# ATM: IMPROVING MODEL MERGING BY ALTERNATING TUNING AND MERGING

Anonymous authors

Paper under double-blind review

### Abstract

Model merging has recently emerged as a cost-efficient paradigm for Multi-task Learning (MTL). Among merging solutions, Task Arithmetic (Ilharco et al., 2022) stands out for its simplicity and effectiveness. In this paper, we start by motivating the effectiveness of task vectors with their relation to multi-task gradients. We show that in the single epoch scenario, task vectors are exactly equivalent to gradients obtained by performing gradient descent in a multi-task setting, and still approximate the latter with further epochs. We further strengthen the explanation by showing that task vectors work best when equality is maintained and motivate their effectiveness in the general case by showing that most of the contribution in the total update is determined by the gradient of the first epoch. Guided by this parallel, we propose viewing model merging as a single step in an iterative process that Alternates between Tuning and Merging (ATM). Acting as a midpoint between model merging and multi-task gradient descent, ATM obtains state-ofthe-art results with the same data and computing requirements. We first extensively evaluate our approach under diverse settings, demonstrating state-of-the-art performance, leading by an accuracy of up to 19% in computer vision and 20% in NLP over the best baselines. We then motivate its effectiveness empirically, showing increased orthogonality between task vectors and, theoretically, proving it to minimize an upper bound to the loss obtained by finetuning jointly on all tasks.

# 030

1

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028 029

031

INTRODUCTION

032 The pretrain-and-finetune paradigm has become the standard approach for numerous deep learning 033 tasks. In this framework, a model pretrained on large-scale unlabeled data is adapted to a specific 034 downstream task with minimal tuning. However, when addressing multiple tasks, a key limitation is the need to store separate finetuned models for each task. Model merging addresses this issue 035 by combining task-specific models into a single model capable of handling all tasks. This approach significantly reduces storage costs, as the unified model's size remains comparable to that of a single 037 task model, regardless of task count. Among numerous model merging methods, task arithmetic (Ilharco et al., 2022) stands out for its simplicity and effectiveness. Given a pretrained model  $\theta_0$ and a model  $\theta_i$  finetuned on task  $t_i$ , the task vector  $\tau_i$  is the subtraction of the pretrained weights 040 from the finetuned ones (i.e.  $\tau_i = \theta_i - \theta_{init}$ ). For multi-task learning with n tasks, task arithmetic 041 computes the sum of the n task vectors, properly scales it with a coefficient  $\alpha$ , and adds the resulting 042 vector to the pretrained model. 043

In this paper, we motivate the effectiveness of task arithmetic by relating task vectors with gradients 044 of the average loss over all the tasks. In particular, we show that when a model is finetuned for a single epoch using Gradient Descent, the corresponding task vector is exactly the additive inverse of 046 the gradient of the loss, rescaled by the learning rate. Analogously, the multi-task vector obtained 047 by summing the task vectors is equivalent to the additive inverse of the gradient of the average loss 048 across all tasks. With this perspective, task addition is equivalent to a step of gradient descent on the sum of the average losses of the tasks. When the finetuning is performed for several epochs, the equality is replaced by an approximation with an error dependent on the learning rate. We 051 further show that, while the single-epoch-finetuning assumption is violated in practice, the analogy to gradients can still explain why task vectors work in the first place. In fact, we show that, if we 052 consider the trajectory in the parameter space followed by the model during fine-tuning, the greatest contribution in terms of gradient norm is given by the first epoch. When this does not hold, we often



Figure 1: An illustration of the ATM framework up to the iteration h with |T| = 2 tasks (A and B). Starting with the pretrained model  $\theta_{base}^{(0)}$  as the base model, the *FT* step consists of finetuning it separately on the |T| tasks and the *Merge* step performs task vector aggregation and applies the multi-task task vector to the current base model, resulting in the next-iteration multi-task model. The process iterates with the resulting model at each iteration as the new base model for the next, until predefined iteration h or when some predefined condition is met. We observe increased task vector orthogonality as the ATM iteration grows.

find the gradients from the subsequent epochs to be aligned with the previous one, confirming thatthe direction is indeed dictated by the latter.

080 In this view, aggregation and merging in task arithmetic correspond to a noisy step of gradient de-081 scent when finetuning on the union of tasks, using as loss the sum of the average losses of the tasks. 082 In practice, this means that the one-step nature of these techniques likely results in overshooting the 083 multi-task minimum, as they would be actually tackling gradient descent over a multi-task dataset 084 with a single noisy step, where the multiplicative factor that is optimized over the validation set is, in fact, the learning rate. Building on these insights, we overcome the limitations of traditional one-step 085 task vector applications by introducing Alternating Tuning and Merging (ATM) - a novel multi-task model merging framework that generalizes task arithmetic by iteratively alternating between fine-087 tuning and merging, enabling a more gradual and refined integration of task-specific knowledge. 088 Given a compute budget of b epochs per task, traditional methods finetune each task for b epochs in a single pass. In contrast, ATM distributes the budget across k iterations, with each iteration 090 performing  $\frac{\partial}{\partial t}$  epochs of finetuning followed by task vector aggregation. The unified model from 091 each iteration serves as the starting point for the next. After k iterations, the final unified model is 092 deployed and evaluated. Notably, ATM is agnostic to the merging framework, allowing the integration of any interference-resolution techniques during the merge step to enhance performance. In 094 general, ATM significantly reduces time overhead compared to current baselines, and extensive experiments in computer vision and NLP show that it achieves state-of-the-art results without requiring hyperparameter tuning. 096

Our analysis of ATM unveils intriguing properties that shed light on its effectiveness. First, ATM task vectors exhibit a higher degree of orthogonality compared to baseline methods. We further prove that, in the simplified case where gradient descent is used to update the model, ATM minimizes the loss obtained by training jointly on all tasks. Our code, together with experiment configurations and checkpoints, is available for reproducibility purposes<sup>1</sup>

- <sup>102</sup> Our contribution is four-fold and is summarized as follows:
- 103 104
- 105
- We expose that task vectors, under certain conditions, are either equivalent to or approximate the gradients of loss of the corresponding tasks.
- 107

<sup>&</sup>lt;sup>1</sup>Link concealed to preserve anonymity.

111

112

113

114

115

116

117

- We point out that existing merging frameworks adopt one-shot merging, which likely overshoots the multi-task optimum, especially when task vectors exhibit large norms.
- We propose Alternating Tuning and Merging (ATM), a novel state-of-the-art model merging framework that generalizes task vector arithmetic. Thanks to its flexibility, ATM can readily integrate any interference-resolution framework, requiring no additional overhead.
- We empirically and mathematically motivate the effectiveness of ATM by showing increased orthogonality between task vectors obtained through ATM as compared to standard ones. We additionally prove that ATM reduces the loss of a multi-task model finetuned on the union of the datasets.
- 118 119 2 Related work

120 Mode connectivity and model merging Mode connectivity studies the weights that characterize 121 local minima on the loss landscape. Frankle et al. (2020) investigated linear mode connectivity in 122 models briefly trained from identical starting points, while Entezari et al. (2022) speculated that 123 all models converge to a shared basin once neuron permutations are resolved. Leveraging these 124 insights, permutation-based model merging aims to combine diverse models into a unified one, 125 aiming to inherit their capabilities without the overhead associated with ensembling. Singh & Jaggi 126 (2020) proposed an optimal-transport weight-matching method, while Git Re-Basin (Ainsworth 127 et al., 2022) proposed three novel matching methods acting both on weights and activations, with 128 REPAIR (Jordan et al., 2023) demonstrating significant barrier reduction through activation renormalization. Most recently, Navon et al. (2023) proposed merging models in the embedding space 129 of deep neural networks, while Crisostomi et al. (2024) proposed a cycle-consistent matching pro-130 cedure for improved merging. When models share the same pretrained initialization, Wortsman 131 et al. (2022) propose fusing them via a simple average. Jolicoeur-Martineau et al. (2023) propose to merge models by pushing them towards the population mean to ensure stability. RegMean Jin et al. 133 (2022) and Fisher-weighted averaging Matena & Raffel (2021) fall under the regime of weighted 134 averaging, where the weights are optimized according to some criteria. Daheim et al. shed light 135 on the positive relation between post-averaging multi-task performance and the gradient mismatch 136 between the constituent models. Finally, Choshen et al. (2022) even proposed model merging in re-137 placement of pertaining. They argue that pretrained checkpoints are not always the optimal starting 138 point for further finetuning, and a model resulting from merging finetuned models can be a better 139 starting point than any of its constituents.

140

141 **Task vectors** Task vector-based merging (Ilharco et al., 2022) finetunes a pretrained model on dif-142 ferent tasks to obtain task vectors (differences between finetuned and original checkpoints). Arithmetic operations on these vectors enable forgetting, analogy learning, and multi-task learning. Sev-143 eral works aim to improve task vector merging by reducing task interference (Deep et al., 2024; 144 Wang et al., 2024; Huang et al., 2024). Some methods include sparsifying task vectors or finetuning 145 only lottery tickets (Panda et al., 2024). TIES-merging (Yadav et al., 2023) merges vectors by prun-146 ing, selecting a unified sign vector, and merging disjointly, while Model Breadcrumbs (Davari & 147 Belilovsky, 2023) prunes both small and large-magnitude weights. DARE Merging (Yu et al., 2023) 148 randomly masks out a portion of weights and scales up the rest. AdaMerging (Yang et al., 2023) op-149 timizes aggregation coefficients, while Yang et al. (2024) propose task-specific modules for test-time 150 adaptation. Ortiz-Jimenez et al. (2024) introduce the concept of weight disentanglement and pro-151 pose finetuning in the tangent space. Unlike these one-shot methods, we introduce the perspective 152 of iterative model merging, evolving a base model towards a better multi-task performance.

153 154

155

# **3** TASK VECTORS AS GRADIENTS

In this section, we show that task vectors are tightly related to the gradients of the loss over the union of the tasks.

**Theorem 3.1.** Let  $\left\{\theta_t^{(k)}\right\}_{t=1}^{|T|}$  be a set of models obtained by finetuning  $\theta_{base}$  for k epochs over the set of tasks T via Gradient Descent (GD) with a learning rate of  $\eta$ , where fine-tuning over task  $t \in T$  minimizes  $\overline{L}_t(\theta) = \frac{1}{n_t} \sum_{i=1}^{n_t} \ell(x_i, y_i, \theta)$ . Let moreover  $\left\{\tau_t^{(k)}\right\}_{t=1}^{|T|}$  be a set of task vectors, with each  $\tau_t^{(k)} = \theta_t^{(k)} - \theta_{base}$ . Let  $\tau_{MT}^{(k)}$  be the multi-task vector  $\tau_{MT}^{(k)} = \sum_{t \in T} \tau_t^{(k)}$ . Finally let  $\theta_{MT}^{(k)}$ be the model obtained minimizing the loss  $\sum_{i=1}^{|T|} \overline{L}_i$  for k epochs with GD using learning rate equal to  $\alpha \eta$ . It holds that

$$\tau_{MT}^{(1)} = -\eta \nabla \sum_{t \in T} \overline{L}_t(\theta_{base}) \tag{1}$$

$$\tau_{MT}^{(k)} = -\eta \sum_{t \in T} \sum_{j=0}^{k-1} \nabla \overline{L}_i(\theta_{MT}^{(j)}) + \frac{\eta^2}{2} C(\{\theta_{MT}^{(j)}\}_{j=1}^{k-2}) + O(\eta^3)$$
(2)

with

177

178

181

183

186

191

207 208

$$C(\{\theta_{MT}^{(j)}\}_{j=1}^{h}) = \sum_{t \in T} \sum_{\ell=0}^{h} \nabla^{2} \overline{L}_{t}(\theta_{MT}^{(\ell)}) \sum_{m=0}^{\ell} \left[ \alpha \sum_{t' \neq t, t' \in T} \nabla \overline{L}_{t}'(\theta_{MT}^{(m)}) + (\alpha - 1) \nabla \overline{L}_{t}(\theta_{MT}^{(m)}) \right]$$

We report the proof in appendix A.1. To better appreciate the relation between a task vector and a gradient computed over the corresponding task dataset, it is worth focusing on the single task case, in which one is exactly the additive inverse of the other, scaled by the learning rate  $\eta$ .

**Remark 3.1.** From theorem 3.1, it follows that, for a single task t, and after a single epoch of finetuning,

$$\tau_t = -\eta \nabla \overline{L}_t(\theta_{base})$$

182 where  $\eta$  is the learning rate.

This also implies that, under the abovementioned assumptions, adding the task vector to the pre trained model can be seen as an approximation of a finetuning of the latter.

**Corollary 3.1.1.** Let  $\theta_{TA}^{(k)}$  be the model obtained using vanilla task arithmetics i.e.,  $\theta_{TA}^{(k)} = \theta_{base} + \alpha \sum_{t=1}^{T} \tau_t^{(k)}$ . Using the same notation of Theorem 3.1, it holds that

$$\theta_{TA}^{(1)} = \theta_{MT}^{(1)} \tag{3}$$

$$\theta_{TA}^{(k)} = \theta_{MT}^{(k)} + \frac{\eta^2}{2} C(\{\theta_{MT}^{(j)}\}_{j=1}^{k-2}) + O(\eta^3) \quad \text{for } k \ge 2$$
(4)

192 When k > 1,  $\tau_{\rm MT}$  still approximates the gradient of a 193 model finetuned for the same number of epochs via gra-194 dient descent, but with an error that is  $o(\eta^2)$ . In fact, we 195 show in fig. 2 that the multi-task model obtained by merg-196 ing models finetuned for a single epoch works better than 197 the one obtained by merging models finetuned to convergence. In other words, the best resulting merged model is obtained when task arithmetic is exactly a step of gradi-199 ent descent. The effectiveness of task vectors in the k > 1200 case, when the error nullifies the equality with respect 201 to gradients, can instead be motivated by the observation 202 that most of the trajectory followed by the model during 203 the finetuning phase is dictated by the gradient of the first 204 epoch: fig. 3a in fact shows that if we consider the epoch-205 wise normalized gradient norm for epoch k obtained as 206

$$\nabla_k = \frac{\|\nabla_{\theta}^{(k)}L\|}{\sum_{k'=1}^{K} \|\nabla_{\theta}^{(k')}L\|},$$



Figure 2: Accuracy of a model merged using task vectors over the task-specific models fine-tuned for 1 epoch and at convergence.

then the first epoch accounts for the largest contribution, reaching, in some cases, almost 70% of the
total of the gradient norms for all the epochs. When this is not true, e.g. in dataset RESISC45, we
speculate that the direction is still mostly the one dictated by the previous epoch: looking at fig. 3b,
we can see that the gradients of the first 5 epochs maintain high cosine similarity (> 0.8) with that of
the first one. It is worth remarking that, by considering Gradient Descent (GD) instead of Stochastic
Gradient Descent (SGD), the analysis does not apply exactly to the real case. However, considering
SGD as an approximation to GD, we speculate that the intuition is still correct. Under this unifying
lens, reducing task interference is analogous to reducing conflicting gradients in multi-task learning.



Cosine Similarity Matrix of Gradients Across Epochs Cosine Similarity Matrix of Gradients Across Epo

(a) Normalized gradient norms after 5 epochs of finetuning.

(b) Pairwise cosine similarities of the gradients of the first 10 epochs over dataset RESISC45.

#### 4 ATM: Alternating Tuning and Merging

Building upon the insights of section 3, we argue that task arithmetic is an approximation to a single GD step over the union of all the tasks. Following this parallel, we advocate taking further update steps sequentially and iteratively. The overall framework of ATM is depicted in fig. 1. Specifically, starting from a pretrained checkpoint as the base model  $\theta_{\text{base}}^{(0)}$ , we finetune it separately on each task to obtain the first-iteration task vectors  $\tau_1^{(1)}, \ldots, \tau_{|T|}^{(1)}$ . Task vectors are then aggregated and added to the base model to form the first-iteration unified model  $\theta_{\text{base}}^{(1)}$ . The procedure is iterated, and at each iteration k, the unified model  $\theta_{\text{base}}^{(k)}$  becomes the new base model for the next iteration of ATM (eq. (5)). The k-th iteration task vector for task t,  $\tau_t^{(k)}$ , is obtained as  $\theta_t^{(k)} - \theta_{\text{base}}^{(k)}$ , where  $\theta_t^{(k)}$  is a model obtained finetuning the k-th iteration base model  $\theta_{\text{base}}^{(k)}$  on task t.

$$\theta_{\text{base}}^{(k+1)} = \theta_{\text{base}}^{(k)} + \frac{\alpha}{|T|} \sum_{t \in T} \tau_t^{(k)} \quad \forall k = 0, \dots, K-1$$
(5)

After K iterations, we ultimately obtain  $\theta_{\text{base}}^{(K)}$ . The exact value of K can be predefined or based on a stop condition. Using a pretrained model is the standard practice. However, Choshen et al. 250 251 (2022) have suggested that pretrained checkpoints are not always the optimal starting point for fur-252 ther finetuning. Rather, they show that a model resulting from the merging of finetuned models can 253 be a better starting point than any of its constituents. Inspired by this, throughout ATM iterations, 254 we evolve the base model by iteratively merging task-specific finetuned models at each iteration 255 until the compute budget is met. In practice, each iteration of ATM involves finetuning the current 256 base model on all |T| tasks of interest, thereby obtaining |T| task vectors. These task vectors deter-257 mine the task-specific directions the current base model should follow in order to attain enhanced 258 performance on the corresponding tasks. The merging step of ATM at a given iteration simply con-259 sists of summing the mean current-iteration task vectors to the current base model, although any interference-resolution method in the task vector literature can be integrated. This step is intended 260 to pull the base model closer to the multi-task basin on the loss landscape. The averaging step en-261 sures the magnitude of the update remains insensitive to the number of tasks. Note that after each 262 iteration, the task vectors of the previous iterations can be safely discarded. Therefore, at any instant 263 of ATM, we only store one task vector for each of the |T| tasks and the base model, incurring no 264 additional memory requirements. 265

266 267

268

269

#### 5 UPPER BOUNDING THE MULTI-TASK LOSS

In this section we explore the relationship between the the ATM loss, defined as the mean of average losses over all tasks, and the loss of a model trained jointly of all the datasets. Analogously to

232

237

238

239

240 241

242

243 244

245 246

270 section 3, we will conduct this analysis under the assumption that GD, rather than SGD, is used for 271 optimizing the model parameters. This simplifying assumption removes the stochasticity introduced 272 by random sampling, enabling a more straightforward analysis while still providing valuable insights 273 into the underlying dynamics of the optimization process. With a slight abuse of notation, we denote 274 with t both the task and its corresponding dataset with cardinality  $n_t$ . The total number of samples for all tasks is given by  $N = \sum_{t \in T} n_t$ . 275

276 Inspired by the target model introduced by Daheim et al. (2023), we define as target loss for model 277 merging the loss  $L_{\text{target}}(\theta)$  of a model trained jointly on all the datasets. Namely,  $L_{\text{target}}(\theta) =$ 278  $\frac{1}{N}\sum_{i=1}^{N}\ell(x_i, y_i, \theta)$ . Given theorem 3.1, in the case in which we perform the merging after one step 279 of finetuning on each dataset, the ATM update using gradient is 280

281 282 283

287

290 291 292  $\theta_{\text{base}}^{(k+1)} = \theta_{\text{base}}^{(k)} + \frac{\alpha}{|T|} \sum_{t \in T} \tau_t^{(k)} = \theta_{\text{base}}^{(k)} - \alpha \eta \nabla \left( \frac{1}{|T|} \sum_{t \in T} \overline{L_t} \right)$ 

284 Namely, we are performing a gradient descent step minimizing the function  $L_{\text{ATM}} = \frac{1}{|T|} \sum_{t \in T} \overline{L_t}$ . 285 Having established that one step of ATM in gradient descent minimizes LATM, a crucial question 286 arises: under what conditions does minimizing  $L_{\text{ATM}}$  also lead to the minimization of  $L_{\text{target}}$ ? In other words, when can we be certain that optimizing the ATM loss will also minimize the loss 288 associated with training jointly on all datasets? To answer this question, we first note that  $L_{\text{ATM}}$  is 289 an unweighted average of the individual dataset losses, while the target loss is a weighted average:

$$L_{\text{target}}(\theta) = \frac{\sum_{t \in T} n_t \overline{L}_t(\theta)}{\sum_{t \in T} n_t},$$

293 We now analyze the parameter update from  $\theta^{(k)}$  to  $\theta^{(k+1)}$ . For both ATM and target methods, we denote the change in loss,  $L_{\text{method}}$ , as  $\Delta L_{\text{method}} = L_{\text{method}}(\theta^{(k)}) - L_{\text{method}}(\theta^{(k+1)})$ . In the following 295 theorem, we prove that if the drop in ATM loss exceeds a threshold  $\delta$  the target loss will also 296 decrease. The value of  $\delta$  depends on the size of the largest dataset with a decreasing loss and the 297 smallest dataset with an increasing loss. In particular, if the former dataset is larger than the latter, a 298 reduction in ATM loss reduces the target loss. In practice, this is ensured when the loss is reduced 299 on the largest dataset.

300 **Theorem 5.1.** Let D be the set of datasets where the loss decreases after a parameter update, and I 301 be the set of datasets where the loss increases or remains unchanged, defined as  $D = \{t \mid \Delta L_t > 0\}$ 302 and  $I = \{t \mid \Delta \overline{L}_t \leq 0\}$ . If the reduction in the ATM loss satisfies  $\Delta L_{ATM} > \delta$ , where 303

$$\delta = \frac{1}{|T|} \left( 1 - \frac{\min_{t \in I} n_t}{\max_{t \in D} n_t} \right) \sum_{t \in I} |\Delta \overline{L}_t|.$$

then the target loss  $L_{target}$  will also decrease, i.e.,  $\Delta L_{target} > 0$ . 307

308 We redirect the reader to appendix A.2 for the formal proof. 309

**Remark 5.1.** If we choose the target loss to be the maximum of the average loss across all datasets 310  $L_{target1} = \max_{t \in T} L_t$ , by leveraging the equivalence between the  $L_1$ -norm and the max norm, we 311 obtain the bound  $L_{target} \leq T \cdot L_{ATM}$ . 312

313 314

315 316

317

318

304 305 306

In this section, we illustrate the experimental setup and outcomes by comparing ATM against several recent baselines across a number of classification tasks in computer vision and NLP.

319 320

6.1 EXPERIMENTAL SETTING

321 **Datasets and Models** To evaluate ATM, we conduct experiments across multiple neural architectures and datasets in both computer vision and natural language processing (NLP) domains. For 322 computer vision tasks, we test ATM with a ViT-B-16 backbone (Dosovitskiy et al., 2021) and eval-323 uate it on a diverse set of datasets: CIFAR100 (Krizhevsky et al., 2009), DTD (Cimpoi et al., 2014),

340

341

342

343

353

354

372



Figure 4: Comparisons as computational budget K varies for ViT-B-16.

*EuroSAT* (Helber et al., 2019), *GTSRB* (Houben et al., 2013), *MNIST* (Lecun et al., 1998), *RE-SISC45* (Cheng et al., 2017), and *SVHN* (Netzer et al., 2011). For NLP tasks, we instead employ RoBERTa-base (Liu, 2019) and BERT-base-uncased (Devlin et al., 2019), evaluating them on the *GLUE benchmark* (Wang et al., 2019).

**Baselines and Metrics** To gauge the performance of ATM, we consider several model merging 344 baselines, including task arithmetic (TA) (Ilharco et al., 2022), TIES-merging (Yadav et al., 2023), 345 and model breadcrumbs (Davari & Belilovsky, 2023) for both computer vision and NLP tasks, and 346 DARE merging (Yu et al., 2023) for NLP tasks only. We adhere to author-recommended hyper-347 parameters, whenever needed or available, in order to ensure fair comparisons across experiments. 348 Specifically, for TIES-merging, we retain the top 15% of weights based on magnitude ranking. For 349 model breadcrumbs, we set  $\beta = 0.85$  and  $\gamma = 0.993$ . For DARE merging, we use a drop rate of 350 0.9. In all settings, we adopt mean aggregation of task vectors and use a scaling factor of 1 when 351 applying them to the base model. 352

#### 6.2 IMPACT OF EPOCH DISTRIBUTION ON PERFORMANCE

In this experiment, we first establish a fixed compute budget of 10 finetuning epochs for each task. Then, we seek
the optimal distribution of epochs among different numbers of ATM iterations. To exemplify, if 10 epochs are
distributed among 5 iterations, then in each iteration a
task is finetuned for 2 epochs.

As shown in fig. 5, when a fixed compute budget is assumed, maximizing the number of iterations while minimizing the number of epochs per iteration achieves the best results for ATM. This suggests that applying more fine-grained updates to the base model is preferred to the application of abrupt updates. Distributing the budget of 10 epochs into 10 iterations achieves the highest average accuracy of 89%, surpassing the extreme of 1 iteration of



Figure 5: Multi-task accuracy for different budget distributions (ViT-B-16)

10 epochs by 21%. Note that this latter setting is analogous to task arithmetic, which performs one-shot merging. Following this finding, we use this 1 epoch, 10 iterations setting for most of the presented experiments.

#### 6.3 EFFECT OF COMPUTE BUDGET



Figure 6: Average multi-task accuracy as budget varies.

We extend the comparison of ATM against baseline methods under different compute budgets. Specifically, we vary the per-task number of finetuning epochs (K) within {2, 4, 10}, and compare the average test accuracy across tasks. As shown in 4, ATM's overall performance is remarkably superior, regardless of the budget.



Figure 7: Comparisons as computational budget K varies for RoBERTa-base.

Moreover, as the budget increases up to 10 epochs, the proportional advantage of ATM over baselines is also increased. ATM leads with an edge of 7%, 11%, and 21% over the best-performing baseline, for 2, 4, and 10 epochs of budget, respectively. For detailed results on both vision and NLP benchmarks, we refer the reader to B.2.

Interestingly, we observe that while the accuracy of ATM increases with more epochs of finetuning, the opposite is true for the baseline methods, see 6. In other words, the more specialized the taskspecific models, the lower the performance of the unified multi-task model. ATM takes a different approach by gradually specializing the intermediate multi-task models, eschewing this issue.

#### 6.4 COMPARISONS IN ORIGINAL SETTINGS

Assuming the availability of training data, we compare the following three ATM settings against var-406 ious baselines under their original settings: (i) ATM finetuned on validation data (valFTATM) for 10 407 iterations of 1 epoch, (ii) ATM finetuned on training data for 10 iterations of 1 epoch, and (iii) ATM 408 finetuned on training data until convergence for 30 iterations of 1 epoch. As shown in table 1, all 409 ATM variants outperform the baselines by a large margin. When the training data is available and 410 the budget is limited to 10 epochs, ATM achieves an average accuracy of 89%, leading by 17% over 411 the best-performing baseline. Assuming no compute budget restriction, ATM converges after 30 412 iterations, achieving a remarkable average test accuracy of 91%.

413 414

415

392 393

396

397

399

400

401

402 403

404 405

> **TRAINING-DATA FREE SETTING** 6.5

One realistic constraint in practice is the absence of per-task training data, as finetuned models may 416 be downloaded from an online repository. In this section, we assume only the availability of the 417 validation data split, which is commonly exploited for hyperparameter tuning by baseline methods. 418 In contrast to the baselines, in this setting, ATM uses the validation data for finetuning the vari-419 ous tasks, leaving hyperparameters untuned. We call this textitATM variant *valFT ATM*, and we 420 compare its performance of against ATM and the baselines, adopting the author-suggested hyper-421 parameters to ensure fair comparisons. As shown in Table 1, *valFT ATM* significantly outperforms 422 the best-performing baseline by 10%. Note that the gap of 7% between valFT ATM and ATM is 423 admittedly due to the amount of data for finetuning, but on all tasks (except for EuroSAT), valFT 424 ATM performs comparably to ATM. We observe a similar phenomenon in the NLP domain as well.

425 426

427

#### 6.6 TIME AND MEMORY COMPLEXITIES

428 This experiment compares ATM against baselines in terms of time and memory consumption, assuming the 429 parameter count of the backbone to be d. Memory-wise, 430 ATM introduces no additional requirements compared to 431 task arithmetic, as the task vectors of each iteration can

Method	Time	Memory
Task Arithmetic	O(n * d)	O(d)
TIES	O(n * d * log(d))	O(d)
Breadcrumbs	O(n * d * log(d))	O(d)
ATM	O(k * n * d)	O(d)

CIFAR100	DTD	EuroSAT	GTSRB	MNIST	RESISC45	SVHN	Average
0.59	0.46	0.78	0.61	0.96	0.63	0.77	0.69
0.70	0.51	0.76	0.60	0.94	0.73	0.78	0.72
0.60	0.47	0.77	0.60	0.95	0.63	0.75	0.68
0.78	0.61	0.53	0.97	0.99	0.88	0.95	0.82
0.79	0.61	0.98	0.97	0.99	0.89	0.96	0.89
0.83	0.68	0.99	0.99	1.00	0.94	0.97	0.91
	CIFAR100 0.59 0.70 0.60 0.78 0.79 <b>0.83</b>	CIFAR100         DTD           0.59         0.46           0.70         0.51           0.60         0.47           0.78         0.61           0.79         0.61 <b>0.83 0.68</b>	CIFAR100         DTD         EuroSAT           0.59         0.46         0.78           0.70         0.51         0.76           0.60         0.47         0.77           0.78         0.61         0.53           0.79         0.61         0.98 <b>0.83 0.68 0.99</b>	CIFAR100         DTD         EuroSAT         GTSRB           0.59         0.46         0.78         0.61           0.70         0.51         0.76         0.60           0.60         0.47         0.77         0.60           0.78         0.61         0.53         0.97           0.79         0.61         0.98         0.97           0.83         0.68         0.99         0.99	CIFAR100         DTD         EuroSAT         GTSRB         MNIST           0.59         0.46         0.78         0.61         0.96           0.70         0.51         0.76         0.60         0.94           0.60         0.47         0.77         0.60         0.95           0.78         0.61         0.53         0.97         0.99           0.79         0.61         0.98         0.97         0.99           0.83         0.68         0.99         0.99         1.00	CIFAR100         DTD         EuroSAT         GTSRB         MNIST         RESISC45           0.59         0.46         0.78         0.61         0.96         0.63           0.70         0.51         0.76         0.60         0.94         0.73           0.60         0.47         0.77         0.60         0.95         0.63           0.78         0.61         0.53         0.97         0.99         0.88           0.79         0.61         0.98         0.97         0.99         0.89 <b>0.83 0.68 0.99 0.99 0.94</b>	CIFAR100         DTD         EuroSAT         GTSRB         MNIST         RESISC45         SVHN           0.59         0.46         0.78         0.61         0.96         0.63         0.77           0.70         0.51         0.76         0.60         0.94         0.73         0.78           0.60         0.47         0.77         0.60         0.95         0.63         0.77           0.78         0.61         0.53         0.97         0.99         0.88         0.95           0.79         0.61         0.98         0.97         0.99         0.89         0.96           0.83         0.68         0.99         0.99         0.89         0.96

Table 1: Accuracy comparison under original baseline settings (ViT-B-16)

be deleted after the merge step. The time complexity

443 of iteration-k ATM is equivalent to that of k times task

arithmetic, as it iterates finetuning and merging for k it-

erations. We argue that this time complexity is generally

asymptotically negligible compared to those of TIES and breadcrumbs, of which the dominant time overhead stems from pruning each task vector post-hoc according to the magnitude ranking of all weights, incurring O(d \* log(d)) time complexity. We compare the time and space complexities in table 2. Note that as long as log(d) is asymptotically greater than k, which is generally the case since k is usually a predefined constant integer (e.g. 10), ATM is faster than both TIES and breadcrumbs.

451

439

440 441

# 7 DISCUSSION

452 453 454

455

# 7.1 Orthogonality

Orthogonality between task vectors has been recommended as a desirable property for multi-task merging (Ilharco et al., 2022). Davari & Belilovsky (2023) adopts pairwise cosine similarity between task vectors as a proxy for task interference. Following this line, we again back the validity of this observation by identifying a positive correlation between ATM performance and task vector orthogonality. As shown in Figure fig. 9, as the number of ATM iterations increases, the magnitude of cosine similarity between task vectors tends to shrink, suggesting greater orthogonality as performance improves. Furthermore, as shown in 10, we find that ATM task vectors exhibit lower-magnitude average cosine similarity compared to the baseline methods.

463 464 465

### 7.2 TASK PROFICIENCY IS NOT MERGEABILITY

466 As shown in fig. 6, task-specific expertise does not imply multi-task performance. We observe that 467 better-performing task-specific models result in worse multi-task models when adopting baseline 468 methods, hinting that downstream performance is not a predictor of post-merging performance. We speculate that specialized models end up in highly dispersed locations in the parameter space, and 469 merging them abruptly in a one-shot fashion generates a suboptimal multi-task model; this can be 470 observed in fig. 8, where baseline methods end up all in the same (suboptimal) loss basin. On 471 the contrary, a lower degree of specialization ensures the task-specific models remain closer to the 472 pretrained checkpoint in the parameter space, leading to less aggressive updates when merging. The 473 above insights translate to less aggressive updates (shorter-norm task vectors) being more amenable 474 to merging. Capitalizing on this, ATM gradually aggregates task-specific models and updates the 475 base model accordingly, merging less aggressively but over multiple iterations. At each iteration, 476 the ATM task vectors represent the best nudges for the current base model without referring back to 477 the initial pretrained model as all baseline methods do.

478 479

480

### 7.3 EDUCATED TRAJECTORY

Task arithmetic performs the aggregation step abruptly in
a one-shot fashion over the initial pretrained checkpoint,
likely overshooting the multi-task optimum. In ATM,
however, the loss landscape is traversed iteration by iteration as the base model updates, leading to more informed
nudges toward the multi-task optimum. Figure 8 depicts





Figure 9: Average magnitudes of pairwise cosine similarity between ATM task vectors



Figure 10: Average magnitudes of cosine similarity between task vectors

the 2D projection of various checkpoints via PCA. Notably, TIES and breadcrumbs, being post-hoc enhancements of task arithmetic, end up around the same basin, whereas ATM takes gradual steps toward a different and better basin, signaling the effectiveness of our novel iterative merging paradigm.

#### 505 7.4 LIMITATIONS OF ITERATIVE MERGING 506

It is important to note that, while we have extensively shown the benefits of ATM in all the con-507 sidered settings, its iterative nature also has its drawbacks. In particular, iterative merging does not 508 yield a task representation that can be used for immediate adaptation to new tasks, as we instead se-509 quentially obtain a series of K vectors per task. storage. However, the approach is still advisable for 510 obtaining a single model that performs best overall in all the tasks. In this regard, we made sure to 511 maintain a fair comparison by using the same computing budget and data requirements, showing that 512 a small validation set is sufficient to effectively employ the framework; such a set is always assumed 513 to be presented in the literature and is usually used to optimize for the merging hyperparameters. 514 Finally, while the approach is similar to performing gradient descent on the union of all the tasks, 515 it still inherits one important property of task arithmetic: by obtaining the task vectors separately, 516 the approach does not require centralizing the data on a computing node, making it amenable to federated settings where data privacy is a key requirement. 517

518 519

495

496 497

# 8 CONCLUSIONS

520 521

In conclusion, this paper identifies a key limitation of task arithmetic—overshooting due to abrupt 522 model merging—and establishes its connection to gradient descent, forming the basis for our pro-523 posed model merging framework. We present Alternating Tuning and Merging (ATM), an iterative 524 framework that addresses the shortcomings of one-shot merging techniques. By alternating between 525 finetuning and merging, ATM effectively prevents overshooting and enhances multi-task perfor-526 mance. Extensive experiments on computer vision and NLP benchmarks demonstrate that ATM achieves state-of-the-art accuracy while maintaining computational efficiency comparable to exist-527 ing baselines. Additionally, our theoretical analysis reveals that ATM optimizes the upper bound 528 of the loss over the union of all the tasks and improves task vector orthogonality. The flexibility of 529 ATM opens numerous future research directions, including the integration of interference-mitigation 530 techniques and further refinement through advancements from the gradient descent literature. 531

532

# 533 ETHICS STATEMENT

534

This research was conducted with a strong commitment to ethical standards in both data usage and
experimental methodology. All datasets utilized in this study are publicly available. No personally
identifiable information was accessed or used during the course of this research. Additionally, the
experiments were designed to ensure fair comparisons across methods. We encourage future work
that adheres to these same ethical principles and addresses broader societal impacts of machine
learning technologies.

# 540 REPRODUCIBILITY STATEMENT

541 542 543

544

546

547 548

549 550

551

555

556

561

566

567

568

569

580

581

582

We are committed to ensuring the reproducibility of our results and have taken steps to facilitate this for the broader research community. The code, datasets, and configurations used for the experiments in this paper are made available via a public repository. We are open to providing further instructions on the usage of our code. The hyperparameters, frameworks, and evaluation metrics have been thoroughly documented, and we provide clear descriptions of our experimental setup to allow for straightforward replication of our findings. We encourage the community to utilize these resources and provide feedback for further improvements.

# References

- Samuel Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. Git Re-Basin: Merging models
   modulo permutation symmetries. In *The Eleventh International Conference on Learning Representations*, 2022.
  - Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10), 2017.
- Leshem Choshen, Elad Venezian, Noam Slonim, and Yoav Katz. Fusing finetuned models for better pretraining. *ArXiv preprint*, abs/2204.03044, 2022. URL https://arxiv.org/abs/ 2204.03044.
- Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In 2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014. IEEE Computer Society, 2014. doi: 10.1109/CVPR.2014.461. URL https://doi.org/10.1109/CVPR.2014.461.
  - Donato Crisostomi, Marco Fumero, Daniele Baieri, Florian Bernard, and Emanuele Rodolà.  $c^2m^3$ : Cycle-consistent multi-model merging. In *Advances in Neural Information Processing Systems*, volume 37, 2024.
- 570 Nico Daheim, Thomas Möllenhoff, Edoardo Ponti, Iryna Gurevych, and Mohammad Emtiyaz Khan.
  571 Model merging by uncertainty-based gradient matching. In *The Twelfth International Conference* 572 *on Learning Representations*.
- 573
  574 Nico Daheim, Thomas Möllenhoff, Edoardo Ponti, Iryna Gurevych, and Mohammad Emtiyaz Khan.
  575 In *The Twelfth International Conference on Learning Representations*, 2023.
- 576
  577
  578
  578
  579
  MohammadReza Davari and Eugene Belilovsky. Model breadcrumbs: Scaling multi-task model merging with sparse masks. ArXiv preprint, abs/2312.06795, 2023. URL https://arxiv. org/abs/2312.06795.
  - Pala Tej Deep, Rishabh Bhardwaj, and Soujanya Poria. Della-merging: Reducing interference in model merging through magnitude-based sampling. *ArXiv preprint*, abs/2406.11617, 2024. URL https://arxiv.org/abs/2406.11617.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, 2019. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/forum?
   id=YicbFdNTTy.

594 595 596 597	Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. The role of permutation invariance in linear mode connectivity of neural networks. In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.</i> OpenReview.net, 2022. URL https://openreview.net/forum?id=dNigytemkL.
598 599 600 601 602 603	Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode con- nectivity and the lottery ticket hypothesis. In <i>Proceedings of the 37th International Conference</i> <i>on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event</i> , volume 119 of <i>Proceedings</i> <i>of Machine Learning Research</i> . PMLR, 2020. URL http://proceedings.mlr.press/ v119/frankle20a.html.
604 605 606	Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. <i>IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing</i> , 12(7), 2019.
607 608 609	Sebastian Houben, Johannes Stallkamp, Jan Salmen, Marc Schlipsing, and Christian Igel. Detection of traffic signs in real-world images: The German Traffic Sign Detection Benchmark. In <i>International Joint Conference on Neural Networks</i> , number 1288, 2013.
610 611 612 613	Chenyu Huang, Peng Ye, Tao Chen, Tong He, Xiangyu Yue, and Wanli Ouyang. Emr-merging: Tuning-free high-performance model merging. <i>ArXiv preprint</i> , abs/2405.17461, 2024. URL https://arxiv.org/abs/2405.17461.
614 615 616	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. <i>The Eleventh International Conference on Learning Representations</i> , 2022.
617 618 619 620	Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. Dataless knowledge fusion by merging weights of language models. <i>ArXiv preprint</i> , abs/2212.09849, 2022. URL https: //arxiv.org/abs/2212.09849.
621 622 623	Alexia Jolicoeur-Martineau, Emy Gervais, Kilian Fatras, Yan Zhang, and Simon Lacoste-Julien. Population parameter averaging (papa). <i>ArXiv preprint</i> , abs/2304.03094, 2023. URL https: //arxiv.org/abs/2304.03094.
624 625 626	Keller Jordan, Hanie Sedghi, Olga Saukh, Rahim Entezari, and Behnam Neyshabur. REPAIR: REnormalizing permuted activations for interpolation repair. In <i>The Eleventh International Conference on Learning Representations</i> , 2023.
627 628 629	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
630 631 632	Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recog- nition. <i>Proceedings of the IEEE</i> , 86(11), 1998. doi: 10.1109/5.726791.
633 634	Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. ArXiv preprint, abs/1907.11692, 2019. URL https://arxiv.org/abs/1907.11692.
635 636	Michael Matena and Colin Raffel. Merging models with fisher-weighted averaging. 2021.
637 638	Aviv Navon, Aviv Shamsian, Ethan Fetaya, Gal Chechik, Nadav Dym, and Haggai Maron. Equivariant deep weight space alignment, 2023.
639 640 641	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 <i>NIPS workshop on deep</i> <i>learning and unsupervised feature learning</i> , volume 2011. Granada, 2011.
642 643 644 645	Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. Task arithmetic in the tangent space: Improved editing of pre-trained models. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
646 647	Ashwinee Panda, Berivan Isik, Xiangyu Qi, Sanmi Koyejo, Tsachy Weissman, and Prateek Mit- tal. Lottery ticket adaptation: Mitigating destructive interference in llms. <i>ArXiv preprint</i> , abs/2406.16797, 2024. URL https://arxiv.org/abs/2406.16797.

- 648 Model fusion via optimal transport. Sidak Pal Singh and Martin Jaggi. In Hugo 649 Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien 650 Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference 651 on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, 652 URL https://proceedings.neurips.cc/paper/2020/hash/ *virtual*, 2020. fb2697869f56484404c8ceee2985b01d-Abstract.html. 653
- 654 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 655 Glue a multi-task benchmark and analysis platform for natural language understanding. In 7th 656 International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, 657 May 6-9, 2019. OpenReview.net, 2019. URL https://openreview.net/forum?id= 658 rJ4km2R5t7.
- Ke Wang, Nikolaos Dimitriadis, Guillermo Ortiz-Jimenez, François Fleuret, and Pascal Frossard. 660 Localizing task information for improved model merging and compression. In *Forty-first Inter*national Conference on Machine Learning, 2024.
- Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo 663 Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and 664 Ludwig Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accu-665 racy without increasing inference time. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba 666 Szepesvári, Gang Niu, and Sivan Sabato (eds.), International Conference on Machine Learn-667 ing, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA, volume 162 of Proceedings of 668 Machine Learning Research. PMLR, 2022. URL https://proceedings.mlr.press/ 669 v162/wortsman22a.html. 670
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. Ties-merging: Re-671 solving interference when merging models. 2023. 672
- 673 Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guibing Guo, Xingwei Wang, and Dacheng 674 Tao. Adamerging: Adaptive model merging for multi-task learning. In The Twelfth International 675 Conference on Learning Representations, 2023.
- 676 Enneng Yang, Li Shen, Zhenyi Wang, Guibing Guo, Xiaojun Chen, Xingwei Wang, and Dacheng 677 Tao. Representation surgery for multi-task model merging. ArXiv preprint, abs/2402.02705, 2024. 678 URL https://arxiv.org/abs/2402.02705. 679
  - Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Language models are super mario: Absorbing abilities from homologous models as a free lunch. In Forty-first International Conference on Machine Learning, 2023.

#### PROOFS А

#### PROOF THEOREM 3.1 AND COROLLARY 3.1.1 A.1

In this section we will prove the Theorem 3.1 and Corollary 3.1.1. We repeat the statement of the theorem and of the corollary for easiness of reading.

**Theorem.** 3.1 Let  $\{\theta_t^{(k)}\}_{t=1}^{|T|}$  be a set of models obtained by finetuning  $\theta_{base}$  for k epochs over the set of tasks T via Gradient Descent (GD) with a learning rate of  $\eta$ , where finetuning over task  $t \in T$  minimizes  $\overline{L}_t(\theta) = \frac{1}{n_t} \sum_{i=1}^{n_t} \ell(x_i, y_i, \theta)$ . Let moreover  $\{\tau_t^{(k)}\}_{t=1}^{|T|}$  be a set of task vectors, with each  $\tau_t^{(k)} = \theta_t^{(k)} - \theta_{base}$ . Let  $\tau_{MT}^{(k)}$  be the multi-task vector  $\tau_{MT}^{(k)} = \sum_{t \in T} \tau_t^{(k)}$ . Finally, let  $\theta_{MT}^k$  be the 690 691 692 693 694 695 model obtained minimizing the loss  $\sum_{i=1}^{|T|} L_i$  for k epochs. It holds that 696

$$\tau_{MT}^{1} = -\eta \nabla \sum_{t \in T} \overline{L}_{t}(\theta_{base}) \tag{6}$$

699

697 698

659

661

662

680

681

682 683 684

685 686

687

688

700  
701 
$$\tau_{MT}^{k} = -\eta \sum_{t \in T} \sum_{j=0}^{k-1} \nabla \overline{L}_{i}(\theta_{MT}^{j}) + \frac{\eta^{2}}{2} C(\{\theta_{MT}^{j}\}_{j=1}^{k}) + O(\eta^{3})$$
(7)

with  $C(\{\theta_{MT}^{j}\}_{j=1}^{k}) = \sum_{t \in T} \sum_{\ell=0}^{k} \nabla^{2} \overline{L}_{t}(\theta_{MT}^{\ell}) \sum_{m=0}^{\ell} \left[ \alpha \sum_{t' \neq t, t' \in T} \nabla \overline{L}_{t}'(\theta_{MT}^{m}) + (\alpha - 1) \nabla \overline{L}_{t}(\theta_{MT}^{m}) \right]$ 

**Corollary.** 3.1.1 Let  $\theta_{TA}^k$  be the model obtained using vanilla task arithmetics i.e.,  $\theta_{TA}^k = \theta_{base} + \alpha \sum_{t=1}^{|T|} \tau_t^k$ . It holds that

$$\theta_{TA}^1 = \theta_{MT}^1 \tag{8}$$

$$\theta_{TA}^k = \theta_{MT}^k + \eta^2 o(1) \quad for \ k > 1 \tag{9}$$

where  $\eta$  is the learning rate used for the finetuning of the model on the single datasets.

We recall that  $\theta_i^k$  is the model obtained finetuing on task *i* for *k* epochs and that both the finetuing on different tasks and the training on the average loss start from a pretrained model  $\theta^0$ .

To prove the statement of the theorem and of the corollary we need a intermediate result. We introduce the following notation:

$$r_i(\theta) = \alpha \sum_{j \neq i} \nabla \overline{L}_j(\theta) + (\alpha - 1) \nabla \overline{L}_i(\theta) = \alpha \sum_{j=1}^{|T|} \nabla \overline{L}_j(\theta) - \nabla \overline{L}_i(\theta^{(0)})$$
(10)

$$p_{i}^{k}(\theta^{0}, \theta_{\text{MT}}^{1}, \dots, \theta_{\text{MT}}^{k}) = \sum_{j=0}^{k} r_{i}(\theta_{\text{MT}}^{k})$$
(11)

$$s_i^k(\theta^0, \dots, \theta_{\mathrm{MT}}^k) = \sum_{j=0}^k \nabla^2 \overline{L}_i(\theta_{\mathrm{MT}}^j) [p_i^j(\theta^0, \dots, \theta_{\mathrm{MT}}^{j-1})].$$
(12)

**Lemma A.1.** Using the notation introduced in theorem 3.1, it holds that

$$\theta_i^{(1)} = \theta_{MT}^{(1)} + \eta p_i^0(\theta^0) \tag{13}$$

and for  $m \ge 2$ 

$$\theta_i^{(m+1)} = \theta_{MT}^{(m+1)} + \eta p_i^m(\theta^0, \dots, \theta_{MT}^m) - \frac{\eta^2}{2} s_i^{m-1}(\theta^0, \dots, \theta_{MT}^{m-1}) + O(\eta^3)$$
(14)

*Proof.* We first show that the statement is true for m = 1, and then prove the results for  $m \ge 2$  by induction. In this case, the base case is given for m = 2. In the induction step, instead, we prove that if the statement holds for any given case m then it must also hold for the next case m + 1.

m = 1. First epoch For each task  $i = 1, \ldots, |T|$ 

$$\theta_i^{(1)} = \theta^{(0)} - \eta \nabla \overline{L}_i(\theta^{(0)}) \text{ while } \theta_{\mathrm{MT}}^{(1)} = \theta^{(0)} - \alpha \eta \sum_{i \in T} \nabla \overline{L}_i(\theta^{(0)}).$$

Consequently, it holds that

754  $\theta_i^1 = \theta_{\mathrm{MT}}^{(1)} + \eta \left[ \alpha \sum_{j \neq i} \nabla \overline{L}_j(\theta^{(0)}) + (\alpha - 1) \nabla \overline{L}_i(\theta^{(0)}) \right]$ 

$$= \theta_{\rm MT}^{(1)} + \eta r_i(\theta^{(0)}) = \theta_{\rm MT}^{(1)} + \eta p_i^0(\theta^0).$$

#### m = 2. Second epoch $\theta_i^{(2)} = \theta_i^1 - \eta \nabla \overline{L}_i(\theta_i^1)$ $= \boldsymbol{\theta}_{\mathrm{MT}}^{(1)} + \eta \boldsymbol{r}_i(\boldsymbol{\theta}^{(0)}) - \eta \nabla \overline{L}_i \left(\boldsymbol{\theta}_{\mathrm{MT}}^{(1)} + \eta \boldsymbol{r}_i(\boldsymbol{\theta}^{(0)})\right)$ $\overset{Taylor}{\approx} \theta_{\mathrm{MT}}^{(1)} + \eta r_i(\theta^{(0)}) - \eta \nabla \overline{L}_i(\theta_{\mathrm{MT}}^{(1)}) - \frac{\eta^2}{2} \nabla^2 \overline{L}_i(\theta_{\mathrm{MT}}^1) r_i(\theta^{(0)}) + O(\eta^3)$ $=\theta_{\rm MT}^{(1)} - \eta \nabla \overline{L}_i(\theta_{\rm MT}^{(1)}) + \eta r_i(\theta^{(0)}) - \frac{\eta^2}{2} \nabla^2 \overline{L}_i(\theta_{\rm MT}^1) r_i(\theta^{(0)}) + O(\eta^3)$ $= \theta_{\mathrm{MT}}^{(1)} - \eta \nabla \overline{L}_i(\theta_{\mathrm{MT}}^{(1)}) + \eta \alpha \sum_{t \in T} \nabla \overline{L}_i(\theta_{\mathrm{MT}}^{(1)}) - \eta \alpha \sum_{t \in T} \nabla \overline{L}_i(\theta_{\mathrm{MT}}^{(1)}) + \eta r_i(\theta^{(0)})$ $-\frac{\eta^2}{2}\nabla^2\overline{L}_i(\theta_{\mathrm{MT}}^1)r_i(\theta^{(0)})+O(\eta^3)$ $=\theta_{\rm MT}^{(1)} + \eta r_i(\theta_{\rm MT}^1) + \eta r_i(\theta^{(0)}) - \frac{\eta^2}{2} \nabla^2 \overline{L}_i(\theta_{\rm MT}^1) r_i(\theta^{(0)}) + O(\eta^3)$ $=\theta_{MT}^{1}+\eta p_{i}^{1}(\theta^{0},\ldots,\theta_{MT}^{1})-\frac{\eta^{2}}{2}s_{i}^{0}(\theta^{0})+O(\eta^{3})$

### **Inductive step** Let us assume that

$$\theta_i^m = \theta_{MT}^m + \eta p_i^{m-1}(\theta^0, \dots, \theta_{MT}^{m-1}) - \frac{\eta^2}{2} s_i^{m-2}(\theta^0, \dots, \theta_{MT}^{m-2}) + O(\eta^3)$$

We can derive that

$$\begin{array}{ll} \theta_{i}^{m+1} = \theta_{i}^{m} - \eta \nabla \overline{L}_{i}(\theta_{i}^{m}) \\ = \theta_{MT}^{m} + \eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-2}(\theta^{0}, \dots, \theta_{MT}^{m-2}) - \eta \nabla \overline{L}_{i}(\theta_{i}^{m}) + O(\eta^{3}) \\ = \theta_{MT}^{m} + \eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-2}(\theta^{0}, \dots, \theta_{MT}^{m-2}) \\ - \eta \nabla \overline{L}_{i}\left(\theta_{MT}^{m} + \eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-2}(\theta^{0}, \dots, \theta_{MT}^{m-2})\right) + O(\eta^{3}) \\ = \theta_{MT}^{m} + \eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-2}(\theta^{0}, \dots, \theta_{MT}^{m-2}) \\ - \eta \nabla \overline{L}_{i}(\theta_{MT}^{m}) - \frac{\eta}{2} \nabla^{2} \overline{L}_{i}(\theta_{MT}^{m}) \left(\eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-2}(\theta^{0}, \dots, \theta_{MT}^{m-2})\right) \\ = \theta_{MT}^{m} + \eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-2}(\theta^{0}, \dots, \theta_{MT}^{m-2}) \\ = \theta_{MT}^{m} + \eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-2}(\theta^{0}, \dots, \theta_{MT}^{m-2}) \\ = \theta_{MT}^{m} + \eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-2}(\theta^{0}, \dots, \theta_{MT}^{m-1}) + O(\eta^{3}) \\ = \theta_{MT}^{m} + \eta p_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) + O(\eta^{3}) \\ = \theta_{MT}^{(m+1)} + \eta p_{i}^{m}(\theta^{0}, \dots, \theta_{MT}^{m}) - \frac{\eta^{2}}{2} s_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) + O(\eta^{3}) \\ = \theta_{MT}^{(m+1)} + \eta p_{i}^{m}(\theta^{0}, \dots, \theta_{MT}^{m}) - \frac{\eta^{2}}{2} s_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) + O(\eta^{3}) \\ = \theta_{MT}^{(m+1)} + \eta p_{i}^{m}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) + O(\eta^{3}) \\ = \theta_{MT}^{(m+1)} + \eta p_{i}^{m}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) + O(\eta^{3}) \\ = \theta_{MT}^{(m+1)} + \eta p_{i}^{m}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) + O(\eta^{3}) \\ = \theta_{MT}^{(m+1)} + \eta p_{i}^{m}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \frac{\eta^{2}}{2} s_{i}^{m-1}(\theta^{0}, \dots, \theta_{MT}^{m-1}) + O(\eta^{3}) \\ = \theta_{MT}^{(m+1)} + \eta p_{i}^{m}(\theta^{0}, \dots, \theta_{MT}^{m-1}) - \theta_{MT}^{(m+1)} + \theta_{M$$

Proof Theorem and Corollary. For the first epoch

$$\theta_{\mathrm{TA}}^{1} = \theta^{0} + \alpha \sum_{i \in T} \tau_{i}^{1} = \theta^{0} - \eta \alpha \sum_{i \in T} \nabla \overline{L}_{i}(\theta^{0})$$

while, choosing  $\alpha \eta$  as learning rate for the loss  $\sum_{i \in T} \overline{L}_i$ :

$$\theta_{\mathrm{MT}}^{(1)} = \theta^{(0)} - \alpha \eta \sum_{i \in T} \nabla \overline{L}_i(\theta^{(0)}).$$

$$\begin{aligned} -\alpha\eta\sum_{j=0}^{k-1}\nabla\sum_{t\in T}\overline{L}_{t}(\theta_{MT}^{j}) + \eta p_{i}^{k-1}(\eta^{0},\ldots,\eta_{MT}^{k-1}) \\ &= -\alpha\eta\sum_{j=0}^{k-1}\nabla\sum_{t\in T}\overline{L}_{t}(\theta_{MT}^{j}) + \sum_{j=0}^{k-1}r_{i}(\theta_{MT}^{k}) \\ &= -\alpha\eta\sum_{j=0}^{k-1}\nabla\sum_{t\in T}\overline{L}_{i}(\theta_{MT}^{j}) + \sum_{j=0}^{k-1}\alpha\sum_{j\in T}\nabla\overline{L}_{j}(\theta_{MT}^{j}) - \nabla\overline{L}_{i}(\theta_{MT}^{j}) \\ &= -\eta\sum_{j=0}^{k-1}\nabla\overline{L}_{i}(\theta_{MT}^{j}) \end{aligned}$$

# Namely

$$\begin{aligned} \theta_i^{m+1} &= \theta_{MT}^{(m+1)} + \eta p_i^m(\theta^0, \dots, \theta_{MT}^m) - \frac{\eta^2}{2} s_i^{m-1}(\theta^0, \dots, \theta_{MT}^{m-1}) + O(\eta^3) \\ &= \theta^0 - \alpha \eta \sum_{j=0}^m \nabla \sum_{t \in T} \overline{L}_i(\theta_{MT}^j) + \eta p_i^m(\theta^0, \dots, \theta_{MT}^m) - \frac{\eta^2}{2} s_i^{m-1}(\theta^0, \dots, \theta_{MT}^{m-1}) + O(\eta^3) \\ &= \theta^0 - \eta \sum_{j=0}^m \nabla \overline{L}_i(\theta_{MT}^j) - \frac{\eta^2}{2} s_i^{m-1}(\theta^0, \dots, \theta_{MT}^{m-1}) + O(\eta^3) \end{aligned}$$

$$au_i^k = heta_i^k - heta^0$$

$$= -\eta \sum_{j=0}^{k-1} \nabla \overline{L}_i(\theta_{MT}^j) - \frac{\eta^2}{2} s_i^{k-2}(\theta^0, \dots, \theta_{MT}^{k-2}) + O(\eta^3)$$
(17)

(16)

Consequently the model obtain with TA is

we can rewrite the tasks vectors as

$$\begin{aligned} \theta_{\text{TA}}^{k} &= \theta^{0} + \alpha \sum_{i \in T} \tau_{i}^{k} \\ &= \theta^{0} - \eta \alpha \sum_{j=0}^{k-1} \sum_{i \in T} \nabla \overline{L}_{i}(\theta_{MT}^{j}) - \alpha \sum_{i \in T} \frac{\eta^{2}}{2} s_{i}^{k-2}(\theta^{0}, \dots, \theta_{\text{MT}}^{k-2}) + O(\eta^{3}) \\ &= \theta_{\text{MT}}^{k} - \alpha \sum_{i \in T} \frac{\eta^{2}}{2} s_{i}^{k-2}(\theta^{0}, \dots, \theta_{\text{MT}}^{k-2}) + O(\eta^{3}) \end{aligned}$$

#### A.2 PROOF THEOREM 5.1

In this section, we provide a rigorous proof of Theorem 5.1, establishing the conditions under which ATM implicitly optimizes the loss of the model trained on the union of datasets.

$$\sum_{j \in N} |D_j| \Delta \overline{L_j} = -\sum_{j \in N} |D_j| |\Delta \overline{L_j}|.$$

For hypothesis  $\Delta L_{\text{ATM}} > \delta$ , this implies  $\sum_{j=1}^{n} \Delta_j > n\delta$ . We have that  $\sum_{j \in P} \Delta \overline{L_j} + \sum_{j \in N} \Delta \overline{L_j} = \sum_{j \in P} \Delta \overline{L_j} - \sum_{j \in N} |\Delta \overline{L_j}| > n\delta$ , namely  $\sum_{j \in P} \Delta \overline{L_j} > n\delta + \sum_{j \in N} |\Delta \overline{L_j}|$ 

We want to prove that  $\Delta L_{\text{target}} > 0$ . Let us now consider  $\Delta L_{\text{target}} > 0$  iff  $\sum_{j=1}^{n} |D_j| \Delta_j > 0$ .

$$\sum_{j=1}^{n} |D_j| \Delta_j = \sum_{j \in P} |D_j| \Delta \overline{L_j} + \sum_{j \in N} |D_j| \Delta \overline{L_j}$$
$$= \sum_{j \in P} |D_j| \Delta \overline{L_j} - \sum_{j \in N} |D_j| |\Delta \overline{L_j}|$$
$$> \min_{j \in P} |D_j| \sum_{j \in P} \Delta \overline{L_j} - \max_{j \in N} |D_j| \sum_{j \in N} |\Delta \overline{L_j}|$$

$$> \min_{j \in P} |D_j| \left\lfloor n\delta + \sum_{j \in N} |\Delta \overline{L_j}| \right\rfloor - \max_{j \in N} |D_j| \sum_{j \in N} |\Delta \overline{L_j}|$$
$$= n \min_{j \in P} |D_j|\delta + (\min_{j \in P} |D_j| - \max_{j \in N} |D_j| \sum_{j \in N}) |\Delta \overline{L_j}|$$

The last line of the previous equation is positive by hypothesis since we assumed  $\delta > \frac{1}{n}(1 - \frac{\min_{j \in N} |D_j|}{\max_{j \in P} |D_j|}) \sum_{j \in N} |\Delta \overline{L_j}|$ .

# **B** ADDITIONAL RESULTS

#### B.1 FULL RESULTS OVER VARYING COMPUTATIONAL BUDGET

In the main paper, we pictorially depicted the multi-task accuracies of the baselines and variants of ATM, in the form of radar plots. For deeper analysis, here we report the full results of ATM compared to the baselines, as the computational budget varies for 2, 4, 7, and 10 epochs.

#### B.2 COSINE SIMILARITY OF EPOCH-WISE GRADIENTS

We report in fig. 11 the cosine similarity of gradients for the first 10 epochs over DTD and EuroSAT, as these do not have a marked difference in gradient norms between the first epoch and the remaining ones in fig. 3a. We see that the alignment of subsequent gradients observed in RESISC45 still holds in DTD, even if with less marked similarities. On the other hand, this does not seem to hold for EuroSAT.

		CIFAR100	DTD	EuroSAT	GTSRB	MNIST	RESISC45	SVHN	average
	Task Arithmetic	0.58	0.39	0.78	0.59	0.98	0.58	0.83	0.68
5	TIES	0.76	0.46	0.71	0.55	0.96	0.71	0.81	0.71
II	Breadcrumbs	0.72	0.48	0.80	0.61	0.97	0.70	0.81	0.72
K	ValFT ATM	0.70	0.53	0.84	0.84	0.98	0.78	0.90	0.80
	ATM	0.74	0.51	0.89	0.76	0.98	0.74	0.89	0.79
	Task Arithmetic	0.71	0.48	0.80	0.62	0.97	0.71	0.82	0.73
4	TIES	0.75	0.45	0.69	0.56	0.97	0.69	0.84	0.70
	Breadcrumbs	0.68	0.43	0.79	0.61	0.97	0.67	0.81	0.71
K	ValFT ATM	0.66	0.55	0.73	0.85	0.98	0.83	0.94	0.79
	ATM	0.75	0.55	0.95	0.90	0.99	0.82	0.94	0.84
	Task Arithmetic	0.67	0.44	0.79	0.62	0.97	0.67	0.83	0.71
1	TIES	0.72	0.43	0.66	0.56	0.97	0.66	0.85	0.69
	Breadcrumbs	0.65	0.41	0.76	0.62	0.97	0.62	0.81	0.69
K	ValFT ATM	0.75	0.62	0.66	0.96	0.99	0.85	0.96	0.83
	ATM	0.77	0.58	0.98	0.96	0.99	0.87	0.95	0.87
	Task Arithmetic	0.63	0.41	0.76	0.63	0.97	0.62	0.83	0.69
10	TIES	0.68	0.41	0.70	0.56	0.98	0.62	0.86	0.69
	Breadcrumbs	0.60	0.40	0.78	0.58	0.98	0.59	0.81	0.68
K	ValFT ATM	0.78	0.61	0.53	0.97	0.99	0.88	0.95	0.82
	ATM	0.79	0.61	0.98	0.97	0.99	0.89	0.96	0.89

Table 3: ATM vs. Baselines as budget varies (ViT-B-16)

		CoLA	SST2	MRPC	QQP	MNLI	QNLI	RTE	Average
K = 2	Task Arithmetic TIES Breadcrumbs DARE	0.70 0.69 0.69 <b>0.84</b> 0.71	0.56 0.51 0.53 0.53	0.66 0.68 0.66 0.68	0.37 0.37 0.37 0.41	0.47 0.37 0.45 0.32	0.54 0.51 0.52 0.50	0.51 0.47 0.48 0.47	0.54 0.51 0.53 0.54
	ATM	0.71	0.78	0.08 <b>0.69</b>	0.40	0.64 0.65	0.68	0.62	0.66
K = 4	Task Arithmetic	0.71	0.60	0.66	0.37	0.51	0.54	0.51	0.56
	TIES	0.69	0.51	0.68	0.37	0.42	0.51	0.47	0.52
	Breadcrumbs	0.70	0.57	0.66	0.37	0.49	0.53	0.51	0.55
	DARE	<b>0.83</b>	0.50	0.65	0.41	0.33	0.50	0.46	0.53
	valFT ATM	0.70	<b>0.83</b>	<b>0.70</b>	<b>0.68</b>	0.67	0.73	0.62	0.70
	ATM	0.73	0.81	<b>0.70</b>	0.58	<b>0.70</b>	<b>0.80</b>	<b>0.67</b>	<b>0.71</b>
K = 7	Task Arithmetic	0.71	0.65	0.65	0.38	0.52	0.56	0.55	0.57
	TIES	0.69	0.51	0.68	0.37	0.42	0.51	0.47	0.52
	Breadcrumbs	0.71	0.61	0.65	0.37	0.50	0.54	0.52	0.56
	DARE	<b>0.84</b>	0.52	0.65	0.43	0.33	0.49	0.44	0.52
	valFT ATM	0.69	<b>0.85</b>	0.69	<b>0.75</b>	0.68	0.72	0.60	0.71
	ATM	0.72	0.79	<b>0.70</b>	0.73	<b>0.73</b>	<b>0.82</b>	<b>0.67</b>	<b>0.74</b>
K = 10	Task Arithmetic	0.71	0.66	0.60	0.38	0.56	0.59	0.60	0.59
	TIES	0.69	0.51	0.68	0.37	0.43	0.51	0.47	0.52
	Breadcrumbs	0.71	0.61	0.62	0.38	0.54	0.57	0.57	0.57
	DARE	<b>0.85</b>	0.51	0.64	0.42	0.33	0.49	0.44	0.52
	valFT ATM	0.68	<b>0.86</b>	0.66	<b>0.76</b>	0.70	0.73	0.63	0.72
	ATM	0.72	0.83	<b>0.69</b>	<b>0.76</b>	<b>0.74</b>	<b>0.81</b>	<b>0.66</b>	<b>0.74</b>

Table 4: ATM vs Baselines as budget varies (RoBERTa-base)

		CoLA	SST2	MRPC	QQP	MNLI	QNLI	RTE	Average
	Task Arithmetic	0.61	0.70	0.39	0.65	0.41	0.57	0.54	0.55
$\sim$ 1	TIES	0.32	0.50	0.32	0.63	0.37	0.50	0.54	0.45
	Breadcrumbs	0.60	0.68	0.39	0.65	0.40	0.57	0.54	0.55
X	DARE	0.82	0.54	0.69	0.40	0.33	0.50	0.46	0.53
	ValFT ATM	0.66	0.71	0.47	0.65	0.44	0.57	0.56	0.58
	ATM	0.68	0.68	0.39	0.64	0.46	0.56	0.61	0.58
	Task Arithmetic	0.59	0.66	0.39	0.67	0.41	0.55	0.57	0.55
	TIES	0.32	0.50	0.32	0.63	0.38	0.51	0.53	0.46
7	Breadcrumbs	0.58	0.65	0.40	0.67	0.41	0.55	0.56	0.54
5	DARE	0.82	0.53	0.66	0.43	0.33	0.50	0.47	0.54
ł	ValFT ATM	0.67	0.75	0.47	0.70	0.54	0.64	0.62	0.62
	ATM	0.68	0.76	0.44	0.74	0.59	0.66	0.67	0.65
	Task Arithmetic	0.58	0.65	0.38	0.66	0.42	0.55	0.58	0.55
	TIES	0.33	0.50	0.33	0.63	0.38	0.51	0.53	0.46
1	Breadcrumbs	0.58	0.65	0.38	0.66	0.41	0.55	0.59	0.55
	DARE	0.83	0.52	0.66	0.43	0.33	0.51	0.46	0.53
ł	ValFT ATM	0.66	0.75	0.47	0.70	0.59	0.67	0.64	0.64
	ATM	0.69	0.81	0.51	0.78	0.63	0.70	0.66	0.68
	Task Arithmetic	0.59	0.66	0.39	0.64	0.42	0.55	0.59	0.55
0	TIES	0.35	0.51	0.34	0.64	0.38	0.51	0.53	0.47
Ξ	Breadcrumbs	0.59	0.65	0.39	0.64	0.42	0.55	0.60	0.55
	DARE	0.82	0.51	0.69	0.42	0.33	0.51	0.47	0.54
K	ValFT ATM	0.66	0.73	0.50	0.71	0.59	0.69	0.63	0.64
	ATM	0.68	0.81	0.53	0.78	0.65	0.72	0.66	0.69

Table 5: ATM vs. Baselines as budget varies (BERT-base-uncased)



Figure 11: Pairwise cosine similarities of the gradients of the first 10 epochs over datasets that do not exhibit most of the gradient norm localized in the first epoch.