

# Improving Generalization of Meta-learning with Inverted Regularization at Inner-level

Lianzhe Wang<sup>1\*</sup> Shiji Zhou<sup>1\*</sup> Shanghang Zhang<sup>2†</sup> Xu Chu<sup>1</sup> Heng Chang<sup>1</sup> Wenwu Zhu<sup>1†</sup>

<sup>1</sup>Tsinghua University <sup>2</sup>National Key Laboratory for Multimedia Information Processing, Peking University.

{wanglz20, zsj17, changh17}@mails.tsinghua.edu.cn,  
shanghang@pku.edu.cn, {chu.xu, wwzhu}@tsinghua.edu.cn

## Abstract

*Despite the broad interest in meta-learning, the generalization problem remains one of the significant challenges in this field. Existing works focus on meta-generalization to unseen tasks at the meta-level by regularizing the meta-loss, while ignoring that adapted models may not generalize to the task domains at the adaptation level. In this paper, we propose a new regularization mechanism for meta-learning – Minimax-Meta Regularization, which employs inverted regularization at the inner loop and ordinary regularization at the outer loop during training. In particular, the inner inverted regularization makes the adapted model more difficult to generalize to task domains; thus, optimizing the outer-loop loss forces the meta-model to learn meta-knowledge with better generalization. Theoretically, we prove that inverted regularization improves the meta-testing performance by reducing generalization errors. We conduct extensive experiments on the representative scenarios, and the results show that our method consistently improves the performance of meta-learning algorithms.*

## 1. Introduction

Meta-learning has been proven to be a powerful paradigm for extracting well-generalized knowledge from previous tasks and quickly learning new tasks [47]. It has received increasing attention in many machine learning settings such as few-shot learning [10, 45, 46, 50] and robust learning [27, 39, 42], and can be deployed in many practical applications [7, 21, 29, 54]. The key idea of meta-learning is to improve the learning ability of agents through a learning-to-learn process. In recent years, optimization-based algorithms have emerged as a popular approach for realizing the learning-to-learn process in meta-learning [10, 28]. These methods formulate the problem as a bi-level optimization problem and have demonstrated impressive per-

formance across various domains, leading to significant attention from the research community. The primary focus of our paper is to further advance this line of research.

The training process of meta-learning takes place at two levels [10, 19]. At the inner-level, a base model, which is initialized using the meta-model’s parameters, adapts to each task by taking gradient descent steps over the support set. At the outer-level, a meta-training objective is optimized to evaluate the generalization capability of the initialization on all meta-training tasks over the query set, helping to ensure that the model is effectively optimized for the desired goal. With this learning-to-learn process, the final trained meta-model could be regarded as the model with good initialization to adapt to new tasks.

Despite the success of meta-learning, the additional level of learning also introduces a new source of potential overfitting [36], which poses a significant challenge to the generalization of the learned initialization. This generalization challenge is twofold: first, the meta-model must generalize to unseen tasks (*meta-generalization*); and second, the adapted model must generalize to the domain of a specific task, which we refer to as *adaptation-generalization*. As the primary objective of meta-learning is to achieve strong performance when adapting to new tasks, the ability of the meta-model to generalize well is critical. Recent works aim to address the meta-generalization problem by meta-regularizations, such as constraining the meta-initialization space [52], enforcing the performance similarity of the meta-model on different tasks [20], and augmenting meta-training data [33, 36, 51]. These approaches are verified to enhance generalization to unseen tasks. However, they do not address the problem of adaptation-generalization to the data distribution of meta-testing tasks.

To address this issue, we propose Minimax-Meta Regularization, a novel regularization mechanism that improves both adaptation-generalization and meta-generalization. Specifically, our approach particularly employs inverted regularization at the inner-level to hinder the adapted model’s generalizability to the task domain. This forces the

\*Equal contributions †Corresponding authors

meta-model to learn hypotheses that better generalize to the task domains, which improves adaptation-generalization. Meanwhile, we use ordinary regularization at the outer-level to optimize the meta-model’s generalization to new tasks, which helps meta-generalization. By improving both adaptation-generalization and meta-generalization simultaneously, our method results in a more robust and effective meta-learning regularization mechanism.

Theoretically, we prove that under certain assumptions, if we add L2-Norm as the regularization term to the inner-level loss function, the *inverted regularization* will reduce the generalization bound of MAML, while the *ordinary regularization* will increase the generalization bound. In terms of total test error, which includes both generalization error and training bias caused by regularization, the inverted L2-Norm also reduces the total test error when the reg parameter is selected within a negative interval. These results suggest that the regularization at the inner-level should be inverted. As it has been verified that ordinary regularization at the outer-level helps the meta-generalization, our theory implies that the proposed Minimax-Meta Regularization helps both meta-generalization and adaptation-generalization.

We conduct experiments on the few-shot classification problem for MAML [10] with different regularization types (ordinary/inverted) at the inner- and outer-level. The results demonstrate the efficacy of Minimax-Meta Regularization, and support the theoretical results that regularization at the inner-level improves test performance only when it’s inverted. Additionally, we empirically verify that Minimax-Meta regularization can be applied with different types of regularization terms (norm/entropy), implying the flexibility for applying the proposed method in practice.

## 2. Related Work

**Meta-learning.** A line of meta-learning methods has sought to train recurrent neural networks that ingest entire datasets [8, 41]. However, they need to place constraints on the model architecture. Another line aims to learn a transferable metric space between samples from previous tasks [31, 34, 44, 49]. However, it is mainly limited to classification problems. In this paper, we focus on optimization-based meta-learning methods that learn a meta-initialization [10–13, 18, 26, 28, 35], which are well-generalized for meta-training tasks, being agnostic to both model architecture and problems. However, these approaches are shown to be overfitting the meta-training tasks [6, 40, 51, 53].

**Meta-Regularization.** Standard regularizations such as weight decay [22], dropout [15], and incorporating noise [1, 2, 48], which can significantly enhance the generality of single-level machine learning. However, it limits the flexibility of fast adaptation in the inner-level [51]. MR-MAML [52] constrains the search space of the meta-model and allows the adaptation to be sufficient at the inner-

level. Jamal *et al.* [20] proposed TAML to enforce the meta-model to perform similarly across tasks. Rajenran *et al.* [37] explored an information-theoretic framework of meta-augmentation. Yao *et al.* [51] proposed two task augmentation methods – MetaMix and Channel Shuffle, which is theoretically proven to be generalized to unseen tasks. Ni *et al.* [33] investigated the distinct ways where data augmentation can be integrated at both the image and class levels. Rothfuss *et al.* [40] addressed the meta-generalization problem using the PAC-Bayesian framework. However, these works focus on meta-generalization, while adaptation-generalization is merely considered.

## 3. Preliminary

Model-Agnostic Meta-Learning (MAML) [10] with a single inner-step is adopted as the representative algorithm to derive the theoretical results in this paper. We follow the framework proposed by Fallah *et al.* [9] to make problem formulation for MAML with a single inner-step. We denote each data point by  $z = (x, y) \in \mathcal{Z}$  and evaluate the performance of a model parameterized by  $w \in \mathcal{W}$  using loss function  $\ell(w, z)$ . Tasks  $\{\mathcal{T}_i\}_{i=1}^m$  are drawn from distributions  $\{\mathcal{P}_i\}_{i=1}^m$ , with corresponded *population loss* for model  $w$  defined as  $\mathcal{L}_i(w) := \mathbb{E}_{z \sim p_i}[\ell(w, z)]$ . Throughout the paper, we adopt the hat notation to denote empirical losses, i.e.,  $\hat{\mathcal{L}}(w; \mathcal{D}) := \frac{1}{|\mathcal{D}|} \sum_{z \in \mathcal{D}} \ell(w, z)$  means the *empirical loss* of model  $w$  with dataset  $\mathcal{D}$ .

$F_i(w)$  is defined to evaluate the performance of the model updated by one single stochastic gradient descent (SGD) from  $w$ , on task  $\mathcal{T}_i$ .  $\mathcal{D}_i$  denotes a data batch consisting of  $K$  samples drawn from  $\mathcal{P}_i$ . The goal of MAML is to find a good model parameter  $w$  that generally performs well across different tasks after taking the SGD step:

$$\begin{aligned} \min_{w \in \mathcal{W}} F(w) &:= \frac{1}{m} \sum_{i=1}^m F_i(w) \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathcal{D}_i} \mathbb{E}_{z \sim p_i} \left[ \ell \left( w - \frac{\alpha}{K} \sum_{z' \in \mathcal{D}_i} \nabla \ell(w, z'), z \right) \right] \end{aligned} \quad (1)$$

However, directly solving (1) is usually impractical since the true task distributions  $\{\mathcal{P}_i\}_{i=1}^m$  are usually unknown. Instead, the common practice is to approximate  $F_i$  by the empirical loss. For simplicity, suppose we have access to totally  $2n$  training samples from each task  $\mathcal{T}_i$ , and we further group the samples into two distinct sets of size  $n$ :  $\mathcal{S}_i^{\text{in}}$  for meta-training(support) at inner-level and  $\mathcal{S}_i^{\text{out}}$  for meta-validation(query) at outer-level. Then, for each task  $\mathcal{T}_i$ , we have one corresponding training set  $\mathcal{S}_i := \{\mathcal{S}_i^{\text{in}}, \mathcal{S}_i^{\text{out}}\}$ . During training, each distinct  $K$ -shot data batch  $\mathcal{D}_i$  is sampled from each  $\mathcal{S}_i^{\text{in}}$  to serve as a meta-training(support) set.

The approximation of (1) is given by

$$\arg \min_{w \in \mathcal{W}} \hat{F}(w, \mathcal{S}) := \frac{1}{m} \sum_{i=1}^m \hat{F}_i(w, \mathcal{S}_i) \quad (2)$$

where  $\mathcal{S} := \{\mathcal{S}_i\}_{i=1}^m$ . And  $\hat{F}_i$  stands for empirical loss that estimates  $F_i$  by

$$\hat{F}_i(w, \mathcal{S}_i) := \frac{1}{\binom{n}{k}} \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}|=k}} \frac{1}{n} \sum_{z \in \mathcal{S}_i^{\text{out}}} \ell \left( w - \frac{\alpha}{K} \sum_{z' \in \mathcal{D}_i^{\text{in}}} \nabla \ell(w, z'), z \right)$$

MAML solves the minimization problem in (2) by using each per-task gradient  $\nabla \hat{F}_i(w, \mathcal{S}_i)$  to take SGD step at meta-level. Specifically, at each iteration  $t$ , for each sampled task data  $\{\mathcal{D}_i^{t, \text{in}}, \mathcal{D}_i^{t, \text{out}}\}$ , MAML calculates

$$w_i^{t+1} := w^t - \beta_t \nabla_{w^t} \hat{\mathcal{L}} \left( w^t - \alpha \nabla \hat{\mathcal{L}} \left( w^t, \mathcal{D}_i^{t, \text{in}} \right), \mathcal{D}_i^{t, \text{out}} \right) \quad (3)$$

and update the model at the end of each iteration by

$$w^{t+1} := \frac{1}{r} \sum_{i \in \mathcal{B}_t} w_i^{t+1}$$

where  $\mathcal{B}_t$  is the set of indices of  $r$  randomly chosen tasks at iteration  $t$ . When referring to the per-task adapted model in the paper, we denote it as  $w_i^{t+1}$  and its calculation is in fact embedded within (3), that is,  $w_i^{t+1} := w^t - \alpha \nabla \hat{\mathcal{L}} \left( w^t, \mathcal{D}_i^{t, \text{in}} \right)$ .

In the context of evaluating the performance of meta-learning algorithms, the test error is generally considered the most critical metric. This error represents the population loss of a meta-model, denoted as  $\mathcal{A}(\mathcal{S})$ , obtained by algorithm  $\mathcal{A}$  with a given dataset  $\mathcal{S}$ . The test error can be decomposed into three distinct terms:

$$\begin{aligned} \mathbb{E}_{\mathcal{A}, \mathcal{S}} \left[ F(\mathcal{A}(\mathcal{S})) - \min_{\mathcal{W}} F \right] & \quad (\text{test error}) = \\ \underbrace{\mathbb{E}_{\mathcal{A}, \mathcal{S}} \left[ \hat{F}(\mathcal{A}(\mathcal{S}), \mathcal{S}) - \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) \right]}_{\text{training error}} & \\ + \underbrace{\mathbb{E}_{\mathcal{A}, \mathcal{S}} \left[ F(\mathcal{A}(\mathcal{S})) - \hat{F}(\mathcal{A}(\mathcal{S}), \mathcal{S}) \right]}_{\text{generalization error}} & + \underbrace{\mathbb{E}_{\mathcal{S}} \left[ \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) \right] - \min_{\mathcal{W}} F}_{\leq 0} \end{aligned} \quad (4)$$

Fallah *et al.* [9] have shown that the first training error term will converge to zero as the number of training steps  $T$  increases, given that the loss function  $\ell(w, z)$  satisfies certain assumptions, and that the third term is non-positive. Therefore, to improve the performance of the obtained model on the test error, we aim to apply regularization to reduce the generalization error term.

## 4. Method

In this section, we introduce the Minimax-Meta Regularization method for bi-level meta-learning and its application to the popular MAML algorithm. We also provide an intuitive explanation of the effectiveness of the inner-level inverted regularization.

### 4.1. Minimax-Meta Regularization

Our Minimax-Meta Regularization method is designed to improve the generalization performance of bi-level meta-learning by combining two types of regularizations: one at the outer-level and the other at the inner-level. In particular, we propose to use an ordinary regularization at the outer-level to encourage the meta-model to learn more generalized hypotheses, and an inverted regularization at the inner-level to increase the adaptation difficulty and help the meta-model improve generalization during training.

Specifically, when the regularizations involved can be achieved in the loss function, the Minimax-Meta Regularization shifts the learning objective of the inner level from  $\hat{\mathcal{L}} \left( w^t, \mathcal{D}_i^{t, \text{in}} \right)$  to

$$\hat{\mathcal{L}} \left( w^t, \mathcal{D}_i^{t, \text{in}} \right) + \sigma^{\text{in}} \text{Inverted\_Reg} \left( w^t, \mathcal{D}_i^{t, \text{in}} \right),$$

and the learning objective of the outer level from  $\hat{\mathcal{L}} \left( w_i^t, \mathcal{D}_i^{t, \text{out}} \right)$  to

$$\hat{\mathcal{L}} \left( w_i^t, \mathcal{D}_i^{t, \text{out}} \right) + \sigma^{\text{out}} \text{Ordinary\_Reg} \left( w_i^t, \mathcal{D}_i^{t, \text{out}} \right),$$

where  $\sigma^{\text{in}}$  and  $\sigma^{\text{out}}$  are regularization coefficients.

The outer-level regularization term  $\text{Ordinary\_Reg}(w, \mathcal{D})$  can be any classic ordinary regularization term, such as L1/L2-Norm or information entropy regularization, which encourages the meta-model to learn more generalized hypotheses. In contrast, the inner-level regularization term  $\text{Inverted\_Reg}(w, \mathcal{D})$  should be an inverted regularization term, which could typically be achieved by changing the sign of an ordinary regularization term (e.g., negative L1/L2-Norm, inverted entropy regularization), and this increases the adaptation difficulty and forces the meta-model to learn better-generalized hypotheses.

It is worth noting that the inner-level inverted regularization is only added during the training phase, and we do not use it for the meta-testing phase. Specifically, during the meta-testing phase, which evaluates the performance of the learned meta-model on new tasks, we only adapt the model without any additional regularization to avoid influencing its task-specific performance.

**Intuition for Inverted Regularization at Inner-level.** The intuition behind using inverted regularization at the inner-level is that it can help the meta-model learn better-generalized hypotheses (meta-knowledge) by increasing the

---

**Algorithm 1** Minimax-MAML
 

---

**Require:** Datasets  $\mathcal{S} = \{\mathcal{S}_i^{\text{in}}, \mathcal{S}_i^{\text{out}}\}_{i=1}^m$ ; total number of iterations  $T$ ; regularization coefficients  $\sigma^{\text{in}}$  and  $\sigma^{\text{out}}$ .

- 1: Initialize the meta-model  $w^0$
- 2: **for**  $t = 0$  to  $T - 1$  **do**
- 3:   Randomly sample  $r$  tasks with indices stored in  $\mathcal{B}_t$ ;
- 4:   **for** each sampled task  $\mathcal{T}_i$  **do**
- 5:     Sample a support data batch  $\mathcal{D}_i^{t, \text{in}}$  from  $\mathcal{S}_i^{\text{in}}$ ;
- 6:     Sample a query data batch  $\mathcal{D}_i^{t, \text{out}}$  from  $\mathcal{S}_i^{\text{out}}$ ;
- 7:     (Inner-level) Compute per-task adapted parameters with gradient descent:

$$w_i^t := w^t - \alpha \nabla_{w^t} \left( \hat{\mathcal{L}}(w^t, \mathcal{D}_i^{t, \text{in}}) + \sigma^{\text{in}} \text{Inverted\_Reg}(w^t, \mathcal{D}_i^{t, \text{in}}) \right);$$

- 8:     (Outer-level) SGD step for meta-model, save per-task meta-weight for meta-update:

$$w_i^{t+1} := w^t - \beta_t \nabla_{w^t} \left( \hat{\mathcal{L}}(w_i^t, \mathcal{D}_i^{t, \text{out}}) + \sigma^{\text{out}} \text{Ordinary\_Reg}(w_i^t, \mathcal{D}_i^{t, \text{out}}) \right);$$

- 9:   **end for**
  - 10:   Meta-update  $w^{t+1} := \frac{1}{r} \sum_{i \in \mathcal{B}_t} w_i^{t+1}$
  - 11: **end for**
  - 12: **Return:**  $w^T$
- 

adaptation difficulty during training. Specifically, by making the adapted model more difficult to learn a generalized hypothesis by fitting the meta-support set, the meta-model is forced to learn better-generalized meta-knowledge to achieve good performance on the meta-query set. In this sense, we can think of the Minimax-Meta Regularization as a form of ‘‘adversarial training’’ for the meta-model, which can improve its generalization performance during training. Importantly, the ‘‘adversarial training’’ is only applied during the training phase and is not used during meta-testing. Thus, the meta-model does not carry the ‘‘adversarial training’’ burden in the actual deployment after learning better-generalized meta-knowledge, which can lead to better generalization in the new environment.

While the concept of using inverted regularization at the inner-level to improve generalization may seem too intuitional or counterintuitive to some, we provide a theoretical analysis in the next section to support its utility.

## 4.2. Application to MAML

To apply Minimax-Meta Regularization to MAML, we modify the MAML algorithm by adding the regularization to the inner- and outer-level training objective. The modified algorithm, which we refer to as Minimax-MAML, is shown in Algorithm 1. Note that this modification for Minimax-Meta Regularization is also generally applicable to other MAML variants.

## 5. Theoretical Analysis

In this section, we provide an analysis of the effectiveness of inverted regularization in meta-learning by taking L2-Norm regularization at the inner-level of the single-step MAML algorithm as a typical example, which is very possible to generalize to other regularization.

It is important to note that the process of adding regularization often involves changes to the loss function during training. This means if the model is obtained by a new regularized algorithm  $\tilde{\mathcal{A}}$ , it is usually optimized for a different function  $\tilde{F}(\cdot)$  instead of the original  $F(\cdot)$  (e.g., added weight-norm in the inner-level). However, in the meta-testing phase, the model’s test error is still calculated using  $F(\cdot)$ . As a result, to evaluate the test error change with a new regularized method  $\tilde{\mathcal{A}}$ , instead of directly adopting (4)’s decomposition in Preliminary, we need to further decompose the test error by

$$\begin{aligned}
 & \mathbb{E}_{\tilde{\mathcal{A}}, \mathcal{S}} \left[ F(\tilde{\mathcal{A}}(\mathcal{S})) - \min_{\mathcal{W}} F \right] \quad (\text{test error}) = \\
 & \underbrace{\mathbb{E}_{\tilde{\mathcal{A}}, \mathcal{S}} \left[ \hat{F}(\tilde{\mathcal{A}}(\mathcal{S}), \mathcal{S}) - \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) \right]}_{\text{training error}} \\
 & + \underbrace{\mathbb{E}_{\tilde{\mathcal{A}}, \mathcal{S}} [F(\tilde{\mathcal{A}}(\mathcal{S})) - \hat{F}(\tilde{\mathcal{A}}(\mathcal{S}), \mathcal{S})]}_{\text{generalization error}} \\
 & + \underbrace{\mathbb{E}_{\mathcal{S}} \left[ \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) \right] - \min_{\mathcal{W}} F}_{\leq 0} \\
 & + \underbrace{\mathbb{E}_{\mathcal{S}} \left[ \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) \right]}_{\text{training bias}}
 \end{aligned} \tag{5}$$

where  $\hat{F}(\cdot)$  refers to the regularized empirical loss function. (5) has one more training bias term compared to (4), which is caused by the changing of the objective function. Usually, regularization would reduce the expected generalization error while increasing the training bias. The goal of regularization is to decrease test error by reducing generalization error while trading off training bias.

Adding L2-Norm regularization at the inner-level for MAML could be obtained by changing the inner-level training objective from  $\hat{\mathcal{L}}(w^t, \mathcal{D}_i^{t, \text{in}})$  to  $(\hat{\mathcal{L}}(w^t, \mathcal{D}_i^{t, \text{in}}) + \frac{\delta}{2} \|w^t\|^2)$ , where  $\delta$  is the regularisation parameter. The meta updating rule would be accordingly changed from (3) to:

$$\begin{aligned}
 & w_i^{t+1} := \\
 & w^t - \beta_t \nabla_{w^t} \hat{\mathcal{L}} \left( w^t - \alpha \nabla_{w^t} \left( \hat{\mathcal{L}}(w^t, \mathcal{D}_i^{t, \text{in}}) + \frac{\delta}{2} \|w^t\|^2 \right), \mathcal{D}_i^{t, \text{out}} \right)
 \end{aligned} \tag{6}$$

Here  $\delta$  can be either positive or negative to represent the ordinary and inverted regularization, respectively. We treat

$\delta$  as a variable and analyze how its value would influence the generalization error and the training bias of the total error introduced in (5).

The analysis of generalization error closely follows the work of [9], and holds the same assumptions about function  $\ell(\cdot, z)$  and task distribution as follows.

**Assumption 1.** We assume the function  $\ell(\cdot, z)$  satisfies the following properties for any  $z \in \mathcal{Z}$ :

1. (Strong convexity)  $\ell(\cdot, z)$  is  $\mu$ -strongly convex, i.e.,  $(\nabla\ell(w, z) - \nabla\ell(u, z))^T(w - u) \geq \mu\|w - u\|^2$ ;
2. (Lipschitz in function value)  $\ell(\cdot, z)$  has gradients with norm bounded by  $G$ , i.e.,  $\|\nabla\ell(w, z)\| \leq G$ ;
3. (Lipschitz gradient)  $\ell(\cdot, z)$  is  $L$ -smooth, i.e.,  $\|\nabla\ell(w, z) - \nabla\ell(u, z)\| \leq L\|w - u\|$ ;
4. (Lipschitz Hessian)  $\ell(\cdot, z)$  has  $\rho$ -Lipschitz Hessian, i.e.,  $\|\nabla^2\ell(w, z) - \nabla^2\ell(u, z)\| \leq \rho\|w - u\|$

**Assumption 2.** We assume  $\mathcal{F}_{\mathcal{Z}}$  is the Borel  $\sigma$ -algebra over  $\mathcal{Z}$  and  $\mathcal{Z}$  is a Polish space. And each  $p_i$  is a non-atomic distribution over  $(\mathcal{Z}, \mathcal{F}_{\mathcal{Z}})$

## 5.1. Generalization Error

We derive our generalization bound for MAML with L2 regularization at the inner-level through the theoretical framework proposed by [9], which mainly adopts an algorithmic stability approach for the derivation. We denote the algorithm combines MAML with inner-level regularization as  $\tilde{\mathcal{A}}$ , and the below generalization bound could be obtained. We provide detailed proof in Appendix.

**Theorem 1** (generalization bound). *If Assumption 1 and 2 hold. With  $\alpha \leq \frac{1}{2L}$ ,  $\beta_t \leq \frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L}$ ,  $\delta < \frac{1}{2\alpha}$  and  $\frac{\alpha\rho G}{\mu} < (\frac{1}{2} - \alpha L)^2$ . The model  $\tilde{\mathcal{A}}(\mathcal{S})$  generated by the last iterate of MAML with regularized updating rule introduced in (6) satisfies*

$$\mathbb{E}_{\tilde{\mathcal{A}}, \mathcal{S}}[F(\tilde{\mathcal{A}}(\mathcal{S})) - \hat{F}(\tilde{\mathcal{A}}(\mathcal{S}), \mathcal{S})] \leq \frac{2G^2(1 + \alpha L)(1 - \alpha\mu - \alpha\delta + (2 + \alpha L - \alpha\delta)\alpha LK)}{mn} \left( \frac{1}{\alpha\rho G + (1 - \alpha\mu - \alpha\delta)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2 \mu} \right)$$

where the expectation is taken over the randomness of  $\tilde{\mathcal{A}}$  and sampling of  $\mathcal{S}$ .

The generalization bound could be regarded as a function  $GB(\delta)$ , and its derivative  $GB'(\delta)$  is positive  $\forall \delta \in (-\infty, \frac{1}{2\alpha})^1$ . It suggests that  $GB(\delta)$  is monotonically increasing if  $\delta \in (-\infty, \frac{1}{2\alpha})$ , implying that L2 regularization at the inner-level decreases the generalization bound of MAML only when it's inverted (i.e.  $\delta < 0$ ). And ordinary regularization (i.e.  $\delta \in (0, \frac{1}{2\alpha})$ ) at the inner-level would increase the generalization bound.

<sup>1</sup> Derivation is included in Appendix A.3.1.  $\delta \geq \frac{1}{2\alpha}$  are excluded from discussion because they may break the convexity of the meta loss function.

## 5.2. Training Bias

**Theorem 2** (training bias bound). *If Assumption 1 and 2 hold. With  $\alpha \leq \frac{1}{2L}$ ,  $\delta < \frac{1}{2\alpha}$  and  $\frac{\alpha\rho G}{\mu} < (\frac{1}{2} - \alpha L)^2$ . The training bias from MAML with inner-level L2 regularization to the original MAML is bounded by*

$$\mathbb{E}_{\mathcal{S}} \left[ \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) \right] \leq \frac{\alpha^2(\alpha\rho G + (1 - \alpha\mu)^2 L)((1 - \alpha\mu - \alpha\delta)L\|w^*\| + G)^2 \delta^2}{2(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2 \mu)^2}$$

where  $\|w^*\| := \max_{\mathcal{S}} \|\arg \min_{\mathcal{W}} \hat{F}(w, \mathcal{S})\|$ , the maximum is taken over sampling of  $\mathcal{S}$ .

The training bias bound could also be regarded as a function  $TB(\delta)$ . We could observe that  $TB(\delta) > TB(0) = 0$  for  $\delta \neq 0$ , which suggests that training bias is inevitable when regularization is adopted. Another important finding is that for any legal choice of  $\delta_0 > 0$ , we have  $TB(-\delta_0) < TB(\delta_0)$ <sup>2</sup>, which suggests that the inverted regularization has less corruption to training bias bound at the inner-level than the ordinary regularization with the same coefficient.

## 5.3. Test Error

Since the training error term in the test error (5) vanishes with iteration  $T$  as long as the outer-level loss is strongly-convex [9], the training error term could be negligible for  $\delta < \frac{1}{2\alpha}$ . So we could just consider the training bias and generalization error for bounding the test error, i.e.,

$$\mathbb{E}_{\tilde{\mathcal{A}}, \mathcal{S}} \left[ F(\tilde{\mathcal{A}}(\mathcal{S})) - \min_{\mathcal{W}} F \right] \leq \frac{2G^2(1 + \alpha L)(1 - \alpha\mu - \alpha\delta + (2 + \alpha L - \alpha\delta)\alpha LK)}{mn} \underbrace{\left( \frac{1}{\alpha\rho G + (1 - \alpha\mu - \alpha\delta)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2 \mu} \right)}_{\text{generalization error bound } GB(\delta)} + \underbrace{\frac{\alpha^2(\alpha\rho G + (1 - \alpha\mu)^2 L)((1 - \alpha\mu - \alpha\delta)L\|w^*\| + G)^2 \delta^2}{2(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2 \mu)^2}}_{\text{training bias bound } TB(\delta)}$$

The test error bound could be described by  $TE(\delta) := TB(\delta) + GB(\delta)$ . When  $\delta$  is positive, we have  $TB(\delta) > TB(0)$  and  $GB(\delta) > GB(0)$  (since  $GB'(\delta) > 0 \forall \delta \in (-\infty, \frac{1}{2\alpha})$ ), which suggests ordinary regularization at the inner-level worsens the model's test error bound. Instead, for inverted regularization, since  $TE'(0) = TB'(0) + GB'(0) = 0 + GB'(0) > 0$ , there must be an interval  $[\delta^*, 0)$  in which all values can be used as the inverted regularization parameter to decrease the test error bound.

<sup>2</sup> Derivation is included in Appendix A.3.2

## 6. Experiments

We conduct extensive experiments on three types of classical meta-learning tasks: few-shot classification, few-shot regression, and robust reweighting. The experiments include: **i)** an empirical verification of the regularization at inner- and outer-level on the Mini-Imagenet few-shot classification task, which demonstrates the effectiveness of both the inverted regularization at inner-level and the ordinary regularization at outer-level; **ii)** further experiments on few-shot classification and regression benchmarks to compare our Minimax-Meta regularized algorithms with other representative methods; **iii)** a few-shot learning experiment on a limited number of tasks evaluating generalization of different regularization strategies; and **iv)** an experiment on meta-reweighting for robust learning, which demonstrates the broad applicability of our method to different meta-learning problems. (*Due to page-size limitations, the experiments on limited tasks, Meta-Dataset with larger backbones, and meta-reweighting are included in the Appendix*)

### 6.1. Few-shot Classification

We first conduct experiments on the few-shot classification task, one of the most popular tasks to evaluate meta-learning algorithms. To verify the effectiveness of our approach, we adapt Minimax-Meta Regularization into bi-level optimization meta-learning algorithms and make a benchmark to compare with other methods.

#### 6.1.1 Experimental Setup

**Datasets.** For the few-shot classification task, we experiment on the Mini-Imagenet [38, 49] and Omniglot [23] datasets. The Mini-Imagenet [38] is sampled from ImageNet with 600 instances of 100 classes. In the experiment, the Mini-Imagenet dataset is split into 64 classes for training, 12 classes for validation, and 24 classes for testing. The Omniglot dataset is a collection of 1623 character classes with different alphabets. Each class in the dataset contains 20 instances. The classes are shuffled and divided into the training, validation, and test sets, with 1150, 50, and 423 instances in the experiment.

**Experimental details.** We select MAML [10] as the representative bi-level optimization meta-learning algorithm to conduct the experiment. The few-shot benchmark settings for Omniglot and Mini-ImageNet experiments provided in [4] are adopted for our experiment build. Details about the experiment can be found in Appendix.

To verify the theoretical results and show the effectiveness of our regularization design, we first conduct an empirical verification experiment on Mini-ImageNet using L2-Norm as the regularizer.

In other few-shot classification experiments, we use a combination of L2-Norm and output-entropy as the regu-

larizer to further improve the generalization. (Although we only use L2-Norm as the sample regularizer to derive the theoretical results in Section 5, the use of inverted regularization can cover many other regularizers in practice, including the entropy regularizer.) That is, in this part of the experiment, when we say that "adding ordinary regularization" at a certain level, its corresponding learning objective will include minimizing the L2-Norm of model weights and maximizing the entropy of the model's output prediction (improves generalization); when we say that "adding inverted regularization" at a certain level, its corresponding learning objective will include maximizing of the L2-Norm of model weights and minimizing the entropy of the model's output prediction (hinders generalization). And we keep the magnitude of the L2-Norm parameter =  $5e-4$  and the magnitude of the information entropy parameter = 2.0 across the experiments, i.e., the difference between ordinary and inverted regularization in this group of few-shot learning experiments is only the sign of the regularization term. Note that we only add regularization at the training phase, so the inner-levels are not regularized in meta-testing time.

#### 6.1.2 Empirical Verification for regularization at inner- and outer-level.

To verify our view that the regularization at the inner- and outer-level should respectively be inverted and ordinary, we conduct two experiments for MAML [10] with different regularization methods on Mini-Imagenet 5-way few-shot problem. There are five regularization methods being compared: *no regularization*, *regularize the outer-level*, *regularize the inner-level*, *invertedly regularize the inner-level*, and *Minimax-Meta Regularization*. In the first experiment, We only use L2-Norm regularization to match the setting of theoretical analysis. In the second experiment, We use L2-Norm & entropy combined regularization to verify whether inverted inner-level regularization is suitable for different types of regularizers and whether a combination of multiple regularizers leads to better generalization. We follow [4]'s setting to build the experiment with 48-48-48-48 conv backbone and use the ensemble of per-epoch models to generate more stable results (MAML baseline achieves higher performance under this setting compared to classic 32-32-32-32 conv backbone implementations [10]), The results are respectively presented in Table 1 and 2. Based on the results, we make the following observations:

*Inner-level inverted regularization enhances the generalization performance.* Compare the results from "no regularization" and "invertedly regularize the inner-level", we observe that adding inner inverted regularization achieves accuracy improvements in both 1-shot and 5-shot experiments, which verifies the efficacy of the inner inverted regularization. This is aligned with our intuition and theoretical

Table 1. Test accuracy of MAML with different types of regularization in the Mini-Imagenet 5-way MAML Few-shot Classification experiment (*L2-Norm as regularization objective only*). Backbone: 48-48-48-48 conv. We report the test accuracy with a 95% confidence interval for the mean.

Mini-Imagenet 5-way Few-shot Classification for MAML ( <i>Reg Objective: L2-Norm</i> )				
Regularization Type	Outer Reg	Inner Reg	1-Shot	5-Shot
<i>no regularization</i>	-	-	49.58±0.45%	65.39±0.50%
<i>regularize the outer-level</i>	<i>Ordinary</i>	-	49.90±0.54%	66.47±1.21%
<i>regularize the inner-level</i>	-	<i>Ordinary</i>	49.28±0.37%	64.80±0.25%
<i>invertedly regularize the inner-level</i>	-	<i>Inverted</i>	49.92±0.42%	66.05±0.68%
<i>Minimax-Meta Regularization</i>	<i>Ordinary</i>	<i>Inverted</i>	<b>50.25±0.38%</b>	<b>68.17±0.92%</b>

Table 2. Test accuracy of MAML with different types of regularization in the Mini-Imagenet 5-way MAML Few-Shot Classification experiment (*Combining L2-Norm and output entropy as regularization objective*). Backbone: 48-48-48-48 conv. We report the test accuracy with a 95% confidence interval for the mean.

Mini-Imagenet 5-way Few-Shot Classification for MAML ( <i>Reg Objective: L2-Norm &amp; Entropy</i> )				
Regularization Type	Outer Reg	Inner Reg	1-Shot	5-Shot
<i>no regularization</i>	-	-	49.58±0.45%	65.39±0.50%
<i>regularize the outer-level</i>	<i>Ordinary</i>	-	50.23±0.67%	67.18±0.88%
<i>regularize the inner-level</i>	-	<i>Ordinary</i>	48.07±1.01%	64.32±0.35%
<i>invertedly regularize the inner-level</i>	-	<i>Inverted</i>	49.96±0.33%	65.91±0.41%
<i>Minimax-Meta Regularization</i>	<i>Ordinary</i>	<i>Inverted</i>	<b>50.85±0.37%</b>	<b>69.36±0.34%</b>

result.

*Inner-level ordinary regularization impairs the generalization performance.* Compare the results from “no regularization” and “regularize the inner-level”, we observe that adding inner ordinary regularization suffers from accuracy impairments. This observation is also consistent with our intuition and theoretical findings.

*Outer-level ordinary regularization enhances the generalization performance.* Compare the results from “no regularization” and “regularize the outer-level”, we observe that adding outer regularization can get accuracy improvements, which verifies the efficacy of adding ordinary regularization at the outer-level.

*The outer-level ordinary regularization and inner-level inverted regularization are compatible.* We observe that Minimax-Meta Regularization outperforms solely outer-level or inverted inner-level regularization, indicating compatibility between the regularizations at the two distinct levels. This aligns with the intuition that meta and adaptation generalization are not conflicting.

*Inner-level inverted regularization and the outer-level ordinary regularization are suitable for combined regularizer* We observe consistent effects across L2-Norm regularizer and L2-Norm & entropy combined regularizer when using different regularization strategies. Furthermore, combining the L2-Norm and entropy regularizer led to improved performance compared to using L2-Norm regularizer alone.

### 6.1.3 Minimax-Meta Regularization for Few-shot Classification

So far, we have proved that Minimax-Meta Regularization is a promising regularization strategy for bi-level meta-learning. Here, we do experiments to further test the effectiveness of Minimax-Meta Regularization.

The experiments are conducted on Omniglot and Mini-ImageNet datasets. We implement Minimax-Meta Regularization for bi-level meta-learning algorithms: *MAML* [10], which is the most representative bi-level meta-learning algorithm; *MAML++* [4], which is an adapted version of MAML with additional techniques for performance improvements. L2-Norm & entropy combined regularizer is adopted in this experiment.

Representative algorithms with comparable backbone structures are selected for making the comparison. We use the 64-64-64-64 conv backbone for the Mini-ImageNet experiment to make a fairer comparison with other methods. The results are shown in Table 3 and 4.

The results suggest that Minimax-Meta Regularization generally improves test performances. Minimax-MAML++ achieves the best performance on both datasets.

## 6.2. Minimax-Meta Regularization for Few-shot Regression

We then conduct experiments on the few-shot regression task to test the efficacy of Minimax-Meta Regularization.

Table 3. Omniglot 20-way 1-shot experiment. We report the test accuracy with a 95% confidence interval for the mean.

*the \* indicates result generated in our experiment.*

Omniglot 20-way 1-Shot Classification	
	Accuracy
Meta-SGD [28]	95.93±0.38%
Prototypical Net [44]	96.00%
Meta-Networks [32]	97.00%
GNN [16]	97.40%
Relation Network [46]	97.60±0.20%
R2-D2 [5]	96.24±0.05%
SNAIL [31]	97.64±0.30%
TAML(Entropy) [20]	95.62±0.50%
MAML [10]*	94.20±0.41%
<b>Minimax-MAML(ours)*</b>	<b>95.76±0.39%</b>
MAML++ [4]*	97.21±0.51%
<b>Minimax-MAML++(ours)*</b>	<b>97.77±0.06%</b>

## 6.2.1 Experimental Setup

**Datasets.** We follow the few-shot regression experiment setting proposed in [40] to build the experiment. One synthetic and three real-world few-shot regression datasets are included. The synthetic dataset is created by a 2-dimensional mixture of Cauchy distributions plus random GP functions. One real-world dataset is SwissFEL [30] which corresponds to Swiss Free Electron Laser’s calibration sessions. Another two datasets are from the PhysioNet 2012 challenge [43], which contains time-series data related to patients’ health metrics, in particular, the Glasgow Coma Scale (GCS) and the hematocrit value (HCT).

**Experimental details.** We implement Minimax-MAML for the regression task by adding inverted and ordinary L2-Norm at the inner-level and outer-level of MAML, respectively. To obtain optimal results, unlike the single-inner-step MAML implemented in [40], we perform three inner update steps for the meta-training of Minimax-MAML. In order to verify the effect of minimax, we also compared the results of unregularized MAML with three inner steps.

## 6.2.2 Experimental Results

As shown in Table 5, the Minimax-Meta Regularization improved the performance in all four datasets. And Minimax-MAML achieves near-best performance on the synthetic Cauchy datasets and outperforms other algorithms on the two Physionet datasets. The results suggest that the Minimax-Meta Regularization could improve the performance of the few-shot regression task for meta-learning.

## 7. Conclusion

This paper studies the generalization problem of bi-level optimization-based meta-learning. While most of the exist-

Table 4. Mini-Imagenet 5-way few-shot experiment. We report the test accuracy with a 95% confidence interval for the mean.

*the \* indicates result generated in our experiment.*

Mini-Imagenet 5-way Few-Shot Classification			
Approach	Backbone	1-Shot Accuracy	5-Shot Accuracy
Meta-SGD [28]	64-64-64-64	50.47±1.87%	64.03±0.94%
Prototypical Nets [44]	64-64-64-64	49.42±0.78%	68.20±0.66%
GNN [16]	64-96-128-256	50.33±0.36%	66.41±0.63%
R2-D2 [5]	64-64-64-64	49.50±0.20%	65.40±0.20%
LR-D2 [5]	96-192-384-512	51.90±0.20%	68.70±0.20%
MetaOptNet [25]	64-64-64-64	53.23±0.59%	69.51±0.48%
TAML(Entropy) [20]	64-64-64-64	51.73±1.88%	66.05±0.85%
MAML-Meta Dropout [24]	32-32-32-32	51.93±0.67%	67.42±0.52%
MAML-MMCF [51]	32-32-32-32	50.35±1.82%	64.91±0.96%
MAML [10]*	64-64-64-64	50.20±1.65%	65.86±0.61%
<b>Minimax-MAM(ours)*</b>	64-64-64-64	<b>51.70±0.42%</b>	<b>68.41±1.28%</b>
MAML++ [4]*	64-64-64-64	52.96±0.78%	70.02±0.55%
<b>Minimax-MAML++(ours)*</b>	64-64-64-64	<b>53.28±0.35%</b>	<b>71.70±0.23%</b>

Table 5. Test RMSE comparison of algorithms in four meta-learning environments for few-shot regression.

*the \* indicates the result generated in our experiment, other results are reported from [40]*

	Cauchy	SwissFel	Physionet-GCS	Physionet-HCT
MLL-GP [14]	0.216±0.003	0.974±0.093	1.654±0.094	2.634±0.144
MLAP [3]	0.219±0.004	0.486±0.026	2.009±0.248	2.470±0.039
NP [17]	0.224±0.008	0.471±0.053	2.056±0.209	2.594±0.107
PACOH-GP [40]	0.209±0.008	0.376±0.024	1.498±0.081	2.361±0.047
PACOH-NN [40]	<b>0.195±0.001</b>	<b>0.372±0.002</b>	1.561±0.061	2.405±0.017
MAML [10](1 inner step)	0.219±0.004	0.730±0.057	1.895±0.141	2.413±0.113
MAML [10](3 inner steps)*	0.212±0.003	0.535±0.042	1.532±0.074	2.396±0.047
<b>Minimax-MAML*</b>	0.201±0.002	0.477±0.026	<b>1.483±0.052</b>	<b>2.343±0.019</b>

ing works focus on meta-generalization to unseen tasks at the meta-level, they leave out that adapted models may not be generalized to the task domain at the adaptation-level. We give an intuitive explanation of why the inverted regularization at the inner-level could improve the adaptation generalization of meta-learning. We provide theoretical support for this intuition by deriving generalization error and training bias bound. We empirically verify that both *inverted regularization at inner-level* and *ordinary regularization at outer-level* improve the test performance of meta-learning. Based on the aligned theoretical and empirical results, we propose meta-learning with Minimax-Meta Regularization, combining regularization at inner- and outer-level. Finally, we conduct experiments on multiple meta-learning tasks to show the efficacy of the proposed method.

## Acknowledgements

This work was supported in part by the National Key Research and Development Program of China No. 2020AAA0106300, National Natural Science Foundation of China (No. 62250008, 62222209, 62102222, 61936011). Shanghang would like to thank the support from CCF-DiDi GAIA Collaborative Research Funds for Young Scholars. We gratefully thank Jiaming Liu, Yuzhou Cao and Wenpeng Zhang for their valuable discussions.



## References

- [1] Alessandro Achille and Stefano Soatto. Information dropout: Learning optimal representations through noisy computation. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2897–2905, 2018. [2](#)
- [2] Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. Deep variational information bottleneck. *arXiv preprint arXiv:1612.00410*, 2016. [2](#)
- [3] Ron Amit and Ron Meir. Meta-learning by adjusting priors based on extended pac-bayes theory. In *International Conference on Machine Learning*, pages 205–214. PMLR, 2018. [8](#)
- [4] Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *arXiv preprint arXiv:1810.09502*, 2018. [6](#), [7](#), [8](#)
- [5] Luca Bertinetto, Joao F Henriques, Philip Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. In *International Conference on Learning Representations*, 2019. [8](#)
- [6] Liam Collins, Aryan Mokhtari, and Sanjay Shakkottai. Task-robust model-agnostic meta-learning. *arXiv preprint arXiv:2002.04766*, 2020. [2](#)
- [7] Zi-Yi Dou, Keyi Yu, and Antonios Anastasopoulos. Investigating meta-learning algorithms for low-resource natural language understanding tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1192–1197, 2019. [1](#)
- [8] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel.  $R1^2$ : Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016. [2](#)
- [9] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Generalization of model-agnostic meta-learning algorithms: Recurring and unseen tasks. *Advances in Neural Information Processing Systems*, 34, 2021. [2](#), [3](#), [5](#)
- [10] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, pages 1126–1135. PMLR, 2017. [1](#), [2](#), [6](#), [7](#), [8](#)
- [11] Chelsea Finn and Sergey Levine. Meta-learning and universality: Deep representations and gradient descent can approximate any learning algorithm. In *International Conference on Learning Representations*, 2018. [2](#)
- [12] Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 9537–9548, 2018. [2](#)
- [13] Sebastian Flennerhag, Andrei Rusu, Razvan Pascanu, Francisco Visin, Hujun Yin, and Raia Hadsell. Meta-learning with warped gradient descent. In *International Conference on Learning Representations 2020*, 2020. [2](#)
- [14] Vincent Fortuin and Gunnar Rätsch. Deep mean functions for meta-learning in gaussian processes. *arXiv preprint arXiv:1901.08098*, 2019. [8](#)
- [15] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016. [2](#)
- [16] Victor Garcia and Joan Bruna. Few-shot learning with graph neural networks. In *6th International Conference on Learning Representations, ICLR 2018*, 2018. [8](#)
- [17] Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye Teh. Neural processes. *arXiv preprint arXiv:1807.01622*, 2018. [8](#)
- [18] Erin Grant, Chelsea Finn, Sergey Levine, Trevor Darrell, and Thomas Griffiths. Recasting gradient-based meta-learning as hierarchical bayes. In *International Conference on Learning Representations*, 2018. [2](#)
- [19] Mike Huisman, Jan N Van Rijn, and Aske Plaat. A survey of deep meta-learning. *Artificial Intelligence Review*, 54(6):4483–4541, 2021. [1](#)
- [20] Muhammad Abdullah Jamal and Guo-Jun Qi. Task agnostic meta-learning for few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11719–11727, 2019. [1](#), [2](#), [8](#)
- [21] Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, and Trevor Darrell. Few-shot object detection via feature reweighting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8420–8429, 2019. [1](#)
- [22] Anders Krogh and John A Hertz. A simple weight decay can improve generalization. In *Advances in neural information processing systems*, pages 950–957, 1992. [2](#)
- [23] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015. [6](#)
- [24] Hae Beom Lee, Taewook Nam, Eunho Yang, and Sung Ju Hwang. Meta dropout: Learning to perturb latent features for generalization. In *International Conference on Learning Representations*, 2020. [8](#)
- [25] Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex optimization. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10649–10657. IEEE Computer Society, 2019. [8](#)
- [26] Yoonho Lee and Seungjin Choi. Gradient-based meta-learning with learned layerwise metric and subspace. In *International Conference on Machine Learning*, pages 2927–2936. PMLR, 2018. [2](#)
- [27] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S Kankanhalli. Learning to learn from noisy labeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5051–5059, 2019. [1](#)
- [28] Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few-shot learning. *arXiv preprint arXiv:1707.09835*, 2017. [1](#), [2](#), [8](#)
- [29] Andrea Madotto, Zhaojiang Lin, Chien-Sheng Wu, and Pascale Fung. Personalizing dialogue agents via meta-learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5459, 2019. [1](#)

- [30] Christopher J Milne, Thomas Schietinger, Masamitsu Aiba, Arturo Alarcon, Jürgen Alex, Alexander Anghel, Vladimir Arsov, Carl Beard, Paul Beaud, Simona Bettoni, et al. Swissfel: the swiss x-ray free electron laser. *Applied Sciences*, 7(7):720, 2017. 8
- [31] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. In *International Conference on Learning Representations*, 2018. 2, 8
- [32] Tsendsuren Munkhdalai and Hong Yu. Meta networks. In *International Conference on Machine Learning*, pages 2554–2563. PMLR, 2017. 8
- [33] Renkun Ni, Micah Goldblum, Amr Sharaf, Kezhi Kong, and Tom Goldstein. Data augmentation for meta-learning. In *International Conference on Machine Learning*, pages 8152–8161. PMLR, 2021. 1, 2
- [34] Boris N Oreshkin, Pau Rodriguez, and Alexandre Lacoste. Tadam: task dependent adaptive metric for improved few-shot learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 719–729, 2018. 2
- [35] Eunbyung Park and Junier B Oliva. Meta-curvature. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 3314–3324, 2019. 2
- [36] Janarthanan Rajendran, Alexander Irpan, and Eric Jang. Meta-learning requires meta-augmentation. In *Advances in Neural Information Processing Systems*, pages 5705–5715, 2020. 1
- [37] Janarthanan Rajendran, Alex Irpan, and Eric Jang. Meta-learning requires meta-augmentation. *arXiv preprint arXiv:2007.05549*, 2020. 2
- [38] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017. 6
- [39] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *International Conference on Machine Learning*, pages 4334–4343. PMLR, 2018. 1
- [40] Jonas Rothfuss, Vincent Fortuin, Martin Josifoski, and Andreas Krause. Pacoh: Bayes-optimal meta-learning with pac-guarantees. In *International Conference on Machine Learning*, pages 9116–9126. PMLR, 2021. 2, 8
- [41] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*, pages 1842–1850. PMLR, 2016. 2
- [42] Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, and Deyu Meng. Meta-weight-net: Learning an explicit mapping for sample weighting. *Advances in Neural Information Processing Systems*, 32:1919–1930, 2019. 1
- [43] Ikaro Silva, George Moody, Daniel J Scott, Leo A Celi, and Roger G Mark. Predicting in-hospital mortality of icu patients: The physionet/computing in cardiology challenge 2012. In *2012 Computing in Cardiology*, pages 245–248. IEEE, 2012. 8
- [44] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 4080–4090, 2017. 2, 8
- [45] Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning for few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 403–412, 2019. 1
- [46] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1199–1208, 2018. 1, 8
- [47] Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer Science & Business Media, 2012. 1
- [48] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015. 2
- [49] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29:3630–3638, 2016. 2, 6
- [50] Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM Computing Surveys (CSUR)*, 53(3):1–34, 2020. 1
- [51] Huaxiu Yao, Long-Kai Huang, Linjun Zhang, Ying Wei, Li Tian, James Zou, Junzhou Huang, et al. Improving generalization in meta-learning via task augmentation. In *International Conference on Machine Learning*, pages 11887–11897. PMLR, 2021. 1, 2, 8
- [52] Mingzhang Yin, George Tucker, Mingyuan Zhou, Sergey Levine, and Chelsea Finn. Meta-learning without memorization. *arXiv preprint arXiv:1912.03820*, 2019. 1, 2
- [53] Jaesik Yoon, Taesup Kim, Ousmane Dia, Sungwoong Kim, Yoshua Bengio, and Sungjin Ahn. Bayesian model-agnostic meta-learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 7343–7353, 2018. 2
- [54] Tianhe Yu, Chelsea Finn, Annie Xie, Sudeep Dasari, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot imitation from observing humans via domain-adaptive meta-learning. *arXiv preprint arXiv:1802.01557*, 2018. 1

# Supplementary Material for “Improving Generalization of Meta-learning with Inverted Regularization at Inner-level”

Due to the space limitation of the main paper, we provide supplementary theoretical proof and supplementary experimental results in this Appendix, including: more detailed theoretical analyses, more experiment details, an additional experiment on Mini-ImageNet few-shot classification with limited tasks, an additional experiment on Meta-dataset with the first-order method and larger backbone, and an additional experiment on meta-reweighting with Minimax-Meta Regularization for robust learning.

## A. Theoretical Analysis

In this section, we provide detailed proof derivations of the theoretical results in the main paper.

### A.1. Lemmas

This section lists the Lemmas that help prove our main results.

**Lemma 1.** (from [6]) Let  $\phi$  be a  $\lambda$ -strongly convex and  $\eta$ -smooth function. Then, for any  $\beta \leq \frac{2}{\lambda+\eta}$ , we have

$$\|(u - \beta \nabla \phi(u)) - (v - \beta \nabla \phi(v))\| \leq \left(1 - \frac{\beta \lambda \eta}{\lambda + \eta}\right) \|u - v\|$$

for any  $u$  and  $v$ .

**Lemma 2.** (based on [5]) Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be an function that is  $L$ -smooth,  $\mu$ -strongly convex, and has gradient bounded by  $G$ . Consider a function  $U(\cdot)$  that describes the MAML inner-level update rule, with L2-Norm regularization parameterized by  $\frac{\delta}{2} : U(\mathbf{w}) = \mathbf{w} - \alpha \nabla_{\mathbf{w}}(f(\mathbf{w}) + \frac{\delta}{2} \|\mathbf{w}\|^2)$ , with  $\alpha \leq \frac{1}{2L}$ ,  $\delta < \frac{1}{2\alpha}$ . Then,

$$\|U(\mathbf{w}) - U(\mathbf{v})\| \leq (1 - \alpha\delta - \alpha\mu) \|\mathbf{w} - \mathbf{v}\| \quad \forall \mathbf{w}, \mathbf{v} \in \mathbb{R}^d$$

*Proof.* Firstly, note that

$$U(\mathbf{w}) = \mathbf{w} - \alpha(\nabla f(\mathbf{w}) + \delta \mathbf{w}) = (1 - \alpha\delta)\mathbf{w} - \alpha \nabla f(\mathbf{w})$$

The Jacobian of  $U(\cdot)$  is given by  $\nabla U(\mathbf{w}) = (1 - \alpha\delta)\mathbf{I} - \alpha \nabla^2 f(\mathbf{w})$ .

Like in [5], we use  $\mathbf{A} \succeq \mathbf{0}$  to denote the positive semi-definite nature of the matrix. Similarly,  $\mathbf{A} \succeq \mathbf{B}$  means that  $\boldsymbol{\theta}^T(\mathbf{A} - \mathbf{B})\boldsymbol{\theta} \geq 0 \forall \boldsymbol{\theta}$ . Since  $f$  is  $\mu$ -strongly convex, and  $L$ -smooth, we could have  $\mu\mathbf{I} \preceq \nabla^2 f(\mathbf{w}) \preceq L\mathbf{I} \quad \forall \mathbf{w} \in \mathbb{R}^d$ . Then the Jacobian can be bounded by

$$(1 - \alpha\delta - \alpha L)\mathbf{I} \preceq \nabla U(\mathbf{w}) \preceq (1 - \alpha\delta - \alpha\mu)\mathbf{I} \quad \forall \mathbf{w} \in \mathbb{R}^d$$

The upper bound implies  $\|\nabla U(\mathbf{w})\| \leq (1 - \alpha\mu - \alpha\delta) \|\mathbf{w}\| \quad \forall \mathbf{w} \in \mathbb{R}^d$ .

Let  $\psi(t) = \mathbf{v} + t(\mathbf{w} - \mathbf{v})$ ,  $t \in [0, 1]$  be the line function connecting  $\mathbf{w}$  and  $\mathbf{v}$ . Taking the line integral, we have

$$\begin{aligned} U(\mathbf{w}) - U(\mathbf{v}) &= \int_{t=0}^{t=1} \nabla U(\psi(t)) d\psi(t) \\ &= \int_{t=0}^{t=1} \nabla U(\psi(t)) \frac{d\psi(t)}{dt} dt \\ &= \int_{t=0}^{t=1} \nabla U(\psi(t)) (\mathbf{w} - \mathbf{v}) dt \\ &= \left( \int_{t=0}^{t=1} \nabla U(\psi(t)) dt \right) (\mathbf{w} - \mathbf{v}) \end{aligned}$$

Using the Cauchy-Schwartz inequality and  $\|\nabla U(\mathbf{w})\| \leq (1 - \alpha\mu - \alpha\delta) \forall \mathbf{w}$ , we have

$$\begin{aligned}
\|U(\mathbf{w}) - U(\mathbf{v})\| &= \left\| \int_{t=0}^{t=1} \nabla U(\psi(t))(\mathbf{w} - \mathbf{v}) dt \right\| \\
&\leq \int_{t=0}^{t=1} \|\nabla U(\psi(t))(\mathbf{w} - \mathbf{v})\| dt \\
&\leq \int_{t=0}^{t=1} \|\nabla U(\psi(t))\| \|\mathbf{w} - \mathbf{v}\| dt \\
&\leq \int_{t=0}^{t=1} (1 - \alpha\mu - \alpha\delta) \|\mathbf{w} - \mathbf{v}\| dt \\
&= (1 - \alpha\mu - \alpha\delta) \|\mathbf{w} - \mathbf{v}\| \int_{t=0}^{t=1} dt \\
&= (1 - \alpha\mu - \alpha\delta) \|\mathbf{w} - \mathbf{v}\|
\end{aligned}$$

□

## A.2. Main Results

We start by restating the assumptions we use to derive the results, and then we move on to prove our main results.

**Assumption 1.** We assume the function  $\ell(\cdot, z)$  satisfies the following properties for any  $z \in \mathcal{Z}$ :

1. (Strong convexity)  $\ell(\cdot, z)$  is  $\mu$ -strongly convex, i.e.,  $(\nabla \ell(w, z) - \nabla \ell(u, z))^T(w - u) \geq \mu \|w - u\|^2$ ;
2. (Lipschitz in function value)  $\ell(\cdot, z)$  has gradients with norm bounded by  $G$ , i.e.,  $\|\nabla \ell(w, z)\| \leq G$ ;
3. (Lipschitz gradient)  $\ell(\cdot, z)$  is  $L$ -smooth, i.e.,  $\|\nabla \ell(w, z) - \nabla \ell(u, z)\| \leq L \|w - u\|$ ;
4. (Lipschitz Hessian)  $\ell(\cdot, z)$  has  $\rho$ -Lipschitz Hessian, i.e.,  $\|\nabla^2 \ell(w, z) - \nabla^2 \ell(u, z)\| \leq \rho \|w - u\|$

**Assumption 2.** We assume  $\mathcal{F}_{\mathcal{Z}}$  is the Borel  $\sigma$ -algebra over  $\mathcal{Z}$  and  $\mathcal{Z}$  is a Polish space. And each  $p_i$  is a non-atomic distribution over  $(\mathcal{Z}, \mathcal{F}_{\mathcal{Z}})$ .

### A.2.1 Strongly Convexity and Smoothness

**Lemma 3.** (based on [5]) Suppose  $f$  and  $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfy assumptions 1. We formulate the MAML's outer-level evaluation function with inner-level L2-Norm regularization parameterized by  $\frac{\delta}{2}$  with  $f$  and  $\hat{f}$ , and let  $\tilde{f}$  be the function evaluated after a one-step gradient update procedure, i.e.,

$$\tilde{f}(w) := f(w - \alpha \nabla_w(\hat{f}(w) + \frac{\delta}{2} \|w\|^2))$$

then, with  $\alpha < \frac{1}{2L}$ ,  $\delta < \frac{1}{2\alpha}$  and  $\frac{\alpha\rho G}{\mu} < (\frac{1}{2} - \alpha L)^2$ ,  $\tilde{f}(\cdot)$  is  $(-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2\mu)$  strongly convex and  $(\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2L)$  smooth.

*Proof.* Let  $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$  be two arbitrary points. Let  $U(\mathbf{w}) = w - \alpha \nabla_w(\hat{f}(w) + \frac{\delta}{2} \|w\|^2)$ . Note that

$$\begin{aligned}
U(\mathbf{w}) &= w - \alpha(\nabla \hat{f}(w) + \delta w) \\
&= (1 - \alpha\delta)w - \alpha \nabla \hat{f}(w)
\end{aligned}$$

We use shorthand of  $\tilde{\mathbf{w}} \equiv U(\mathbf{w})$ ,  $\tilde{\mathbf{v}} \equiv U(\mathbf{v})$ . Using the chain rule we could have

$$\begin{aligned}
\nabla \tilde{f}(\mathbf{w}) - \nabla \tilde{f}(\mathbf{v}) &= \nabla U(\mathbf{w}) \nabla f(\tilde{\mathbf{w}}) - \nabla U(\mathbf{v}) \nabla f(\tilde{\mathbf{v}}) \\
&= (\nabla U(\mathbf{w}) - \nabla U(\mathbf{v})) \nabla f(\tilde{\mathbf{w}}) + \nabla U(\mathbf{v}) (\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))
\end{aligned}$$

We first move towards the smoothness property. Taking the norm on both sides, based on triangle inequality, we have:

$$\begin{aligned}
\|\nabla \tilde{f}(\mathbf{w}) - \nabla \tilde{f}(\mathbf{v})\| &= \|(\nabla U(\mathbf{w}) - \nabla U(\mathbf{v})) \nabla f(\tilde{\mathbf{w}}) + \nabla U(\mathbf{v}) (\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))\| \\
&\leq \|(\nabla U(\mathbf{w}) - \nabla U(\mathbf{v})) \nabla f(\tilde{\mathbf{w}})\| + \|\nabla U(\mathbf{v}) (\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))\|
\end{aligned} \tag{7}$$

We could bound the first term on the RHS by

$$\begin{aligned}
\|(\nabla U(\mathbf{w}) - \nabla U(\mathbf{v}))\nabla f(\tilde{\mathbf{w}})\| &\stackrel{(a)}{\leq} \|\nabla U(\mathbf{w}) - \nabla U(\mathbf{v})\| \|\nabla f(\tilde{\mathbf{w}})\| \\
&= \left\| \left( (1 - \alpha\delta)\mathbf{I} - \alpha\nabla^2 \hat{f}(\mathbf{w}) \right) - \left( (1 - \alpha\delta)\mathbf{I} - \alpha\nabla^2 \hat{f}(\mathbf{v}) \right) \right\| \|\nabla f(\tilde{\mathbf{w}})\| \\
&= \alpha \left\| \nabla^2 \hat{f}(\mathbf{w}) - \nabla^2 \hat{f}(\mathbf{v}) \right\| \|\nabla f(\tilde{\mathbf{w}})\| \\
&\stackrel{(b)}{\leq} \alpha\rho \|\mathbf{w} - \mathbf{v}\| \|\nabla f(\tilde{\mathbf{w}})\| \\
&\stackrel{(c)}{\leq} \alpha\rho G \|\mathbf{w} - \mathbf{v}\|
\end{aligned} \tag{8}$$

where (a) is due to Cauchy-Schwarz inequality, (b) is due to the Hessian Lipschitz property, and (c) is due to bounded gradient assumption. Similarly, we could bound the second term on (7)'s RHS by

$$\begin{aligned}
\|\nabla U(\mathbf{v})(\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))\| &= \left\| \left( (1 - \alpha\delta)\mathbf{I} - \alpha\nabla^2 \hat{f}(\mathbf{v}) \right) (\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}})) \right\| \\
&\stackrel{(a)}{\leq} (1 - \alpha\delta - \alpha\mu) \|\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}})\| \\
&\stackrel{(b)}{\leq} (1 - \alpha\delta - \alpha\mu)L \|\tilde{\mathbf{w}} - \tilde{\mathbf{v}}\| \\
&\stackrel{(c)}{=} (1 - \alpha\delta - \alpha\mu)L \|\mathbf{U}(\mathbf{w}) - \mathbf{U}(\mathbf{v})\| \\
&\stackrel{(d)}{\leq} (1 - \alpha\delta - \alpha\mu)L(1 - \alpha\delta - \alpha\mu) \|\mathbf{w} - \mathbf{v}\| \\
&= (1 - \alpha\delta - \alpha\mu)^2 L \|\mathbf{w} - \mathbf{v}\|
\end{aligned} \tag{9}$$

Here, (a) is due to  $(1 - \alpha\delta)\mathbf{I} - \alpha\nabla^2 \hat{f}(\mathbf{v})$  being symmetric, semi-positive definite, and  $\lambda_{\max} \left( (1 - \alpha\delta)\mathbf{I} - \alpha\nabla^2 \hat{f}(\mathbf{v}) \right) \leq (1 - \alpha\delta) - \alpha\mu$  (see Lemma 2). Step (b) is due to  $f(\cdot)$  is L-smooth. Step (c) is the use of short hand  $\tilde{\mathbf{w}} \equiv \mathbf{U}(\mathbf{w})$ ,  $\tilde{\mathbf{v}} \equiv \mathbf{U}(\mathbf{v})$ . Finally, step (d) is achieved by using Lemma 2 on  $\mathbf{U}(\cdot)$ . Put the result of (8) and (9) into (7), we have

$$\begin{aligned}
\|\nabla \tilde{f}(\mathbf{w}) - \nabla \tilde{f}(\mathbf{v})\| &\leq \|(\nabla U(\mathbf{w}) - \nabla U(\mathbf{v}))\nabla f(\tilde{\mathbf{w}})\| + \|\nabla U(\mathbf{v})(\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))\| \\
&\leq \alpha\rho G \|\mathbf{w} - \mathbf{v}\| + (1 - \alpha\delta - \alpha\mu)^2 L \|\mathbf{w} - \mathbf{v}\| \\
&= (\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L) \|\mathbf{w} - \mathbf{v}\|
\end{aligned}$$

and thus  $\tilde{f}(\cdot)$  is  $\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L$  smooth.

Similarly, we first use triangle inequality to find the lower bound.

$$\begin{aligned}
\|\nabla \tilde{f}(\mathbf{w}) - \nabla \tilde{f}(\mathbf{v})\| &= \|(\nabla U(\mathbf{w}) - \nabla U(\mathbf{v}))\nabla f(\tilde{\mathbf{w}}) + \nabla U(\mathbf{v})(\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))\| \\
&\geq \|\nabla U(\mathbf{v})(\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))\| - \|(\nabla U(\mathbf{w}) - \nabla U(\mathbf{v}))\nabla f(\tilde{\mathbf{w}})\|
\end{aligned}$$

The second term on RHS has already been derived in (8). For the first term, we could bound it by

$$\begin{aligned}
\|\nabla U(\mathbf{v})(\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))\| &= \left\| \left( (1 - \alpha\delta)\mathbf{I} - \alpha\nabla^2 \hat{f}(\mathbf{v}) \right) (\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}})) \right\| \\
&\stackrel{(a)}{\geq} (1 - \alpha\delta - \alpha L) \|\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}})\| \\
&\stackrel{(b)}{\geq} (1 - \alpha\delta - \alpha L)\mu \|\tilde{\mathbf{w}} - \tilde{\mathbf{v}}\| \\
&= (1 - \alpha\delta - \alpha L)\mu \|(1 - \alpha\delta)\mathbf{w} - \alpha\nabla \hat{f}(\mathbf{w}) - (1 - \alpha\delta)\mathbf{v} + \alpha\nabla \hat{f}(\mathbf{v})\| \\
&\geq \mu(1 - \alpha\delta - \alpha L)((1 - \alpha\delta)\|\mathbf{w} - \mathbf{v}\| - \alpha\|\nabla \hat{f}(\mathbf{w}) - \nabla \hat{f}(\mathbf{v})\|) \\
&\stackrel{(c)}{\geq} \mu(1 - \alpha\delta - \alpha L)((1 - \alpha\delta)\|\mathbf{w} - \mathbf{v}\| - \alpha L\|\mathbf{w} - \mathbf{v}\|) \\
&\geq \mu(1 - \alpha\delta - \alpha L)^2 \|\mathbf{w} - \mathbf{v}\|
\end{aligned}$$

Here (a) is due to  $\lambda_{\min} \left( I - \alpha\delta - \alpha\nabla^2\hat{f}(v) \right) \geq 1 - \alpha\delta - \alpha L$ , (b) is due to  $f(\cdot)$  being  $\mu$ -strongly convex, and (c) is due to  $\hat{f}(\cdot)$  being  $L$ -smooth. Put the results together, we have that

$$\begin{aligned} \|\nabla\tilde{f}(\mathbf{w}) - \nabla\tilde{f}(\mathbf{v})\| &\geq \|\nabla U(\mathbf{v})(\nabla f(\tilde{\mathbf{w}}) - \nabla f(\tilde{\mathbf{v}}))\| - \|(\nabla U(\mathbf{w}) - \nabla U(\mathbf{v}))\nabla f(\tilde{\mathbf{w}})\| \\ &\geq (\mu(1 - \alpha\delta - \alpha L)^2 - \alpha\rho G) \|\mathbf{w} - \mathbf{v}\| \end{aligned}$$

Thus the function  $\tilde{f}(\cdot)$  is  $\mu(1 - \alpha\delta - \alpha L)^2 - \alpha\rho G$  strongly convex.  $\mu(1 - \alpha\delta - \alpha L)^2 - \alpha\rho G$  is positive since  $\frac{\alpha\rho G}{\mu} < (\frac{1}{2} - \alpha L)^2$ .  $\square$

## A.2.2 Generalization Bound

---

**Algorithm 2** MAML [4] (the Original Algorithm without Regularization)

---

**Require:** Datasets  $\mathcal{S} = \{\mathcal{S}_i^{\text{in}}, \mathcal{S}_i^{\text{out}}\}_{i=1}^m$ ; few-shot meta-query batch size  $K$ ; the number of training tasks sampled at each round  $r$ ; the total number of iterations  $T$ .

- 1: Initialize the model parameters  $w^0$ .
  - 2: **for**  $t = 0$  to  $T - 1$  **do**
  - 3:   Randomly sample  $r$  tasks from the set of  $m$  available tasks with indices stored in  $\mathcal{B}_t$ .
  - 4:   **for** each sampled task  $\mathcal{T}_i$  **do**
  - 5:     Sample a size  $K$  support data batch  $\mathcal{D}_i^{t, \text{in}}$  from  $\mathcal{S}_i^{\text{in}}$ ;
  - 6:     Sample a size  $b$  query data batch  $\mathcal{D}_i^{t, \text{out}}$  from  $\mathcal{S}_i^{\text{out}}$ ;
  - 7:     Calculate  $w_i^{t+1} := w^t - \beta_t \nabla_{w^t} \hat{\mathcal{L}} \left( w^t - \alpha \nabla \hat{\mathcal{L}} \left( w^t, \mathcal{D}_i^{t, \text{in}} \right), \mathcal{D}_i^{t, \text{out}} \right)$ ;
  - 8:   **end for**
  - 9:   Meta-update  $w^{t+1} := \frac{1}{r} \sum_{i \in \mathcal{B}_t} w_i^{t+1}$
  - 10: **end for**
  - 11: **Return:**  $w^T$
- 

This section provides the derivation of the generalization bound of MAML with inner-level L2-Norm regularization. The proofs are based on the derivation framework based on algorithm stability proposed by [3]. The framework's preliminary is consistent with this work and shares the same assumptions as this work. To include the notations used in this section, we provide a restatement of the unregularized MAML steps in Algorithm 2.

We restate some notations for a clearer explanation. In the following of this appendix, we use the *hat* superscript to distinguish *empirical losses* from *population losses*. And we use the *tilda* superscript to denote the functions, algorithms, or processes *involving the inner-level regularization*, e.g.,

$$\begin{aligned} \hat{\tilde{F}}_i(w, \mathcal{S}_i) &:= \frac{1}{\binom{n}{K}} \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \hat{\mathcal{L}} \left( w - \alpha \nabla_w \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}) + \frac{\delta}{2} \|w\|^2, \mathcal{S}_i^{\text{out}} \right) \\ &= \frac{1}{\binom{n}{K}} \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \frac{1}{n} \sum_{z \in \mathcal{S}_i^{\text{out}}} \ell \left( w - \frac{\alpha}{K} \sum_{z' \in \mathcal{D}_i^{\text{in}}} \nabla_w \left( \ell(w, z') + \frac{\delta}{2} \|w\|^2 \right), z \right) \end{aligned}$$

Similarly,  $\tilde{F}(\cdot)$  and  $\tilde{F}_i(\cdot, \mathcal{S}_i)$  are corresponding to functions with inner-level regularization, distinguished from  $F(\cdot)$  and  $F_i(\cdot, \mathcal{S}_i)$  corresponding to unregularized MAML.  $\tilde{\mathcal{A}}$  also refers to the algorithm (MAML) with inner-level regularization, with output  $\tilde{\mathcal{A}}(\mathcal{S})$ .

Our goal is to bound

$$\mathbb{E}_{\tilde{\mathcal{A}}, \mathcal{S}} [F(\tilde{\mathcal{A}}(\mathcal{S})) - \hat{F}(\tilde{\mathcal{A}}(\mathcal{S}), \mathcal{S})] \quad (10)$$

Which is the expected discrepancy between population loss and empirical loss. Note that the loss is evaluated using the original MAML's unregularized inner-updating rule at the test time, while the model  $\tilde{\mathcal{A}}(\mathcal{S})$  is generated by MAML with inner-level regularization.

Then, we are going to bound (10) using the algorithm-stability-based framework proposed by [3]. We first include definitions and the key lemma of the framework.

**Definition 1. (symmetric algorithm)** We define an algorithm  $\mathcal{A} : \mathcal{Z}^n \rightarrow \mathbb{R}^d$  to be symmetric if its output distribution, denoted by  $\mathcal{A}(\mathcal{S})$ , remains unchanged under any permutation of the input set  $\mathcal{S} \subset \mathcal{Z}^n$ . In other words, if we take another set  $\mathcal{S}'$  that is a permutation of  $\mathcal{S}$ , the distributions of  $\mathcal{A}(\mathcal{S})$  and  $\mathcal{A}(\mathcal{S}')$  would be similar.

**Definition 2.  $(\gamma, K)$ -uniformly stability, from [3])** Consider the problem in (2) of the main paper, and let  $\mathcal{A}$  be a randomized algorithm that produces output  $w_{\mathcal{S}}$  given dataset  $\mathcal{S}$ . We say that  $\mathcal{A}$  is  $(\gamma, K)$ -uniformly stable if the following condition holds: for any  $i \in \{1, \dots, m\}$ , let  $\tilde{\mathcal{S}}$  be a dataset that is identical to  $\mathcal{S}$  except that  $\tilde{\mathcal{S}}_i^{\text{in}}$  and  $\tilde{\mathcal{S}}_i^{\text{out}}$  differ from  $\mathcal{S}_i^{\text{in}}$  and  $\mathcal{S}_i^{\text{out}}$ , respectively, in at most  $K$  and one data points. For any  $\bar{z} \in \mathcal{Z}$  and any set of  $K$  distinct points  $\{z_1, \dots, z_K\}$  in  $\mathcal{Z}$ ,

$$\mathbb{E}_{\mathcal{A}} \left[ \left| \ell \left( w_{\mathcal{S}} - \alpha \nabla \hat{\mathcal{L}} \left( w_{\mathcal{S}}, \{z_j\}_{j=1}^K \right), \bar{z} \right) - \ell \left( w_{\tilde{\mathcal{S}}} - \alpha \nabla \hat{\mathcal{L}} \left( w_{\tilde{\mathcal{S}}}, \{z_j\}_{j=1}^K \right), \bar{z} \right) \right| \right] \leq \gamma$$

the expectation is with respect to the randomness of  $\mathcal{A}$ .

**Lemma 4. (stability and generalization error, from [3])** Let  $F$  and  $\hat{F}$  be the population and empirical losses defined in Equations (1) and (2) of the main paper, respectively. Suppose Assumption 2 holds and let  $\mathcal{A}$  be a (possibly randomized) symmetric and  $(\gamma, K)$ -uniformly stable algorithm that produces output  $w_{\mathcal{S}} \in \mathcal{W}$  given input dataset  $\mathcal{S}$ . Then, the expected difference between the population loss and the empirical loss of  $w_{\mathcal{S}}$  is bounded by  $\gamma$ , i.e.,  $\mathbb{E}_{\mathcal{A}, \mathcal{S}} \left[ F(w_{\mathcal{S}}) - \hat{F}(w_{\mathcal{S}}, \mathcal{S}) \right] \leq \gamma$ .

Lemma 4 shows that if a symmetric algorithm could be proven to be  $(\gamma, K)$ -uniformly stable, we could bound its generalization error by capturing its stability parameter  $\gamma$ . Then we show the proof of generalization bound for (10) following this idea.

### Proof of Theorem 1.

**Theorem 1. (generalization bound)** If Assumption 1 and 2 hold. With  $\alpha \leq \frac{1}{2L}$ ,  $\beta_t \leq \frac{1}{\alpha \rho G + (1 - \alpha \delta - \alpha \mu)^2 L}$ ,  $\delta < \frac{1}{2\alpha}$  and  $\frac{\alpha \rho G}{\mu} < (\frac{1}{2} - \alpha L)^2$ . The model  $\tilde{\mathcal{A}}(\mathcal{S})$  generated by the last iterate of MAML with regularized updating rule introduced in (6) of the main paper satisfies

$$\mathbb{E}_{\tilde{\mathcal{A}}, \mathcal{S}} [F(\tilde{\mathcal{A}}(\mathcal{S})) - \hat{F}(\tilde{\mathcal{A}}(\mathcal{S}), \mathcal{S})] \leq \frac{2G^2(1 + \alpha L)(1 - \alpha \mu - \alpha \delta + (2 + \alpha L - \alpha \delta)\alpha L K)}{mn} \left( \frac{1}{\alpha \rho G + (1 - \alpha \delta - \alpha \mu)^2 L} + \frac{1}{-\alpha \rho G + (1 - \alpha \delta - \alpha L)^2 \mu} \right)$$

where the expectation is taken over the randomness of  $\tilde{\mathcal{A}}$  and sampling of  $\mathcal{S}$ .

*Proof.* The result in Lemma 4 means that we could bound the generalization error of MAML with inner-level regularization  $\tilde{\mathcal{A}}$  by proving its  $(\gamma, K)$ -uniformly stable as defined in Definition 2 and capture the  $\gamma$  parameter.

The Definition 2 of  $(\gamma, K)$ -uniformly stability means that there is a dataset  $\tilde{\mathcal{S}}$  which is the same as  $\mathcal{S}$  except one  $i$  such that:

- $\tilde{\mathcal{S}}_i^{\text{in}}$  has at most  $K$  data points different from  $\mathcal{S}_i^{\text{in}}$ . We denote the  $K$  samples in each dataset by  $\{z_j\}_{j=1}^K$  and  $\{\bar{z}_j\}_{j=1}^K$ .
- $\tilde{\mathcal{S}}_i^{\text{out}}$  has at most 1 data points different from  $\mathcal{S}_i^{\text{out}}$ . They are denoted by  $\zeta$  and  $\bar{\zeta}$ .

We consider the two parallel processes of training  $\{w^t\}$  and  $\{\bar{w}^t\}$  using datasets  $\mathcal{S}$  and  $\tilde{\mathcal{S}}$ . The bar superscript is used to denote the process using  $\tilde{\mathcal{S}}$ .  $D_i^{t, \text{out}}$  and  $D_i^{t, \text{in}}$  are referring to indices of samples in  $\mathcal{D}_i^{t, \text{out}}$  and  $\mathcal{D}_i^{t, \text{in}}$ , respectively. We could assume the parallel using a same random machine to sample batches for generating  $\{w^t\}$  and  $\{\bar{w}^t\}$ , i.e.,  $\mathcal{B}_t = \bar{\mathcal{B}}_t$ ,  $D_i^{t, \text{out}} = \bar{D}_i^{t, \text{out}}$ , and  $D_i^{t, \text{in}} = \bar{D}_i^{t, \text{in}}$ .

We use  $v_t$  to denote the number of indices corresponding to  $\{z_j\}_{j=1}^K$  (or  $\{\bar{z}_j\}_{j=1}^K$ ) is chosen in  $D_i^{t, \text{in}}$ . And  $u_t$  is denoting the number of times that the index of sample  $\zeta$  (or  $\bar{\zeta}$ ) is chosen in  $D_i^{t, \text{out}}$ , respectively. As shown in [3], recalling the definition of  $b$  and  $r$  from Algorithm 2, for each  $t$ , the expectations of  $v_t$  and  $u_t$  are given by

$$\mathbb{E}[v_t] = \frac{K^2 r}{nm}, \quad \mathbb{E}[u_t] = \frac{br}{nm} \tag{11}$$

Then we are coming to the main proof. We first claim that

$$\mathbb{E}_{\tilde{\mathcal{A}}} [\|w^T - \bar{w}^T\|] \leq \frac{2G(1 - \alpha\mu - \alpha\delta + (2 + \alpha L - \alpha\delta)\alpha LK)}{mn} \left( \frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu} \right) \quad (12)$$

when conditions in Assumption 1 are satisfied and  $\{w^T, \bar{w}^T\}$  is generated using  $\tilde{\mathcal{A}}$ . The next is to prove this claim. To simplify the notation, let us define  $\psi(w; \mathcal{D}, z) := \ell(w - \alpha \nabla_w (\hat{\mathcal{L}}(w, \mathcal{D}) + \frac{\delta}{2} \|w\|^2), z)$ . Note that

$$\begin{aligned} \psi(w; \mathcal{D}, z) &:= \ell(w - \alpha \nabla_w (\hat{\mathcal{L}}(w, \mathcal{D}) + \frac{\delta}{2} \|w\|^2), z) \\ &= \ell((1 - \alpha\delta)w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}), z) \end{aligned}$$

and

$$\nabla \psi(w; \mathcal{D}, z) = \left( (1 - \alpha\delta)I - \alpha \nabla^2 \hat{\mathcal{L}}(w, \mathcal{D}) \right) \nabla \ell \left( (1 - \alpha\delta)w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}), z \right).$$

Recalling from Lemma 1, we know for a  $\lambda$ -strongly convex and  $\eta$ -smooth function  $\phi$ , we have

$$\|(u - \beta \nabla \phi(u)) - (v - \beta \nabla \phi(v))\| \leq \left( 1 - \frac{\beta \lambda \eta}{\lambda + \eta} \right) \|u - v\|$$

for any  $u$  and  $v$ .

And Lemma 3 shows that for any batch  $\mathcal{D}$  and any  $z \in \mathcal{Z}$ ,  $\psi(w; \mathcal{D}, z)$  is  $\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L$  smooth and  $-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu$  strongly convex. Hence, using Lemma 1, for any  $j \in \mathcal{B}_t$  that  $j \neq i$ , we have

$$\|w_j^{t+1} - \bar{w}_j^{t+1}\| \leq \left( 1 - \beta_t \frac{1}{\frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu}} \right) \|w^t - \bar{w}^t\|. \quad (13)$$

Next, For the case that  $i \in \mathcal{B}_t$ , we could have

$$\begin{aligned} \|w_i^{t+1} - \bar{w}_i^{t+1}\| &\leq \frac{1}{b} \sum_{z \in \mathcal{D}_i^{t, \text{out}}} \left\| \left( w^t - \beta_t \nabla \psi \left( w^t; \mathcal{D}_i^{t, \text{in}}, z \right) \right) - \left( \bar{w}^t - \beta_t \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^{t, \text{in}}, z \right) \right) \right\| \\ &\quad + \frac{1}{b} \beta_t \sum_{z \in \bar{\mathcal{D}}_i^{t, \text{out}} / \mathcal{D}_i^{t, \text{out}}} \left\| \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^{t, \text{in}}, z \right) - \nabla \psi \left( w^t; \mathcal{D}_i^{t, \text{in}}, z \right) \right\|. \end{aligned} \quad (14)$$

To bound the second term on RHS of (14), we first consider that

$$\begin{aligned} \|\nabla \psi(w; \mathcal{D}, z)\| &= \left\| \left( (1 - \alpha\delta)I - \alpha \nabla^2 \hat{\mathcal{L}}(w, \mathcal{D}) \right) \nabla \ell \left( (1 - \alpha\delta)w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}), z \right) \right\| \\ &\leq \left\| \left( (1 - \alpha\delta)I - \alpha \nabla^2 \hat{\mathcal{L}}(w, \mathcal{D}) \right) \right\| \|\nabla \ell \left( (1 - \alpha\delta)w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}), z \right)\| \\ &\leq (1 - \alpha\mu - \alpha\delta)G. \end{aligned} \quad (15)$$

The last inequality is due to bounded gradient assumption for  $\ell(\cdot, z)$  and  $(1 - \alpha\delta)I - \alpha \nabla^2 \hat{\mathcal{L}}(w, \mathcal{D})$  being symmetric, semi-positive definite, and  $\lambda_{\max} \left( (1 - \alpha\delta)I - \alpha \nabla^2 \hat{\mathcal{L}}(w, \mathcal{D}) \right) \leq (1 - \alpha\delta) - \alpha\mu$  (see Lemma 2). Then, we have

$$\begin{aligned} \left\| \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^{t, \text{in}}, z \right) - \nabla \psi \left( w^t; \mathcal{D}_i^{t, \text{in}}, z \right) \right\| &\leq \left\| \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^{t, \text{in}}, z \right) \right\| + \left\| \nabla \psi \left( w^t; \mathcal{D}_i^{t, \text{in}}, z \right) \right\| \\ &\leq 2(1 - \alpha\mu - \alpha\delta)G. \end{aligned} \quad (16)$$

Since  $|\bar{\mathcal{D}}_i^{t, \text{out}} / \mathcal{D}_i^{t, \text{out}}| = u_t$ , using the above result, the second term of (14)'s RHS could be bounded by  $2\beta_t u_t G(1 - \alpha\mu - \alpha\delta)/b$ , i.e.,

$$\begin{aligned} \|w_i^{t+1} - \bar{w}_i^{t+1}\| &\leq 2\beta_t G(1 - \alpha\mu - \alpha\delta) \frac{u_t}{b} \\ &\quad + \frac{1}{b} \sum_{z \in \mathcal{D}_i^{t, \text{out}}} \left\| \left( w^t - \beta_t \nabla \psi \left( w^t; \mathcal{D}_i^{t, \text{in}}, z \right) \right) - \left( \bar{w}^t - \beta_t \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^{t, \text{in}}, z \right) \right) \right\|. \end{aligned} \quad (17)$$



For the second term, note that

$$\begin{aligned}
& \left\| \left( w^t - \beta_t \nabla \psi \left( w^t; \mathcal{D}_i^t, \text{in}, z \right) \right) - \left( \bar{w}^t - \beta_t \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^t, \text{in}, z \right) \right) \right\| \\
& \leq \left\| \left( w^t - \beta_t \nabla \psi \left( w^t; \mathcal{D}_i^t, \text{in}, z \right) \right) - \left( \bar{w}^t - \beta_t \nabla \psi \left( \bar{w}^t; \mathcal{D}_i^t, \text{in}, z \right) \right) \right\| \\
& \quad + \beta_t \left\| \nabla \psi \left( \bar{w}^t; \mathcal{D}_i^t, \text{in}, z \right) - \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^t, \text{in}, z \right) \right\|.
\end{aligned} \tag{18}$$

For the first term on the RHS of (18), we could bound it similarly to the derivation of (13) by

$$\begin{aligned}
& \left\| \left( w^t - \beta_t \nabla \psi \left( w^t; \mathcal{D}_i^t, \text{in}, z \right) \right) - \left( \bar{w}^t - \beta_t \nabla \psi \left( \bar{w}^t; \mathcal{D}_i^t, \text{in}, z \right) \right) \right\| \\
& \leq \left( 1 - \beta_t \frac{1}{\frac{1}{\alpha \rho G + (1 - \alpha \delta - \alpha \mu)^2 L} + \frac{1}{-\alpha \rho G + (1 - \alpha \delta - \alpha L)^2 \mu}} \right) \|w^t - \bar{w}^t\|.
\end{aligned} \tag{19}$$

And for the second term on the RHS of (18), we have

$$\begin{aligned}
& \left\| \nabla \psi \left( \bar{w}^t; \mathcal{D}_i^t, \text{in}, z \right) - \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^t, \text{in}, z \right) \right\| \\
& = \left\| \left( (1 - \alpha \delta) I - \alpha \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) \right) \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) \right. \\
& \quad \left. - \left( (1 - \alpha \delta) I - \alpha \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right) \right) \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| \\
& \leq (1 - \alpha \delta) \left\| \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) - \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| + \\
& \alpha \left\| \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) \right. \\
& \quad \left. - \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right) \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| \\
& \leq (1 - \alpha \delta) \left\| \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) - \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| + \\
& \alpha \left\| \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) - \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right. \\
& \quad \left. + \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) - \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right) \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| \\
& \leq (1 - \alpha \delta) \left\| \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) - \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| + \\
& \alpha \left\| \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) \right\| \left\| \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) - \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| + \\
& \alpha \left\| \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) - \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right) \right\| \left\| \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| \\
& \leq (1 - \alpha \delta + \alpha L) \left\| \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) - \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| + \\
& \alpha G \left\| \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) - \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right) \right\|.
\end{aligned} \tag{20}$$

The last inequality is given by the smoothness and bounded gradient assumption for  $\ell(\cdot, z)$ . Next, we are going to bound the terms in (20). Note that

$$\begin{aligned}
& \left\| \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right), z \right) - \nabla \ell \left( (1 - \alpha \delta) \bar{w}^t - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right), z \right) \right\| \\
& \leq \alpha L \left\| \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) - \nabla \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right) \right\| \leq 2\alpha L G \frac{v_t}{K}
\end{aligned}$$

and

$$\left\| \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \mathcal{D}_i^t, \text{in} \right) - \nabla^2 \hat{\mathcal{L}} \left( \bar{w}^t, \bar{\mathcal{D}}_i^t, \text{in} \right) \right\| \leq 2L \frac{v_t}{K}.$$

Putting the above two results into (20), we have

$$\begin{aligned} \left\| \nabla \psi \left( \bar{w}^t; \mathcal{D}_i^{t, \text{in}}, z \right) - \nabla \psi \left( \bar{w}^t; \bar{\mathcal{D}}_i^{t, \text{in}}, z \right) \right\| &\leq 2(1 - \alpha\delta + \alpha L)\alpha LG \frac{v_t}{K} + 2\alpha LG \frac{v_t}{K} \\ &= 2(2 + \alpha L - \alpha\delta)\alpha LG \frac{v_t}{K}. \end{aligned} \quad (21)$$

Putting the result in (19) and (18) into (17), we have

$$\begin{aligned} \|w_i^{t+1} - \bar{w}_i^{t+1}\| &\leq \left( 1 - \beta_t \frac{1}{\frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu}} \right) \|w^t - \bar{w}^t\| \\ &\quad + 2\beta_t G \left( (1 - \alpha\mu - \alpha\delta) \frac{u_t}{b} + \alpha L (2 + \alpha L - \alpha\delta) \frac{v_t}{K} \right). \end{aligned}$$

Along with (13), we have

$$\begin{aligned} \left\| \frac{1}{r} \sum_{j \in \mathcal{B}_t} w_j^{t+1} - \frac{1}{r} \sum_{j \in \mathcal{B}_t} \bar{w}_j^{t+1} \right\| &\leq \left( 1 - \beta_t \frac{1}{\frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu}} \right) \|w^t - \bar{w}^t\| \\ &\quad + 2\beta_t G \left( (1 - \alpha\mu - \alpha\delta) \frac{u_t}{rb} + \alpha L (2 + \alpha L - \alpha\delta) \frac{v_t}{rK} \right), \end{aligned}$$

which indicates

$$\begin{aligned} \|w^{t+1} - \bar{w}^{t+1}\| &\leq \left( 1 - \beta_t \frac{1}{\frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu}} \right) \|w^t - \bar{w}^t\| \\ &\quad + 2\beta_t G \left( (1 - \alpha\mu - \alpha\delta) \frac{u_t}{rb} + \alpha L (2 + \alpha L - \alpha\delta) \frac{v_t}{rK} \right). \end{aligned}$$

Using (11), we could take the expectation for both sides and get

$$\begin{aligned} \mathbb{E}_{\bar{\mathcal{A}}} [\|w^{t+1} - \bar{w}^{t+1}\|] &\leq \left( 1 - \beta_t \frac{1}{\frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu}} \right) \mathbb{E}_{\bar{\mathcal{A}}} [\|w^t - \bar{w}^t\|] \\ &\quad + 2 \frac{\beta_t G}{mn} (1 - \alpha\mu - \alpha\delta + (2 + \alpha L - \alpha\delta)\alpha LK). \end{aligned}$$

The bound could be rewritten as

$$\mathbb{E}_{\bar{\mathcal{A}}} [\|w^{t+1} - \bar{w}^{t+1}\|] \leq (1 - \beta_t \lambda) \mathbb{E}_{\bar{\mathcal{A}}} [\|w^t - \bar{w}^t\|] + \beta_t \eta,$$

where the  $\lambda$  and  $\eta$  are given by

$$\lambda := \frac{1}{\frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu}}, \quad \eta := \frac{2G}{mn} (1 - \alpha\mu - \alpha\delta + (2 + \alpha L - \alpha\delta)\alpha LK).$$

In fact, the main claim (12) is equivalent to

$$\mathbb{E}_{\bar{\mathcal{A}}} [\|w^t - \bar{w}^t\|] \leq \frac{\eta}{\lambda}.$$

For  $t = 1$ , this is true because of  $\beta_0 \leq \frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} \leq \frac{1}{\lambda}$ . Then the result could be easily obtained by induction. The result could be written as

$$\mathbb{E}_{\bar{\mathcal{A}}} [\|w^T - \bar{w}^T\|] \leq \frac{2G(1 - \alpha\mu - \alpha\delta + (2 + \alpha L - \alpha\delta)\alpha LK)}{mn} \left( \frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu} \right).$$

Having proved the above result, we are ready to finish proving Theorem 1. We have

$$\begin{aligned} &\left| \ell \left( w^T - \alpha \nabla \hat{\mathcal{L}} \left( w^T, \{z_j\}_{j=1}^K \right), \bar{z} \right) - \ell \left( \bar{w}^T - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^T, \{z_j\}_{j=1}^K \right), \bar{z} \right) \right| \\ &\leq G \left\| \left( w^T - \alpha \nabla \hat{\mathcal{L}} \left( w^T, \{z_j\}_{j=1}^K \right), \bar{z} \right) - \left( \bar{w}^T - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^T, \{z_j\}_{j=1}^K \right), \bar{z} \right) \right\| \\ &\leq G \|w^T - \bar{w}^T\| + \alpha G \left\| \nabla \hat{\mathcal{L}} \left( w^T, \{z_j\}_{j=1}^K \right) - \nabla \hat{\mathcal{L}} \left( \bar{w}^T, \{z_j\}_{j=1}^K \right) \right\| \\ &\leq (1 + \alpha L) G \|w^T - \bar{w}^T\|. \end{aligned}$$

Then,

$$\mathbb{E}_{\tilde{\mathcal{A}}} \left[ \ell \left( w^T - \alpha \nabla \hat{\mathcal{L}} \left( w^T, \{z_j\}_{j=1}^K \right), \bar{z} \right) - \ell \left( \bar{w}^T - \alpha \nabla \hat{\mathcal{L}} \left( \bar{w}^T, \{z_j\}_{j=1}^K \right), \bar{z} \right) \right] \leq \frac{2G^2(1 + \alpha L)(1 - \alpha\mu - \alpha\delta + (2 + \alpha L - \alpha\delta)\alpha LK)}{mn} \left( \frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu} \right).$$

This means the algorithm is  $(\gamma, K)$ -uniformly stable with RHS as the  $\gamma$  parameter.

Finally, for the meta-testing phase learning objective of the original MAML (unregularized), by Lemma 4, the generalization bound in (10) is given by

$$\mathbb{E}_{\tilde{\mathcal{A}}, \mathcal{S}} [F(\tilde{\mathcal{A}}(\mathcal{S})) - \hat{F}(\tilde{\mathcal{A}}(\mathcal{S}), \mathcal{S})] \leq \frac{2G^2(1 + \alpha L)(1 - \alpha\mu - \alpha\delta + (2 + \alpha L - \alpha\delta)\alpha LK)}{mn} \left( \frac{1}{\alpha\rho G + (1 - \alpha\delta - \alpha\mu)^2 L} + \frac{1}{-\alpha\rho G + (1 - \alpha\delta - \alpha L)^2 \mu} \right),$$

where  $F(\cdot)$  and  $\hat{F}(\cdot, \mathcal{S})$  are population loss and empirical loss for unregularized MAML, respectively. The proof is complete.  $\square$

### A.2.3 Training Bias

In this section, we give the proof of Theorem 2 on training bias bound in the paper.

**Theorem 2.** (training bias bound) *If Assumption 1 and 2 hold. With  $\alpha \leq \frac{1}{2L}$ ,  $\delta < \frac{1}{2\alpha}$  and  $\frac{\alpha\rho G}{\mu} < (\frac{1}{2} - \alpha L)^2$ . The training bias from MAML with inner-level L2 regularization to the original MAML is bounded by*

$$\mathbb{E}_{\mathcal{S}} \left[ \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) \right] \leq \frac{\alpha^2(\alpha\rho G + (1 - \alpha\mu)^2 L)((1 - \alpha\mu - \alpha\delta)L\|w^*\| + G)^2 \delta^2}{2(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2 \mu)^2}$$

where  $\|w^*\| := \max_{\mathcal{S}} \|\arg \min_w \hat{F}(w, \mathcal{S})\|$ , the expectation is taken over sampling of  $\mathcal{S}$ .

*Proof.* The empirical loss of unregularized MAML is defined by

$$\hat{F}(w, \mathcal{S}) := \frac{1}{m} \sum_{i=1}^m \hat{F}_i(w, \mathcal{S}_i), \quad (22)$$

where  $\hat{F}_i(\cdot, \mathcal{S}_i)$  is given by

$$\begin{aligned} \hat{F}_i(w, \mathcal{S}_i) &:= \frac{1}{\binom{n}{k}} \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = k}} \frac{1}{n} \sum_{z \in \mathcal{S}_i^{\text{out}}} \ell \left( w - \frac{\alpha}{K} \sum_{z' \in \mathcal{D}_i^{\text{in}}} \nabla \ell(w, z'), z \right) \\ &= \frac{1}{\binom{n}{k}} \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = k}} \hat{\mathcal{L}} \left( w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}), \mathcal{S}_i^{\text{out}} \right). \end{aligned} \quad (23)$$

So (22) could also be written as

$$\hat{F}(w, \mathcal{S}) := \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = k}} \hat{\mathcal{L}} \left( w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}), \mathcal{S}_i^{\text{out}} \right). \quad (24)$$

For MAML with inner-level L2-Norm regularization, the corresponding empirical loss is given by

$$\hat{\hat{F}}(w, \mathcal{S}) := \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = k}} \hat{\mathcal{L}} \left( w - \alpha \nabla_w (\hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}) + \frac{\delta}{2} \|w\|^2), \mathcal{S}_i^{\text{out}} \right), \quad (25)$$

where  $\delta$  is the parameter for regularization. Our goal is to bound the training bias

$$\hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) = \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}). \quad (26)$$

To bound (26), we could first bound

$$\|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) - \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\| \quad (27)$$

We denote the two model parameters by

$$u := \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})$$

and

$$v := \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}).$$

Lemma 3 shows the strongly convexity of both  $\hat{\mathcal{L}}(w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}), \mathcal{S}_i^{\text{out}})$  and  $\hat{\mathcal{L}}(w - \alpha \nabla_w (\hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}) + \frac{\delta}{2} \|w\|^2), \mathcal{S}_i^{\text{out}})$ , which also indicates the strongly convexity of  $\hat{F}(\cdot, \mathcal{S})$  and  $\hat{F}(\cdot, \mathcal{S})$ . Suppose the optimal solution lies within  $\mathcal{W}$  for  $\hat{F}(\cdot, \mathcal{S})$  and  $\hat{F}(\cdot, \mathcal{S})$ , we have

$$\nabla \hat{F}(u, \mathcal{S}) = \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = k}} \nabla_w \hat{\mathcal{L}}(w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}), \mathcal{S}_i^{\text{out}}) \Bigg|_{w=u} = 0 \quad (28)$$

and

$$\nabla \hat{F}(v, \mathcal{S}) = \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = k}} \nabla_w \hat{\mathcal{L}}\left(w - \alpha \nabla_w (\hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}) + \frac{\delta}{2} \|w\|^2), \mathcal{S}_i^{\text{out}}\right) \Bigg|_{w=v} = 0. \quad (29)$$

Note that

$$\begin{aligned} \nabla_w \hat{\mathcal{L}}(w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}), \mathcal{S}_i^{\text{out}}) &= \\ (I_d - \alpha \nabla^2 \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}})) \nabla \hat{\mathcal{L}}(w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}}) & \end{aligned} \quad (30)$$

and

$$\begin{aligned} \nabla_w \hat{\mathcal{L}}\left(w - \alpha \nabla_w (\hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}) + \frac{\delta}{2} \|w\|^2), \mathcal{S}_i^{\text{out}}\right) &= \\ \left((1 - \alpha \delta) I_d - \alpha \nabla^2 \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}})\right) \nabla \hat{\mathcal{L}}\left((1 - \alpha \delta) w - \alpha \nabla \hat{\mathcal{L}}(w, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}}\right). & \end{aligned} \quad (31)$$

By plugging (30) and (31) into (28) and (29) respectively, we have

$$\nabla \hat{F}(u, \mathcal{S}) = \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = k}} \left(I_d - \alpha \nabla^2 \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}})\right) \nabla \hat{\mathcal{L}}(u - \alpha \nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}}) = 0 \quad (32)$$

and

$$\nabla \hat{F}(v, \mathcal{S}) = \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = k}} \left((1 - \alpha \delta) I_d - \alpha \nabla^2 \hat{\mathcal{L}}(v, \mathcal{D}_i^{\text{in}})\right) \nabla \hat{\mathcal{L}}\left((1 - \alpha \delta) v - \alpha \nabla \hat{\mathcal{L}}(v, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}}\right) = 0. \quad (33)$$

Then, we could bound  $\|\nabla \hat{F}(u, \mathcal{S}) - \nabla \hat{F}(v, \mathcal{S})\|$  by

$$\begin{aligned}
& \|\nabla \hat{F}(u, \mathcal{S}) - \nabla \hat{F}(v, \mathcal{S})\| \\
&= \|\nabla \hat{F}(u, \mathcal{S}) - 0\| = \\
&= \|\nabla \hat{F}(u, \mathcal{S}) - \nabla \hat{F}(u, \mathcal{S})\| \\
&= \left\| \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \left( (1 - \alpha\delta)I_d - \alpha\nabla^2 \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) \nabla \hat{\mathcal{L}}\left( (1 - \alpha\delta)u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right. \\
&\quad \left. - \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \left( I_d - \alpha\nabla^2 \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) \nabla \hat{\mathcal{L}}\left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right\| \\
&\stackrel{(a)}{=} \left\| \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \left[ \left( (1 - \alpha\delta)I_d - \alpha\nabla^2 \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) \nabla \hat{\mathcal{L}}\left( (1 - \alpha\delta)u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right. \right. \\
&\quad \left. \left. - \left( I_d - \alpha\nabla^2 \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) \nabla \hat{\mathcal{L}}\left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right] \right\| \\
&= \left\| \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \left[ \left( (1 - \alpha\delta)I_d - \alpha\nabla^2 \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) \left( \nabla \hat{\mathcal{L}}\left( (1 - \alpha\delta)u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) - \nabla \hat{\mathcal{L}}\left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right) \right. \right. \\
&\quad \left. \left. - \alpha\delta \nabla \hat{\mathcal{L}}\left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right] \right\| \\
&\stackrel{(b)}{\leq} \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \left\| \left( (1 - \alpha\delta)I_d - \alpha\nabla^2 \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) \left( \nabla \hat{\mathcal{L}}\left( (1 - \alpha\delta)u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) - \nabla \hat{\mathcal{L}}\left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right) \right. \\
&\quad \left. - \alpha\delta \nabla \hat{\mathcal{L}}\left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right\| \\
&\stackrel{(c)}{\leq} \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \left( \left\| \left( (1 - \alpha\delta)I_d - \alpha\nabla^2 \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) \left( \nabla \hat{\mathcal{L}}\left( (1 - \alpha\delta)u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) - \nabla \hat{\mathcal{L}}\left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right) \right\| \right. \\
&\quad \left. + \alpha|\delta| \left\| \nabla \hat{\mathcal{L}}\left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}), \mathcal{D}_i^{\text{out}} \right) \right\| \right) \\
&\stackrel{(d)}{\leq} \frac{1}{m} \frac{1}{\binom{n}{k}} \sum_{i=1}^m \sum_{\substack{\mathcal{D}_i^{\text{in}} \subset \mathcal{S}_i^{\text{in}} \\ |\mathcal{D}_i^{\text{in}}| = K}} \left( (1 - \alpha\mu - \alpha\delta)L \left\| \left( (1 - \alpha\delta)u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) - \left( u - \alpha\nabla \hat{\mathcal{L}}(u, \mathcal{D}_i^{\text{in}}) \right) \right\| \right. \\
&\quad \left. + \alpha|\delta|G \right) \\
&= \frac{1}{m} \frac{1}{\binom{n}{k}} m \binom{n}{k} \left( (1 - \alpha\mu - \alpha\delta)L \|\alpha\delta u\| + \alpha|\delta|G \right) \\
&= (\alpha|\delta|(1 - \alpha\mu - \alpha\delta)L\|u\| + \alpha|\delta|G) \\
&= \alpha|\delta|((1 - \alpha\mu - \alpha\delta)L\|u\| + G).
\end{aligned}$$

Here, (a) is by combining each term with the same index within the two summations, (b) and (c) are due to triangle inequality, and (d) is due to the strongly convex, smooth and bounded gradient property of  $\hat{\mathcal{L}}(\cdot, \mathcal{D}_i^{\text{in}})$ .

By the definition of  $u$  and  $v$ , (34) actually shows

$$\|\nabla \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \nabla \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S})\| \leq \alpha|\delta|((1 - \alpha\mu - \alpha\delta)L\|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\| + G). \quad (35)$$

Lemma 3 also indicates that  $\hat{F}(\cdot, \mathcal{S})$  is  $(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2\mu)$  strongly-convex, so we could bound  $\|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) - \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\|$  by

$$\begin{aligned} \|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) - \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\| &\leq \frac{1}{(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2\mu)} \|\nabla \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \nabla \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S})\| \\ &\leq \frac{1}{(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2\mu)} \alpha|\delta|((1 - \alpha\mu - \alpha\delta)L\|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\| + G). \end{aligned} \quad (36)$$

Take square for both sides of (36), we have

$$\|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) - \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\|^2 \leq \left( \frac{1}{(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2\mu)} \alpha|\delta|((1 - \alpha\mu - \alpha\delta)L\|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\| + G) \right)^2, \quad (37)$$

and then take the expectation over the sampling of  $\mathcal{S}$ , we further have

$$\mathbb{E}_{\mathcal{S}} \left[ \|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) - \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\|^2 \right] \leq \left( \frac{1}{(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2\mu)} \alpha|\delta|((1 - \alpha\mu - \alpha\delta)L\|w^*\| + G) \right)^2, \quad (38)$$

where  $\|w^*\| := \max_{\mathcal{S}} \|\arg \min_w \hat{F}(w, \mathcal{S})\|$ , the maximum is taken over sampling of  $\mathcal{S}$ .

Recall from Lemma 3 that  $\hat{F}(\cdot, \mathcal{S})$  is  $(\alpha\rho G + (1 - \alpha\mu)^2L)$  smooth, and note that  $\nabla \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) = 0$ , we could bound (26) by

$$\begin{aligned} \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) &= \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) \\ &\leq \frac{1}{2}(\alpha\rho G + (1 - \alpha\mu)^2L)\|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) - \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\|^2. \end{aligned} \quad (39)$$

Finally, by taking the expectation, we have

$$\begin{aligned} &\mathbb{E}_{\mathcal{S}} \left[ \hat{F}(\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}), \mathcal{S}) - \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) \right] \\ &\leq \mathbb{E}_{\mathcal{S}} \left[ \frac{1}{2}(\alpha\rho G + (1 - \alpha\mu)^2L)\|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) - \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\|^2 \right] \\ &= \frac{1}{2}(\alpha\rho G + (1 - \alpha\mu)^2L)\mathbb{E}_{\mathcal{S}} \left[ \|\arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S}) - \arg \min_{\mathcal{W}} \hat{F}(\cdot, \mathcal{S})\|^2 \right] \\ &\stackrel{(a)}{\leq} \frac{1}{2}(\alpha\rho G + (1 - \alpha\mu)^2L) \left( \frac{1}{(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2\mu)} \alpha|\delta|((1 - \alpha\mu - \alpha\delta)L\|w^*\| + G) \right)^2 \\ &= \frac{\alpha^2(\alpha\rho G + (1 - \alpha\mu)^2L)((1 - \alpha\mu - \alpha\delta)L\|w^*\| + G)^2\delta^2}{2(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2\mu)^2}. \end{aligned} \quad (40)$$

(a) is because of (38). The proof is complete.  $\square$

### A.3. Further Analysis

#### A.3.1 Property of Generalization Error Bound

In 4.1, we claimed that if we regard the generalization bound as a function  $GB(\delta)$ , its derivative  $GB'(\delta)$  would be positive for  $\delta \in (-\infty, \frac{1}{2\alpha})$ , i.e.,

$$GB'(\delta) > 0 \quad \forall \delta \in (-\infty, \frac{1}{2\alpha}) \quad (41)$$

In this section, we provide proof of this claim.

*Proof.* Based on the result of Theorem 1, the function  $GB(\delta)$  is given by

$$GB(\delta) = \frac{2G^2(1+\alpha L)(1-\alpha\mu-\alpha\delta+(2+\alpha L-\alpha\delta)\alpha LK)}{mn} \left( \frac{1}{\alpha\rho G+(1-\alpha\delta-\alpha\mu)^2L} + \frac{1}{-\alpha\rho G+(1-\alpha\delta-\alpha L)^2\mu} \right).$$

Taking its derivative, we have

$$\begin{aligned} GB'(\delta) = & \\ & \frac{2\alpha G^2(1+\alpha L)}{mn(-\alpha\rho G+(1-\alpha\delta-\alpha L)^2\mu)} \cdot \\ & (-\alpha LK+1)((1-\alpha\delta-\alpha L)^2\mu-\alpha\rho G)+2\mu(1-\alpha\delta-\alpha L)(\alpha LK(2+\alpha L-\alpha\delta)+(1-\alpha\delta-\alpha\mu)) \\ & + \\ & \frac{2\alpha G^2(1+\alpha L)}{mn(\alpha\rho G+(1-\alpha\delta-\alpha\mu)^2L)} \cdot \\ & (-\alpha LK+1)((1-\alpha\delta-\alpha\mu)^2L+\alpha\rho G)+2L(1-\alpha\delta-\alpha\mu)(\alpha Lk(2+\alpha L-\alpha\delta)+(1-\alpha\delta-\alpha\mu)). \end{aligned} \quad (42)$$

To prove (41), we are going to prove both terms on RHS of (42) are greater than 0 for  $\delta \in (-\infty, \frac{1}{2\alpha})$ .

For the first term on RHS of (42), having  $\frac{2\alpha G^2(1+\alpha L)}{mn(-\alpha\rho G+(1-\alpha\delta-\alpha L)^2\mu)} > 0$ , we only need to prove

$$-\alpha LK+1)((1-\alpha\delta-\alpha L)^2\mu-\alpha\rho G)+2\mu(1-\alpha\delta-\alpha L)(\alpha LK(2+\alpha L-\alpha\delta)+(1-\alpha\delta-\alpha\mu)) > 0. \quad (43)$$

(43) could be re-written by

$$-\alpha LK+1)((1-\alpha\delta-\alpha L)^2\mu-\alpha\rho G)+2\mu(1-\alpha\delta-\alpha L)((\alpha LK+1)(2+\alpha L-\alpha\delta)-(1+\alpha L+\alpha\mu)) > 0. \quad (44)$$

Then, since  $\mu > 0$ ,  $(\alpha LK+1) > 0$  and  $(1-\alpha\delta-\alpha L) > 0$ , when  $\delta < \frac{1}{2\alpha}$ , by dividing both sides of (44) by  $\mu(1-\alpha\delta-\alpha L)(\alpha LK+1)$ , we find this inequality is equivalent to

$$-(1-\alpha\delta-\alpha L-\frac{\alpha\rho G}{\mu(1-\alpha\delta-\alpha L)})+2(2+\alpha L-\alpha\delta-\frac{1+\alpha L+\alpha\mu}{\alpha LK+1}) > 0. \quad (45)$$

(45) is equivalent to

$$3+3\alpha L-\alpha\delta+\frac{\alpha\rho G}{\mu(1-\alpha\delta-\alpha L)}-\frac{2(1+\alpha L+\alpha\mu)}{\alpha LK+1} > 0. \quad (46)$$

Since  $\alpha LK+1 > 1$ , (46) is true if

$$3+3\alpha L-\alpha\delta+\frac{\alpha\rho G}{\mu(1-\alpha\delta-\alpha L)}-2(1+\alpha L+\alpha\mu) > 0. \quad (47)$$

By recombining the LHS of (47), we obtain its equivalent form

$$(1-\alpha\mu-\alpha\delta)+\alpha(L-\mu)+\frac{\alpha\rho G}{\mu(1-\alpha\delta-\alpha L)} > 0. \quad (48)$$

When  $\delta < \frac{1}{2\alpha}$ , (48) is true since  $(1-\alpha\mu-\alpha\delta) > 0$ ,  $\frac{\alpha\rho G}{\mu(1-\alpha\delta-\alpha L)} > 0$ , and  $\alpha(L-\mu)$  is non-negative. This proves (43) to be true and shows the first term of (42)'s RHS is greater than 0 for  $\delta \in (-\infty, \frac{1}{2\alpha})$ .

Then, we move on to prove the second term of (42)'s RHS is greater than 0 for  $\delta \in (-\infty, \frac{1}{2\alpha})$ . Having  $\frac{2\alpha G^2(1+\alpha L)}{mn(-\alpha\rho G+(1-\alpha\delta-\alpha\mu)^2L)} > 0$ , we only need to prove

$$-(\alpha LK + 1)((1 - \alpha\delta - \alpha\mu)^2L + \alpha\rho G) + 2L(1 - \alpha\delta - \alpha\mu)(\alpha LK(2 + \alpha L - \alpha\delta) + (1 - \alpha\delta - \alpha\mu)) > 0. \quad (49)$$

Similar to the proof of the first term, when  $\delta < \frac{1}{2\alpha}$  we could obtain an equivalent inequality for (49) by dividing its both sides by  $L(1 - \alpha\delta - \alpha\mu)(\alpha LK + 1)$  and reordering:

$$3 + 2\alpha L + \alpha\mu - \alpha\delta - \frac{\alpha\rho G}{L(1 - \alpha\delta - \alpha\mu)} - 2\frac{1 + \alpha L + \alpha\mu}{\alpha LK + 1} > 0. \quad (50)$$

Since  $\frac{\alpha\rho G}{\mu} < (\frac{1}{2} - \alpha L)^2$  and  $\alpha Lk + 1 > 1$ , (50) is true if

$$3 + 2\alpha L + \alpha\mu - \alpha\delta - \frac{\mu(\frac{1}{2} - \alpha L)^2}{L(1 - \alpha\delta - \alpha\mu)} - 2(1 + \alpha L + \alpha\mu) \geq 0. \quad (51)$$

Note that we have

$$\frac{\mu(\frac{1}{2} - \alpha L)^2}{L(1 - \alpha\delta - \alpha\mu)} \stackrel{(a)}{<} \frac{\mu(1 - \alpha\delta - \alpha L)^2}{L(1 - \alpha\delta - \alpha\mu)} \stackrel{(b)}{\leq} \frac{L(1 - \alpha\delta - \alpha L)^2}{L(1 - \alpha\delta - \alpha L)} = 1 - \alpha\delta - \alpha L, \quad (52)$$

where (a) is true when  $\delta < \frac{1}{2\alpha}$  and (b) is because of  $\mu \leq L$ . By plugging (52) into (51), we have (51) being true when

$$3 + 2\alpha L + \alpha\mu - \alpha\delta - (1 - \alpha\delta - \alpha L) - 2(1 + \alpha L + \alpha\mu) \geq 0 \quad (53)$$

(53) is equivalent to

$$\alpha L - \alpha\mu \geq 0. \quad (54)$$

This is obviously true since we have  $\alpha > 0$  and  $L \geq \mu$ . This proves (49) to be true and shows the second term of (42)'s RHS is greater than 0 for  $\delta \in (-\infty, \frac{1}{2\alpha})$ .

Since both terms of (42)'s RHS being greater than 0 for  $\delta < \frac{1}{2\alpha}$  has been proved, the conclusion  $GB'(\delta) > 0 \forall \delta \in (-\infty, \frac{1}{2\alpha})$  in (41) is obtained. The proof is complete.  $\square$

### A.3.2 Property of Training Bias Bound

In 4.2, we claimed that if we regard the training bias bound as a function  $TB(\delta)$ , for a legal positive choice of  $\delta$ -value  $\delta_0$ , we would always have  $TB(\delta_0) > TB(-\delta_0)$ , i.e.,

$$TB(\delta_0) > TB(-\delta_0) \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}) \quad (55)$$

In this section, we provide proof of this claim.

*Proof.* From the result of Theorem 2, the training bias bound function of  $\delta$  is given by

$$TB(\delta) = \frac{\alpha^2(\alpha\rho G + (1 - \alpha\mu)^2L)((1 - \alpha\mu - \alpha\delta)L\|w^*\| + G)^2\delta^2}{2(-\alpha\rho G + (1 - \alpha L - \alpha\delta)^2\mu)^2}. \quad (56)$$

It's easy to find that  $TB(\delta) > 0$  for any  $\delta \neq 0$ . So the conclusion (55) is equivalent to

$$\frac{TB(\delta_0)}{TB(-\delta_0)} > 1 \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \quad (57)$$

By plugging (56) into (57), (57) is equivalent to

$$\frac{((1 - \alpha\mu - \alpha\delta_0)L\|w^*\| + G)^2 (-\alpha\rho G + (1 - \alpha L + \alpha\delta_0)^2\mu)^2}{((1 - \alpha\mu + \alpha\delta_0)L\|w^*\| + G)^2 (-\alpha\rho G + (1 - \alpha L - \alpha\delta_0)^2\mu)^2} > 1 \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \quad (58)$$



By taking the square root for both sides, (58) is equivalent to

$$\frac{((1 - \alpha\mu - \alpha\delta_0)L\|w^*\| + G)(-\alpha\rho G + (1 - \alpha L + \alpha\delta_0)^2\mu)}{((1 - \alpha\mu + \alpha\delta_0)L\|w^*\| + G)(-\alpha\rho G + (1 - \alpha L - \alpha\delta_0)^2\mu)} > 1 \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \quad (59)$$

A sequence of equivalent transformations on (59) can be performed as follows:

$$\begin{aligned} (59) &\iff \frac{((1 - \alpha\mu - \alpha\delta_0) + \frac{G}{L\|w^*\|})(-\frac{\alpha\rho G}{\mu} + (1 - \alpha L + \alpha\delta_0)^2)}{((1 - \alpha\mu + \alpha\delta_0) + \frac{G}{L\|w^*\|})(-\frac{\alpha\rho G}{\mu} + (1 - \alpha L - \alpha\delta_0)^2)} > 1 \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \\ &\iff \frac{((1 - \alpha\mu + \frac{G}{L\|w^*\|}) - \alpha\delta_0)((-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2) + 2\alpha\delta_0(1 - \alpha L))}{((1 - \alpha\mu + \frac{G}{L\|w^*\|}) + \alpha\delta_0)((-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2) - 2\alpha\delta_0(1 - \alpha L))} > 1 \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \\ &\iff ((1 - \alpha\mu + \frac{G}{L\|w^*\|}) - \alpha\delta_0)((-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2) + 2\alpha\delta_0(1 - \alpha L)) > \\ &\iff ((1 - \alpha\mu + \frac{G}{L\|w^*\|}) + \alpha\delta_0)((-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2) - 2\alpha\delta_0(1 - \alpha L)) \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \\ &\iff ((1 - \alpha\mu + \frac{G}{L\|w^*\|})(-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2) - 2\alpha^2\delta_0^2(1 - \alpha L)) \\ &\iff + (2\alpha\delta_0(1 - \alpha L)(1 - \alpha\mu + \frac{G}{L\|w^*\|}) - \alpha\delta_0(-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2)) > \\ &\iff ((1 - \alpha\mu + \frac{G}{L\|w^*\|})(-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2) - 2\alpha^2\delta_0^2(1 - \alpha L)) \\ &\iff - (2\alpha\delta_0(1 - \alpha L)(1 - \alpha\mu + \frac{G}{L\|w^*\|}) - \alpha\delta_0(-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2)) \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \\ &\iff (2\alpha\delta_0(1 - \alpha L)(1 - \alpha\mu + \frac{G}{L\|w^*\|}) - \alpha\delta_0(-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2)) > 0 \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \\ &\iff 2(1 - \alpha L)(1 - \alpha\mu + \frac{G}{L\|w^*\|}) > -\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2 \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \quad (60) \end{aligned}$$

For (60)'s LHS, since  $\mu \leq L$  and  $\alpha \leq \frac{1}{2L}$ , we have

$$2(1 - \alpha L)(1 - \alpha\mu + \frac{G}{L\|w^*\|}) > 2(1 - \alpha L)(1 - \alpha L + 0) \geq 2(1 - \frac{1}{2L} \cdot L)(1 - \frac{1}{2L} \cdot L + 0) = \frac{1}{2} \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \quad (61)$$

For (60)'s RHS, since  $\alpha \leq \frac{1}{2L}$  and  $\delta_0 < \frac{1}{2\alpha}$ , we have

$$-\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2 < 0 + (1 - \frac{1}{2L} \cdot L)^2 + \alpha^2(\frac{1}{2\alpha})^2 = \frac{1}{2} \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}). \quad (62)$$

Taking the result of (61) and (62) altogether, (60) is true since

$$2(1 - \alpha L)(1 - \alpha\mu + \frac{G}{L\|w^*\|}) > \frac{1}{2} > -\frac{\alpha\rho G}{\mu} + (1 - \alpha L)^2 + \alpha^2\delta_0^2 \quad \forall \delta_0 \in (0, \frac{1}{2\alpha}).$$

The proven (55) is equivalent to the desired conclusion. The proof is complete.  $\square$

## B. Supplementary Experiment Details

This section provides more details about the experimental settings and hyper-parameter choices.

### B.1. Few-shot Classification

---

**Algorithm 3** MAML with inner- and outer-level regularization

---

**Require:** Datasets  $\mathcal{S} = \{\mathcal{S}_i^{\text{in}}, \mathcal{S}_i^{\text{out}}\}_{i=1}^m$ ; few-shot meta-query batch size  $K$ ; the number of training tasks sampled at each round  $r$ ; the total number of iterations  $T$ .

**Require:** Regularization term  $Reg(w, \mathcal{D})$ ; Inner-level regularization selector  $\sigma^{\text{in}} \in \{-1, 0, 1\}$ , Outer-level regularization selector  $\sigma^{\text{out}} \in \{-1, 0, 1\}$ .

- 1: Initialize the model parameters  $w^0$  randomly.
  - 2: **for**  $t = 0$  to  $T - 1$  **do**
  - 3:   Randomly select  $r$  tasks from the set of  $m$  available tasks with indices stored in  $\mathcal{B}_t$ .
  - 4:   **for** each sampled task  $\mathcal{T}_i$  **do**
  - 5:     Sample a size  $K$  support data batch  $\mathcal{D}_i^{t, \text{in}}$  from  $\mathcal{S}_i^{\text{in}}$ ;
  - 6:     Sample a size  $b$  query data batch  $\mathcal{D}_i^{t, \text{out}}$  from  $\mathcal{S}_i^{\text{out}}$ ;
  - 7:     (Inner-level) Compute adapted parameters with gradient descent:
  - 8:          $w_i^t := w^t - \alpha \nabla_{w^t} \left( \hat{\mathcal{L}}(w^t, \mathcal{D}_i^{t, \text{in}}) + \sigma^{\text{in}} Reg(w^t, \mathcal{D}_i^{t, \text{in}}) \right)$ ;
  - 9:     (Outer-level) SGD step for meta-model, save the per-task weight for meta-update:
  - 10:          $w_i^{t+1} := w^t - \beta_t \nabla_{w^t} \left( \hat{\mathcal{L}}(w_i^t, \mathcal{D}_i^{t, \text{out}}) + \sigma^{\text{out}} Reg(w_i^t, \mathcal{D}_i^{t, \text{out}}) \right)$ ;
  - 11:   **end for**
  - 12:   Meta-update  $w^{t+1} := \frac{1}{r} \sum_{i \in \mathcal{B}_t} w_i^{t+1}$
  - 13: **end for**
  - 14: **Return:**  $w^T$
- 

The experiment setup for Omniglot and Mini-ImageNet follows [1]. **Datasets.** For the few-shot classification task, we experiment on the Mini-Imagenet [10, 15] and Omniglot [8] datasets. The Mini-Imagenet [10] is sampled from ImageNet with 600 instances of 100 classes. Each image is resized into  $84 \times 84$ . In the experiment, the Mini-Imagenet dataset is split into 64 classes for training, 12 classes for validation, and 24 classes for testing. The Omniglot dataset is a collection of 1623 character classes with different alphabets. Each class in the dataset contains 20 instances. The classes are shuffled and divided into the training, validation, and test sets, with 1150, 50, and 423 instances in the experiment. **Models.** We use the classic 4-layer convolution backbone models [1, 4] in the experiments. Each convolution layer has conv-filters of  $3 \times 3$  size and is followed by batchnorm and max-pooling. For the Omniglot dataset, we use the backbone model with 64 filters in each convolition layer (i.e., the backbone is 64-64-64-64 conv model). For the empirical verification experiment on Mini-ImageNet, the 48-48-48-48 conv backbone model is adopted. And for the experiment that comparing Minimax-MAML and Minimax-MAML++ with other baseline methods on Mini-ImageNet, we use the 64-64-64-64 conv backbone model to make a fairer comparison with other methods. **Training.** All the MAML experiments take 5 inner-steps. In one experiment, the training takes 150 epochs for 64-64-64-64 conv model and 120 epochs for 48-48-48-48 conv model, and each epoch consists of 500 iterations. The task batch size for all Omniglot experiments is 16. Mini-Imagenet experiments use task batch sizes of 4 and 2 for 1-shot and 5-shot experiments, respectively. After each epoch, the model’s performance is evaluated on the validation set. When the training is complete, a prediction of the test set is made by the ensemble of the best 5 per-epoch-models on the validation set (following [1], all the MAML-type methods’ results are generated under this paradigm). The Adam optimizer is adopted for the model training, with a scheduled learning rate starting from 0.001,  $\beta_1 = 0.9$ , and  $\beta_2 = 0.99$ . Cross-entropy loss is adopted as the loss function for all the models in the experiments.

**Regularization:**

The pseudo-code for implementing MAML with inner/outer-level regularization in the experiment is shown in Algorithm 3. For the Few-shot image classification experiments, the MAML-type methods are sharing the same form of regularization objective (except the first verification experiment isolating the L2-Norm regularizer). The regularization is achieved by combining the L2-Norm regularization and output entropy regularization, i.e., the regularization term  $Reg(w, \mathcal{D})$  in Algorithm 3

is given by

$$Reg(w, \mathcal{D}) = -\gamma^{entropy} H(w, \mathcal{D}) + \gamma^{norm} \frac{1}{2} \|w\|^2, \quad (63)$$

where  $H(w, \mathcal{D})$  denotes information entropy of the output generated by model  $w$  for data batch  $\mathcal{D}$ :

$$H(w, \mathcal{D}) = -\mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}} \sum_{i=1}^K p_w(y = i | \mathbf{x}) \log p_w(y = i | \mathbf{x}),$$

where  $K$  is the number of classes and  $p_w(y = i|x)$  represents the probability of prediction to class  $i$  generated by model  $w$ . Entropy represents the diversity of model output, and using negative entropy as a regularization objective can encourage the model to make more conservative outputs and suppress overconfident outputs, thereby avoiding overfit. Entropy regularization is also referred to as label-smooth sometimes. (*Since entropy needs to be maximized in order to serve as a regularizer, it is necessary to include a negative sign when incorporating it into the loss function.*)

In order to explain entropy regularizer more specifically, we provide a PyTorch-based sample implementation here:

```

1 class Self_Entropy(torch.nn.Module):
2     def __init__(self, reduction = True):
3         super(Self_Entropy, self).__init__()
4         self.reduction = reduction
5
6     def forward(self, x):
7         b = F.softmax(x, dim=1) * F.log_softmax(x, dim=1)
8         if self.reduction:
9             b = -1.0 * b.mean()
10        else:
11            b = -1.0 * b.sum()
12        return b
13
14 self_entropy = Self_Entropy()

```

$\gamma^{entropy}$  and  $\gamma^{norm}$  in (63) are positive hyper-parameters controlling the regularization rate. We use  $\gamma^{entropy} = 2.0$  and  $\gamma^{norm} = 5e - 4$  for all of the few-shot classification experiments.

In Algorithm 3, the selectors  $\sigma^{in}$  and  $\sigma^{out}$  respectively determine the type of regularization for the inner-level and outer-level. The values of  $\sigma^{in}$  and  $\sigma^{out}$  can be 1, 0 or -1, corresponding to ordinary regularization, non-regularization, and inverted regularization respectively. For instance, in the empirical verification experiment that uses a combined regularizer, to evaluate the effect of inverted inner-level regularization, we set  $\{\sigma^{in}, \sigma^{out}\}$  as  $\{-1, 0\}$ . Similarly, we set  $\{\sigma^{in}, \sigma^{out}\}$  as  $\{1, 0\}$  to evaluate the effect of ordinary regularization at inner-level.

In terms of other regularization types, we have  $\{\sigma^{in}, \sigma^{out}\} = \{0, 1\}$  for *regularize the outer-level*,  $\{\sigma^{in}, \sigma^{out}\} = \{1, 1\}$  for *regularize the loss function* (in the limited-tasks experiment), and  $\{\sigma^{in}, \sigma^{out}\} = \{-1, 1\}$  for *minimax-meta regularization*. Since the original MAML doesn't have any regularization, it is equivalent to having  $\delta^{in} = 0$  and  $\delta^{out} = 0$ . (see Table 1 and 2 in the main paper.)

It is worth noting that we only add the inner-level inverted regularization during the training phase, and we do not use it for the meta-testing phase. Specifically, during the meta-testing phase, which evaluates the performance of the learned meta-model on new tasks, we only adapt the model without any additional regularization to avoid influencing its task-specific performance.

### Implementation.

The implementation of inner- and outer-level regularizations is simple and straightforward, and often involves only modifications to loss functions. Assuming that cross-entropy is used as the classification loss, and the combined regularization term  $Reg(w, \mathcal{D}) = -\gamma^{entropy} H(w, \mathcal{D}) + \gamma^{norm} \frac{1}{2} \|w\|^2$  is adopted, we provide a PyTorch implementation example of Minimax-Meta Regularization to further explain the regularizations and demonstrate the simplicity of implementation.

During training, at the inner-level, the invertedly regularized loss  $\hat{\mathcal{L}}(w^t, \mathcal{D}_i^{t, in}) - Reg(w^t, \mathcal{D}_i^{t, in})$  now can be expressed by  $\hat{\mathcal{L}}(w^t, \mathcal{D}_i^{t, in}) - (-\gamma^{entropy} H(w^t, \mathcal{D}_i^{t, in}) + \gamma^{norm} \frac{1}{2} \|w^t\|^2)$ , which could be implemented as:

```

1 # inner-loop training loss of MAML, with inverted regularization.
2 loss = F.cross_entropy(preds, y) - (- gamma_e * self_entropy(preds) + gamma_n * l2_norm(weights))

```

where  $preds$  is the model’s prediction for the input data batch,  $y$  is the true label batch and  $weights$  stores the weight values of the model.

Similarly, at the outer-level, the ordinarily regularized loss  $\hat{\mathcal{L}}(w_i^t, \mathcal{D}_i^{t, out}) + Reg(w_i^t, \mathcal{D}_i^{t, out})$  now can be expressed by  $\hat{\mathcal{L}}(w_i^t, \mathcal{D}_i^{t, out}) + (-\gamma^{entropy} H(w_i^t, \mathcal{D}_i^{t, out}) + \gamma^{norm} \frac{1}{2} \|w_i^t\|^2)$ , which could be implemented as:

```
1 # outer-loop training loss of MAML, with ordinary regularization.
2 loss = F.cross_entropy(preds, y) + (- gamma_e * self_entropy(preds) + gamma_n * l2_norm(weights))
```

Since the built-in norm/weight-decay methods in popular libraries usually do not support negative parameters, the  $l2\_norm$  function may require manual implementation, but it is also easy to accomplish.

When training is complete, during the testing phase, which evaluates the performance of the learned meta-model on new tasks, we adapt the model to each new task without any additional regularization:

```
1 # meta-testing phase loss of MAML, without additional regularization
2 loss = F.cross_entropy(preds, y)
```

The aforementioned implementation example can be readily incorporated into the widely-used open-source MAML directory “How to train your MAML in Pytorch” proposed in [1], which serves as the basis for our experimental setup.

## B.2. Few-shot Regression

The experiment setting follows the few-shot regression experiment in [12]. **Datasets** One synthetic and three real-world few-shot regression datasets are considered. The synthetic dataset is created by a 2-dimensional mixture of Cauchy distributions plus random GP functions. One real-world dataset is SwissFEL [9] which corresponds to Swiss Free Electron Laser’s calibration sessions. Another two datasets are from the PhysioNet 2012 challenge [13], which contains time-series data related to patients’ health metrics, in particular, the Glasgow Coma Scale (GCS) and the hematocrit value (HCT). Cauchy contains 20 tasks, and each task consists of 20 samples. SwissFel contains 5 tasks, and each task consists of 200 samples. Each Physionet dataset contains 100 tasks, and each task consists of 4 ~ 24 samples. **Models** We use a fully-connected neural network with 4 layers with each 32 neurons as the base-learner model, aligning with the base-learner structure adopted by other baseline methods. *ReLU* is used as activation. MAML takes 3 inner steps in our experiment.

**Regularization** For regression problems, the output entropy used in the classification experiments cannot again be used as the regularization objective. So we adopt L2-Norm as the only regularization objective. Let  $\gamma$  be the parameter controlling the magnitude and direction of the L2-Norm regularization, the inner-level regularization rate parameter  $\gamma^{in}$  is set to negative (inverted regularization) and the outer-level regularization rate parameter  $\gamma^{out}$  is set to positive (ordinary regularization). We use hyper-parameter search to select the value of parameters. Specifically, we use  $\{\gamma^{in}, \gamma^{out}\} = \{-1e-3, 1e-3\}$ ,  $\{\gamma^{in}, \gamma^{out}\} = \{-1e-2, 1e-2\}$ ,  $\{\gamma^{in}, \gamma^{out}\} = \{-5e-3, 5e-3\}$ , and  $\{\gamma^{in}, \gamma^{out}\} = \{-5e-2, 5e-2\}$  for Cauchy, SwissFel, Physionet-GCS, and Physionet-HCT experiment respectively. The number of iterations in each experiment is determined using the validation set.

## C. Supplementary Experimental Analysis

Due to the space limitation of the main paper, we provide supplementary experimental results in this section, including an additional experiment on Mini-ImageNet few-shot classification with limited tasks, an additional experiment on Meta-dataset with the first-order method and larger backbone, and an additional experiment on meta-reweighting with Minimax-Meta Regularization for robust learning.

### C.1. Mini-ImageNet Few-shot Classification with Limited Tasks

To further illustrate the generalization ability of Meta-Minimax regularization, we conduct an experiment to compare it with other common regularization strategies on meta-learning with the limited number of training tasks. The fewer the task number is, the easier the meta-model would overfit.

In the implementation of N-way few-shot classification experiments, in the training phase, each task is sampled by combining N training classes as one N-way classification task. That is, for a dataset with M training classes available, there would be accordingly  $\binom{M}{N}$  training tasks available. So we could restrict the number of training tasks by restricting the number of training classes.

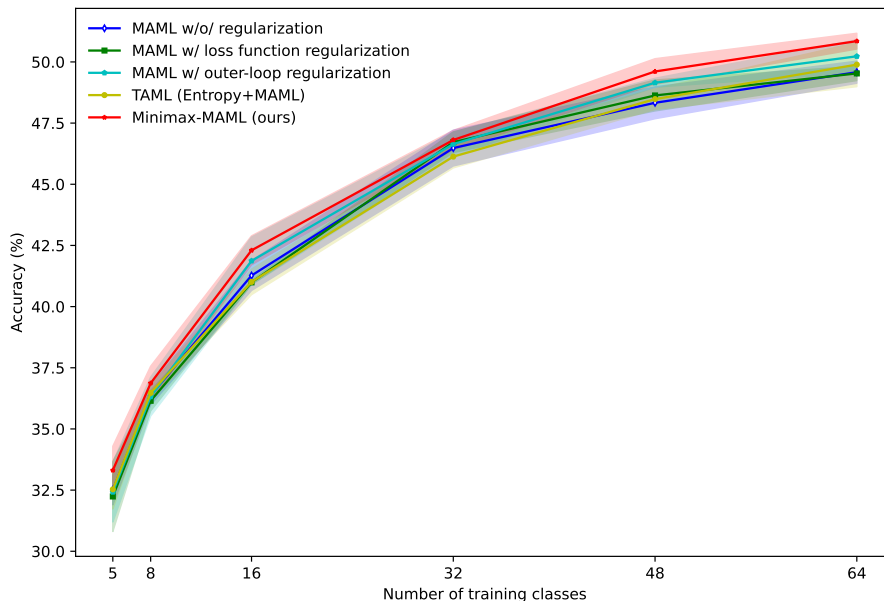


Figure 1. Test accuracies (%) with varying training classes number. The shaded region denotes the 95% confidence interval.

We take Mini-ImageNet 5-way 1-shot as the experiment scenario. The experiment setting follows the same setting of Mini-ImageNet empirical verification experiments. In the original experiment, there are 64 classes available for training. We restrict the number of training classes to 48/32/16/8/5 in this limited classes experiment. And Meta-Minimax regularization is compared with original MAML and MAML with common regularization: *MAML w/ outer-loop regularization*, *MAML w/ loss function regularization* (*MAML w/ loss function regularization* means simply adding ordinary reg term in loss function at both inner- and outer-level during the training). We also implemented TAML(Entropy+MAML) proposed by [7] for comparison.

Figure 1 shows the accuracy curve of experiment results with the varying number of training classes. The result suggests Meta-Minimax regularization continuously outperforms other methods under the limited task number scenario and improves the accuracy to a certain margin even under a very small task number.

## C.2. Meta-Dataset Few-shot Classification Experiment

To test the effectiveness of Minimax-Regularization on larger backbones and to validate if it is fitful for first-order meta-learning methods, We further conduct an experiment using first-order MAML (fo-MAML) and ResNet-12 backbone on Meta-Dataset [14]. Meta-Dataset creates a dataset of datasets benchmark for meta-learning. In our experiment, we only train the model on the ILSVRC training set and test on the ILSVRC testing test and other 8 datasets (Omniglot, Aircraft, Birds, Textures, QuickDraw, VGG Flower, Traffic, MSCOCO). The experiment settings follow the benchmark proposed in [14], and we use the open PyTorch repository provided by [2] to build the experiment. We implement Minimax-Meta Regularization for fo-MAML to compare it with the original version.

Each model takes 6 inner-loop steps in the training phase, and additional 5 inner-loop steps are adopted for the testing phase. The inner stepsize is set to 0.3 for each model. Adam optimizer is adopted for the meta-model updating with a learning rate of 0.00025. Each model completed 10000 training updates. We repeat the experiment for 3 independent runs and report the mean accuracies and 95% confidence intervals.

Like in the Mini-ImageNet and Omniglot experiments, we adopt the entropy & L2-Norm combined regularization term to achieve the Minimax-Meta Regularization. We use  $\gamma^{entropy} = 2.0$  and  $\gamma^{norm} = 3e-5$  as the reg magnitude coefficients for both inner- and outer-level regularizations.

Table 6 shows the experiment results. The results suggest that the Minimax-fo-MAML generalizes better on all 9 testing

datasets.

### C.3. Meta-reweighting with Minimax Regularization for Robust Learning

To verify the general effectiveness of our proposed methods, we further conduct experiments on the meta-learning problem of meta-reweighting for robust learning.

#### C.3.1 Experimental Setup

For this experiment, we evaluate the performance of our proposed method and baselines on a robust-learning task: the noisy MNIST dataset. The dataset is created by randomly flipping the labels of 40% of the training images, resulting in 10000 training images with 40% incorrectly labeled data. Each image has a dimension of 28x28, and the task is to classify them into ten handwritten digits (0 ~ 9). There is also a clean validation set consisting of 100 correctly labeled images with balanced categories available for helping the training process on the noisy set.

---

#### Algorithm 4 Minimax Meta-Reweighting.

---

**Require:** model  $\theta_0$ , noisy training set  $D_f$ , clean validation set  $D_g$ , training batch size  $n$ , validation batch size  $m$ , inner-level regularization parameter  $\gamma^{in}$ , outer-level regularization parameter  $\gamma^{out}$

**Ensure:**  $\theta_T$

- 1: **for**  $t = 0$  to  $T - 1$  **do**
  - 2:   Sample a  $n$ -size mini-Batch data  $\{X_f, y_f\}$  from  $D_f$ ;
  - 3:   Sample a  $m$ -size mini-Batch data  $\{X_g, y_g\}$  from  $D_g$ ;
  - 4:   Forward  $X_f$  using model  $\theta_t$ , get predicted labels  $\hat{y}_f$ ;
  - 5:   Set temporary example weights to zero:  $\epsilon = 0$ ;
  - 6:   Calculate weighted loss on noisy data batch:  $l_f = \sum_{i=1}^n \epsilon_i C(y_{f,i}, \hat{y}_{f,i})$ ;
  - 7:   Calculate  $\hat{\theta}_t = \theta_t - \alpha \nabla_{\theta_t} l_f$ ;
  - 8:   Forward  $X_g$  using model  $\hat{\theta}_t$ , get predicted labels  $\hat{y}_g$ ;
  - 9:   Evaluate loss on clean data batch, with inverted entropy reg:  
 $l_g = \frac{1}{m} \sum_{i=1}^m (C(y_{g,i}, \hat{y}_{g,i}) + \gamma^{in} Entropy(\hat{y}_{g,i}))$ ;
  - 10:   Calculate new example weights  $\tilde{w} = \max(-\nabla_{\epsilon} l_g, 0)$ , and normalize  $w = \frac{\tilde{w}}{\sum_j \tilde{w} + \delta (\sum_j \tilde{w})}$ ;
  - 11:   Calculate new weighted loss on noisy data batch, with ordinary entropy reg:  
 $\hat{l}_f = \sum_{i=1}^n w_i (C(y_{f,i}, \hat{y}_{f,i}) - \gamma^{out} Entropy(\hat{y}_{f,i}))$ ;
  - 12:    $\theta_{t+1} \leftarrow \text{OptimizerStep}(\theta_t, \nabla_{\theta_t} \hat{l}_f)$ ;
  - 13: **end for**
- 

The basic robust-learning baseline we evaluate here is Meta-Reweighting introduced in [11]. The Meta-Reweighting algorithm learns to assign weights to training examples for robust learning. To determine the example weights, Meta-Reweighting performs a meta gradient descent step on the mini-batch example weights (which are initialized from zero) to minimize the loss on a clean, unbiased validation set.

Our method adds the Minimax-Meta Regularization on top of Meta-Reweighting. We add ordinary regularization at the outer-level, where the optimal weights are adopted for meta-update. And inverted regularization is added at the inner-level, where the weighted inner-model fits the clean unbiased validation set for optimal weight calculation. Intuitively, through the meta-weighted learning process, such a regularization method makes the model become more conservative when updating based on the noisy training data in the outer loop and values the diversity of predictions more, thereby resisting overfit.

Table 6. Few-shot classification results on **Meta-Dataset** using models trained on ILSVRC. Backbone: **ResNet-12**

Datasets except ILSVRC are only used for testing, and we report the test accuracy with a 95% confidence interval. Each model completed 10000 training updates.

Method	ILSVRC (test)	Omniglot	Aircraft	Birds	Textures	QuickDraw	VGG Flower	Traffic	MSCOCO
fo-MAML	38.24±2.30	44.75±6.26	28.06±2.43	37.64±3.56	39.41±4.50	42.57±3.79	58.55±5.20	36.62±2.85	42.38±5.09
<b>Minimax-fo-MAML(ours)</b>	<b>40.53±1.54</b>	<b>68.43±3.53</b>	<b>30.95±2.97</b>	<b>41.09±0.40</b>	<b>45.12±1.41</b>	<b>51.57±2.68</b>	<b>66.23±0.89</b>	<b>38.83±2.71</b>	<b>45.15±0.85</b>

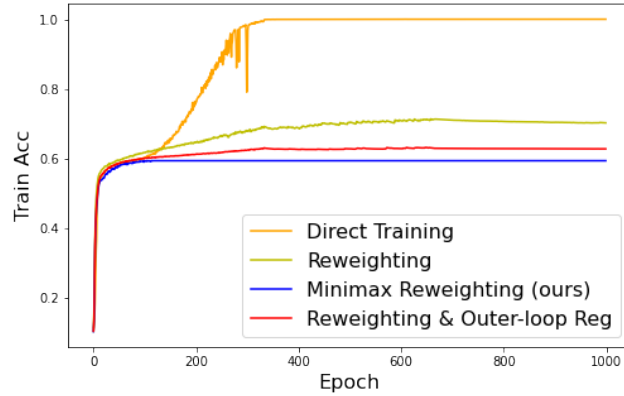


Figure 2. Training accuracy curve. Since 40% of training samples are incorrectly labeled, the model keeps a training accuracy of around 60% would be considered resistant to overfitting during the training.

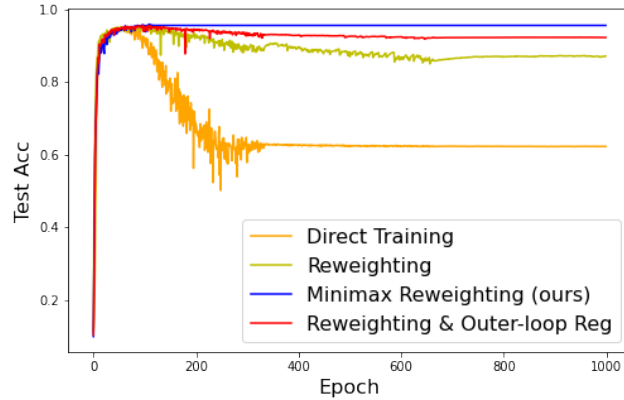


Figure 3. Test accuracy curve. Since the test dataset is clean, the model that can maintain a higher test accuracy is considered with better learning robustness (less affected by the noise in the training set).

At the same time, the inner model was encouraged to make sharper predictions on the clean validation set by the inverted regularization.

The regularization objective used in our method is maximizing output entropy (minimizing output entropy at the inner-level). We call our method Minimax Meta-Reweighting. The pseudo-code for implementation is shown in Algorithm 4. In our experiment, we use a  $\gamma^{in} = 0.25$  and  $\gamma^{out} = 2.0$ .

For each method in the experiment, we use the LeNet-5 as the backbone model and train the model for 1000 epochs. The learning rates for the first 1/3, the middle 1/3, and the last 1/3 training epochs are set to 1e-2, 1e-3, and 1e-4, respectively.

### C.3.2 Results and Analysis

Under this setting, models are extremely prone to overfit the noisy dataset during the training phase. To understand the models' performance, we could look at the training and testing curves in Figure 2 and 3. Since the training set is noisy, models overfitted to the train set would show significant performance deduction on the clean test set.

From the perspective of robust learning, the *direct training* method sets the lower performance bound to some extent. Since it does not have any denoising ability, it quickly overfits the training set during the training. It reaches peak accuracy on the clean test set around the 80th epoch, and starts to overfit after that. We could identify the overfitting characteristic from the training and testing accuracy curve. Since 40% of the labels in the training set are incorrect, once the model starts

to predict the training data with an accuracy larger than 60%, it fits the distribution of the noisy training data instead of the ground truth distribution. At the same time, the performance deduction on the clean test set would also start. Finally, we could observe the training accuracy and testing accuracy of the directly trained model to converge to nearly 100% and 60%, respectively, which indicates a complete overfit. On the contrary, the model with optimal learning robustness should not overfit the train set, keep a train accuracy value close to 60% and maintain the optimal performance on the clean test set.

Compared to direct training, the training curve of Meta-Reweighting baseline [11] shows a significant improvement in the learning robustness. However, it still suffers from overfitting. It neither completely overfits the training dataset nor ignores all the noises; its training accuracy converges to around 70%. After around the 100th epoch, the Meta-Reweighting model experienced continual test accuracy deduction and finally maintained test accuracy at around 87.5%.

Minimax-Reweighting nearly reached the optimal learning robustness under this setting. The training accuracy of Minimax-Reweighting stuck at around 60% with rarely any change throughout the training phase. And the testing accuracy maintained a peak value of around 95.5% without observable deduction.

To further evaluate the effectiveness of Minimax-Reweighting, we implemented the outer-loop-only regularization on top of the Meta-Reweighting algorithm to make comparisons. While this approach did show improvement from the baseline method, it was unable to achieve the same level of performance as Minimax-Meta Regularization, as shown in Figure 2 and 3.



## References

- [1] Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *arXiv preprint arXiv:1810.09502*, 2018. 16, 18
- [2] Malik Boudiaf, Ziko Imtiaz Masud, Jérôme Rony, Jose Dolz, Ismail Ben Ayed, and Pablo Piantanida. Mutual-information based few-shot classification, 2021. 19
- [3] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Generalization of model-agnostic meta-learning algorithms: Recurring and unseen tasks. *Advances in Neural Information Processing Systems*, 34, 2021. 4, 5
- [4] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, pages 1126–1135. PMLR, 2017. 4, 16
- [5] Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *International Conference on Machine Learning*, pages 1920–1930. PMLR, 2019. 1, 2
- [6] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR, 2016. 1
- [7] Muhammad Abdullah Jamal and Guo-Jun Qi. Task agnostic meta-learning for few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11719–11727, 2019. 19
- [8] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015. 16
- [9] Christopher J Milne, Thomas Schietinger, Masamitsu Aiba, Arturo Alarcon, Jürgen Alex, Alexander Anghel, Vladimir Arsov, Carl Beard, Paul Beaud, Simona Bettoni, et al. Swissfel: the swiss x-ray free electron laser. *Applied Sciences*, 7(7):720, 2017. 18
- [10] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017. 16
- [11] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *International Conference on Machine Learning*, pages 4334–4343. PMLR, 2018. 20, 22
- [12] Jonas Rothfuss, Vincent Fortuin, Martin Josifoski, and Andreas Krause. Pacoh: Bayes-optimal meta-learning with pac-guarantees. In *International Conference on Machine Learning*, pages 9116–9126. PMLR, 2021. 18
- [13] Ikaro Silva, George Moody, Daniel J Scott, Leo A Celi, and Roger G Mark. Predicting in-hospital mortality of icu patients: The physionet/computing in cardiology challenge 2012. In *2012 Computing in Cardiology*, pages 245–248. IEEE, 2012. 18
- [14] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, et al. Meta-dataset: A dataset of datasets for learning to learn from few examples. *arXiv preprint arXiv:1903.03096*, 2019. 19
- [15] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29:3630–3638, 2016. 16