WHAT MATTERS FOR MODEL MERGING AT SCALE?

Anonymous authors

Paper under double-blind review

ABSTRACT

Model merging aims to combine multiple expert models into a more capable single model, offering benefits such as reduced storage and serving costs, improved generalization, and support for decentralized model development. Despite its promise, previous studies have primarily focused on merging a few small models. This leaves many unanswered questions about the effect of scaling model size and how it interplays with other key factors—like the base model quality and number of expert models—, to affect the merged model's performance. This work systematically evaluates the utility of model merging at scale, examining the impact of these different factors. We experiment with merging fully fine-tuned models using four popular merging methods—Averaging, Task Arithmetic, Dare-TIES, and TIES-Merging—across model sizes ranging from 1B to 64B parameters and merging up to 8 different expert models. We evaluate the merged models on both held-in tasks, i.e., the expert's training tasks, and zero-shot generalization to unseen held-out tasks. Our wide range of experiments provide several new insights about model merging at scale and the interplay between different factors. *First*, we find that merging is more effective when experts are created from strong base models, i.e., models with good zero-shot performance, compared to pre-trained ones. *Second*, larger models facilitate easier merging. *Third* merging consistently improves generalization capabilities. Notably, when merging eight large expert models, the merged models often generalize better compared to the multitask trained models. *Fourth*, we can better merge more expert models when working with larger models. Fifth, different merging methods behave very similarly at larger scales. Overall, our findings shed light on some interesting properties of model merging while also highlighting some limitations. We hope that this study will serve as a reference point on large-scale merging for upcoming research.

032 033

000

001 002 003

004

005 006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

023

025

026

028

029

031

034

1 INTRODUCTION

Model merging (Raffel, 2021) refers to the process of combining two or more *constituent* (*expert*) models to produce a new, and potentially more powerful model. The appeal of this technique is rooted in several benefits it can confer: *first*, it dramatically reduces storage and serving costs by reusing a single model across tasks; *second*, it enables compositional combination of capabilities from expert models, which can improve generalization to novel tasks; and *third*, merging supports decentralized and modular model development by allowing multiple contributors to independently build models and later combine them together.

044 These characteristics have led to a great deal of recent efforts in developing cost-effective model merging methods (Matena & Raffel, 2022b; Ilharco et al., 2022; Jin et al., 2022; Yadav et al., 2024b; 046 Yang et al., 2023; Yu et al., 2024a; Shah et al., 2023; Tam et al., 2023; Zhao et al., 2024), often using 047 simple *arithmetic operations*, such as averaging the parameters of the constituent models. However, 048 most of these studies are limited to small-scale experiments with relatively small models (typically < 7B parameters) and merging 2 or 3 experts (Yu et al., 2024a;b), and mainly focus on improving benchmark performance on *held-in* tasks that the expert models were trained on (Yu et al., 2024a; Yadav et al., 2024b). Despite the promises that model merging holds, the research community still lacks a comprehensive study to evaluate its effectiveness as we scale the model size. Moreover, it is 052 not clear how scale interplays with other factors like number of expert models and base model quality to affect the merged model's held-in performance and zero-shot generalization. This is of paramount



066 Figure 1: Held-In performance results from our large scale model merging experiments con-067 ducted over keys factors like base models, model sizes, merging methods, and number of ex-068 perts being merged. We present results for two base models, PaLM-2 and an instruction tuned 069 version of it, PaLM-2-IT, four different models sizes (1B, 8B, 24B, 64B), four merging methods (Averaging, Task Arithmetic, Dare-TIES, and TIES-Merging), when merging either 2 or 8 ex-071 pert models. We report the performance normalized with the oracle expert's performance which is denoted by the bold black circle of radius 1. We also present the performance of multitask base-073 line train on the held-in tasks. We find merging expert models created from the instruction tuned 074 PaLM-2-IT model always performs better than merging PaLM-2 based experts. Moreover, the gap between these model increase when we merge more experts. Larger experts (64B) merge better and 075 show the best held-in performance. 076

077

090

091

092

093

094

importance, as models are rapidly growing in size, and more open-weight models and datasets are becoming available,¹ driving the need for practical and scalable merging methods.

Our primary goal in this paper is to provide insights into the scalability of model merging. While a few studies have explored merging at the 13B parameter scale (Huang et al., 2024a; Yu et al., 2024a;b), they primarily leverage increased model size and combine only 2-3 models to attain better performance on held-in tasks. As such, the interplay of factors like model size, base model quality, number of constituent models—and their effect on both held-in and zero-shot generalization performance (*held-out*)—remains largely *unexplored*. Hence, we aim to address the following *four* research questions (RQ):

- **RQ1:** What is the effect of using *pretrained* vs. *instruction-tuned* base models for creating expert models for merging?
- **RQ2:** Does model merging become *easier* or *harder* as the model size increases?
 - **RQ3:** How does merging affect *zero-shot generalization* to held-out tasks, and how is this influenced by model size?
 - **RQ4:** How *many* expert models can be merged without performance loss, and how does this depend on model size?

To answer these question, we systematically evaluate the effectiveness of current state-of-the-art 096 merging methods through empirical experiments. Specifically, we utilize the PaLM-2 model (Anil 097 et al., 2023) and its instruction-tuned variant, PaLM-2-IT, while scaling the model sizes up to 64B 098 parameters. We experiment with *four* popular merging methods, namely, Averaging (Wortsman et al., 099 2022a; Choshen et al., 2022b), Task Arithmetic (Ilharco et al., 2022), TIES-Merging (Yadav et al., 100 2024b), and Dare-TIES (Yu et al., 2024a). We conduct a series of sensitivity and ablation experiments 101 to understand the relative importance of several factors like model size (1B, 8B, 24B, 64B parameters), 102 base model quality (pretrained vs. instruction-tuned), and number of constituent models (2, 4, 6, 8) 103 being merged. Additionally, we consider two axes of evaluation using the T0 data collection (Sanh et al., 2021a): held-in evaluation with tasks the expert models were trained on, and held-out, for 104 zero-shot generalization to unseen tasks. 105

106 107

¹As of writing, the largest open-weight AI model is Llama 3.1 405B parameters, and Hugging Face hosts a plethora of community-contributed resources, with 1M+ models and 200K+ datasets.



Figure 2: Merged experts created from big and strong base models generalize better than multitask models. We find that for strong base models as we merge more experts (x-axis, \rightarrow), the merged model's generalization performance (y-axis, \uparrow) monotonically increases to approach and eventually surpasses multitask baseline. (yellow line). More details in Section 4.3.

122 Our experiment results shed light on the promises of model merging and reveal interesting insights 123 into the behaviors of different factors at scale. First, we find that the model initialization plays a 124 crucial role in enhancing the performance of the merged model. Specifically, across all evaluation 125 settings, using strong zero-shot instruction-tuned base models to create expert models leads to 126 improved performance compared to using pretrained models (see §4.1). Second, larger models are 127 consistently easier to merge. This holds true regardless of the base model used (instruction-tuned or not), number of models merged, or merging method (see §4.2). *Third*, our results demonstrate that 128 merging significantly enhances zero-shot generalization, consistently improving the ability to adapt 129 to new tasks. *Notably*, when using strong base models as the number of merged experts increases, 130 our merged model either matches or exceeds the performance of a strong multi-task training baseline 131 (see $\S4.3$). Fourth, larger models are better at merging a larger number of expert models (see $\S4.4$). 132 Finally, our numerous experiments identify specific settings where we expect model merging to be 133 much more useful. From this we provide general recommendations for practitioners (see $\S5$). Taken 134 as a whole, our findings are a powerful testament to the potential of model merging at scale for 135 creating highly generalizable language models, which we hope will spur more fundamental research 136 into the development of practical and scalable merging methods.

137 138 139

140

117

118

119

120 121

2 BACKGROUND

Model merging (Yang et al., 2024; Goddard et al., 2024) has emerged as a cost-effective method for developing improved models. Two common use cases of merging are: (1) combining model checkpoints from different data versions, hyperparameters, or training stages to enhance distributional robustness (Team et al., 2024; Dubey et al., 2024), and (2) combining multiple expert models trained on different datasets to leverage their complementary capabilities. In both scenarios, the expert models generally share a common architecture and a base model from which the expert models are created via fine-tuning.

This work focuses on merging specialized, fine-tuned versions (experts) of a single base model to enhance its capabilities. Each expert model is trained on distinct datasets covering different tasks, domains, and/or capabilities. We refer to the tasks/datasets used for training the expert models as "held-in", while those that are new and unseen are called "held-out". Our goal is to create a unified model that retains the individual expert models' capabilities on held-in tasks while improving zeroshot generalization on held-out tasks. This merging approach provides a flexible, modular method for post-training large language models, facilitating the addition of new features and capabilities to top-performing models.

155 156

157

2.1 MODEL MERGING METHODS

We denote the set of N expert tasks as t_1, \ldots, t_N and the base model weights, representing the common ancestor of all expert models as θ_{base} . The weights of the corresponding specialized expert models, each obtained by fully fine-tuning the base model on a specific expert task, are denoted as $\theta_1, \ldots, \theta_N$, respectively. We focus on "open vocabulary" models which utilize natural language as input and output for both classification and generation tasks, eliminating the need for task-specific classification heads making the merging process simpler. Given this, model merging methods can be defined as a function $\mathcal{M}(.)$. This function takes as input the base model, the set of N expert models, and potentially additional information, denoted by Φ . This additional information may include activation statistics, Fisher matrices, or other method-specific data. The output of the function is the merged model, represented by its parameters θ_m . Formally, $\theta_m = \mathcal{M}(\{\theta_i\}_{i=1}^{N}, \theta_{base}, \Phi)$, where Φ is method specific data.

168 Given our focus on studying model merging with large models, we select four merging methods 169 based on their popularity and simplicity. We only study merging methods that can scale to tens 170 of billions of model weight parameters and do not require any additional information to perform 171 merging, i.e., $\Phi = \{\}$, as these techniques are efficient for even larger models. Other more complex 172 methods that require computing fisher matrices (Matena & Raffel, 2022a), backward passes (Yang et al., 2023), or additional information like model activation (Jin et al., 2023) are skipped because of 173 their computational complexities for large scale model merging that we focus on in this work. Next, 174 we describe the four selected model merging methods in detail. 175

176 177

186

187

2.1.1 AVERAGING

Parameter averaging (Choshen et al., 2022b; Wortsman et al., 2022a) is a well-established technique in federated learning (McMahan et al., 2017) and recent applications extend its utility to merge models for enhancing model robustness against out-of-distribution data (Wortsman et al., 2022b; Ramé et al., 2022), refine pre-trained models (Yu et al., 2024a), develop multimodal models (Sung et al., 2023), and create multitask models by combining capabilities (Yadav et al., 2024b; Ilharco et al., 2022). Parameter averaging is achieved by taking a mean of all the expert model weights together without using the base model which can be formally described as, $\mathcal{M}(\{\theta_i\}_{i=1}^{\mathbb{N}}, \theta_{base}) = \frac{1}{\mathbb{N}} \sum_{i=1}^{\mathbb{N}} \theta_i$.

2.1.2 TASK ARITHMETIC

Task Arithmetic (Ilharco et al., 2022) introduces a novel concept of "*task vectors*" for model merging. For task t_i , the task vector is denoted as $\tau_i = \theta_i - \theta_{base}$ which captures task-specific knowledge by quantifying the difference between the fine-tuned expert parameters (θ_i) and the original base model parameters (θ_{base}). A scaling hyperparameter λ controls the contribution of the aggregated taskspecific knowledge to the final model. The merged model is then constructed by linearly combining the base model parameters with a scaled sum of all task vectors. Formally, task arithmetic can be described as, $\mathcal{M}(\{\theta_i\}_{i=1}^{\mathbb{N}}, \theta_{base}; \lambda) = \theta_{base} + \lambda * \sum_{i=1}^{\mathbb{N}} (\theta_i - \theta_{base})$.

2.1.3 TIES MERGING

197 TIES-Merging (Yadav et al., 2024b) identifies two main challenges with model merging: **1** during finetuning expert models accumulate a lot of noise in the parameters, and *Q* different experts might want to change the same parameter in different directions leading to interference/conflict between 199 the expert models. They demonstrate that both of these factors hurt model merging and propose a 200 three steps process to remove redundant parameters, followed by resolving sign conflicts, and finally 201 aggregating only the parameters that are not conflicting. Specifically, in TIES Merging they first 202 zero out the values in each task vector that have low magnitudes to obtain the trimmed task vector 203 $\hat{\tau}_i$ for each task. Next, they chose the aggregate sign (γ_m) for each parameter based on whether the 204 parameter has a higher total magnitude in the positive or the negative direction across all trimmed 205 task vector, formally, $\gamma_{m} = \operatorname{sgn}(\sum_{i=1}^{N} \hat{\tau}_{i})$. Finally, for each parameters p the models whose sign 206 matches the aggregate sign are averaged to obtain the merged task vector. Finally, the merged model 207 is obtained by scaling the merged task vector using a hyperparameter λ and then added back to the base model as, $\theta_{m}^{p} = \theta_{base} + \lambda * \frac{1}{|\mathcal{A}^{p}|} \sum_{i \in \mathcal{A}^{p}} \hat{\tau}_{i}^{p}$, where $\mathcal{A}^{p} = \{i \in [N] \mid \hat{\gamma}_{i}^{p} = \gamma_{m}^{p}\}$. 208

209 210 211

2.1.4 DARE MERGING

212 Dare (Yu et al., 2024a) extends the idea of TIES merging by proposing to use a dropout-like pruning 213 stage to remove noise before merging. Specifically, a Bernoulli mask M_i with drop probability p 214 is applied to each task vector to obtain the pruned task vector $\hat{\tau}_i = (1 - M_i) \odot \tau_i / (1 - p)$. This 215 stochastic process randomly zeroes out elements within the task vector while preserving its expected value. These pruned task vectors are then used along with either TIES Merging or Task Arithmetic. Due to the popularity of the Dare variant that uses TIES Merging, we use that to represent the Dare method and call it *Dare-TIES*.

218 219

220

2.2 CHALLENGES/LIMITATIONS

221 Model Merging has been utilized at a growing rate in practice as it has recently been applied to building modern language models like Llama-3 (Dubey et al., 2024) and Gemma-2 (Team et al., 2024). 222 Recently, Lu et al. (2024) showed that merged models can lead to the emergence of capabilities that 223 surpass the individual contributions of the parent models. However, most formal studies on model 224 merging have been performed with relatively small models. There are a few studies that look at larger 225 models with 7B and 13B parameters. However, those studies mostly focus on merging 2-3 models to 226 improve benchmark numbers as opposed to better understanding how the size of the model affects 227 the model merging process and the resultant model. To motivate our work, we present some of the 228 limitations of the existing studies and highlight their difference with our work. 229

230 **Most Studies on Small Models** (< 7B **parameters**): Almost all existing model merging papers 231 rarely use large models (> 7B). For example past works (He et al., 2024; Daheim et al., 2023; 232 Ortiz-Jimenez et al., 2024; Jang et al., 2024), including popular methods like ModelSoup (Wortsman 233 et al., 2022a), Task Arithmetic (Ilharco et al., 2023) and TIES-Merging (Yadav et al., 2024b), 234 RegMean (Jin et al., 2023), Fisher-Merging (Matena & Raffel, 2022a) Ada-Merging (Yang et al., 2023), MatS (Tam et al., 2024) perform experiments with model families like CLIP (Radford 235 et al., 2021), ViT (Dosovitskiy et al., 2021), T5 (Raffel et al., 2020a), DeBERTa (He et al., 2021), 236 Roberta (Liu et al., 2019), BERT (Devlin et al., 2018) with less than 1B parameters. Hence, it is 237 unclear how well model merging works for large models, what factors play an important role, the 238 effect of model size, number of tasks being merged, and its effect on both held-in performance and 239 generalization of the model. Some studies hypothesize that bigger models might be easier to merge 240 however there are no concrete large scale studies to thoroughly assess such claims at large scale. 241

242 Model Merging Studies with Large Models are Shallow: Some recent works like DARE (Yu 243 et al., 2024a), WIDEN (Yu et al., 2024b), Chat-Vector (Huang et al., 2024b) demonstrate merging 244 results for larger models with up to 13B parameters, however these studies have a few limitations: **0** 245 They primarily focus on using model merging to improve model quality and hence their experiments 246 do not provide concrete insights on how model size interplays with merging, ^(a) They only merge a maximum of two or three models at once, ⁽³⁾ They primarily focus on held-in tasks and do not 247 provide any insights on the effect of merging on a model's generalization abilities. Other works like 248 RewardSoup (Rame et al., 2024), WARM (Rame et al.), WARP (Ramé et al., 2024), FuseLLM (Wan 249 et al., 2024a), FuseChat (Wan et al., 2024b) also work with \sim 7B sized models and focus on specific 250 applications of model merging without providing any deeper insight about how merging performance 251 changes for large models.

252 253

Varied Evaluation Setups: Most previous works rarely share their experimental setup where 254 both the expert datasets and the objective vary. For example, RegMean (Jin et al., 2023), Task 255 Arithmetic (Ilharco et al., 2023), TIES (Yadav et al., 2024b), MaTS (Tam et al., 2024) uses GLUE 256 tasks (Wang et al., 2018), Vision tasks, T0 held-out, and T0 held-in (Sanh et al., 2021b) tasks 257 respectively. Moreover, different works evaluate for different use cases like intermediate task training 258 in Fisher merging (Matena & Raffel, 2022a), robustness in modelsoups (Wortsman et al., 2022a), and held-in performance for Dare (Yu et al., 2024a), both held-in and held-out performance in TIES 259 Merging (Yadav et al., 2024b). Tang et al. (2024) attempts to unify these various problem setting 260 and different evaluation setup and compare many method on similar setups. Given our focus on 261 combining model capabilities in the post training phase, we focus on evaluating on both held-in tasks 262 and generalization to unseen held-out tasks. 263

264 265

266

3 LARGE SCALE EVALUATION OF MODEL MERGING

In this work, we address the limitations mentioned above by systematically understanding the effect
 of various factors like model size, base model quality, merging method, and the number of models
 being merged on both the held-in and generalization performance of the final merged model. Next,
 we describe our experimental design.

270 **Data:** Sanh et al. (2021a) found that explicit multitask training of T5 (Raffel et al., 2020b) on a 271 collection of prompted datasets produces a model with strong zero-shot performance on unseen tasks. 272 This has become a common experimental setting for benchmarking zero-shot generalization (e.g. 273 (Longpre et al., 2023; Jang et al., 2023; Zhou et al., 2022; Chung et al., 2024; Muqeeth et al., 2024). 274 Hence, we adopt the experimental setting from the T0 mixture (Sanh et al., 2021a) which contains 8 held-in and 4 held-out task categories. For each of these categories there are multiple datasets 275 in the T0 mixture (Sanh et al., 2021b) and hence to reduce evaluation costs, we select 2 datasets 276 from each category based on the popularity and the train dataset size. Specifically, the 8 held-in task 277 categories (with a total of 16 datasets) include Multiple-choice QA, Extractive Qa, Closed-Book 278 QA, Sentiment Analysis, Topic Classification, Structure-to-text, Summarization, and Paraphrase 279 Identification. Similary, the 4 held-out task categories (with a total of 7 datasets) are Sentence 280 Completion, Natural Language Inference, Coreference Resolution, and Word Sense Disambiguation. 281 For more details see Section B. 282

283

Expert Model Creation: Recognizing the significance of post-training for LLMs where models 284 are typically fully fine-tuned, we perform full fine-tuning to create our expert models to better mimic 285 the post-training setting. Moreover, in post-training phases it is common to first perform Instruction 286 Tuning (IT) on the model before moving on to other steps. Hence, we examine the effect of using 287 strong instruction-tuned base models on the process and outcome of model merging. Given this, we 288 utilize the PaLM-2 models (Anil et al., 2023) with sizes 1B, 8B, 24B, and 64B as our base models 289 (θ_{base}) . To obtain the instruction tuned base model, we further fine-tuned the PaLM-2 models on the 290 FLAN-v2 dataset (Longpre et al., 2023) while excluding the T0-mixture tasks (Sanh et al., 2021a). 291 These instruction-tuned variants are denoted as PaLM-2-IT. For each of the 2 base model types (non-IT vs IT) and 4 model sizes, we perform full fine-tuning on the 8 held-in task categories resulting 292 64 specialized experts models which are then used further in our experiments. Comprehensive details 293 regarding hyper parameters and computational requirements are provided in Appendix C. 294

295

296 **Experimental Setting:** Given our collection of expert models, for each merging experiment we 297 select a subset of expert models which we call the *constituent models*. We create a large merging 298 experiment grid with 2 base models (PaLM-2 and PaLM-2-IT), four model sizes (1B, 8B, 24B, 64B), 299 four Merging methods (Averaging, Task Arithmetic, Dare-TIES, and TIES), the number of constituent models (2, 4, 6, 8), and 3 seeds to randomly select the constituent tasks for the experiment resulting 300 in a total of 384 merging experiments. These seeds are shared across different experimental settings 301 to ensure the same tasks are selected across base models, model sizes and merging methods to ensure 302 fair comparison. For example, in an experiment we merged 2 expert models, derived from the 64B 303 PaLM-2 base model with the constituent models being MCQ and Summarization experts while the 304 same experiment with a different seed resulted in Closed Book QA and Sentiment Analysis experts 305 as the constituent models.

306 307

Evaluation: For each of the experiments above, we assess the merged model's performance by evaluating it on both the held-in tasks – i.e., the training tasks of the constituent expert models – and all 4 held-out task categories. For example, if the constituent models are MCQ and Summarization experts, then for held-in tasks we evaluate on the MCQ datasets (DREAM and Cosmos QA) and Summarization datasets (CNN Daily Mail and XSum) resulting a total of 4 held-in evaluation datasets. Moreover, all merging experiments are also evaluated on the 4 held-out tasks categories consisting of 7 datasets listed in Appendix B. There we perform approximately ~ 9000 model evaluations across all of our experiments.

315

316 Metric: Given that different datasets use different metrics, we normalize the performance metrics to 317 make them unitless so that they can be aggregated. For held-in tasks, the merged model's performance 318 is normalized against the corresponding task expert model's performance. However, for held-out 319 tasks, the normalization was performed relative to the base model's performance. We denote this 320 metric as *normalized performance* throughout the paper. Importantly, we want to emphasise that 321 this metric is relative, with a value of 1 indicating performance comparable to the reference model. Hence, for held-in tasks a value of 1 means performance similar to the domain expert model while 322 for held-out tasks it means performance is similar to the base model. We mark this line in most of 323 our figures and specify the models that are used for normalization. Finally, to generate aggregated



Figure 3: Instruction-tuned models facilitate easier merging. $PaLM-2-IT(\bullet)$ consistently outperforms $PaLM-2(\bullet)$ as shown by the huge gap between the green point (•) being higher than red points (•), across various merging methods, model sizes, and numbers of constituent models, indicating that stronger instruction-tuned base models enhance the performance of merged models. The dashed lines denoted the performance of the experts trained on the held-in tasks as defined in § 3. For more details see Section 4.1.

results, we compute the mean of normalized performance across all datasets within each category, then across all categories and then over the three seeds.

341 342 343

344

332

333

334

335

336

337

338 339

340

4 EXPERIMENTAL RESULTS

345 In this section, we explore the interplay between model size and key factors such as base model quality, 346 merging method, and the number of constituent (expert) model, along with their effect on both held-in 347 and zero-shot generalization (held-out) performance. Our findings are: 1 Merging is more effective 348 when the constituent models are derived from instruction-tuned base models rather than pretrained 349 ones (see $\S4.1$); **\Theta** Larger models facilitate easier merging ($\S4.2$); **\Theta** Merging significantly improves 350 zero-shot generalization, with instruction-tuned models benefiting from increased constituent models, and larger model sizes allowing the merged model to match or exceed multi-task training ($\{4.3\}$; 351 We can merge more models effectively when using larger models ($\S4.4$); and **\Theta** Different merging 352 methods perform similarly when applied to large-scale instruction-tuned models. Below, we outline 353 the experimental setup and discuss these findings in detail. 354

- 355
- 356

357

4.1 INSTRUCTION-TUNED MODELS FACILITATE EASIER MERGING

Experimental Setup: Prior research suggests a connection between robust zero-shot models and effective model merging. Wortsman et al. (2022a) demonstrate that averaging strong zero-shot models improves out-of-distribution robustness. Ortiz-Jimenez et al. (2024) indicate that effective pretraining allows for weight disentanglement, and thus enhancing merging. Other studies (Yadav et al., 2024b; Ilharco et al., 2023) propose that strong base models could aid in model merging, though this hypothesis remains largely untested.

To assess how base model quality affects the held-in performance of merged models, we perform merging experiments with fully fine-tuned experts from PaLM-2 and PaLM-2-IT. We vary model sizes in {1B, 8B, 24B, 64B} and the number of constituent models in {2, 4, 6, 8}. Held-in performance is measured over three trials to minimize the impact of selected expert models and their data distributions. A consistent seed is used across different base models, model sizes, and merging methods to ensure fair task comparisons. We evaluate four merging methods: averaging, task arithmetic, TIES, and Dare-TIES, and also compare against the performance of task-specific expert models.

371

Findings: Our results, presented in Figure 3, indicate that PaLM-2-IT models denoted by green color (•), consistently outperforms PaLM-2 models (•) across various merging methods (•, \bigstar , \diamondsuit , \star), model sizes (x-axis \rightarrow), and numbers of constituent models (subplots). This supports our hypothesis that stronger instruction-tuned base models enhance the performance of merged models. Similar to the findings of Ortiz-Jimenez et al. (2024), we believe that large-scale instruction tuning further disentangles model weights, facilitating effective model merging and improving the base model's zero-shot performance.



Figure 4: **Bigger models merge better.** On Held-In evaluations, we find that bigger models always perform better compared to smaller models, barring a few outliers. We find that large instruction tuned models like 64B PaLM-2-IT are the easiest to merge. For more details see Section 4.2.



Figure 5: Merged models at scale generalize better. We plot the held-out generalization of the merged model for two merging methods. We also include the performance of base model (dashed line) and the multitask baseline (yellow line) which trained on a mixture of held-in tasks. We find that the number of constituent expert models (x-axis, \rightarrow) had little effect on zero-shot generalization as shown in the left and center plots. However, increasing model size significantly to 64B improved the merged model's performance over the base model (right plot). For more details see Section 4.3.

4.2 MODEL MERGING BECOMES EASIER WITH BIGGER MODELS

Experimental Setup: In this section, we explore the effect of model size on the held-in performance
 of merged models. We run experiments using different model sizes, base models, merging methods,
 and numbers of constituent models. As in the previous experiment, we report the average results over
 three random seeds and compare the performance of the merged models to that of the task-specific
 expert models.

Findings: Figure 4 illustrates how increasing base model size impacts merging effectiveness. As model size grows (denoted by colors, \blacksquare , \blacksquare , \blacksquare , \blacksquare), merged model performance generally improves. This positive trend is consistent across all base models (different subplots), merging methods (x-axis \rightarrow), and numbers of constituent models (subplots). For large instruction-tuned PaLM-2-IT models, the merged models perform nearly as well as task-specific expert models denoted by dashed line. These results demonstrate that larger models facilitate merging. This suggests a promising approach for developing adaptive, modular post-training recipes. If the remaining performance gap can be further reduced, model merging could become a cost-effective alternative to multitask training. Our full results across settings are available in the Appendix D.

4.3 MERGED MODELS AT SCALE GENERALIZE BETTER

Experimental Setup: Expert models are created by fine-tuning our base model on specialized tasks, which can lead to a decrease in its generalization capabilities. This raises the question: *How well, if at all, can the merged model generalize to held-out tasks?* Ideally, the merged model should perform at least as well as the base model on these tasks. To explore this, we evaluate the merged



Figure 6: **Bigger model sizes can merge more experts.** We merge experts of various sizes created from PaLM-2 and PaLM-2-IT models and plot the held-in (left) and held-out (right) performance of merged models. While PaLM-2's held-in performance degrades with more experts, PaLM-2-IT's performance stabilizes at a much higher level. Both PaLM-2 and PaLM-2-IT models consistently improve held-out generalization, particularly at 24B and 64B scales with increasing expert count. For more details see Section 4.4.

model's performance on unseen tasks across various model sizes, merging methods, and numbers of constituent models. Additionally, we compare our merging approach to a traditional multitask baseline, where a single model is trained on a mixture of all eight held-in task categories. As detailed in Section 3, we normalize the performance of both the merged and multitask model against the base model to assess relative gains or losses in generalization abilities.

452 **Findings:** Figure 2 and Figure 5 show the zero-shot generalization performance of the merged 453 model using PaLM-2-IT and PaLM-2, respectively. Overall, we find that: **1** The merged models 454 outperform their corresponding base models in zero-shot generalization to held-out tasks, as indicated 455 by performance values greater than 1 in most cases; ⁽²⁾ This improvement is consistent across 456 various model sizes (denoted by subplot), base models (different figures), merging methods (different 457 colors \blacksquare , \blacksquare), and numbers of constituent models (on x-axis \rightarrow), suggesting that merging generally improves generalization; ⁽³⁾ For weak base models (i.e., PaLM-2) illustrated in Figure 5, the number 458 of constituent expert models had little effect on zero-shot generalization (Left and Center plots). 459 However, increasing model size significantly improved the merged model's performance over the base 460 model (Right plot); ④ In contrast, strong base models (PaLM-2-IT) show a different trend, zero-shot 461 generalization monotonically improves with the addition of more expert models as shown in Figure 2. 462 We hypothesize this positive correlation arises from reduced model noise through the inclusion of 463 multiple experts, resulting in better generalization; and **6** Notably, our merged model outperforms 464 the multitask baseline when combining more than 6 large instruction-tuned expert models (over 465 24B). This indicates that models developed through merging can generalize even better than those 466 trained on a multitask mixture, offering a promising approach for developing highly capable language 467 models. Our full results on other merging methods and model size are available in Appendix D.

468 469

470

440

441

442

443

444

445 446

447

448

449

450

451

4.4 BIGGER MODEL SIZES CAN MERGE MORE EXPERTS

Experimental Setup: When creating multitask models, datasets for different tasks or domains are
typically combined. In contrast, model merging involves developing separate expert models for each
task or domain before combining them. Previous work has shown that merging multiple models can
reduce the quality of the resulting model (Yadav et al., 2024b; Ilharco et al., 2022). In this study, we
experiment with merging up to 8 expert models from various base models, model sizes, and merging
methods to assess their impact on successful merges.

477

Findings: Figure 6 shows the held-in and held-out performance of the merged models using 478 Task Arithmetic as the number of constituent models increases shown on x-axis. Results for other 479 methods can be found in Appendix D. Overall, we observe that: **1** Unlike merging with PaLM-2, 480 where held-in performance typically declines with more model merges, merging with stronger zero-481 shot PaLM-2-IT initially drops slightly in performance before stabilizing as number of constituent 482 models increase. For example, merging eight 8B PaLM-2 models decreases performance from 0.66 483 to 0.39 when increasing the number of experts from 2 to 8, whereas PaLM-2-IT's performance only slightly drops from 0.91 to 0.86; @ In the held-out evaluations, the merged experts based on PaLM-2 484 models generally outperform the base PaLM-2 models by a small margin. However, with larger model 485 sizes (64B), the performance improvement increases significantly, achieving about 30 percentage

relative improvement. We attribute this substantial gain to the base PaLM-2 model's weak zero-shot
 performance; and ⁽³⁾ The merged models based on PaLM-2-IT show improved generalization over
 PaLM-2-IT across all settings. Additionally, for the 24B and 64B models, we observe a consistent
 increase in generalization capabilities with the addition of more constituent expert models.

490 491

492

4.5 MERGING METHODS BECOME SIMILAR AT SCALE

493 We find that all merging methods exhibit similar 494 performance when merging large instruction-495 tuned models. This suggests that simpler meth-496 ods, such as Averaging, can be sufficient for merging powerful large expert models. Figure 7 497 shows the held-in and held-out performance of 498 the 64B experts derived from PaLM-2-IT. All 499 merging methods yield comparable results on 500 both held-in and held-out tasks for any number 501 of constituent models (shown on x-axis). We 502 hypothesize that as model size increases, expert 503 models are highly over-parameterized due to 504 limited training data. Consequently, the subtle 505 advantages of certain merging techniques - such 506 as highlighting information via task vectors (IIharco et al., 2022), resolving interference (Ya-507 dav et al., 2024b), or pruning (Yu et al., 2024a) 508 - which benefit smaller models, become less rel-509 evant. This indicates a need for more practical 510 and scalable methods to improve merging at scale.



Figure 7: **Different merging methods become similar at scale.** We plot the held-in and held-out performances of merged 64B PaLM-2-IT models across different merging methods and numbers of constituent models. For more details see Section 4.5.

- 511 512
- 513
- 514

5 CONCLUSION, FUTURE WORK AND TAKEAWAYS

515 We summarize key insights from our study and provide practical recommendations for model 516 merging practitioners. Overall, we find that: **0** Creating expert models from the best available 517 base model is always beneficial. The quality of the base model can be gauged by its zero-shot 518 generalization capabilities. We hypothesize that better generalization leads to improved weight 519 disentanglement (Ortiz-Jimenez et al., 2024) and a flatter loss landscape, enhancing linear mode 520 connectivity and facilitating model merging; ⁽²⁾ Merged models often underperform compared to task-specific expert models, indicating a potential loss in performance. Despite this, specialized 521 expert models generally outperform general-purpose multitask models (Liu et al., 2022; Roziere et al., 522 2023; Luo et al., 2023), suggesting that the performance loss may not be significant when compared 523 to multitask models trained on specific tasks; and **③** Our findings indicate that large-scale merging 524 can accommodate more models and significantly improve generalization, outperforming multitask 525 training when a powerful zero-shot base model is employed. • Surprisingly, we find that when 526 working with large instruction tuned models, different merging method perform very similary. This 527 implies that using simple merging methods like averaging will result in models that are comparable 528 in quality with the models obtained from more advanced merging method. 529

Based on our findings, the following future directions are worth exploring. First, exploring strategies 530 to mitigate held-in performance loss during merging. Some sort of iterative branched training and 531 merging of expert models could be fruitful to mitigate performance held-in performance loss. Second, 532 investigating weight disentanglement and its relationship with zero-shot capabilities could yield 533 deeper insights for improving both held-in and held-out performance. Third, diving deeper into the 534 insight as to why different merging methods perform similarly at scale. It might be valuable to assess the generality of this finding either theoretically or empirically across different models. Lastly, there 536 are some minor aberrations in our findings which we believe are impacted by the choice of the expert 537 data distributions. A thorough study analyzing the impact of data distribution on the merging process would be very useful. Answering questions like which models are easier to merge and which ones are 538 harder can result in insights to build better merging methods. We hope our research inspires further fundamental studies on developing more practical and scalable merging methods.

540 REFERENCES 541

341	
542	David Ifeoluwa Adelani, Jade Abbott, Graham Neubig, Daniel D'souza, Julia Kreutzer, Constan-
543	tine Lignos, Chester Palen-Michel, Happy Buzaaba, Shruti Rijhwani, Sebastian Ruder, et al.
544	Masakhaner: Named entity recognition for african languages. Transactions of the Association for
545	Computational Linguistics, 9:1116–1131, 2021.
546	Samuel K Ainsworth Ionathan Hayase and Siddhartha Sriniyasa, Git re-basin: Merging models
547	modulo permutation symmetries, 2022. https://arxiv.org/abs/2209.04836.
548	
549	Takuya Akiba, Makoto Shing, Yujin Tang, Qi Sun, and David Ha. Evolutionary optimization of model merging regipes arXiv preprint arXiv:2403.13187.2024
550	model merging recipes. <i>urxiv preprut urxiv.2405.15107</i> , 2024.
551	Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,
552	Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. arXiv
553	preprint arXiv:2305.10403, 2023.
555	Gregory Benton, Wesley Maddox, Sanae Lotfi, and Andrew Gordon Gordon Wilson. Loss surface
556	simplexes for mode connecting volumes and fast ensembling. In International Conference on
557	Machine Learning (ICML), 2021. https://arxiv.org/abs/2102.13042.
558	Leshem Choshen, Elad Venezian, Shachar Don-Yehia, Noam Slonim, and Yoay Katz. Where to
559	start? analyzing the potential value of intermediate models. 2022a. https://arxiv.org/
560	abs/2211.00107.
561	
562	Leshem Choshen, Elad Venezian, Noam Slonim, and Yoav Katz. Fusing finetuned models for better
563	pretraining. <i>urxiv preprint urxiv.2204.03044</i> , 20220.
564	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,
565	Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language
566	models. <i>Journal of Machine Learning Research</i> , 25(70):1–53, 2024.
567	Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment
500	challenge. In Machine Learning Challenges Workshop, 2005. https://link.springer.
570	com/chapter/10.1007/11736790_9.
571	Nice Dahaim Thomas Möllenhoff Edearde Maria Ponti, Irvna Gurayych, and Mohammad Emtiyaz
572	Khan Model merging by uncertainty-based gradient matching arXiv preprint arXiv:2310.12808
573	2023.
574	
575	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Ioutanova. BERI: Pre-training of deep hidiractional transformers for language understanding arXiv preprint arXiv:1810.04805, 2018
576	bidirectional transformers for language understanding. <i>arXiv preprint arXiv</i> .1810.04805, 2018.
577	William B. Dolan and Chris Brockett. Automatically constructing a corpus of sentential para-
578	phrases. In International Workshop on Paraphrasing, 2005. https://aclanthology.org/
579	105-5002.
580	Shachar Don-Yehiya, Elad Venezian, Colin Raffel, Noam Slonim, Yoay Katz, and Leshem Choshen.
581	Cold fusion: Collaborative descent for distributed multitask finetuning, 2022.
582	-
584	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner Mostefa Dehghani, Matthias Minderer Georg Heigold, Sylvain Celly, Jakob Uszkorait
585	and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In
586	International Conference on Learning Representations (ICLR), 2021. https://openreview.
587	net/forum?id=YicbFdNTTy.
588	Ealiy Draylor Kampia Vasahaini Manfrad Salmhafar and Frad Hampracht Essantially as harrise
589	in neural network energy landscape. In International Conference on Machine Learning (ICML)
590	2018. https://arxiv.org/abs/1803.00885.
591	
592	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
593	arXiv preprint arXiv:2407.21783, 2024.

594 595 596	Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. The role of permutation invariance in linear mode connectivity of neural networks. <i>arXiv preprint arXiv:2110.06296</i> , 2021.
597 598 599	Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In <i>International Conference on Machine Learning (ICML)</i> , 2020. https://proceedings.mlr.press/v119/frankle20a.html.
600 601 602	C Daniel Freeman and Joan Bruna. Topology and geometry of half-rectified network optimization. <i>arXiv preprint arXiv:1611.01540</i> , 2016.
603 604 605	Timur Garipov, Pavel Izmailov, Dmitrii Podoprikhin, Dmitry Vetrov, and Andrew Gordon Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. In <i>Advances in Neural Information</i> <i>Processing Systems (NeurIPS)</i> , 2018. https://arxiv.org/abs/1802.10026.
607 608 609	Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vlad Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. Arcee's mergekit: A toolkit for merging large language models. <i>arXiv preprint arXiv:2403.13257</i> , 2024.
610 611 612	Almog Gueta, Elad Venezian, Colin Raffel, Noam Slonim, Yoav Katz, and Leshem Choshen. Knowl- edge is a region in weight space for fine-tuned language models. <i>arXiv preprint arXiv:2302.04863</i> , 2023.
613 614 615 616	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. In <i>International Conference on Learning Representations</i> , 2021. URL https://openreview.net/forum?id=XPZIaotutsD.
617 618 619	Yifei He, Yuzheng Hu, Yong Lin, Tong Zhang, and Han Zhao. Localize-and-stitch: Efficient model merging via sparse task arithmetic. <i>arXiv preprint arXiv:2408.13656</i> , 2024.
620 621 622	Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Cosmos qa: Machine reading comprehension with contextual commonsense reasoning. <i>arXiv preprint arXiv:1909.00277</i> , 2019.
623 624 625 626 627	Shih-Cheng Huang, Pin-Zu Li, Yu-chi Hsu, Kuang-Ming Chen, Yu Tung Lin, Shih-Kai Hsiao, Richard Tsai, and Hung-yi Lee. Chat vector: A simple approach to equip LLMs with instruction following and model alignment in new languages. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL 2024)</i> , pp. 10943–10959, 2024a. URL https://aclanthology.org/2024.acl-long.590.
628 629 630 631 632 633 634	Shih-Cheng Huang, Pin-Zu Li, Yu-chi Hsu, Kuang-Ming Chen, Yu Tung Lin, Shih-Kai Hsiao, Richard Tsai, and Hung-yi Lee. Chat vector: A simple approach to equip LLMs with instruction following and model alignment in new languages. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 10943–10959, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.590. URL https: //aclanthology.org/2024.acl-long.590.
635 636 637 638	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. <i>arXiv preprint</i> <i>arXiv:2212.04089</i> , 2022.
639 640 641 642	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In <i>The Eleventh International Confer-</i> <i>ence on Learning Representations</i> , 2023. URL https://openreview.net/forum?id= 6t0Kwf8-jrj.
643 644 645 646	Shankar Iyer, Nikhil Dandekar, and Kornel Csernai. First quora dataset release: Question pairs, 2017. URL https://data.quora.com/ First-Quora-Dataset-Release-Question-Pairs.

647 Dong-Hwan Jang, Sangdoo Yun, and Dongyoon Han. Model stock: All we need is just a few fine-tuned models, 2024. URL https://arxiv.org/abs/2403.19522.

648 649 650	Joel Jang, Seungone Kim, Seonghyeon Ye, Doyoung Kim, Lajanugen Logeswaran, Moontae Lee, Kyungjae Lee, and Minjoon Seo. Exploring the benefits of training expert language models over instruction tuning. <i>arXiv preprint arXiv:2302.03202</i> , 2023.
652 653	Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. Dataless knowledge fusion by merging weights of language models. <i>arXiv preprint arXiv:2212.09849</i> , 2022.
654 655 656 657	Xisen Jin, Xiang Ren, Daniel Preotiuc-Pietro, and Pengxiang Cheng. Dataless knowledge fusion by merging weights of language models. In <i>The Eleventh International Conference on Learning</i> <i>Representations</i> , 2023. URL https://openreview.net/forum?id=FCnohuR6AnM.
658 659	Keller Jordan, Hanie Sedghi, Olga Saukh, Rahim Entezari, and Behnam Neyshabur. Repair: Renor- malizing permuted activations for interpolation repair. <i>arXiv preprint arXiv:2211.08403</i> , 2022.
660 661 662 663 664	Keller Jordan, Hanie Sedghi, Olga Saukh, Rahim Entezari, and Behnam Neyshabur. REPAIR: REnormalizing permuted activations for interpolation repair. In <i>The Eleventh International</i> <i>Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum? id=gU5sJ6ZggcX.
665 666 667 668	Rohith Kuditipudi, Xiang Wang, Holden Lee, Yi Zhang, Zhiyuan Li, Wei Hu, Rong Ge, and Sanjeev Arora. Explaining landscape connectivity of low-cost solutions for multilayer nets. <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2019. https://arxiv.org/abs/1906.06247.
669 670 671	Rémi Lebret, David Grangier, and Michael Auli. Neural text generation from structured data with application to the biography domain. <i>arXiv preprint arXiv:1603.07771</i> , 2016.
672 673 674 675	Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. Dbpedia–a large- scale, multilingual knowledge base extracted from wikipedia. <i>Semantic web</i> , 6(2):167–195, 2015.
676 677 678	Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In <i>Thirteenth international conference on the principles of knowledge representation and reasoning</i> , 2012a.
679 680 681	Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. <i>Thirteenth</i> <i>International Conference on the Principles of Knowledge Representation and Reasoning</i> , 2012b.
682 683 684 685	Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. CommonGen: A constrained text generation challenge for generative commonsense reasoning. In <i>Findings of the Association for Computational Linguistics: EMNLP</i> , 2020. https://www.aclweb.org/anthology/2020.findings-emnlp.165.
686 687 688	Kevin Lin, Oyvind Tafjord, Peter Clark, and Matt Gardner. Reasoning over paragraph effects in situations. <i>arXiv preprint arXiv:1908.05852</i> , 2019.
689 690 691	Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. <i>Advances in Neural Information Processing Systems</i> , 35:1950–1965, 2022.
692 693 694 695	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019. https://arxiv.org/abs/1907.11692.
696 697 698	Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. <i>arXiv preprint arXiv:2301.13688</i> , 2023.
700 701	Wei Lu, Rachel K Luu, and Markus J Buehler. Fine-tuning large language models for domain adaptation: Exploration of training strategies, scaling, model merging and synergistic capabilities. <i>arXiv preprint arXiv:2409.03444</i> , 2024.

702 703	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Oinguyai Lin, Shifang Chan, and Dangmai Zhang, Wigardmath, Empowering methamatical
704	reasoning for large language models via reinforced evol instruct. arXiv preprint arXiv:2308.00583
705	2023
706	
707	Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts.
708	Learning word vectors for sentiment analysis. In Proceedings of the 49th annual meeting of the
709	association for computational linguistics: Human language technologies, pp. 142–150, 2011.
710	Michael S Matons and Colin A Doffel Marging models with fisher weighted everyging Advances in
711	Neural Information Processing Systems 35:17703–17716 2022a
712	Neura Information 1 rocessing Systems, 55.17765 17716, 2022a.
713	Michael S Matena and Colin A Raffel. Merging models with fisher-weighted averaging. Advances in
714	Neural Information Processing Systems, 35:17703–17716, 2022b.
715	Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
716	Communication-efficient learning of deep networks from decentralized data. In Artificial intelli-
717	gence and statistics, pp. 1273–1282. PMLR, 2017.
718	
719	wonammed Muqeeth, Haokun Liu, Yutan Liu, and Colin Kattel. Learning to route among specialized
720	Adrian Waller, Nuria Oliver, Ionathan Scarlett, and Ealiy Barkankamp (ads.), Proceedings of the
721	41st International Conference on Machine Learning, volume 235 of Proceedings of Machine
722	Learning Research pp 36829–36846 PMLR 21–27 Jul 2024 URL https://proceedings
723	mlr.press/v235/mugeeth24a.html.
724	
725	Shashi Narayan, Shay B Cohen, and Mirella Lapata. Don't give me the details, just the sum-
726	mary! topic-aware convolutional neural networks for extreme summarization. arXiv preprint
727	arXiv:1808.08/45, 2018.
728	Behnam Nevshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learning?
729	Advances in neural information processing systems, 33:512–523, 2020.
730	
731	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial
732	nli: A new benchmark for natural language understanding. arXiv preprint arXiv:1910.14599, 2019.
734	Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. Task arithmetic in the tangent
735	space: Improved editing of pre-trained models. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 36, 2024.
730	
729	Fidel A Guerrero Peña, Heitor Rapela Medeiros, Thomas Dubail, Masih Aminbeidokhti, Eric
730	Granger, and Marco Pedersoli. Re-basin via implicit sinkhorn differentiation. In <i>Proceedings of</i>
739	the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 20237–20246, 2023.
740	Mohammad Taher Pilehvar and Jose Camacho-Collados. Wic: the word-in-context dataset for
742	evaluating context-sensitive meaning representations. <i>arXiv preprint arXiv:1808.09121</i> , 2018.
743	
744	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Gon, Sandhini Agarwal,
745	models from natural language supervision. In International conference on machine learning, pp
746	8748–8763 PMLR 2021
747	0/10/0/00/11/1ER, 2021.
748	Colin Raffel. A call to build models like we build open-source software, 2021.
749	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi
750	Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text
/51	transformer. Journal of Machine Learning Research (JMLR), 2020a. http://jmlr.org/
752	papers/v21/20-074.html.
753	Colin Doffal Noom M. Shazaar, Adam Daharta, Katharing Lao, Sharan Narang, Mishari Matara
754 755	Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>ArXiv</i> , abs/1910.10683, 2020b.

756 757 758 750	Alexandre Rame, Nino Vieillard, Leonard Hussenot, Robert Dadashi, Geoffrey Cideron, Olivier Bachem, and Johan Ferret. Warm: On the benefits of weight averaged reward models. In <i>Forty-first International Conference on Machine Learning</i> .
759 760 761 762	Alexandre Ramé, Kartik Ahuja, Jianyu Zhang, Matthieu Cord, Léon Bottou, and David Lopez-Paz. Model ratatouille: Recycling diverse models for out-of-distribution generalization. <i>arXiv preprint</i> <i>arXiv:2212.10445</i> , 2022.
763 764 765 766	Alexandre Rame, Guillaume Couairon, Corentin Dancette, Jean-Baptiste Gaya, Mustafa Shukor, Laure Soulier, and Matthieu Cord. Rewarded soups: towards pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
767 768 769	Alexandre Ramé, Johan Ferret, Nino Vieillard, Robert Dadashi, Léonard Hussenot, Pierre-Louis Cedoz, Pier Giuseppe Sessa, Sertan Girgin, Arthur Douillard, and Olivier Bachem. Warp: On the benefits of weight averaged rewarded policies. <i>arXiv preprint arXiv:2406.16768</i> , 2024.
771 772	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In 2011 AAAI spring symposium series, 2011.
773 774 775	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> , 2023.
776 777 778 770	Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training enables zero-shot task generalization. <i>arXiv preprint arXiv:2110.08207</i> , 2021a.
780 781 782 783	Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training enables zero-shot task generalization. In <i>International Conference on Learning Representations (ICLR)</i> , 2021b. https://arxiv.org/abs/2110.08207.
784 785	Abigail See, Peter J Liu, and Christopher D Manning. Get to the point: Summarization with pointer-generator networks. <i>arXiv preprint arXiv:1704.04368</i> , 2017.
786 787 788 789	Viraj Shah, Nataniel Ruiz, Forrester Cole, Erika Lu, Svetlana Lazebnik, Yuanzhen Li, and Varun Jampani. Ziplora: Any subject in any style by effectively merging loras. <i>arXiv preprint arXiv:2311.13600</i> , 2023.
790 791	Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In <i>International Conference on Machine Learning</i> . PMLR, 2018.
792 793 794	Ken Shoemake. Animating rotation with quaternion curves. In <i>Proceedings of the 12th annual conference on Computer graphics and interactive techniques</i> , pp. 245–254, 1985.
795 796	Sidak Pal Singh and Martin Jaggi. Model fusion via optimal transport. <i>Advances in Neural Information</i> <i>Processing Systems</i> , 33:22045–22055, 2020.
797 798 799	Kai Sun, Dian Yu, Jianshu Chen, Dong Yu, Yejin Choi, and Claire Cardie. Dream: A challenge data set and models for dialogue-based reading comprehension. <i>Transactions of the Association for Computational Linguistics</i> , 7:217–231, 2019.
800 801 802	Yi-Lin Sung, Linjie Li, Kevin Lin, Zhe Gan, Mohit Bansal, and Lijuan Wang. An empirical study of multimodal model merging. <i>Empirical Methods in Natural Language Processing (Findings)</i> , 2023.
803 804	Derek Tam, Mohit Bansal, and Colin Raffel. Merging by matching models in task subspaces. <i>arXiv</i> preprint arXiv:2312.04339, 2023.
805 806 807	Derek Tam, Mohit Bansal, and Colin Raffel. Merging by matching models in task parameter subspaces. <i>Transactions on Machine Learning Research</i> , 2024.
808 809	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.

823

824

825

830

831

832

833

847

853

859

861

- 810 Anke Tang, Li Shen, Yong Luo, Han Hu, Bo Du, and Dacheng Tao. Fusionbench: A comprehensive 811 benchmark of deep model fusion. arXiv preprint arXiv:2406.03280, 2024. 812
- Norman Tatro, Pin-Yu Chen, Payel Das, Igor Melnyk, Prasanna Sattigeri, and Rongjie Lai. Optimizing 813 mode connectivity via neuron alignment. Advances in Neural Information Processing Systems, 33: 814 15300-15311, 2020. 815
- 816 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya 817 Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 818 Gemma 2: Improving open language models at a practical size. arXiv preprint arXiv:2408.00118, 819 2024. 820
- Fanqi Wan, Xinting Huang, Deng Cai, Xiaojun Quan, Wei Bi, and Shuming Shi. Knowledge fusion 821 of large language models. arXiv preprint arXiv:2401.10491, 2024a. 822
 - Fanqi Wan, Ziyi Yang, Longguang Zhong, Xiaojun Quan, Xinting Huang, and Wei Bi. Fusechat: Knowledge fusion of chat models. arXiv preprint arXiv:2402.16107, 2024b.
- 826 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: 827 A multi-task benchmark and analysis platform for natural language understanding. In International 828 Conference on Learning Representations (ICLR), 2018. https://arxiv.org/abs/1804. 829 07461.
 - Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. In International Conference on Learning Representations, 2020.
- 834 Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, 835 Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig 836 Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In Proceedings of the 39th International Conference on Machine 837 Learning, volume 162 of Proceedings of Machine Learning Research, pp. 23965–23998. PMLR, 838 2022a. URL https://proceedings.mlr.press/v162/wortsman22a.html. 839
- 840 Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, 841 Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust 842 fine-tuning of zero-shot models. In Proceedings of the IEEE/CVF conference on computer vision 843 and pattern recognition, pp. 7959-7971, 2022b. 844
- 845 Prateek Yadav, Leshem Choshen, Colin Raffel, and Mohit Bansal. Compett: Compression for communicating parameter efficient updates via sparsification and quantization, 2023a. 846
- Prateek Yadav, Qing Sun, Hantian Ding, Xiaopeng Li, Dejiao Zhang, Ming Tan, Parminder Bhatia, 848 Xiaofei Ma, Ramesh Nallapati, Murali Krishna Ramanathan, Mohit Bansal, and Bing Xiang. 849 Exploring continual learning for code generation models. In Proceedings of the 61st Annual 850 Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pp. 782–792, 851 Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023. 852 acl-short.68. URL https://aclanthology.org/2023.acl-short.68.
- Prateek Yadav, Colin Raffel, Mohammed Muqeeth, Lucas Caccia, Haokun Liu, Tianlong Chen, Mohit 854 Bansal, Leshem Choshen, and Alessandro Sordoni. A survey on model moerging: Recycling and 855 routing among specialized experts for collaborative learning, 2024a. URL https://arxiv. 856 org/abs/2408.07057.
- 858 Prateek Yaday, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. Ties-merging: Resolving interference when merging models. Advances in Neural Information Processing Systems, 860 36, 2024b.
- Enneng Yang, Zhenyi Wang, Li Shen, Shiwei Liu, Guibing Guo, Xingwei Wang, and Dacheng Tao. 862 Adamerging: Adaptive model merging for multi-task learning. arXiv preprint arXiv:2310.02575, 863 2023.

- Enneng Yang, Li Shen, Guibing Guo, Xingwei Wang, Xiaochun Cao, Jie Zhang, and Dacheng Tao.
 Model merging in llms, mllms, and beyond: Methods, theories, applications and opportunities.
 arXiv preprint arXiv:2408.07666, 2024.
- Yi Yang, Wen-tau Yih, and Christopher Meek. WikiQA: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 2013–2018, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1237. URL https://aclanthology.org/D15–1237.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
- Peng Ye, Chenyu Huang, Mingzhu Shen, Tao Chen, Yongqi Huang, Yuning Zhang, and Wanli Ouyang.
 Merging vision transformers from different tasks and domains. *arXiv preprint arXiv:2312.16240*, 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, 2024a.
- Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Extend model merging from fine-tuned to pre-trained large language models via weight disentanglement. *arXiv preprint arXiv:2408.03092*, 2024b.
 - Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
 - Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In Advances in Neural Information Processing Systems (NeurIPS), 2015. https://proceedings.neurips.cc/paper/2015/file/ 250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf.
 - Yiran Zhao, Wenxuan Zhang, Huiming Wang, Kenji Kawaguchi, and Lidong Bing. Adamergex: Cross-lingual transfer with large language models via adaptive adapter merging. *arXiv preprint arXiv:2402.18913*, 2024.
 - Jing Zhou, Zongyu Lin, Yanan Zheng, Jian Li, and Zhilin Yang. Not all tasks are born equal: Understanding zero-shot generalization. In *The Eleventh International Conference on Learning Representations*, 2022.
 - A EXTENDED RELATED WORK

885

887

889

890

891

892

893

894 895

896

897 898 899

900 901

902

A.1 LOSS LANDSCAPE AND WEIGHT INTERPOLATION

903 While the loss function of a neural network is generally non-convex, recent work (Draxler et al., 2018; Freeman & Bruna, 2016; Garipov et al., 2018; Jordan et al., 2023; Gueta et al., 2023) has 904 demonstrated that the parameter values from different training runs can sometimes be interpolated 905 without increasing the loss (i.e. they are mode-connected). Many methods (Kuditipudi et al., 2019; 906 Tatro et al., 2020; Benton et al., 2021) have explored finding these low-loss paths between models, 907 focusing on simple (not necessarily linear) interpolations. For example, Frankle et al. (2020) showed 908 that if a part of the optimization trajectory is shared between two neural networks then they can be 909 interpolated without lowering accuracy. On the other hand, Neyshabur et al. (2020) showed that 910 naively interpolating two neural networks with completely disjoint optimization trajectories can result 911 in a catastrophic drop in their accuracies. Entezari et al. (2021) hypothesized that if we account 912 for the permutation symmetry of neural networks, then all neural networks of a given architecture 913 trained on the same dataset are linear mode connected. This assumption of the existence of a low-loss 914 "basin" in parameter space encompassing the models is critical for model merging (Ilharco et al., 915 2023). Ainsworth et al. (2022); Singh & Jaggi (2020); Wang et al. (2020); Jordan et al. (2022); Peña et al. (2023) therefore used techniques based on finding permutations (Wang et al., 2020; Ainsworth 916 et al., 2022) and optimal transport (Singh & Jaggi, 2020) to better align neural networks trained from 917 scratch so that they can be merged or interpolated without increasing the loss.

918 A.2 MODEL MERGING 919

920

Section 2.1 discusses the merging methods that we use for our experiments, however, the popularity of 921 model merging has led to a ever-growing number of methods and applications of model merging (He 922 et al., 2024; Daheim et al., 2023; Yadav et al., 2023a;b; 2024b; Matena & Raffel, 2022a; Jin et al., 923 2023). Next, we discuss some of these methods which were omitted due to large scale practical 924 considerations. Tangent Task Arithmetic (Ortiz-Jimenez et al., 2024) fine-tune models in the tangent 925 space for better weight disentanglement when using Task Arithmetic. Akiba et al. (2024) explore 926 using evolutionary algorithms to choose which layers to merge. SLERP (Shoemake, 1985) and Model Stock (Jang et al., 2024) consider the geometric properties in weight space where SLERP performs 927 spherical interpolation of model weights while Model Stock approximates a center-close weight based 928 on several FT models, utilizing their backbone as an anchor point. Tang et al. (2023) train a mask that 929 learns which parameters are important for the merged model. Ye et al. (2023) train a gating network to 930 predict a weight that is then used to compute a weighted average of examples during inference. Yaday 931 et al. (2024a) provides a comprehensive survey of methods that train a router to route between the 932 different models to merge. Moreover, other applications of model merging include intermediate-task 933 training (Ramé et al., 2022; Choshen et al., 2022a;b), continual learning (Don-Yehiya et al., 2022), 934 model alignment (Rame et al., 2024; Rame et al.; Ramé et al., 2024), merging pretrained models Yu 935 et al. (2024b), or merging models in different modalities (Sung et al., 2023). 936

937

940

938 939

В DETAILED TASK DESCRIPTIONS.

941 We adopt the experimental setting from the T0 mixture (Sanh et al., 2021a) which contains 8 held-in 942 and 4 held-out task categories Specifically, the 8 held-in task categories include Multiple-choice QA 943 (with selected datasets DREAM (Sun et al., 2019), Cosmos OA (Huang et al., 2019)), Extractive Oa 944 (Adversarial QA (Adelani et al., 2021), ROPES (Lin et al., 2019)), Closed-Book QA (Hotpot QA (Yang 945 et al., 2018), Wiki QA (Yang et al., 2015)), Sentiment Analysis (App Reviews (), IMDB (Maas et al., 946 2011)), Topic Classification (AG News (Zhang et al., 2015), DBPedia (Lehmann et al., 2015)), 947 Structure-to-text (Common Gen (Lin et al., 2020), Wiki Bio (Lebret et al., 2016)), Summarization (CNN Daily Mail (See et al., 2017), XSum (Narayan et al., 2018)) and Paraphrase Identification 948 (MRPC (Dolan & Brockett, 2005), QQP (Iyer et al., 2017)). Similary, the 4 held-out task categories 949 are Sentence Completion (with selected dataset COPA (Roemmele et al., 2011), HellaSwag (Zellers 950 et al., 2019)), Natural Language Inference (ANLI (Nie et al., 2019), RTE (Dagan et al., 2005)), Coreference Resolution (WSC (Levesque et al., 2012b), Winogrande (Levesque et al., 2012a)) and 952 Word Sense Disambiguation (WiC (Pilehvar & Camacho-Collados, 2018)).

953 954

951

955 956

С **EXPERT TRAINING DETAILS**

957 958

In our research, we utilized two base models, namely PaLM-2 and PaLM-2-IT to create specialized 959 expert models. We train the PaLM-2model for an additional 60000 steps on the Flan-v2 dataset (Long-960 pre et al., 2023) to obtain the PaLM-2-IT model. We removed the T0 tasks from the flan mixture in 961 order to training experts on them in future. Many of these training jobs were early stopped due to 962 convergence. We used Sharded Adafactor (Shazeer & Stern, 2018) optimizer along with a cosine 963 decay and a learning rate of 1e-4 for 1B, 24B, and 64B model sizes and 3e-5 for 8B model. We use a 964 dropout value of 0.05. Following Chung et al. (2024), we used an input length of 2048 and output 965 length of 512. To create expert models we perform full finetuning with the following hyperparameters. 966 For training the experts model, for all model size, we train by default for 2000 steps with a learning 967 rate of 3e-5 and dropout of 0.05. For some task we adjust the number of steps depending upon the 968 convergence. For the purpose of evaluating classification tasks (Raffel et al., 2020b), we perform rank 969 *classification*. In this method, the model's log probabilities for all potential label strings are ranked. The model's prediction is deemed accurate if the choice ranked highest aligns with the correct answer. 970 It should be noted that rank classification evaluation can accommodate both classification tasks and 971 multiple-choice tasks.

FULL RESULT TABLES D

In this section, we provide the result for the full grid of experiments that we performed. The results contain information about any of the plots that are not provided in the main paper. Table 3 and 4 present the held-in and held-out performance of PaLM-2 model across all model sizes, base models, merging methods, and the number of experts being merged. Similarly, Table 1 and 2 present the held-in and held-out performance of PaLM-2-IT model.

Table 1: The table reports the average normalized performance for the held-in tasks when merging experts created from PaLM-2-IT base models.

Merging Method (\downarrow)		1	в			8	в			24	4B		64B				
# of Experts (\rightarrow)	2	4	6	8	2	4	6	8	2	4	6	8	2	4	6	8	
Average	0.85	0.78	0.81	0.83	0.90	0.82	0.82	0.85	0.94	0.84	0.80	0.77	0.97	0.91	0.89	0.93	
Task Arithmetic	0.91	0.82	0.84	0.86	0.95	0.86	0.85	0.88	0.96	0.90	0.91	0.92	1.00	0.91	0.90	0.93	
Dare-TIES	0.90	0.81	0.83	0.86	0.93	0.86	0.84	0.88	0.94	0.89	0.87	0.88	0.97	0.91	0.89	0.93	
TIES	0.89	0.81	0.82	0.85	0.93	0.86	0.84	0.88	0.95	0.88	0.86	0.86	0.97	0.90	0.89	0.93	
Multitask	0.97	0.96	0.96	0.96	0.96	0.96	0.97	0.96	0.99	0.97	0.98	0.98	0.99	0.98	0.98	0.99	

Table 2: The table reports the average normalized performance on the held-out tasks when merging experts created from PaLM-2-IT base models.

Merging Method (\downarrow)		1	в			8	в		24B					64B				
# of Experts (\rightarrow)	2	4	6	8	2	4	6	8	2	4	6	8	2	4	6	8		
Average	0.99	1.00	1.04	1.05	1.03	1.02	1.03	1.02	1.05	1.10	1.11	1.16	1.00	1.03	1.06	1.09		
Task Arithmetic	1.03	1.03	1.04	1.05	1.06	1.05	1.05	1.03	1.05	1.10	1.13	1.18	1.00	1.03	1.06	1.09		
Dare-TIES	1.02	1.03	1.04	1.05	1.05	1.04	1.04	1.03	1.05	1.10	1.12	1.17	1.00	1.03	1.06	1.09		
TIES	1.02	1.03	1.04	1.05	1.06	1.05	1.06	1.04	1.04	1.09	1.11	1.16	1.00	1.03	1.06	1.10		
Multitask	1.11	1.11	1.11	1.11	1.12	1.12	1.12	1.12	1.18	1.18	1.18	1.18	1.05	1.05	1.05	1.05		

Table 3: The table reports the average normalized performance on the held-in tasks when merging experts created from PaLM-2 base models.

Merging Method (\downarrow)	1B				8B					24	4B		64B			
# of Experts (\rightarrow)	2	4	6	8	2	4	6	8	2	4	6	8	2	4	6	8
Average	0.63	0.44	0.36	0.26	0.66	0.53	0.50	0.32	0.70	0.48	0.51	0.27	0.80	0.74	0.69	0.67
Task Arithmetic	0.66	0.52	0.44	0.39	0.68	0.54	0.54	0.42	0.72	0.56	0.60	0.46	0.80	0.74	0.69	0.67
Dare-TIES	0.65	0.51	0.42	0.37	0.66	0.51	0.51	0.32	0.67	0.44	0.51	0.27	0.80	0.74	0.69	0.67
TIES	0.66	0.50	0.41	0.33	0.67	0.52	0.48	0.29	0.68	0.49	0.52	0.27	0.80	0.71	0.65	0.56
Multitask	0.88	0.88	0.88	0.87	1.06	1.04	1.04	1.06	1.25	1.15	1.11	1.20	0.97	0.96	0.96	0.96

Table 4: The table reports the average normalized performance on the held-out tasks when merging experts created from PaLM-2 base models.

Merging Method (\downarrow)	1B				8B					24	4B		64B				
# of Experts (\rightarrow)	2	4	6	8	2	4	6	8	2	4	6	8	2	4	6	8	
Average	0.98	1.00	1.02	1.04	1.01	0.97	1.02	0.98	0.95	0.85	0.93	0.83	1.28	1.24	1.29	1.25	
Task Arithmetic	1.01	1.03	1.05	1.07	1.06	1.03	1.04	1.00	1.05	1.03	1.10	1.08	1.29	1.28	1.36	1.35	
Dare-TIES	0.99	1.01	1.04	1.05	1.02	1.00	1.05	1.01	0.97	0.89	0.99	0.90	1.28	1.24	1.28	1.24	
TIES	1.05	1.06	1.03	1.04	1.07	1.04	1.02	0.99	1.01	0.93	0.98	0.90	1.31	1.22	1.24	1.15	
Multitask	1.10	1.10	1.10	1.10	1.62	1.62	1.62	1.62	1.51	1.51	1.51	1.51	1.73	1.73	1.73	1.72	