# Continual Optimization with Symmetry Teleportation for Multi-Task Learning

Zhipeng Zhou[1],    Ziqiao Meng[2],    Pengcheng Wu[1],    Peilin Zhao[3],    Chunyan Miao[1]

[1]Nanyang Technological University,    [2]The Chinese University of Hong Kong,    [3]Tencent AI Lab

## Abstract

*Multi-task learning (MTL) is a widely explored paradigm that enables the simultaneous learning of multiple tasks using a single model. Despite numerous solutions, the key issues of optimization conflict and task imbalance remain under-addressed, limiting performance. Unlike existing optimization-based approaches that typically reweight task losses or gradients to mitigate conflicts or promote progress, we propose a novel approach based on **C**ontinual **O**ptimization with **S**ymmetry **T**eleportation (COST). During MTL optimization, when an optimization conflict arises, we seek an alternative loss-equivalent point on the loss landscape to reduce conflict. Specifically, we utilize a low-rank adapter (LoRA) to facilitate this practical teleportation by designing convergent, loss-invariant objectives. Additionally, we introduce a historical trajectory reuse strategy to continually leverage the benefits of advanced optimizers. Extensive experiments on multiple mainstream datasets demonstrate the effectiveness of our approach. COST is a plug-and-play solution that enhances a wide range of existing MTL methods. When integrated with state-of-the-art methods, COST achieves superior performance.*

## 1. Introduction

Traditional machine learning typically requires separate models for each task, leading to higher computational and storage demands as the number of tasks increases. To overcome this issue, multi-task learning (MTL) offers an efficient approach, enabling the simultaneous learning of multiple tasks using a single model. Recent developments in MTL methods can be broadly divided into two categories: structure-based [6, 11, 27] and optimization-based [17, 24, 28]. Structure-based methods focus on designing architectures that enhance task learning by utilizing task relationships and promoting individual progress. On the other hand, optimization-based methods prioritize the learning process by addressing challenges such as gradient conflicts and task imbalances. Since this paper concentrates on optimization-based methods, our analysis and comparisons will primarily focus on these approaches.
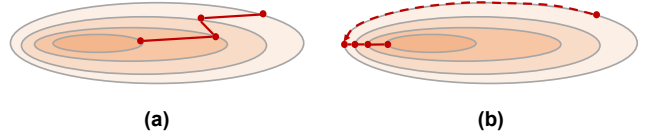


Figure 1. The illustration of symmetry teleportation. (a) is the original gradient descent. (b) is the gradient descent with a faster convergence rate after teleporting the start point from (a).

Optimization-based MTL aims to resolve the aforementioned issues by re-weighting on various aspects. For example, a series of studies [16, 24, 28] explore different gradient combinations to prevent improving some tasks while sacrificing others. Another group of works [4, 17] re-weight task loss to ensure fair progress for individual tasks and thereby address the imbalance issue. While the former group of works endeavors to balance conflict and imbalance issues, the latter focuses on attaining balanced individual progress with minimal concern for the conflict issue. As a result, the latter is generally less robust in different scenarios compared to the former according to empirical observations [4]. However, the former also struggles to achieve a proper balance. In this paper, distinct from these two paradigms and based on the definition of Pareto dominance, we approach MTL from a new perspective, i.e., seeking the less conflicting and more convergent point through symmetry teleportation during MTL optimization.

Unlike the traditional gradient descent regime, symmetry teleportation aims to accelerate the optimization process by seeking another point within the same loss level set, as depicted in Figure 1. Several recent works [1, 30, 31] have explored optimization through symmetry teleportation. For example, [30] introduces a simple teleportation algorithm for non-linear neural networks, based on the assumption that activation functions are bijective, and seeks the point of maximal gradient magnitude using gradient ascent. However, these methods do not provide practical algorithms for larger, modern neural networks, primarily due to their reliance on strict assumptions about non-linearity and computational intensity (Section 3.3). These limitations are especially pronounced in more complex tasks, e.g., MTL.

Therefore, in this paper, we aim to develop a practical

symmetry teleportation method that is applicable for modern deep models, and addressing MTL issues. Specifically, we leverage the low-rank adapter (LoRA) to realize teleportation when encountered with the conflicts issue. By designing the objectives to ensure the invariant task loss and promote progress, we are able to further extend the boundaries of individual task learning for MTL models in a balanced manner. Besides, we also design a historical trajectory reuse strategy to continually benefit from advanced optimizer (e.g., Adam). In a nutshell, our contribution can be summarized as follows:

- We approach MTL from a new angle, i.e., symmetry teleportation, and empirically verify its applicability for MTL (Section 3.2).
- A new practical teleportation method COST is proposed for mitigating the conflict and imbalance issue. To the best of our knowledge, we are the first to develop a practical teleportation method for non-small deep models, specifically for MTL.
- By proposing a historical trajectory reuse strategy, we can continually benefit from the advanced optimizer (e.g., Adam and its variants).
- Taking the advanced method as the baseline, our COST can well augment it to achieve state-of-the-art (SOTA) performance across diverse evaluations. Besides, we also equip mainstream MTL methods with COST, and showing its plug-and-play property.

## 2. Related Work

### 2.1. Optimization-based MTL

Optimization-based methods aim to optimize multiple tasks simultaneously by enhancing the gradient-based learning process itself. For example, MGDA [24] reduces conflicts between task gradients by combining them using the Frank-Wolfe algorithm [12] to generate a gradient with minimal norm. PCGrad [28] addresses gradient conflicts by projecting gradients from different tasks onto directions that minimize interference. CAGrad [16] attempts to balance global optimization and task-specific performance, maintaining both Pareto efficiency and overall optimization with the assistance of a hyperparameter. Nash-MTL [22] introduces a game-theoretic approach where tasks negotiate to update parameters in a manner that enables balanced progression across tasks. Additionally, MoCo [9] focuses on correcting biases in gradient direction by tracking parameters during the learning process, improving gradient alignment and task performance. FairGrad [3] is a pioneering MTL algorithm that puts forward fairness measurements to facilitate maximal loss reduction. It can be considered as an advanced version of Nash-MTL, being capable of balancing task progress in a more fine-grained manner.

### 2.2. Symmetry Teleportation for Deep Model

Before presenting some recent works on symmetry teleportation, we first provide its definition here as per [30]. Let $\mathcal{L}(\boldsymbol{\theta})$ be the loss function. Here, $\mathbb{R}^d$ denotes the model's parameter space, and $A$ represents the acting space on the parameters that leaves the loss value unchanged. Subsequently, we have the following definition:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathcal{L}(a \cdot \boldsymbol{\theta}), \quad \forall a \in A, \quad \forall \boldsymbol{\theta} \in \mathbb{R}^d. \quad (1)$$

$$\boldsymbol{\theta}' = a \cdot \boldsymbol{\theta}, \quad a = \underset{a \in A}{\operatorname{argmax}} \left\| \nabla \mathcal{L}(a \cdot \boldsymbol{\theta}) \right\|^2. \quad (2)$$

we can observe that symmetry teleportation aims to find a loss-invariant point (Eqn. 1) with a maximum gradient norm (Eqn. 2) on the loss level set by acting with a group element.

As a recent research topic, symmetry teleportation has been explored in only a few works [1, 30, 31]. [1] first introduced the concept of 'neural teleportation' and investigated its impact on optimization. [30] proposed a gradient ascent-based teleportation algorithm for small neural networks (e.g., three-layer MLPs). And [31] established the connection between symmetry teleportation and generalization through a series of theoretical analyses and provided an alternative for enhancing the meta optimizer.

### 2.3. Low-Rank Adapter

LoRA is gaining increasing popularity in tandem with the rapid advancement of foundation models and parameter-efficient fine-tuning (PEFT). It operates by maintaining the pre-trained weights of a large model in a fixed state and incorporating small, trainable rank decomposition matrices. During fine-tuning, rather than modifying all the parameters of the model, only these low-rank matrices are subject to update.

Moreover, LoRA has several variants that can attain dynamic rank [29], or quantization [8]. For instance, AdaLoRA [29] adaptively assigns dynamic rank to different parameters, thereby enabling the capture of important updates while preserving efficiency. In contrast, QLoRA [8] introduces 4-bit NormalFloat, double quantization, and paged optimizers to more effectively optimize LoRA, while significantly reducing the required memory.

**Connection and Difference**: Our work tackles conflict and imbalance issues in optimization-based MTL through symmetry teleportation. Specifically, we utilize LoRA to implement practical teleportation. In contrast to previous studies, we explore MTL from a novel perspective and introduce a new teleportation algorithm for modern deep models. This algorithm is scalable, easily integratable, and compatible with both PEFT and MTL.

## 3. Motivation and Empirical Observation

### 3.1. Preliminary

As mentioned, optimization-based MTL approaches operate under the assumption that the model consists of a task-shared backbone network alongside task-specific branches. Consequently, the primary objective of these approaches is to devise gradient combination strategies that optimize the backbone network to yield benefits across all tasks. Let us consider a scenario where there are $K \geq 2$ tasks available, each associated with a differentiable loss function $\mathcal{L}_i(\boldsymbol{\theta})$, where $\boldsymbol{\theta}$ represents the task-shared parameters. The goal of optimization-based MTL is to search for the optimal $\boldsymbol{\theta}^* \in \mathbb{R}^m$ that minimizes the losses for all tasks.

**Definition 1 (Gradient Similarity)** *Denote $\phi_{ij}$ as the angle between two task gradients $\boldsymbol{g_i}$ and $\boldsymbol{g_j}$, and assume $\|\boldsymbol{g_i}\|_2 \leq \|\boldsymbol{g_j}\|_2$, then we define the gradient similarity as $\cos \phi_{ij}$ and the gradients as conflicting when $\cos \phi_{ij} < 0$ (referred as **Weak Conflict**). When the mean gradient $\boldsymbol{g_0}$ is conflicting with $\boldsymbol{g_i}$, we call it as **Dominated Conflict** (see Figure 2).*
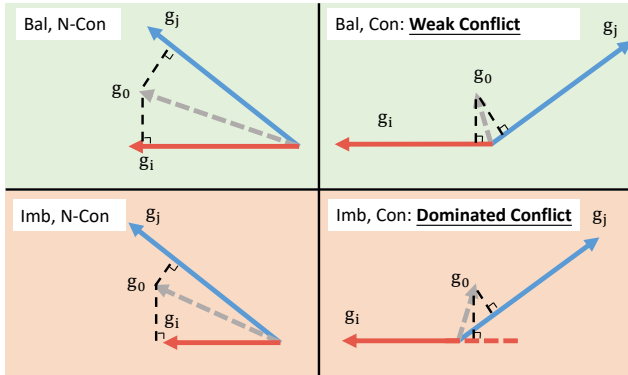


Figure 2. Illustration of conflict and imbalance issues in MTL. 'Bal' and 'Imb' represent balanced and imbalanced, while 'N-Con' and 'Con' represent non-conflicting and conflicting.

### 3.2. Applicability of Symmetry Teleportation

Before delving into the principal design of our method, it is necessary to verify the existence of parameter symmetries with differing conflict statuses. To this end, we examine the optimization process of mainstream MTL approaches. We analyze the mean loss across all tasks and the associated conflict status during optimization from various initial points, with the results presented in Figure 5.

As shown in Figure 5, it is often possible to identify a non-conflict alternative at the same loss level when encountering conflict, demonstrating the potential of symmetry teleportation. Additional statistical results from other MTL approaches are provided in the **Appendix** (Sec. 1.1).
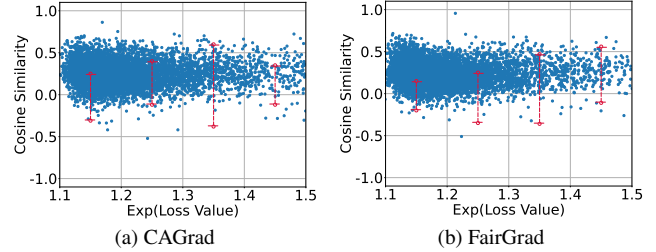


Figure 3. Dominated conflict vs. loss examination. The pink backdrop designates the conflicting area, whereas the green backdrop indicates the non-conflicting area. The blue scatter points are the individual recorded points throughout the optimization process. The red dashed line symbolizes the teleportation occurring from a conflict point to a non-conflict point. An exponential amplification has been applied to the loss values to enhance visual clarity.

### 3.3. Pitfall of Current Paradigms

While several works [1, 30, 31] have proposed symmetry teleportation algorithms for neural network-based models, we demonstrate their limitations with current deep models. First, these algorithms require activation functions to be bijective, which poses a significant challenge for widely-used deep models (e.g., ResNet-50) that use non-bijective activation functions, e.g., ReLU and Sigmoid. Second, they require calculating the pseudo-inverse of inputs layer by layer to ensure output and loss invariance. This process is computationally intensive and may be impractical for modern deep models. As a result, these approaches have only been tested on simple three-layer MLP networks and small-scale datasets (e.g. MNIST) for verification.

## 4. Principal Design

In this section, we present the detailed design of COST, incorporating the symmetry teleportation paradigm and a historical trajectory reuse strategy. We also provide an analysis of convergence.

### 4.1. Continual Optimization with Symmetry Teleportation

The overall framework of COST is depicted in Figure 4. At a certain training stage $t$, we utilize LoRA to teleport the weight of the shared backbone to the non-conflict point (merge the trained LoRA into the backbone's weight) with the same loss level. Subsequently, the model (including both the backbone and branches) is continuously optimized by other MTL algorithms. In this framework, there are two questions that need to be answered:

**When**: The first question is, when should teleportation be triggered? Unfortunately, the previous solutions presented in [1, 30, 31] did not offer a clear answer to this question. They merely triggered it in a random or intuitive manner. In contrast, our goal is to address two key challenges in MTL:
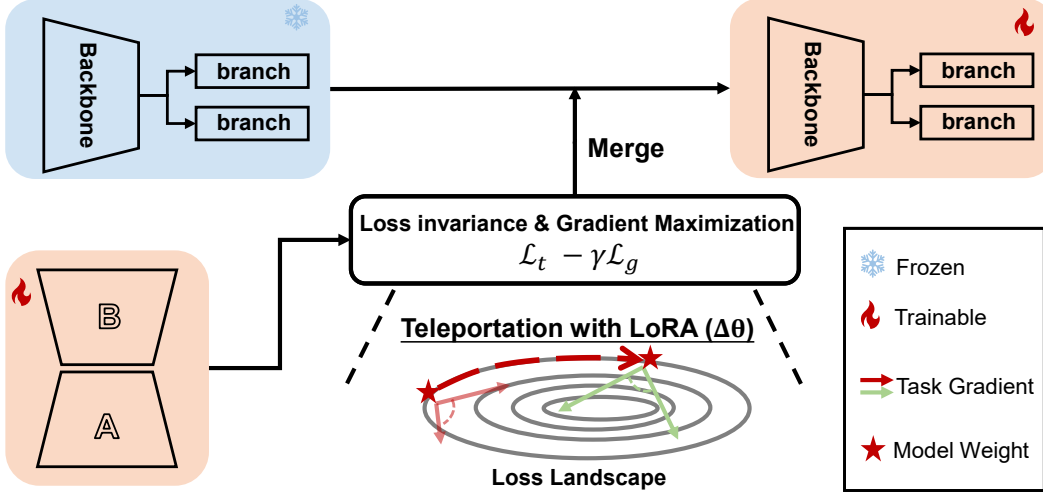
Figure 4. The Illustration of COST. Here, we depict a one-time teleportation procedure by using a 2-task example for the sake of illustration. It is worth noting that LoRA is only applied to the shared backbone.

conflict and imbalance, challenges that are not concurrently addressed by existing solutions. Moreover, a naïve linear scalarization (LS) strategy can effectively promote all tasks, as illustrated in Figure 2 and has been empirically verified in [26]. Thus, the primary challenge lies in resolving conflict arising from imbalance, i.e., dominated conflict. Therefore, we establish the teleportation trigger condition based on the occurrence of dominated conflict [1]:

$$\cos\phi_{i0} < 0, \quad \phi_{i0} = \angle(\boldsymbol{g_i}, \boldsymbol{g_0}) \tag{3}$$

where $\boldsymbol{g_i}$ and $\boldsymbol{g_0}$ represent the task gradient with the smallest norm and the mean gradient, respectively. However, when handling a large number of tasks, dominated conflicts become inevitable, reaching a 97% conflict ratio per epoch on CelebA [20], as shown in Figure 5(a). Then if we still employ dominated conflict as the trigger condition, frequent teleportation would occurs and results in inefficiency. Therefore, our objective shifts to mitigating dominant conflicts, balancing efficiency and effectiveness. To achieve this, we adopt the following condition:

$$\sum_i^K \mathbb{1}[\cos\phi_{i0} < 0] \geq \left\lceil \frac{K}{2} \right\rceil \tag{4}$$

Under this condition, the trigger frequency is significantly reduced (see Figure 5(a)) while maintaining effectiveness, as demonstrated in the evaluation. Additionally, we analyze the trade-off between effectiveness and efficiency for this condition in the **Appendix** (Section 1.2).

---

[1] We have provided a comparison between dominated conflict and weak conflict in the **Appendix** (Sec. 1.2).



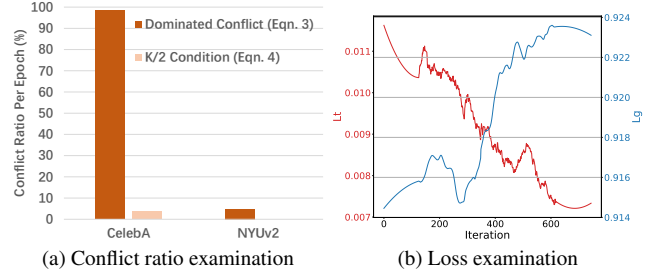(a) Conflict ratio examination     (b) Loss examination

Figure 5. (a) Conflict ratio per epoch on CelebA (40-task) and NYUv2 (3-task) and (b) loss examinations during a single teleportation.

**How**: In the symmetry teleportation paradigm, there are two key objectives: loss invariance and gradient maximization, as outlined in Eqn. 1 and Eqn. 2. Since finding a group action $g$ is infeasible for deep models, we instead use LoRA ($\boldsymbol{\Delta\theta}$) as an alternative, reformulating it as:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta} + \Delta\boldsymbol{\theta}) \tag{5}$$

$$\Delta\boldsymbol{\theta} = \underset{\Delta\boldsymbol{\theta}}{\arg\max} \|\nabla\mathcal{L}(\boldsymbol{\theta} + \Delta\boldsymbol{\theta})\|^2 . \tag{6}$$

With respect to the specific symmetry teleportation taking place during the optimization process, in order to ensure the task loss remains invariant, we undertake the minimization of the loss fluctuation in the following way:

$$\mathcal{L}_t = \frac{1}{K} \sum_i^K |\mathcal{L}_i - \mathcal{L}_i^*| \tag{7}$$

where $\mathcal{L}_i$ represents the individual task loss, and $\mathcal{L}_i^*$ is its loss before starting teleportation, which is unchanged during teleportation.

To maximize the gradient of the target point, a simplistic solution would be to incorporate it into the objective of LoRA optimization. However, two challenges arise when attempting to do so: (1) It is difficult to explicitly incorporate gradient maximization into the objective design. (2) Even if it were possible, the computation of the Hessian matrix would be overly burdensome for LoRA optimization. On the other hand, upon closely examining Eqn. 6, we can observe that our intention is merely to identify the point with the maximal gradient, rather than precisely attaining the maximal gradient itself. Consequently, we choose to select another metric for measuring the gradient norm, i.e., *Sharpness* [10]. Since the negative direction of the gradient is locally the fastest direction of descent, thus we can estimate the gradient by seeking the sharpest direction at $\theta^*$ as follow:

$$\text{Sharpness} = \max_{\|\epsilon\| \leq \delta} |\mathcal{L}(\theta^* + \epsilon) - \mathcal{L}(\theta^*)| \quad (8)$$

where $\delta$ is the radius of the sphere. We further implement Eqn. 8 by randomly sampling $\epsilon$ $\tilde{n}$ times from the sphere, and estimating sharpness by selecting the maximum one. Since $\mathcal{L}(\theta^*)$ remains unchanged for each sampling operation, we obtain the following objective:

$$\mathcal{L}_g = \max \left\{ \left\| \frac{1}{K} \sum_i^K R_i \cdot \mathcal{L}_i(\theta + \Delta\theta + \epsilon_j) \right\| \right\}_{j=1}^{\tilde{n}} \quad (9)$$

$$\mathbf{R} = K \cdot \text{softmax} \left( \left[ \frac{\sum_{j=1}^K \|g_j\|}{\|g_i\|} \right]_{i=1}^K \right) \quad (10)$$

where $\mathbf{R}$ is computed to facilitate the search for more balanced alternatives, mitigating imbalance issues. Consequently, the overall objective for LoRA optimization can be formulated as follows:

$$\mathcal{L}_{lora} = \mathcal{L}_t - \gamma \mathcal{L}_g \quad (11)$$

where $\gamma$ is the hyper-parameter. As depicted in Figure 5(b), $\mathcal{L}_t$ largely decreases while $\mathcal{L}_g$ increases as expected during the teleportation.

To enhance understanding of our approach, we provide a trajectory illustration on toy examples [16] in Figure 6. As shown, LS may fail to reach the Pareto front from certain initializations due to conflict issues. However, with the augmentation of COST, it successfully explores alternative paths for continuous optimization rather than getting stuck.

## 4.2. Convergence Analysis

In this section, we present a convergence analysis to further enhance the understanding of the applicability of our proposed method. By formulating a theorem, it has been proven that our method converges to the Pareto stationary point with guarantee.
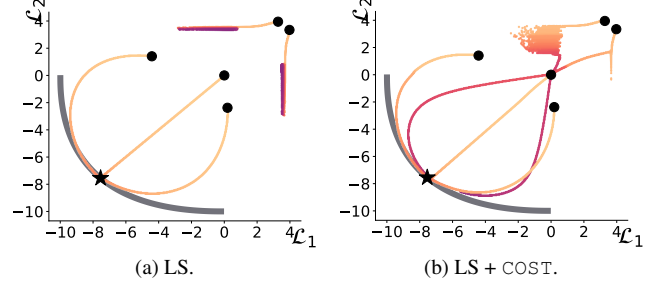


(a) LS.　　　(b) LS + COST.

Figure 6. Trajectory illustration on toy examples.

**Theorem 1** *Assume task loss functions $\mathcal{L}_1, ..., \mathcal{L}_K$ are differentiable and $\Lambda$-smooth ($\Lambda > 0$) such that $\|\nabla\mathcal{L}_i(\theta_1) - \nabla\mathcal{L}_i(\theta_2)\| \leq \Lambda \|\theta_1 - \theta_2\|$ for any two points $\theta_1$, $\theta_2$, and our symmetry teleportation property holds. Set the step size as $\eta = \frac{1}{\Lambda\sqrt{T-1}}$, $T$ is the training iteration. Then, there exists a subsequence $\{\theta^{t_j}\}$ of the output sequence $\{\theta^t\}$ that converges to a Pareto stationary point $\theta^*$.*

The proof of this theorem is provided in the **Appendix** (Sec. 4).

## 4.3. Historical Trajectory Reuse Strategy

When training MTL models with advanced optimizers, a minor issue arises after reaching the loss-invariant point through our symmetry teleportation. Specifically, the teleportation process disrupts the continuous optimization flow, preventing the MTL model from leveraging its historical trajectory—one of the key advantages of advanced optimizers (e.g., Adam [14]).
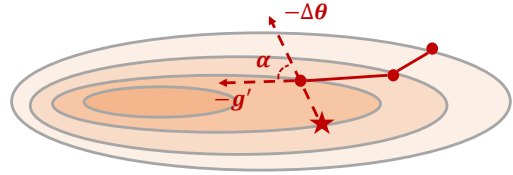


Figure 7. The illustration of HTR strategy. The red star represents the teleported point. $g'$ is the gradient at the pre-teleportation point, and $\alpha$ is the angle between $g'$ and $\Delta\theta$.

Taking the Adam optimizer as an example, which is commonly utilized in mainstream MTL approaches [16, 17, 22]. It employs an exponentially weighted moving average to estimate the momentum ($v_t$) and quadratic moments ($s_t$) of the gradient (historical trajectory). However, when the model is teleported to another point, the stored historical trajectory would supply misleading information for the current model optimization. To tackle this issue, we partially preserve the historical trajectory by computing the correlation between teleportation and previous updating (as de-

picted in Figure 7):

$$\sigma = \cos\_{\text{sim}}(\Delta\boldsymbol{\theta}, \boldsymbol{g}') \tag{12}$$

where $\cos\_{\text{sim}}$ is the function of computing cosine similarity, and $\boldsymbol{g}'$ is the gradient at the pre-teleportation point. In this way, the historical trajectory can be modulated and reused according to Eqn. 13. After that, $\sigma$ is still set to 1, as is typically the case with the Adam optimizer.

$$\boldsymbol{v_t} = \sigma\beta_1\boldsymbol{v_{t-1}} + (1 - \sigma\beta_1)\boldsymbol{g_t} \tag{13}$$
$$\boldsymbol{s_t} = \sigma\beta_2\boldsymbol{s_{t-1}} + (1 - \sigma\beta_2)\boldsymbol{g_t^2}$$

The overall training algorithm is concluded in the **Appendix** (Section 3).

## 5. Performance Evaluation

In this section, we initially evaluate our method using mainstream MTL benchmarks and compare it with the following baselines: Linear Scalarization (LS), Scale-Invariant (SI), Random Loss Weighting (RLW) as described in [15], Dynamic Weight Average (DWA) from [19], Uncertainty Weighting (UW) detailed in [13], MGDA from [24], GradDrop presented in [5], PCGrad as in [28], CAGrad from [16], IMTL detailed in [18], Nash-MTL from [22], FAMO described in [17], and FairGrad from [3]. Subsequently, we offer some additional analyses regarding conflict and gradient examinations, ablation studies, and plug-and-play verification, etc., to further enhance understanding. We also provide additional analysis on alternatives to PEFT (Sec. 1.3), and time cost (Sec. 1.4) in the **Appendix**. All experiments are carried out on a single Tesla V100 GPU. For more experimental details, please refer to the **Appendix** (Sec. 2). Code will be released once this paper is accepted.

**Evaluation Metric**. In addition to reporting individual performance, we also incorporate a widely used metric, $\Delta\mathbf{m}\%$ [21], which evaluates the overall degradation compared to independently trained models that are considered as the reference oracles. The formal definition of $\Delta\mathbf{m}\%$ is given as:

$$\Delta\mathbf{m}\% = \frac{1}{K}\sum_{k=1}^{K}(-1)^{\delta_k}\frac{M_{m,k} - M_{b,k}}{M_{b,k}} \times 100 \tag{14}$$

where $M_{m,k}$ and $M_{b,k}$ represent the metric $M_k$ for the compared method and the independent model, respectively. The value of $\delta_k$ is assigned as 1 if a higher value is better for $M_k$, and 0 otherwise. Besides, we also report another popular metric named **Mean Rank** (**MR**), which computes the average ranks of each methods across all tasks.

### 5.1. Overall Evaluation

**Dense Prediction.** CityScapes [7] and NYUv2 [25] are two widely-used scene understanding datasets, which are employed for the evaluation of MTL. NYUv2 comprises 1449

Table 1. **Scene understanding** (*CityScapes*, 2 tasks).

| Method | Segmentation (Higher Better) | | Depth (Lower Better) | | MR ↓ | Δm% ↓ |
|---|---|---|---|---|---|---|
| | mIoU | Pix. Acc. | Abs. Err. | Rel. Err. | | |
| Independent | 74.01 | 93.16 | 0.0125 | 27.77 | - | - |
| LS | 75.18 | 93.49 | 0.0155 | 46.77 | 8.25 | 22.60 |
| RLW [15] | 74.57 | 93.41 | 0.0158 | 47.79 | 11.00 | 24.37 |
| DWA [19] | 75.24 | 93.52 | 0.0160 | 44.37 | 8.25 | 21.43 |
| Uncertainty [13] | 72.02 | 92.85 | 0.0140 | 30.13 | 7.50 | 5.88 |
| MGDA [24] | 68.84 | 91.54 | 0.0309 | 33.50 | 11.00 | 44.14 |
| GradDrop [5] | 75.27 | 93.53 | 0.0157 | 47.54 | 7.75 | 23.67 |
| PCGrad [28] | 75.13 | 93.48 | 0.0154 | 42.07 | 8.50 | 18.21 |
| CAGrad [16] | 75.16 | 93.48 | 0.0141 | 37.60 | 7.50 | 11.58 |
| IMTL [18] | 75.33 | 93.49 | 0.0135 | 38.41 | 5.75 | 11.04 |
| Nash-MTL [22] | 75.41 | 93.66 | 0.0129 | 35.02 | 3.00 | 6.82 |
| FAMO [17] | 74.54 | 93.29 | 0.0145 | 32.59 | 8.25 | 8.13 |
| FairGrad [3] | 75.72 | 93.68 | 0.0134 | 32.25 | 2.25 | 5.18 |
| COST | 75.73 | 93.53 | 0.0133 | 31.53 | 2.00 | 4.30 |

annotated images and is utilized for three fine-grained tasks, i.e., semantic segmentation, depth estimation, and surface normal prediction. CityScapes consists of 5000 annotated scene images, which are readied for two fine-grained tasks: semantic segmentation and depth estimation.

In line with the implementation and training strategy of FairGrad [3], we construct our model using SegNet [2] and employ MTAN [19] as the backbone within it. We train our model with the Adam optimizer for a total of 200 epochs, setting the initial learning rate to 1.0e-4 and reducing it to half after 100 epochs. The batch size is set to 2 for NYUv2 and 8 for CityScapes, respectively.

The results obtained on these two datasets are presented in Table 1 and Table 2, respectively. With FairGrad serving as the baseline, our method not only successfully surpasses it but also attains the SOTA performance in terms of **MR** and $\Delta$m%. Specifically, upon closely examining the performance of each individual task, we can note that COST significantly enhances FairGrad on the CityScapes dataset and considerably improves the surface normal prediction task, while also showing some promise on the other tasks on the NYUv2 dataset. These observations clearly demonstrate the effectiveness of our design, which aids in alleviating conflict and facilitating convergence.

**Image Classification.** CelebA [20] is a commonly utilized face attributes dataset that contains over 200,000 images and is annotated with 40 attributes. Recently, it has been adopted as a 40-task MTL benchmark to assess the model's capacity to handle a large number of tasks. In accordance with the setup of FairGrad, we utilize a 9-layer convolutional neural network (CNN) as the backbone and linear layers as the task-specific heads on top of it. We train our model with the Adam optimizer for a total of 15 epochs, setting the initial learning rate to 3.0e-4. Moreover, the batch

Table 2. **Scene understanding** (*NYUv2*, 3 tasks). We report MTAN model performance averaged over 3 random seeds. The best scores are provided in gray , and the second scores are underlined.

| Method | Segmentation | | Depth | | Surface Normal | | | | | MR ↓ | Δm% ↓ |
| | (Higher Better) | | (Lower Better) | | Angle Distance | | Within $t°$ | | | | |
| | | | | | (Lower Better) | | (Higher Better) | | | | |
| | mIoU | Pix. Acc. | Abs Err | Rel Err | Mean | Median | 11.25 | 22.5 | 30 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Independent | 38.30 | 63.76 | 0.68 | 0.28 | 25.01 | 19.21 | 30.14 | 57.20 | 69.15 | - | - |
| LS | 39.29 | 65.33 | 0.55 | 0.23 | 28.15 | 23.96 | 22.09 | 47.50 | 61.08 | 9.44 | 5.46 |
| RLW [15] | 37.17 | 63.77 | 0.58 | 0.24 | 28.27 | 24.18 | 22.26 | 47.05 | 60.62 | 12.22 | 7.67 |
| DWA [19] | 39.11 | 65.31 | 0.55 | 0.23 | 27.61 | 23.18 | 24.17 | 50.18 | 62.39 | 8.56 | 3.49 |
| Uncertainty [13] | 36.87 | 63.17 | 0.54 | 0.23 | 27.04 | 22.61 | 23.54 | 49.05 | 63.65 | 8.78 | 4.01 |
| MGDA [24] | 30.47 | 59.90 | 0.61 | 0.26 | 24.88 | <u>19.45</u> | 29.18 | <u>56.88</u> | <u>69.36</u> | 7.11 | 1.47 |
| GradDrop [5] | 39.39 | 65.12 | 0.55 | 0.23 | 27.48 | 22.96 | 23.38 | 49.44 | 62.87 | 8.89 | 3.61 |
| PCGrad [28] | 38.06 | 64.64 | 0.56 | 0.23 | 27.41 | 22.80 | 23.86 | 49.83 | 63.14 | 9.33 | 3.83 |
| CAGrad [16] | 39.79 | 65.49 | 0.55 | 0.23 | 26.31 | 21.58 | 25.61 | 52.36 | 65.58 | 6.33 | 0.29 |
| IMTL [18] | 39.35 | 65.60 | 0.54 | 0.23 | 26.02 | 21.19 | 26.20 | 53.13 | 66.24 | 5.56 | -0.59 |
| Nash-MTL [22] | <u>40.13</u> | 65.93 | 0.53 | 0.22 | 25.26 | 20.08 | 28.40 | 55.47 | 68.15 | <u>3.11</u> | -4.04 |
| FAMO [17] | 40.30 | 66.07 | 0.56 | 0.21 | 26.67 | 21.83 | 25.61 | 51.78 | 64.85 | 5.44 | 0.16 |
| FairGrad [3] | 39.74 | <u>66.01</u> | 0.54 | 0.22 | <u>24.84</u> | 19.60 | <u>29.26</u> | 56.58 | 69.16 | 3.00 | <u>-4.66</u> |
| COST | 38.06 | 64.71 | <u>0.54</u> | 0.23 | 24.47 | 18.80 | 30.84 | 58.25 | 70.30 | 3.22 | -5.39 |

Table 3. Comparison of methods on *CelebA* and *QM9* datasets with MR and Δm%. The results of FairGrad-R are reported according to the official implementation of FairGrad.

| Method | CelebA | | QM9 | |
| | MR ↓ | Δm% ↓ | MR ↓ | Δm% ↓ |
|---|---|---|---|---|
| LS | 7.08 | 4.15 | 9.09 | 177.6 |
| SI | 8.80 | 7.20 | 5.64 | 77.8 |
| RLW | 5.98 | 1.46 | 10.64 | 203.8 |
| DWA | 7.78 | 2.40 | 8.91 | 175.3 |
| UW | 6.65 | 3.23 | 7.00 | 108.0 |
| MGDA | 11.98 | 14.85 | 8.91 | 120.5 |
| PCGrad | 7.58 | 3.17 | 7.36 | 125.7 |
| CAGrad | 7.13 | 2.48 | 8.09 | 112.8 |
| IMTL-G | <u>5.53</u> | 0.84 | 6.91 | 77.2 |
| Nash-MTL | 5.73 | 2.84 | <u>4.27</u> | 62.0 |
| FAMO | 5.65 | 1.21 | 5.18 | <u>58.5</u> |
| FairGrad-R | 6.35 | 1.15 | 4.82 | 59.9 |
| COST | 4.80 | <u>0.93</u> | 4.18 | 58.3 |

size is set to 256.

The evaluation results are shown in Table 3. Given that our method is mainly developed based on FairGrad, our performance is thus highly associated with it. We conscientiously re-implemented FairGrad using the official code they provided and were able to achieve the reported performance on CityScapes, NYUv2. However, we were unable to do so on CelebA and QM9. Consequently, we only report our re-implemented performance of FairGrad here (referred to as FairGrad-R). As can be observed, COST still

significantly enhances its baseline and attains the SOTA performance according to MR, ranking second according to Δm%. These results demonstrate COST's remarkable ability to handle numerous tasks simultaneously.

**Regression.** QM9 [23] is another widely used MTL dataset specifically for regression tasks. It contains 130,000 organic molecules that are organized as graphs with node and edge features. This task is designed to predict 11 properties having different measurement scales and can also be considered as an evaluation scenario for MTL involving a large number of tasks. Our approach is trained for 300 epochs with a batch size of 120. The initial learning rate is set to 1.0e-3, and a learning rate scheduler is applied to reduce the rate when the validation performance shows no further improvement.

According to Table 3, our method still achieves competitive performance on this specific dataset. However, in comparison to other datasets, it exhibits fewer enhancements over its baseline. One crucial reason for these relatively less satisfactory results is that this task adopts a graph model with only two layers supporting LoRA in current PEFT package, which reduces its effectiveness.

## 5.2. Conflict and Gradient Examinations

Although our method achieves competitive performance, it remains unclear whether it effectively resolves the targeted issues, i.e., conflict mitigation and greater gradient norm discovery. To investigate this, we analyze the training process by recording the results before and after tele-

portation, as shown in Figure 8. The findings indicate that conflict is significantly alleviated, with task gradients becoming positively correlated in most cases after teleportation. Besides, teleportation consistently yields greater gradient norms, confirming the effectiveness of COST's design.
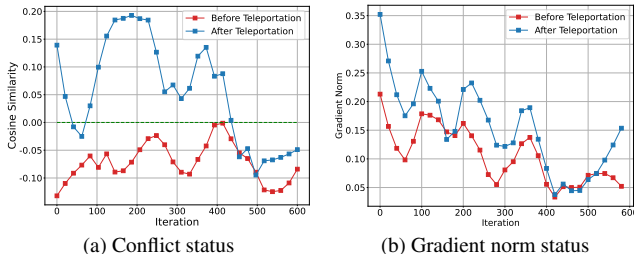


(a) Conflict status          (b) Gradient norm status

Figure 8. Examinations before and after teleportation.

## 5.3. Ablation Study

We consider COST as an integrated system, and thus each component ought to be evaluated to showcase its effectiveness. In our design, there are primarily three key components: the loss invariance objective ($\mathcal{L}_t$), the gradient maximization objective ($\mathcal{L}_g$), and the HTR strategy. Consequently, we carry out ablation studies for verification purposes and present the results in Table 4. In the context of symmetry teleportation, $\mathcal{L}_t$ and $\mathcal{L}_g$ serve as the foundation for seeking alternatives within the orbit. Hence, we exclude $\mathcal{L}_t$ and $\mathcal{L}_g$ during the LoRA optimization process, respectively. The results indicate that without $\mathcal{L}_t$ or $\mathcal{L}_g$, COST performs worse than its baseline (which is FairGrad in this case). More specifically, COST would experience a severe deterioration without $\mathcal{L}_t$, thereby underlining the crucial importance of loss invariance. These results might address another concern regarding COST, namely: Are the improvements brought about by COST rooted in the capability expansion facilitated by LoRA? Without an appropriate objective design, LoRA is unable to effectively augment the base models.

Table 4. Ablation study of COST on *CityScapes* (2 tasks).

| $\mathcal{L}_t$ | $\mathcal{L}_g$ | HTR | $\Delta\mathbf{m}\% \downarrow$ |
|---|---|---|---|
| | | | 5.18 |
| ✓ | | | 7.90 |
| - | ✓ | | **381.86** |
| ✓ | ✓ | | 4.65 |
| ✓ | ✓ | ✓ | 4.30 |

On the other hand, it should be noted that LoRA is incorporated into the base model after each teleportation. Thus,

the model's capability remains unchanged during the inference time. All that we are doing is assisting in finding a better convergence point. Furthermore, when the HTR strategy is excluded, $\Delta$m% decreases from 4.30 to 4.65, which demonstrates the significance of benefiting from advanced optimizers.

## 5.4. Plug-and-Play Verification

Intuitively, our method is orthogonal to existing MTL approaches and is therefore plug-and-play, enabling augmentation when integrated. Here, we take three baselines (i.e., CAGrad, Nash-MTL, and FairGrad) to demonstrate the effectiveness of COST, and present the results in Table 5. As anticipated, our method successfully brings considerable augmentation to its baselines, with improvements ranging from 0.88 to 3.21 according to $\Delta$m%. Specifically, CAGrad and FairGrad receive improvements on almost each individual metric. More results please see the **Appendix** (Sec. 1.5).

Table 5. Plug-and-play verification on *CityScapes* (2 tasks) dataset. We adopt FAMO's implementation for Nash-MTL (denoted as Nash-MTL-R) and augment it with COST, since Nash-MTL does not provide the official implementation on CityScapes.

| Method | Segmentation | | Depth | | $\Delta\mathbf{m}\% \downarrow$ |
|---|---|---|---|---|---|
| | (Higher Better) | | (Lower Better) | | |
| | mIoU | Pix. Acc. | Abs. Err. | Rel. Err. | |
| CAGrad | 75.16 | 93.48 | 0.0141 | 37.60 | 11.58 |
| CAGrad + COST | 75.46 | 93.57 | 0.0134 | 35.68 | 8.37 |
| Nash-MTL-R | 75.87 | 93.57 | 0.0135 | 37.29 | 9.89 |
| Nash-MTL-R + COST | 75.70 | 93.56 | 0.0134 | 34.34 | 7.15 |
| FairGrad | 75.72 | 93.68 | 0.0134 | 32.25 | 5.18 |
| FairGrad + COST | 75.73 | 93.53 | 0.0133 | 31.53 | 4.30 |

## 6. Conclusion

This paper explores the MTL problem from a brand new perspective, i.e., alleviating the conflict issue through symmetry teleportation. Specifically, we utilize LoRA to achieve practical symmetry teleportation for contemporary deep models. Additionally, we design loss-invariant and gradient maximization objectives to assist in identifying non-conflict and more convergent points. We also devise a historical trajectory reuse strategy to continuously benefit from advanced optimizers. Extensive experiments have demonstrated the effectiveness of our proposed method as well as its plug-and-play characteristic. As a scalable framework, we anticipate that our method can offer some valuable insights to researchers engaged in optimization-based MTL. Currently, there are still rooms for improvement within this system, and our future work will focus on these aspects.

# References

[1] Marco Armenta, Thierry Judge, Nathan Painchaud, Youssef Skandarani, Carl Lemaire, Gabriel Gibeau Sanchez, Philippe Spino, and Pierre-Marc Jodoin. Neural teleportation. *Mathematics*, 11(2):480, 2023. 1, 2, 3

[2] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017. 6

[3] Hao Ban and Kaiyi Ji. Fair resource allocation in multi-task learning. *arXiv preprint arXiv:2402.15638*, 2024. 2, 6, 7

[4] Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, and Andrew Rabinovich. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *International conference on machine learning*, pages 794–803. PMLR, 2018. 1

[5] Zhao Chen, Jiquan Ngiam, Yanping Huang, Thang Luong, Henrik Kretzschmar, Yuning Chai, and Dragomir Anguelov. Just pick a sign: Optimizing deep multitask models with gradient sign dropout. *Advances in Neural Information Processing Systems*, 33:2039–2050, 2020. 6, 7

[6] Zitian Chen, Yikang Shen, Mingyu Ding, Zhenfang Chen, Hengshuang Zhao, Erik G Learned-Miller, and Chuang Gan. Mod-squad: Designing mixtures of experts as modular multi-task learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11828–11837, 2023. 1

[7] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 6

[8] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024. 2

[9] Heshan Devaka Fernando, Han Shen, Miao Liu, Subhajit Chaudhury, Keerthiram Murugesan, and Tianyi Chen. Mitigating gradient bias in multi-objective learning: A provably convergent approach. In *The Eleventh International Conference on Learning Representations*, 2023. 2

[10] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. 2021. 5

[11] Falk Heuer, Sven Mantowsky, Saqib Bukhari, and Georg Schneider. Multitask-centernet (mcn): Efficient and diverse multitask learning using an anchor free approach. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 997–1005, 2021. 1

[12] Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *International conference on machine learning*, pages 427–435. PMLR, 2013. 2

[13] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018. 6, 7

[14] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[15] Baijiong Lin, Feiyang Ye, Yu Zhang, and Ivor W Tsang. Reasonable effectiveness of random weighting: A litmus test for multi-task learning. *arXiv preprint arXiv:2111.10603*, 2021. 6, 7

[16] Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. *Advances in Neural Information Processing Systems*, 34:18878–18890, 2021. 1, 2, 5, 6, 7

[17] Bo Liu, Yihao Feng, Peter Stone, and Qiang Liu. Famo: Fast adaptive multitask optimization. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 5, 6, 7

[18] Liyang Liu, Yi Li, Zhanghui Kuang, J Xue, Yimin Chen, Wenming Yang, Qingmin Liao, and Wayne Zhang. Towards impartial multi-task learning. In *International Conference on Learning Representations*, 2021. 6, 7

[19] Shikun Liu, Edward Johns, and Andrew J Davison. End-to-end multi-task learning with attention. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1871–1880, 2019. 6, 7

[20] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. 4, 6

[21] Kevis-Kokitsi Maninis, Ilija Radosavovic, and Iasonas Kokkinos. Attentive single-tasking of multiple tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1851–1860, 2019. 6

[22] Aviv Navon, Aviv Shamsian, Idan Achituve, Haggai Maron, Kenji Kawaguchi, Gal Chechik, and Ethan Fetaya. Multi-task learning as a bargaining game. In *International Conference on Machine Learning*, pages 16428–16446. PMLR, 2022. 2, 5, 6, 7

[23] Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014. 7

[24] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018. 1, 2, 6, 7

[25] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part V 12*, pages 746–760. Springer, 2012. 6

[26] Derrick Xin, Behrooz Ghorbani, Justin Gilmer, Ankush Garg, and Orhan Firat. Do current multi-task optimization methods in deep learning even help? *Advances in neural information processing systems*, 35:13597–13609, 2022. 4

[27] Hanrong Ye and Dan Xu. Taskprompter: Spatial-channel multi-task prompting for dense scene understanding. In *The Eleventh International Conference on Learning Representations*, 2023. 1

[28] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*, 33:5824–5836, 2020. 1, 2, 6, 7

[29] Q Zhang, M Chen, A Bukharin, P He, Y Cheng, W Chen, and T Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. preprint (2023). *arXiv preprint arXiv:2303.10512*. 2

[30] Bo Zhao, Nima Dehmamy, Robin Walters, and Rose Yu. Symmetry teleportation for accelerated optimization. *Advances in neural information processing systems*, 35:16679–16690, 2022. 1, 2, 3

[31] Bo Zhao, Robert M Gower, Robin Walters, and Rose Yu. Improving convergence and generalization using parameter symmetries. *arXiv preprint arXiv:2305.13404*, 2023. 1, 2, 3