INTERFERING WITH INTERFERENCE: BLIND SHUF FLING AND SUPERPOSITION FOR BETTER MULTI MODEL COMPRESSION

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

023

Paper under double-blind review

ABSTRACT

We present two complementary random mechanisms to significantly reduce interference when eliminating cross-model redundancy for efficient multi-model serving: *Layer Shuffling* and *Task Vector Superposition*. They work together to increase the orthogonality among interfering task vectors, forcing them into selfdestruction without requiring any post-training learning or optimization. *Layer Shuffling* randomly reorders layers of each individual models to reduce the alignment between interfering task vectors. While *Task Vector Superposition* leverages random orthogonal transformations to decorrelate task vectors further. Together, these techniques drastically minimize interference, yielding improved performance across multiple tasks with effectively zero incremental memory cost when incorporating new models. Their data and model-independent nature also allows for seamless on-the-fly addition or removal of models, without requiring any re-computation, making them highly practical for real-world deployment scenarios.

$$\begin{split} M_i^{ta} & \hat{M}_i^{ta} = (1 - \lambda) M_0 + \lambda M_i + \lambda \sum_{j=1}^2 T_j \\ M_i^{T_1} & \hat{M}_i^{our} = (1 - \lambda) M_0 + \lambda M_i + \lambda \sum_{j=1}^2 T_j' \\ T_1' & Cos(T_1', T_2') < Cos(T_1, T_2) \\ \hat{M}_i^{our} & T_i \\ \hat{M}_i^{our} & M_0 & \left\| \lambda \sum_j^2 T_j' \right\|_F < \left\| \lambda \sum_j^2 T_j \right\|_F \end{split}$$

Figure 1: Illustration of interference reduction in multi-model compression. M_0 is the pre-trained checkpoint, and M_i the *i*-th fine-tuned checkpoint, with task vectors T_i , T_1 , and T_2 . Standard task arithmetic (\hat{M}_i^{ta} , red) sums aligned task vectors, causing interference. With S_2 (\hat{M}_i^{our} , blue), layer shuffling and superposition decorrelate the interfering task vectors into T'_1 and T'_2 , lowering the Frobenius norm of interference. This allows better retrieval of M_i with a higher merging coefficient λ .

042 043

040

041

037

044

1

INTRODUCTION

045

Contemporary advances in machine learning are fueled by ever larger models. Language models and multimodal language models now run into billions of parameters (Cohen & Gokaslan, 2020;
Radford et al., 2019; 2021; Workshop et al., 2022). These models are often finetuned into task-specific models to capture the intricacies of individual tasks. As the number of these large models proliferates, serving them becomes a challenge. It is no longer possible to even store multiple models, even in high-end GPUs, significantly impacting downstream applications. Approaches to compress these models without losing accuracy are thus becoming increasingly important. When there are multiple models finetuned from the same pre-trained checkpoint on potentially related tasks, one would expect that the models have a lot of redundant information and can be compressed

together. In fact, a line of work has shown that all these models can be *merged* together into a single model that can tackle all the tasks involved (Ainsworth et al., 2022; Frankle et al., 2020; Wortsman et al., 2020; 2022; Li et al., 2024; Yadav et al., 2024; Tang et al., 2024c; Ilharco et al., 2022; Yang et al., 2023). Of these, a popular framework is task arithmetic (Ilharco et al., 2022), which computes the difference between finetuned model weights and pre-trained model weights to produce task vectors for each task, and add these task vectors together with the pretrained model weights to yield a merged model.

However, the accuracy of these merged models still lags behind the accuracy of the original fine-tuned models. Prior work has identified as a potential reason the interference between the different tasks, which may not be perfectly correlated with each other (Yadav et al., 2024; Wang et al., 2024;
While many techniques have been proposed to limit interference, it has generally been difficult to reduce this interference.

066 In this paper, we take a renewed look at this interference, and find that the cause of this interference 067 is not that the models involved are too different, but *that they are too similar*. Concretely, we find that 068 task arithmetic works best when the task vectors are as orthogonal to each other as possible. Armed 069 with this insight, we propose two new ways of improving upon task arithmetic. Our first approach is to *shuffle* task vectors across the layers of each model before combining them, with an inverse 071 shuffle applied at test time. Our second approach is to apply a random sign flip or rotation to the task vectors before merging them, again inverting the transformation at test time. Both approaches 072 significantly reduce interference between the task vectors. They also have the advantage of being 073 simple, efficient and requiring no training or optimization. 074

We test our approach on three different benchmarks involving large models and their finetuned versions: CLIP-ViT-B/32 and CLIP-ViT-L/14 for zero-shot image classification (Radford et al., 2021), Flan-T5-base for text generation (Longpre et al., 2023), and GPT-2 for text classification (Radford et al., 2019). We find that across all of these benchmarks, our approach substantially improves in terms of accuracy over prior model merging-based approaches. When compared to the original finetuned models, in two of the three benchmarks our approach yields near-identical accuracy to the individual models while reducing the storage costs by $4 \times .$ In sum, our contributions are:

- 1. We provide an analysis of the interference between tasks in task arithmetic, which suggests that similarity between the task vectors may be a problem.
- 2. We propose two complementary strategies for reducing interference. Our first strategy randomly shuffles parameter matrices across layers. Our second strategy applies a random rotation or a sign flip to the task vectors before merging.
- 3. We demonstrate through experiments on three benchmarks that our approach compresses multiple models together and achieves much higher accuracy than prior model merging based approaches.

2 PROBLEM SETUP

082

084

085

090 091

092 093

We are given T models $\{\Theta_i\}_{i=1}^T$ fine-tuned from a single pre-trained model Θ_0 on tasks $i = 1, \ldots, T$. Each model Θ_i as a set of parameter matrices:

$$\mathbf{\Theta}_{i} = \left\{ \mathbb{I}_{i}, \left(\boldsymbol{M}_{i}^{k,1}, \boldsymbol{M}_{i}^{k,2}, \dots, \boldsymbol{M}_{i}^{k,m_{k}}
ight)_{k=1}^{K}, \mathbb{O}_{i}
ight\}.$$

Here \mathbb{I}_i and \mathbb{O}_i are the input and output layers. Each model has K blocks, with the k-th block containing m_k matrices $M_i^{k,1}, M_i^{k,2}, \ldots, M_i^{k,m_k}$.

101 102 Our goal is to compress $\{\Theta_i\}_{i=1}^T$ into a compact representation $\Theta_* = \text{compress}(\{\Theta_i\}_{i=1}^T)$ with 103 minimal memory usage, so that at test time, given a task *i*, we can retrieve an *approximate* model $\hat{\Theta}_i = \text{retrieve}(\Theta_*, i)$ for the task that achieves high accuracy on this task.

- One way to address this problem is obviously to compress each individual model using strategies such as pruning or quantization. However, here we are interested in techniques that can leverage the structure of the problem (namely, T models finetuned from the same source) to yield storage that is *sub-linear* in T.
 - 2

¹⁰⁸ 3 TASK ARITHMETIC

110 A promising line of approaches for this problem derives from the observation that models fine-111 tuned for different tasks could be *merged* into a single model that gives reasonable accuracy for 112 all tasks (Ilharco et al., 2022). Concretely, their framework, called task arithmetic, first computes 113 the difference between the finetuned model weights for each layer k, M_i^k , and the corresponding 114 pre-trained model weights, M_0^k , to produce *task vectors* $T_i^k = M_i^k - M_0^k$. Task arithmetic then 115 computes a weighted average of the pretrained model weights and the task vectors:

116 117

118 119

120

123 124

$$\boldsymbol{M}_{\star}^{k} \leftarrow \boldsymbol{M}_{0}^{k} + \lambda \sum_{i=1}^{T} \boldsymbol{T}_{i}^{k} = \boldsymbol{M}_{0}^{k} + \lambda \sum_{i=1}^{T} (\boldsymbol{M}_{i}^{k} - \boldsymbol{M}_{0}^{k}), \tag{1}$$

$$\boldsymbol{\Theta}_{\star}^{TA} \leftarrow \{\boldsymbol{M}_{\star}^{k}\}_{k=1}^{K}, \tag{2}$$

where $\lambda \in \mathbb{R}^+$ is the merging coefficient. At test time, this compressed model is directly applied no matter what the task:

$$\hat{\boldsymbol{\Theta}}_i \leftarrow \boldsymbol{\Theta}_\star^{TA} = \{\boldsymbol{M}_\star^k\}_{k=1}^K.$$
(3)

This approach can be used to reduce model storage by a factor of T since we only need to store one model instead of T different models. However, as we show later, this yields much lower accuracy than the original fine-tuned models.

One reason that has been put forward for the low accuracies offered by task arithmetic is *task inter-ference*: different tasks may want to set particular parameters differently, and merging these parameters naively will cause one task to harm another task's accuracy (Yadav et al., 2024; Wang et al., 2024; Tang et al., 2023). However, a mathematical analysis of this interference is missing. Below, we delve deeper into this interference, and find a counter-intuitive solution.

134 3.1 INTERFERENCE IN TASK ARITHMETIC

To understand the interference term, let us consider what happens when we apply the merged model to task i:

$$\hat{\boldsymbol{M}}_{i}^{k} = \boldsymbol{M}_{\star}^{k}, \qquad \text{(from equation 3)}$$
$$= \boldsymbol{M}_{0}^{k} + \lambda \sum_{i=1}^{T} (\boldsymbol{M}_{i}^{k} - \boldsymbol{M}_{0}^{k}),$$
$$= (1 - \lambda) \boldsymbol{M}_{0}^{k} + \lambda \boldsymbol{M}_{i}^{k} + \lambda \sum_{j \neq i} \boldsymbol{T}_{j}^{k}. \qquad (4)$$

143 144 145

The last line suggests that the model being applied to the *i*-th task is interpolating between the pretrained model M_0^k and the task-specific finetuned model M_i^k , but with an additional interference coming from other merged models : $\lambda \sum_{j \neq i} T_j^k$. To retrieve the finetuned model, the first two terms suggest that we should set λ to 1. However this will *increase* the interference term, $\lambda \sum_{j \neq i} T_j^k$. Some prior work has tried to achieve a good balance by optimizing λ for each model and layer using test-time adaptation(Yang et al., 2023), but attaining this balance has often been challenging.

Instead of focusing on λ , let us look at this interference term in greater detail by analyzing its Frobenius norm:

154

$$\left\| \lambda \sum_{j \neq i} \mathbf{T}_{j}^{k} \right\|_{F}^{2} = \lambda^{2} \left(\sum_{j \neq i} \|\mathbf{T}_{i}^{k}\|_{F}^{2} + 2 \sum_{\substack{1 \leq l < j \leq n \\ l, j \neq i}} \|\mathbf{T}_{l}^{k}\|_{F} \|\mathbf{T}_{j}^{k}\|_{F} \cos(\mathbf{T}_{l}^{k}, \mathbf{T}_{j}^{k}) \right).$$
(5)

160 We observe that the interference term is directly correlated with two quantities: the magnitude of 161 the task vectors T_i^k (which is out of our control since it depends on the task-specific finetuning), and 162 the cosine products between them. Interestingly, the interference term is maximum when the task



Figure 2: Average pairwise cosine similarity of three out of eight CLIP-ViT-B/32 task vectors during model retrieval for SUN397 across three repetitions. Both random layer shuffling and superposition increase mutual orthogonality, with an additive effect when combined.

176

177

178

179 180 181

182

172

vectors are very closely aligned with each other. Thus, the problem with task arithmetic is not that the individual task vectors are very different from each other, but that they are *too similar*.

Our goal, therefore, should be to make the task vectors as different from each other as possible. Below, we propose two strategies for doing this.

4 Methodology

As described above, to minimize interference, we want the task vectors to be as orthogonal to each other as possible. We propose two complementary random algorithms to achieve this: *random layer shuffling* and *task vector superposition*.



Figure 3: (a) Layer shuffling illustration in a three-layer model with three checkpoints. Different task vectors (distinct arrowheads) initially align across layers, causing interference when directly merged. Shuffling layers followed by inverse transformation retain the target task's orientation (standard arrowhead) while reducing interference through increased orthogonality among other vectors. (b) shows cosine similarity distributions between task vector layers within and across models for CLIP-ViT-B/32, Flan-T5, and GPT-2, with standard error of mean (SEM) as error bars.

207 208

209

4.1 RANDOM LAYER SHUFFLING

Across several model architectures (CLIP-ViT-B/32, CLIP-ViT-L/14 (Radford et al., 2021), Flan-T5 (Longpre et al., 2023), and GPT-2 (Radford et al., 2019)), we observed that task vector layers within the same model exhibit greater variability compared to corresponding layers across finetuned models (see Figure 3 (b)). This insight suggests that by randomly shuffling layers across different task vectors, we can reduce the pairwise cosine similarity of interfering task vectors and thus minimize their contribution to the interference. To that end, we propose *random layer shuffling* as a simple fix to the problem. Method Description. The models we consider are made of multiple parameter blocks of similar structure. For example, transformers have multiple MLP layers and multiple attention layers. Many of the MLP and Attention layers have weight matrices of the same size.

Our proposal is that when merging the task vectors, for each model, we first *randomly permute* the task vectors across layers of the same type and with the same dimensions of parameter matrices (as illustrated in Figure 3 (a)). Then, when we want to perform inference on a particular task, we perform the inverse of the corresponding permutation to obtain the model for the task.

Concretely, for each task *i*, we produce a random permutation of the layers σ_i , taking care to only permute across layers of the same type and same dimensionality. We then produce merged task vectors by adding up these shuffled task vectors across models. The merged task vector for the *k*-th layer is:

$$\boldsymbol{T}_{\star}^{k} \leftarrow \sum_{i=1}^{T} \boldsymbol{T}_{i}^{\sigma_{i}(k)}.$$
(6)

The number of such merged task vectors is equal to the total number of layers K. We then store both the pretrained model Θ_0 and the merged task vectors $\{T_{\star}^k\}_{k=1}^K$:

$$\boldsymbol{\Theta}_{*}^{Shuffle} \leftarrow \left(\boldsymbol{\Theta}_{0}, \{\boldsymbol{T}_{\star}^{k}\}_{k=1}^{K}\right).$$

$$\tag{7}$$

Reduction in Interference: Because parameter vectors from different layers are less likely to align, we effectively reduce the cosine product between the task vectors being merged: instead of the term $\cos(T_i^k, T_j^k)$ in the interference term (equation 5), we now have the product $\cos(T_i^{\sigma_i(k)}, T_j^{\sigma_j(k)})$, which is expected to be significantly lower due to the reduced alignment of parameter vectors from different layers. We thus expect smaller interference and thus more faithful retrieval of model weights for each task. The first two parts of Figure 2 shows this effect in action, with layer shuffling reducing the pairwise cosine similarity among interfering vectors unanimously.

4.2 TASK VECTOR SUPERPOSITION

We can also leverage the *blessing of dimensionality* (Gorban & Tyukin, 2018) to promote orthogo nality among high dimensional vectors. We take inspiration from Cheung et al. (2019) on continual
 learning and introduce *superposition* as a complementary approach to increase the mutual orthogo nality among interfering task vectors.

Method Description. Considering merging the parameters of layer k, we sample a random binary diagonal matrices whose diagonal entries have equal probability to be +1 or -1 to each of the Ttask vectors and apply them to the vectors before summation:

$$T^k_{\star} \leftarrow \sum_{i=1}^T T^k_i C^k_i.$$
 (8)

We call them context matrices and $\forall i \in [1, \dots, T]$, $C_i^k C_i^{k(T)} = C_i^k C_i^{k(-1)} = I$.

When performing task *i*, we apply the inverse transformation $C_i^{k(-1)}$ to retrieve task vector T_i^k from the superposition:

$$\hat{\boldsymbol{T}}_{i}^{k} = \boldsymbol{T}_{\star}^{k} \boldsymbol{C}_{i}^{k(-1)}, \tag{9}$$

$$=\sum_{i=1}^{I} [T_i^k C_i^k] C_i^{k(-1)},$$
(10)

$$= T_i^k + \sum_{j \neq i} [T_j^k C_j^k C_i^{k(-1)}].$$
(11)

We store both the pretrained model Θ_0 , the merged task vectors $\{T^k_{\star}\}_{k=1}^K$, as well as the context matrices $\{C^k_{\star}\}_{k=1}^K$:

$$\boldsymbol{\Theta}_{*}^{Superpose} \leftarrow \left(\boldsymbol{\Theta}_{0}, \{\boldsymbol{T}_{\star}^{k}\}_{k=1}^{K}, \{\boldsymbol{C}_{\star}^{k}\}_{k=1}^{K}\right).$$
(12)

262 263 264

265 266

260 261

253 254

228 229 230

231

236

237

238 239

240

241 242

Reduction in Interference. Two random vectors in high dimensional space is very likely to be nearly orthogonal with each other. Now the cosine similarity in equation 5 changes from $\cos(T_i^k, T_j^k)$ to $\cos(T_i^k C_i C_j^{l(-1)}, T_j^k C_i C_j^{l(-1)})$ when performing task *l*. The randomly sampled diagonal binary matrices $\{C_{\star}^k\}_{k=1}^K$ will randomize the task vectors, leading to more orthogonal task vectors, smaller interference, and thus better retrieval of model parameters for the task at hand. Again, Figure 2 confirmed the cosine similarity reduction when task vectors are superposed together.

5 EXPERIMENTS

We evaluate our methods on FusionBench (Tang et al., 2024a) across vision and language tasks, showing comparable performance for both discriminative and generative models. Through ablation studies, we analyze component importance, merging coefficient effects, and context matrix designs.
Finally, we demonstrate broader applications including PEFT model compression and large-scale merging of twenty CLIP-ViT-L/14 models.

5.1 EXPERIMENT SETUP

Datasets and Models. We follow Tang et al. (2024a) and select three representative scenarios to evaluate our methods. This includes i) CLIP-ViT-B/32 fine-tuned on eight image classification datasets (adopted from (Ilharco et al., 2022)); ii) Flan-T5-base fine-tuned on eight text generation datasets; and iii) GPT-2 fine-tuned on seven text classification datasets. Detailed information on the datasets and models is in Appendix B.

292

277

278

285

286

Baselines and Metrics. We evaluate baselines from model merging/compression literature, grouped by memory requirements: methods using the original footprint (pre-trained model, standard merging techniques) and those requiring additional memory (fine-tuned models, newer merging baselines). The pre-trained and fine-tuned models provide lower and upper performance bounds respectively. Following Ilharco et al. (2022), we optimize merging coefficient λ via validation set grid search. We report accuracy and memory usage across three runs per experiment with random operations. See Appendix B for details.

300 301

5.2 PERFORMANCE ANALYSIS

Table 1: Performance and memory comparison of CLIP-ViT-B/32 models across eight image classification tasks, showing absolute and normalized accuracy (%), as well as memory footprint (Gb). Results averaged over three runs where applicable. Variances smaller than 0.1% are omitted.

Method	Avg.(%) \uparrow	$Bits(Gb)\downarrow$	SUN397	Cars	RESISC45	EuroSAT	SVHN	GTSRB	MNIST	DTD
Pre-trained	48.2 (53.4)	0.564 (1.00)	63.2	59.8	60.7	46.0	31.6	32.5	48.3	43.9
Weight Averaging	66.5 (73.6)	0.564 (1.00)	65.4	62.6	70.8	76.9	64.5	54.9	86.3	50.9
Fisher Merging	70.6 (78.2)	0.564 (1.00)	66.7	64.0	72.2	91.6	69.0	64.3	83.5	53.7
RegMean	80.5 (89.1)	0.564 (1.00)	67.8	68.9	82.5	94.4	90.6	79.2	94.7	63.2
Task Arithmetic	69.8 (77.2)	0.564 (1.00)	64.4	61.5	70.5	80.4	73.9	62.8	93.0	51.6
Ties-Merging	72.2 (80.0)	0.564 (1.00)	67.1	64.2	74.1	91.6	77.7	69.4	94.1	54.0
Layerwise AdaMerging	82.6 (91.5)	0.564 (1.00)	67.9	71.3	83.5	92.7	87.4	92.9	98.2	67.0
PSP	4.5 (5.0)	0.564 (1.00)	0.3	0.5	1.9	10.4	8.8	2.3	9.8	1.9
Fine-tuned	90.3 (100)	2.84 (5.03)1	75.0	78.3	95.2	99.0	97.3	98.9	99.6	79.7
WEMoE	89.2 (98.8)	2.27 (4.03)	73.7	76.8	93.4	98.2	96.8	98.2	99.6	76.6
SMILE	89.3 (98.9)	1.23 (2.20)	73.6	77.8	92.0	98.3	96.9	98.1	99.6	78.1
TA+Shuffle (Ours)	81.3 (90.0)	0.89 (1.58)	65.6	58.5	86.8	94.5	93.2	91.4	98.5	62.2
STA (Ours)	<u>89.6</u> (99.2)	0.89 (1.58)	74.4	75.6	<u>94.6</u>	99.0	<u>97.1</u>	<u>98.5</u>	<u>99.5</u>	77.8
STA+Shuffle (Ours)	89.9 (99.6)	0.89 (1.58)	74.8	76.7	94.8	99.0	97.2	98.6	99.5	78.7

316 317

318

Superior MTL Performance. Our approach achieves significant accuracy gains across bench marks (Tables 1, 2 and 6), with *STA+Shuffle* nearly matching individual fine-tuned models. We
 outperform WEMOE (Tang et al., 2024c) and SMILE (Tang et al., 2024b) on image classification
 while using only 40% and 72% of their respective memory footprints, and surpass SMILE's text generation performance at benchmark saturation. Though Task Arithmetic (Ilharco et al., 2022) uses
 of our storage, its performance is substantially lower. Parameter Superposition (Cheung et al.,

Method	Avg.(%) \uparrow	$Bits(Gb) \downarrow$	CoLA	MNLI	MRPC	QNLI	QQP	RTE	SST2	STSB
Pre-trained	75.7 (87.6)	1.19 (1.00)	69.1	56.5	76.2	88.4	82.1	80.1	91.2	62.2
Weight Averaging	78.9 (91.3)	1.19 (1.00)	69.1	62.6	79.4	89.8	83.9	81.2	91.7	73.2
Task Arithmetic	79.6 (92.1)	1.19 (1.00)	69.7	64.1	79.2	90.2	83.9	81.6	92.1	76.4
Ties-Merging	79.9 (92.5)	1.19 (1.00)	70.3	65.0	78.9	90.2	83.5	81.6	91.7	78.3
PSP	0.0 (0.0)	1.19 (1.00)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	N/A
Fine-tuned	86.4 (100)	9.52 (8.00)	75.0	83.4	87.5	91.5	85.4	85.9	93.6	88.7
SMILE	85.5 (99.0)	1.81 (1.52)	73.2	84.2	85.0	91.3	84.9	84.8	<u>93.5</u>	87.3
TA+Shuffle (Ours)	85.7 (99.0)	2.38 (2.00)	75.5	82.0	87.5	91.1	83.9	83.8	93.6	88.4
STA (Ours)	86.5 (100)	2.38 (2.00)	77.2	82.1	87.6	<u>91.6</u>	85.3	85.7	93.2	89.0
STA+Shuffle (Ours)	<u>86.4</u> (100)	2.38 (2.00)	<u>75.6</u>	<u>82.8</u>	88.2	91.7	85.3	85.7	<u>93.5</u>	88.9

Table 2: Performance and memory comparison of Flan-T5-base models across eight GLUE text generation
 tasks, showing absolute and normalized accuracy (%), as well as memory footprint (Gb). Results averaged
 over three runs where applicable. Variances smaller than 0.1% are omitted.

2019), while effective for continual learning, underperforms here, demonstrating the importance of task vector superposition for offline compression.

Amortizable Memory Overhead. Our method requires only 2x memory, mainly from storing merged task vectors $\{T_{\star}^k\}_{k=1}^K$ and binary context matrices $\{C_{\star}^k\}_{k=1}^K$ (equations 7, 12). By storing random seeds and regenerating context matrices on-the-fly with minimal overhead (292.70 ms for CLIP-ViT-B/32, 658.19 ms for CLIP-ViT-L/14 on Intel Xeon Gold 6448Y CPU), we achieve effectively zero additional memory per model. This enables efficient scaling, demonstrated by merging 20 CLIP-ViT-L/14 models with state-of-the-art performance and 9x memory reduction (Sec. 5.9).

5.3 Key Components Ablation

In this section, we ablate both the *random layer shuffling* and *superposition* to show their individual contribution. Specifically, we derive two variants:

- **TA+Shuffle**: we randomly shuffle the layers before task arithmetic without performing superposition.
- STA: we superpose task vectors without layer shuffling.

As shown in Tables 1, 2, 4, 5 and 6, both methods significantly outperform task arithmetic consistently. In some benchmarks, shuffling works better (Flan-T5-base and GPT-2) while superposition works better in others (CLIP). This difference may be because of the nature of task vectors themselves: how they vary across the model layers and across different models. We find that the combined approach is able to combine gains from both components, yielding consistently the best result across all the benchmarks. This complementary effect is further manifested in Figure 2, where introducing both mechanism gets the smallest pairwise cosine similarity among interfering task vectors. We provide detailed analysis of the interplay between shuffling and superposition in Section D, including how different implementation choices affect model performance.

366 367

368

327 328

340

341 342

343 344

345

346

347 348 349

350

351

352 353

354

355

356

5.4 IMPACT OF MERGING COEFFICIENT λ

Here we examine the interplay between the merging coefficient λ and the average performance across different setup. For each variant derived in section 4, we perform a grid search on $\lambda = \{0.1, 0.2, \dots, 1.0\}$ when compressing eight CLIP-ViT-B/32 models for image classification. Figure 4 shows the change of optimal model performance and the coefficient λ when *layer shuffling* and *superposition* are introduced to task arithmetic.

We observe that when shuffling and superposition are introduced, the best performance increases along with the value of λ . This shows the effectiveness of our method in reducing interference, allowing larger λ to be selected for more authentic model retrieval as according to equation 4.



Figure 4: The impact of λ on average accuracy over eight image classification tasks.

378 5.5 IMPACT OF CONTEXT MATRIX DESIGN

To further examine how *task vector superposition* works and shad light on better context matrix design, we make a comparison between three types of context matrices: random binary diagonal matrix with $\{-1, +1\}$ entries (RBD), identity matrix (Identity), and random diagonal matrix with entries draw from Normal distribution (RD). We use random layer shuffling when compressing the 8 CLIP-ViT-B/32 models on the image classification tasks.

The average accuracy and its variance with the optimal merging coefficient is shown in Figure 5. RBD receives higher accuracy than Identity due to the randomness it introduces, which reduces interference as discussed in section 4.2. Despite being random, RD's accuracy is much lower than RBD. We think this happens because RD is not an orthogonal matrix. It fails to preserve the Frobenius norm of $T_j^k C_j^k C_i^{k(-1)}$ and thus disturb this self-cancellation process.



Figure 5: (a) Impact of context matrix design to the average accuracy. RBD stands for random binary diagonal; RD stands for random diagonal. (b) Impact of target layer selection to the average accuracy. ALL stands for choosing all layers; MLP means only MLP layers are selected; and ATTN stands for attention layers.

5.6 TARGET LAYER SELECTION

By default we apply the random operations on all layers within the models. In this section, we evaluate the benefits of targeting specific types of layers. To do this, we create two variants: MLP (which selects only the MLP layers) and ATTN (which selects only the attention layers), in addition to the default setup (ALL). Figure 5 shows the average accuracy for each setup across eight image classification tasks using CLIP-ViT-B/32. The ALL configuration achieves the highest accuracy, followed by MLP and ATTN. Note that the total number of parameters in MLP is twice that of ATTN, explaining the gradual decline in performance as fewer parameters are selected.

5.7 MODEL HOT SWAPPING

416 The ability to hot-swap models in real-world applications is crucial, especially in dynamic environ-417 ments like model serving, where new models need to be integrated into the system regularly, and deprecated ones need to be removed in a timely fashion. As mentioned, the STA+Shuffle method 418 allows for this by shuffling layers and sampling diagonal binary matrices *independently* of data or 419 model parameters, thus enabling the on-the-fly addition of new models without the need for re-420 computation. This provides our a method a big advantage over methods like WEMoE which require 421 recomputation of the router when new models are added (Tang et al., 2024c), or TALL-masks, which 422 also needs to recompute binary masks when new models are added (Wang et al., 2024). We dub this 423 feature hot swapping to borrow a term from the hardware literature.

424 425 426

414

415

390 391

392

393

396 397

398 399

400

401

402

403 404 405

5.8 PARAMETER EFFICIENT FINETUNING (PEFT) MODEL COMPRESSION

427 We also apply our method on PEFT adapter weights. Consider a LoRA (Hu et al., 2021), where we 428 have a fixed pre-trained model Θ_0 , along with LoRA weights L_i . We merge the LoRA weights to 429 get the fine-tuned model: $\Theta_i = \Theta_0 + \lambda L_i$. Similar to section 4.1 and 4.2, we apply *random layer* 430 *shuffling* and *superposition* on these LoRA weight vectors before retrieval.

¹CLIP models' text encoder is frozen and shared by all fine-tuned models.

Table 3: Comparison of selected methods with hot adding and recomputation requirements when new models are added to the pool.

434		H (C	D (()
435	Method	Hot Swap	Recomputation
436	Task Arithmetic	\checkmark	X
407	WEMoE	×	\checkmark
437	TALL-masks	×	\checkmark
438	STA+Shuffle	\checkmark	X
400			

Experiments on Flan-T5-base LoRA fine-tunes (Longpre et al., 2023; Tang et al., 2024a;b) demonstrate that our method is performative in PEFT compression settings as well (Table 4). With 99.8% normalized average accuracy compare to the fine-tuned baseline, and 1.20 Gb memory usage, our method presents a better trade-off point between performance and storage usage than the state-ofthe-art model SMILE (Tang et al., 2024b).

Table 4: Performance and memory comparison of Flan-T5-base LoRA models across eight GLUE text generation tasks, showing absolute and normalized accuracy (%), as well as memory footprint (Gb). Results averaged over three runs where applicable. Variances smaller than 0.1% are omitted.

Method	Avg.(%) \uparrow	$Bits(Gb) \downarrow$	CoLA	MNLI	MRPC	QNLI	QQP	RTE	SST2	STSB
Pre-trained	75.7 (87.6)	1.19 (1.00)	69.1	56.5	76.2	88.4	82.1	80.1	91.2	62.2
Weight Averaging	78.2 (92.4)	1.19 (1.00)	69.7	59.7	78.9	90.1	83.8	80.5	91.2	72.0
Task Arithmetic	77.4 (91.5)	1.19 (1.00)	68.8	55.2	78.7	89.8	83.7	79.1	91.5	72.4
Ties-Merging	77.5 (91.6)	1.19 (1.00)	68.3	56.3	79.4	89.8	83.7	79.4	91.6	71.2
Fine-tuned	84.6 (100)	1.25 (1.05)	69.1	82.7	85.5	90.9	84.0	84.4	92.9	87.4
SMILE	<u>84.0</u> (<u>99.3</u>)	1.21 (1.02)	69.3	82.9	83.8	<u>90.6</u>	83.9	83.4	93.1	85.1
TA+Shuffle (Ours)	83.9 (99.2)	1.20 (1.01)	69.2	79.0 ± 0.3	84.2	90.4	84.1	85.0	92.9	86.5
STA (Ours)	83.0 (98.1)	1.20 (1.01)	69.1	81.3	82.2	90.5	83.2	79.1	92.7	85.6
STA+Shuffle (Ours)	84.4 (99.8)	1.20 (1.01)	69.1	<u>82.7</u>	85.0	90.9	83.8	<u>84.2</u>	92.7	86.9

5.9 SCALABILITY ANALYSIS

Our method scales effectively to merging larger models and more tasks, as demonstrated on CLIP ViT-L/14 with 8, 14, and 20 image classification tasks (Table 5). STA+Shuffle achieves near finetuned performance (93.5% vs 94.2% for 20 tasks) while maintaining constant 2.87GB storage regardless of task count. In contrast, TALL-masks+TA (Wang et al., 2024) requires progressively more storage (5.42GB to 9.25GB) as tasks increase. Though Task Arithmetic uses only 1.59GB, its performance drops significantly with more tasks. Model retrieval remains efficient, requiring just 658.19ms per CLIP ViT-L/14 model on an Intel Xeon Gold 6448Y CPU.

Table 5: Performance and memory comparison of CLIP ViT-L/14 models across three test scenarios with 8, 14, and 20 image classification tasks, showing absolute and normalized accuracy (%), as well as memory footprint (Gb). Results averaged over three runs where applicable. Variances smaller than 0.1% are omitted.

Method	8 ta	asks	14 1	tasks	20 tasks		
	Acc.(%) ↑	Bits(Gb)↓	Acc.(%) ↑	Bits(Gb) \downarrow	Acc.(%) ↑	Bits(Gb)↓	
Pre-trained	64.5 (68.3)	1.59 (1.00)	68.1 (72.8)	1.59 (1.00)	65.2 (69.2)	1.59 (1.00)	
Task Arithmetic	84.0 (88.7)	1.59 (1.00)	79.1 (84.2)	1.59 (1.00)	73.8 (78.3)	1.59 (1.00)	
Fine-tuned	94.4 (100)	10.53 (6.62)	93.5 (100)	18.18 (11.43)	94.2 (100)	25.84 (16.25	
Magnitude Masking	92.8 (98.2)	5.42 (3.41)	90.6 (96.7)	7.34 (4.62)	90.9 (96.4)	9.25 (5.82)	
TALL Mask+TA	94.2 (99.7)	5.42 (3.41)	92.4 (98.8)	7.34 (4.62)	93.2 (98.9)	9.25 (5.82)	
TA+Shuffle (Ours)	93.0 (98.4)	2.87 (1.81)	88.8 (94.6)	2.87 (1.81)	87.1 ±0.1(92.2)	2.87 (1.81)	
STA (Ours)	94.2 (99.8)	2.87 (1.81)	92.8 (99.2)	2.87 (1.81)	93.4 (99.1)	2.87 (1.81)	
STA+Shuffle (Ours)	94.3 (99.9)	2.87 (1.81)	93.0 (99.5)	2.87 (1.81)	93.5 (99.3)	2.87 (1.81)	

486 6 RELATED WORK

488 489

490 **Model Merging.** Recent research on model merging is largely founded on *linear mode connectiv*-491 ity (LMC) (Frankle et al., 2020; Neyshabur et al., 2020), which posits that models fine-tuned from 492 the same pre-trained model are connected by a linear path along which performance remains con-493 stant. Building upon this concept, Wortsman et al. (2022) and Li et al. (2024) demonstrated that a set 494 of specialist models can be directly interpolated to obtain a multi-task model. Ilharco et al. (2022) proposed interpolating the parameter deltas (referred to as "task vectors") instead. However, these 495 methods suffer from *task interference*: when different models adjust the same parameters in conflict-496 ing ways, summing these adjustments leads to interference and degraded performance on individual 497 tasks (Yadav et al., 2024; Tang et al., 2024b; Wang et al., 2024). To mitigate this interference, various 498 strategies have been proposed. Yang et al. (2023) optimized the merging coefficients for different 499 tasks and layers to reduce interference. Yadav et al. (2024) addressed the conflict by removing re-500 dundant parameters and resolving sign disagreements. Tang et al. (2024c) reduced interference by 501 upscaling the multilayer perceptron (MLP) layers. Tang et al. (2024b) compressed task vectors us-502 ing singular value decomposition (SVD) and performed routing between them to further diminish 503 interference. Both Wang et al. (2024) and Yu et al. (2024) sparsified the task vectors to prevent 504 task conflicts. Additionally, Ortiz-Jimenez et al. (2024) proposed fine-tuning the linearized model 505 along the tangent space of the pre-trained model to promote weight disentanglement and avoid interference. In contrast to the above mentioned methods that aim to avoid conflicts, we intentionally 506 accumulate interference among conflicting task vectors to facilitate their mutual cancellation. 507

508 509

510 Model Compression. Model compression techniques aim to reduce the memory footprint of mod-511 els while maintaining their performance. Model pruning compresses neural networks by removing inessential parameters in either a structured (Anwar et al., 2017; Fang et al., 2023; He & Xiao, 2023; 512 Wang et al., 2019) or unstructured (Liao et al., 2023; Kwon et al., 2020) manner. Parameter quan-513 tization saves memory and speeds up inference by converting the weights and activation values of 514 a neural network from high precision to low precision (Gholami et al., 2022; Liu et al., 2021; Yuan 515 et al., 2022). Knowledge distillation reduces the memory footprint by training a smaller network to 516 mimic a larger network's behavior (Gou et al., 2021; Cho & Hariharan, 2019; Park et al., 2019; Zhao 517 et al., 2022). Leveraging the low-rank nature of model parameters, many works decompose weight 518 matrices into low-rank matrices for memory reduction (Yu et al., 2017; Li et al., 2023; Guo et al., 519 2024). Ryu et al. (2023) observed the low-rank nature of weight residuals in overparameterized 520 models and proposed reducing storage demands for fine-tuned models through low-rank approxi-521 mation of these residuals. Similarly, Tang et al. (2024b) compresses individual task vectors using 522 SVD and routes through a set of them conditioned on input. Our work differs from these works in 523 that we try to reduce redundancy across a set of aligned models rather than within them.

524 525

526

7 DISCUSSION AND FUTURE WORK

527 In this work, we introduce random layer shuffling and task vector superposition to enhance or-528 thogonality between task vectors, thereby significantly reduce task interference during multi-model 529 merging and compression. These data- and model-agnostic random operations enable users to i) 530 efficiently modify the model merging combinations without the need for additional training or op-531 timization; ii) merge additional models without increasing memory usage by saving random seeds. Evaluation on diverse model and task sets demonstrates that our method maintains high performance 532 while keeping a constant memory footprint as more and larger models are merged. These attributes 533 make our approach highly practical for real-world multi-model serving environments. 534

An interesting future direction is to further improve performance by increasing orthogonality, po tentially through alternative random operations or more systematic approaches. Our method relies
 on specific properties of model parameters that emerge from fine-tuning. Identifying these proper ties and enhancing fine-tuning strategies could lead to better merging and compression performance.
 Since we reduce cross-model redundancy, applying model compression algorithms could potentially further decrease memory footprint.

540 REFERENCES

547

553

554

555

570

571

572

573

579

- Samuel K Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. Git re-basin: Merging models
 modulo permutation symmetries. *arXiv preprint arXiv:2209.04836*, 2022.
- Sajid Anwar, Kyuyeon Hwang, and Wonyong Sung. Structured pruning of deep convolutional neural networks. ACM Journal on Emerging Technologies in Computing Systems (JETC), 13(3):1–18, 2017.
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101–mining discriminative components with random forests. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part VI 13*, pp. 446–461. Springer, 2014.
- Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Bench mark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017.
 - Brian Cheung, Alexander Terekhov, Yubei Chen, Pulkit Agrawal, and Bruno Olshausen. Superposition of many models into one. *Advances in neural information processing systems*, 32, 2019.
- Jang Hyun Cho and Bharath Hariharan. On the efficacy of knowledge distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4794–4802, 2019.
- M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), 2014.
- Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David
 Ha. Deep learning for classical japanese literature. *arXiv preprint arXiv:1812.01718*, 2018.
- Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised
 feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 215–223. JMLR Workshop and Conference Proceedings, 2011.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist
 to handwritten letters. In 2017 international joint conference on neural networks (IJCNN), pp.
 2921–2926. IEEE, 2017.
 - Vanya Cohen and Aaron Gokaslan. Opengpt-2: open language models and implications of generated text. XRDS, 27(1):26–30, September 2020. ISSN 1528-4972. doi: 10.1145/3416063. URL https://doi.org/10.1145/3416063.
- Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- Gongfan Fang, Xinyin Ma, Mingli Song, Michael Bi Mi, and Xinchao Wang. Depgraph: Towards any structural pruning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16091–16101, 2023.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode con nectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pp. 3259–3269. PMLR, 2020.
- Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A
 survey of quantization methods for efficient neural network inference. In *Low-Power Computer Vision*, pp. 291–326. Chapman and Hall/CRC, 2022.
- Jan J Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, et al. Challenges in representation learning: A report on three machine learning contests. In *Neural information processing: 20th international conference, ICONIP 2013, daegu, korea, november 3-7, 2013. Proceedings, Part III* 20, pp. 117–124. Springer, 2013.
- Alexander N Gorban and Ivan Yu Tyukin. Blessing of dimensionality: mathematical foundations of the statistical physics of data. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 376(2118):20170237, 2018.

594 Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. International Journal of Computer Vision, 129(6):1789–1819, 2021. 596 Yangyang Guo, Guangzhi Wang, and Mohan Kankanhalli. Pela: Learning parameter-efficient mod-597 els with low-rank approximation. In Proceedings of the IEEE/CVF Conference on Computer 598 Vision and Pattern Recognition, pp. 15699–15709, 2024. 600 Yang He and Lingao Xiao. Structured pruning for deep convolutional neural networks: A survey. 601 IEEE transactions on pattern analysis and machine intelligence, 2023. 602 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset 603 and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected* 604 *Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019. 605 606 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 607 and Weizhu Chen. Lora: Low-rank adaptation of large language models. arXiv preprint 608 arXiv:2106.09685, 2021. 609 Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, 610 Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. arXiv preprint 611 arXiv:2212.04089, 2022. 612 Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. Dataless knowledge fusion by 613 merging weights of language models. arXiv preprint arXiv:2212.09849, 2022. 614 615 Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained 616 categorization. In Proceedings of the IEEE international conference on computer vision work-617 shops, pp. 554-561, 2013. 618 Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009. URL https: 619 //api.semanticscholar.org/CorpusID:18268744. 620 621 Se Jung Kwon, Dongsoo Lee, Byeongwook Kim, Parichay Kapoor, Baeseong Park, and Gu-Yeon 622 Wei. Structured compression by weight encryption for unstructured pruning and quantization. 623 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 624 1909–1918, 2020. 625 Tao Li, Weisen Jiang, Fanghui Liu, Xiaolin Huang, and James T Kwok. Scalable learned model 626 soup on a single gpu: An efficient subspace training strategy. arXiv preprint arXiv:2407.03641, 627 2024. 628 Yixiao Li, Yifan Yu, Qingru Zhang, Chen Liang, Pengcheng He, Weizhu Chen, and Tuo Zhao. 629 Losparse: Structured compression of large language models based on low-rank and sparse ap-630 proximation. In International Conference on Machine Learning, pp. 20336–20350. PMLR, 2023. 631 632 Zhu Liao, Victor Quétu, Van-Tam Nguyen, and Enzo Tartaglione. Can unstructured pruning reduce 633 the depth in deep neural networks? In Proceedings of the IEEE/CVF International Conference 634 on Computer Vision, pp. 1402–1406, 2023. 635 Zhenhua Liu, Yunhe Wang, Kai Han, Wei Zhang, Siwei Ma, and Wen Gao. Post-training quanti-636 zation for vision transformer. Advances in Neural Information Processing Systems, 34:28092-637 28103, 2021. 638 639 Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V 640 Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. In International Conference on Machine Learning, pp. 22631–22648. PMLR, 641 2023. 642 643 Michael S Matena and Colin A Raffel. Merging models with fisher-weighted averaging. Advances 644 in Neural Information Processing Systems, 35:17703–17716, 2022. 645 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. 646 Reading digits in natural images with unsupervised feature learning. In NIPS workshop on deep 647 learning and unsupervised feature learning, volume 2011, pp. 4. Granada, 2011.

648 649 650	Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learn- ing? <i>Advances in neural information processing systems</i> , 33:512–523, 2020.
651 652 653	Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In 2008 Sixth Indian conference on computer vision, graphics & image processing, pp. 722–729. IEEE, 2008.
654 655 656	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.
658 659 660	Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 3967–3976, 2019.
661 662	Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In 2012 IEEE conference on computer vision and pattern recognition, pp. 3498–3505. IEEE, 2012.
664 665	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.
666 667 668 669	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021.
670 671 672	Simo Ryu, Seunghyun Seo, and Jaejun Yoo. Efficient storage of fine-tuned models via low-rank approximation of weight residuals. <i>arXiv preprint arXiv:2305.18425</i> , 2023.
673 674 675 676	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 conference on empirical methods in natural language processing</i> , pp. 1631–1642, 2013.
677 678 679 680	Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. Man vs. computer: Bench- marking machine learning algorithms for traffic sign recognition. <i>Neural networks</i> , 32:323–332, 2012.
681 682 683	Anke Tang, Li Shen, Yong Luo, Liang Ding, Han Hu, Bo Du, and Dacheng Tao. Concrete subspace learning based interference elimination for multi-task model fusion. <i>arXiv preprint arXiv:2312.06173</i> , 2023.
684 685 686	Anke Tang, Li Shen, Yong Luo, Han Hu, Bo Do, and Dacheng Tao. Fusionbench: A comprehensive benchmark of deep model fusion. <i>arXiv preprint arXiv:2406.03280</i> , 2024a.
687 688 689	Anke Tang, Li Shen, Yong Luo, Shuai Xie, Han Hu, Lefei Zhang, Bo Du, and Dacheng Tao. Smile: Zero-shot sparse mixture of low-rank experts construction from pre-trained foundation models. <i>arXiv preprint arXiv:2408.10174</i> , 2024b.
690 691 692	Anke Tang, Li Shen, Yong Luo, Nan Yin, Lefei Zhang, and Dacheng Tao. Merging multi-task models via weight-ensembling mixture of experts. <i>arXiv preprint arXiv:2402.00433</i> , 2024c.
693 694 695 696	Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant cnns for digital pathology. In <i>Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part II 11</i> , pp. 210–218. Springer, 2018.
697 698 699	Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understand- ing. <i>arXiv preprint arXiv:1804.07461</i> , 2018.
700 701	Ke Wang, Nikolaos Dimitriadis, Guillermo Ortiz-Jimenez, François Fleuret, and Pascal Frossard. Localizing task information for improved model merging and compression. <i>arXiv preprint</i> <i>arXiv:2405.07813</i> , 2024.

Ziheng Wang, Jeremy Wohlwend, and Tao Lei. Structured pruning of large language models. *arXiv* preprint arXiv:1910.04732, 2019.

- BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- Mitchell Wortsman, Vivek Ramanujan, Rosanne Liu, Aniruddha Kembhavi, Mohammad Rastegari,
 Jason Yosinski, and Ali Farhadi. Supermasks in superposition. *Advances in Neural Information Processing Systems*, 33:15173–15184, 2020.
- Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International conference on machine learning*, pp. 23965–23998. PMLR, 2022.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark ing machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 3485–3492, 2010. doi: 10.1109/CVPR.2010.
 5539970.
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. Ties-merging: Re solving interference when merging models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Finneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guibing Guo, Xingwei Wang, and Dacheng Tao.
 Adamerging: Adaptive model merging for multi-task learning. *arXiv preprint arXiv:2310.02575*, 2023.
- 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*, 2024.
- Xiyu Yu, Tongliang Liu, Xinchao Wang, and Dacheng Tao. On compressing deep models by low
 rank and sparse decomposition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7370–7379, 2017.
- Zhihang Yuan, Chenhao Xue, Yiqi Chen, Qiang Wu, and Guangyu Sun. Ptq4vit: Post-training quantization for vision transformers with twin uniform quantization. In *European conference on computer vision*, pp. 191–207. Springer, 2022.
- Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation.
 In Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition, pp. 11953–11962, 2022.
- 744 745
- 746
- 747 748
- 749
- 750
- 751
- 752
- 753
- 754 755

757		
758	Α	OVERVIEW
759		
700	In t	his appendix we present more information about the experiment settings and analysis that could
701	not	fit in the main paper. In Sec. B we present more details about the datasets, models, and baseline
762	met	hods used in evaluation. In Sec. C we derive the squared Frobenius norm of the interference
763	tern	used in equation 5. In Sec. D we include additional analysis and results.
764		
765 766	В	EXPERIMENT SETUP
767		
768	Thi	s section provides detailed descriptions for the datasets, baselines, and model fine-tuning settings.
770	Dat	asets Details. Evaluation are performed on two sets of datasets with different type of tasks.
771		1. Image Classification Datasets : For image classification, following Ilharco et al. (2022):
772		Tang et al. (2024a); Wang et al. (2024), we use twenty tasks from CLIP's (Radford
773		et al., 2021) test set: SUN397 (Xiao et al., 2010), Cars (Krause et al., 2013), RE-
774		SISC45 (Cheng et al., 2017), EuroSAT (Helber et al., 2019), SVHN (Netzer et al., 2011),
775		GTSRB (Stallkamp et al., 2012), MNIST (Deng, 2012), DTD (Cimpoi et al., 2014),
776		CIFAR100 (Krizhevsky, 2009), STL10 (Coates et al., 2011), Flowers102 (Nilsback &
777		Zisserman, 2008), OxfordIIITPet (Parkhi et al., 2012), PCAM (Veeling et al., 2018),
778		FER2013 (Goodfellow et al., 2013), EMNIST (Cohen et al., 2017), CIFAR10 (Krizhevsky, 2000), Feed101, (Deserved et al., 2014), Feedier MOHET, (Xiag et al., 2017), Ber
779		2009), FOODIOI (Bossard et al., 2014), Fashionivinisi (Alao et al., 2017), Ren- deredSST2 (Seeber et al. 2012; Bedford et al. 2010), and KMNIST (Clanuwet et al.
780		2018) For experiments on K tasks, the first K datasets from this list are select
781		2010). For experiments on K tasks, the first K datasets from this list are select.
782		2. Text Classification and Generation Datasets : For text classification and generation, fol-
783		2018): Col A MNUL MRPC ONUL OOP RTE SST2 and STSR
784		2010). COLA, MINLI, MIXIC, QIVLI, QQI, KIE, 3512 , and $513D$.
785 786	Bas	eline Details. Our experiments compare the following baselines and our methods:
787		• Pre-trained : Pre-trained model used across all tasks (performance lower bound).
788		• Fine-tuned: Individual fine-tuned models (performance upper bound).
789		• Weight Averaging (Wortsman et al., 2022): Merge models by directly averaging their
790		parameters.
791		• Fisher Merging (Matena & Raffel 2022): Fisher Merging uses the Fisher information as
792		a weight for each parameter during weight averaging.
793		• RegMean (Jin et al. 2022): RegMean introduces a constraint in model merging by mini-
794		mizing the L2 distance between the merged model and each individual model
795		• Task Arithmatic (Ilbarco et al. 2022): Task Arithmatic computes the delta peremeters
796		between fine-tuned models and the base model (known as "task vectors") and aggregates
/97		them before adding into a pre-trained model.
798		• Ties-Marging (Vaday et al. 2024): Ties-Marging addresses task conflict issues found in
799		Task Arithmetic by eliminating redundant parameters and resolving symbol conflicts
008		• Leven wise Ado Monoing (Vong et al. 2022): Leven wise Ado Monoing Ende active large
801 000		• Layer-wise Aualyterging (Tang et al., 2023): Layer-wise Aualyterging initias optimal merg- ing coefficients for each layer of each task vector in Task Arithmetic using test time adap
802 000		tation
003		Barometer Supernecition (BSB) (Chaung at al. 2010), DSD applies render arthogonal
004		• 1 at ameter Superposition (FSF) (Cliculing et al., 2019): FSF applies random orthogonal matrices periodically during training to store many models into one model. We adopted
000		it in our offline setting by treating fine-tuned models as different model instances during
000		training to provide some contexts to our task vector superposition approach.
007		• WFMoF (Tang et al. 2024c): WFMoF only margas the layer norm and attention layers
809		while keeping the multi-layer perceptron layers unmerged, with a router to dynamically allocate weights to each MLP conditioned on the input.

810 811 812	• SMILE (Tang et al., 2024b): SMILE compresses task vectors with singular value decomposition (SVD). It then determines the routing weights based on the alignment between input and each low-rank matrix.
813	
814	• TALL-masks+TA (Wang et al., 2024): TALL-masks+TA finds a binary parameter mask
815	for each task vector by finding task-specific parameters with values deviate a lot from the
816	aggregated multi-task vector. The corresponding mask for each task is applied to the multi-
817	task vector before adding to a pre-trained model.
818	
819	• Magnitude Masking (Wang et al., 2024): Magnitude Masking differ from TALL-masks in
820	that it determines per-task masks by keeping the top $k\%$ of each task vector's parameters.
821	
822	• TA+Shuffle (Ours): TA+Shuffle performs random layer shuffling among the repetitive
823	layers in each task vector before merging them with Task Arithmetic.
824	
825	• STA (Ours) : STA applies random orthogonal transformations to each layer in each task
826	vector in Task Arithmetic.
827	
828	• STA+Shuffle (Ours) : STA+Shuffle combines layer shuffling and superposition.
829	
830	
831	Model Details. We utilize fine-tuned models from Tang et al. (2024a) and Wang et al. (2024).
832	Here we describe the experimental setup for fine-tuning these models.
833	
834	• CLIP-ViT-B/32 Models: The CLIP-ViT-B/32 models are fine-tuned by Tang et al. (2024a).
835	The Adam optimizer is employed with a fixed learning rate of $1e^{-5}$ for a total of 4,000
836	training steps with the batch size of 32. The zero-shot classification layer is computed
837	on-the-fly with a frozen text encoder.
838	
839	• CLIP-ViT-L/14 Models: Different from CLIP-ViT-B/32 models, these models are fine-
840	tuned by Wang et al. (2024) with the training procedure described in Ilharco et al. (2022).
841	The AdamW optimizer is employed with a fixed learning rate of $1e^{-5}$ for a total of 2,000
842	training steps with the batch size of 128, and a cosine annealing learning rate schedule
843	with 200 warm-up steps. The zero-shot classification heads are pre-computed and freezed
844	during fine-tuning process, following Ilharco et al. (2022) and Ortiz-Jimenez et al. (2024).
845	
846	• GPT-2 Models : These models are fine-tuned by Tang et al. (2024a) with a constant learning
847	rate of $5e^{-6}$ for 3 epochs.
848	
849	• Flan-T5-base and LoRA Models: These models come from Tang et al. (2024a), with
850	unspecified fine-tuning settings.
851	
852	Evolution Metrica We measure performance using eveness tests economy and normalized economy
853	Evaluation interfects. We measure performance using average task accuracy and normalized ac- curacy (relative to Fine typed baseline). For STSB (Wang 2018), we use Spearman's correlation
854	Memory efficiency is evaluated by estimated memory footprint in Gb and normalized footprint (rel-
855	ative to Pre-trained baseline).
856	
857	
858	Default Experimental Setup. We use global random seeds 42, 43, 44 for three runs per exper-
859	iment on random approaches. Each model's specific seed is generated by adding its index to the
860	global seed, and is used consistently for layer shuffling and binary diagonal matrices across target
861	layers. Following Ilharco et al. (2022), we apply uniform merging coefficients across models, op-
000	timized via grid search on validation sets (10% of training data may 1,000 samples (Wang et a)

timized via grid search on validation sets (10% of training data, max 1,000 samples (Wang et al., 2024)). The search space is 0.1, 0.2, ..., 1.0, extended to 0.1, 0.2, ..., 2.0 for Flan-T5-base LoRA experiments.

C DERIVATION OF EQUATION 5

Here we derive the squared Frobenius norm of the interference $\lambda \sum_{i \neq k} T_i^k$ in more details:

$$\left\|\lambda\sum_{i\neq k}\boldsymbol{T}_{i}^{k}\right\|_{F}^{2} = \left\langle\lambda\sum_{i\neq k}\boldsymbol{T}_{i}^{k},\lambda\sum_{i\neq k}\boldsymbol{T}_{i}^{k}\right\rangle_{F},$$
(13)

$$=\lambda^2 \left(\sum_{i\neq k} \sum_{j\neq k} \langle \mathbf{T}_i^k, \mathbf{T}_j^k \rangle_F \right),\tag{14}$$

$$=\lambda^{2}\left(\sum_{i\neq k}\langle \boldsymbol{T}_{i}^{k},\boldsymbol{T}_{i}^{k}\rangle_{F}+\sum_{i,j\neq k}\langle \boldsymbol{T}_{i},\boldsymbol{T}_{j}^{k}\rangle_{F}\right),\tag{15}$$

$$=\lambda^2 \left(\sum_{i \neq k} \|\boldsymbol{T}_i^k\|_F^2 + 2 \sum_{\substack{1 \leq i < j \leq n \\ i, j \neq k}} \langle \boldsymbol{T}_i^k, \boldsymbol{T}_j^k \rangle_F \right),$$
(16)

$$= \lambda^{2} \left(\sum_{i \neq k} \|\boldsymbol{T}_{i}^{k}\|_{F}^{2} + 2 \sum_{\substack{1 \leq i < j \leq n \\ i, j \neq k}} \|\boldsymbol{T}_{i}^{k}\|_{F} \|\boldsymbol{T}_{j}^{k}\|_{F} \cos(\boldsymbol{T}_{i}^{k}, \boldsymbol{T}_{j}^{k}) \right).$$
(17)

\

D ADDITIONAL ANALYSIS

D.1 GPT-2 TEXT CLASSIFICATION EXPERIMENTS

We evaluated our proposed methods against established baselines by merging seven independently trained GPT-2 models on text classification tasks. As shown in Table 6, both our *STA+Shuffle* algorithm and its *TA+Shuffle* variants achieved significantly higher classification accuracy while doubling the memory footprint, consistent with the performance on other benchmarks.

Table 6: Performance and memory comparison of GPT-2 models across seven GLUE text classification tasks, showing absolute and normalized accuracy (%), as well as memory footprint (Gb). Results averaged over three runs where applicable. Variances smaller than 0.1% are omitted.

$\begin{array}{ c c c c c c c c c c c c c c c c c c c$										
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	Method	Avg.(%) \uparrow	$Bits(Gb) \downarrow$	CoLA	MNLI	MRPC	QNLI	QQP	RTE	SST-2
$ \begin{array}{l c c c c c c c c c c c c c c c c c c c$	Pre-trained	44.5 (54.3)	0.498 (1.00)	30.9	33.0	31.4	49.2	63.2	52.7	50.9
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Weight Averaging	56.1 (63.3)	0.498 (1.00)	55.0	55.1	51.0	57.6	76.7	44.8	52.5
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Fisher Merging	58.7 (64.7)	0.498 (1.00)	54.8	58.0	39.5	63.3	81.5	49.1	64.7
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	RegMean	68.8 (79.7)	0.498 (1.00)	61.7	70.4	65.4	69.7	78.8	56.0	79.7
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Task Arithmetic	70.0 (85.4)	0.498 (1.00)	68.7	68.6	69.6	70.5	81.8	47.3	83.6
PSP 44.5 (54.3) 0.498 (1.00) 30.9 33.6 31.6 49.5 63.2 52.5 50 Fine-tuned 82.0 (100) 3.49 (7.00) 76.8 82.1 80.4 88.3 89.6 65.3 91 TA+Shuffle (Ours) 76.7 (93.5) 0.997 (2.00) 71.6 80.3 73.9 85.8 88.5 47.5 89 STA (Ours) 71.3 ±0.6 (87.0) 0.997 (2.00) 70.3 ±0.1 81.0 61.0 ±1.3 87.2 89.3 57.5 ±0.3 90	Ties-Merging	70.0 (82.4)	0.498 (1.00)	68.4	71.4	68.4	69.6	82.4	47.7	81.8
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	PSP	44.5 (54.3)	0.498(1.00)	30.9	33.6	31.6	49.5	63.2	52.5	50.3
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	Fine-tuned	82.0 (100)	3.49 (7.00)	76.8	82.1	80.4	88.3	89.6	65.3	91.2
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	TA+Shuffle (Ours)	76.7 (93.5)	0.997 (2.00)	71.6	80.3	73.9	85.8	88.5	47.5	89.3
$\textbf{STA+Shuffle (Ours)} \boxed{76.6 \pm 0.2}_{(93.4)} 0.997_{(2.00)} \boxed{70.3 \pm 0.1}_{(2.00)} \textbf{81.0}_{(2.00)} \textbf{81.0}_{(2.00)$	STA (Ours)	71.3 ±0.6 (87.0)	0.997 (2.00)	62.3 ± 0.3	78.2	46.1 ± 4	82.6	88.4	52.7	88.9
	STA+Shuffle (Ours)	$\underline{76.6} \pm \underline{0.2} \ \underline{(93.4)}$	0.997 (2.00)	$\underline{70.3} \pm \underline{0.1}$	81.0	61.0 ± 1.3	87.2	89.3	$57.5 \ \overline{\pm 0.3}$	90.2

D.2 DYNAMICS BETWEEN LAYER SHUFFLING AND SUPERPOSITION

915 We investigate how varying layer shuffling and superposition parameters affects model performance.
916 We test target layer skip rates of 1, 2, 3, 4, where every k-th target layer within repetitive layer sets
917 is shuffled and superposed. We also introduce *layer shifting* – a deterministic alternative to shuffling that shifts layers one position deeper with wrap-around – to study how different decorrelation



Figure 6: Average accuracy and cosine similarity among interfering task vectors when retrieving SUN397 and GTSRB models from 8 merged CLIP-ViT-B/32 models with various target layer skipping rates and shuf-fling/superposition setups.

approaches affect performance.

Experiments on eight CLIP-ViT-B/32 benchmarks (Figure 6) show averaged results across
three repetitions, focusing on overall benchmark accuracy and two specific tasks: SUN397 (Xiao
et al., 2010) and GTSRB (Stallkamp et al., 2012). We analyze both task performance and average
pairwise cosine similarity among task vectors - both original and among interfering vectors during
model retrieval.

956
957As skip rate increases, accuracy declines while cosine similarity rises. Performance remains
stable up to skip rate 2, suggesting potential memory savings through selective layer manipulation.
For GTSRB, TA+Shuffle outperforms TA+Shift despite higher cosine similarity, indicating the
method of achieving orthogonality matters beyond decorrelation levels. This pattern reverses for
SUN397, revealing task-dependent variations and opportunities for task-specific optimization.

The correlation between interfering task vectors' cosine similarity and accuracy shows a negative trend (Figure 7), most pronounced in EuroSAT and MNIST. While patterns vary across tasks, peak accuracy consistently occurs near zero cosine similarity. This observation, combined with STA+Shift's strong performance at low skip rates (Figure 6), suggests a cosine similarity threshold may exist above which method selection becomes less critical.

