# AdaTask: A Task-aware Adaptive Learning Rate Approach to Multi-task Learning

**Enneng Yang** [1*†], **Junwei Pan** [2*], **Ximei Wang** [2], **Haibin Yu** [2], **Li Shen** [3]
**Xihua Chen** [2], **Lei Xiao** [2], **Jie Jiang** [2], **Guibing Guo** [1‡]

[1] Northeastern University, China [2] Tencent Inc, China [3] JD Explore Academy, China
ennengyang@stumail.neu.edu.cn, mathshenli@gmail.com, guogb@swc.neu.edu.cn
{jonaspan,messixmwang,nathanhbyu,tinychen,shawnxiao,zeus}@tencent.com

## Abstract

Multi-task learning (MTL) models have demonstrated impressive results in computer vision, natural language processing, and recommender systems. Even though many approaches have been proposed, how well these approaches balance different tasks on each parameter still remains unclear. In this paper, we propose to measure the task dominance degree of a parameter by the total updates of each task on this parameter. Specifically, we compute the total updates by the *exponentially decaying Average of the squared Updates* (AU) on a parameter from the corresponding task. Based on this novel metric, we observe that many parameters in existing MTL methods, especially those in the higher shared layers, are still *dominated* by one or several tasks. The dominance of AU is mainly due to the dominance of accumulative gradients from one or several tasks. Motivated by this, we propose a Task-wise Adaptive learning rate approach, AdaTask in short, to separate the *accumulative gradients* and hence the learning rate of each task for each parameter in adaptive learning rate approaches (e.g., AdaGrad, RMSProp, and Adam). Comprehensive experiments on computer vision and recommender system MTL datasets demonstrate that AdaTask significantly improves the performance of dominated tasks, resulting SOTA average task-wise performance. Analysis on both synthetic and real-world datasets shows AdaTask balance parameters in every shared layer well.

## Introduction

Multi-task learning (MTL) has emerged as a promising approach to jointly learning multiple tasks together by sharing structure across tasks in computer vision (Misra et al. 2016; Zamir et al. 2018; Liu, Johns, and Davison 2019), natural language processing (Dong et al. 2015; Radford et al. 2019), and recommender systems (Ma et al. 2018a; Pan et al. 2019; Tang et al. 2020). However, due to optimizing different tasks simultaneously, MTL is prone to be dominated by one or several tasks, leading to performance deterioration of other tasks (Tang et al. 2020).

While there has been a significant number of approaches to MTL, it is still unclear *to what extent these approaches*

---

*resolve the task dominance of parameter optimization during MTL training*. This is due to the fact that there are currently no metrics to quantify the degree of task dominance in MTL. To fill this gap, in this paper, we propose to calculate the total updates of a parameter from a specific task by the exponentially decaying Average of the squared parameter Updates (AU in short) from this task. The dominance of a task on a given parameter is then measured by the ratio of AU from this task against the overall AU from all tasks. Intuitively, for a parameter, if a task has a much larger AU ratio than other tasks, then this parameter can be regarded as being dominated by this task. Based on this metric, we observe that all existing MTL methods still suffer from task dominance on shared parameters, especially on those in higher shared layers.

We suspect that the dominance of accumulative gradients may be one of the main reasons of the task dominance issue in MTL. Modern adaptive learning rate optimizers commonly used in deep learning adaptively adjust the learning rate of each parameter during the model training process, which is usually better than the gradient descent optimizer with a global learning rate (see Appendix F for details). The adaptive learning rate is usually proportional to the inverse of the accumulative gradients, and both the learning rate and accumulative gradients are computed across all tasks for shared parameters in MTL. Therefore, for a given parameter, if a task provides much larger gradients than other tasks, then this task would dominate both overall accumulative gradients and the learning rate.

Motivated by this, we propose AdaTask to separate the accumulative gradients of each task for each shared parameter in MTL optimizers. By doing this, no task would dominate the overall accumulative gradients as well as the learning rate anymore since each task accumulates and computes its own ones. Almost all modern adaptive learning rate optimizer and its variants (Duchi, Hazan, and Singer 2011; Zeiler 2012; Kingma and Ba 2015; Reddi, Kale, and Kumar 2018; Zhuang et al. 2020; Zou et al. 2019; Chen et al. 2022a,b, 2021; Zou et al. 2018), such as AdaGrad, Weighted AdaGrad, AdaDelta, RMSProp, Adam, AMSGrad, etc., can utilize AdaTask to resolve the task dominance issue in MTL training. The contribution of this work is summarized as follows:

(1) We quantify the task dominance of a parameter by the task-wise exponentially decaying average of the squared updates(AU) in MTL and observe that existing MTL meth-

ods still suffer from task dominance.

(2) We propose a Task-aware Adaptive Learning Rate method (AdaTask) to separate the accumulative gradients of each task during optimization in the MTL models.

(3) Extensive experiments on the synthetic and three public datasets from the CV and recommendation demonstrate that AdaTask significantly improves the performance of the dominated task(s), while achieving SOTA average task-wise performance. Analysis shows that AdaTask balances parameters in all shared layers well.

## Related Work

**Multi-Task Learning Architecture.** Shared Bottom (Caruana 1997) is the most basic and commonly used MTL structure. With hard parameters shared across tasks, it suffers from negative transfer. Cross-stitch network (Misra et al. 2016) and Sluice network (Ruder et al. 2019) were proposed to learn weights of linear combinations to fuse representations from different tasks selectively. However, the weights for different tasks are static for all samples. MOE (Shazeer et al. 2017) first proposed to share some experts at the bottom and combine experts through a gating network. MMOE (Ma et al. 2018a) extends MOE to utilize different gates for each task to obtain different fusing weights in MTL. Similarly, MRAN (Zhao et al. 2019) applied multi-head self-attention to learn different representation subspaces at different feature sets. PLE (Tang et al. 2020) develops both shared experts and task-specific experts, together with a progressive routing mechanism to further improve learning efficiency. In addition, there have been some studies of adaptive selection of layers or neurons instead of experts. AdaShare (Sun et al. 2020) selects some layers for each task from multiple shared layers. TAAN (Liu et al. 2020) adopts task adaptive activation network to enable flexible and low-cost MTL. MTAN (Liu, Johns, and Davison 2019) uses the attention mechanism to extract the feature representation corresponding to the task from the shared backbone. In contrast to the above, which focuses on the shared module design, some works explicitly accomplish knowledge transfer between tasks at the output layer. ESSM (Ma et al. 2018b) model the sequential pattern impression → click → conversion. CrossDistil (Yang et al. 2022) proposes to enhance knowledge transfer between tasks through distillation.

**Multi-Task and Multi-Objective Optimization.** There are also many MTL works to revolve negative transfer by either balancing the gradients of all tasks (Kendall, Gal, and Cipolla 2018; Liu, Johns, and Davison 2019; Liu, Liang, and Gitter 2019; Chen et al. 2018) or avoiding gradient conflicts (Yu et al. 2020; Chen et al. 2020; Liu et al. 2021). UW (Kendall, Gal, and Cipolla 2018) calculates the weights based on homoscedastic uncertainty, while DWA (Liu, Johns, and Davison 2019) and LBTW (Liu, Liang, and Gitter 2019) were based on the rate of change of loss for each task. Grad-Norm (Chen et al. 2018) proposed to add a regularization term to the loss function to balance the training rates of all tasks. (Yu et al. 2020) proposes gradient surgery to avoid interference between task gradients. GradDrop (Chen et al. 2020) samples gradients based on their level of consistency.

CAGrad (Liu et al. 2021) seeks the gradient update direction by maximizing the task with the least loss reduction. In addition, some works formulate multi-task learning as a multi-objective optimization problem and seek discrete or continuous Pareto optimal solutions (Sener and Koltun 2018; Ma, Du, and Matusik 2020; Xie et al. 2021). Finally, some works in MTL attempt to balance the interference of auxiliary tasks with the main task (Lin et al. 2019a; Navon et al. 2021; He et al. 2022).

## Rethinking Task Dominance in MTL

In this section, we want to answer the following questions: (RQ1) How can we quantify the task dominance of a parameter in MTL models? (RQ2) To what extent do exist MTL approaches tackle the task dominance issue? (RQ3) How does task dominance impact the training of MTL models?

### Synthetic Dataset Setting

To answer the above research questions, we generate a synthetic MTL dataset with two regression tasks, following Grad-Norm (Chen et al. 2018). Specifically, the ground truth for task $k \in \{A, B\}$ is obtained as follows:

$$\mathbf{y}^k(\mathbf{x}) = \mathbf{W}_1^k \mathbf{x} + \mathbf{W}_2^k \mathbf{x}^2 + \mathbf{W}_3 \mathbf{x}^3 + \epsilon, \qquad (1)$$

where $\mathbf{W}_1^A$ and $\mathbf{W}_2^A$ are constant matrices with elements sampled from $\mathcal{N}(1, 1)$, $\mathbf{W}_1^B$, $\mathbf{W}_2^B$ and $\mathbf{W}_3$ are constant matrices with elements sampled from $\mathcal{N}(10, 10)$, and $\epsilon$ is a Gaussian noise sampled from $\mathcal{N}(0, 0.1)$. Note that $\mathbf{W}_1^B$, $\mathbf{W}_2^B$ is roughly ten times as much as $\mathbf{W}_1^A$ and $\mathbf{W}_2^A$, making the dataset dominated by task $B$. Moreover, $\mathbf{x} \in \mathbb{R}^{250}$ and $\mathbf{y}^k(\mathbf{x}) \in \mathbb{R}^{100}$. The MTL model is a 4-layer fully-connected neural network with ELU (Clevert, Unterthiner, and Hochreiter 2016) as the activation function. We use the RMSProp optimizer for parameter optimization.

### (RQ1) How can we quantify the task dominance of parameters in MTL model?

We quantify the task dominance of a parameter by the accumulative updates from the corresponding task divided by the overall updates from all tasks for this parameter. Here we choose the exponentially decaying Average of the squared Updates (AU) to measure the accumulative updates, another choice may also be feasible. Note that we mainly focus on the shared parameters rather than the task-specific parameters since the later ones are designed to be only updated and therefore dominated by the corresponding task.

Given a shared parameter $\theta_i$, the AU at step $t$ from task $k$ is defined as:

$$\mathrm{AU}(i, t, k) := G\left[\Delta\theta^2\right]_{t,i}^k = \gamma G\left[\Delta\theta^2\right]_{t-1,i}^k + (1-\gamma)[\Delta\theta_{t,i}^k]^2 \quad (2)$$

where $\Delta\theta_{t,i}^k$ denotes the update of parameter $\theta_i$ at setp $t$ from task $k$. In RMSProp, $\Delta\theta_{t,i}^k = -\frac{\eta}{\sqrt{G_{t,i}+\epsilon}} g_{t,i}^k$.

The dominance of task $k$ on a parameter $i$ at step $t$ can therefore be defined as the ratio of its AU over the total AU of all tasks, i.e.,

$$\mathrm{rAU}(i, t, k) := \frac{\mathrm{AU}(i, t, k)}{\sum_{k'} \mathrm{AU}(i, t, k')} = \frac{G\left[\Delta\theta^2\right]_{t,i}^k}{\sum_{k'} G\left[\Delta\theta^2\right]_{t,i}^{k'}} \quad (3)$$
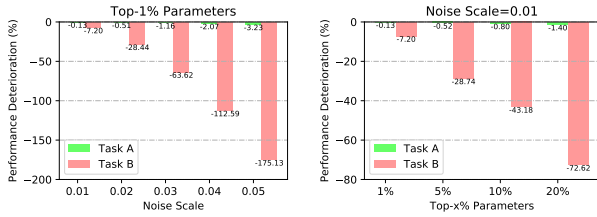
Figure 1: Illustration of rAU$(i, T, B)$ as a metric to measure the task dominance for shared parameters in MTL.

We denote this AU ratio w.r.t. task $k$, or rAU$(i, t, k)$ in short.

In order to verify whether AU indeed measures the dominance of a specific task, we conduct the following analysis. On the one hand, we pick up the Top-1% parameters w.r.t. rAU$(i, T, B)$, i.e., the Top-1% parameters that are treated as being dominated by task $B$ at the last step $T$ according to our definition. Random Gaussian noises with mean 0 and various variances of $\{0.01, 0.02, 0.03, 0.04, 0.05\}$ are added to these parameters; on the other hand, we select the Top-$x$% ($x \in \{1, 5, 10, 20\}$) w.r.t. rAU$(i, T, B)$, and add Gaussian noise with a mean of 0 and a variance of 0.01. We then evaluate the performance of these models on task $A$ and task $B$, and check their performance deterioration, compared with the original model. If a parameter is dominated by task $B$, then adding Gaussian noise may lead to a more significant performance deterioration on task $B$ than on task $A$.

We plot the relative root mean square error (RMSE) increment of these models in Fig. 1. We observe that (1) With Gaussian noise added to the Top-1% parameters w.r.t. rAU$(i, T, B)$, the loss increases much more on task $B$ than task $A$. Specifically, the loss is increased by 175.13% on task $B$, but only 3.23% on task $A$, with a $\mathcal{N}(0, 0.05)$ Gaussian noise. (2) By adding a $\mathcal{N}(0, 0.01)$ Gaussian noise to the the Top-1% to the Top-20% parameters w.r.t. rAU$(i, T, B)$, the performance deteriorates much more on task $B$ than on task $A$. With more parameters with high rAU$(i, T, B)$ polluted by the noise, the performance deteriorates more. For example, the loss is increased by 7.20% when only the Top-1% parameters w.r.t. rAU$(i, T, B)$ is polluted, while it's increased by 72.62% with the Top-20% polluted. This concludes that the rAU$(i, t, k)$ is a reasonable metric to measure the dominance of a task on a shared parameter: a parameter with a high rAU$(i, t, k)$ for one task becomes much more important for that task, while less important for the remaining tasks.

### (RQ2) To what extent do exist MTL approaches tackle the task dominance issue?

In order to quantify to what extent existing MTL approaches tackle the task dominance issue (taking the dominance of task $B$ as an example here), we calculate the percentage of shared parameters based on their rAU$(i, T, B)$ using the following pre-defined thresholds: $\{0\%, 20\%, 40\%, 60\%, 80\%, 100\%\}$. If rAU$(i, T, B) \in (80\%, 100\%]$, we regard the parameter $i$ as being dominated by task $B$. In addition, if both rAU$(i, T, A)$ and rAU$(i, T, B)$ are in $(40\%, 60\%]$, we can say that parameter $i$ is roughly balanced. Without loss of generality, we pick

up Shared Bottom as the backbone MTL model architecture. Please note other MTL structures such as MMOE, PLE, or MTAN are also feasible. Following (Chen et al. 2018), we then pick up UW, CAGrad, and GradNorm as representative MTL balancing methods. We name the vanilla Shared Bottom method as EqualWeight since it corresponds to weighting all tasks equally.

The percentage of shared parameters of each model w.r.t. rAU$(i, T, B)$ in different buckets is shown in Fig. 2. We observe that: (1) Almost all shared parameters in Equal-Weight are dominated by task $B$, in the sense that the red area (rAU$(i, t, B) \in (80\%, 100\%]$) is always close to 100%. This result is because these two methods either use the original gradient directly or use the projected gradient direction, and they have a serious gradient dominance problem when the magnitudes of the gradients of the two tasks are different. (2) In UW and CAGrad, some of the parameters are dominated by task $B$. This can be explained as both methods adjust the magnitude of the gradient by weighting the loss or weighting the gradient, thus alleviating the problem of accumulative gradient dominance. However, their weights are designed in terms of learning speed or gradient direction optimization and are not directly related to the magnitude of the gradient, thus only slightly mitigating task dominance. (3) In contrast, around 78% percent of shared parameters in GradNorm are not dominated by any task. This is because the weights of the loss in GradNorm are designed to balance the overall gradient magnitude of the task.

We further split the shared parameters into layers, and calculate the rAU$(i, t, B)$ for parameters from each layer, respectively. The detailed result is shown in Fig. 6 of Appendix D. In EqualWeight, PCGrad and UW, parameters from all layers are dominated by task $B$ in the sense parameters from all layers are with rAU$(i, t, B) \in (80\%, 100\%]$. CAGrad and GradNorm are both balanced at some level. Specifically, In GradNorm, only the first layer is balanced well in the sense that around 90% of the parameters with rAU$(i, t, B) \in (40\%, 60\%]$. However, this percentage drops to around 62%, 61%, and 40%, respectively for the 2nd, 3rd, and 4th layers. Such failure in balancing shared parameters from higher layers may be due to that the lower layers learn more general representations while the higher layers learn task-specific representations (Yosinski et al. 2014).

### (RQ3) How does task dominance impact the training of MTL models?

We have observed that some parameters, especially those in higher layers, are still dominated by task $B$. We then wonder what would happen due to such dominance? Or more specifically, what will happen if the Accumulative Updates (AU) are dominated?

In modern optimizers, the accumulative gradients are used to adjust the learning rate of each parameter during model training. Then dominance of accumulative gradients may lead to a dominance of learning rate(LR) in MTL. For example, in a MTL setting with only two tasks $A$ and task $B$, the accumulative gradients of task $A$ is small, and it would have a large LR if it's trained on its own. However, if task $B$ has larger accumulative gradients, then the overall accumulative
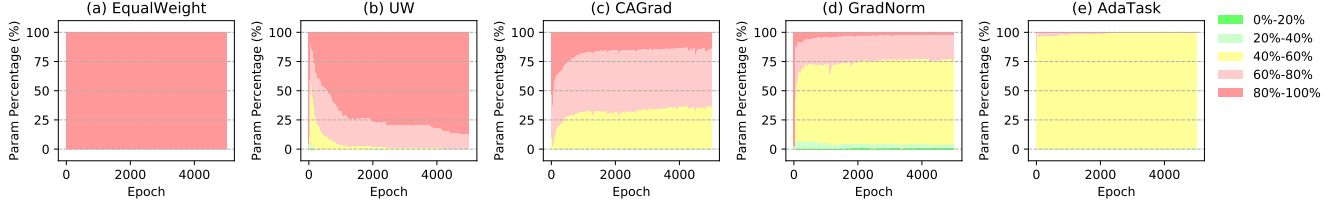
Figure 2: rAU$(i, T, B)$ of all shared parameters on the synthetic dataset for five MTL models: (a) EqualWeight (PCGrad is close to EqualWeight, it was removed due to page limitations.), (b) UW, (c) CAGrad, (d) GradNorm and (e) our AdaTask. The green area denotes the percentage of parameters dominated by task $A$, the red area denotes the percentage of parameters dominated by task $B$, and the yellow area denotes the percentage of balanced parameters.

gradients in MTL would be also large, leading to a small overall learning rate. With a much smaller learning rate in MTL than in training $A$ alone, task $A$ may be optimized much slower, leading to performance deterioration.

To verify this, we randomly select multiple parameters from each layer. We then calculate the separate and whole accumulative gradients, as well as the overall LR when training two tasks simultaneously in MTL, as well as the separate LR for each task if it's trained independently. The results are shown in Fig. 5(a) of Appendix D. Each line represents the average learning rate of multiple parameters. It is obvious that the whole LR is dominated by task $B$ in the sense it's very close to the separate LR of task $B$.

## Our Proposed Method: AdaTask

It is observed in the above section existing MTL optimization methods still suffer from task dominance, and such dominance leads to the dominance of learning rate in optimizers. In order to tackle such dominance, in this section, we propose a Task-wise Adaptive learning rate approach, to separate the accumulative gradients of each task for each shared parameter. We name this approach as *AdaTask*.

**Preliminaries** Adaptive learning rate optimizers adjust the learning rate of each parameter based on some statics $G_{t,i}$ of history gradients w.r.t. this parameter: $\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,i}+\epsilon}} g_{t,i}$, where $\eta$ is a hyper-parameter representing the initial global learning rate, $\epsilon$ is a small number used to avoid a zero denominator. In AdaGrad, $G_{t,i}$ is defined as the sum of all history gradient squares, i.e., $G_{t,i} = G_{t-1,i} + (g_{t,i})^2$, where $g_{t,i}$ denotes the gradient of this parameter at step $t$. In AdaDelta, Adam or RMSProp, it's defined as the exponential moving average, i.e., $G_{t,i} = \gamma G_{t-1,i} + (1-\gamma)(g_{t,i})^2$, where $\gamma$ is a decaying coefficient.

When optimizing shared parameters in MTL models, these statics are aggregated over all tasks, i.e.:

$$G_{t,i} = G_{t-1,i} + (g_{t,i})^2 \text{ AdaGrad in MTL}$$
$$G_{t,i} = \gamma G_{t-1,i} + (1-\gamma)(g_{t,i})^2 \text{ RMSProp in MTL} \quad (4)$$

where $g_{t,i} = g_{t-1,i} + \sum (g_{t,i}^k)^2$ and $k \in \{1, 2, \ldots, K\}$ represents the $k$-th task.

**AdaGrad, RMSProp with AdaTask.** As discussed in the above section, the accumulative gradients across tasks in Eq. 4 tend to be dominated by some tasks, which in turn

leads to the dominance of the adaptive learning rate. We propose to separate the accumulative gradients of different tasks for each shared parameter. Formally, we maintain task-wise accumulative gradients $G_{t,i}^k$ for parameter $i$ and task $k$ at step $t$. Each such variable only aggregates gradients from the corresponding task, i.e.,

$$G_{t,i}^k = f(G_{t-1,i}^k, (g_{t,i}^k)^2) \quad (5)$$

where $f(\cdot)$ denotes the corresponding aggregation function in each method, i.e., average in AdaGrad, and exponentially decaying average in RMSProp. During MTL model training, when optimizing a sample from one task, only the accumulative gradients w.r.t. that task is used to calculate the learning rate. Therefore, AdaGrad with AdaTask and RMSProp with AdaTask in MTL can then be formulated as:

$$\theta_{t+1,i} = \theta_{t,i} - \sum_k \frac{\eta}{\sqrt{G_{t,i}^k}+\epsilon} g_{t,i}^k$$

$$G_{t,i}^k = G_{t-1,i}^k + (g_{t,i}^k)^2 \text{ AdaGrad with AdaTask in MTL}$$

$$G_{t,i}^k = \gamma G_{t-1,i}^k + (1-\gamma)(g_{t,i}^k)^2 \text{ RMSProp with AdaTask in MTL}$$

**Adam with AdaTask.** The Adam (Kingma and Ba 2015) includes the decayed average of both past gradient squares $G_{t,i}$ as AdaDetla (Zeiler 2012) (or RMSProp) and past gradients $m_{t,i}$ as momentum (Qian 1999). Adam's update rule in MTL is as follows:

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,i}}+\epsilon} m_{t,i}$$
$$g_{t,i} = g_{t-1,i} + \sum (g_{t,i}^k)^2$$
$$G_{t,i} = \gamma_1 G_{t-1,i} + (1-\gamma_1)(g_{t,i})^2 \quad (6)$$
$$m_{t,i} = \gamma_2 m_{t-1,i} + (1-\gamma_2)(g_{t,i})$$

When applying AdaTask to Adam, if we only separate the $G_{t,i}$ per task, the accumulation of gradients in $m_{t,i}$ can still be dominated by one or several tasks due to the different magnitude of $g_{t,i}^k$. In order to prevent task dominance in $m_{t,i}$, we propose to separate it for each task. Adam with AdaTask can be formulated as:

$$\theta_{t+1,i}^k = \theta_{t,i} - \sum_k \frac{\eta}{\sqrt{G_{t,i}^k}+\epsilon} m_{t,i}^k$$
$$G_{t,i}^k = \gamma_1 G_{t-1,i}^k + (1-\gamma_1)(g_{t,i}^k)^2 \quad (7)$$
$$m_{t,i}^k = \gamma_2 m_{t-1,i}^k + (1-\gamma_2)(g_{t,i}^k)$$

**Algorithm 1:** AdaGrad, RMSProp and Adam in MTL

1: **Require** $\gamma, \gamma_1, \gamma_2$: Exponential decay factors
2: **Require** $\eta$: Initial learning rate
3: **Require** $\epsilon$: A smoothing factor to avoid division by zero
4: **Require** $f^k(\theta)$: Objective function of task $k$, $\forall k \in \{1, 2, \ldots, K\}$
5: **Require** Initialize: $G_0 = 0, m_0 = 0$
6: **For** step $t = 1 : T$ **do**
7:     $g_{t,i} = 0$
8:     **For** task $k = 1 : K$ **do**
9:        $g_{t,i}^k = \nabla f_t^k(\theta_{t,i})$
10:        $g_{t,i} = g_{t,i} + g_{t,i}^k$
11:     **End For**
12:     **If** AdaGrad
13:        $G_{t,i} = G_{t-1,i} + {g_{t,i}}^2$
14:        $\Delta\theta_{t,i} = \frac{\eta}{\sqrt{G_{t,i}+\epsilon}} g_{t,i}$
15:     **Else If** RMSProp
16:        $G_{t,i} = \gamma G_{t-1,i} + (1-\gamma){g_{t,i}}^2$
17:        $\Delta\theta_{t,i} = \frac{\eta}{\sqrt{G_{t,i}+\epsilon}} g_{t,i}$
18:     **Else If** Adam
19:        $G_{t,i} = \gamma_1 G_{t-1,i} + (1-\gamma_1){g_{t,i}}^2$
20:        $m_{t,i} = \gamma_2 m_{t-1,i} + (1-\gamma_2)g_{t,i}$
21:        $\hat{G}_{t,i} = G_{t,i} / (1-\gamma_1^t)$
22:        $\hat{m}_{t,i} = m_{t,i} / (1-\gamma_2^t)$
23:        $\Delta\theta_{t,i} = \frac{\eta}{\sqrt{\hat{G}_{t,i}+\epsilon}} \hat{m}_{t,i}$
24:     **End If**
25:     $\theta_{t+1,i} = \theta_{t,i} - \Delta\theta_{t,i}$
26: **End For**

---

**Algorithm 2:** AdaGrad, RMSProp and Adam **with AdaTask** in MTL

1: **Require** $\gamma, \gamma_1, \gamma_2$: Exponential decay factors
2: **Require** $\eta$: Initial learning rate
3: **Require** $\epsilon$: A smoothing factor to avoid division by zero
4: **Require** $f^k(\theta)$: Objective function of task $k$, $\forall k \in \{1, 2, \ldots, K\}$
5: **Require** Initialize: $G_0^k = 0, m_0^k = 0, \forall k \in \{1, 2, \ldots, K\}$
6: **For** step $t = 1 : T$ **do**
7:     **If** AdaGrad with AdaTask
8:        **For** task $k = 1 : K$ **do**
9:           $g_{t,i}^k = \nabla f_t^k(\theta_{t,i})$
10:           $G_{t,i}^k = G_{t-1,i}^k + {g_{t,i}^k}^2$
11:           $\Delta\theta_{t,i}^k = \frac{\eta}{\sqrt{G_{t,i}^k+\epsilon}} g_{t,i}^k$
12:        **End For**
13:     **Else If** RMSProp with AdaTask
14:        **For** task $k = 1 : K$ **do**
15:           $g_{t,i}^k = \nabla f_t^k(\theta_{t,i})$
16:           $G_{t,i}^k = \gamma G_{t-1,i}^k + (1-\gamma){g_{t,i}^k}^2$
17:           $\Delta\theta_{t,i}^k = \frac{\eta}{\sqrt{G_{t,i}^k+\epsilon}} g_{t,i}^k$
18:        **End For**
19:     **Else If** Adam with AdaTask
20:        **For** task $k = 1 : K$ **do**
21:           $g_{t,i}^k = \nabla f_t^k(\theta_{t,i})$
22:           $G_{t,i}^k = \gamma_1 G_{t-1,i}^k + (1-\gamma_1){g_{t,i}^k}^2$
23:           $m_{t,i}^k = \gamma_2 m_{t-1,i}^k + (1-\gamma_2)g_{t,i}^k$
24:           $\hat{G}_{t,i}^k = G_{t,i}^k / (1-\gamma_1^t)$
25:           $\hat{m}_{t,i}^k = m_{t,i}^k / (1-\gamma_2^t)$
26:           $\Delta\theta_{t,i}^k = \frac{\eta}{\sqrt{\hat{G}_{t,i}^k+\epsilon}} \hat{m}_{t,i}^k$
27:        **End For**
28:     **End If**
29:     $\theta_{t+1,i} = \theta_{t,i} - \sum_k \Delta\theta_{t,i}^k$
30: **End For**

---

In Alg. 1 / Alg. 2, we shown the AdaGrad, RMSProp, and Adam without / with AdaTask in MTL, respectively.

**Efficient AdaTask.** In order to separate the Accumulative Gradients in AdaTask, $K-1$ times more intermediate variables are needed during optimization. This is negligible when there are many shared parameters. We can reduce the number of variables to store accumulative gradients by grouping parameters by some dimensions, e.g., layers or modules. Without loss of generality, we discuss grouping accumulative gradients by layers here. That is, all parameters in each layer ($l \in \{1, \ldots, L\}$) share the same accumulative gradients per task and formalize it as follows:

$$C_{t,l}^k = \gamma C_{t-1,l}^k + (1-\gamma)h\left(\{(g_{t,i}^k)^2, \forall i \in l(\theta)\}\right), \quad (8)$$

where $h(\cdot)$ represents an aggregate function, such as mean or sum operation, which is used to aggregate the gradient of all parameters at the $l$-th layer, $\forall i \in l(\theta)$ represents all parameters in the $l$-th shared layer. Therefore, $C_{t,l}^k$ can be regarded as the overall contribution of task $k$ to the shared layer $l$ at step $t$. The accumulative gradient of task $k$ w.r.t parameter $i$ is expressed as follows:

$$G_{t,i}^k = \frac{C_{t,l(i)}^k}{\sum_k C_{t,l(i)}^k} G_{t,i}, G_{t,i} = \gamma G_{t-1,i} + (1-\gamma)(g_{t,i}^k)^2, \quad (9)$$

where $l(i)$ which layer parameter $i$ belongs to. Due to limited space, we provide a detailed comparison of AdaTask with the efficient version (LAdaTask) in Appendix A.

## Experiments

In this section, we experimentally verify the effectiveness of the proposed AdaTask. We conduct experiments on a synthetic dataset, and several real-world datasets: the CityScapes dataset in computer vision, and the TikTok and WeChat datasets in recommender systems. For the baseline methods, we choose two loss weighting approaches GradNorm (Chen et al. 2018), UW (Kendall, Gal, and Cipolla 2018), and two gradient conflict-avoiding approaches PCGrad (Yu et al. 2020) and CAGrad (Liu et al. 2021). More details of our experiment, including baselines, datasets, and implementation details, are provided in Appendix B.

We adopt the commonly used *SharedBottom* architecture (mainly based on fully connected layers) and *RMSProp* optimizer in the synthetic and recommendation datasets. In the computer vision dataset, we choose the commonly used *MTAN* architecture (including complex operations such as convolution and attention mechanism, etc.) and the *Adam* optimizer. Each experiment is repeated over 3 random seeds and the mean is reported. Validation on different network architectures (such as MMoE and SegNet) is presented in Appendix C. In addition, we provide a comparison of the

| Method | Task A | | Task B | | $\Delta p \uparrow$ |
| | Abs Err↓ | Rel Err↓ | mIoU↑ | Pix Acc↑ | |
|---|---|---|---|---|---|
| EqualWeight | 0.0152 | 47.00 | 75.01 | 93.40 | 0.00% |
| GridSearch | 0.0132 | 41.88 | 74.13 | 93.16 | 5.65% |
| MGDA | 0.0166 | 33.50 | 69.42 | 91.22 | 2.40% |
| GradNorm | 0.0145 | 44.16 | 75.18 | 93.38 | 2.71% |
| UW | 0.0134 | 36.66 | 74.42 | 93.05 | 8.05% |
| PCGrad $^{\dagger}$ | 0.0150 | 42.52 | **75.46** | **93.49** | 2.88% |
| CAGrad $^{\dagger}$ | 0.0141 | 38.40 | 75.30 | 93.48 | 6.50% |
| AdaTask | 0.0128 | 36.85 | 75.02 | 93.40 | 9.34% |
| + PCGrad | **0.0127** | 34.46 | 74.76 | 93.25 | 10.52% |
| + CAGrad | 0.0133 | **31.50** | 73.17 | 92.88 | **10.59%** |

Table 1: Performance evaluation on the CityScapes dataset. The $^{\dagger}$ symbol represents the gradient direction modification methods, and they are orthogonal to our work. Task $A$ is the dominated task, and Task $B$ is the dominant task.

training times in Appendix E.

## Performance Evaluation

**CityScapes Dataset Results.** For the semantic segmentation and depth estimation tasks on CityScapes datasets, we follow MTAN (Liu, Johns, and Davison 2019) and CA-Grad (Liu et al. 2021) to use Absolute Error (Abs Err) and Relative Error (Rel Err) to evaluate the performance of the Depth Estimation task (Task $A$) and use Mean Intersection over Union (mIoU) and Pixel Accuracy (Pix Acc) to evaluate the performance of Semantic Segmentation task (Task $B$). Similar to CAGrad, we also compute the average per-task performance improvement $\Delta p$ of an MTL method $m$ with respect to the EqualWeight baseline $b$: $\Delta p = \frac{1}{4} \sum_p (-1)^{I_p} (M_{m,p} - M_{b,p}) / M_{b,p}$, where $I_p = 0$ if a higher value is better for a criterion $M_p$ on performance $p \in \{$Abs Err, Rel Err, mIoU and Pix Acc$\}$ and 1 otherwise.

As shown in Table 1, GridSearch delivers higher performance than Equalweight, but manual weighting may require efforts to find the proper weights. AdaTask achieved a significant improvement in two metrics of task $A$ (the dominated task in this dataset). Specifically, Abs Err is improved by 15.78% compared to EqualWeight, and Rel Err is improved by 21.59% compared to EqualWeight. Meanwhile, AdaTask achieves a result close to EqualWeight on task $B$. In addition, in terms of the overall performance, AdaTask has achieved an improvement of 1.29% and 2.84% respectively compared with the best baseline (UW) and the second-best baseline (CAGrad), respectively. Notably, AdaTask is orthogonal to the methods of modifying gradient direction, such as CA-Grad and PCGrad, so it can be further combined with these methods to improve performance. When combined with PC-Grad and CAGrad (The implementation details are in Alg. 4 and Alg. 3 of Appendix G), the performance can be further improved by 1.18% and 1.25%, respectively. Finally, by comparing MGDA (which seeks Pareto optimal solutions), we can find that the solutions of GradNorm, UW, PCGrad, CA-Grad, AdaTask, and MGDA do not dominate each other, that is, none of them performs better than the other in all metrics.

**Synthetic Dataset Results.** For the regression tasks on the Synthetic dataset, we follow GradNorm (Chen et al. 2018) to use Task-Normalized RMSE (Root Mean Square Error) to evaluate the performance of all methods. As shown in

| Method | Task A ↓ | Task B ↓ | Average ↓ |
|---|---|---|---|
| EqualWeight | 0.2172 | 0.0113 | 0.1142 |
| GradNorm | 0.1134 | 0.0226 | 0.0680 |
| UW | 0.2032 | **0.0100** | 0.1066 |
| PCGrad | 0.2221 | 0.0155 | 0.1188 |
| CAGrad | 0.1495 | 0.0140 | 0.0817 |
| AdaTask | **0.0704** | 0.0417 | **0.0560** |

Table 2: Performance on the synthetic dataset. Task $A$ is the dominated task, and task $B$ is the dominant task.

Table 2, we observe that the proposed AdaTask achieves a huge RMSE improvement in task $A$, with an absolute improvement of 0.14 and 0.04 compared with the EqualWeight and the best baseline (GradNorm), respectively. Such a huge improvement is due to the fact that task $A$ is dominated dramatically but task $B$ and our proposed method resolve such dominance explicitly. The RMSE absolute reduction of task $B$ is 0.03 compared with EqualWeight. Task $B$ showed some decline since it already dominates baseline models and is more inclined by them. For the average results of the two tasks, AdaTask achieved the best performance.

**TikTok Dataset Results.** For the MTL classification tasks on the TikTok dataset, we evaluate Area under the ROC curve (AUC) (Fawcett 2006) of each task, as well as the average and weighed AUC, with weights specified by the provider. As shown in Table 3, similar to the observation on the synthetic dataset, AdaTask gets a significant AUC lift on the dominated task ($A$ in this case). The improvement on task $A$ is 1.51% and 0.81%, compared with EqualWeight and the best baseline (CAGrad). While on task $B$, there is a slight performance drop of 0.51%, compared with EqualWeight. Such performance drop may be due to a strong preference for task $B$ in these baseline methods. Compared with EqualWeight, AdaTask improves the Average AUC and Weighted AUC by 0.61% and 0.99%, respectively. While against the best baseline model, AdaTask improves the Average AUC and Weighted AUC by 0.59% (against GradNorm), 0.6% (against PCGrad), and 0.43% (against CAGrad).

**WeChat Dataset Results.** For the WeChat dataset, similar to Tiktok, we used average AUC and weighted AUC to measure the overall performance. As shown in Table 4, we observed a similar phenomenon with the above three datasets, that is, AdaTask improved significantly in the dominated tasks (Task $A$, $B$, and $C$), while its performance decreased slightly in the dominate task (Task $D$). Compared to the best baseline method, i.e., CAGrad, AdaTask slightly outperformed it by 0.22% on average AUC. However, w.r.t. the weighted AUC, CAGrad outperforms AdaTask by 0.28%. This is because AdaTask usually performs worse on the dominant task, but in the WeChat dataset the dominant one, i.e., Task D, gets a high weight (The weights for task $A$, $B$, $C$, and $D$ is 1:2:3:4, please refer to Appendix B for details).

## Overall Results

**Multi-task overall performance perspective.** Experiments on four datasets demonstrate that our proposed AdaTask model achieves state-of-the-art performance w.r.t. Average and Weighted metrics. Such improvement mainly comes from

| Method | Task A ↑ | Task B ↑ | Average AUC ↑ | Weighted AUC ↑ |
|---|---|---|---|---|
| EqualWeight | 0.9173 | <u>0.7437</u> | 0.8305 (+0.00%) | 0.8652 (+0.00%) |
| GradNorm (Chen et al. 2018) | 0.9182 | 0.7431 | 0.8307 (+0.02%) | 0.8657 (+0.05%) |
| PCGrad (Yu et al. 2020) | 0.9199 | 0.7413 | 0.8306 (+0.01%) | 0.8663 (+0.12%) |
| UW (Kendall, Gal, and Cipolla 2018) | <u>0.9237</u> | 0.7351 | 0.8294 (-0.13%) | 0.8671 (+0.21%) |
| CAGrad (Liu et al. 2021) | 0.9185 | **0.7455** | <u>0.8320</u> (+0.18%) | <u>0.8666</u> (+0.16%) |
| AdaTask | **0.9312** | 0.7399 | **0.8356** (+0.61%) | **0.8738** (+0.99%) |

Table 3: Performance evaluation on the TikTok. Task $A$ is the dominated task, and task $B$ is the dominant task.

| Method | Task A↑ | Task B↑ | Task C↑ | Task D↑ | Average AUC↑ | Weighted AUC↑ |
|---|---|---|---|---|---|---|
| EqualWeight | 0.9080 | 0.8912 | 0.8881 | 0.9514 | 0.9097(+0.00%) | 0.9160(+0.00%) |
| GradNorm (Chen et al. 2018) | 0.9408 | 0.9072 | 0.8840 | 0.9462 | 0.9196(+1.08%) | 0.9192(+0.34%) |
| UW (Kendall, Gal, and Cipolla 2018) | **0.9645** | 0.9187 | 0.8654 | 0.9451 | <u>0.9234</u>(+1.51%) | 0.9179(+0.19%) |
| PCGrad (Yu et al. 2020) | 0.9182 | 0.8987 | <u>0.8942</u> | <u>0.9520</u> | 0.9158(+0.67%) | 0.9206(+0.50%) |
| CAGrad (Liu et al. 2021) | 0.9462 | <u>0.9313</u> | **0.9040** | **0.9538** | 0.9338(+2.65%) | **0.9336**(+1.91%) |
| AdaTask | <u>0.9643</u> | **0.9391** | 0.8920 | 0.9480 | **0.9359**(+2.87%) | <u>0.9311</u>(+1.63%) |

Table 4: Performance evaluation on the WeChat. Task $A$, $B$, and $C$ are the dominated tasks, and task $D$ is the dominant task.

the dominated task, which is task $A$ in Synthetic, CityScapes and TikTok cases, task $A$, $B$, and $C$ in WeChat case. Compared with baselines, our results show a certain decrease in the dominant tasks. This is because baseline methods are heavily skewed toward the dominant task, resulting in better performance for the dominant task. However, our method restricts the shared parameters to be skewed to the dominant task to a certain extent, so the performance of the dominant task will be slightly degraded, but we achieve the overall optimal performance.

**Multi-objective Pareto solution perspective.** In addition to overall performance, we are also interested in reviewing these results from a MOO perspective. According to the definition of the Pareto stationary solution, solution $\theta^*$ is said to belong to the Pareto stationary solution set if there is no solution $\theta$ that outperforms $\theta^*$ in all tasks. So we can consider that the solution found by AdaTask also belongs to Pareto stationary solution because there is no method (such as MGDA, GradNorm, UW, PCGrad, or CAGrad) that outperforms AdaTask on all task metrics. We provide more detailed comparison and discussion with the MOO approach in Appendix C.

**Study on Task Dominance**

**Task Dominance on AU.** In this part, we study whether the proposed AdaTask approach indeed resolves the task dominance w.r.t. AU. We calculate the percentage of shared parameters from all layers w.r.t. their rAU$(i, T, B)$ using some pre-defined thresholds: $\{0\%, 20\%, 40\%, 60\%, 80\%, 100\%\}$. The result of AdaTask is shown in Fig. 2(e) for the synthetic dataset. Very few parameters in AdaTask are dominated by task $B$ (yellow area in (e)). We then further split the parameters into four layers in the synthetic dataset. The results are shown in Fig. 6. The proposed AdaTask achieves a decent balance of AU: after the first several epochs, more than 98% of shared parameters in all four layers have rAU$(i, T, B)$ in (40%, 60%]. We also show that task dominance is a real problem in the real-world dataset (CityScapes). As shown in Fig. 4 of Appendix D, we calculate the AU of EqualWeight(using MTAN) and GradNorm(using MTAN) on the CityScapes dataset and do observe that 99% and 95% of shared parameters in these two methods are dominated by task $B$. After

applying AdaTask(on MTAN), only 7% of parameters are slightly dominated by task $B$, while 93% of shared parameters are balanced well. It shows that AdaTask solves the task dominance issue well.

**Task Dominance on Learning Rate.** In this part, we study the task dominance of the learning rate in AdaTask. The learning rate of the baseline methods dominates, which we analyze in Task-dominance analysis section and our AdaTask is shown in Fig. 5(b) of Appendix D. There is no such learning rate dominance problem in our proposed AdaTask since we separate learning rates between tasks for each parameter. Therefore a parameter can be optimized toward each task's objective based on its own learning rate, without any intervention from other tasks.

## Conclusion and Future Works

In this work, we quantify the dominance of parameter training and demonstrates that such dominance still occurs in existing works and it leads to serious optimization problem using adaptive learning rate methods for optimization. We propose a Task-wise Adaptive Learning Rate Method, named AdaTask, to use task-specific accumulative gradients when adjusting the learning rate of each parameter. Comprehensive experiments on synthetic and three real-world datasets demonstrate that AdaTask achieves comparable performance with state-of-the-art MTL approaches.

There are several potential directions for future research. First, when AdaTask significantly improves the performance of dominated task, it could damage the performance of the dominant task, so further research is needed to ensure the performance of the dominant task. Second, deep multi-task learning is a highly nonconvex problem, making convergence analysis quite challenging. Therefore, we expect to perform a theoretical convergence analysis for AdaTask in the future.

## Acknowledgement

# References

Badrinarayanan, V.; Kendall, A.; and Cipolla, R. 2017. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *TPAMI*, 39(12): 2481–2495.

Caruana, R. 1997. Multitask learning. *Machine learning*, 28(1): 41–75.

Chen, C.; Shen, L.; Huang, H.; and Liu, W. 2021. Quantized adam with error feedback. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 12(5): 1–26.

Chen, C.; Shen, L.; Liu, W.; and Luo, Z.-Q. 2022a. Efficient-Adam: Communication-efficient distributed Adam with complexity analysis. *arXiv preprint arXiv:2205.14473*.

Chen, C.; Shen, L.; Zou, F.; and Liu, W. 2022b. Towards practical adam: Non-convexity, convergence theory, and mini-batch acceleration. *Journal of Machine Learning Research*, 23: 1–47.

Chen, Z.; Badrinarayanan, V.; Lee, C.-Y.; and Rabinovich, A. 2018. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *ICML*, 794–803. PMLR.

Chen, Z.; Ngiam, J.; Huang, Y.; Luong, T.; Kretzschmar, H.; Chai, Y.; and Anguelov, D. 2020. Just Pick a Sign: Optimizing Deep Multitask Models with Gradient Sign Dropout. In *NeurIPS*.

Clevert, D.; Unterthiner, T.; and Hochreiter, S. 2016. Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs). In *ICLR*, 1–14.

Dong, D.; Wu, H.; He, W.; Yu, D.; and Wang, H. 2015. Multi-task learning for multiple language translation. In *ACL*, 1723–1732.

Duchi, J. C.; Hazan, E.; and Singer, Y. 2011. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *Journal of Machine Learning Research*, 12: 2121–2159.

Fawcett, T. 2006. An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8): 861–874.

He, Y.; Feng, X.; Cheng, C.; Ji, G.; Guo, Y.; and Caverlee, J. 2022. MetaBalance: Improving Multi-Task Recommendations via Adapting Gradient Magnitudes of Auxiliary Tasks. *WWW*, 2205–2215.

Kendall, A.; Gal, Y.; and Cipolla, R. 2018. Multi-Task Learning Using Uncertainty to Weigh Losses for Scene Geometry and Semantics. In *CVPR*, 7482–7491. IEEE Computer Society.

Kingma, D. P.; and Ba, J. 2015. Adam: A Method for Stochastic Optimization. In *ICLR*.

Lin, X.; Baweja, H. S.; Kantor, G.; and Held, D. 2019a. Adaptive Auxiliary Task Weighting for Reinforcement Learning. In *NeurIPS*, 4773–4784.

Lin, X.; Zhen, H.-L.; Li, Z.; Zhang, Q.-F.; and Kwong, S. 2019b. Pareto multi-task learning. *NeurIPS*, 32.

Liu, B.; Liu, X.; Jin, X.; Stone, P.; and Liu, Q. 2021. Conflict-Averse Gradient Descent for Multi-task Learning. *NeurIPS*, 34.

Liu, S.; Johns, E.; and Davison, A. J. 2019. End-To-End Multi-Task Learning With Attention. In *CVPR*, 1871–1880. Computer Vision Foundation / IEEE.

Liu, S.; Liang, Y.; and Gitter, A. 2019. Loss-Balanced Task Weighting to Reduce Negative Transfer in Multi-Task Learning. In *AAAI*, 9977–9978. AAAI Press.

Liu, Y.; Yang, X.; Xie, D.; Wang, X.; Shen, L.; Huang, H.; and Balasubramanian, N. 2020. Adaptive activation network and functional regularization for efficient and flexible deep multi-task learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 4924–4931.

Ma, J.; Zhao, Z.; Yi, X.; Chen, J.; Hong, L.; and Chi, E. H. 2018a. Modeling Task Relationships in Multi-task Learning with Multi-gate Mixture-of-Experts. In *KDD*, 1930–1939. ACM.

Ma, P.; Du, T.; and Matusik, W. 2020. Efficient continuous pareto exploration in multi-task learning. In *ICML*, 6522–6531. PMLR.

Ma, X.; Zhao, L.; Huang, G.; Wang, Z.; Hu, Z.; Zhu, X.; and Gai, K. 2018b. Entire Space Multi-Task Model: An Effective Approach for Estimating Post-Click Conversion Rate. In *SIGIR*, 1137–1140. ACM.

Misra, I.; Shrivastava, A.; Gupta, A.; and Hebert, M. 2016. Cross-Stitch Networks for Multi-task Learning. In *CVPR*, 3994–4003. IEEE Computer Society.

Navon, A.; Achituve, I.; Maron, H.; Chechik, G.; and Fetaya, E. 2021. Auxiliary Learning by Implicit Differentiation. In *ICLR*. OpenReview.net.

Pan, J.; Mao, Y.; Ruiz, A. L.; Sun, Y.; and Flores, A. 2019. Predicting different types of conversions with multi-task learning in online advertising. In *KDD*, 2689–2697.

Qian, N. 1999. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1): 145–151.

Radford, A.; Wu, J.; Child, R.; Luan, D.; Amodei, D.; Sutskever, I.; et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8): 9.

Reddi, S. J.; Kale, S.; and Kumar, S. 2018. On the Convergence of Adam and Beyond. In *ICLR*.

Ruder, S.; Bingel, J.; Augenstein, I.; and Søgaard, A. 2019. Latent Multi-Task Architecture Learning. In *AAAI*, 4822–4829. AAAI Press.

Sener, O.; and Koltun, V. 2018. Multi-Task Learning as Multi-Objective Optimization. In *NeurIPS*, 525–536.

Shazeer, N.; Mirhoseini, A.; Maziarz, K.; Davis, A.; Le, Q. V.; Hinton, G. E.; and Dean, J. 2017. Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer. In *ICLR*. OpenReview.net.

Sun, X.; Panda, R.; Feris, R.; and Saenko, K. 2020. Adashare: Learning what to share for efficient deep multi-task learning. *NeurIPS*, 33: 8728–8740.

Tang, H.; Liu, J.; Zhao, M.; and Gong, X. 2020. Progressive Layered Extraction (PLE): A Novel Multi-Task Learning

(MTL) Model for Personalized Recommendations. In *RecSys*, 269–278. ACM.

Xi, D.; Chen, Z.; Yan, P.; Zhang, Y.; Zhu, Y.; Zhuang, F.; and Chen, Y. 2021. Modeling the Sequential Dependence among Audience Multi-step Conversions with Multi-task Learning in Targeted Display Advertising. *KDD*.

Xie, R.; Liu, Y.; Zhang, S.; Wang, R.; Xia, F.; and Lin, L. 2021. Personalized Approximate Pareto-Efficient Recommendation. In *WWW*, 3839–3849. ACM / IW3C2.

Yang, C.; Pan, J.; Gao, X.; Jiang, T.; Liu, D.; and Chen, G. 2022. Cross-Task Knowledge Distillation in Multi-Task Recommendation. *AAAI*.

Yosinski, J.; Clune, J.; Bengio, Y.; and Lipson, H. 2014. How transferable are features in deep neural networks? In *NeurIPS*, 3320–3328.

Yu, T.; Kumar, S.; Gupta, A.; Levine, S.; Hausman, K.; and Finn, C. 2020. Gradient Surgery for Multi-Task Learning. *NeurIPS*, 33: 5824–5836.

Zamir, A. R.; Sax, A.; Shen, W.; Guibas, L. J.; Malik, J.; and Savarese, S. 2018. Taskonomy: Disentangling task transfer learning. In *CVPR*, 3712–3722.

Zeiler, M. D. 2012. ADADELTA: An Adaptive Learning Rate Method. *CoRR*, abs/1212.5701.

Zhao, J.; Du, B.; Sun, L.; Zhuang, F.; Lv, W.; and Xiong, H. 2019. Multiple relational attention network for multi-task learning. In *KDD*, 1123–1131.

Zhuang, J.; Tang, T.; Ding, Y.; Tatikonda, S. C.; Dvornek, N.; Papademetris, X.; and Duncan, J. 2020. Adabelief optimizer: Adapting stepsizes by the belief in observed gradients. *NeurIPS*, 33: 18795–18806.

Zou, F.; Shen, L.; Jie, Z.; Sun, J.; and Liu, W. 2018. Weighted adagrad with unified momentum. *arXiv preprint arXiv:1808.03408*.

Zou, F.; Shen, L.; Jie, Z.; Zhang, W.; and Liu, W. 2019. A sufficient condition for convergences of adam and rmsprop. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, 11127–11135.

# Appendix A: Efficient AdaTask

In this section, we will discuss and verify the effect of LA-daTask in more detail. Specifically, the number of parameters, including intermediate variables, and the communication cost for each method is summarized in Table 5. For the number of parameters or intermediate variables, besides the $N$ parameters, RMSProp (as well as other optimizers) needs to maintain $N$ accumulative gradients $G$, one for each parameter. Parameter wise task-aware adaptive learning rate methods, i.e., AdaTask, as defined in Alg. 2, need to maintain $NK$ accumulative gradients $G^k$ during model training, and there is $NK$ additional communication cost to transfer these accumulative gradients from workers to parameter servers. For the Layer wise task-aware adaptive learning rate methods, i.e., LAdaTask as defined in Eq. 9, only $LK$ accumulative gradients need to be maintained.

| Method | #Intermediate Variables | #Communication Cost |
|--------|------------------------|---------------------|
| RMSProp | $O(N)$ | $O(N)$ |
| AdaTask | $O(NK)$ | $O(NK)$ |
| LAdaTask | $O(N + LK)$ | $O(N + LK)$ |

Table 5: The number of parameters and communication cost of each method. Here AdaTask, and LAdaTask denote the Parameter-wise and Layer-wise Task-aware Adaptive Learning Rate methods, respectively.

| Dataset | Method | Task A↓ | Task B↓ | Average↓ | Weighted |
|---------|--------|---------|---------|----------|----------|
| Synthetic | EqualWeight | 0.2172 | 0.0113 | 0.1142 | - |
| | AdaTask | 0.0704 | 0.0417 | 0.0560 | - |
| | LAdaTask | 0.0989 | 0.0285 | 0.0637 | - |

| Dataset | Method | Task A↑ | Task B↑ | Average↑ | Weighted↑ |
|---------|--------|---------|---------|----------|-----------|
| TikTok | EqualWeight | 0.9173 | 0.7437 | 0.8305 | 0.8652 |
| | AdaTask | 0.9312 | 0.7399 | 0.8356 | 0.8738 |
| | LAdaTask | 0.9196 | 0.7446 | 0.8321 | 0.8671 |

Table 6: Performance evaluation of efficient AdaTask (i.e., LAdaTask) on the synthetic dataset (Lower Better) and Tik-Tok dataset (Higher Better). Task $A$ is the dominated task, and task $B$ is the dominant task.

We conduct experiments on Synthetic and TikTok datasets to evaluate AdaTask and LAdaTask, and the results are shown in Table 6. We observe that LAdaTask is better than Equal-Weight on the dominated task (Task $A$), but the improvement is not as significant as AdaTask. Specifically, on the synthetic dataset, LAdaTask achieves an average Task-Normalized RMSE of 0.0637 (lower better), while EqualWeight and AdaTask are 0.1142 and 0.0560, respectively. On the TikTok dataset, the average AUC of LAdaTask is 0.8321 (weighted AUC is 0.8671), and the higher the two evaluation metrics, the better the performance. The EqualWeight and AdaTask are 0.8305(0.8652) and 0.8356(0.8738), respectively.

In addition, we also study LAdaTask's dominance in the synthetic dataset. As shown in Fig. 3(b) and (c), about $85\%$ of LAdaTask's parameters are not dominated by any task, while almost all parameters in AdaTask are not dominated. However, in EqualWeight, almost all parameters are dominated by task $B$ (Fig. 3(a)). This shows that LAdaTask can alleviate the task dominance problem to a certain extent.
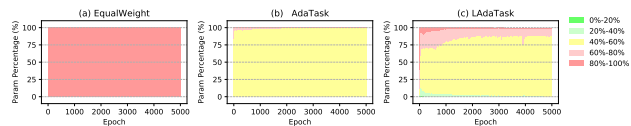


Figure 3: rAU$(i, T, B)$ of shared parameters on the synthetic dataset for EqualWeight(RMSprop), AdaTask, and LAdaTask methods.

# Appendix B: Experiment Details

This section describes the details of our experimental setting, including baseline methods, datasets, and parameter settings.

## Baselines

We present the performance mainly comparison of AdaTask with the following MTL approaches:

- **EqualWeight** (Caruana 1997): A basic MTL model that learns shared bottom across tasks as well as task-specific towers.

- **GradNorm** (Chen et al. 2018): A MTL optimization approach that automatically balances training in deep multi-task models by normalizing gradients across tasks.

- **Uncertainty Weighting** (UW) (Kendall, Gal, and Cipolla 2018): A MTL optimization approach that weighs multiple loss functions by considering the homoscedastic uncertainty of each task.

- **PCGrad/Gradient Surgery** (Yu et al. 2020): A MTL optimization technique that projects the gradient of a task onto the normal plane of the gradient of any other task to avoid gradient conflicts.

- **CAGrad** (Liu et al. 2021) A MTL optimization technique that seeks a gradient update direction in the mean gradient direction neighborhood, to maximize the sub-task with the smallest lift in all tasks.

## Datasets

We conducted experiments on three large-scale real-world datasets, namely CityScapes, TikTok, and WeChat.

- **CityScapes** dataset is a multitask learning benchmark for depth estimation (Task $A$) and semantic segmentation (Task $B$). Cityscapes have 5,000 high-resolution images of street scenes from 50 different cities[1]. In this paper, we use the dataset provided by MTAN[2].

- **TikTok** dataset comes from the ICME 2019 Short Video Content Understanding and Recommendation Competition[3]. TikTok consists of two tasks (Like, Finish), named Task $A$ and Task $B$, with official weights of 7:3. The complete dataset contains 19,622,340 samples and 4,691,488 features. Since the dataset does not provide timing information in logs, we randomly divided it into training, validation and test sets with a ratio of 8:1:1.

---

[1]https://www.cityscapes-dataset.com/

[2]https://github.com/lorenmt/mtan

[3]https://www.biendata.xyz/competition/icmechallenge2019/

| Method | Task A ↑ | Task B ↑ | Task C ↑ | Task D ↑ | Average AUC ↑ | Weighted AUC ↑ |
|---|---|---|---|---|---|---|
| EqualWeight (SharedBottom) | 0.9080 | 0.8912 | 0.8881 | 0.9514 | 0.9097 (+0.00%) | 0.9160 (+0.00%) |
| AdaTask (SharedBottom) | 0.9643 | 0.9391 | 0.8920 | 0.9480 | 0.9359 (+2.87%) | 0.9311 (+1.63%) |
| EqualWeight (MMoE: 4 experts) | 0.9126 | 0.9052 | 0.8898 | 0.9512 | 0.9147 (+0.00%) | 0.9096 (+0.00%) |
| AdaTask (MMoE: 4 experts) | 0.9568 | 0.9263 | 0.8949 | 0.9499 | 0.9320 (+1.89%) | 0.9345 (+2.73%) |

Table 7: Performance evaluation of different architectures (SharedBottom v.s. MMoE) on the WeChat dataset. Task $A$, $B$, and $C$ are the dominated tasks, and task $B$ is the dominant task.

- **WeChat** dataset is from the 2021 WeChat Big Data Challenge[4]. It consists of four tasks (named Task $A$-$D$), with official weights of 1:2:3:4. It consists of 8,373,653 interaction data samples across 14 days. To simulate the real-world setting, we split the WeChat data by dates, i.e., the first 12 days as the training set, the next one day as the validation set, and the last day as the testing set.

## Implementation Details

To evaluate the methods on the **CityScapes** dataset, we follow the same setting as MTAN (Liu, Johns, and Davison 2019) and CAGrad (Liu et al. 2021). Specifically, we use the MTAN (Liu, Johns, and Davison 2019) as the backbone. After combining various baseline methods into the backbone, we use the Adam (Kingma and Ba 2015) optimizer with an initial learning rate of 0.0001 to train 200 epochs, and halve the learning rate at the 100-th epoch. The batch size is set to 8. As MTAN and CAGrad, we do not create a separate validation set, and report the average performance of the last 10 epochs of the test set.

To evaluate the methods on **TikTok and WeChat** datasets, we choose Shared Bottom as the backbone, which is widely used in multi-task recommendation systems (Ma et al. 2018b; Xi et al. 2021; He et al. 2022). The Shared Bottom consists of an embedding layer and two hidden layers. In addition, each task has a private tower composed of three layers of MLPs with ELU, with the number of neurons as [128, 64, 32], respectively. The embedding size is set as 10 for all methods on both datasets. We used RMSprop with an initial learning rate of 0.001 and a batch size of 4,096. For TikTok and WeChat datasets, we implement low-frequency filtering for the original features, so that features that occur less than ten times are uniformly set to 'Other'. In addition, we use the early-stop strategy and end the training when the result does not improve on the validation set for five consecutive epochs.

In all datasets, we use three random seeds for all methods and report average results. For all baseline methods, we use the parameter search range recommended in the original paper (Chen et al. 2018; Kendall, Gal, and Cipolla 2018; Yu et al. 2020; Liu et al. 2021) for hyper-parameters tuning. All experiments are conducted on an NVIDIA V100 GPU with 16 GB memory.

## Appendix C: Compared with MOO methods and More Model Architectures&Optimizers

In this section, we conduct experiments in the following aspects: 1) Performance comparison and discussion between

| | Task A | | Task B | | Rel Avg |
|---|---|---|---|---|---|
| Method | AbsErr↓ | RelErr↓ | mIoU↑ | PixAcc↑ | $\Delta p$ ↑ |
| MGDA | 0.0166 | **33.50** | 69.42 | 91.22 | -13.26% |
| ParetoMTL(0.25, 0.75) | 0.0154 | 38.89 | 74.97 | **93.52** | -6.56% |
| ParetoMTL(0.75, 0.25) | 0.0149 | 36.30 | **75.18** | 93.45 | -3.66% |
| AdaTask | **0.0128** | 36.85 | 75.02 | 93.40 | 0.00% |

Table 8: Performance evaluation of MOO methods (MGDA and ParetoMTL) on the CityScapes dataset. Task $A$ is the dominated task, and task $B$ is the dominant task.

AdaTask and MOO methods; 2) Validation of AdaTask on more architectures and optimizers.

## Comparison with MOO methods

Recent advances formulate MTL as a Multi-Objective Optimization (MOO) problem (Sener and Koltun 2018; Lin et al. 2019b), and these approaches achieve a Pareto stationary solution, or a set of Pareto solutions with various trade-off preferences. In this section, we compare AdaTask with ParetoMTL (Lin et al. 2019b), a representative MOO approach to MTL. As shown in Table 8, we compare our AdaTask with MGDA (Sener and Koltun 2018) and ParetoMTL (Lin et al. 2019b) with different preferences. Specifically, ParetoMTL(0.75, 0.25) and ParetoMTL(0.25, 0.75) represent the Pareto solution inclined to task $A$ and $B$, respectively. We have the following observations. From the perspective of overall performance (i.e., $\Delta p$), AdaTask is better than Pareto MTL with either (0.25, 0.75) or (0.75, 0.25) preference. From the MOO perspective, AdaTask also achieves a Pareto solution in the sense that Pareto MTL with various preferences doesn't dominate AdaTask.

## More Model Architectures and Optimizers

In the main text, we apply the MTAN (Liu, Johns, and Davison 2019) model on the CityScapes dataset and the Share-Bottom (Caruana 1997) model on the WeChat dataset. In this subsection, we apply AdaTask to more multi-task network architectures and optimizer. Specifically, we use the SegNet (Badrinarayanan, Kendall, and Cipolla 2017) on the CityScapes dataset, and the MMoE (Ma et al. 2018a) on the WeChat dataset. In addition, we also tested the RMSProp optimizer on the CityScapes dataset. The results are presented in Table 7 and Table 9.

On the WeChat dataset, we used average AUC and weighted AUC to measure the overall performance. As shown in Table 7, by applying AdaTask on MMoE, we found that it can further improve the average AUC by +1.89% compared with applying EqualWeight on MMoE. In addition, on the MMoE model, AdaTask improves the weighted AUC by 2.73% compared to EqualWeight. Similar to MTAN, MMoE can also allocate capacity among tasks via attention.

Similar to CAGrad (Liu et al. 2021), we compute the average per-task performance improvement $\Delta p$ of AdaTask to the EqualWeight baseline on the CityScapes dataset. As shown in Table 9, we found that by applying AdaTask on SegNet, overall performance resulted in a $23.95\%$ improvement. We also observe that MTAN outperforms SegNet in all task evaluation metrics (i.e., Abs Err, Rel Err, mIoU, and Pix Acc), because it designs task-specific attention mechanisms to share the network backbone flexibly.

Finally, we add the experiments of RMSprop with AdaTask on the CityScapes dataset. We observe that it brings a $6.84\%$ average improvement compared to the traditional RMSProp (EqualWeight baseline). These results are consistent with the one using Adam($9.34\%$ average improvement) and validate the effectiveness of AdaTask for different optimizers on the same dataset.
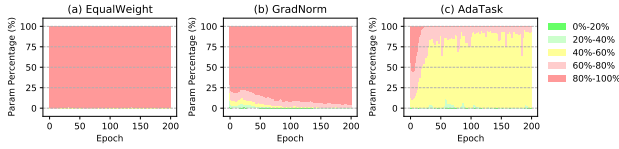


Figure 4: $\mathrm{rAU}(i, T, B)$ of shared parameters on the CityScapes dataset for EqualWeight, GradNorm, and AdaTask methods.

## Appendix D: Task Dominance Analysis
### Task Dominance on the CityScapes

In the main paper, we analyze the task dominance of multiple MTL methods on the synthetic dataset (i.e., Fig. 2). Here we present the task dominance analysis on real-world datasets. We choose EqualWeight, GradNorm and AdaTask in the following analysis, and present the results in Fig. 4.

We observe that $99\%$ parameters in EqualWeight are dominated by task $B$; that is, $\mathrm{rAU}(i, T, B)$ of almost all parameters belongs to $(80\%, 100\%]$, which is the red area in Fig. 4(a). In addition, the GradNorm (Fig. 4(b)) has $95\%$ of the parameters dominated by task $B$. However, only $7\%$ of shared parameters in AdaTask(Fig. 4(c)) are dominated, i.e., $93\%$ of $\mathrm{rAU}(i, T, B)$ belongs to $(40\%, 60\%]$. These results verify that most parameters of EqualWeight and GradNorm are dominated by one task in real-world datasets, while AdaTask solves it decently.

### Layer-wise Task Dominance on the Synthetic Dataset

In this section, we present evaluation results of AdaTask. As shown in Fig. 6, PCGrad is similar to EqualWeight in the sense that they are completely dominated by task $B$ at each layer. UW, CAGrad, and GradNorm are balanced better, especially at the bottom layer (1st layer), while they are still dominated at the top layers. In contrast, AdaTask achieves a good balance across all layers.

### Dominance of Learning Rates on the Synthetic Dataset

In this section we demonstrate the dominance of learning rates in MTL. Similar to the parameter cumulative gradient-

---

**Algorithm 3: CAGrad(Adam) with AdaTask**

1: **Require** Model parameters $\theta$, a constant $c \in [0, 1)$
2: **For** step $t \in \{1, 2, \dots, T\}$ **do**
3:    $(\mathbf{g}_t^k)^{\mathrm{PC}} \leftarrow \mathbf{g}_t^k \ \ \forall k$
4:    $g_0 = \frac{1}{K} \sum_k g_t^k$
5:    $\phi = c^2 \|g_0\|^2$
6:    Solve $\min_{w \in \mathcal{W}} F(w) := g_w^\top g_0 + \sqrt{\phi} \|g_w\|$, where $g_w = \frac{1}{K} \sum_k w^k g^k$
7:    **For** task $k \in \{1, 2, \dots, K\}$ **do**
8:      $(\mathbf{g}_t^k)^{\mathrm{CA}} \leftarrow \frac{1}{K}(1 + \frac{\phi^{0.5}}{\|g_w\|} w^k) g^k$
9:      $\Delta \theta_t^k \leftarrow -\frac{\eta}{\sqrt{\hat{\mathbf{G}}_t^k} + \epsilon} \hat{\mathbf{m}}_t^k$
10:      where
11:      $\hat{\mathbf{G}}_t^k = \mathbf{G}_t^k / (1 - \gamma_1^t), \quad \hat{\mathbf{m}}_t^k = \mathbf{m}_t^k / (1 - \gamma_2^t)$
12:      $\mathbf{G}_t^k = \gamma_1 \mathbf{G}_{t-1}^k + (1 - \gamma_1)((\mathbf{g}_t^k)^{\mathrm{CA}})^2$
13:      $\mathbf{m}_t^k = \gamma_2 \mathbf{m}_{t-1}^k + (1 - \gamma_2)(\mathbf{g}_t^k)^{\mathrm{CA}}$
14:    **End For**
15:    $\theta_t = \theta_{t-1} + \sum_k \Delta \theta_t^k$
16: **End For**

---

**Algorithm 4: PCGrad(Adam) with AdaTask**

1: **Require** Model parameters $\theta$
2: **For** step $t \in \{1, 2, \dots, T\}$ **do**
3:    **For** task $k \in \{1, 2, \dots, K\}$ **do**
4:      $(\mathbf{g}_t^k)^{\mathrm{PC}} \leftarrow \mathbf{g}_t^k \ \ \forall k$
5:      **For** task $j \sim \{1, 2, \dots, K\} \setminus k$ in random order **do**
6:        **If** $(\mathbf{g}_t^k)^{\mathrm{PC}} \cdot \mathbf{g}_t^j < 0$ **then**
7:        // *Subtract the projection of* $(\mathbf{g}_t^k)^{\mathrm{PC}}$ *into* $\mathbf{g}^j$
8:        Set $(\mathbf{g}_t^k)^{\mathrm{PC}} = (\mathbf{g}_t^k)^{\mathrm{PC}} - \frac{(\mathbf{g}_t^k)^{\mathrm{PC}} \odot \mathbf{g}_t^j}{\|\mathbf{g}_t^j\|^2} \mathbf{g}_t^j$
9:        **End If**
10:      **End For**
11:      $\Delta \theta_t^k \leftarrow -\frac{\eta}{\sqrt{\hat{\mathbf{G}}_t^k} + \epsilon} \hat{\mathbf{m}}_t^k$
12:      where
13:      $\hat{\mathbf{G}}_t^k = \mathbf{G}_t^k / (1 - \gamma_1^t), \quad \hat{\mathbf{m}}_t^k = \mathbf{m}_t^k / (1 - \gamma_2^t)$
14:      $\mathbf{G}_t^k = \gamma_1 \mathbf{G}_{t-1}^k + (1 - \gamma_1)((\mathbf{g}_t^k)^{\mathrm{PC}})^2$
15:      $\mathbf{m}_t^k = \gamma_2 \mathbf{m}_{t-1}^k + (1 - \gamma_2)(\mathbf{g}_t^k)^{\mathrm{PC}}$
16:    **End For**
17:    $\theta_t = \theta_{t-1} + \sum_k \Delta \theta_t^k$
18: **End For**

---

dominated analysis in the main text, we analyze the learning rate-dominated problem on synthetic dataset (Eq. 1) based on the SharedBottom model (containing four shared fully-connected layers). We randomly select ten parameters from each shared layer and calculate the learning rates calculated using their own accumulative gradients (i.e., `LR-A` and `LR-B`) and the learning rates calculated using the overall gradients accumulated across both tasks (i.e., `LR-Whole`).

Fig. 5(a) presents the learning rate of parameters of each layer in the training process. We observe that: `LR-B` is much closer to `LR-Whole` than `LR-A` because the overall accumulative gradients are dominated by task $B$. Task $A$ should have a larger learning rate based on its own accumulative gradients.

As shown in Fig. 5(b), our proposed AdaTask adjusts the learning rate by using their own accumulative gradients, there-

| | Task A | | Task B | | Rel Avg |
| Method | Abs Err↓ | Rel Err↓ | mIoU↑ | Pix Acc↑ | $\Delta p$ ↑ |
|---|---|---|---|---|---|
| EqualWeight (SegNet model, Adam optimizer) | 0.0198 | 121.81 | 70.23 | 91.92 | +0.00% |
| AdaTask (SegNet model, Adam optimizer) | 0.0134 | 45.53 | 70.76 | 92.03 | +23.95% |
| EqualWeight (MTAN model, Adam optimizer) | 0.0152 | 47.00 | 75.01 | 93.40 | +0.00% |
| AdaTask (MTAN model, Adam optimizer) | 0.0128 | 36.85 | 75.02 | 93.40 | +9.34% |
| EqualWeight (MTAN model, RMSProp optimizer) | 0.0152 | 44.63 | 74.86 | 93.43 | +0.00% |
| AdaTask (MTAN model, RMSProp optimizer) | 0.0132 | 37.46 | 74.05 | 93.02 | +6.84% |

Table 9: Performance evaluation of different architectures (SegNet v.s. MTAN) on the CityScapes dataset. Task $A$ is the dominated task, and task $B$ is the dominant task.
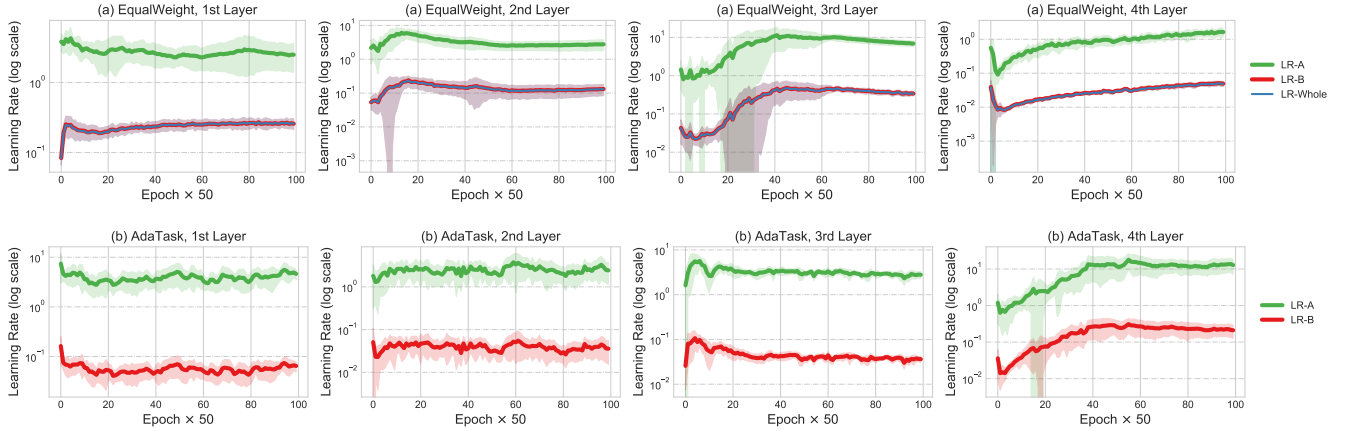


Figure 5: Illustration of learning rate dominance of randomly sampled parameters from the 1st to 4th layer for two MTL methods: (a) EqualWeight and (b) our AdaTask. The learning rate of Task $A$ (i.e., `LR-A`) or Task $B$ (i.e., `LR-B`) is based on Task $A$'s or Task $B$'s own accumulative gradient, respectively. The whole learning rate (i.e., `LR-Whole`) is based on the sum of accumulative gradients from both tasks.

fore there is no dominance on learning rates at all.

## Appendix E: Training Time

In this section, we evaluate the training time cost by running each method for one epoch on the CityScapes dataset. As shown in Table 10, it can be observed that the training time of AdaTask is less than all other MTL methods except Equal-Weight. This is because, compared to other MTL approaches, AdaTask only needs to separate the accumulative gradients, which introduces little additional computation cost.

| Method | EqualWeight | PCGrad | GradNorm |
|---|---|---|---|
| Time cost | 249.20s | 343.01s | 346.51s |
| Method | CAGrad | MGDA | AdaTask (ours) |
| Time cost | 353.64s | 442.62s | 332.14s |

Table 10: Each method performs the time cost comparison of a single Epoch on the CityScapes dataset.

## Appendix F: Adaptive Learning Rate Methods (RMSProp and Adam) v.s. SGD Method

In this section, we compare the performance of the adaptive learning rate methods and the global learning rate method on the CityScapes dataset. Like CAGrad (Liu et al. 2021), we use average per-task performance improvement to measure the overall performance improvement of RMSProp and

| | Task A | | Task B | | Rel Avg |
| Optimizer | Abs↓ | Rel↓ | mIoU↑ | Pix↑ | $\Delta p$↑ |
|---|---|---|---|---|---|
| SGD | 0.0242 | 115.35 | 66.51 | 90.89 | 0.00% |
| RMSProp | 0.0152 | 44.63 | 74.86 | 93.43 | +28.33% |
| Adam | 0.0152 | 47.00 | 75.01 | 93.40 | +27.86% |

Table 11: Performance comparison of adaptive learning rate approaches and SGD on CityScapes dataset.

Adam over SGD. As shown in Table 11, we find that adaptive learning rate approaches (RMSProp and Adam) perform better than the SGD. Specifically, in the CityScapes dataset, we observe a 28.33% increase in RMSProp relative to SGD and a 27.86% increase in Adam relative to SGD. However, as pointed out in this paper, these traditional adaptive learning rate optimizers accumulate gradients across tasks in multi-task learning. Therefore they still suffer from accumulative gradient dominance.

## Appendix G: Algorithm Pseudocode

PCGrad (Yu et al. 2020) and CAGrad (Liu et al. 2021) are two gradient direction modification methods, and they are orthogonal to our proposed AdaTask and therefore they can be combined with AdaTask. We present the pseudocode of PC-Grad+AdaTask and CAGrad+AdaTask in Alg. 4 and Alg. 3, respectively. Their performances are reported in Table 1.

Figure 6: rAU$(i, T, B)$ of shared parameters in each layer on the synthetic dataset for six MTL methods: EqualWeight, PCGrad, UW, CAGrad, GradNorm, and our AdaTask.