

---

# A Simple Yet Effective Strategy to Robustify the Meta Learning Paradigm

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Meta learning is a promising paradigm to enable skill transfer across tasks. Most  
2        previous methods employ the empirical risk minimization principle in optimization.  
3        However, the resulting worst fast adaptation to a subset of tasks can be catastrophic  
4        in risk-sensitive scenarios. To robustify fast adaptation, this paper optimizes meta  
5        learning pipelines from a distributionally robust perspective and meta trains models  
6        with the measure of expected tail risk. We take the two-stage strategy as heuristics  
7        to solve the robust meta learning problem, controlling the worst fast adaptation  
8        cases at a certain probabilistic level. Experimental results show that our simple  
9        method can improve the robustness of meta learning to task distributions and reduce  
10       the conditional expectation of the worst fast adaptation risk.

## 11    1 Introduction

12       The past decade has witnessed the remarkable  
13       progress of deep learning in real-world appli-  
14       cations (LeCun et al., 2015). However, train-  
15       ing deep learning models requires an enormous  
16       dataset and intensive computational power. At  
17       the same time, these pre-trained models can fre-  
18       quently encounter deployment difficulties when  
19       the dataset’s distribution drifts in testing time  
20       (Lesort et al., 2021).

21       As a result, the paradigm of *meta learning* or  
22       *learning to learn* is proposed and impacts the ma-  
23       chine learning scheme (Finn et al., 2017), which  
24       leverages past experiences to enable fast adapta-  
25       tion to unseen tasks. Moreover, in the past few  
26       years, there has grown a large body of meta learning methods to find plausible strategies to distill  
27       common knowledge into separate tasks (Finn et al., 2017; Duan et al., 2016; Garnelo et al., 2018a).

28       Notably, most previous work concentrates merely on the fast adaptation strategies and employs the  
29       standard risk minimization principle, *e.g.* the empirical risk minimization, ignoring the difference  
30       between tasks in fast adaptation. Given the sampled batch from the task distribution, the standard meta  
31       learning methods weight tasks equally in fast adaptation. Such an implementation raises concerns  
32       in some real-world scenarios, when worst fast adaptation is catastrophic in a range of risk-sensitive  
33       applications (Johannsmeier et al., 2019; Jaafra et al., 2019). For example, in robotic manipulations,  
34       humanoid robots (Duan, 2017) can quickly leverage past motor primitives to walk on plain roads but  
35       might suffer from tribulation doing this on rough roads.

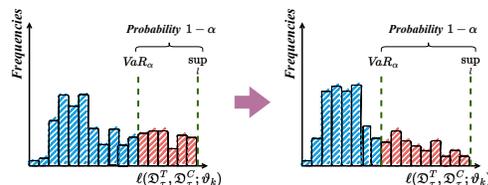


Figure 1: **Illustrations of Distributionally Robust Fast Adaptation.** Shown are histograms of meta risk function values  $\ell(\mathcal{D}_\tau^T, \mathcal{D}_\tau^C; \vartheta)$  in the task distribution  $p(\tau)$ . Given a probability  $\alpha$ , we optimize meta learning model parameters  $\vartheta$  to decrease the risk quantity CVaR<sub>α</sub> in Definition (3).

36 **Research Motivations.** Instead of seeking novel fast adaptation strategies, we take more interest  
37 in optimization principles for meta learning. Given the meta trained model, this paper stresses the  
38 performance difference in fast adaptation to various tasks as an indispensable consideration. As *the*  
39 *concept of robustness in fast adaptation* has not been sufficiently explored from the task distribution  
40 perspective, researching this topic has more practical significance and deserves more attention in  
41 meta learning. Naturally, we raise the question below:

42 *Can we reconsider the meta learning paradigm through the lens of risk distribution, and are there*  
43 *plausible measures to enhance the fast adaptation robustness in some vulnerable scenarios?*

44 **Developed Methods.** In an effort to address the above concerns and answer these questions, *we*  
45 *reduce robust fast adaptation in meta learning to a stochastic optimization problem within the*  
46 *principle of minimizing the expected tail risk, e.g., conditional value-at-risk (CVaR) (Rockafellar*  
47 *et al., 2000).* To tractably solve the problem, we adopt a two-stage heuristic strategy for optimization  
48 with the help of crude Monte Carlo methods (Kroese and Rubinstein, 2012) and give some theoretical  
49 analysis. In each optimization step, the algorithm estimates the value-at-risk (VaR) (Rockafellar et al.,  
50 2000) from a meta batch of tasks and screens and optimizes a percentile of task samples vulnerable  
51 to fast adaptation. As illustrated in Fig. (1), such an operation is equivalent to iteratively reshaping  
52 the task risk distribution to increase robustness. The consequence of optimizing the risk function  
53 distributions  $\vartheta_k \rightarrow \vartheta_{k+1}$  is to transport the probability mass in high-risk regions to the left side  
54 gradually. In this manner, the distribution of risk functions in the task domain can be optimized  
55 toward the anticipated direction that controls the worst-case fast adaptation at a certain probabilistic  
56 level.

57 **Outline & Primary Contributions.** We overview related meta learning and robust optimization  
58 work in Section (2). Section (3) introduces general notations and describes meta learning optimization  
59 objectives together with typical models. The distributionally robust meta learning problem is presented  
60 together with a heuristic optimization strategy in Section (4). We report experimental results and  
61 analysis in Section (5), followed by conclusions and limitations in Section (6). **Our primary**  
62 **contribution** is two-fold:

- 63 1. We recast the robustification of meta learning to a distributional optimization problem. The  
64 resulting framework minimizes the conditional expectation of task risks, namely the tail risk,  
65 which unifies vanilla meta-learning and worst-case meta learning frameworks.
- 66 2. To resolve the robust meta learning problem, we adopt the heuristic two-stage strategy and  
67 demonstrate its improvement guarantee. Experimental results show the effectiveness of our  
68 method, enhancing fast adaptation robustness and mitigating the worst-case performance.

## 69 2 Literature Review

70 **Meta Learning Methods.** In practice, meta learning enables fast learning (adaptation to unseen  
71 tasks) via slow learning (meta training in a collection of tasks). There exist different families of meta  
72 learning methods. The optimization-based methods, such as model agnostic meta learning (MAML)  
73 (Finn et al., 2017) and its variants (Finn et al., 2018; Rajeswaran et al., 2019; Grant et al., 2018;  
74 Vuorio et al., 2019; Abbas et al., 2022), try to find the optimal initial parameters of models and then  
75 execute gradient updates over them to achieve adaptation with a few examples. The context-based  
76 methods, e.g. conditional neural processes (CNPs) (Garnelo et al., 2018a), neural processes (NPs)  
77 (Garnelo et al., 2018b) and extensions (Gordon et al., 2019; Foong et al., 2020; Gondal et al., 2021;  
78 Wang and van Hoof, 2022; Wang et al., 2023), learn the representation of tasks in the function  
79 space and formulate meta learning models as exchangeable stochastic processes. The metrics-based  
80 methods (Snell et al., 2017; Allen et al., 2019; Bartunov and Vetrov, 2018) embed tasks in a metric  
81 space and can achieve competitive performance in few-shot classification tasks. Other methods like  
82 memory-augmented models (Santoro et al., 2016), recurrent models (Duan et al., 2016) and hyper  
83 networks (Zhao et al., 2020; Beck et al., 2023) are also modified for meta learning purposes.

84 **Robust Optimization.** When performing robust optimization for downstream tasks in deep learning,  
85 we can find massive work concerning the adversarial input noise (Goodfellow et al., 2018; Goel et al.,  
86 2020; Ren et al., 2021), or the perturbation on the model parameters (Goodfellow et al., 2014; Kurakin  
87 et al., 2016; Liu et al., 2018; Silva and Najafirad, 2020). In contrast, this paper studies the robustness  
88 of fast adaptation in meta learning. In terms of robust principles, the commonly considered one is

89 the worst-case optimization (Olds, 2015; Zhang et al., 2020; Tay et al., 2022). For example, Collins  
 90 et al. (2020) conducts the worst-case optimization in MAML to obtain the robust meta initialization.  
 91 Considering the influence of adversarial examples, Goldblum et al. (2019) propose to adversarially  
 92 meta train the model for few-shot image classification. Wang et al. (2020) adopt the worst-case  
 93 optimization in MAML to increase the model robustness by injecting adversarial noise to the input.  
 94 However, distributionally robust optimization (Rahimian and Mehrotra, 2019) is rarely examined in  
 95 the presence of the meta learning task distribution.

### 96 3 Preliminaries

97 **Notations.** Consider the distribution of tasks  $p(\tau)$  for meta learning and denote the task space by  $\mathcal{T}$ .  
 98 Let  $\tau$  be a task sampled from  $p(\tau)$  with  $\mathcal{T}$  the set of all tasks. We denote the meta dataset by  $\mathcal{D}$ . For  
 99 example, in few-shot regression problems,  $\mathcal{D}$  refers to a set of data points  $f(x_i, y_i)g_{i=1}^m$  to fit.

100 Generally,  $\mathcal{D}$  are processed into the context set  $\mathcal{D}^C$  for fast adaptation, and the target set  $\mathcal{D}^T$  for  
 101 evaluating adaptation performance. As an instance, we process the dataset  $\mathcal{D} = \mathcal{D}^C \cup \mathcal{D}^T$  with a  
 102 fixed partition in MAML (Finn et al., 2017).  $\mathcal{D}^C$  and  $\mathcal{D}^T$  are respectively used for the inner loop and  
 103 the outer loop in model optimization.

104 **Definition 1 (Meta Risk Function)** With the task  $\tau \in \mathcal{T}$  and the pre-processed dataset  $\mathcal{D}$  and the  
 105 model parameter  $\vartheta \in \mathcal{V}$ , the meta risk function is a map  $\ell : \mathcal{D} \times \mathcal{V} \rightarrow \mathbb{R}^+$ .

106 In meta learning, the meta risk function  $\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta)$ , e.g. instantiations in Example (1)/(2), is  
 107 to evaluate the model performance after fast adaptation. Now we turn to the commonly-used risk  
 108 minimization principle, which plays a crucial role in fast adaptation. To summarize, we include the  
 109 vanilla and worst-case optimization objectives as follows.

110 **Expected Risk Minimization for Meta Learning.** The objective that applies to most vanilla meta  
 111 learning methods can be formulated in Eq. (1), and the optimization executes in a distribution over  
 112 tasks  $p(\tau)$ . The Monte Carlo estimate corresponds to the *empirical risk minimization* principle.

$$113 \quad \min_{\vartheta} E(\vartheta) := \mathbb{E}_{p(\tau)} \left[ \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \right] \quad (1)$$

114 Here  $\vartheta$  is the parameter of meta learning models, which includes parameters for common knowledge  
 115 shared across all tasks and for fast adaptation. Furthermore, the task distribution heavily influences  
 116 the direction of optimization in meta training.

117 **Worst-case Risk Minimization for Meta Learning.** This also considers meta learning in the task  
 118 distribution  $p(\tau)$ , but the worst case in fast adaptation is the top priority in optimization.

$$119 \quad \min_{\vartheta} \max_{\tau \in \mathcal{T}} \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \quad (2)$$

120 The optimization objective is built upon the min-max framework, advancing the robustness of meta  
 121 learning to the worst case. Approaches like TR-MAML (Collins et al., 2020) sample the worst  
 122 task in a batch to meta train with gradient updates. Nevertheless, this setup might result in a highly  
 123 conservative solution where the worst case only happens with an incredibly lower chance.

### 124 4 Distributionally Robust Fast Adaptation

125 This section starts with the concept of risk measures and the derived meta learning optimization  
 126 objective. Then a heuristic strategy is designed to approximately solve the problem. Finally, we  
 127 provide two examples of distributionally robust meta learning methods.

128 **4.1 Meta Risk Functions as Random Variables**

**Assumption 1** The meta risk function  $\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)$  is  $\beta$ -Lipschitz continuous w.r.t.  $\vartheta$ , which suggests: there exists a positive constant  $\beta$  such that  $\forall \vartheta, \vartheta^0 \in \mathcal{G}$ :

$$|\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta^0)| \leq \beta \|\vartheta - \vartheta^0\|.$$

Let  $(\mathcal{T}, F, P)$  denote a probability measure over the task space, where  $F$  corresponds to a  $\sigma$ -algebra on the subsets of  $\mathcal{T}$ . And we have  $(\mathbb{R}^+, B)$  a probability measure over the non-negative real domain for the previously mentioned meta risk function  $\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)$  with  $B$  a Borel  $\sigma$ -algebra. For any  $\vartheta \in \mathcal{G}$ , the meta learning operator  $\mathcal{M}_\# : \mathcal{T} \rightarrow \mathbb{R}^+$  is defined as:

$$\mathcal{M}_\# : \tau \mapsto \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta).$$

129 In this way,  $\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)$  can be viewed as a random variable to induce the distribution  
 130  $p(\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta))$ . Further, the cumulative distribution function can be formulated as  $F_\cdot(l; \vartheta) =$   
 131  $P(\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \leq l; \tau \in \mathcal{T}, l \in \mathbb{R}^+)$  w.r.t. the task space. Note that  $F_\cdot(l; \vartheta)$  implicitly depends  
 132 on the model parameter  $\vartheta$ , and we cannot access a closed-form in practice.

**Definition 2 (Value-at-Risk)** Given the confidence level  $\alpha \in [0, 1]$ , the task distribution  $p(\tau)$  and the model parameter  $\vartheta$ , the VaR (Rockafellar et al., 2000) of the meta risk function is defined as:

$$\text{VaR}_\alpha[\ell(T, \vartheta)] = \inf_{l \in \mathbb{R}^+} l \text{ s.t. } F_\cdot(l; \vartheta) \geq 1 - \alpha, \tau \in \mathcal{T}.$$

**Definition 3 (Conditional Value-at-Risk)** Given the confidence level  $\alpha \in [0, 1]$ , the task distribution  $p(\tau)$  and the model parameter  $\vartheta$ , we focus on the constrained domain of the random variable  $\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)$  with  $\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \leq \text{VaR}_\alpha[\ell(T, \vartheta)]$ . The conditional expectation of this is termed as conditional value-at-risk (Rockafellar et al., 2000):

$$\text{CVaR}_\alpha[\ell(T, \vartheta)] = \int_0^{1-\alpha} l dF_\cdot(l; \vartheta),$$

where the normalized cumulative distribution is as follows:

$$F_\cdot(l; \vartheta) = \begin{cases} 0, & l < \text{VaR}_\alpha[\ell(T, \vartheta)] \\ \frac{F_\cdot(l; \vartheta)}{1 - \alpha}, & l \geq \text{VaR}_\alpha[\ell(T, \vartheta)]. \end{cases}$$

133 This results in the normalized probability measure  $(\mathcal{T}_\alpha, F_\cdot, P_\cdot)$  over the task space, where  
 134  $\mathcal{T}_\alpha := \bigcup_{\vartheta \in \mathcal{G}} \text{VaR}_\alpha[\ell(T; \vartheta)] \cap \mathcal{M}_\#^{-1}(\ell)$ . For ease of presentation, we denote the corresponding task  
 135 distribution constrained in  $\mathcal{T}_\alpha$  by  $p_\cdot(\tau; \vartheta)$ .

136 Rather than optimizing VaR, a quantile, in meta learning, we take more interest in CVaR optimization,  
 137 a type of the expected tail risk. Such risk measure regards the conditional expectation and has  
 138 more desirable properties for meta learning: more adequate in handling adaptation risks in extreme  
 139 tails, more accessible sensitivity analysis w.r.t.  $\alpha$ , and more efficient optimization.

140 **Remark 1**  $\text{CVaR}_\alpha[\ell(T, \vartheta)]$  to minimize is respectively equivalent with the vanilla meta learning  
 141 optimization objective in Eq. (1) when  $\alpha = 0$  and the worst-case meta learning optimization objective  
 142 in Eq. (3) when  $\alpha$  is sufficiently close to 1.

143 **4.2 Meta Learning via Controlling the Expected Tail Risk**

144 As mentioned in Remark (1), the previous two meta learning objectives can be viewed as special  
 145 cases within the CVaR principle. Furthermore, we turn to a particular distributionally robust fast  
 146 adaptation with the adjustable confidence level  $\alpha$  to control the expected tail risk in optimization as  
 147 follows.

148 **Distributionally Robust Meta Learning Objective.** With the previously introduced normalized  
 149 probability density function  $p_\cdot(\tau; \vartheta)$ , minimizing  $\text{CVaR}_\alpha[\ell(T, \vartheta)]$  can be rewritten as Eq. (3).

$$\min_{\vartheta \in \mathcal{G}} E_\cdot(\vartheta) := E_{p_\alpha(\cdot; \vartheta)}[\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)] \quad (3)$$

151 Even though CVaR  $[\ell(T, \vartheta)]$  is a function of the model parameter  $\vartheta$ , the integral in Eq. (3) is  
 152 intractable due to the involvement of  $p(\tau; \vartheta)$  in a non-closed form.

153 **Assumption 2** For meta risk function values, the cumulative distribution function  $F(l; \vartheta)$  is  $\beta$ -  
 154 Lipschitz continuous w.r.t.  $l$ , and the implicit normalized probability density function of tasks  $p(\tau; \vartheta)$   
 155 is  $\beta$ -Lipschitz continuous w.r.t.  $\vartheta$ .

**Assumption 3** For any valid  $\vartheta \in \mathcal{V}$  and corresponding implicit normalized probability density  
 function of tasks  $p(\tau; \vartheta)$ , the meta risk function value can be bounded by a positive constant  $L_{\max}$ :

$$\sup_{\mathcal{D}^T, \mathcal{D}^C} \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \leq L_{\max}.$$

156 **Proposition 1** Under assumptions (1)/(2)/(3), the meta learning optimization objective  $E(\vartheta)$  in Eq.  
 157 (3) is continuous w.r.t.  $\vartheta$ .

Further, we use  $\xi(\vartheta)$  to denote the VaR  $[\ell(T, \vartheta)]$  for simple notations. The same as that in  
 (Rockafellar et al., 2000), we introduce a slack variable  $\xi \in \mathbb{R}$  and the auxiliary risk function  
 $[\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) - \xi]^+ := \max\{\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) - \xi, 0\}$ . To circumvent directly optimizing the non-  
 analytical  $p(\tau; \vartheta)$ , we can convert the probability constrained function  $E(\vartheta)$  to the below uncon-  
 strained one after optimizing  $\xi$ :

$$\varphi(\xi; \vartheta) = \xi + \frac{1}{1 - \alpha} \mathbb{E}_{p(\cdot)} \left[ [\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) - \xi]^+ \right].$$

158 It is demonstrated that  $E(\vartheta) = \min_{\xi \in \mathbb{R}} \varphi(\xi; \vartheta)$  and  $\xi \in \arg \min_{\xi \in \mathbb{R}} \varphi(\xi; \vartheta)$  in (Rockafellar  
 159 et al., 2000), and also note that CVaR is the upper bound of  $\xi$ , implying

$$\xi \leq \varphi(\xi; \vartheta) \leq \varphi(\xi; \vartheta), \quad \forall \xi \in \mathbb{R} \text{ and } \vartheta \in \mathcal{V}. \quad (4)$$

160 With the deductions from Eq.s (3)/(4), we can resort the distributionally robust meta learning  
 161 optimization objective with the probability constraint into a unconstrained optimization objective as  
 162 Eq.(5).

$$\min_{\xi \in \mathbb{R}} \xi + \frac{1}{1 - \alpha} \mathbb{E}_{p(\cdot)} \left[ [\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) - \xi]^+ \right] \quad (5)$$

164 **Sample Average Approximation.** For the stochastic programming problem above, it is mostly  
 165 intractable to derive the analytical form of the integral. Hence, we need to perform Monte Carlo  
 166 estimates of Eq. (5) to obtain Eq. (6) for optimization.

$$\min_{\xi \in \mathbb{R}} \xi + \frac{1}{(1 - \alpha)B} \sum_{i=1}^B [\ell(\mathcal{D}_i^T, \mathcal{D}_i^C; \vartheta) - \xi]^+ \quad (6)$$

168 **Remark 2** If  $\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta)$  is convex w.r.t.  $\vartheta$ , then Eq.s (5)/(6) are also convex functions. In this  
 169 case, the optimization objective Eq. (6) of our interest can be resolved with the help of several convex  
 170 programming algorithms (Fan et al., 2017; Meng et al., 2020; Levy et al., 2020).

### 171 4.3 Heuristic Algorithms for Optimization

172 Unfortunately, most existing meta learning models' risk functions (Finn et al., 2017; Garnelo et al.,  
 173 2018a; Santoro et al., 2016; Li et al., 2017; Duan et al., 2016),  $\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta)$  are non-convex w.r.t.  $\vartheta$ ,  
 174 bringing difficulties in optimization of Eq.s (5)/(6).

175 To this end, we propose a simple yet effective optimization strategy, where the sampled task batch is  
 176 used to approximate the VaR and the integral in Eq. (5) for deriving Eq. (6). In detail, two stages  
 177 are involved in iterations: (i) approximate VaR  $[\ell(T, \vartheta)] \approx \hat{\xi}$  with the meta batch values, which can  
 178 be achieved via either a quantile estimator (Dong and Nakayama, 2018) or other density estimators;  
 179 (ii) optimize  $\vartheta$  in Eq. (6) via stochastic updates after replacing  $\xi$  by the estimated  $\hat{\xi}$ .

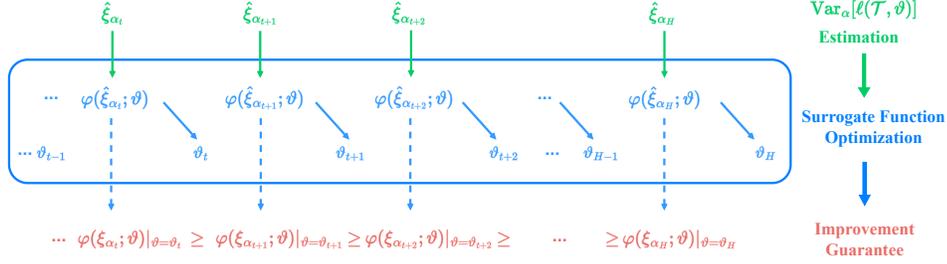


Figure 2: **Optimization Diagram of Distributionally Robust Meta Learning with Surrogate Functions.** From left to right: the meta model parameters  $\vartheta$  in the middle block are optimized w.r.t. the constructed surrogate function  $\varphi(\hat{\xi}_{\alpha_t}; \vartheta)$  marked in blue in  $t$ -th iteration. Under certain conditions in Theorem (1), the distributionally robust meta learning objective  $\varphi(\xi_{\alpha}; \vartheta)$  marked in pink can be decreased monotonically until it reaches the convergence in the  $H$ -th iteration.

**Proposition 2** Suppose there exists  $\delta \geq \mathbb{R}^+$  such that  $j\xi(\vartheta) - \hat{\xi}(\vartheta)j < \delta$  with  $\hat{\xi}(\vartheta)$  an estimate of  $\xi(\vartheta)$ . Then there exists a constant  $\kappa = \max\{f_{\frac{2}{1-\gamma}}, \frac{1}{1-\gamma}g\}$  such that

$$\varphi(\hat{\xi}(\vartheta); \vartheta) - \kappa \delta < E(\vartheta) < \varphi(\hat{\xi}(\vartheta); \vartheta).$$

180 The performance gap resulting from VaR approximation error is estimated in Proposition (2). For  
 181 ease of implementation, we adopt crude Monte Carlo methods (Kroese and Rubinstein, 2012) to  
 182 obtain a consistent estimator of  $\xi$ .

183 **Theorem 1 (Improvement Guarantee)** Under assumptions (1)/(2)/(3), suppose that the estimate  
 184 error with the crude Monte Carlo holds:  $j\hat{\xi}_t - \xi_j \leq \frac{\delta}{\sqrt{t}}$ ,  $\forall t \geq 2N^+$ , with the subscript  $t$  the  
 185 iteration number,  $\lambda$  the learning rate in stochastic gradient descent,  $\beta$  the Lipschitz constant of  
 186 the risk cumulative distribution function, and  $\alpha$  the confidence level. Then the proposed heuristic  
 187 algorithm with the crude Monte Carlo can produce at least a local optimum for distributionally  
 188 robust fast adaptation.

189 Note that Theorem (1) indicates that under certain conditions, using the above heuristic algorithm has  
 190 the performance improvement guarantee, which corresponds to Fig. (2). The error resulting from the  
 191 approximate algorithm is further estimated in Appendix Theorem (2).

#### 192 4.4 Instantiations & Implementations

193 Our proposed optimization strategy applies to all meta learning methods and has the improvement guarantee  
 194 in Theorem (1). Here we take two representative meta learning methods, MAML (Finn et al., 2017) and CNP  
 195 (Garnelo et al., 2018a), as examples and show how to robustify them through the lens of risk distributions. Note that the forms  
 196 of  $\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta)$  sometime differ in these methods. Also, the detailed implementation of the strategy relies on specific  
 197 meta learning algorithms and models.  
 198  
 199  
 200  
 201

202 **Example 1 (DR-MAML)** With the task distribution  $p(\tau)$   
 203 and model agnostic meta learning (Finn et al., 2017), the  
 204 distributionally robust MAML treats the meta learning  
 205 problem as a bi-level optimization with a VaR relevant  
 206 constraint.

$$\min_{\vartheta \in \mathbb{R}^2} \xi + \frac{1}{1-\alpha} E_{p(\cdot)} \left[ \left[ \ell(\mathcal{D}^T; \vartheta) - \lambda r \# \ell(\mathcal{D}^C; \vartheta) \right] \xi^+ \right] \quad (7)$$

207 The gradient operation  $r \# \ell(\mathcal{D}^C; \vartheta)$  corresponds to the inner loop with the learning rate  $\lambda$ .

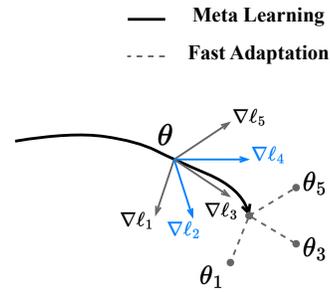


Figure 3: **Diagram of Distributionally Robust Fast Adaptation for Model Agnostic Meta Learning (Finn et al., 2017).** For example, with the size of the meta batch 5 and  $\alpha = 40\%$ , 5 (1 -  $\alpha$ ) tasks in gray with the worst fast adaptation performance are screened for updating meta initialization.

208 The resulting distributionally robust MAML (DR-MAML) is still a optimization-based method,  
 209 where a fixed percentage of tasks are screened for the outer loop. As shown in Eq. (7), the objective  
 210 is to obtain a robust meta initialization of the model parameter  $\vartheta$ .

211 **Example 2 (DR-CNP)** With the task distribution  $p(\tau)$  and the conditional neural process (Gar-  
 212 nelo et al., 2018a), the distributionally robust conditional neural process learns the functional  
 213 representations with a CVaR constraint.

$$\min_{\xi} \xi + \frac{1}{\alpha} \mathbb{E}_{p(\cdot)} \left[ \left[ \ell(\mathcal{D}^T; z, \theta_2) - \xi \right]^+ \right] \quad (8)$$

s.t.  $z = h_{\vartheta_1}(\mathcal{D}^C)$  with  $\vartheta = \{\theta_1, \theta_2\}$

214 Here  $h_{\vartheta_1}$  is a set encoder network with  $\theta_2$  the parameter of the decoder network.

215 The resulting distributionally robust CNP (DR-CNP) is to find a robust functional embedding to  
 216 induce underlying stochastic processes. Still, in Eq. (8), a proportion of tasks with the worst functional  
 217 representation performance are used in meta training.

218 Moreover, we convey the pipelines of optimizing these developed distributionally robust models in  
 219 Appendix Algorithms (1)/(2).

## 220 5 Experimental Results and Analysis

221 This section presents experimental results and examines fast adaptation performance in a distributional  
 222 sense. Without loss of generality, we take DR-MAML in Example (1) to run experiments.

223 **Benchmarks.** The same as work in (Collins et al., 2020), we use two commonly-seen downstream  
 224 tasks for meta learning experiments: few-shot regression and image classification. Besides, ablation  
 225 studies are included to assess other factors’ influence or the proposed strategy’s scalability.

226 **Baselines & Evaluations.** Since the primary investigation is regarding risk minimization principles,  
 227 we consider the previously mentioned *expected risk minimization*, *worst-case minimization*, and  
 228 *expected tail risk minimization* for meta learning. Hence, MAML (empirical risk), TR-MAML  
 229 (worst-case risk), and DR-MAML (expected tail risk) serve as examined methods. We evaluate these  
 230 methods’ performance based on the *Average*, *Worst-case*, and *CVaR* metrics. For the confidence  
 231 level to meta train DR-MAML, we empirically set  $\alpha = 0.7$  for few-shot regression tasks and  $\alpha = 0.5$   
 232 image classification tasks without external configurations.

### 233 5.1 Sinusoid Regression

234 Following (Finn et al., 2017; Collins et al., 2020),  
 235 we conduct experiments in  
 236 sinusoid regression tasks.  
 237 The mission is to approxi-  
 238 mate the function  $f(x) =$   
 239  $a \sin(x - b)$  with K-shot  
 240 randomly sampled function  
 241 points, where the task is de-  
 242 fined by  $a$  and  $b$ . In the sine function family, the target range, amplitude range, and phase range are  
 243 respectively  $[ -5.0, 5.0] \in \mathbb{R}$ ,  $a \in [0.1, 5.0]$  and  $b \in [0, 2\pi]$ . In the setup of meta training and testing  
 244 datasets, a distributional shift exists: numerous easy tasks and several difficult tasks are generated to  
 245 formulate the training dataset with all tasks in the space as the testing one. Please refer to Appendix  
 246 (J) for a detailed partition of meta-training, testing tasks, and neural architectures.

Table 1: Test average mean square errors (MSEs) with reported standard deviations for sinusoid regression (5 runs). We respectively consider 5-shot and 10-shot cases with  $\alpha = 0.7$ . The results are evaluated across the 490 meta-test tasks, as in (Collins et al., 2020). The best results are in bold.

Method	5-shot						10-shot					
	Average	Worst	CVaR	Average	Worst	CVaR	Average	Worst	CVaR	Average	Worst	CVaR
MAML (Finn et al., 2017)	1.02	0.10	3.89	0.83	2.25	0.15	0.66	0.16	2.57	0.70	1.15	0.19
TR-MAML (Collins et al., 2020)	1.09	0.08	<b>2.28</b>	0.35	1.79	0.06	0.77	0.11	<b>1.68</b>	0.43	1.27	0.28
DR-MAML (Ours)	<b>0.89</b>	<b>0.04</b>	2.91	0.46	<b>1.76</b>	<b>0.02</b>	<b>0.54</b>	<b>0.01</b>	1.70	0.17	<b>0.96</b>	<b>0.01</b>

248 **Result Analysis.** We list meta-testing MSEs in sinusoid regression in Table (1). As expected, the  
 249 tail risk minimization principle in DR-MAML can lead to an intermediate performance in the worst-  
 250 case. In both cases, the comparison between MAML and DR-MAML in MSEs indicates that such  
 251 probabilistic-constrained optimization in the task space even has the potential to advance average fast  
 252 adaptation performance. In contrast, TR-MAML has to sacrifice more average performance for the  
 253 worst-case improvement.

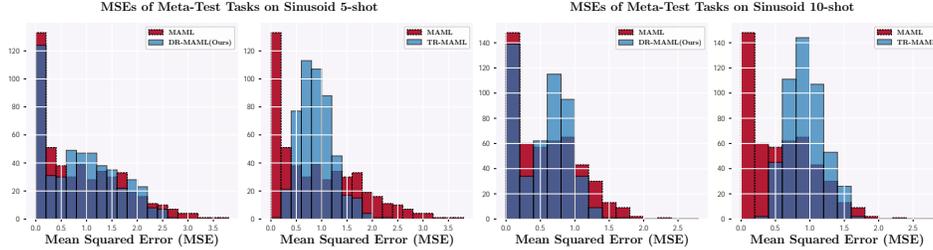


Figure 4: **Histograms of Meta-Testing Performance in Sinusoid Regression Problems.** With  $\alpha = 0.7$ , we respectively visualize the comparison results, DR-MAML-vs-MAML and TR-MAML-vs-MAML in 5-shot (Two Sub-figures Left Side) and 10-shot (Two Sub-figures Right Side) cases, for a sample trial.

254 More intuitively, Fig. (4) illustrates MSE statistics on the testing task distribution and further verifies  
 255 the effect of the CVaR principle in decreasing the proportional worst-case errors. In 5-shot  
 256 cases, the difference in MSE statistics is more significant: *the worst-case method tends to increase*  
 257 *the skewness in risk distributions* with many task risk values gathered in regions of intermediate  
 258 performance, which is unfavored in general cases. As for why DR-MAML surpasses MAML in terms of  
 259 average performance, we attribute it to the benefits of external robustness in several scenarios, *e.g.*,  
 260 the drift of training/testing task distributions.

Table 2: **Average N-way K-shot classification accuracies in Omniglot with reported standard deviations (3 runs).** With  $\alpha = 0.5$ , the best results are in bold.

Method	Meta-Training Alphabets						Meta-Testing Alphabets																	
	5-way 1-shot			20-way 1-shot			5-way 1-shot			20-way 1-shot														
	Average	Worst	CVaR	Average	Worst	CVaR	Average	Worst	CVaR	Average	Worst	CVaR												
MAML (Finn et al., 2017)	<b>98.4</b>	0.2	82.4	1.1	<b>96.9</b>	0.5	99.2	0.1	33.9	3.0	80.9	0.7	93.5	0.2	82.5	0.2	91.6	0.6	67.6	2.0	49.7	3.5	60.4	1.7
TR-MAML (Collins et al., 2020)	97.4	0.6	<b>95.0</b>	0.3	96.5	0.4	92.2	0.8	<b>82.4</b>	2.1	<b>87.2</b>	0.9	93.1	1.1	<b>85.3</b>	1.9	91.3	0.9	74.3	1.4	<b>58.4</b>	1.8	<b>68.5</b>	1.2
DR-MAML (Ours)	97.1	0.3	84.0	0.4	95.1	0.3	<b>99.6</b>	0.6	57.9	2.4	84.8	0.7	<b>93.7</b>	0.4	84.1	0.8	<b>92.1</b>	0.5	<b>74.6</b>	1.2	51.0	2.3	66.4	1.4

## 261 5.2 Few-Shot Image Classification

262 Here we do investigations in few-shot image classification. Each task is an N-way K-shot clas-  
 263 sification with N the number of classes and K the number of labeled examples in one class. The  
 264 Omniglot (Lake et al., 2015) and mini-ImageNet (Vinyals et al., 2016) datasets work as benchmarks for examination. We retain the setup of datasets in work  
 265 (Collins et al., 2020).

Table 3: **Average 5-way 1-shot classification accuracies in mini-ImageNet with reported standard deviations (3 runs).** With  $\alpha = 0.5$ , the best results are in bold.

Method	Eight Meta-Training Tasks			Four Meta-Testing Tasks								
	Average	Worst	CVaR	Average	Worst	CVaR						
MAML (Finn et al., 2017)	70.1	2.2	48.0	4.5	63.2	2.6	46.6	0.4	44.7	0.7	44.6	0.7
TR-MAML (Collins et al., 2020)	63.2	1.3	60.7	1.6	62.1	1.2	48.5	0.6	45.9	0.8	46.6	0.5
DR-MAML (Ours)	<b>70.2</b>	<b>0.2</b>	<b>63.4</b>	<b>0.2</b>	<b>67.2</b>	<b>0.1</b>	<b>49.4</b>	<b>0.1</b>	<b>47.1</b>	<b>0.1</b>	<b>47.5</b>	<b>0.1</b>

272 **Result Analysis.** The classification accuracies in Omniglot are illustrated in Table (2): In 5-way  
 273 1-shot cases, DR-MAML obtains the best average and CVaR in meta-testing datasets, while TR-  
 274 MAML achieves the best worst-case performance with a slight degradation of average performance  
 275 compared to MAML in both training/testing datasets. In 20-way 1-shot cases, for training/testing  
 276 datasets, we surprisingly notice that the expected tail risk is not well optimized with DR-MAML, but  
 277 there is an average performance gain; while TR-MAML works best in worst-case/CVaR metrics.

278 When it comes to mini-ImageNet, findings are distinguished a lot from the above one: in Table (3),  
 279 DR-MAML yields the best result in all evaluation metrics and cases, even regarding the worst-case.  
 280 TR-MAML can also improve all metrics in meta-testing cases. Overall, challenging meta-learning  
 281 tasks can reveal more advantages of DR-MAML over others.

## 282 5.3 Ablation Studies

283 This part mainly checks the influence of the confidence level  $\alpha$  and the meta batch size towards the  
 284 distribution of fast adapted risk values. Apart from examining these factors of interest, additional  
 285 experiments are also conducted in this paper; please refer to Appendix (K) for more details.

286 **Sensitivity to Confidence Levels  $\alpha$ .** To deepen understanding of the confidence level  $\alpha$ 's effect in  
 287 fast adaptation performance, we vary  $\alpha$  to train DR-MAML and evaluate models under previously

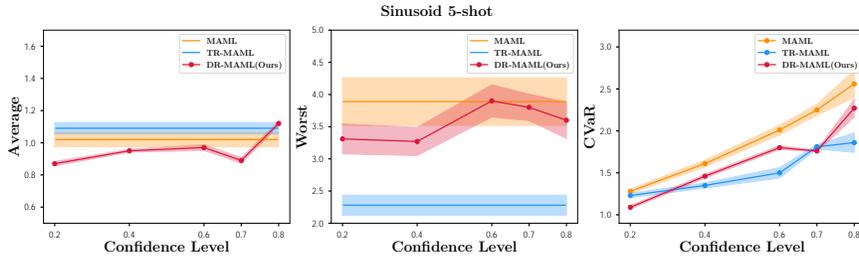


Figure 5: **Meta Testing MSEs of Meta-Trained DR-MAML with Various Confidence Levels  $\alpha$ .** MAML and TR-MAML are irrelevant with the variation of  $\alpha$  in meta-training. The plots report testing MSEs with standard error bars in shadow regions.

288 mentioned metrics. Taking the sinusoid 5-shot regression as an example, we calculate MSEs and  
 289 visualize the fluctuations with additional  $\alpha$ -values in Fig. (5). The results in the worst-case exhibit  
 290 higher deviations. The trend in Average/CVaR metrics shows that with increasing  $\alpha$ , DR-MAML  
 291 gradually approaches TR-MAML in average and CVaR MSEs. With  $\alpha = 0.8$ , ours is less sensitive  
 292 to the confidence level and mostly beats MAML/TR-MAML in average performance. Though ours  
 293 aims at optimizing CVaR, it cannot always ensure such a metric to surpass TR-MAML in all  
 294 confidence levels due to rigorous assumptions in Theorem (1).

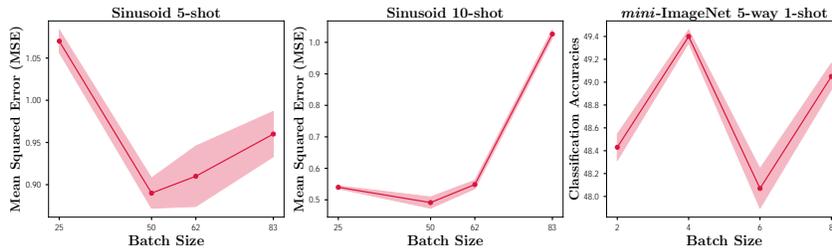


Figure 6: **Meta Testing Average Performance of Meta-Trained DR-MAML with Various Sizes of the Meta Batch.** The plots report average results with standard error bars in shadow regions.

295 **Influence of the Task Batch Size.** Note that our optimization strategy relies on the estimate of  
 296 VaR, and the improvement guarantee relates to the estimation error. Theoretically, the meta batch in  
 297 training can directly influence the optimization result. Here we vary the meta batch size in training,  
 298 and evaluated results are visualized in Fig. (6). In regression scenarios, the average MSEs can  
 299 decrease to a certain level with an increase of meta batch, and then performance degrades. We  
 300 attribute the performance gain with a batch increase to more accurate VaR estimates; however, a  
 301 meta batch larger than some threshold can worsen the first-order meta learning algorithms' efficiency,  
 302 similarly observed in (Nichol et al., 2018). As for classification scenarios, there appears no clear  
 303 trend since the meta batch is smaller enough.

## 304 6 Conclusion and Limitations

305 **Technical Discussions.** This work contributes more insights into robustifying fast adaptation in meta  
 306 learning. Our utilized expected tail risk trades off the expected risk minimization and worst-case risk  
 307 minimization, and the two-stage strategy works as the heuristic to approximately solve the problem  
 308 with an improvement guarantee. Our strategy can empirically alleviate the worst-case fast adaptation  
 309 and sometimes even improve average performance.

310 **Existing Limitations.** Though our robustification strategy is simple yet effective in implementations,  
 311 empirical selection of the optimal meta batch size is challenging, especially for first-order optimization  
 312 methods. Meanwhile, the theoretical analysis only applies to a fraction of meta learning tasks when  
 313 risk function values are in a compact continuous domain.

314 **Future Extensions.** Designing a heuristic algorithm with an improvement guarantee is non-trivial  
 315 and relies on the properties of risk functions. This research direction has practical meaning in the era  
 316 of large models and deserves more investigations in terms of optimization methods. Also, establishing  
 317 connections between the optimal meta batch size and specific stochastic optimization algorithms can  
 318 be a promising theoretical research issue in this domain.

319 **References**

- 320 Abbas, M., Xiao, Q., Chen, L., Chen, P.-Y., and Chen, T. (2022). Sharp-maml: Sharpness-aware  
321 model-agnostic meta learning. In *International Conference on Machine Learning*, pages 10–32.  
322 PMLR.
- 323 Allen, K., Shelhamer, E., Shin, H., and Tenenbaum, J. (2019). Infinite mixture prototypes for few-shot  
324 learning. In *International Conference on Machine Learning*, pages 232–241. PMLR.
- 325 Antoniou, A., Edwards, H., and Storkey, A. (2019). How to train your maml. In *Seventh International  
326 Conference on Learning Representations*.
- 327 Bahadur, R. R. (1966). A note on quantiles in large samples. *The Annals of Mathematical Statistics*,  
328 37(3):577–580.
- 329 Barabás, B. (1987). Estimation of density functions by order statistics. *Periodica Mathematica  
330 Hungarica*, 18(2):115–122.
- 331 Bartunov, S. and Vetrov, D. (2018). Few-shot generative modelling with generative matching  
332 networks. In *International Conference on Artificial Intelligence and Statistics*, pages 670–678.  
333 PMLR.
- 334 Beck, J., Jackson, M. T., Vuorio, R., and Whiteson, S. (2023). Hypernetworks in meta-reinforcement  
335 learning. In *Conference on Robot Learning*, pages 1478–1487. PMLR.
- 336 Chen, R. S., Lucier, B., Singer, Y., and Syrgkanis, V. (2017). Robust optimization for non-convex  
337 objectives. *Advances in Neural Information Processing Systems*, 30.
- 338 Collins, L., Mokhtari, A., and Shakkottai, S. (2020). Task-robust model-agnostic meta-learning.  
339 *Advances in Neural Information Processing Systems*, 33:18860–18871.
- 340 Dong, H. and Nakayama, M. K. (2018). A tutorial on quantile estimation via monte carlo. In  
341 *International Conference on Monte Carlo and Quasi-Monte Carlo Methods in Scientific Computing*,  
342 pages 3–30. Springer.
- 343 Duan, Y. (2017). *Meta learning for control*. University of California, Berkeley.
- 344 Duan, Y., Schulman, J., Chen, X., Bartlett, P. L., Sutskever, I., and Abbeel, P. (2016). R12: Fast  
345 reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*.
- 346 Duchi, J. C. and Namkoong, H. (2021). Learning models with uniform performance via distribution-  
347 ally robust optimization. *The Annals of Statistics*, 49(3):1378–1406.
- 348 Fallah, A., Mokhtari, A., and Ozdaglar, A. (2021). Generalization of model-agnostic meta-learning  
349 algorithms: Recurring and unseen tasks. *Advances in Neural Information Processing Systems*,  
350 34:5469–5480.
- 351 Fan, Y., Lyu, S., Ying, Y., and Hu, B. (2017). Learning with average top-k loss. *Advances in neural  
352 information processing systems*, 30.
- 353 Finn, C., Abbeel, P., and Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep  
354 networks. In *International conference on machine learning*, pages 1126–1135. PMLR.
- 355 Finn, C., Xu, K., and Levine, S. (2018). Probabilistic model-agnostic meta-learning. *Advances in  
356 neural information processing systems*, 31.
- 357 Foong, A., Bruinsma, W., Gordon, J., Dubois, Y., Requeima, J., and Turner, R. (2020). Meta-learning  
358 stationary stochastic process prediction with convolutional neural processes. *Advances in Neural  
359 Information Processing Systems*, 33:8284–8295.
- 360 Gagne, C. and Dayan, P. (2021). Two steps to risk sensitivity. *Advances in Neural Information  
361 Processing Systems*, 34:22209–22220.
- 362 Garnelo, M., Rosenbaum, D., Maddison, C., Ramalho, T., Saxton, D., Shanahan, M., Teh, Y. W.,  
363 Rezende, D., and Eslami, S. A. (2018a). Conditional neural processes. In *International Conference  
364 on Machine Learning*, pages 1704–1713. PMLR.

- 365 Garnelo, M., Schwarz, J., Rosenbaum, D., Viola, F., Rezende, D. J., Eslami, S., and Teh, Y. W.  
366 (2018b). Neural processes. *arXiv preprint arXiv:1807.01622*.
- 367 Goel, A., Agarwal, A., Vatsa, M., Singh, R., and Ratha, N. K. (2020). Dndnet: Reconfiguring cnn  
368 for adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
369 *Pattern Recognition Workshops*, pages 22–23.
- 370 Goldblum, M., Fowl, L., and Goldstein, T. (2019). Robust few-shot learning with adversarially  
371 queried meta-learners.
- 372 Gondal, M. W., Joshi, S., Rahaman, N., Bauer, S., Wuthrich, M., and Scholkopf, B. (2021). Function  
373 contrastive learning of transferable meta-representations. In *Proceedings of the 38th International*  
374 *Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of  
375 *Proceedings of Machine Learning Research*, pages 3755–3765. PMLR.
- 376 Goodfellow, I., McDaniel, P., and Papernot, N. (2018). Making machine learning robust against  
377 adversarial inputs. *Communications of the ACM*, 61(7):56–66.
- 378 Goodfellow, I. J., Shlens, J., and Szegedy, C. (2014). Explaining and harnessing adversarial examples.  
379 *arXiv preprint arXiv:1412.6572*.
- 380 Gordon, J., Bronskill, J., Bauer, M., Nowozin, S., and Turner, R. (2019). Meta-learning probabilistic  
381 inference for prediction. In *International Conference on Learning Representations*.
- 382 Grant, E., Finn, C., Levine, S., Darrell, T., and Griffiths, T. (2018). Recasting gradient-based  
383 meta-learning as hierarchical bayes. *arXiv preprint arXiv:1801.08930*.
- 384 Hiraoka, T., Imagawa, T., Mori, T., Onishi, T., and Tsuruoka, Y. (2019). Learning robust options by  
385 conditional value at risk optimization. *Advances in Neural Information Processing Systems*, 32.
- 386 Hsieh, Y.-P., Mertikopoulos, P., and Cevher, V. (2021). The limits of min-max optimization algorithms:  
387 Convergence to spurious non-critical sets. In *International Conference on Machine Learning*,  
388 pages 4337–4348. PMLR.
- 389 Jaafra, Y., Laurent, J. L., Deruyver, A., and Naceur, M. S. (2019). Robust reinforcement learning for  
390 autonomous driving.
- 391 Jiang, Y., Li, C., Dai, W., Zou, J., and Xiong, H. (2021). Monotonic robust policy optimization with  
392 model discrepancy. In *International Conference on Machine Learning*, pages 4951–4960. PMLR.
- 393 Johannsmeier, L., Gerchow, M., and Haddadin, S. (2019). A framework for robot manipulation: Skill  
394 formalism, meta learning and adaptive control. In *2019 International Conference on Robotics and*  
395 *Automation (ICRA)*, pages 5844–5850. IEEE.
- 396 Juditsky, A., Nemirovski, A., and Tauvel, C. (2011). Solving variational inequalities with stochastic  
397 mirror-prox algorithm. *Stochastic Systems*, 1(1):17–58.
- 398 Kroese, D. P. and Rubinstein, R. Y. (2012). Monte carlo methods. *Wiley Interdisciplinary Reviews:*  
399 *Computational Statistics*, 4(1):48–58.
- 400 Kurakin, A., Goodfellow, I., and Bengio, S. (2016). Adversarial machine learning at scale. *arXiv*  
401 *preprint arXiv:1611.01236*.
- 402 Lake, B. M., Salakhutdinov, R., and Tenenbaum, J. B. (2015). Human-level concept learning through  
403 probabilistic program induction. *Science*, 350(6266):1332–1338.
- 404 Larochelle, S. (2017). Optimization as a model for few-shot learning. In *International Conference on*  
405 *Learning Representations*.
- 406 LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *nature*, 521(7553):436–444.
- 407 Lesort, T., Caccia, M., and Rish, I. (2021). Understanding continual learning settings with data  
408 distribution drift analysis. *arXiv preprint arXiv:2104.01678*.

- 409 Levy, D., Carmon, Y., Duchi, J. C., and Sidford, A. (2020). Large-scale methods for distributionally  
410 robust optimization. *Advances in Neural Information Processing Systems*, 33:8847–8860.
- 411 Li, Z., Zhou, F., Chen, F., and Li, H. (2017). Meta-sgd: Learning to learn quickly for few-shot  
412 learning. *arXiv preprint arXiv:1707.09835*.
- 413 Liu, Q., Liu, T., Liu, Z., Wang, Y., Jin, Y., and Wen, W. (2018). Security analysis and enhancement of  
414 model compressed deep learning systems under adversarial attacks. In *2018 23rd Asia and South  
415 Pacific Design Automation Conference (ASP-DAC)*, pages 721–726. IEEE.
- 416 Meng, S. Y., Charisopoulos, V., and Gower, R. M. (2020). A stochastic prox-linear method for cvar  
417 minimization.
- 418 Nemirovski, A., Juditsky, A., Lan, G., and Shapiro, A. (2009). Stochastic approximation approach to  
419 stochastic programming. *SIAM J. Optim. CiteSeer*.
- 420 Nichol, A., Achiam, J., and Schulman, J. (2018). On first-order meta-learning algorithms. *arXiv  
421 preprint arXiv:1803.02999*.
- 422 Olds, K. C. (2015). Global indices for kinematic and force transmission performance in parallel  
423 robots. *IEEE Transactions on Robotics*, 31(2):494–500.
- 424 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein,  
425 N., Antiga, L., Desmaison, A., Köpf, A., Yang, E. Z., DeVito, Z., Raison, M., Tejani, A., Chil-  
426 amkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. (2019). Pytorch: An imperative style,  
427 high-performance deep learning library. In *Advances in Neural Information Processing Systems 32:  
428 Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December  
429 8-14, 2019, Vancouver, BC, Canada*, pages 8024–8035.
- 430 Pinto, L., Davidson, J., Sukthankar, R., and Gupta, A. (2017). Robust adversarial reinforcement  
431 learning. In *International Conference on Machine Learning*, pages 2817–2826. PMLR.
- 432 Quaranta, A. G. and Zaffaroni, A. (2008). Robust optimization of conditional value at risk and  
433 portfolio selection. *Journal of Banking & Finance*, 32(10):2046–2056.
- 434 Rahimian, H. and Mehrotra, S. (2019). Distributionally robust optimization: A review.
- 435 Rajeswaran, A., Finn, C., Kakade, S. M., and Levine, S. (2019). Meta-learning with implicit gradients.  
436 *Advances in neural information processing systems*, 32.
- 437 Ren, J., Zhang, D., Wang, Y., Chen, L., Zhou, Z., Chen, Y., Cheng, X., Wang, X., Zhou, M., Shi, J.,  
438 et al. (2021). Towards a unified game-theoretic view of adversarial perturbations and robustness.  
439 *Advances in Neural Information Processing Systems*, 34:3797–3810.
- 440 Rockafellar, R. T., Uryasev, S., et al. (2000). Optimization of conditional value-at-risk. *Journal of  
441 risk*, 2:21–42.
- 442 Rosenblatt, M. (1956). A central limit theorem and a strong mixing condition. *Proceedings of the  
443 national Academy of Sciences*, 42(1):43–47.
- 444 Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. (2020). Distributionally robust neural  
445 networks. In *International Conference on Learning Representations*.
- 446 Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., and Lillicrap, T. (2016). Meta-learning with  
447 memory-augmented neural networks. In *International conference on machine learning*, pages  
448 1842–1850. PMLR.
- 449 Shalev-Shwartz, S. and Wexler, Y. (2016). Minimizing the maximal loss: How and why. In  
450 *International Conference on Machine Learning*, pages 793–801. PMLR.
- 451 Silva, S. H. and Najafirad, P. (2020). Opportunities and challenges in deep learning adversarial  
452 robustness: A survey. *arXiv preprint arXiv:2007.00753*.
- 453 Snell, J., Swersky, K., and Zemel, R. (2017). Prototypical networks for few-shot learning. *Advances  
454 in neural information processing systems*, 30.

- 455 Tay, S. S., Foo, C. S., Daisuke, U., Leong, R., and Low, B. K. H. (2022). Efficient distributionally  
456 robust bayesian optimization with worst-case sensitivity. In *International Conference on Machine*  
457 *Learning*, pages 21180–21204. PMLR.
- 458 Triantafillou, E., Zhu, T., Dumoulin, V., Lamblin, P., Evcı, U., Xu, K., Goroshin, R., Gelada, C.,  
459 Swersky, K., and Manzagol, P. A. (2019). Meta-dataset: A dataset of datasets for learning to learn  
460 from few examples.
- 461 Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al. (2016). Matching networks for one shot  
462 learning. *Advances in neural information processing systems*, 29.
- 463 Vuorio, R., Sun, S.-H., Hu, H., and Lim, J. J. (2019). Multimodal model-agnostic meta-learning via  
464 task-aware modulation. *Advances in neural information processing systems*, 32.
- 465 Wang, Q., Federici, M., and van Hoof, H. (2023). Bridge the inference gaps of neural processes via  
466 expectation maximization. In *The Eleventh International Conference on Learning Representations*.
- 467 Wang, Q. and van Hoof, H. (2022). Learning expressive meta-representations with mixture of expert  
468 neural processes. In *Advances in neural information processing systems*.
- 469 Wang, R., Xu, K., Liu, S., Chen, P.-Y., Weng, T.-W., Gan, C., and Wang, M. (2020). On fast  
470 adversarial robustness adaptation in model-agnostic meta-learning. In *International Conference on*  
471 *Learning Representations*.
- 472 Wang, Z., Shen, Y., and Zavlanos, M. (2022). Risk-averse no-regret learning in online convex games.  
473 In *International Conference on Machine Learning*, pages 22999–23017. PMLR.
- 474 Wilder, B. (2018). Risk-sensitive submodular optimization. In *Proceedings of the AAAI Conference*  
475 *on Artificial Intelligence*, volume 32.
- 476 Zhang, J., Cheung, B., Finn, C., Levine, S., and Jayaraman, D. (2020). Cautious adaptation for  
477 reinforcement learning in safety-critical settings. In *International Conference on Machine Learning*,  
478 pages 11055–11065. PMLR.
- 479 Zhao, D., von Oswald, J., Kobayashi, S., Sacramento, J., and Grewe, B. F. (2020). Meta-learning via  
480 hypernetworks.

481	<b>Contents</b>	
482	<b>1 Introduction</b>	<b>1</b>
483	<b>2 Literature Review</b>	<b>2</b>
484	<b>3 Preliminaries</b>	<b>3</b>
485	<b>4 Distributionally Robust Fast Adaptation</b>	<b>3</b>
486	4.1 Meta Risk Functions as Random Variables . . . . .	4
487	4.2 Meta Learning via Controlling the Expected Tail Risk . . . . .	4
488	4.3 Heuristic Algorithms for Optimization . . . . .	5
489	4.4 Instantiations & Implementations . . . . .	6
490	<b>5 Experimental Results and Analysis</b>	<b>7</b>
491	5.1 Sinusoid Regression . . . . .	7
492	5.2 Few-Shot Image Classification . . . . .	8
493	5.3 Ablation Studies . . . . .	8
494	<b>6 Conclusion and Limitations</b>	<b>9</b>
495	<b>A Frequently Asked Question</b>	<b>16</b>
496	<b>B Pseudo Algorithms of DR-MAML &amp; DR-CNPs</b>	<b>17</b>
497	<b>C Properties of Risk Minimization Principles</b>	<b>18</b>
498	C.1 Stochastic Optimization in the Constrained Form . . . . .	18
499	C.2 SGD Intractability of Meta Learning CVaR Optimization Objective . . . . .	18
500	C.3 Risk-Sensitive Applications & Optimization Strategies . . . . .	18
501	<b>D Computational Complexity</b>	<b>19</b>
502	<b>E Proof of Remark (2)</b>	<b>19</b>
503	<b>F Proof of Proposition (1)</b>	<b>20</b>
504	<b>G Proof of Proposition (2)</b>	<b>21</b>
505	<b>H Proof of Improvement Guarantee in Theorem (1)</b>	<b>22</b>
506	<b>I Proof of Approximation Error in Theorem (2)</b>	<b>25</b>
507	<b>J Experimental Set-up &amp; Implementation Details</b>	<b>25</b>
508	J.1 Meta Learning Datasets & Tasks . . . . .	26
509	J.2 Neural Architectures & Optimization . . . . .	27

510	<b>K Additional Experimental Results</b>	<b>27</b>
511	<b>L Platforms &amp; Computational Tools</b>	<b>29</b>

## 512 A Frequently Asked Question

513 Here we collected technical questions and suggestions from researchers who helped check out the  
514 manuscript. We thank these researchers for precious questions and provide more details.

515 **Novelty & primary findings of this work.** Here we mainly summarize two points of novelty in this  
516 work:

- 517 • *Meta Learning Robustification Framework.* Though the concept of the expected tail risk  
518 has emerged for several decades and has been widely employed in financial domains, the  
519 application to fast adaptation or robustification of meta learning remains limited in literature  
520 as far as we know.
- 521 • *Optimization Strategy & Theoretical Analysis.* We theoretically analyze two-stage opti-  
522 mization strategies as the heuristic algorithm in optimizing the distributionally robust meta  
523 learning models and demonstrate the improvement guarantee under certain conditions.

524 As verified in experimental results in the main paper, placing a probabilistic constraint in the task space  
525 is meaningful. It circumvents the effect of over-pessimistic consideration (worst-case optimization),  
526 increases robustness in proportional cases, and mostly retains or even improves average performance.

527 Apart from the novelty in the framework and algorithm parts, we have several findings which bring  
528 crucial insights into meta learning: (i) Not all tasks are necessary to perform fast adaptation. (ii)  
529 Additional focus on the tail risk has the potential to enhance models' generalization capability. (iii)  
530 The tail risk instead of extreme worst-case risk can better advance robustness in challenging datasets.

531 **Meta risk function values as random variables.** In some few-shot learning related work, the context  
532 and the target dataset are equivalently called the support and query datasets. The definition of a task  
533 in meta learning is up to application scenarios and specific meta learning algorithms or models. The  
534 commonly-used sinusoid regression using model-agnostic meta learning (Finn et al., 2017) considers  
535 the fixed number of context points to induce tasks. In contrast, conditional neural processes (Garnelo  
536 et al., 2018a) for few-shot regression vary the number of context points to induce tasks. Once the  
537 context and the target are partitioned and the model parameter is specified, we can obtain a meta risk  
538 function value. However, the meta risk function values are not in a compact Euclidean subspace in  
539 several cases.

540 **The continuity of the meta risk probability density function.** This relates to the task distribution  
541 and meta learning problems. Throughout the few-shot regression task, the meta risk function value  
542 can be approximately viewed as a continuous random variable. However, when it comes to the  
543 few-shot image classification mission, the meta risk function value is impossible to cover the entire  
544 continuous interval. In these scenarios, the probable values of accuracies are finite. This makes  
545 the theoretical analysis, e.g. Theorem (1), no longer holds. For example, Assumption (2) will be  
546 unrealistic when there exists a constant gap between two accuracy values. Hence, we leave this part a  
547 future research direction in theoretical analysis. Regarding the meta risk function, ours shares the  
548 same setup as that in (Fallah et al., 2021).

549 **The selection of baselines in robust meta learning.** A large body of prior work considers the  
550 robustness of meta learning in scenarios when the input of a data point, the model parameter, and the  
551 number of context points are trembled or modified. Robustness in presence of tasks distributions is  
552 seldom investigated except for the worst-case optimization in (Collins et al., 2020). Hence, we retain  
553 most of the setups in (Collins et al., 2020) for a fair comparison.

554 **Influence of risk minimization principles.** This paper is primarily devoted to studying the influence  
555 of the risk minimization principle on meta learning. The empirical risk minimization principle  
556 corresponds to reducing the Monte Carlo estimate of Average-case Meta Learning in this work.  
557 In comparison to TR-MAML (Collins et al., 2020), our approach has a couple of advantages as  
558 follows: (i) easier implementations. Note that min-max optimization is numerically unstable and  
559 requires a relaxation method for computationally intensive convex optimization, e.g., robust stochastic  
560 mirror-prox algorithm used in TR-MAML (Collins et al., 2020). (ii) more flexible in terms of  
561 robustness concept. Theoretically, the worst-case meta learning corresponds to the extreme case  
562 of the distributionally robust risk minimization principle. (iii) empirically better performance in  
563 most cases. The expected tail risk minimization preserves a particular property that minimizing

564 proportional worst-case fast adaptation seldom sacrifices the average performance. These advantages  
 565 are also why we call our approach *simple yet effective*.

566 **Comparison with Other Heuristic Optimization Strategies.** To sum up the proposed optimization  
 567 strategy, we empirically highlight the following points regarding the adopted heuristic strategies.  
 568 There exist a couple of approximate algorithms for CVaR optimization. In comparison, the crude  
 569 Monte Carlo one is the simplest for VaR estimates. It has an improvement guarantee under certain  
 570 conditions, leaving it easier to analyze. Besides, we also compare ours to the risk reweighted method  
 571 (Sagawa et al., 2020) in previously used benchmarks, and please take a closer look at that in Section  
 572 (K).

## 573 B Pseudo Algorithms of DR-MAML & DR-CNPs

---

### Algorithm 1: DR-MAML

---

**Input** : Task distribution  $p(\tau)$ ; Confidence level  $\alpha$ ; Task batch size  $B$ ; Learning rates:  $\lambda_1$  and  $\lambda_2$ .

**Output** : Meta-trained model parameter  $\vartheta$ .

Randomly initialize the model parameter  $\vartheta$ ;

**while** *not converged* **do**

Sample a batch of tasks  $f_{\tau_i} g_{i=1}^B \sim p(\tau)$ ;

*// inner loop via gradient descent*

**for**  $i = 1$  **to**  $B$  **do**

Evaluate the gradient:  $r \nabla \ell(\mathcal{D}_i^C; \vartheta)$  in Eq. (7);

Perform task-specific gradient updates:

$\vartheta_i \leftarrow \vartheta - \lambda_1 r \nabla \ell(\mathcal{D}_i^C; \vartheta)$ ;

**end**

*// estimate VaR*  $[\ell(T, \vartheta)] \hat{\xi}$

Evaluate performance  $L_B = f\ell(\mathcal{D}_{i=1}^T; \vartheta_i) g_{i=1}^B$ ;

Estimate VaR  $[\ell(T, \vartheta)]$  and set  $\xi = \hat{\xi}$  in Eq. (7) with either percentile rank or density estimators;

*// outer loop via gradient descent*

Screen the subset  $L_{\hat{\beta}} = f\ell(\mathcal{D}_{\hat{\lambda}_i}^T; \vartheta_i) g_{i=1}^K$  with  $\hat{\xi}$  for meta initialization updates;

$\vartheta \leftarrow \vartheta - \lambda_2 r \nabla \sum_{i=1}^K \ell(\mathcal{D}_{\hat{\lambda}_i}^T; \vartheta)$  in Eq. (7);

**end**

---

### Algorithm 2: DR-CNP

---

**Input** : Task distribution  $p(\tau)$ ; Confidence level  $\alpha$ ; Task batch size  $B$ ; Learning rate  $\lambda$ .

**Output** : Meta-trained model parameter  $\vartheta$ .

Randomly initialize the model parameter  $\vartheta$ ;

**while** *not converged* **do**

Sample a batch of tasks  $f_{\tau_i} g_{i=1}^B \sim p(\tau)$ ;

*// estimate VaR*  $[\ell(T, \vartheta)] \hat{\xi}$

Evaluate performance  $L_B = f\ell(\mathcal{D}_{i=1}^T; z, \vartheta_i) g_{i=1}^B$ ;

Estimate VaR  $[\ell(T, \vartheta)] \hat{\xi}$  with either percentile rank or density estimators;

*// execute gradient descent*

Screen the subset  $L_{\hat{\beta}} = f\ell(\mathcal{D}_{\hat{\lambda}_i}^T; z, \vartheta) g_{i=1}^K$  with  $\hat{\xi}$  for meta initialization updates;

$\vartheta \leftarrow \vartheta - \lambda r \nabla \sum_{i=1}^K \ell(\mathcal{D}_{\hat{\lambda}_i}^T; z, \vartheta)$  in Eq. (8);

**end**

---

## 576 C Properties of Risk Minimization Principles

### 577 C.1 Stochastic Optimization in the Constrained Form

578 In the section above, meta learning optimization objectives are discussed within three different  
 579 principles, respectively the average-case in Eq. (1), the worst-case in Eq. (2) and CVaR worst-case  
 580 in Eq. (5). This subsection continues this topic and introduces the relaxation variable  $\xi$  in the  
 581 optimization objective. In this way, the robust fast adaptation can be reframed in the case of stochastic  
 582 optimization with the probabilistic constraint.

583 We can equivalently express the minimization of the worst-case problem (Shalev-Shwartz and Wexler,  
 584 2016) within the following constrained stochastic optimization framework as follows.

$$\begin{aligned} & \min_{\tau \in \mathcal{T}; \xi \in \mathbb{R}^+} \xi \\ \text{s.t. } & \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \leq \xi \quad \text{with } \tau \in \mathcal{T} \end{aligned} \quad (9)$$

### 585 C.2 SGD Intractability of Meta Learning CVaR Optimization Objective

586 This subsection shows that directly stochastic gradient descent is intractable for meta learning to  
 587 optimize CVaR. Note that the normalized density distribution function  $p(\tau; \vartheta)$  implicitly depends  
 588 on  $\vartheta$  and  $\alpha$ , so we cannot access the exact form of such a distribution.

$$r_{\#} E(\vartheta) = \int p(\tau; \vartheta) \left[ \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) r_{\#} \ln p(\tau; \vartheta) + r_{\#} \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \right] d\tau \quad (10a)$$

$$\frac{1}{K} \sum_{k=1}^K \left[ \underbrace{\ell(\mathcal{D}_{k}^T, \mathcal{D}_{k}^C; \vartheta) r_{\#} \ln p(\tau_k; \vartheta) + r_{\#} \ell(\mathcal{D}_{k}^T, \mathcal{D}_{k}^C; \vartheta)}_{\text{Score Function}} \right] \quad (10b)$$

589 As illustrated in Eq. (10), the stochastic gradient estimate is not plausible since  $p(\tau; \vartheta)$  has no  
 590 closed form. Hence, heuristic or convex programming algorithms for specific cases are mostly used  
 591 as this domain's optimization strategy. However, designing optimization strategies for non-convex  
 592 cases is non-trivial in this domain and requires more consideration of theoretical guarantees.

### 593 C.3 Risk-Sensitive Applications & Optimization Strategies

594 **Related Applications.** There are a number of applications concerning robust optimization with  
 595 probabilistic constraints. Most of them are for the sake of safety. The risk principle CVaR firstly  
 596 occurs in the financial domain as a coherent principle (Rockafellar et al., 2000). It enjoys much  
 597 popularity in portfolio optimization (Quaranta and Zaffaroni, 2008). Gagne and Dayan (2021) adopt  
 598 CVaR principles to improve the distributional reinforcement learning performance. To robustify  
 599 robotic control, Pinto et al. (2017) varies hyper-parameters of Markov decision processes and  
 600 optimizes proportional worst-case trajectories in policy optimization within the principle of CVaR.  
 601 In work (Wilder, 2018), a CVaR related strategy is devised to solve submodular optimization  
 602 problems. Such a risk measure is also included in producing robust options (Hiraoka et al., 2019).  
 603 *Regarding probabilistic robust meta learning within CVaR principles, there exists scarce related*  
 604 *work until now.*

605 **Optimization Strategies.** Concerning the constrained stochastic optimization problem, we partic-  
 606 ularly overview related work in this subsection. In the past few decades, there emerge substantial  
 607 optimization strategies together with theoretical analysis for convex risk functions in CVaR op-  
 608 timization (Nemirovski et al., 2009; Fan et al., 2017; Meng et al., 2020; Levy et al., 2020; Duchi  
 609 and Namkoong, 2021; Wang et al., 2022). As for the min-max risk minimization principle, which  
 610 focuses on the worst case instead of a propositional worst cases, some researchers have designed  
 611 relaxation or other heuristic algorithms to handle convex or non-convex risk functions (Chen et al.,  
 612 2017; Collins et al., 2020; Jiang et al., 2021; Hsieh et al., 2021). Nevertheless, it remains challenging  
 613 to design algorithms with convergence guarantees for non-convex risk function cases. One of the  
 614 latest CVaR work on non-convex risk functions is (Sagawa et al., 2020), where a risk reweighted  
 615 algorithm for robust neural networks is proposed to handle distributional shifts. Moreover, most  
 616 CVaR optimization in non-convex risk functions follows this type of risk reweighted strategies in  
 617 applications.

618 **Remarks on Literature Work:** Risk functions for robust fast adaptation cases are mostly non-convex,  
619 and it is non-trivial to design optimization strategies with a convergence guarantee. Meanwhile,  
620 we also conduct comparison experiments between work in (Sagawa et al., 2020) and ours in meta  
621 learning downstream tasks. We refer the reader to Appendix (K) for more details and analysis.

## 622 D Computational Complexity

623 The primary difference between distributionally robust meta learning and expected risk based meta  
624 learning lies in that a fixed probabilistic portion of task gradients are considered in meta updates. Such  
625 an operation drives the optimization procedure to focus more on proportional vulnerable scenarios,  
626 increasing the robustness in worst cases.

627 However, the iteration of the surrogate function  $\varphi(\hat{\xi}_t; \vartheta)$  involves the screening of proportional  
628 worst cases since our strategies require extra evaluation of fast adaptation. This brings additional  
629 computational cost with complexity  $O(B \log(B))$ , where  $B$  is the number of tasks in each batch. On  
630 the other hand, the proposed strategy performs sub-gradient updates instead of complete gradient  
631 updates. This helps reduce computational cost with complexity  $O(\alpha B j M j)$  in each iteration, where  
632  $j M j$  corresponds to the scale of parameters of meta learning models and  $\alpha$  is the confidence level in  
633 CVaR optimization.

634 Universally analyzing the computational complexity in meta learning is intractable since various  
635 meta learning methods exist. Some are gradient-based ones, while some are non-parametric ones.  
636 The exact number of iterations for convergence heavily relies on specific methods. In summary, given  
637 the same number of iterations and a fixed confidence level, the computational complexity difference  
638 for CVaR optimization in meta learning scenarios is  $O(B(\alpha j M j \log(B)))$ .

639 To deepen understanding of our method, we explain more via specific examples. The main idea is to  
640 execute the sub-gradient operation over the batch of task gradients. Here we take the DR-MAML in  
641 Example (1) to show the operation and how the distributionally robust meta initialization is obtained:

$$\begin{aligned} \vartheta_{t+1}^{\text{meta}} &= \vartheta_t^{\text{meta}} - \lambda_1 \left[ \sum_{i=1}^B r_{\#} [\delta(\tau_i) \ell(\mathcal{D}_i^T; \vartheta_t^i)] \right], \\ \text{with } \vartheta_t^i &= \vartheta_t^{\text{meta}} - \lambda_2 r_{\#} \ell(\mathcal{D}_i^C; \vartheta), \tau_i = p(\tau), \end{aligned} \quad (11)$$

642 where  $\lambda_1$  is the outer loop learning rate,  $\lambda_2$  is the inner loop learning rate, and  $\delta(\tau_i)$  is the indicator  
643 variable. Here  $\delta(\tau_i) = 1$  in the case when the  $\ell(\mathcal{D}_i^T; \vartheta_t^i)$  falls into the  $(1 - \alpha)$ -probabilistic  
644 worst-case region otherwise  $\delta(\tau_i) = 0$ .

645 As for the optimal rate for convergence or the generalization bound, it is up to specific meta learning  
646 methods and risk minimization principles. For worst-case risk minimization for meta learning,  
647 there already exists theoretical analysis in previous work, e.g., optimization-based meta learning  
648 (Collins et al., 2020) when fast adaptation functions hold the convexity property. When it comes to  
649 more universal cases, considering diverse meta learning methods and optimization strategies, it is  
650 still challenging to estimate the optimal rate for convergence. Also, our considered scenarios exist  
651 distributional drift between meta training and meta testing task distributions, which makes it tough to  
652 derive the generalization bound.

653 Finally, note that our developed optimization strategies for meta learning are regardless of meta  
654 learning methods, and DR-MAMLs and DR-CNPs are merely two examples. For the sake of  
655 convenience, we only implement DR-MAML to compare with TR-MAML in this work.

## 656 E Proof of Remark (2)

657 *Remark (2).* If  $\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta)$  is convex w.r.t.  $\vartheta$ , then Eq.s (5)/(6) are also convex functions. In  
658 this case, the optimization objective Eq. (6) of our interest can be resolved with the help of convex  
659 programming algorithms (Fan et al., 2017; Meng et al., 2020; Levy et al., 2020).

660 **Proof:** We at first show that  $[\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \xi]^+ := \max_{\xi \in \mathcal{G}} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \xi$  is convex w.r.t.  $\vartheta$   
661 if  $\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)$  is convex w.r.t.  $\vartheta$ : For ease of derivation, let us redenote two functions  $f_1(\xi; \vartheta) :=$   
662  $\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \xi$  and  $f_2(\xi; \vartheta) := 0$ .

663 With any constant  $\lambda \in [0, 1]$  and any two parameters  $\vartheta_1 \in \mathcal{G}$  and  $\vartheta_2 \in \mathcal{G}$ , we can have:

$$[\ell(\mathfrak{D}^T, \mathfrak{D}^C; \lambda\vartheta_1 + (1 - \lambda)\vartheta_2) - \xi]^+ = f_i(\xi; \lambda\vartheta_1 + (1 - \lambda)\vartheta_2) \text{ [for some } i \in \{1, 2\}] \quad (12a)$$

$$\lambda f_1(\xi; \vartheta_1) + (1 - \lambda)f_1(\xi; \vartheta_2) \quad (12b)$$

$$\lambda \max_{\xi \in \mathcal{G}} f_1(\xi; \vartheta_1) + (1 - \lambda) \max_{\xi \in \mathcal{G}} f_1(\xi; \vartheta_2) \quad (12c)$$

$$\lambda [\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta_1) - \xi]^+ + (1 - \lambda) [\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta_2) - \xi]^+, \quad (12d)$$

664 which shows that the risk function with slack variables is convex w.r.t.  $\vartheta$ .

665 As  $\varphi(\xi; \vartheta) = \xi + \frac{1}{\tau} \mathbb{E}_{p(\cdot)} \left[ [\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \xi]^+ \right]$  is the convex combination of the above  
666 mentioned convex function, the resulting  $\varphi(\xi; \vartheta)$  is naturally convex w.r.t.  $\vartheta$ .

## 667 F Proof of Proposition (1)

**Assumption (1):** The meta risk function  $\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)$  is  $\beta$ -Lipschitz continuous w.r.t.  $\vartheta$ , which suggests: there exists a positive constant  $\beta$  such that  $\|\partial \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)\| \leq \beta$ :

$$|\ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta^0)| \leq \beta \|\vartheta - \vartheta^0\|.$$

668 **Assumption (2):** For meta risk function values, the risk cumulative distribution function  $F(\cdot; \vartheta)$  is  
669  $\beta$ -Lipschitz continuous w.r.t.  $\vartheta$ , and the implicit normalized probability density function of tasks  
670  $p(\tau; \vartheta)$  is  $\beta$ -Lipschitz continuous w.r.t.  $\vartheta$ .

**Assumption (3):** For any valid  $\vartheta \in \mathcal{G}$  and corresponding implicit normalized probability density function of tasks  $p(\tau; \vartheta)$ , the meta risk function value can be bounded by a positive constant  $L_{\max}$ :

$$\sup_{\vartheta \in \mathcal{G}} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \leq L_{\max}.$$

671 **Proposition (1):** Under Assumptions (1)/(2)/(3), the meta learning optimization objective  $E(\vartheta)$  in  
672 Eq. (3) is continuous w.r.t.  $\vartheta$ .

**Proof:** Suppose that  $\|\partial \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)\| \leq \beta$  and  $\|\partial p(\tau; \vartheta)\| \leq \beta$ , we can have the following inequality based on Assumption (2):

$$\left| p(\tau; \vartheta) - p(\tau; \vartheta^0) \right| \leq \beta \|\vartheta - \vartheta^0\|.$$

Meanwhile, we can have the following inequality based on Assumption (1):

$$\left| \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta^0) \right| \leq \beta \|\vartheta - \vartheta^0\|.$$

Together with the boundness Assumption (3) of meta risk function value:

$$\sup_{\vartheta \in \mathcal{G}} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \leq L_{\max},$$

673 we can roughly estimate the probabilistic constrained expected meta risk function values as follows:

$$\left| E(\vartheta) - E(\vartheta^\theta) \right| = \left| E_{p_\alpha(\cdot; \#)} \left[ \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \right] - E_{p_\alpha(\cdot; \#^\theta)} \left[ \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta^\theta) \right] \right| \quad (13a)$$

$$= \left| E_{p_\alpha(\cdot; \#)} \left[ \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \right] - E_{p_\alpha(\cdot; \#^\theta)} \left[ \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \right] \right| \quad (13b)$$

$$+ E_{p_\alpha(\cdot; \#^\theta)} \left[ \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta^\theta) \right] \quad (13c)$$

$$\int \left| p(\tau; \vartheta) - p(\tau; \vartheta^\theta) \right| \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) d\tau + E_{p_\alpha(\cdot; \#^\theta)} \left[ \left| \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta^\theta) \right| \right] \quad (13d)$$

$$\beta \int \vartheta^j \sup_{\frac{2}{\alpha}} f \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) g + \sup_{\frac{2}{\alpha}} f \beta \int g \vartheta^j \quad (13e)$$

$$= \left( \beta L_{\max} + \beta_{\max} \right) \int \vartheta^j \quad (13f)$$

674 From the above inequality, we can see that  $E(\vartheta) - E(\vartheta^\theta)$  is a  $(\beta L_{\max} + \beta_{\max})$ -Lipschitz  
675 continuous w.r.t.  $\vartheta$ . As a result, we demonstrate Proposition (1).

## 676 G Proof of Proposition (2)

*Proposition (2):* Suppose there exists  $\delta \geq 2R^+$  such that  $j\xi(\vartheta) - \hat{\xi}(\vartheta)^j < \delta$  with  $\hat{\xi}(\vartheta)$  an estimate of  $\xi(\vartheta)$ . Then there exists a constant  $\kappa = \max\left\{\frac{2}{1-\alpha}, \frac{1}{1-\alpha}\right\} g$  such that

$$\varphi(\hat{\xi}(\vartheta); \vartheta) - \kappa \delta < E(\vartheta) - \varphi(\hat{\xi}(\vartheta); \vartheta).$$

677 *Proof:* Based on the direct deduction in work (Rockafellar et al., 2000), we know that for any  $\vartheta \geq \cdot$ ,  
678 the inequality holds:  $\varphi(\xi; \vartheta) - \varphi(\hat{\xi}; \vartheta)$ .

679 In the case when  $\hat{\xi} = \xi + \delta$  with  $\delta \geq 2R^+$ , the probability space of the task can be respectively  
680 partitioned into three disjoint probability space:

$$P(\tau) = P\left( FM_{\#}^{-1}(\ell) j \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \leq \xi, \tau \geq g \right) \quad (14a)$$

$$P^+(\tau) = P\left( FM_{\#}^{-1}(\ell) j \xi < \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) \leq \xi + \delta, \tau \geq g \right) \quad (14b)$$

$$P^{++}(\tau) = P\left( FM_{\#}^{-1}(\ell) j \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) > \xi + \delta, \tau \geq g \right). \quad (14c)$$

681

682 As a result, we can estimate the difference between the two terms as follows:

$$\varphi(\hat{\xi}; \vartheta) - \varphi(\xi; \vartheta) = \hat{\xi} - \xi \quad (15a)$$

$$+ \frac{1}{1-\alpha} E_{P^+(\cdot)} \left[ \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta) - \xi \right] + \frac{1}{1-\alpha} E_{P^{++}(\cdot)} \left[ \xi - \hat{\xi} \right] \quad (15b)$$

$$\delta + \frac{1}{1-\alpha} E_{P^+(\cdot)} [\delta] - \frac{1}{1-\alpha} E_{P^{++}(\cdot)} [\delta] = \delta + \frac{1}{1-\alpha} \delta = \frac{2}{1-\alpha} \delta. \quad (15c)$$

683

684 Similarly, in the case when  $\hat{\xi} = \xi + \delta$  with  $\delta \geq \mathbb{R}^+$ , the probability space of the task can be  
 685 partitioned into three disjoint space with the probability respectively:

$$P^-(\tau) = P\left(\mathcal{FM}_{\#}^{-1}(\ell)j\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \mid \xi + \delta, \tau \geq g\right) \quad (16a)$$

$$P^0(\tau) = P\left(\mathcal{FM}_{\#}^{-1}(\ell)j \mid \delta < \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \mid \xi < 0, \tau \geq g\right) \quad (16b)$$

$$P^+(\tau) = P\left(\mathcal{FM}_{\#}^{-1}(\ell)j\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \mid \xi, \tau \geq g\right). \quad (16c)$$

686

687 As a result, we can estimate the difference between the two terms as follows:

$$\varphi(\hat{\xi}; \vartheta) - \varphi(\xi; \vartheta) = \xi \quad (17a)$$

$$+ \frac{1}{1-\alpha} E_{P^-(\cdot)}[\ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \mid \hat{\xi}] + \frac{1}{1-\alpha} E_{P^+(\cdot)}[\xi \mid \hat{\xi}] \quad (17b)$$

$$\delta + \frac{1}{1-\alpha} E_{P^-(\cdot)}[\delta] + \frac{1}{1-\alpha} E_{P^+(\cdot)}[\delta] \quad \delta + \frac{1}{1-\alpha} \delta = \frac{\alpha}{1-\alpha} \delta. \quad (17c)$$

688 Based on the inequalities (15)/(17), we can have  $\kappa = \max\{f_1^2, \frac{1}{1-\alpha}g\}$  such that the proposition is  
 689 verified as  $\varphi(\hat{\xi}; \vartheta) - \kappa \delta < \varphi(\xi; \vartheta)$ .

## 690 H Proof of Improvement Guarantee in Theorem (1)

691 **Theorem (1):** Under assumptions (1)/(2)/(3), suppose that the estimate error with the crude Monte  
 692 Carlo holds:  $j\hat{\xi}_t - \xi_t \leq \frac{1}{\epsilon(1-\alpha)^2}$ ,  $\forall t \geq \mathbb{N}^+$ , with the subscript  $t$  the iteration number,  $\lambda$   
 693 the learning rate in stochastic gradient descent,  $\beta$  the Lipschitz constant of the risk cumulative  
 694 distribution function, and  $\alpha$  the confidence level. Then the proposed heuristic algorithm with the  
 695 crude Monte Carlo can produce at least a local optimum for distributionally robust fast adaptation.

696 **Proof:** In the main paper, Fig. (2) provides the outline of the improvement guarantee proof. Note  
 697 that  $\hat{\xi}_t$  is an estimate of  $\xi_t$  with the help of Monte Carlo samples, and this depends on the model  
 698 parameters  $\vartheta_t$  in optimization. Performing the gradient updates w.r.t. the surrogate function  $\varphi(\hat{\xi}; \vartheta)$ ,  
 699 we can have the following equation with a small step-size learning rate  $\lambda$ :

$$\text{Gradient Descent : } \vartheta_{t+1} = \vartheta_t - \lambda r_{\#} \varphi(\hat{\xi}; \vartheta) \quad (18)$$

$$\text{) Monotonic Sequence : } \varphi(\hat{\xi}_t; \vartheta_{t+1}) \leq \varphi(\hat{\xi}_t; \vartheta_t).$$

700 To verify the improvement guarantee, we need to show that with the meta model parameters derived  
 701 from the surrogate function:

$$\varphi(\xi_{t+1}; \vartheta_{t+1}) \leq \varphi(\xi_t; \vartheta_t). \quad (19)$$

702 With the previous deduction  $\varphi(\xi_{t+1}; \vartheta_{t+1}) \leq \varphi(\xi_t; \vartheta_{t+1})$  from the property of CVaR, the  
 703 demonstration is equivalently reduced to show that:

$$\varphi(\xi_t; \vartheta_{t+1}) \leq \varphi(\xi_t; \vartheta_t). \quad (20)$$

704 We can perform the one order Taylor expansion with Peano's form of remainders w.r.t.  $\varphi(\xi_t; \vartheta)$   
 705 around the point  $\vartheta_t$ :

$$\varphi(\xi_t; \vartheta_{t+1}) = \varphi(\xi_t; \vartheta_t) + \lambda \left[ r_{\#} \varphi(\hat{\xi}_t; \vartheta) \Big|_{j_{\#}=\#_t} - r_{\#} \varphi(\xi_t; \vartheta) \Big|_{j_{\#}=\#_t} \right] + O(j_{\#} \vartheta_{t+1} - \vartheta_t) \varphi(\xi_t; \vartheta_t). \quad (21)$$

706 In the case when  $\hat{\xi} = \xi + \delta$  with  $\delta \geq \mathbb{R}^+$ , we use the partitioned task probability space in Eq. (14)  
 707 and can derive the gradient estimate:

$$r_{\#} \varphi(\hat{\xi}_t; \vartheta) \Big|_{j_{\#}=\#_t} = \frac{1}{1-\alpha} \left[ E_{P^-(\cdot)} \left[ r_{\#} \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \Big|_{j_{\#}=\#_t} \right] + \frac{1}{1-\alpha} \left[ E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \Big|_{j_{\#}=\#_t} \right] \right]. \quad (22)$$

708 With  $jF \cdot (\hat{\xi}; \vartheta) = F \cdot (\xi; \vartheta)j - \beta \cdot \delta$  in the Assumption (2),

$$E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]^T = E_{P^{++}(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] \quad (23a)$$

$$jj E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] jj = jj E_{P^{++}(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] jj \quad (23b)$$

$$E_{P^+(\cdot)} \left[ \sup_{\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj \right] = E_{P^{++}(\cdot)} \left[ \sup_{\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj \right] \quad (23c)$$

$$\beta \cdot \delta (1 - \alpha) \left( \sup_{\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj \right)^2 = \beta \cdot \delta (1 - \alpha) \mu^2, \quad (23d)$$

709 where  $\mu$  defines the  $\sup_{\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj$ .

710 Then we can derive the following inequality with the deduction from Eq. (23):

$$\varphi(\xi_{t+1}; \vartheta_{t+1}) = \varphi(\xi_t; \vartheta_t) \quad (24a)$$

$$\lambda \left[ r_{\#} \varphi(\hat{\xi}_{t+1}; \vartheta)_{j \neq \#_t}^T - r_{\#} \varphi(\xi_t; \vartheta)_{j \neq \#_t} \right] + O(j\vartheta_{t+1} - \vartheta_t) \quad (24b)$$

$$= \varphi(\xi_t; \vartheta_t) \quad (24c)$$

$$\frac{\lambda}{(1 - \alpha)^2} \left[ E_{P^{++}(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]^T - E_{P^{++}(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] \right] \quad (24d)$$

$$\frac{\lambda}{(1 - \alpha)^2} \left[ E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]^T - E_{P^{++}(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] \right] \quad (24e)$$

$$+ O(j\vartheta_{t+1} - \vartheta_t) \quad (24f)$$

$$\varphi(\xi_t; \vartheta_t) - \frac{\lambda j \nu_1 j_2^2}{(1 - \alpha)^2} + \beta \cdot \delta (1 - \alpha) \mu^2 + O(j\vartheta_{t+1} - \vartheta_t), \quad (24g)$$

711 where  $\nu_1 = E_{P^{++}(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]$ .

712 To ensure the existence of improvement guarantee, we need that the following inequality holds:

$$\varphi(\xi_t; \vartheta_t) - \frac{\lambda j \nu_1 j_2^2}{(1 - \alpha)^2} + \beta \cdot \delta (1 - \alpha) \mu^2 + O(j\vartheta_{t+1} - \vartheta_t) \geq \varphi(\xi_t; \vartheta_t). \quad (25)$$

713 Similarly, in the case when  $\hat{\xi} = \xi - \delta$  with  $\delta \geq R^+$ , we use the probability space of partitioned  
714 tasks in Eq. (16) and can derive the gradient estimate:

$$r_{\#} \varphi(\hat{\xi}_{t+1}; \vartheta)_{j \neq \#_t}^T = \frac{1}{1 - \alpha} \left[ E_{P(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] + E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] \right]. \quad (26)$$

715 With  $jF(\xi; \vartheta) = F(\xi; \vartheta) - \beta \cdot \delta$  in the Assumption (2), the following formula can be easily  
716 verified:

$$E_{P(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]^T = E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] \quad (27a)$$

$$jj E_{P(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] jj = jj E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] jj \quad (27b)$$

$$E_{P(\cdot)} \left[ \sup_{\mathcal{Z}_\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj \right] = E_{P^+(\cdot)} \left[ \sup_{\mathcal{Z}_\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj \right] \quad (27c)$$

$$\beta \cdot \delta (1 - \alpha) \left( \sup_{\mathcal{Z}_\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj \right)^2 = \beta \cdot \delta (1 - \alpha) \mu^2. \quad (27d)$$

717

718 Once again, we can derive the following inequality with the deduction from Eq. (27):

$$\varphi(\xi_{t+1}; \vartheta_{t+1}) = \varphi(\xi_t; \vartheta_t) - \lambda \left[ r_{\#} \varphi(\xi_t; \vartheta_t)_{j \neq \#_t}^T - r_{\#} \varphi(\xi_t; \vartheta_t)_{j \neq \#_t} \right] + O(j\vartheta_{t+1} - \vartheta_t) \quad (28a)$$

$$= \varphi(\xi_t; \vartheta_t) \quad (28b)$$

$$\frac{\lambda}{(1 - \alpha)^2} \left[ E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]^T - E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] \right] \quad (28c)$$

$$\frac{\lambda}{(1 - \alpha)^2} \left[ E_{P(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]^T - E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] \right] \quad (28d)$$

$$+ O(j\vartheta_{t+1} - \vartheta_t) \quad (28e)$$

$$\varphi(\xi_t; \vartheta_t) - \frac{\lambda jj \nu_2 jj^2}{(1 - \alpha)^2} + \beta \cdot \delta (1 - \alpha) \mu^2 + O(j\vartheta_{t+1} - \vartheta_t), \quad (28f)$$

719 where  $\nu_2 = E_{P^+(\cdot)} \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]$ .

720 To ensure the existence of improvement guarantee, we need that the following holds:

$$\varphi(\xi_t; \vartheta_t) - \frac{\lambda jj \nu_2 jj^2}{(1 - \alpha)^2} + \beta \cdot \delta (1 - \alpha) \mu^2 + O(j\vartheta_{t+1} - \vartheta_t) \geq \varphi(\xi_t; \vartheta_t). \quad (29)$$

721 Considering the above two cases and estimated bounds Eq. (25)/(29), we can roughly estimate the  
722 upper bound of the required  $\delta$  to guarantee performance improvement using our developed strategy  
723 in optimization:

$$\delta \leq \frac{\lambda jj \nu_m jj^2}{\beta \cdot (1 - \alpha)^3 \mu^2}, \quad (30)$$

724 where  $\lambda$  is the formerly mentioned learning rate,  $jj \nu_m jj^2$  is  $\max_{f, j} jj \nu_1 jj^2, jj \nu_2 jj^2 g$ , and  $\mu$  is supremum  
725 of the meta risk function derivatives in the task domain  $\sup_{\mathcal{Z}_\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj$ .

726 With the help of Jensen inequality,  $jj \nu_i jj^2$  can be roughly bounded as:

$$jj \nu_i jj^2 = E_{P^+(\cdot)} \left[ \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right]^T \left[ r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} \right] \right] \quad (31)$$

$$(1 - \alpha) \left( \sup_{\mathcal{Z}_\tau} jj r_{\#} \ell(\mathfrak{D}^T, \mathfrak{D}^C; \vartheta)_{j \neq \#_t} jj \right)^2 = (1 - \alpha) \mu^2.$$

727 With Eq.s (30)/(31), we can finally obtain the necessary condition for improvement guarantee:

$$\delta > \frac{\lambda}{\beta \cdot (1 - \alpha)^2}. \quad (32)$$

## 728 I Proof of Approximation Error in Theorem (2)

729 This section is to build up connections between the number of Monte Carlo samples in estimating  
730 VaR and the gap of solutions between the approximately derived one and the theoretical one.

**Theorem 2 (Gaps of Optimized Solutions)** Suppose  $F(\cdot; l; \vartheta) \in \mathcal{C}^2$  in  $l$ -domain. With meta trained  $\vartheta$ , the crude Monte Carlo estimate of  $\xi$ , the constant  $\kappa$  in Proposition (2), and the sufficiently large number of the task batch  $B$ , and  $R_B = O(B^{-3/4} \ln B)$ , we have the expected error between the exact optimum and the approximate optimum:

$$E(\vartheta) - \varphi(\hat{\xi}; \vartheta) + \kappa \left[ \frac{\alpha \hat{F}(\hat{\xi}; B, \vartheta)}{\frac{dF_\ell(\cdot; \#)}{d} \Big|_{j=\alpha}} + R_B \right].$$

731 *Proof:* As noted in (Bahadur, 1966), with assumptions that the cumulative distribution function  
732  $F(\cdot; l; \vartheta)$ 's second order derivative is continuous in  $l$ -domain, namely  $F(\cdot; l; \vartheta) \in \mathcal{C}^2$ , and  $\frac{dF_\ell(\cdot; \#)}{d} \Big|_{j=\alpha}$ ,  
733 the quantile estimate with crude Monte Carlo can be asymptotically written in the form with the help  
734 of central limit theory (Rosenblatt, 1956):

$$\hat{\xi} - \xi = \frac{\alpha \hat{F}(\hat{\xi}; B, \vartheta)}{\frac{dF_\ell(\cdot; \#)}{d} \Big|_{j=\alpha}} + R_B, \quad \text{with } R_B = O(B^{-3/4} \ln B) \text{ when } B \rightarrow \infty, \quad (33)$$

735 where the empirical cumulative distribution function is computed as follows.

736 With the sampled meta risk function values  $L_B = \{f_\ell(\mathcal{D}_i^T, \mathcal{D}_i^C; \vartheta) g_{i=1}^B\}$ , we rank them by values as  
737  $\hat{L}_B = \{f_\ell(\mathcal{D}_{\lambda_i}^T, \mathcal{D}_{\lambda_i}^C; \vartheta) g_{i=1}^B\}$ , which means  $\ell(\mathcal{D}_{\lambda_{i-1}}^T, \mathcal{D}_{\lambda_{i-1}}^C; \vartheta) \leq \ell(\mathcal{D}_{\lambda_i}^T, \mathcal{D}_{\lambda_i}^C; \vartheta)$ . Then the empirical  
738 cumulative distribution function with these order statistics (Barabás, 1987) can be written as:

$$\hat{F}(\xi; B, \vartheta) = \begin{cases} 0, & \text{if } \xi \leq \ell(\mathcal{D}_{\lambda_1}^T, \mathcal{D}_{\lambda_1}^C; \vartheta) \\ \frac{k}{B}, & \text{if } \ell(\mathcal{D}_{\lambda_k}^T, \mathcal{D}_{\lambda_k}^C; \vartheta) < \xi \leq \ell(\mathcal{D}_{\lambda_{k+1}}^T, \mathcal{D}_{\lambda_{k+1}}^C; \vartheta) \quad (k = 1, 2, \dots, B) \\ 1, & \text{if } \ell(\mathcal{D}_{\lambda_B}^T, \mathcal{D}_{\lambda_B}^C; \vartheta) < \xi. \end{cases} \quad (34)$$

739 Based on Proposition (2) and  $\kappa = \max_{\tau} f_{\tau}^2 - \tau - g$ , we know the inequality holds:

$$\varphi(\hat{\xi}; \vartheta) - E(\vartheta) < \kappa \delta. \quad (35)$$

With Eq. (33)/(35), the following inequality naturally holds when  $B$  is large enough:

$$\varphi(\hat{\xi}; \vartheta) - E(\vartheta) + \kappa \left[ \frac{\alpha \hat{F}(\hat{\xi}; B, \vartheta)}{\frac{dF_\ell(\cdot; \#)}{d} \Big|_{j=\alpha}} + R_B \right].$$

## 740 J Experimental Set-up & Implementation Details

741 This section is to provide experimental details in this paper. For the implementation of few-shot  
742 sinusoid regression and few-shot image classification, we respectively refer the reader to TR-MAML's  
743 codes (<https://github.com/lgcollins/tr-maml>) in (Collins et al., 2020) and vanilla MAML's  
744 codes (<https://github.com/AntreasAntoniou/HowToTrainYourMAMLPytorch>) in (Antoniu  
745 et al., 2019). And ours is built on top of the above codes except for simple modification of loss  
746 functions. The learning rates for the inner loop and the outer loop of all methods are the same as the  
747 above ones.

748 To facilitate the use of our heuristic optimization strategy, we leave the pytorch version of loss  
749 functions within the expected tail risk minimization. The example is provided in the case of mean  
750 square errors after MAML's inner loop, which is *simple to implement yet effective in robustifying fast  
751 adaptation*, as follows:

```

752 1 import torch
753 2 from torch.nn import MSELoss
754 3
755 4 def cvar_mse(y_pred, y_true, conf_level=0.5):
756 5
757 6     batch_MSE=MSELoss(reduction='none')
758 7     batch_loss=batch_MSE(y_pred, y_true)
759 8
760 9     # average risk values over non-task dimensions
76110    batch_avg_loss=torch.mean(batch_loss, dim=-1)
76211
76312    # crude Monte Carlo to estimate VaR and sub-tasks
76413    topk_mse, topk_idxs=torch.topk(batch_avg_loss, int((1-conf_level)*
765    y_true.size()[0]))
76614
76715    return torch.mean(topk_mse)

```

Listing 1: Loss Functions in Two-Stage Heuristic Algorithm with Crude Monte Carlo

## 768 J.1 Meta Learning Datasets & Tasks

769 **Sinusoid Regression.** MAML (Finn et al., 2017), TR-MAML (Collins et al., 2020), and DR-MAML  
770 are considered in this experiment. We retain the setup in task generation and partition in (Collins  
771 et al., 2020). More specifically, there exists a distribution drift between the meta-training and the  
772 meta-testing function families  $\tilde{f}_m(x) = a_m \sin(x - b_m) g_{m=1}^M$ .

773 Numerous easy tasks and a small proportion of difficult tasks are available in meta-training, while  
774 all tasks in the space are used in the evaluation. The range of the phase parameter is  $b \in [0, \pi]$ ,  
775 and the amplitude range of the parameter is  $a \in [0.1, 1.05]$  for easy tasks and  $a \in [4.95, 5.0]$  for  
776 difficult tasks. It is noted that the sinusoid task is more challenging to fit with larger amplitudes as the  
777 resulting function is less smooth. The loss function corresponds to the mean-squared error between  
778 the predicted value  $f(x)$  and the ground truth value. The number of task batches is 50 for 5-shot and  
779 25 for 10-shot. The optimal selection of the confidence level  $\alpha$  is difficult since we need to trade off  
780 the worst and average performance. Our setup is to minimize CVaR, which already considers the  
781 worst-case at some degree, so we watch the average performance in meta training results and set  
782  $\alpha = 0.7$  for all few-shot regression tasks. There is no external configuration for this hyper-parameter.  
783 The maximum number of iterations in meta-training is 70000.

784 **Few-shot Image Classification.** MAML (Finn et al., 2017), TR-MAML (Collins et al., 2020), and  
785 DR-MAML are considered in this experiment. The N-way K-shot classification corresponds to an  
786 N-classification problem with K-labeled examples available to the meta learner.

787 The Omniglot dataset consists of 1623 handwritten characters from 50 alphabets, with each 20  
788 examples. The task distribution is uniform for all task instances consisting of characters from one  
789 specific alphabet. The dataset split follows procedures in (Triantafillou et al., 2019). Finally, 25  
790 alphabets are used for meta-training, with 20 alphabets for meta-testing. The number of task batches  
791 is 16. The confidence level  $\alpha = 0.5$  is selected with the same criteria as that in few-shot regression  
792 tasks. The maximum number of iterations is 60000 in meta-training. As the construction of the  
793 Omniglot meta dataset is related to specific alphabets and the scale of combination for tasks is huge,  
794 this indicates randomly sampled meta-training tasks in the evaluation of the main paper Tables may  
795 not be used in meta-training.

796 The *mini*-ImageNet dataset is pre-processed according to (Larochelle, 2017). In detail, 64 classes  
797 are used for meta-training, with the remaining 36 classes for meta-testing. Tasks are generated as  
798 follows: 64 meta-training classes are randomly grouped into 8 meta-train tasks with the class numbers  
799  $f6, 7, 7, 8, 8, 9, 9, 10g$ , and the 36 meta-testing classes are processed in the same way. Finally, each  
800 task is built by sampling 1 image from 5 different classes within one task, resulting in a 5-way  
801 1-shot problem. The number of task batches is 4. The maximum number of iterations is 60000 in  
802 meta-training.

803 **J.2 Neural Architectures & Optimization**

804 **Sinusoid Regression.** MAML (Finn et al., 2017), TR-MAML (Collins et al., 2020), and DR-MAML  
 805 are considered in this experiment. We retain the neural architecture for regression problems in (Finn  
 806 et al., 2017; Collins et al., 2020). That is, we deploy a fully-connected neural network with two  
 807 hidden layers of 40 ReLU nonlinear activation units. All methods use one stochastic gradient descent  
 808 step as the inner loop.

809 **Few-shot Image Classification.** We retain the neural architecture for few-shot image classification  
 810 problems in (Finn et al., 2017; Collins et al., 2020). In detail, a four-layer convolutional neural  
 811 network is used for both Omniglot and *mini*-ImageNet datasets. All methods use one stochastic  
 812 gradient descent step as the inner loop.

813 **K Additional Experimental Results**

814 **More Quantitative Analysis.** Due to the page limit in the main paper, we include the  $\alpha$ 's sensitivity  
 815 experimental result in sinusoid 10-shot regression. As illustrated in Fig. (7), the trend is similar to  
 816 that in sinusoid 5-shot regression. Worst-case optimization degrades the average performance of TR-  
 817 MAML. DR-MAML is entangled with MAML in the average performance, while the performance  
 gap between them is significant in the worst-case.

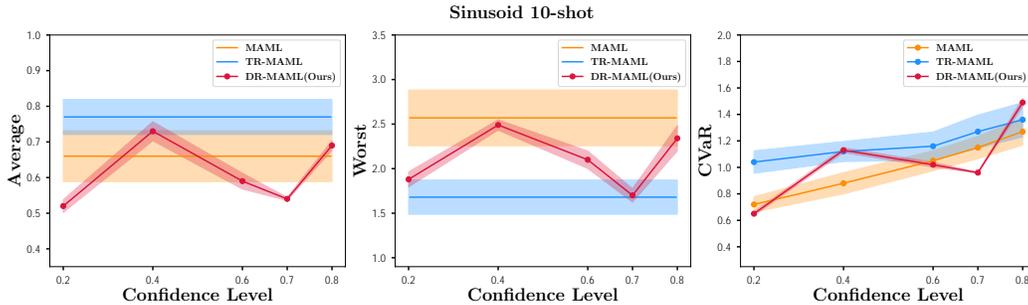


Figure 7: **Meta Testing Performance of Meta-Trained DR-MAML with Various Confidence Levels  $\alpha$ .** MAML and TR-MAML are irrelevant with the variation of  $\alpha$  in meta-training. The plots report testing MSEs with standard error bars in shadow regions.

818

819 Regarding few-shot image classification in the *mini*-ImageNet dataset, we can observe that in Fig.  
 820 (8), the standard error is relatively smaller than in previous regression cases. When the confidence  
 821 level is over a particular value, *e.g.*  $\alpha > 0.5$ , there occurs a significant decline of performance in all  
 822 metrics. Note that when  $\alpha \nearrow 1.0$ , the optimization objective approaches the worst-case optimization  
 823 objective. Here we attach two possible reasons for the performance degradation phenomenon: (i) The  
 824 adopted base optimization technique matters in nearly worst-case optimization. The stochastic mirror  
 825 descent-ascent (Juditsky et al., 2011) is utilized in TR-MAML, which is more stable in deriving the  
 826 optimal solution. In comparison, the stochastic gradient descent with sub-gradient operations works  
 827 as the optimization method, and this method can be unstable when the scale of worst-case examples  
 828 is small in the update. (ii) For few-shot image classification, estimates of VaR can be less precise  
 829 with limited batch sizes and higher  $\alpha$  values since the meta risk function value is discontinuous.  
 830 Consequently, we can also attribute the severe degradation of fast adaptation performance in higher  $\alpha$   
 831 value cases to the approximation errors of quantile estimates.

832 **More Visualization Results.** Further, we explore the influence of the expected tail risk minimization  
 833 principle in meta learning. Here the landscape of meta risk values, namely fast adaptation losses, is  
 834 presented in the sinusoid regression problem.

835 As exhibited in Fig. (9), the evaluated meta risk values from one random trial are associated with  
 836 hyper-parameters of tasks. The final optimized results can discover some tasks difficult in fast  
 837 adaptations. Meta learning methods are difficult to adapt in task regions with higher amplitudes.  
 838 MAML exhibits higher MSEs in regions with the amplitude  $a > 2$ , while DR-MAML minimizes  
 839 a proportion of risks in these regions. In contrast, TR-MAML reduces the risk around task regions  
 840 with  $a > 2$  to a certain extent; however, it shows relatively higher risk values in easy regions. Such

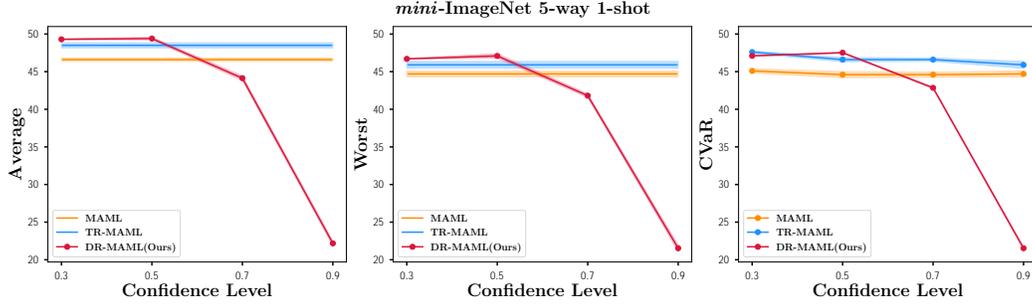


Figure 8: **Meta Testing Classification Accuracies of Meta-Trained DR-MAML with Various Confidence Levels  $\alpha$ .** MAML and TR-MAML are irrelevant with the variation of  $\alpha$  in meta-training. The plots report testing accuracies with standard error bars in shadow regions.

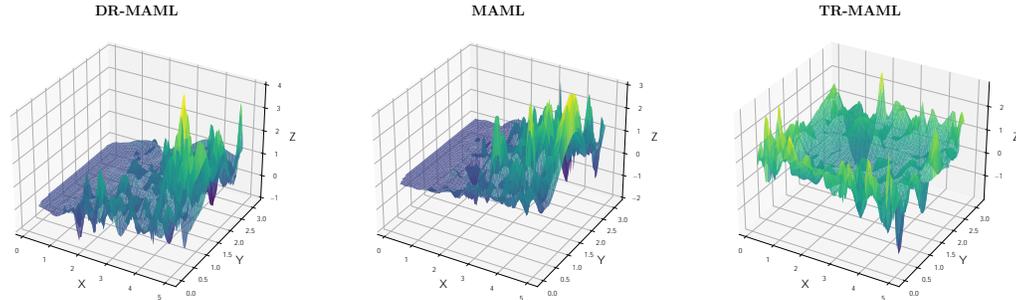


Figure 9: **The Fast Adaptation Risk Landscape of Meta-Trained DR-MAML, TR-MAML and MAML.** Shown is an example of sinusoid 5-shot regression, which corresponds to the function space  $f(x) = a \sin(x + b)$ . The  $X$ -axis denotes the amplitude parameter  $a$ , and the  $Y$ -axis is the phase parameter  $b$ . The confidence level is  $\alpha = 0.7$  in meta-training. The plots report testing MSEs in the  $Z$ -axis with a random trial of generating tasks.

841 evidence reflects the interpretability in optimization within the expected tail risk minimization, and  
 842 the landscape of meta risk values is relatively flat and smooth than others.

843 **Comparison with Other Optimization Strategies.** Note that instantiations of distributionally robust  
 844 meta learning methods, such as DR-MAML and DR-CNPs in Examples (1)/(2) are regardless of  
 845 optimization strategies and can be optimized via any heuristic algorithms for CVaR objectives.

846 Additionally, we use DR-MAML as the example and perform the comparison between our two-stage  
 847 algorithm and the risk reweighted algorithm (Sagawa et al., 2020). The intuition of the risk reweighted  
 848 algorithm is to relax the weights of tasks and assign more weights to the gradient of worst cases. The  
 849 normalization of risk weights is achieved via the softmax operator. Though there is an improvement  
 850 guarantee *w.r.t.* the probabilistic worst group of tasks, the algorithm is not specially designed for  
 851 meta learning or CVaR objective.

$$\min_{\#2} E(\vartheta) := E_{p(\cdot; \vartheta)} \left[ \frac{p(\tau; \vartheta)}{p(\tau; \vartheta)} \ell(\mathcal{D}^T, \mathcal{D}^C; \vartheta) \right] = \frac{1}{B} \sum_{b=1}^B \omega_b(\tau_b; \vartheta) \ell(\mathcal{D}_b^T, \mathcal{D}_b^C; \vartheta) \quad (36)$$

852 Meanwhile, the weight of task gradients after normalization is a biased estimator *w.r.t.* the constrained  
 853 probability  $p(\tau; \vartheta)$  in the task space. In other words, the risk reweighted method can be viewed  
 854 as approximation *w.r.t.* the importance weighted method in Eq. (36). In the importance weighted  
 855 method, for tasks out of the region of  $(1 - \alpha)$ -proportional worst, the probability of sampling such  
 856 tasks  $\tau_b$  is zero, indicating  $\omega_b(\tau_b; \vartheta) = 0$ . While in risk reweighted methods, the approximate weight  
 857 is assumed to satisfy  $\omega_b(\tau_b; \vartheta) \propto \exp\left(-\frac{\ell(\mathcal{D}_{\tau_b}^T, \mathcal{D}_{\tau_b}^C; \vartheta)}{\tau_b}\right)$ , where  $\vartheta$  means last time updated meta model  
 858 parameters and the risk function value is evaluated after fast adaptation.

859 In implementations, we keep the setup the same as Risk-Rweighted methods in (Paszke et al.,  
 860 2019) for meta-training. As illustrated in Table (4)/(5), DR-MAML with the two-stage optimization

Table 4: Test average mean square errors (MSEs) with reported standard deviations for sinusoid regression (5 runs). We mainly compare DR-MAML with different optimization algorithms. 5-shot and 10-shot cases are respectively considered here. The results are evaluated across the 490 meta-test tasks, which is the same as in (Collins et al., 2020). With  $\alpha = 0.7$  for meta training, the best testing results are in bold.

Method	5-shot						10-shot					
	Average		Worst		CVaR		Average		Worst		CVaR	
DR-MAML (Risk-Rewighted, (Sagawa et al., 2020))	0.91	0.06	3.57	0.56	1.83	0.03	0.61	0.02	1.90	0.11	1.13	0.02
DR-MAML (Two-Stage, Ours)	<b>0.89</b>	<b>0.04</b>	<b>2.91</b>	<b>0.46</b>	<b>1.76</b>	<b>0.02</b>	<b>0.54</b>	<b>0.01</b>	<b>1.70</b>	<b>0.17</b>	<b>0.96</b>	<b>0.01</b>

Table 5: Average 5-way 1-shot classification accuracies in *mini-ImageNet* with reported standard deviations (3 runs). We mainly compare DR-MAML with different optimization algorithms. With  $\alpha = 0.5$  for meta training, the best testing results are in bold.

Method	Eight Meta-Training Tasks						Four Meta-Testing Tasks					
	Average		Worst		CVaR		Average		Worst		CVaR	
DR-MAML (Risk-Rewighted, (Sagawa et al., 2020))	67.0	0.2	56.6	0.4	61.6	0.2	49.1	0.2	46.6	0.1	47.2	0.2
DR-MAML (Two-Stage, Ours)	<b>70.2</b>	<b>0.2</b>	<b>63.4</b>	<b>0.2</b>	<b>67.2</b>	<b>0.1</b>	<b>49.4</b>	<b>0.1</b>	<b>47.1</b>	<b>0.1</b>	<b>47.5</b>	<b>0.1</b>

861 strategies consistently outperform that with the risk-weighted ones in both 5-shot and 10-shot  
 862 sinusoid cases regarding all metrics. The performance advantage of using the two-stage ones is not  
 863 significant in *mini-ImageNet* scenarios. We can hypothesize that the estimate of VaR in continuous  
 864 task domains, *e.g.*, sinusoid regression, is more accurate, and this probabilistically ensures the  
 865 improvement guarantee with two-stage strategies. Both the VaR estimate in two-stage strategies  
 866 and the importance weight estimate in the risk-reweighted ones may have a lot of biases in few-shot  
 867 image classification, which lead to comparable performance.

## 868 L Platforms & Computational Tools

869 In this research project, we use NVIDIA 1080-Ti GPUs in computation. Pytorch (Paszke et al.,  
 870 2019) works as the deep learning toolkit in implementing few-shot image classification experiments.  
 871 Meanwhile, Tensorflow is the deep learning toolkit for implementing sinusoid few-shot regression  
 872 experiments.