4	72.8

Table 6: Comparison with different sensitivities of TATR when merging ViT-B/32 models on eight tasks. 

1053											
1054	Method	Sensitivity	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD	Avg Acc
1054	Pre-trained	-	63.8	64.6	65.7	54.5	52.0	43.3	51.7	45.1	55.0
1055	TATR positive	$\frac{1}{K(K-1)} \sum_{i \neq j} \nabla_{\theta} L_j(\theta_{pre}) \odot \Delta_i$	60.0	54.3	44.8	9.3	18.9	14.3	17.1	39.4	32.3
1050	TATR negative	$-\frac{1}{K(K-1)}\sum_{i\neq j}\nabla_{\theta}L_{j}\left(\theta_{\text{pre}}\right)\odot\Delta_{i}$	29.5	11.5	22.6	30.0	65.5	40.0	83.3	30.5	39.1
1050	TATR ntk	$\frac{1}{K(K-1)}\sum_{i\neq j}  \nabla_{\theta}L_{j}(\theta_{\rm pre})  \odot  \nabla_{\theta}L_{i}(\theta_{\rm pre}) $	61.8	59.0	71.5	81.3	81.3	72.9	97.2	55.3	72.5
1057	TATR zero-shot	$\frac{1}{K(K-1)}\sum_{i\neq j}  \Delta_j  \odot  \Delta_i $	59.0	56.6	69.2	80.2	79.0	70.5	97.0	53.5	70.6
1058	TATR	$\frac{1}{K(K-1)}\sum_{i\neq j}  \nabla_{\theta}L_j(\theta_{\text{pre}})  \odot  \Delta_i $	62.7	59.3	72.3	82.3	80.5	72.6	97.0	55.4	72.8

## C.4 ANALYSIS OF SENSITIVITY FOR KNOWLEDGE CONFLICT

In this section, we explore various forms of conflict sensitivity in knowledge sharing. Specifically, we examine the following five approaches:

• **TATR positive**: This is calculated as the product between the task vector and the gradient. It promotes the merging of the components in the task vector that align with the gradient's ascent direction. The sensitivity of TATR positive is formulated as:

$$\frac{1}{K(K-1)}\sum_{i\neq j}\nabla_{\theta}L_{j}\left(\theta_{\text{pre}}\right)\odot\Delta_{i}.$$

• **TATR negative**: This is computed as the product between the negative task vector and the gradient. It encourages the merging of the components in the task vector that follow the gradient's descent direction. The sensitivity of TATR negative is formulated as:

$$-\frac{1}{K(K-1)}\sum_{i\neq j}\nabla_{\theta}L_{j}\left(\theta_{\text{pre}}\right)\odot\Delta_{i}.$$

• TATR ntk: This approach computes the product of the absolute values of gradients from different tasks, analogous to the Neural Tangent Kernel (NTK) (Jacot et al., 2018). It

 measures the influence of model updates for one task on another. Specifically, TATR ntk utilizes the following sensitivity for knowledge conflict:

$$\frac{1}{K(K-1)}\sum_{i\neq j}\left|\nabla_{\theta}L_{j}\left(\theta_{\text{pre}}\right)\right|\odot\left|\nabla_{\theta}L_{i}\left(\theta_{\text{pre}}\right)\right|.$$

• **TATR Zero-shot**: The zero-shot variant calculates the product of the absolute values between task vectors from different tasks:

$$\frac{1}{K(K-1)}\sum_{i\neq j} |\Delta_j| \odot |\Delta_i|$$

• **TATR** (Standard Version): The standard version computes the product of the absolute values of the task vector and the gradient. It encourages the fusion of the components in the task vector that are orthogonal to the gradient. The sensitivity of TATR is calculated as follows:

$$\frac{1}{K(K-1)}\sum_{i\neq j}\left|\nabla_{\theta}L_{j}\left(\theta_{\text{pre}}\right)\right|\odot\left|\Delta_{i}\right|.$$

Table 6 reports the multi-task performance of these methods when merging ViT-B/32 models across eight tasks. It is evident that both TATR Negative and TATR Positive result in a significant performance drop, indicating severe knowledge conflicts. In contrast, the standard TATR method effectively improves the performance of the pre-trained model across tasks, significantly mitigating knowledge conflicts. While the zero-shot and NTK variants exhibit slight performance degradation compared to TATR, they also demonstrate the ability to alleviate knowledge conflicts, suggesting that task vectors can approximate gradients to some extent.