

Robust Fast Adaptation from Adversarially Explicit Task Distribution Generation

Qi (Cheems) Wang^{†*}
Tsinghua University
Beijing, China

Yiqin Lv^{*}
Tsinghua University
Beijing, China

Yixiu Mao^{*}
Tsinghua University
Beijing, China

Yun Qu
Tsinghua University
Beijing, China

Yi Xu
Dalian University of Technology
Dalian, China

Xiangyang Ji[†]
Tsinghua University
Beijing, China

ABSTRACT

Meta-learning is a practical learning paradigm to transfer skills across tasks from a few examples. Nevertheless, the existence of *task distribution shifts* tends to weaken meta-learners' generalization capability, particularly when the training task distribution is naively hand-crafted or based on simple priors that fail to cover critical scenarios sufficiently. Here, we consider explicitly generative modeling task distributions placed over task identifiers and propose robustifying fast adaptation from adversarial training. Our approach, which can be interpreted as a model of a Stackelberg game, not only uncovers the task structure during problem-solving from an explicit generative model but also theoretically increases the adaptation robustness in worst cases. This work has practical implications, particularly in dealing with task distribution shifts in meta-learning, and contributes to theoretical insights in the field. Our method demonstrates its robustness in the presence of task subpopulation shifts and improved performance over SOTA baselines in extensive experiments. The code is available at the project site (<https://sites.google.com/view/ar-metalearn>).

CCS CONCEPTS

• **Computing methodologies** → *Meta learning; Generative modeling.*

KEYWORDS

Meta Learning, Generative Models, Game Theory

ACM Reference Format:

Qi (Cheems) Wang[†], Yiqin Lv, Yixiu Mao, Yun Qu, Yi Xu, and Xiangyang Ji. 2025. Robust Fast Adaptation from Adversarially Explicit Task Distribution Generation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V.1 (KDD '25)*, August 3–7, 2025, Toronto, ON, Canada. ACM, New York, NY, USA, 27 pages. <https://doi.org/10.1145/3690624.3709337>

^{*}These authors contributed equally to this research.

[†]Correspondence: cheemswang@mail.tsinghua.edu.cn; xyji@tsinghua.edu.cn

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
KDD '25, August 3–7, 2025, Toronto, ON, Canada
© 2025 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-1245-6/25/08
<https://doi.org/10.1145/3690624.3709337>

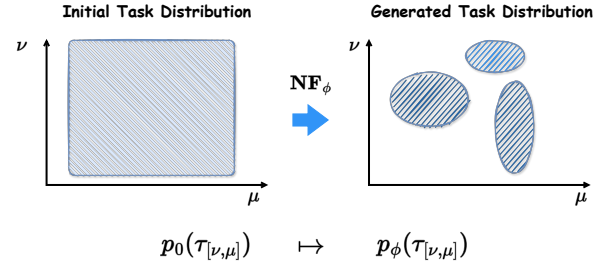


Figure 1: Diagram of Generating Task Distribution as the Adversary in Meta-Learning. Here, the initial task distribution $p_0(\tau)$ is a uniform distribution governed by two task identifiers $[\nu, \mu]$. Then, it is transformed into an explicit distribution $p_\phi(\tau)$ with the help of normalizing flows NF_ϕ .

1 INTRODUCTION

Deep learning has made remarkable progress in the past decade, ranging from academics to industry [36]. However, training deep learning models is generally time-consuming, and the previously trained model on one task might perform poorly in deployment when faced with unseen scenarios [39].

Fortunately, meta-learning, or learning to learn, offers a scheme to generalize learned knowledge to unseen scenarios [13, 17, 26, 27]. The strategy is to leverage past experience, extract meta knowledge as the prior, and utilize a few shot examples to transfer skills across tasks. This way, we can avoid learning from scratch and quickly adapt the model to unseen but similar tasks, catering to practical demands, such as fast autonomous driving in diverse scenarios. Due to these desirable properties, such a learning paradigm is playing an increasingly crucial role in building foundation models [3, 31, 48, 79].

Literature Challenges: Despite the promising adaptation performance in meta-learning, several concerns remain. Among them, the automatically task distribution design is under-explored and challenging in the field, which closely relates to the model's generalization evaluation [10, 93].

Overall, task identifiers configure the task, such as the topic type in the corpus for large language models [5, 79], the amplitude and phase in sinusoid functions, or the degree of freedom in robotic manipulators [2, 15]. Most existing studies adopt simple prior, such as uniform distributions over task identifiers [17, 19, 63], or hand-crafted distributions, which heavily rely on domain-specific knowledge difficult to acquire.

Some scenarios even pose more realistic demands for task distributions. In testing an autopilot system, an ideal task distribution

deserves more attention on traffic accidents or even generates some while covering typical cases [59, 71]. Similar circumstances also occur during domain randomization for embodied robots [47]. These imply that the shift between commonly used task distributions, such as uniform, and the expected testing distributions raises robustness issues and probably causes catastrophic failures when adapting to risk-sensitive scenarios [44].

Proposed Solutions: Rather than exploring fast adaptation strategies, we turn to *explicitly create task distribution shifts at a certain level and characterize robust fast adaptation with a Stackelberg game* [53]. To this end, we utilize normalizing flows to parameterize the distribution adversary in Figure 1 for task distribution generation and the meta learner for fast adaptation in the presence of distribution shifts.

Importantly, we constitute the solution concept, adopt the alternative gradient descent ascent to approximately compute the equilibrium [34], and conduct theoretical analysis. The optimization process can be translated as *fast adaptation robustification through adversarially explicit task distribution generation*.

Outline & Primary Contributions: The remainder starts with related work in Section 2. We define the notation and recap fundamentals in Section 3. Then, we present the game-theoretical framework to handle constrained task distribution shifts and robustify fast adaptation in Section 4. The quantitative analysis is conducted in Section 5, followed by conclusions and limitations. In primary, our contributions are:

- This work translates the robust fast adaptation under distribution shifts into a Stackelberg game [76]. To reveal task structures during problem-solving, we explicitly generate the task distribution with normalizing flows over task identifiers and optimize the meta-learner in an adversarial way.
- In theoretical analysis and tractable optimization, we constitute the solution concept *w.r.t.* fast adaptation, approximately solve the game using alternating stochastic gradient descent, and perform convergence and generalization analysis under certain conditions.

Extensive experimental results show that our approach can reveal adaptation-related structures in the task space and achieve robustness improvement in task subpopulation shifts.

2 LITERATURE REVIEW

The past few years have developed a large body of work on skill transfer across tasks or domain generalization in different ways [27, 87–89]. This section overviews the field regarding meta-learning and adaptation robustness.

Meta Learning. Meta learning is a learning paradigm that considers a distribution over tasks. The key is to pursue strategies for leveraging past experiences and distilling extracted knowledge into unseen tasks with a few shots of examples [8, 27, 49]. Currently, there are various families of meta-learning methods. The optimization-based ones, like model agnostic meta-learning (MAML) [17] and its extensions [14, 24, 58, 77], aim at finding a good meta-initialization of model parameters for adapting to all tasks via gradient descent. The deep metrics methods optimize the task representation in a metric space and are superior in few-shot image classification tasks [1, 28, 40, 69, 90]. Typical context-based methods, e.g.,

neural processes (NPs) and variants [19, 20, 23, 32, 60, 68, 78, 81–83], constitute the deep latent variable model as the stochastic process to accomplish tasks. Besides, memory-augmented networks [65], hyper-networks [25], and so forth are designed for meta-learning purposes.

Robustness in Meta Learning. In most previous work, the task distribution is fixed in the training set-up. In order to robustify the fast adaptation performance, a couple of learning strategies or principles emerge. Increasing the robustness to worst cases is a commonly seen consideration in adaptation, and these scenarios include input noise, parameter perturbation, and task distributions [7, 35, 41, 43, 52, 73, 94]. To alleviate the effects of adversarial examples in few-shot image classification, Goldblum et al. [21] meta-train the model in an adversarial way. To handle the distribution mismatch between training and testing tasks, Zhang et al. [95] adopt the adaptive risk minimization principle to enable fast adaptation. Wang et al. [80] propose to optimize the expected tail risk in meta-learning and witness the increase of robustness in proportional worst cases. Ours is a variant of a distributionally robust framework [84], and we seek equilibrium for fast adaptation.

Task Distribution Studies in Meta Learning. Task distributions are directly related to the generalization capability of meta-learning models, attracting increasing attention recently. Aiming to alleviate task overfitting, Murty et al. [50], Ni et al. [51], Rajendran et al. [57], Yao et al. [91] enrich the task space with augmentation techniques. Task relatedness can improve generalization across tasks, Fifty et al. [16] devise an efficient strategy to group tasks in multi-task training. In [42, 92], neural task samplers are developed to schedule the probability of task sampling in the context of few-shot classification. To increase the fidelity of generated tasks, Wu et al. [86] adopt the task representation model and constructs the up-sampling network for meta-training task augmentation. To reduce the required tasks, [38, 92] take the task interpolation strategy and shows that the interpolation strategy outperforms the standard set-up. Distinguished from the above, this work takes more interest in explicitly understanding task identifier structures concerning learning performance and cares about fast adaptation robustness under subpopulation shift constraints. Optimizing the task distribution might reserve the potential to improve generative performance in large models [6].

3 PRELIMINARIES

Notation. Throughout this paper, we use $p(\tau)$ to denote the task distribution with \mathcal{T} the task domain. Here, \mathcal{D}_τ represents the meta dataset with a sampled task τ . With the model parameter domain Θ and the support/query dataset construction, e.g., $\mathcal{D}_\tau = \mathcal{D}_\tau^S \cup \mathcal{D}_\tau^Q$, the risk function in meta-learning is a real-value function $\mathcal{L} : \mathcal{T} \times \Theta \mapsto \mathbb{R}$.

As an example, \mathcal{D}_τ consists of data points $\{(x_i, y_i)\}_{i=1}^{m+n}$ in few shot regression, and it is mostly split into the support dataset \mathcal{D}_τ^S for fast adaptation and query dataset \mathcal{D}_τ^Q for evaluation.

3.1 Problem Statement

To begin with, we revisit a couple of commonly-used risk minimization principles for meta-learning as follows.

Standard Meta-Learning Optimization Objective. We consider the meta-learning problem within the expected risk minimization principle in the statistical learning theory [75]. This results in the objective as Eq. (1), and we execute optimization in the form of task batches in implementation.

$$\min_{\theta \in \Theta} \mathbb{E}_{p(\tau)} \left[\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \right] \quad (1)$$

Here, θ refers to the meta-learning model parameters for meta knowledge and fast adaptation. The risk function depends on specific meta-learning methods. For example, in MAML, the form can be $\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) := \mathcal{L}(\mathcal{D}_\tau^Q; \theta - \lambda \nabla_\theta \mathcal{L}(\mathcal{D}_\tau^S; \theta))$ in regression, where the gradient update with the learning rate λ in the bracket reflects fast adaptation.

Distributionally Robust Meta Learning Optimization Objective. Recently, tail risk minimization has been adopted for meta-learning, effectively alleviating the effects towards fast adaptation in task distribution shifts [80]. In detail, we can express the optimization objective as Eq. (2) in the presence of the constrained distribution $p_\alpha(\tau; \theta)$, which characterizes the $(1 - \alpha)$ proportional θ -dependent worst cases in the task space.

$$\min_{\theta \in \Theta} \mathbb{E}_{p_\alpha(\tau; \theta)} \left[\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \right] \quad (2)$$

It is worth noting that $p_\alpha(\tau; \theta)$ is non-differentiable and θ -dependent with no closed-form. Meanwhile, the worst-case optimization for meta-learning in Eq. (3) can be treated as a particular instance of Eq. (2) when α sufficiently approaches 1.

$$\min_{\theta \in \Theta} \max_{\tau \in \mathcal{T}} \mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \quad (3)$$

Through tail risk minimization, the model's adaptation robustness can be enhanced *w.r.t.* the proportional worst scenarios [9, 80].

3.2 Two-Player Stackelberg Game

Before detailing our approach, it is necessary to describe elements in a two-player, non-cooperative Stackelberg game [76].

Let us assume two competitive players are involved in the game $\Gamma := \langle \{\mathcal{P}_1, \mathcal{P}_2\}, \{\theta \in \Theta, \phi \in \Phi\}, \mathcal{J}(\theta, \phi) \rangle$, where the meta learner as the leader \mathcal{P}_1 makes a decision first in the domain Θ while the distribution adversary as the follower \mathcal{P}_2 tries to deteriorate the leader decision's utility in the domain Φ . We refer to $\mathcal{J}(\theta, \phi)$ as the continuous risk function of the leader \mathcal{P}_1 , and that of the follower \mathcal{P}_2 corresponds to the negative form $-\mathcal{J}(\theta, \phi)$. Without loss of generality, all the players are rational and try to minimize risk functions in the game.

4 TASK ROBUST META LEARNING UNDER DISTRIBUTION SHIFT CONSTRAINTS

This section starts with the game-theoretic framework for meta-learning, followed by approximate optimization. Figure 2 shows a diagram of the constructed Stackelberg game. Then we perform theoretical analysis *w.r.t.* our approach.

4.1 Generate Task Distribution within A Game-Theoretic Framework

As part of an indispensable element in meta-learning, the task distribution is mostly set to be uniform or manually designed from

the heuristics. Such a setup hardly identifies a subpopulation of tasks that are tough to resolve in practice and fails to handle task distribution shifts.

In contrast, this paper considers an explicit task distribution to capture along with the learning progress and then automatically creates task distribution shifts for the meta-learner to adapt robustly. Our framework can be categorized as curriculum learning [4], but there places a constraint over the distribution shift in optimization.

Adversarially Task Robust Optimization with Distribution Shift Constraints. Now, we translate the meta-learning problem, namely generative task distributions for robust adaptation, into a min-max optimization problem:

$$\begin{aligned} \min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) &:= \mathbb{E}_{p_\phi(\tau)} \left[\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \right], \\ \text{s.t. } D_{KL} \left[p_0(\tau) \parallel p_\phi(\tau) \right] &\leq \delta, \end{aligned} \quad (4)$$

where the constraint term defines the maximum distribution shift to tolerate in meta training.

Equivalently, we can rewrite the above optimization objective in the form of unconstrained one with the help of a Lagrange multiplier $\lambda \in \mathbb{R}^+$:

$$\begin{aligned} \min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) &:= \mathbb{E}_{p_\phi(\tau)} \left[\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \right] \\ &\quad - \lambda \left[D_{KL} \left[p_0(\tau) \parallel p_\phi(\tau) \right] - \delta \right]. \end{aligned} \quad (5)$$

The above can be further simplified as:

$$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) := \mathbb{E}_{p_\phi(\tau)} \left[\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \right] + \lambda \mathbb{E}_{p_0(\tau)} \left[\ln p_\phi(\tau) \right], \quad (6)$$

where the constant terms, e.g., $\lambda \delta \in \mathbb{R}^+$ and $\mathbb{E}_{p_0(\tau)} \left[\ln p_0(\tau) \right]$ are eliminated.

As previously mentioned, the role of the distribution adversary attempts to transform the initial task distribution into one that raises challenging task proposals with higher probability. Such a setup drives the evolution of task distributions via adaptively shifting task sampling chance under constraints, which can be more crucial for generalization across risky scenarios. The term $D_{KL} \left[p_0(\tau) \parallel p_\phi(\tau) \right]$ inside Eq. (5) works as regularization to avoid the mode collapse in the generative task distribution. In Figure 2, the goal of the meta learner retains that of traditional meta-learning, while the distribution adversary continually generates the task distribution shifts along optimization processes.

ASSUMPTION 1 (LIPSCHITZ SMOOTHNESS AND COMPACTNESS). The adversarially task robust meta-learning optimization objective $\mathcal{J}(\theta, \phi)$ is assumed to satisfy

- (1) $\mathcal{J}(\theta, \phi)$ with $\forall [\theta, \phi] \in \Theta \times \Phi$ belongs to the class of twice differentiable functions \mathbb{C}^2 .
- (2) The norm of block terms inside Hessian matrices $\nabla^2 \mathcal{J}(\theta, \phi)$ is bounded, meaning that $\forall [\theta, \phi] \in \Theta \times \Phi$:

$$\sup \{ \|\nabla_{\theta, \theta}^2 \mathcal{J}\|, \|\nabla_{\theta, \phi}^2 \mathcal{J}\|, \|\nabla_{\phi, \phi}^2 \mathcal{J}\| \} \leq L_{\max}.$$

- (3) The parameter spaces $\Theta \subseteq \mathbb{R}^{d_1}$ and $\Phi \subseteq \mathbb{R}^{d_2}$ are compact with d_1 and d_2 respectively dimensions of model parameters for two players.

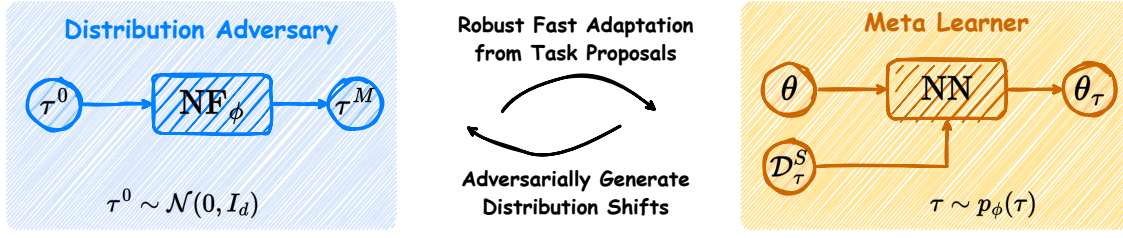


Figure 2: Diagram of Adversarially Task Robust Meta Learning. The proposed framework consists of two players, the distribution adversary and the meta player, in the game of meta-learning. On the left side of the figure: the distribution adversary seeks to transform the distribution from an initial task distribution, e.g., $\mathcal{N}(0, I_d)$ or $\mathcal{U}[a, b]$, via the neural network parameterized by ϕ with the purpose of deteriorating meta player’s fast adaptation performance. On the right side of the figure: the meta player parameterized by θ attempts to learn robust strategies for fast adaptation in sampled worst-case tasks (MAML algorithm [17] as an illustration).

EXAMPLE 1 (ADVERSARIALLY TASK ROBUST MAML, AR-MAML). Given the parameterized task distribution $p_\phi(\tau)$, the risk function \mathcal{L} and the learning rate γ in the inner loop of MAML [17], the adversarially task robust MAML corresponds to the following optimization problem:

$$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathbb{E}_{p_\phi(\tau)} \left[\mathcal{L}(\mathcal{D}_\tau^O; \theta - \gamma \nabla_\theta \mathcal{L}(\mathcal{D}_\tau^S; \theta)) \right] + \lambda \mathbb{E}_{p_0(\tau)} \left[\ln p_\phi(\tau) \right] \quad (7)$$

where \mathcal{D}_τ^S is used for the inner loop with \mathcal{D}_τ^O used for the outer loop.

EXAMPLE 2 (ADVERSARIALLY TASK ROBUST CNP, AR-CNP). Given the parameterized task distribution $p_\phi(\tau)$, the risk function \mathcal{L} , and the conditional neural process [19], the adversarially task robust CNP can be formulated as follows:

$$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathbb{E}_{p_\phi(\tau)} \left[\mathcal{L}(\mathcal{D}_\tau^O; z, \theta_2) \right] + \lambda \mathbb{E}_{p_0(\tau)} \left[\ln p_\phi(\tau) \right], \quad (8)$$

s.t. $z = h_{\theta_1}(\mathcal{D}_\tau^S)$ with $\theta = \{\theta_1, \theta_2\}$,

where θ_1 and θ_2 are respectively a set encoder and the decoder networks.

Here, we take two typical methods, e.g., MAML [17] and CNP [19], to illustrate the meta learner within the adversarially task robust framework, see Examples 1/2 for details.

Explicit Task Distribution Adversary Construction with Normalizing Flows. Learning to transform the task distribution is treated as a generative process: $\Phi: \mathcal{T} \rightarrow \mathcal{T} \subseteq \mathbb{R}^d$ in this paper. Admittedly, there already exist a collection of generative models to achieve the goal of generating task distributions, e.g., variational autoencoders [33, 62], generative adversarial networks [22], and normalizing flows [62].

Among them, we propose to utilize the normalizing flow [62] to achieve due to its tractability of the exact log-likelihood, flexibility in capturing complicated distributions, and a direct understanding of task structures. The basic idea of normalizing flows is to transform a simple distribution into a more flexible distribution with a series of invertible mappings $\mathcal{G} = \{g_i\}_{i=1}^M$, where $g_i: \mathcal{T} \rightarrow \mathcal{T} \subseteq \mathbb{R}^d$ indicates the smooth invertible mapping. We refer to these mappings implemented in the neural networks as NN_ϕ afterward. Specifically, with the base distribution $p_0(\tau)$ and a task sample τ^0 , the model

applies the above mappings to τ^0 to obtain τ^M .

$$\tau^M = g_M \circ \dots \circ g_2 \circ g_1(\tau^0) = \text{NN}_\phi(\tau^0) \quad (9)$$

In this way, the task distribution of interest is adaptive and adversarially exploits information from the shifted task distributions. The density function after transformations can be easily computed with the help of functions’ Jacobians:

$$\ln p_\phi(\tau^M) = \ln p_0(\tau^0) - \sum_{i=1}^M \ln \left| \det \frac{\partial g_i}{\partial \tau^{i-1}} \right|. \quad (10)$$

DEFINITION 1 (((ℓ_1, ℓ_2) -BI-LIPSCHITZ FUNCTION). An invertible function $g: x \subseteq \mathcal{X} \mapsto x \subseteq \mathcal{X}$, is said to be (ℓ_1, ℓ_2) -bi-Lipschitz if $\forall \{x_1, x_2\} \in \mathcal{X}$, the following conditions hold:

$$|g(x_1) - g(x_2)| \leq \ell_2 |x_1 - x_2| \quad \text{and} \quad |g^{-1}(x_1) - g^{-1}(x_2)| \leq \ell_1 |x_1 - x_2|.$$

As the normalizing flow function is invertible, the **Definition 1** is to describe the Lipschitz continuity in bi-directions.

4.2 Solution Concept & Explanations

This work separates players regarding the decision-making order, and the optimization procedure is no longer a simultaneous game. The nature of Stackelberg game enables us to technically express the studied asymmetric bi-level optimization problem as:

$$\min_{\theta \in \Theta} \mathcal{J}(\theta, \phi), \quad \text{s.t. } \phi \in \mathcal{S}(\theta) \quad (11)$$

with the θ -dependent conditional subset $\mathcal{S}(\theta) := \{\phi \in \Phi | \mathcal{J}(\theta, \phi) \geq \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi)\}$. This suggests the variables θ and ϕ are entangled in optimization.

Moreover, we can define the resulting equilibrium as a local minimax point [29] in adversarially task robust meta-learning, due to the non-convex optimization practice.

DEFINITION 2 (LOCAL MINIMAX POINT). The solution $\{\theta_*, \phi_*\}$ is called local Stackelberg equilibrium when satisfying two conditions: (1) $\phi_* \in \Phi' \subset \Phi$ is the maximum of the function $\mathcal{J}(\theta_*, \cdot)$ with Φ' a neighborhood; (2) $\theta_* \in \Theta' \subset \Theta$ is the minimum of the function $\mathcal{J}(\theta, g(\theta))$ with $g(\theta)$ the implicit function of $\nabla_\phi \mathcal{J}(\theta, \phi) = 0$ in the neighborhood Θ' .

Moreover, there exists a clearer interpretation w.r.t. the sequential optimization process and the equilibrium in the **Definition 2**. The meta learner as the leader first optimizes its parameter θ . Then

the distribution adversary as the follower updates the parameter ϕ and explicitly generates the task distribution proposal to challenge adaptation performance. In other words, we expect that meta learners can benefit from generative task distribution shifts regarding the adaptation robustness.

REMARK 1 (ENTROPY OF THE GENERATED TASK DISTRIBUTION). Given the generative task distribution $p_{\phi_*}(\tau)$, we can derive its entropy from the initial task distribution $p_0(\tau)$ and normalizing flows $\mathcal{G} = \{g_i\}_{i=1}^M$:

$$\mathbb{H}[p_{\phi_*}(\tau)] = \mathbb{H}[p_0(\tau)] + \int p_0(\tau) \left[\sum_{i=1}^M \ln \left| \det \frac{\partial g_i}{\partial \tau^{i-1}} \right| \right] d\tau. \quad (12)$$

The above implies that the generated task distribution entropy is governed by the change of task identifiers in the probability measure of the task space.

4.3 Strategies for Finding Equilibrium

Given the previously formulated optimization objective, we propose to approach it with the help of estimated stochastic gradients. As noticed, the involvement of adaptive expectation term $p_{\phi}(\tau)$ requires extra considerations in optimization.

Best Response Approximation. Given two players with completely distinguished purposes, the commonly used strategy to compute the equilibrium is the Best Response (BR), which means:

$$\theta_{t+1} = \arg \min_{\theta \in \Theta} \mathcal{J}(\theta, \phi_t) \quad (13a)$$

$$\phi_{t+1} = \arg \max_{\phi \in \Phi} \mathcal{J}(\theta_{t+1}, \phi). \quad (13b)$$

For implementation convenience, we instead apply the gradient updates to the meta player and the distribution adversary, namely stochastic alternating gradient descent ascent (GDA). The operations are entangled and result in the following iterative equations with the index t :

$$\theta_{t+1} \leftarrow \theta_t - \gamma_1 \nabla_{\theta} \mathcal{J}(\theta_t, \phi_t) \quad (14a)$$

$$\phi_{t+1} \leftarrow \phi_t + \gamma_2 \nabla_{\phi} \mathcal{J}(\theta_{t+1}, \phi_t). \quad (14b)$$

This can be viewed as the gradient approximation for the BR strategy, which leads to at least a local Stackelberg equilibrium for the considered minimax problem [30].

Stochastic Gradient Estimates & Variance Reduction. Addressing the game-theoretic problem is non-trivial especially when it relates to distributions. A commonly-used method is to perform the sample average approximation w.r.t. Eq. (14). It iteratively updates the parameters of the meta player and the distribution adversary to approximate the saddle point.

More specifically, we can have the Monte Carlo estimates of the stochastic gradients for the leader \mathcal{P}_1 :

$$\begin{aligned} \nabla_{\theta} \mathcal{J}(\theta, \phi) &= \int p_{\phi}(\tau) \nabla_{\theta} \mathcal{L}(D_{\tau}^Q, D_{\tau}^S; \theta) d\tau \\ &\approx \frac{1}{K} \sum_{k=1}^K \nabla_{\theta} \mathcal{L}(D_{\tau_k}^Q, D_{\tau_k}^S; \theta). \end{aligned} \quad (15)$$

The form of stochastic gradients w.r.t. the meta player parameter θ is the meta-learning algorithm specific or model dependent. We refer the reader to Algorithm 1/2 as examples.

Now, we can derive the estimates with the help of REINFORCE algorithm [85] for the follower \mathcal{P}_2 and obtain the score function as:

$$\begin{aligned} \nabla_{\phi} \mathcal{J}(\theta, \phi) &\approx \frac{1}{K} \sum_{k=1}^K \mathcal{L}(D_{\tau_k}^Q, D_{\tau_k}^S; \theta) \nabla_{\phi} \ln p_{\phi}(\tau_k) \\ &\quad + \frac{\lambda}{K} \sum_{k=1}^K \nabla_{\phi} \ln p_{\phi}(\tau_k^{-M}), \end{aligned} \quad (16)$$

where the particle $\tau_k \sim p_{\phi}(\tau)$ denotes the task sampled from the generative task distribution, and τ_k^{-M} means the particle sampled from the initial task distribution to enable $\text{NN}_{\phi}(\tau_k^{-M}) = \tau_k$.

As validated in [18], the score estimator is an unbiased estimate of $\nabla_{\phi} \mathcal{J}(\theta, \phi)$. However, such a gradient estimator in Eq. (16) mostly exhibits higher variances, which weakens the stability of training processes. To reduce the variances, we utilize the commonly-used trick by including a constant baseline $\mathcal{V} = \mathbb{E}_{p_{\phi}(\tau)} [\mathcal{L}(D_{\tau}^Q, D_{\tau}^S; \theta)] \approx \frac{1}{K} \sum_{k=1}^K \mathcal{L}(D_{\tau_k}^Q, D_{\tau_k}^S; \theta)$ for the score function, which results in:

$$\begin{aligned} \nabla_{\phi} \mathcal{J}(\theta, \phi) &\approx \frac{1}{K} \sum_{k=1}^K [\mathcal{L}(D_{\tau_k}^Q, D_{\tau_k}^S; \theta) - \mathcal{V}] \nabla_{\phi} \ln p_{\phi}(\tau_k) \\ &\quad + \frac{\lambda}{K} \sum_{k=1}^K \nabla_{\phi} \ln p_{\phi}(\tau_k^{-M}). \end{aligned} \quad (17)$$

Particularly, since the normalizing flow works as the distribution transformation in this work, please refer to Eq. (10) to obtain the derivative of the log-likelihood of the transformed task $\ln p_{\phi}(\tau)$ w.r.t. ϕ inside Eq. (17). For easier analysis, we characterize the iteration sequence in optimization as $\begin{bmatrix} \theta_0 \\ \phi_0 \end{bmatrix} \mapsto \dots \mapsto \begin{bmatrix} \theta_t \\ \phi_t \end{bmatrix} \mapsto \begin{bmatrix} \theta_{t+1} \\ \phi_{t+1} \end{bmatrix} \mapsto \dots$.

REMARK 2 (SOLUTION AS A FIXED POINT). The alternating GDA for solving Eq. (5) results in the fixed point when $\begin{bmatrix} \theta_{H+1} \\ \phi_{H+1} \end{bmatrix} = \begin{bmatrix} \theta_H \\ \phi_H \end{bmatrix}$, or in other words $\begin{bmatrix} \theta_H \\ \phi_H \end{bmatrix}$ is stationary $\nabla \mathcal{J}(\theta_H, \phi_H) = 0$.

4.4 Theoretical Analysis

Built on the deduction of the local Stackelberg equilibrium's existence and the **Remark 2**, we further perform analysis on the considered equilibrium $\begin{bmatrix} \theta_* \\ \phi_* \end{bmatrix}$, in terms of learning dynamics using the alternating GDA. For notation simplicity, we denote the block terms inside the Hessian matrix $\mathbf{H}_* := \nabla^2 \mathcal{J}(\theta_*, \phi_*)$ around

$$[\theta_*, \phi_*]^T \text{ as } \begin{bmatrix} \nabla_{\theta\theta}^2 \mathcal{J} & \nabla_{\theta\phi}^2 \mathcal{J} \\ \nabla_{\phi\theta}^2 \mathcal{J} & \nabla_{\phi\phi}^2 \mathcal{J} \end{bmatrix} \big|_{[\theta_*, \phi_*]^T} := \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix}.$$

THEOREM 1 (CONVERGENCE GUARANTEE). Suppose that the **Assumption 1** and the function condition of the (local) Stackelberg equilibrium $\Delta(\mathbf{A}, \mathbf{B}, \mathbf{C}, \gamma_1, \gamma_2) < \frac{1}{2}$ are satisfied, where norms of the corresponding matrix are involved. Then the following statements hold:

- (1) The resulting iterated parameters $\{\dots \mapsto [\theta_t, \phi_t]^T \mapsto [\theta_{t+1}, \phi_{t+1}]^T \mapsto \dots\}$ are Cauchy sequences;

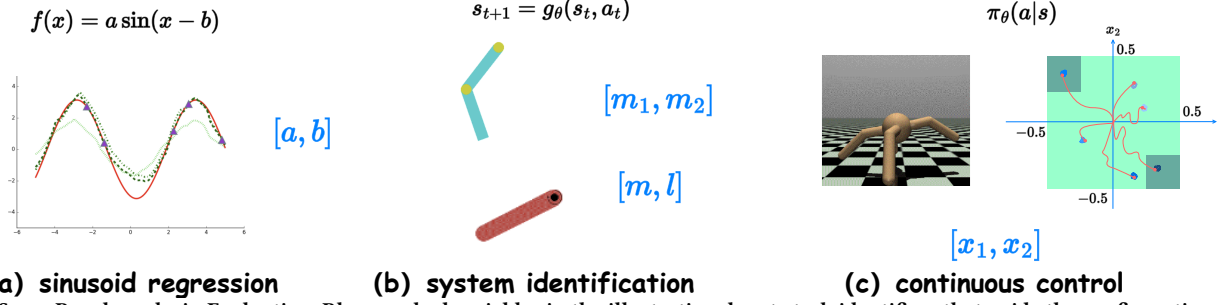


Figure 3: Some Benchmarks in Evaluation. Blue-marked variables in the illustration denote task identifiers that guide the configuration of a specific task. We place distributions over these task identifiers in generating diverse tasks for meta-learning.

- (2) The optimization can guarantee at least the linear convergence to the local Stackelberg equilibrium with the rate $\sqrt{\Delta}$.

The **Theorem 1** clarifies learning rates γ_1 and γ_2 's influence on convergence and the required second-order derivative conditions of the resulting stationary point $[\theta_*, \phi_*]^T$. And when the game arrives at convergence, the local Stackelberg equilibrium is the best response to these two players, which is at least a local min-max solution to Eq. (5).

Next, we estimate the generalization bound of meta learners when confronting the generated task distribution shifts.

THEOREM 2 (GENERALIZATION BOUND WITH THE DISTRIBUTION ADVERSARY). Given the pretrained normalizing flows $\{g_i\}_{i=1}^M$, where g_i is (ℓ_a, ℓ_b) -bi-Lipschitz, and the pretrained meta learner $\theta_* \in \Theta$, we can derive the generalization bound with the initial task distribution p uniform:

$$R_p^\omega(\theta_*) \leq \hat{R}_p^\omega(\theta_*) + \Upsilon(\mathcal{T}) \left(\frac{C \ln \frac{2Ke}{C} + \ln \frac{4}{\delta}}{K} \right)^{\frac{3}{8}}, \quad (18)$$

where $C = \text{Pdim}(\{\mathcal{L}(\cdot; \theta) : \theta \in \Theta\})$ denotes the pseudo-dimension in [56], $R_p^\omega(\theta_*)$ and $\hat{R}_p^\omega(\theta_*)$ are expected and empirical risks.

We refer the reader to Appendix F for formal **Theorem 2** and proofs. It reveals the connection between the bound and task complexity $\Upsilon(\mathcal{T})$, and more training tasks from initial distributions decrease the generalization error in adversarially distribution shifts.

5 EXPERIMENTS

Previous sections recast the adversarially task robust meta-learning to a Stackelberg game, specify the equilibrium, and analyze theoretical properties in distribution generation. This section focuses on the evaluation, and baselines constructed from typical risk minimization principles are reported in **Table 3**. These include vanilla MAML [17], DRO-MAML [64], TR-MAML [9], DR-MAML [80], and AR-MAML (ours).

Technically, we mainly answer the following **Research Questions (RQs)**:

- (1) Does adversarial training help improve few-shot adaptation robustness in case of task distribution shifts?
- (2) How does the type of the initial task distribution influence the performance of resulting solutions?

- (3) Can generative modeling the task distribution discover meaningful task structures and afford interpretability?

Implementation & Examination Setup. As our approach is agnostic to meta-learning methods, we mainly employ AR-MAML as the implementation of this work. Concerning the meta testing distribution, tasks are from the initial task distribution and the adversarial task distribution, respectively. The latter corresponds to the generated task distribution under shift constraints after convergence.

Evaluation Metrics. Here, we use both the average risk and conditional value at risk (CVaR $_\alpha$) in evaluation metrics, where CVaR $_\alpha$ can be viewed as the worst group performance in [64].

5.1 Benchmarks

We consider the few-shot synthetic regression, system identification, and meta reinforcement learning to test fast adaptation robustness with typical baselines. Notably, the task is specified by the generated task identifiers as shown in Figure 3.

Synthetic Regression. The same as that in [17], we conduct experiments in sinusoid functions. The goal is to uncover the function $f(x) = a \sin(x - b)$ with K -shot randomly sampled function points. And the task identifiers are the amplitude a and phase b .

System Identification. Here, we take the Acrobot System [70] and the Pendulum System [37] to perform system identification. In the Acrobot System, we generate different dynamical systems as tasks by varying masses of two pendulums. And the task identifiers are the pendulum mass parameters m_1 and m_2 . In the Pendulum System, the system dynamics are distinguished by varying the mass and the length of the pendulum. And the task identifiers are the mass parameter m and the length parameter l . For both benchmarks, we collect the dataset of state transitions with a complete random policy to interact with sampled environments. The goal is to predict state transitions conditioned on randomly sampled context transitions from an unknown dynamical system.

Meta Reinforcement Learning. We evaluate the role of task distributions in meta-learning continuous control. In detail, the Point Robot in [17] and the Ant-Pos Robot in Mujoco [74] are included as navigation environments. We respectively vary goal/position locations as task identifiers within a designed range to generate diverse tasks. The goal is to seek a policy that guides the robot to the target location with a few episodes derived from an environment.

We refer the reader to Appendix I for set-ups, hyper-parameter configurations and additional experimental results.

Table 1: Average mean square errors in 5-shot sinusoid regression/10-shot Acrobot system identification/10-shot Pendulum system identification with reported standard deviations (5 runs). With $\alpha = 0.5$, the best results are in pink (the lower, the better). U/N in benchmarks denote Uniform/Normal as the initial distribution type.

Benchmark	Meta-Test Distribution	Average					CVaR				
		MAML	TR-MAML	DR-MAML	DRO-MAML	AR-MAML	MAML	TR-MAML	DR-MAML	DRO-MAML	AR-MAML
Sinusoid-U	Initial	0.499±0.01	0.539±0.01	0.479±0.01	0.481±0.01	0.459±0.01	0.858±0.01	0.868±0.02	0.793±0.02	0.816±0.02	0.782±0.03
	Adversarial	0.508±0.01	0.548±0.01	0.499±0.01	0.502±0.02	0.405±0.01	0.883±0.02	0.879±0.02	0.836±0.01	0.826±0.03	0.671±0.01
Sinusoid-N	Initial	0.578±0.03	0.628±0.01	0.556±0.01	0.562±0.02	0.554±0.02	1.017±0.05	1.017±0.02	0.932±0.02	0.983±0.03	0.947±0.03
	Adversarial	0.496±0.01	0.511±0.01	0.492±0.02	0.493±0.01	0.404±0.02	0.838±0.03	0.827±0.02	0.807±0.03	0.835±0.01	0.672±0.03
Acrobot-U	Initial	0.244±0.01	0.233±0.00	0.222±0.00	0.237±0.00	0.219±0.01	0.336±0.01	0.320±0.00	0.303±0.00	0.322±0.01	0.298±0.00
	Adversarial	0.243±0.00	0.238±0.01	0.235±0.01	0.244±0.00	0.230±0.00	0.341±0.01	0.320±0.01	0.325±0.01	0.333±0.01	0.306±0.01
Acrobot-N	Initial	0.231±0.00	0.225±0.00	0.227±0.00	0.222±0.00	0.215±0.00	0.321±0.01	0.311±0.00	0.316±0.01	0.309±0.01	0.301±0.01
	Adversarial	0.246±0.00	0.237±0.00	0.241±0.00	0.242±0.00	0.229±0.00	0.338±0.00	0.327±0.01	0.327±0.00	0.332±0.01	0.314±0.01
Pendulum-U	Initial	0.648±0.02	0.694±0.01	0.634±0.01	0.630±0.02	0.627±0.01	0.799±0.03	0.780±0.02	0.744±0.01	0.751±0.03	0.733±0.02
	Adversarial	0.672±0.01	0.724±0.01	0.669±0.01	0.674±0.00	0.660±0.01	0.845±0.02	0.854±0.02	0.808±0.02	0.826±0.01	0.778±0.01
Pendulum-N	Initial	0.596±0.00	0.637±0.01	0.574±0.01	0.582±0.00	0.586±0.01	0.715±0.01	0.720±0.01	0.685±0.01	0.695±0.01	0.694±0.01
	Adversarial	0.664±0.02	0.702±0.01	0.660±0.02	0.677±0.02	0.635±0.01	0.861±0.03	0.837±0.02	0.817±0.03	0.860±0.04	0.777±0.03

5.2 Empirical Result Analysis

Here, we report the experimental results, perform analysis and answer the raised RQs (1)/(2).

Overall Performance: Table 1 shows that AR-MAML mostly outperforms others in the adversarial distribution, seldom sacrificing performance in the initial distribution. Similar to observations in [80], task distributionally robust optimization methods, like DR-MAML and DRO-MAML, not only retain robustness advantage on shifted distribution but also sometimes boost average performance on the initial distribution. Cases with two types of initial task distributions (Uniform/Normal) come to similar conclusions on average and CVaR $_{\alpha}$ performance. Figures 4/5 show the meta reinforcement learning results for Point Robot and Ant Pos navigation tasks. AR-MAML exhibits similar superiority on both continuous control benchmarks compared to baselines.

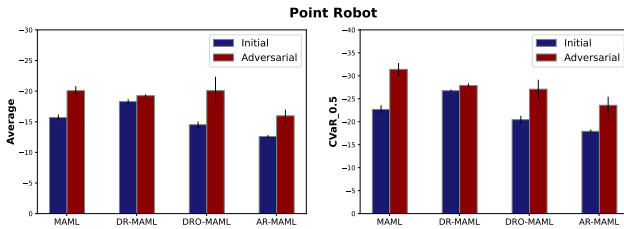


Figure 4: Meta Testing Returns in Point Robot Navigation Tasks (4 runs). The charts report average and CVaR $_{\alpha}$ returns with $\alpha = 0.5$ in initial and adversarial distributions, with standard error bars indicated by black vertical lines. The higher, the better.

Multiple Tail Risk Robustness: Note that CVaR metrics imply the model’s robustness under the subpopulation shift. Figure 6 reports CVaR $_{\alpha}$ values with various confidence values on pendulum system identification. The AR-MAML’s merits in handling the proportional worst cases are consistent across diverse levels. We also illustrate and include these statistics on other benchmarks in Appendix J. Moreover, as suggested in [72], a robust learner seldom encounters a performance gap between a standard (initial) test set and a test set

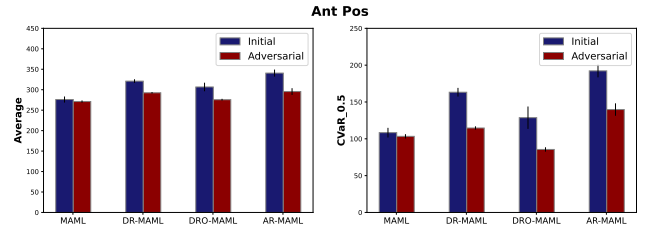


Figure 5: Meta Testing Returns in Ant Pos Tasks (4 runs). The charts report average and CVaR $_{\alpha}$ returns with $\alpha = 0.5$ in initial and adversarial distributions, with standard error bars indicated by black vertical lines. The higher, the better.

with a distribution shift (adversarial). Figure 7 validates the meta-learners’ robustness on sinusoid regression, where AR-MAML’s results are more proximal to the $y = x$ line than other baselines.

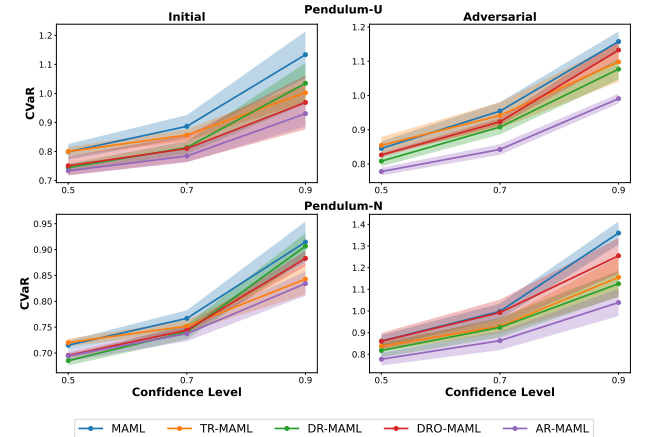


Figure 6: CVaR $_{\alpha}$ MSEs with Various Confidence Level α . Pendulum-U/N denotes Uniform/Normal as the initial distribution type. The plots report meta testing CVaR $_{\alpha}$ MSEs in initial and adversarial distributions with standard error in shadow regions.

Random Perturbation Robustness: We also test meta-learners’ robustness to random noise from the support dataset. To do so, we

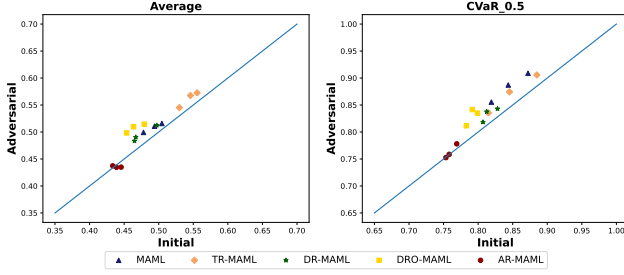


Figure 7: Meta testing MSEs on the initial distribution (x-axis) and on the adversarial distribution (y-axis). The $y = x$ line serves as a baseline for comparison. Models above this line show increased losses when faced with distribution shifts, indicating a decline in performance compared to the standard test set.

take sinusoid regression and inject random noise into the support set, i.e., the noise is drawn from a Gaussian distribution $\mathcal{N}(0, 0.1^2)$ and added to the output y . Figure 8 illustrates that AR-MAML’s performance degradation is somewhat less than others on the adversarial distribution. The noise exhibits similar effects on AR-MAML and DR-MAML on the initial distribution, harming performance severely. AR-MAML and DR-MAML still exhibit lower MSEs than other baselines for all cases. This indicates the adversarial training mechanism can also bring more robustness to challenging test scenarios with random noise.

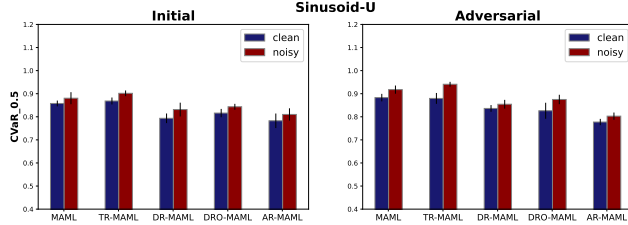


Figure 8: Meta Testing Performance in Clean and Noisy Tasks. The noisy tasks are constructed by adding noise on the outputs of the support dataset. Reported are testing CVar_α MSEs with $\alpha = 0.5$, where black vertical lines indicate standard error bars.

5.3 Task Structure Analysis

In response to RQ (3), we turn to the analysis of the learned distribution adversary. As a result, we visualize the adversarial task probability density.

Explicit Task Distribution: As displayed in Figure 9, our approach enables the discovery of explicit task structures regarding problem-solving. The general learned patterns seem to be regardless of the initial task distributions. In sinusoid regression, more probability mass is allocated in the region with $[3.0, 5.0] \times [0.0, 1.0]$, which reveals more difficulties in adaptation with larger amplitude descriptors. For the Pendulum, the distribution adversary assigns less probability mass to two corner regions, implying that the combination of higher masses and longer pendulums or lower masses and shorter pendulums is easier to predict. Similar phenomena are observed in mass combinations of Acrobat systems. Consistently, the existence of constraint decreases all task distribution entropies to a certain level, which we report in Appendix I. Though such a decrease brings more concentration on some task subsets, AR-MAML

still probably fails to cover other challenging combinations in mode collapse.

Initial Task Distributions’ Influence on Structures: Comparing the top and the bottom of Figure 9, we notice that the uniform and the normal initial distribution results in similar patterns after normalizing flows’ transformations on separate benchmarks. The normal initial distribution can be transformed into smooth ones and captures high-density regions around centroids.

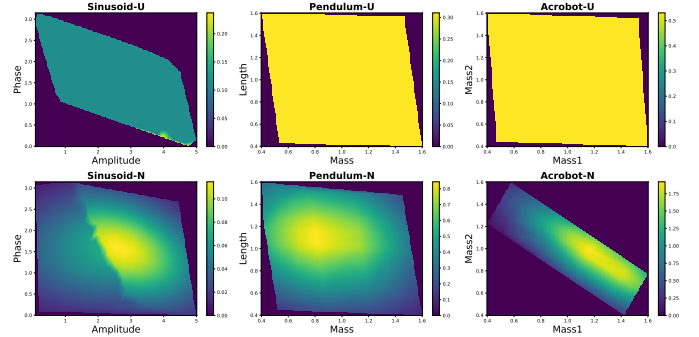


Figure 9: Adversarial Task Probability Distribution. The plots show the adversarial distributions resulting from two different initial distributions: uniform (top row) and normal (bottom row).

5.4 Other Investigations

Here, we conduct additional investigations through the following perspectives.

Impacts of Shift Distribution Constraints: Our studied framework allows the task distribution to shift at a certain level. In Eq. (6), larger λ values tend to cause the generated distribution to collapse into the initial distribution. Consequently, we empirically test the naive and severe adversarial training, e.g., setting $\lambda = \{0.0, 0.1, 0.2\}$ on sinusoid regression. As displayed in Figure 10, the generated distribution with $\lambda = 0.0$ suffers from severe mode collapse, merely covering diagonal regions in the task space. Such a curse is alleviated with increasing λ values. In Figure 11, the meta learner, after heavy distribution shifts, catastrophically fails to generalize well in the initial distribution, illustrating higher adaptation risks in $\lambda = 0.0$.

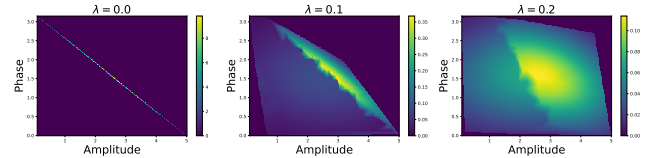


Figure 10: Adversarial Task Probability Distribution on Sinusoid Regression with Various Lagrange Multipliers λ .

Compatibility with Other Meta-learning Methods: Besides the AR-MAML, we also check the effect of adversarially task robust training with other meta-learning methods. Here, AR-CNP in Example 2 is employed in the evaluation. Take the sinusoid regression as an example. Table 2 observes comparable performance between AR-CNP and DR-CNP on the initial task distribution, while results on the adversarial task distribution uncover a significant advantage

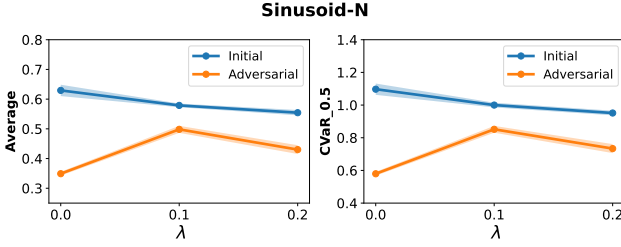


Figure 11: Meta testing MSEs with various lagrange multiplier λ . Reported are testing average and CVaR $_{\alpha}$ MSEs with $\alpha = 0.5$ with standard error in shadow regions.

over others, particularly on robustness metrics, namely CVaR $_{\alpha}$ values.

Table 2: Meta testing MSEs in 5-shot sinusoid regression. With $\alpha = 0.5$, the best results are in pink (the lower, the better).

Method	Average		CVaR	
	Initial	Adversarial	Initial	Adversarial
CNP	0.023 \pm 0.001	0.026 \pm 0.004	0.041 \pm 0.003	0.045 \pm 0.007
TR-CNP	0.048 \pm 0.002	0.050 \pm 0.002	0.076 \pm 0.004	0.079 \pm 0.004
DR-CNP	0.021 \pm 0.001	0.023 \pm 0.002	0.034 \pm 0.003	0.037 \pm 0.003
DRO-CNP	0.023 \pm 0.001	0.025 \pm 0.002	0.039 \pm 0.003	0.041 \pm 0.004
AR-CNP(Ours)	0.019 \pm 0.001	0.018 \pm 0.002	0.033 \pm 0.001	0.029 \pm 0.003

6 CONCLUSIONS

Discussions & Society Impacts. This work develops a game-theoretical approach for generating explicit task distributions in an adversarial way and contributes to theoretical understandings. In extensive scenarios, our approach improves adaptation robustness in constrained distribution shifts and enables the discovery of interpretable task structures in optimization.

Limitations & Future Work. The task distribution in this work relies on the task identifier, which can be inaccessible in some cases, e.g., few-shot classification. Also, the adopted strategy to derive the game solution is approximate, leading to suboptimality in optimization. Hence, future efforts can be made to overcome these limitations and facilitate robust adaptation in applications.

ACKNOWLEDGMENTS

This work is funded by National Natural Science Foundation of China (NSFC) with the Number # 62306326 and # 62495091. And we thank Dong Liang, Yuhang Jiang, Chen Chen, Daming Shi, other anonymous reviewers, and KDD2025 Area Chairs Prof. Yan Liu and Prof. Auroop R Ganguly for suggestions and helpful discussions.

REFERENCES

- [1] Kelsey Allen, Evan Shelhamer, Hanul Shin, and Joshua Tenenbaum. 2019. Infinite mixture prototypes for few-shot learning. In *International Conference on Machine Learning*. PMLR, 232–241.
- [2] Timothée Anne, Jack Wilkinson, and Zhibin Li. 2021. Meta-learning for fast adaptive locomotion with uncertainties in environments and robot dynamics. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 4568–4575.
- [3] Yutong Bai, Xinyang Geng, Karttikeya Mangalam, Amir Bar, Alan Yuille, Trevor Darrell, Jitendra Malik, and Alexei A Efros. 2023. Sequential Modeling Enables Scalable Learning for Large Vision Models. *arXiv:2312.00785 [cs.CV]*
- [4] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*. 41–48.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [6] Haoang Chi, He Li, Wenjing Yang, Feng Liu, Long Lan, Xiaoguang Ren, Tongliang Liu, and Bo Han. 2024. Unveiling Causal Reasoning in Large Language Models: Reality or Mirage?. In *Neural Information Processing Systems*.
- [7] Haoang Chi, Feng Liu, Wenjing Yang, Long Lan, Tongliang Liu, Bo Han, William Cheung, and James Kwok. 2021. TOHAN: A one-step approach towards few-shot hypothesis adaptation. In *Neural Information Processing Systems*, Vol. 34. 20970–20982.
- [8] Haoang Chi, Feng Liu, Wenjing Yang, Long Lan, Tongliang Liu, Bo Han, Gang Niu, Mingyuan Zhou, and Masashi Sugiyama. 2022. Meta Discovery: Learning to Discover Novel Classes given Very Limited Data. In *International Conference on Learning Representations*.
- [9] Liam Collins, Aryan Mokhtari, and Sanjay Shakkottai. 2020. Task-robust model-agnostic meta-learning. *Advances in Neural Information Processing Systems* 33 (2020), 18860–18871.
- [10] Henry Conklin, Bailin Wang, Kenny Smith, and Ivan Titov. 2021. Meta-learning to compositionally generalize. *arXiv preprint arXiv:2106.04252* (2021).
- [11] Corinna Cortes, Yishay Mansour, and Mehryar Mohri. 2010. Learning bounds for importance weighting. *Advances in neural information processing systems* 23 (2010).
- [12] Michael Dennis, Natasha Jaques, Eugene Vinitzky, Alexandre Bayen, Stuart Russell, Andrew Critch, and Sergey Levine. 2020. Emergent complexity and zero-shot transfer via unsupervised environment design. *Advances in neural information processing systems* 33 (2020), 13049–13061.
- [13] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. 2016. RL2: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779* (2016).
- [14] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. 2020. On the convergence theory of gradient-based model-agnostic meta-learning algorithms. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 1082–1092.
- [15] Bernard Faverjon and Pierre Tournassoud. 1987. A local based approach for path planning of manipulators with a high number of degrees of freedom. In *Proceedings. 1987 IEEE international conference on robotics and automation*, Vol. 4. IEEE, 1152–1159.
- [16] Chris Fifty, Ehsan Amid, Zhe Zhao, Tianhe Yu, Rohan Anil, and Chelsea Finn. 2021. Efficiently identifying task groupings for multi-task learning. *Advances in Neural Information Processing Systems* 34 (2021), 27503–27516.
- [17] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*. PMLR, 1126–1135.
- [18] Michael C Fu. 2006. Gradient estimation. *Handbooks in operations research and management science* 13 (2006), 575–616.
- [19] Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. 2018. Conditional neural processes. In *International Conference on Machine Learning*. PMLR, 1704–1713.
- [20] Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye Teh. 2018. Neural processes. *arXiv preprint arXiv:1807.01622* (2018).
- [21] Micah Goldblum, Liam Fowl, and Tom Goldstein. 2019. Robust few-shot learning with adversarially queried meta-learners. (2019).
- [22] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [23] Jonathan Gordon, Wessel P Bruinsma, Andrew YK Foong, James Requeima, Yann Dubois, and Richard E Turner. 2019. Convolutional Conditional Neural Processes. In *International Conference on Learning Representations*.
- [24] Erin Grant, Chelsea Finn, Sergey Levine, Trevor Darrell, and Thomas Griffiths. 2018. Recasting gradient-based meta-learning as hierarchical bayes. *arXiv preprint arXiv:1801.08930* (2018).
- [25] David Ha, Andrew Dai, and Quoc V Le. 2016. Hypernetworks. *arXiv preprint arXiv:1609.09106* (2016).
- [26] Sepp Hochreiter, A Steven Younger, and Peter R Conwell. 2001. Learning to learn using gradient descent. In *International conference on artificial neural networks*. Springer, 87–94.
- [27] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. 2021. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence* 44, 9 (2021), 5149–5169.
- [28] Hongwei Huang, Zhangkai Wu, Wenbin Li, Jing Huo, and Yang Gao. 2021. Local descriptor-based multi-prototype network for few-shot learning. *Pattern Recognition* 116 (2021), 107935.

- [29] Chi Jin, Praneeth Netrapalli, and Michael Jordan. 2020. What is local optimality in nonconvex-nonconcave minimax optimization?. In *International conference on machine learning*. PMLR, 4880–4889.
- [30] Chi Jin, Praneeth Netrapalli, and Michael I Jordan. 2019. Minmax optimization: Stable limit points of gradient descent ascent are locally optimal. *arXiv preprint arXiv:1902.00618* (2019).
- [31] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. 2022. Transformers in vision: A survey. *ACM computing surveys (CSUR)* 54, 10s (2022), 1–41.
- [32] Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. 2018. Attentive Neural Processes. In *International Conference on Learning Representations*.
- [33] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [34] David M Kreps. 1989. Nash equilibrium. In *Game Theory*. Springer, 167–177.
- [35] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. 2016. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236* (2016).
- [36] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521, 7553 (2015), 436–444.
- [37] Kimin Lee, Younggyo Seo, Seunghyun Lee, Honglak Lee, and Jinwoo Shin. 2020. Context-aware dynamics model for generalization in model-based reinforcement learning. In *International Conference on Machine Learning*. PMLR, 5757–5766.
- [38] Seanie Lee, Bruno Andreis, Kenji Kawaguchi, Juho Lee, and Sung Ju Hwang. 2022. Set-based meta-interpolation for few-task meta-learning. *Advances in Neural Information Processing Systems* 35 (2022), 6775–6788.
- [39] Timothée Lesort, Massimo Caccia, and Irina Rish. 2021. Understanding continual learning settings with data distribution drift analysis. *arXiv preprint arXiv:2104.01678* (2021).
- [40] Wenbin Li, Lei Wang, Jinglin Xu, Jing Huo, Yang Gao, and Jiebo Luo. 2019. Revisiting local descriptor based image-to-class measure for few-shot learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 7260–7268.
- [41] Wenbin Li, Lei Wang, Xingxing Zhang, Lei Qi, Jing Huo, Yang Gao, and Jiebo Luo. 2022. Defensive Few-Shot Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 5 (2022), 5649–5667.
- [42] Chenghao Liu, Zhihao Wang, Doyen Sahoo, Yuan Fang, Kun Zhang, and Steven CH Hoi. 2020. Adaptive task sampling for meta-learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII* 16. Springer, 752–769.
- [43] Qi Liu, Tao Liu, Zihao Liu, Yanzhi Wang, Yier Jin, and Wujie Wen. 2018. Security analysis and enhancement of model compressed deep learning systems under adversarial attacks. In *2018 23rd Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 721–726.
- [44] Yixiu Mao, Hongchang Zhang, Chen Chen, Yi Xu, and Xiangyang Ji. 2023. Supported trust region optimization for offline reinforcement learning. In *International Conference on Machine Learning*. PMLR, 23829–23851.
- [45] Yixiu Mao, Hongchang Zhang, Chen Chen, Yi Xu, and Xiangyang Ji. 2024. Supported value regularization for offline reinforcement learning. *Advances in Neural Information Processing Systems* 36 (2024).
- [46] Vladimir G Maz'ya and Tatyana O Shaposhnikova. 1999. *Jacques Hadamard: a universal mathematician*. Number 14. American Mathematical Soc.
- [47] Bhairav Mehta, Manfred Diaz, Florian Golemo, Christopher J Pal, and Liam Paull. 2020. Active domain randomization. In *Conference on Robot Learning*. PMLR, 1162–1176.
- [48] Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2023. Recent advances in natural language processing via large pre-trained language models: A survey. *Comput. Surveys* 56, 2 (2023), 1–40.
- [49] Tsendsuren Munkhdalai and Hong Yu. 2017. Meta networks. In *International Conference on Machine Learning*. PMLR, 2554–2563.
- [50] Shikhar Murty, Tatsunori B Hashimoto, and Christopher D Manning. 2021. Drecta: A general task augmentation strategy for few-shot natural language inference. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 1113–1125.
- [51] Renkun Ni, Micah Goldblum, Amr Sharaf, Kezhi Kong, and Tom Goldstein. 2021. Data augmentation for meta-learning. In *International Conference on Machine Learning*. PMLR, 8152–8161.
- [52] Kevin C Olds. 2015. Global indices for kinematic and force transmission performance in parallel robots. *IEEE Transactions on Robotics* 31, 2 (2015), 494–500.
- [53] Martin J Osborne et al. 2004. *An introduction to game theory*. Vol. 3. Oxford university press New York.
- [54] Jack Parker-Holder, Minqi Jiang, Michael Dennis, Mikayel Samvelyan, Jakob Foerster, Edward Grefenstette, and Tim Rocktäschel. 2022. Evolving curricula with regret-based environment design. In *International Conference on Machine Learning*. PMLR, 17473–17498.
- [55] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019).
- [56] David Pollard. 1984. *Convergence of stochastic processes*. David Pollard.
- [57] Janarthanan Rajendran, Alexander Irpan, and Eric Jang. 2020. Meta-learning requires meta-augmentation. *Advances in Neural Information Processing Systems* 33 (2020), 5705–5715.
- [58] Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. 2019. Meta-learning with implicit gradients. *Advances in neural information processing systems* 32 (2019).
- [59] Davis Rempe, Jonah Philion, Leonidas J Guibas, Sanja Fidler, and Or Litany. 2022. Generating useful accident-prone driving scenarios via a learned traffic prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17305–17315.
- [60] James Requeima, Jonathan Gordon, John Bronskill, Sebastian Nowozin, and Richard E Turner. 2019. Fast and flexible multi-task classification using conditional neural adaptive processes. *Advances in Neural Information Processing Systems* 32 (2019).
- [61] Danilo Rezende and Shakir Mohamed. 2015. Variational inference with normalizing flows. In *International conference on machine learning*. PMLR, 1530–1538.
- [62] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. 2014. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*. PMLR, 1278–1286.
- [63] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. 2018. Meta-Learning with Latent Embedding Optimization. In *International Conference on Learning Representations*.
- [64] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. 2019. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731* (2019).
- [65] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. 2016. Meta-learning with memory-augmented neural networks. In *International conference on machine learning*. PMLR, 1842–1850.
- [66] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust region policy optimization. In *International conference on machine learning*. PMLR, 1889–1897.
- [67] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2016. High-Dimensional Continuous Control Using Generalized Advantage Estimation. In *International Conference on Learning Representations*.
- [68] Jiayi Shen, Xiantong Zhen, Qi Wang, and Marcel Worring. 2023. Episodic Multi-Task Learning with Heterogeneous Neural Processes. *Advances in Neural Information Processing Systems* 36 (2023).
- [69] Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. *Advances in neural information processing systems* 30 (2017).
- [70] Richard S Sutton, Andrew G Barto, et al. 1998. *Introduction to reinforcement learning*. Vol. 135. MIT press Cambridge.
- [71] Shuhan Tan, Kelvin Wong, Shenlong Wang, Sivabalan Manivasagam, Mengye Ren, and Raquel Urtasun. 2021. Scenegen: Learning to generate realistic traffic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 892–901.
- [72] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. 2020. Measuring Robustness to Natural Distribution Shifts in Image Classification. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [73] Sebastian Shenghong Tay, Chuan Sheng Foo, Urano Daisuke, Richalynn Leong, and Bryan Kian Hsiang Low. 2022. Efficient distributionally robust Bayesian optimization with worst-case sensitivity. In *International Conference on Machine Learning*. PMLR, 21180–21204.
- [74] Emanuel Todorov, Tom Erez, and Yuval Tassa. 2012. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 5026–5033.
- [75] Vladimir Vapnik. 1999. *The nature of statistical learning theory*. Springer science & business media.
- [76] Heinrich Von Stackelberg and Stackelberg Heinrich Von. 1952. *The theory of the market economy*. Oxford University Press.
- [77] Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. 2020. On the global optimality of model-agnostic meta-learning. In *International conference on machine learning*. PMLR, 9837–9846.
- [78] Qi Wang, Marco Federici, and Herke van Hoof. 2023. Bridge the Inference Gaps of Neural Processes via Expectation Maximization. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=A7v2DqLjZdq>
- [79] Qi Wang, Yanghe Feng, Jincui Huang, Yiqin Lv, Zheng Xie, and Xiaoshan Gao. 2023. Large-scale generative simulation artificial intelligence: The next hotspot. *The Innovation* 4, 6 (2023).
- [80] Qi Wang, Yiqin Lv, Yanghe Feng, Zheng Xie, and Jincui Huang. 2023. A Simple Yet Effective Strategy to Robustify the Meta Learning Paradigm. *Advances in Neural Information Processing Systems* 36 (2023).
- [81] Qi Wang and Herke Van Hoof. 2020. Doubly stochastic variational inference for neural processes with hierarchical latent variables. In *International Conference*

- on *Machine Learning*. PMLR, 10018–10028.
- [82] Qi Wang and Herke van Hoof. 2022. Learning expressive meta-representations with mixture of expert neural processes. In *Advances in neural information processing systems*.
 - [83] Qi Wang and Herke Van Hoof. 2022. Model-based meta reinforcement learning using graph structured surrogate models and amortized policy search. In *International Conference on Machine Learning*. PMLR, 23055–23077.
 - [84] Wolfram Wiesemann, Daniel Kuhn, and Melvyn Sim. 2014. Distributionally robust convex optimization. *Operations Research* 62, 6 (2014), 1358–1376.
 - [85] Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning* 8, 3 (1992), 229–256.
 - [86] Yichen Wu, Long-Kai Huang, and Ying Wei. 2022. Adversarial task up-sampling for meta-learning. *Advances in Neural Information Processing Systems* 35 (2022), 31102–31115.
 - [87] Zehao Xiao, Jiayi Shen, Mohammad Mahdi Derakhshani, Shengcai Liao, and Cees GM Snoek. 2024. Any-Shift Prompting for Generalization over Distributions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13849–13860.
 - [88] Zehao Xiao, Xiantong Zhen, Shengcai Liao, and Cees GM Snoek. 2023. Energy-based test sample adaptation for domain generalization. *arXiv preprint arXiv:2302.11215* (2023).
 - [89] Zehao Xiao, Xiantong Zhen, Ling Shao, and Cees GM Snoek. 2022. Learning to generalize across domains on single test samples. *arXiv preprint arXiv:2202.08045* (2022).
 - [90] Lihe Yang, Wei Zhuo, Lei Qi, Yinghuan Shi, and Yang Gao. 2021. Mining latent classes for few-shot segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*. 8721–8730.
 - [91] Huaxiu Yao, Long-Kai Huang, Linjun Zhang, Ying Wei, Li Tian, James Zou, Junzhou Huang, et al. 2021. Improving generalization in meta-learning via task augmentation. In *International Conference on Machine Learning*. PMLR, 11887–11897.
 - [92] Huaxiu Yao, Yu Wang, Ying Wei, Peilin Zhao, Mehrdad Mahdavi, Defu Lian, and Chelsea Finn. 2021. Meta-learning with an adaptive task scheduler. *Advances in Neural Information Processing Systems* 34 (2021), 7497–7509.
 - [93] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. 2020. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*. PMLR, 1094–1100.
 - [94] Jesse Zhang, Brian Cheung, Chelsea Finn, Sergey Levine, and Dinesh Jayaraman. 2020. Cautious adaptation for reinforcement learning in safety-critical settings. In *International Conference on Machine Learning*. PMLR, 11055–11065.
 - [95] Marvin Zhang, Henrik Marklund, Nikita Dhawan, Abhishek Gupta, Sergey Levine, and Chelsea Finn. 2021. Adaptive risk minimization: Learning to adapt to domain shift. *Advances in Neural Information Processing Systems* 34 (2021), 23664–23678.

CONTENTS

Abstract	1
1 Introduction	1
2 Literature Review	2
3 Preliminaries	2
3.1 Problem Statement	2
3.2 Two-Player Stackelberg Game	3
4 Task Robust Meta Learning under Distribution Shift Constraints	3
4.1 Generate Task Distribution within A Game-Theoretic Framework	3
4.2 Solution Concept & Explanations	4
4.3 Strategies for Finding Equilibrium	5
4.4 Theoretical Analysis	5
5 Experiments	6
5.1 Benchmarks	6
5.2 Empirical Result Analysis	7
5.3 Task Structure Analysis	8
5.4 Other Investigations	8
6 Conclusions	9
Acknowledgments	9
References	9
Contents	12
A Quick Guideline for This Work	13
A.1 Pseudo Code of the AR-MAML & AR-CNP	13
A.2 Players' Order in Decision-making	14
A.3 Benefits of the Explicit Generative Modeling	14
A.4 Set-up of Base Distributions	14
A.5 Related Work Summary	15
B Robust Fast Adaptation Strategies in the Task Space	15
C Stochastic Gradient Estimates	16
D Theoretical Understanding of Generating Task Distribution	16
E Equilibrium & Convergence Guarantee	17
E.1 Existence of Stackelberg Equilibrium	17
E.2 Convergence Guarantee	17
E.3 Convergence Properties	19
F Generalization Bound	19
F.1 Importance Weighted Generalization Bound	19
F.2 Formal Adversarial Generalization Bounds	20
G Explicit Generative Task Distributions	20
H Evolution of Entropies in Task Distributions	20
I Experimental Set-up & Implementation Details	21
I.1 Meta Learning Benchmarks	21
I.2 Modules in Pytorch	22
I.3 Neural Architectures & Optimization	25
I.4 Distribution Adversary Implementations	25
J Additional Experimental Results	25
K Platforms & Computational Tools	25

A QUICK GUIDELINE FOR THIS WORK

Here, we include some guidelines for this work. Our focus is to explicitly generate the task distribution with adversarial training. The use of normalizing flows enables task structure discovery under the risk minimization principle. The theoretical understanding is from the Stackelberg game, together with some analyses. Allowing the task distribution to shift at an acceptable level is a promising topic in meta-learning and can help improve robustness in the presence of worse cases. The following further complements these points.

A.1 Pseudo Code of the AR-MAML & AR-CNP

This section lists the pseudo-code for implementing AR-MAML and AR-CNPs in practice. The difference between AR-MAML and AR-CNP lies in the inner loop, as CNP does not require additional gradient updates in evaluating the meta-learner.

Algorithm 1: Adversarially Task Robust MAML

Input : Initial task distribution $p_0(\tau)$; Task batch size K ; Update frequency u ; Learning rates of players: $\{\gamma_{1,1}, \gamma_{1,2}, \gamma_2\}$.

Output : Meta-trained model parameter θ .

Randomly initialize the model parameter θ ;

Set the initial iteration number $t \leftarrow 0$;

while *not converged* **do**

 Sample a batch of tasks $\{\tau_i\}_{i=1}^K \sim p_\phi(\tau)$;

 // **Leader Inner Gradient Descent**

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

 Compute the task-specific gradient: $\nabla_\theta \mathcal{L}(\mathcal{D}_{\tau_i}^S; \theta)$;

 Perform gradient updates as fast adaptation:

$\theta_i \leftarrow \theta - \gamma_{1,1} \nabla_\theta \mathcal{L}(\mathcal{D}_{\tau_i}^S; \theta)$;

end

 Update the iteration number: $t \leftarrow t + 1$;

 // **Leader Outer Gradient Descent**

 Perform meta initialization updates:

$\theta \leftarrow \theta - \frac{\gamma_{1,2}}{K} \sum_{i=1}^K \nabla_\theta \mathcal{L}(\mathcal{D}_{\tau_i}^Q; \theta_i)$;

if $t \% u = 0$ **then**

 // **Follower Gradient Ascent**

 Compute the baseline for the follower:

$\mathcal{V} \approx \frac{1}{K} \sum_{k=1}^K \mathcal{L}(\mathcal{D}_{\tau_k}^Q; \theta)$;

 Perform task distribution $p_\phi(\tau)$ updates:

$\phi \leftarrow \phi + \frac{\gamma_2}{K} \sum_{k=1}^K [\mathcal{L}(\mathcal{D}_{\tau_k}^Q; \theta) - \mathcal{V}] \nabla_\phi \ln p_\phi(\tau_k) + \frac{\lambda \gamma_2}{K} \sum_{k=1}^K \nabla_\phi \ln p_\phi(\tau_k^{-M})$.

end

end

Algorithm 2: Adversarially Task Robust CNP

Input : Initial task distribution $p_0(\tau)$; Task batch size K ; Update frequency u ; Learning rates of players: $\{\gamma_1, \gamma_2\}$.

Output: Meta-trained model parameter $\theta = [\theta_1, \theta_2]$.

Randomly initialize the model parameter θ ;

Set the initial iteration number $t \leftarrow 0$;

while not converged **do**

 Sample a batch of tasks $\{\tau_i\}_{i=1}^K \sim p_\phi(\tau)$;

 // **Leader Gradient Descent**

 Compute the task-specific gradient for the leader: $\nabla_\theta \mathcal{L}(\mathcal{D}_{\tau_i}^Q; \mathcal{D}_{\tau_i}^S, \theta)$;

 Perform meta initialization updates:

$\theta \leftarrow \theta - \frac{\gamma_1}{K} \sum_{i=1}^K \nabla_\theta \mathcal{L}(\mathcal{D}_{\tau_i}^Q; \mathcal{D}_{\tau_i}^S, \theta_i)$;

 Update the iteration number: $t \leftarrow t + 1$;

if $t \% u = 0$ **then**

 // **Follower Gradient Ascent**

 Compute the baseline for the follower:

$\mathcal{V} \approx \frac{1}{K} \sum_{k=1}^K \mathcal{L}(\mathcal{D}_{\tau_k}^Q; \theta)$;

 Perform task distribution $p_\phi(\tau)$ updates:

$\phi \leftarrow \phi + \frac{\gamma_2}{K} \sum_{k=1}^K \left[[\mathcal{L}(\mathcal{D}_{\tau_k}^Q; \theta) - \mathcal{V}] \nabla_\phi \ln p_\phi(\tau_k) + \lambda \nabla_\phi \ln p_\phi(\tau_k^{-M}) \right]$.

end

end

A.2 Players' Order in Decision-making

REMARK 3 (OPTIMIZATION ORDER AND SOLUTIONS). *With the risk function $\mathcal{J}(\theta, \phi)$ convex w.r.t. $\forall \theta \in \Theta$ and concave w.r.t. $\forall \phi \in \Phi$, swapping the order of two players $\{\mathcal{P}_1, \mathcal{P}_2\}$ results in the same equilibrium:*

$$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) = \max_{\phi \in \Phi} \min_{\theta \in \Theta} \mathcal{J}(\theta, \phi).$$

When the convex-concave assumption does not hold, the order generally results in completely different solutions:

$$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) \neq \max_{\phi \in \Phi} \min_{\theta \in \Theta} \mathcal{J}(\theta, \phi).$$

The **Remark 3** indicates the game is asymmetric w.r.t. the players' decision-making order. This work abstracts the decision-making process for the distribution adversary and the meta-learner in a sequential game. Hence, the order inevitably influences the solution concept. Actually, we take MAML as an example to implement the adversarially task robust meta-learning framework. As noted in Algorithm 1, updating the distribution adversary's parameter requires the evaluation of task-specific model parameters, which are not available in the initialization. The traditional game theory no longer applies to our studied case, and simple counter-examples such as rock/paper/scissors show no equilibrium in practice. Hence, we pre-determine the order of decision-makers for implementation, as reported in the main paper.

A.3 Benefits of the Explicit Generative Modeling

The use of flow-based generative models is beneficial, as the density distribution function is accessible in numerical analysis, allowing one to understand the behavior of the distribution adversary. Our training style is adversarial, and the equilibrium is approximately obtained in a gradient update way. As the stationary point works for both the distribution adversary and the meta player, the resulting meta player is the direct solution to robust fast adaptation. All of these have been validated from learned probability task densities and entropies. We treat this as a side product of the method.

A.4 Set-up of Base Distributions

This work does not create new synthetic tasks and only considers the regression and continuous control cases. Our suggestion of a comprehensive base distribution is designed to cover a wide range of scenarios, making it a practical solution in most default setups. Through the optimization of the distribution adversary, the meta-learner converges to a distribution that focuses on more challenging scenarios and lowers the importance of trivial cases. Even though enlarging the scope of the scenarios can result in additional computational costs, generative modeling of the task distribution captures more realistic feedback from adaptation performance. In other words, our study is based on the hypothesis that when the task space is vast enough, the subpopulation shift is allowed under a certain constraint.

A.5 Related Work Summary

As far as we know, robust fast adaptation is an underexplored research issue in the field. In principle, we include the latest SOTA [9, 80] and typical SOTA [64] methods to compare in this work as illustrated in Table 3. Note that in [80], the idea of increasing the robustness is to introduce the task selector for risky tasks in implementations. Though there are some curriculum methods [12, 54, 92] to reschedule the task sampling probability along the learning, our setup focuses on (i) the semantics in the task identifiers to explicitly uncover the task structure from a Stackelberg game and (ii) the robust solution search under a distribution shift constraint. That is, ours requires capturing the differentiation information from the task identifiers of tasks, such as masses and lengths in the pendulum system identification. Ours is agnostic to meta-learning methods and avoids task distribution mode collapse, which might happen in adversarial training without constraints. These configuration and optimization differences drive us to include the risk minimization principle in the comparisons.

Table 3: A Summary of Used Methods from Diverse Optimization Principles and Corresponding Properties in Meta-Learning. These are implemented within the empirical/expected risk minimization (ERM), the distributionally robust risk minimization (DRM), and the adversarially task robust risk minimization (ARM). Here, MAML works as the backbone method to illustrate the difference.

Method	Task Distribution	Optimization Objective	Adaptation Robustness
MAML [17]	Fixed	$\min_{\theta \in \Theta} \mathbb{E}_{\tau \sim p(\tau)} [\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)]$	--
DR-MAML [80]	Fixed	$\min_{\theta \in \Theta} \mathbb{E}_{p_\alpha(\tau; \theta)} [\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)]$	Tail Risk
TR-MAML [9]	Fixed	$\min_{\theta \in \Theta} \max_{\tau \in \mathcal{T}} \mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)$	Worst Risk
DRO-MAML [64]	Fixed	$\min_{\theta \in \Theta} \max_{g \in \mathcal{G}} \mathbb{E}_{p_g(\tau)} [\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)]$	Uncertainty Set
AR-MAML (Ours)	Explicit Generation	$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathbb{E}_{p_\phi(\tau)} [\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)] + \lambda \mathbb{E}_{p_0(\tau)} [\ln p_\phi(\tau)]$	Constrained Shift

B ROBUST FAST ADAPTATION STRATEGIES IN THE TASK SPACE

In this section, we recap recent typical strategies for overcoming the negative effect of task distribution shifts.

Group Distributionally Robust Optimization (GDRO) for Meta Learning. This method [64] considers combating the distribution shifts via the group of worst-case optimization. In detail, the general meta training objective can be written as follows:

$$\min_{\theta \in \Theta} \left\{ \mathcal{R}(\theta) := \sup_{g \in \mathcal{G}} \mathbb{E}_{p_g(\tau)} [\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)] \right\}, \quad (19)$$

where \mathcal{G} denotes a collection of uncertainty sets at a group level. When tasks of interest are finite, g induces a probability measure $p_g(\tau)$ over a subset, which we call the group in the background of meta-learning scenarios. The optimization step inside the bracket of Eq. (19) describes the selection of the worst group in fast adaptation performance. By minimizing the worst group risk, the fast adaptation robustness can be enhanced when confronted with the task distribution shift.

Tail Risk Minimization (CVaR_α) for Meta Learning. This method [80] is from the risk-averse perspective, and the tail risk measured CVaR_α by is incorporated in stochastic programming. The induced meta training objective is written as follows:

$$\min_{\theta \in \Theta} \mathbb{E}_{p_\alpha(\tau; \theta)} [\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)], \quad (20)$$

where $p_\alpha(\tau; \theta)$ represents the normalized distribution over the tail task in fast adaptation performance, which implicitly relies on the meta learner θ .

Also note that in [80], the above equation can be equivalently expressed as an importance weighted one:

$$\min_{\theta \in \Theta} \mathbb{E}_{p(\tau)} \left[\frac{p_\alpha(\tau; \theta)}{p(\tau)} \mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \right] \approx \frac{1}{\mathcal{B}} \sum_{b=1}^{\mathcal{B}} \omega(\tau_b; \theta) \mathcal{L}(\mathcal{D}_{\tau_b}^Q, \mathcal{D}_{\tau_b}^S; \theta), \quad (21)$$

where $\omega(\tau_b; \theta)$ is the ratio of the normalized tail task probability and the original task probability. In this way, Eq. (21) can be connected to the group distributionally robust optimization method, as the worst group is estimated from the tail risk and reflected in the importance ratio.

C STOCHASTIC GRADIENT ESTIMATES

Unlike most adversarial generative models whose likelihood is intractable, our design enables the exact computation of task likelihoods, and the stochastic gradient estimate can be achieved as follows:

$$\begin{aligned}
 \nabla_{\phi} \mathcal{J}(\theta, \phi) &= \nabla_{\phi} \int p_{\phi}(\tau) \mathcal{L}(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) d\tau + \lambda \nabla_{\phi} \int p_0(\tau) \ln p_{\phi}(\tau) d\tau \\
 &= \int p_{\phi}(\tau) \left[\nabla_{\phi} \ln p_{\phi}(\tau) \mathcal{L}(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) \right] d\tau + \lambda \nabla_{\phi} \int p_0(\tau) \nabla_{\phi} \ln p_{\phi}(\tau) d\tau \\
 &\approx \frac{1}{K} \sum_{k=1}^K \left[\mathcal{L}(D_{\tau_k}^Q, D_{\tau_k}^S; \theta) \nabla_{\phi} \ln p_{\phi}(\tau_k) + \lambda \nabla_{\phi} \ln p_{\phi}(\tau_k^{-M}) \right]
 \end{aligned} \tag{22}$$

With the previously mentioned invertible mappings $\mathcal{G} = \{g_i\}_{i=1}^M$, where $g_i : \mathcal{T} \rightarrow \mathcal{T} \subseteq \mathbb{R}^d$, we describe the transformation of tasks as $\tau^{-M} \mapsto \dots \mapsto \tau^{-1} \mapsto \tau$. This indicates that the sampled task particle τ_k^{-M} satisfies the equation:

$$\tau_k = g_M \circ \dots \circ g_2 \circ g_1(\tau_k^{-M}) = \text{NN}_{\phi}(\tau_k^{-M}), \quad \text{with } \tau_k \sim p_{\phi}(\tau), \quad k \in \{1, \dots, K\}. \tag{23}$$

D THEORETICAL UNDERSTANDING OF GENERATING TASK DISTRIBUTION

Unlike previous curriculum learning or adversarial training in the task space, this work places a distribution shift constraint over the task space. Also, our setup tends to create a subpopulation shift under a trust region.

Subpopulation Shift Constraint as the Regularization. The generated task distribution corresponds to the best response in deteriorating the meta-learner’s performance under a tolerant level of the task distribution shift. Note that the proposed optimization objective $\mathcal{J}(\theta, \phi)$ includes a KL divergence term *w.r.t.* the generated task distribution.

$$\begin{aligned}
 \min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) &:= \mathbb{E}_{p_{\phi}(\tau)} \left[\mathcal{L}(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) \right] \\
 \text{s.t. } D_{KL}[p_0(\tau) \parallel p_{\phi}(\tau)] &\leq \delta
 \end{aligned} \tag{24}$$

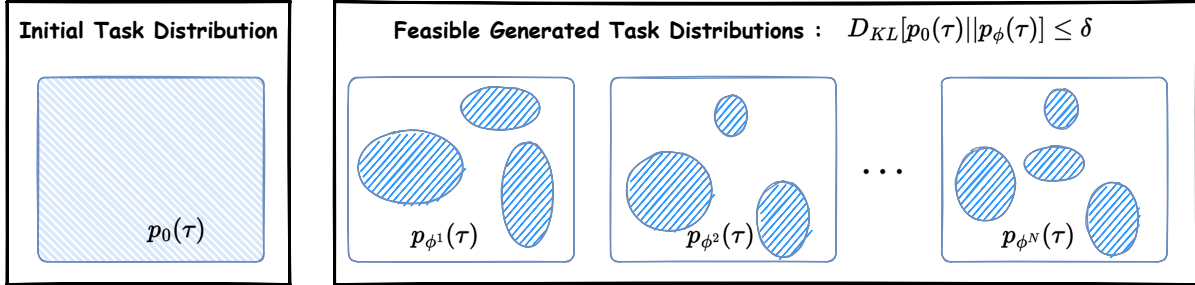


Figure 12: Intuition of distribution shifts and Optimization Steps. Left is the initial task distribution, while the generated task distributions within the constraint of the distribution shift are illustrated on the right.

Tolerant Region in distribution shifts. Next, we interpret the above-induced optimization objective. As exhibited in Eq. (24), the condition constrains the generative task distribution $p_{\phi}(\tau)$ within a neighborhood of the initial task distribution $p_0(\tau)$. This can be viewed as the *tolerant region of task distribution shifts*, which means larger δ allows more severe distribution shifts concerning the initial task distribution. As for the role of the distribution adversary, it attempts to seek the task distribution that can handle the strongest task distribution shift under the allowed region, as displayed in Figure 12.

As mentioned in the main paper, the equivalent unconstrained version can be

$$\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) := \underbrace{\mathbb{E}_{p_{\phi}(\tau)} \left[\mathcal{L}(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) \right]}_{\text{Adversarial Meta Learning}} + \lambda \underbrace{\mathbb{E}_{p_0(\tau)} \left[\ln p_{\phi}(\tau) \right]}_{\text{Distribution Cloning}}, \tag{25}$$

where the second penalty is to prevent the generated task distribution from uncontrollably diverging from the initial one. Hence, it encourages the exploration of crucial tasks in the broader scope while avoiding mode collapse from adversarial training. The larger the Lagrange multiplier λ , the weaker the distribution shift in a generation. In this work, we expect that the tolerant distribution shift can cover a larger scope of shifted task distributions. Hence, λ is configured to be a small value as default.

Best Response in Constrained Optimization. We can also explain the optimization steps in the presence of the Stackelberg game. Here, we can equivalently rewrite the entangled steps to solve the game as follows:

$$\min_{\theta \in \Theta} \mathcal{J}(\theta, \phi_*(\theta)) \quad \text{s.t. } \phi_*(\theta) = \arg \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) \text{ and } \mathcal{D}_{KL}[p_0(\tau) || p_\phi(\tau)] \leq \delta \quad (26a)$$

$$\max_{\phi \in \Phi} \mathcal{J}(\theta, \phi) \text{ and } \mathcal{D}_{KL}[p_0(\tau) || p_\phi(\tau)] \leq \delta, \quad (26b)$$

where Eq. (26.a) indicates the decision-making of the meta learner, while Eq. (26.b) characterizes the decision-making of the distribution adversary. Particularly, Eq. (26.b) is a constrained sub-optimization problem.

E EQUILIBRIUM & CONVERGENCE GUARANTEE

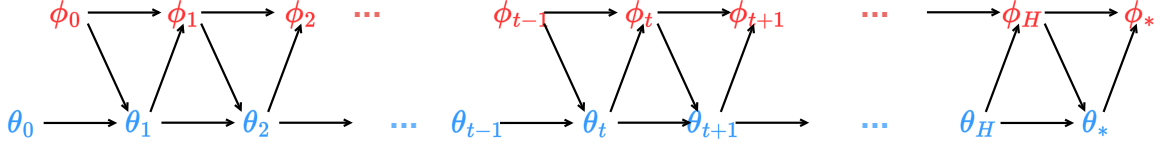


Figure 13: Diagram of Stochastic Alternating Gradient Descent Ascent. In solving adversarially task robust meta-learning, the updates of separate players' parameters are performed repetitively and can be viewed as a bi-level optimization.

E.1 Existence of Stackelberg Equilibrium

DEFINITION 3 (GLOBAL MINIMAX POINT). Once the constructed meta-learning Stackelberg game in Eq. (5) is solved with the optimal solution as the Stackelberg equilibrium $\{\theta_*, \phi_*\}$, we can naturally obtain the optimal expected risk value \mathcal{R}_* and the inequalities $\forall(\theta, \phi) \in \Theta \times \Phi$ as follows:

$$\begin{aligned} \mathcal{J}(\theta_*, \phi) &\leq \mathcal{J}(\theta_*, \phi_*) \leq \max_{\phi \in \Phi} \mathcal{J}(\theta, \hat{\phi}) \\ \mathcal{R}_* = \mathcal{J}(\theta_*, \phi_*) &= \mathbb{E}_{p_{\phi_*}(\tau)} \left[\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta_*) \right] + \lambda \mathbb{E}_{p_0(\tau)} \left[\ln p_{\phi_*}(\tau) \right]. \end{aligned} \quad (27)$$

Note that for the nonconvex-nonconcave min-max optimization problem, there is no general guarantee for the existence of saddle points. The solution concept in the **Definition 3** requires no convexity w.r.t. the optimization objective and provides a weaker equilibrium, also referred to as the global minimax point [29].

With the **Assumption 1** and the **Definition 3**, the global Stackelberg equilibrium always exists for the proposed min-max optimization $\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{J}(\theta, \phi)$. However, its exact search is NP-hard, and the stochastic optimization practically leads to the local Stackelberg equilibrium in the **Definition 2**.

E.2 Convergence Guarantee

Assumption 1 (Lipschitz Smoothness and Compactness) The adversarially task robust meta-learning optimization objective $\mathcal{J}(\theta, \phi)$ is assumed to satisfy

- (1) $\mathcal{J}(\theta, \phi)$ with $\forall[\theta, \phi] \in \Theta \times \Phi$ belongs to the class of twice differentiable functions \mathbb{C}^2 .
- (2) The norm of block terms inside Hessian matrices $\nabla^2 \mathcal{J}(\theta, \phi)$ is bounded, meaning that $\forall[\theta, \phi] \in \Theta \times \Phi$:

$$\sup\{\|\nabla_{\theta, \theta}^2 \mathcal{J}\|, \|\nabla_{\theta, \phi}^2 \mathcal{J}\|, \|\nabla_{\phi, \phi}^2 \mathcal{J}\|\} \leq L_{\max}.$$

- (3) The parameter spaces $\Theta \subseteq \mathbb{R}^{d_1}$ and $\Phi \subseteq \mathbb{R}^{d_2}$ are compact with d_1 and d_2 respectively dimensions of model parameters for two players.

Theorem 1 (Convergence Guarantee) Suppose the **Assumption 1** and the condition of the (local) Stackelberg equilibrium $\Delta(\mathbf{A}, \mathbf{B}, \mathbf{C}, \gamma_1, \gamma_2) := \max \left\{ (1 - \gamma_1 \sigma_{\min})^2 (1 + \gamma_2^2 L_{\max}^2), \left| \gamma_1^2 - 2\gamma_1 \gamma_2 + \gamma_1^2 \gamma_2^2 L_{\max}^2 L_{\max}^2 + (1 + \gamma_2 L_{\max})^2 - 2\gamma_1 \gamma_2^2 \sigma_{\min}(\mathbf{B}^T \mathbf{B} \mathbf{C}) \right| \right\} < \frac{1}{2}$ are satisfied, where $\sigma_{\min}(\cdot)$ denotes the smallest eigenvalues of the corresponding matrix. Then we can have the following statements:

- (1) The resulting iterated parameters $\{\dots \mapsto [\theta_t, \phi_t]^T \mapsto [\theta_{t+1}, \phi_{t+1}]^T \mapsto \dots\}$ are Cauchy sequences.
- (2) The optimization can guarantee at least the linear convergence to the local Stackelberg equilibrium with the rate $\sqrt{\Delta}$.

Proof of Theorem 1:

Remember the learning dynamics as follows,

$$\theta_{t+1} \leftarrow \theta_t - \gamma_1 \nabla_{\theta} \mathcal{J}(\theta_t, \phi_t) \quad (28a)$$

$$\phi_{t+1} \leftarrow \phi_t + \gamma_2 \nabla_{\phi} \mathcal{J}(\theta_{t+1}, \phi_t). \quad (28b)$$

when using the alternating GDA for solving the adaptively robust meta-learning. Figure 13 illustrates the iteration rules and steps.

Let $[\theta_*, \phi_*]^T$ be the obtained (local) Stackelberg equilibrium, we denote the difference between the updated model parameter and the equilibrium by $[\hat{\theta}_t, \hat{\phi}_t]^T = [\theta_t - \theta_*, \phi_t - \phi_*]^T$. As the utility function $\mathcal{J}(\theta, \phi)$ is Lipschitz smooth, we can perform linearization of $\nabla_{\theta} \mathcal{J}(\theta, \phi)$ and $\nabla_{\phi} \mathcal{J}(\theta, \phi)$ round the resulting stationary point $[\theta_*, \phi_*]^T$ respectively as follows.

$$\nabla_{\theta} \mathcal{J}(\theta_t, \phi_t) \approx \nabla_{\theta\theta}^2 \mathcal{J}(\theta_*, \phi_*)(\theta_t - \theta_*) + \nabla_{\theta\phi}^2 \mathcal{J}(\theta_*, \phi_*)(\phi_t - \phi_*) \quad (29a)$$

$$\nabla_{\phi} \mathcal{J}(\theta_{t+1}, \phi_t) \approx \nabla_{\phi\theta}^2 \mathcal{J}(\theta_*, \phi_*)(\theta_{t+1} - \theta_*) + \nabla_{\phi\phi}^2 \mathcal{J}(\theta_*, \phi_*)(\phi_t - \phi_*). \quad (29b)$$

Note that $\nabla_{\phi} \mathcal{J}(\theta_{t+1}, \phi_t)$ in Eq. (29) can be further expressed as

$$\begin{aligned} \nabla_{\phi} \mathcal{J}(\theta_{t+1}, \phi_t) \approx \nabla_{\phi\theta}^2 \mathcal{J}(\theta_*, \phi_*) \left[\theta_t - \theta_* - \gamma_1 \nabla_{\theta\theta}^2 \mathcal{J}(\theta_*, \phi_*)(\theta_t - \theta_*) - \gamma_1 \nabla_{\theta\phi}^2 \mathcal{J}(\theta_*, \phi_*)(\phi_t - \phi_*) \right] \\ + \nabla_{\phi\phi}^2 \mathcal{J}(\theta_*, \phi_*)(\phi_t - \phi_*). \end{aligned} \quad (30)$$

with the help of Eq. (28). Same as that in the main paper, we write the block terms inside the Hessian matrix $\mathbf{H}_* := \nabla^2 \mathcal{J}(\theta_*, \phi_*)$ around

$$[\theta_*, \phi_*]^T \text{ as } \begin{bmatrix} \nabla_{\theta\theta}^2 \mathcal{J} & \nabla_{\theta\phi}^2 \mathcal{J} \\ \nabla_{\phi\theta}^2 \mathcal{J} & \nabla_{\phi\phi}^2 \mathcal{J} \end{bmatrix} \Big|_{[\theta_*, \phi_*]^T} := \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix}.$$

Then we can naturally derive the new form of learning dynamics as follows:

$$\hat{\theta}_{t+1} = \hat{\theta}_t - \gamma_1 \mathbf{A} \hat{\theta}_t - \gamma_1 \mathbf{B} \hat{\phi}_t \quad (31a)$$

$$\hat{\phi}_{t+1} = \hat{\phi}_t + \gamma_2 \mathbf{B}^T \hat{\theta}_t - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{A} \hat{\theta}_t - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B} \hat{\phi}_t + \gamma_2 \mathbf{C} \hat{\phi}_t. \quad (31b)$$

Further, for the sake of simplicity, we rewrite the above mentioned matrices as $\mathbf{P} = \mathbf{I} - \gamma_1 \mathbf{A}$ and $\mathbf{Q} = \mathbf{I} + \gamma_2 \mathbf{C}$. Equivalently, Eq. (31) can be written in the matrix form:

$$\begin{bmatrix} \hat{\theta}_{t+1} \\ \hat{\phi}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{P} & -\gamma_1 \mathbf{B} \\ \gamma_2 \mathbf{B}^T - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{A} & \mathbf{Q} - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B} \end{bmatrix} \cdot \begin{bmatrix} \hat{\theta}_t \\ \hat{\phi}_t \end{bmatrix}. \quad (32)$$

Furthermore, we consider the expression of the updated parameters' norm:

$$\|\hat{\theta}_{t+1}\|_2^2 + \|\hat{\phi}_{t+1}\|_2^2 = \|\mathbf{P} \hat{\theta}_t - \gamma_1 \mathbf{B} \hat{\phi}_t\|_2^2 + \|(\gamma_2 \mathbf{B}^T - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{A}) \hat{\theta}_t + (\mathbf{Q} - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B}) \hat{\phi}_t\|_2^2 \quad (33a)$$

$$= \|\mathbf{P} \hat{\theta}_t - \gamma_1 \mathbf{B} \hat{\phi}_t\|_2^2 + \|\gamma_2 \mathbf{B}^T \mathbf{P} \hat{\theta}_t + (\mathbf{Q} - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B}) \hat{\phi}_t\|_2^2 \quad (33b)$$

$$\leq 2 \left[\|\mathbf{P} \hat{\theta}_t\|_2^2 + \gamma_1^2 \|\mathbf{B} \hat{\phi}_t\|_2^2 \right] + 2 \left[\gamma_2^2 \|\mathbf{B}^T \mathbf{P} \hat{\theta}_t\|_2^2 + \|(\mathbf{Q} - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B}) \hat{\phi}_t\|_2^2 \right] \quad (33c)$$

$$= 2 \left[\|\mathbf{P} \hat{\theta}_t\|_2^2 + \gamma_2^2 \|\mathbf{B}^T \mathbf{P} \hat{\theta}_t\|_2^2 \right] + 2 \left[\gamma_1^2 \|\mathbf{B} \hat{\phi}_t\|_2^2 + \|(\mathbf{Q} - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B}) \hat{\phi}_t\|_2^2 \right], \quad (33d)$$

where the last inequality makes use of the Cauchy-Schwarz inequality trick $\|\mathbf{a} + \mathbf{b}\|_2^2 \leq 2(\|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2)$, $\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^d$.

For terms concerning θ , we can have the following evidence:

$$\|\mathbf{P} \hat{\theta}_t\|_2^2 + \gamma_2^2 \|\mathbf{B}^T \mathbf{P} \hat{\theta}_t\|_2^2 = \hat{\theta}_t^T \left(\mathbf{P}^T \mathbf{P} + \gamma_2^2 \mathbf{P}^T \mathbf{B} \mathbf{B}^T \mathbf{P} \right) \hat{\theta}_t \quad (34a)$$

$$\leq (1 + \gamma_2^2 L_{\max}^2) \hat{\theta}_t^T \mathbf{P}^T \mathbf{P} \hat{\theta}_t \quad (34b)$$

$$\leq (1 + \gamma_2^2 L_{\max}^2) (1 - \gamma_1 \sigma_{\min}(\mathbf{A})) \|\hat{\theta}_t\|_2^2, \quad (34c)$$

where the last two inequalities make use of the assumptions and tricks that (i) the boundness assumption $\mathbf{B} \leq L_{\max}$ and (ii) the trait of the symmetric matrix $\|\mathbf{P}\|_2 = \|\mathbf{I} - \gamma_1 \mathbf{A}\|_2 \leq 1 - \gamma_1 \sigma_{\min}(\mathbf{A})$.

For terms concerning ϕ , we can have the following evidence:

$$\gamma_1^2 \|\mathbf{B} \hat{\phi}_t\|_2^2 + \|(\mathbf{Q} - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B}) \hat{\phi}_t\|_2^2 = \hat{\phi}_t^T \left(\gamma_1^2 \mathbf{B}^T \mathbf{B} + (\mathbf{Q} - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B})^T (\mathbf{Q} - \gamma_1 \gamma_2 \mathbf{B}^T \mathbf{B}) \right) \hat{\phi}_t \quad (35a)$$

$$= \hat{\phi}_t^T \mathbf{B}^T \left(\gamma_1^2 \mathbf{I} + \gamma_1^2 \gamma_2^2 \mathbf{B} \mathbf{B}^T \right) \mathbf{B} \hat{\phi}_t + \hat{\phi}_t^T \mathbf{Q}^T \mathbf{Q} \hat{\phi}_t - 2 \gamma_1 \gamma_2 \hat{\phi}_t^T (\mathbf{B}^T \mathbf{B} \mathbf{Q}) \hat{\phi}_t \quad (35b)$$

$$= \hat{\phi}_t^T \mathbf{B}^T \left(\gamma_1^2 \mathbf{I} + \gamma_1^2 \gamma_2^2 \mathbf{B} \mathbf{B}^T \right) \mathbf{B} \hat{\phi}_t + \hat{\phi}_t^T \mathbf{Q}^T \mathbf{Q} \hat{\phi}_t - 2 \gamma_1 \gamma_2 \hat{\phi}_t^T \mathbf{B}^T \mathbf{B} \hat{\phi}_t - 2 \gamma_2 \gamma_2^2 \hat{\phi}_t^T (\mathbf{B}^T \mathbf{B} \mathbf{C}) \hat{\phi}_t \quad (35c)$$

$$= \hat{\phi}_t^T \mathbf{B}^T \left((\gamma_1^2 - 2 \gamma_1 \gamma_2) \mathbf{I} + \gamma_1^2 \gamma_2^2 \mathbf{B} \mathbf{B}^T \right) \mathbf{B} \hat{\phi}_t + \hat{\phi}_t^T \mathbf{Q}^T \mathbf{Q} \hat{\phi}_t - 2 \gamma_2 \gamma_2^2 \hat{\phi}_t^T (\mathbf{B}^T \mathbf{B} \mathbf{C}) \hat{\phi}_t \quad (35d)$$

$$\leq \left| \gamma_1^2 - 2 \gamma_1 \gamma_2 + \gamma_1^2 \gamma_2^2 L_{\max}^2 \right| L_{\max}^2 + (1 + \gamma_2 L_{\max})^2 - 2 \gamma_1 \gamma_2^2 \sigma_{\min}(\mathbf{B}^T \mathbf{B} \mathbf{C}) \left| \|\hat{\phi}_t\|_2^2 \right|. \quad (35e)$$

E.3 Convergence Properties

Property (1): Given the assumption $\Delta(\mathbf{A}, \mathbf{B}, \mathbf{C}, \gamma_1, \gamma_2) := \max \left\{ (1 - \gamma_1 \sigma_{\min})^2 (1 + \gamma_2^2 L_{\max}^2), \left| \gamma_1^2 - 2\gamma_1 \gamma_2 + \gamma_1^2 \gamma_2^2 L_{\max}^2 \right| L_{\max}^2 + (1 + \gamma_2 L_{\max})^2 - 2\gamma_1 \gamma_2^2 \sigma_{\min}(\mathbf{B}^T \mathbf{B} \mathbf{C}) \right\} < \frac{1}{2}$, we can draw up the deduction that $\|\hat{\boldsymbol{\theta}}_{t+1}\|_2^2 + \|\hat{\boldsymbol{\phi}}_{t+1}\|_2^2 \leq \Delta (\|\hat{\boldsymbol{\theta}}_t\|_2^2 + \|\hat{\boldsymbol{\phi}}_t\|_2^2)$ with $\Delta < 1$ directly from Eq.s (34) and (35).

For ease of simplicity, we rewrite model parameters as $\mathbf{z} = [\boldsymbol{\theta}, \boldsymbol{\phi}]^T \in \mathcal{Z} = \Theta \times \Phi$ and $\hat{\mathbf{z}} = [\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\phi}}]^T$. The deduction can be equivalently expressed as $\|\hat{\mathbf{z}}_{t+1}\|_2 \leq \sqrt{\Delta} \|\hat{\mathbf{z}}_t\|_2 < \|\hat{\mathbf{z}}_t\|_2$.

Then $\forall \epsilon > 0$, there always exists an integer $N \in \mathbb{N}^+$ such that for all $m > n > N$,

$$\|\mathbf{z}_m - \mathbf{z}_n\|_2 = \|\hat{\mathbf{z}}_m - \hat{\mathbf{z}}_n\|_2 \leq \sum_{t=1}^{m-n} \|\hat{\mathbf{z}}_{n+t} - \hat{\mathbf{z}}_{n+t-1}\|_2 \leq \left[\sum_{t=1}^{m-n} \sqrt{\Delta}^{n+k} \right] (\|\hat{\mathbf{z}}_0\|_2) < \epsilon, \quad (36a)$$

where the inequality can be obviously satisfied when N is large enough. This implies the iteration sequence $\{\mathbf{z}_t\}_{t=0}^H$ is Cauchy.

Property (2): Given the assumption $\Delta := \max \left\{ (1 - \gamma_1 \sigma_{\min})^2 (1 + \gamma_2^2 L_{\max}^2), \left| \gamma_1^2 - 2\gamma_1 \gamma_2 + \gamma_1^2 \gamma_2^2 L_{\max}^2 \right| L_{\max}^2 + (1 + \gamma_2 L_{\max})^2 - 2\gamma_1 \gamma_2^2 \sigma_{\min}(\mathbf{B}^T \mathbf{B} \mathbf{C}) \right\} < \frac{1}{2}$ and the property (1), we can derive that $\|\hat{\mathbf{z}}_t\|_2 \leq (\sqrt{\Delta})^t \|\hat{\mathbf{z}}_0\|_2$. After performing the limit operation, we can have:

$$\lim_{t \rightarrow \infty} \|\hat{\mathbf{z}}_t\|_2 \leq \lim_{t \rightarrow \infty} (\sqrt{\Delta})^t \|\hat{\mathbf{z}}_0\|_2 = 0 \implies \lim_{t \rightarrow \infty} [\boldsymbol{\theta}_t, \boldsymbol{\phi}_t]^T = [\boldsymbol{\theta}_*, \boldsymbol{\phi}_*]^T. \quad (37)$$

Hence, the adopted optimization strategy can guarantee at least the linear convergence with the rate $\sqrt{\Delta}$ to the (local) Stackelberg equilibrium.

F GENERALIZATION BOUND

Note that the task distribution is adaptive and learnable, and we take interest in the generalization in the context of generative task distributions. It is challenging to perform direct analysis. Hence, we propose to exploit the importance of the weighting trick. To do so, we first recap the reweighted generalization bound from [11] as **Lemma 1**. The sketch of proofs mainly consists of the importance-weighted generalization bound and estimates of the importance weights' range.

F.1 Importance Weighted Generalization Bound

LEMMA 1 (GENERALIZATION BOUND OF REWEIGHTED RISK [11]). *Given a risk function \mathcal{L} and arbitrary hypothesis $\boldsymbol{\theta}$ inside the hypothesis space $\boldsymbol{\theta} \in \Theta$ together with the pseudo-dimension $C = \text{Pdim}(\{\mathcal{L}(\cdot; \boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta\})$ in [56] and the importance variable $\omega(\tau)$, then the following inequality holds with a probability $1 - \delta$ over samples $\{\tau_1, \tau_2, \dots, \tau_K\}$:*

$$R_p^\omega(\boldsymbol{\theta}) \leq \hat{R}_p^\omega(\boldsymbol{\theta}) + 2^{\frac{5}{4}} V_{p, \hat{p}} \left[\omega(\tau) \mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}) \right] \left(\frac{C \ln \frac{2Ke}{C} + \ln \frac{4}{\delta}}{K} \right)^{\frac{3}{8}}, \quad (38)$$

where $V_{p, \hat{p}} \left[\omega(\tau) \mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}) \right] = \max \left\{ \sqrt{\mathbb{E}_p[\omega^2(\tau) \mathcal{L}^2(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta})]}, \sqrt{\mathbb{E}_{\hat{p}}[\omega^2(\tau) \mathcal{L}^2(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta})]} \right\}$ with p the exact task distribution, \hat{p} the empirical task distribution, and $\mathbb{E}_p[\omega(\tau)] = 1$.

Further we denote the expected risk of interest by $R_p^\omega(\boldsymbol{\theta}_*) = \mathbb{E}_{p_\Phi(\tau)} \left[\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}_*) \right] = \mathbb{E}_{p_0(\tau)} \left[\omega(\tau) \mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}_*) \right]$. The corresponding empirical risk is $\hat{R}_p^\omega(\boldsymbol{\theta}_*) := \frac{1}{K} \sum_{k=1}^K \omega(\tau_k^0) \mathcal{L}(D_{\tau_k^0}^Q, D_{\tau_k^0}^S; \boldsymbol{\theta}_*)$. Details of these terms' relations are attached as below.

$$\mathbb{E}_{p_0(\tau)} \left[\frac{p_\Phi(\tau)}{p_0(\tau)} \mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}_*) \right] = \mathbb{E}_{p_0(\tau)} \left[\omega(\tau) \mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}_*) \right] \approx \frac{1}{K} \sum_{k=1}^K \omega(\tau_k^0) \mathcal{L}(D_{\tau_k^0}^Q, D_{\tau_k^0}^S; \boldsymbol{\theta}_*), \text{ with } \tau_k^0 \sim p_0(\tau). \quad (39)$$

F.2 Formal Adversarial Generalization Bounds

Let us further denote the generated sequence of task identifiers by $\tau_k^{-M} \mapsto \tau_k^{-M+1} \mapsto \dots \mapsto \tau_k^0$ and consider the uniform distribution as the base distribution cases $p_0(\tau_k^{-M}) = p_0(\tau_k^0)$, then we can have the following equation:

$$\omega(\tau_k^0) = \exp \left(\ln p_\phi(\tau_k^0) - \ln p_0(\tau_k^0) \right) \quad (40a)$$

$$= \exp \left(\ln p_0(\tau_k^{-M}) + \sum_{i=1}^M \ln \left| \det \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} \right| - \ln p_0(\tau_k^0) \right) \quad (40b)$$

$$= \exp \left(\ln p_0(\tau_k^0) + \sum_{i=1}^M \ln \left| \det \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} \right| - \ln p_0(\tau_k^0) \right) \quad (40c)$$

$$= \exp \left(\sum_{i=1}^M \ln \left| \det \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} \right| \right). \quad (40d)$$

Note that the invertible functions inside normalizing flows are assumed to hold the bi-Lipschitz property. For the estimate of the term $\left| \det \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} \right|$, we directly apply Hadamard's inequality [46] to obtain:

$$\begin{aligned} \left| \det \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} \right| &\leq \prod_{s=1}^d \left\| \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} v_s \right\| \leq \left\| \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} \right\|^d \leq \ell_a^d, \forall m \in \{1, 2, \dots, M\} \\ &\implies \ln \left| \det \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} \right| \leq d \ln \ell_a \\ &\implies \sum_{i=1}^M \ln \left| \det \frac{\partial g_i^{-1}}{\partial \tau_k^{-M+i}} \right| \leq dM \ln \ell_a, \end{aligned} \quad (41)$$

where v_s is a collection of unit eigenvectors of the considered Jacobian. The above implies that $\omega(\tau_k^0) \leq \ell_a^{Md}$.

Now, we put the above derivations together and present the formal adversarial bound as follows.

$$\max \left\{ \mathbb{E}_{p_0(\tau)} [\omega^2(\tau) \mathcal{L}^2(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta_*)], \mathbb{E}_{\hat{p}_0(\tau)} [\omega^2(\tau) \mathcal{L}^2(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta_*)] \right\} \leq \ell_a^{2Md} \sup_{\tau \in \mathcal{T}} |\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta_*)|^2 \quad (42)$$

This formulates the generalization bound of meta learners under the learned adversarial task distribution.

Theorem 2 (Generalization Bound with the Distribution Adversary) *Given the pretrained normalizing flows $\{g_i\}_{i=1}^M$, where g_i is presumed to be (ℓ_a, ℓ_b) -bi-Lipschitz, the pretrained meta learner θ_* from the hypothesis space $\theta_* \in \Theta$ together with the pseudo-dimension $C = \text{Pdim}(\{\mathcal{L}(\cdot; \theta) : \theta \in \Theta\})$ in [56], we can derive the generalization bound when the initial task distribution p is uniform.*

$$R_p^\omega(\theta) \leq \hat{R}_p^\omega(\theta) + \Upsilon(\mathcal{T}) \left(\frac{C \ln \frac{2Ke}{C} + \ln \frac{4}{\delta}}{K} \right)^{\frac{3}{8}}, \quad (43)$$

the inequality holds with a probability $1 - \delta$ over samples $\{\tau_1, \tau_2, \dots, \tau_K\}$, where the constant means $\Upsilon(\mathcal{T}) = 2^{\frac{5}{4}} \ell_a^{2Md} \sup_{\tau \in \mathcal{T}} |\mathcal{L}(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta_*)|^2$.

G EXPLICIT GENERATIVE TASK DISTRIBUTIONS

Once the meta training process is finished, we can direct access the generative task distributions from the pretrained normalizing flows. For any task τ^0 , we still assume the generated task sequence is $\tau^{-M} \mapsto \tau^{-M+1} \mapsto \dots \mapsto \tau^0$ after the invertible transformations $\{g_i\}_{i=1}^M$. Then, we can obtain the exact likelihood of sampling the task τ as:

$$p_{\phi_*}(\tau^0) = \frac{p_0(\tau^{-M})}{\prod_{i=1}^M \left| \det \frac{\partial g_i^{-1}}{\partial \tau^{-M+i}} \right|}, \quad (44)$$

where the product of the Jacobian scales the task probability density.

H EVOLUTION OF ENTROPIES IN TASK DISTRIBUTIONS

One of the primary benefits of employing normalizing flows in generating task distributions is the tractable likelihood, leaving it possible to analyze the entropy of task distributions. Note that **Remark 1** in the main paper characterizes the way of computing the generative task distribution entropy.

Remark 1 (Entropy of the Generated Task Distribution) Given the generative task distribution $p_{\phi_*}(\tau)$, we can derive its entropy from the initial task distribution $p_0(\tau)$ and normalizing flows $\mathcal{G} = \{g_i\}_{i=1}^M$:

$$\mathbb{H}[p_{\phi_*}(\tau)] = \mathbb{H}[p_0(\tau)] + \int p_0(\tau) \left[\sum_{i=1}^M \ln \left| \det \frac{\partial g_i}{\partial \tau^i} \right| \right] d\tau. \quad (45)$$

The above implies that the generated task distribution entropy is governed by the change of task identifiers in the probability measure of task space.

Next is to provide the proof w.r.t. the above Remark.

Proof: Given the initial distribution $p_0(\tau)$ and the sampled task τ^0 , we know the transformation within the normalizing flows:

$$\tau^M = g_M \circ \dots \circ g_2 \circ g_1(\tau^0) = \text{NN}_{\phi}(\tau^0).$$

Also, the density function after transformations is:

$$\ln p_{\phi}(\tau^M) = \ln p_0(\tau^0) - \sum_{i=1}^M \ln \left| \det \frac{\partial g_i}{\partial \tau^i} \right|.$$

With the above information, we can naturally derive the following equations:

$$\begin{aligned} \mathbb{H}[p_{\phi_*}(\tau^M)] &= - \int p_{\phi_*}(\tau^M) \ln p_{\phi_*}(\tau^M) d\tau^M \\ &= - \int p_0(\tau^0) \left[\ln p_0(\tau^0) - \sum_{i=1}^M \ln \left| \det \frac{\partial g_i}{\partial \tau^i} \right| \right] d\tau^0 \\ &= \mathbb{H}[p_0(\tau^0)] + \int p_0(\tau^0) \left[\sum_{i=1}^M \ln \left| \det \frac{\partial g_i}{\partial \tau^i} \right| \right] d\tau^0. \end{aligned} \quad (46)$$

This completes the proof of **Remark 1**.

I EXPERIMENTAL SET-UP & IMPLEMENTATION DETAILS

I.1 Meta Learning Benchmarks

Table 4: A Summary of Benchmarks in Evaluation. We report task identifiers in configuring tasks, corresponding initial ranges and types of base distributions. \mathcal{U} and \mathcal{N} respectively denote the uniform and the normal distributions, with distribution parameters inside the bracket.

Benchmarks	Task Identifiers	Identifier Range	Initial Distribution
Sinusoid-U	amplitude a and phase b	$(a, b) \sim [0.1, 5.0] \times [0.0, \pi]$	$\mathcal{U}([0.1, 5.0] \times [0.0, \pi])$
Sinusoid-N	amplitude a and phase b	$(a, b) \sim [0.1, 5.0] \times [0.0, \pi]$	$\mathcal{N}([2.5, 1.5], \text{diag}(0.8^2, 0.5^2))$
Acrobot-U	pendulum masses (m_1, m_2)	$(m_1, m_2) \sim [0.4, 1.6] \times [0.4, 1.6]$	$\mathcal{U}([0.4, 1.6] \times [0.4, 1.6])$
Acrobot-N	pendulum masses (m_1, m_2)	$(m_1, m_2) \sim [0.4, 1.6] \times [0.4, 1.6]$	$\mathcal{N}([1.0, 1.0], \text{diag}(0.2^2, 0.2^2))$
Pendulum-U	pendulum mass m and length l	$(m, l) \sim [0.4, 1.6] \times [0.4, 1.6]$	$\mathcal{U}([0.4, 1.6] \times [0.4, 1.6])$
Pendulum-N	pendulum mass m and length l	$(m, l) \sim [0.4, 1.6] \times [0.4, 1.6]$	$\mathcal{N}([1.0, 1.0], \text{diag}(0.2^2, 0.2^2))$
Point Robot	goal location (x_1, x_2)	$(x_1, x_2) \sim [-0.5, 0.5] \times [-0.5, 0.5]$	$\mathcal{U}([-0.5, 0.5] \times [-0.5, 0.5])$
Pos-Ant	target position (x_1, x_2)	$(x_1, x_2) \sim [-3.0, 3.0] \times [-3.0, 3.0]$	$\mathcal{U}([-3.0, 3.0] \times [-3.0, 3.0])$

The following details the meta-learning dataset, and we also refer the reader to the **List (1)** for the preprocessing of data.

Sinusoid Regression. In sinusoid regression, each task is equivalent to mapping the input to the output of a sine wave. Here, the task identifiers are the amplitude a and the phase b . Data points in regression are collected in the way: 10 data points are uniformly sampled from the interval $[-5.0, 5.0]$, coupled with the output $y = a \sin(x - b)$. These data points are divided into the support dataset (5-shot) and the query dataset. As for the range of task identifiers and the types of initial distributions, please refer to Table 4. In meta-testing phases, we randomly sample 500 tasks from the initial task distribution and the generated task distribution to evaluate the performance, and this results in Table 1.

System Identification. In the Acrobot System, angles and angular velocities characterize the state of an environment as $[\theta_1, \theta'_1, \theta_2, \theta'_2]$. The goal is to identify the system dynamics, namely the transited state, after selecting Torque action from $\{-1, 0, +1\}$. In the Pendulum System, environment information can be found in the OpenAI gym. In detail, the observation is $(\cos \theta, \sin \theta, \theta')$ with $\theta \in [-\pi, \pi]$, and the continuous action range is $[-2.0, 2.0]$. The torque is executed on the pendulum body, and the goal is to predict the dynamics given the

observation and action. For both dynamical systems, we use a random policy as the controller to interact with the environment to collect transition samples. For the few-shot purpose, we sample 200 transitions from each Markov Decision Process as one batch and randomly split them into the support dataset (10-shot transitions) and the query dataset. In meta training, the meta-batch in iteration is 16, and the maximum iteration number is 500. In meta-testing phases, we randomly sample 100 tasks respectively from the initial and generated task distributions with each task one batch in evaluation, which results in Table 1.

Meta Learning Continuous Control. The environment of reinforcement learning is mainly treated as a Markov decision process [44, 45]. And the meta RL is about the distribution over environments. We consider the navigation task to examine the performance of methods in reinforcement learning. The mission is to guide the robot, e.g., the point robot and the Ant from Mujoco, to move towards the target goal step by step. Hence, the task identifier is the goal location (x_1, x_2) . The agent performs 20 episode explorations to identify the environment and enable inner policy gradient updates as fast adaptation. The environment information, such as transitions and rewards, is accessible at (<https://github.com/lmzintgraf/cavia/tree/master/rl/envs>). Particularly, for the task distribution like $\mathcal{U}([-0.5, 0.5] \times [-0.5, 0.5])$ in the point robot, we set the dark area in Figure 3c as the sparse reward region, discounting the step-wise reward by 0.6 to area $[-0.25, -0.5] \times [0.25, 0.5]$ and 0.4 to area $[0.25, 0.5] \times [-0.25, -0.5]$. In meta training, the meta-batch in the iteration is 20 for the point robot, and 40 for Ant, and the model is trained for up to 500 meta-iterations. In meta-testing, we randomly sampled 100 MDPs, specifically from the initial and generated task distributions. For each MDP, we run 20 episodes as the support dataset and compute the accumulated returns after fast adaptation.

The general implementation details are retained the same as those in MAML (<https://github.com/tristandeleu/pytorch-maml-rl>) and CAVIA (<https://github.com/lmzintgraf/cavia>).

I.2 Modules in Pytorch

To enable the implementation of our developed framework, we report the neural modules implemented in Pytorch. The whole model consists of the distribution adversary and the meta learner. For the sake of clarity, we attach the separate modules below. Our implementation also relies on the normalizing flow package available at (<https://github.com/VincentStimper/normalizing-flows/tree/master/normflows>).

```

1  """
2  Neural network models for the regression experiments with adaptively robust maml.
3  """
4
5  import math
6  import torch
7  import torch.nn.functional as F
8  from torch import nn
9  import normflows as nf
10 from normflows.flows import Planar, Radial
11
12 #####
13 # This part is to introduce the flow module to transform random variables. --> Distribution Adversary
14 #####
15
16
17 class Distribution_Adversary(nn.Module):
18     def __init__(self,
19                 q0,
20                 latent_size,
21                 num_latent_layers,
22                 flow_type,
23                 hyper_range,
24                 device
25                 ):
26         '''
27         q0: base distribution of task parameters
28         latent_size: the dimension of the latent variable
29         num_latent_layers: number of layers in NFs
30         flow_type: types of NFs
31         hyper_range: range of task hyper-parameters, tensor shape [dim_z, 2] e.g. [[4.9, 0.1], [3.0, -1.0]] -> [
32         param_range, range_min]
33         device: 'cuda' or 'cpu'
34         '''
35
36         super(Distribution_Adversary, self).__init__()
37
38         self.q0 = q0
39         self.latent_size = latent_size
40         self.num_latent_layers = num_latent_layers
41         self.flow_type = flow_type
42         self.hyper_range = hyper_range

```

```

42     self.device = device
43
44     if flow_type == 'Planar_Flow':
45         flows = [Planar(self.latent_size) for k in range(self.num_latent_layers)]
46     elif flow_type == 'Radial_Flow':
47         flows = [Radial(self.latent_size) for k in range(self.num_latent_layers)]
48
49     self.nfm = nf.NormalizingFlow(q0=self.q0, flows=flows)
50     self.nfm.to(device)
51
52     def forward(self, x, train=True):
53         log_det = self.q0.log_prob(x)
54         z, log_det_forward = self.nfm.forward_and_log_det(x)
55
56         min_values = torch.min(z, dim=0).values
57         max_values = torch.max(z, dim=0).values
58         normalized_data = (z - min_values) / (max_values - min_values)
59
60         normalize_tensor = (self.hyper_range.to(self.device)).expand(z.size()[0], -1, -1) # output shape [task_batch,
        dim_z, 2]
61         norm_z = normalize_tensor[:, :, 0] * normalized_data + normalize_tensor[:, :, 1] # normalize the transformed task
        into valid ranges
62
63         log_det_norm = torch.sum(torch.log(normalize_tensor[:, :, 0] / (max_values - min_values)), dim=-1)
64
65         z_reverse, loss = self.nfm.forward_kld(x)
66
67         a, b = self.hyper_range[0][1], torch.sum(self.hyper_range[0])
68         c, d = self.hyper_range[1][1], torch.sum(self.hyper_range[1])
69         condition = (z_reverse[:, 0] >= a) & (z_reverse[:, 0] <= b) & (z_reverse[:, 1] >= c) & (z_reverse[:, 1] <= d)
70         z_reverse = z_reverse[condition]
71
72         if z_reverse.shape[0] == 0:
73             log_det_reverse_total = torch.tensor(0.)
74         else:
75             log_det_z = self.q0.log_prob(z_reverse)
76             log_det_reverse_total = -torch.mean(log_det_z) + loss
77
78         return z, norm_z, log_det, log_det_forward, log_det_norm, log_det_reverse_total
79
80
81 # This part is to introduce the MLP for the implementation of MAML. --> Meta Player
82
83 class Meta_Learner(nn.Module):
84     def __init__(self,
85                 n_inputs,
86                 n_outputs,
87                 n_weights,
88                 task_type,
89                 device
90                 ):
91         '''
92         n_inputs: the number of inputs to the network,
93         n_outputs: the number of outputs of the network,
94         n_weights: for each hidden layer the number of weights, e.g., [128,128,128]
95         device: device to deploy, cpu or cuda
96         '''
97
98         super(Meta_Learner, self).__init__()
99
100         # initialise lists for biases and fully connected layers
101         self.weights = []
102         self.biases = []
103
104         # add one
105         if task_type == 'sine':
106             self.nodes_per_layer = n_weights + [n_outputs]
107         elif task_type == 'acrobot':
108             self.nodes_per_layer = n_weights + [n_outputs-2]

```

```

109     elif task_type == 'pendulum':
110         self.nodes_per_layer = n_weights + [n_outputs-1]
111
112     # additional biases
113     self.task_context = torch.zeros(0).to(device)
114     self.task_context.requires_grad = True
115
116     # set up the shared parts of the layers
117     prev_n_weight = n_inputs
118     for i in range(len(self.nodes_per_layer)):
119         w = torch.Tensor(size=(prev_n_weight, self.nodes_per_layer[i])).to(device)
120         w.requires_grad = True
121         self.weights.append(w)
122         b = torch.Tensor(size=[self.nodes_per_layer[i]]).to(device)
123         b.requires_grad = True
124         self.biases.append(b)
125         prev_n_weight = self.nodes_per_layer[i]
126
127     self._reset_parameters()
128
129     def _reset_parameters(self):
130         for i in range(len(self.nodes_per_layer)):
131             stdv = 1. / math.sqrt(self.nodes_per_layer[i])
132             self.weights[i].data.uniform_(-stdv, stdv)
133             self.biases[i].data.uniform_(-stdv, stdv)
134
135     def forward(self, x, task_type='sine'):
136         x = torch.cat((x, self.task_context))
137
138         for i in range(len(self.weights) - 1):
139             x = F.relu(F.linear(x, self.weights[i].t(), self.biases[i]))
140
141         if task_type == 'sine':
142             y = F.linear(x, self.weights[-1].t(), self.biases[-1])
143         elif task_type == 'acrobot':
144             y = F.linear(x, self.weights[-1].t(), self.biases[-1])
145             y = torch.cat((torch.cos(y[... , 0:1]), torch.sin(y[... , 0:1]), torch.cos(y[... , 1:2]), torch.sin(y[... , 1:2]), y
146             [... , :2]), dim=-1)
147         elif task_type == 'pendulum':
148             y = F.linear(x, self.weights[-1].t(), self.biases[-1])
149             y = torch.cat((torch.cos(y[... , 0:1]), torch.sin(y[... , 0:1]), y[... , :1]), dim=-1)
150
151         return y

```

Listing 1: the Distribution Adversary and the Meta Learner.

As noted in the distribution adversary, the last layer is to normalize the range of the task identifiers into the pre-defined range, where the min-max normalization $\sigma(x) = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$ is utilized. We consider the dimension of the task identifier to be d . Let $a = \text{normalize_tensor}[:, 0] = [a_1, a_2, \dots, a_d]$ and $b = \text{normalize_tensor}[:, 1] = [b_1, b_2, \dots, b_d]$, the final layer of normalizing flows g_m can be expressed as $\tau^m = g_m(\tau^{m-1}) = a \odot \sigma(\tau^{m-1}) + b$, where τ^m indicates the transformed task τ^m and the resulting sequence after $\{g_m\}_{m=1}^M$ is $\tau^0 \mapsto \tau^1 \mapsto \dots \mapsto \tau^M$. Then, the resulting log-probability follows the computation:

$$\begin{aligned}
 \ln p(\tau^m) &= \ln p(\tau^{m-1}) - \ln \left| \det \frac{dg_m(\tau^{m-1})}{d\tau^{m-1}} \right| \\
 &= \ln p(\tau^{m-1}) - \ln \left(\prod_{i=1}^d \frac{a_i}{\tau_{\max}^{m-1,i} - \tau_{\min}^{m-1,i}} \right) \\
 &= \ln p(\tau^{m-1}) - \left(\sum_{i=1}^d \ln a_i - \ln(\tau_{\max}^{m-1,i} - \tau_{\min}^{m-1,i}) \right).
 \end{aligned} \tag{47}$$

$$\mathbb{E}_{p_{\phi_*}(\tau)} \left[\mathcal{L}(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta_*) \right] \approx \frac{1}{K} \sum_{k=1}^K \mathcal{L}(D_{\tau_k^M}^Q, D_{\tau_k^M}^S; \theta_*), \quad \text{with } \tau_k^M = \text{NN}_{\phi_*}(\tau_k^0) \text{ and } \tau_k^0 \sim p_0(\tau) \tag{48}$$

I.3 Neural Architectures & Optimization

Our approach is meta-learning method agnostic, and the implementation is *w.r.t.* MAML [17] and CNP [19] in this work. As a result, we respectively describe the neural architectures in separate implementations and benchmarks. We do not vary the neural architecture of meta learners, and only risk minimization principles are studied in experiments.

In the sinusoid regression and system identification benchmarks, we set the Lagrange multiplier $\lambda = 0.2$ to cover most shifted task distributions. All MAML-like methods employ a multilayer perceptron neural architecture. This architecture comprises three hidden layers, with each layer consisting of 128 hidden units. The activation function utilized in these models is the Rectified Linear Unit (ReLU). The inner loop utilizes the stochastic gradient descent (SGD) algorithm to perform fast adaptation on each task, while the outer loop employs the meta-optimizer Adam to update the initial parameters of the model. The learning rate for both the inner and outer loops is set to $1e-3$. As for CNP-like methods, the encoder uses a three-layer MLP with 128-dimensional hidden units and outputs a 128-dimensional representation. The output representations are averaged to form a single representation. This aggregated representation is concatenated with the query dataset and passed through a two-layer MLP decoder. The optimizer used is Adam, with a learning rate of $1e-3$.

In the continuous control benchmark, we set $\lambda = 0.0$, and MAML serves as the underlying neural network architecture in our research. The agent’s reward is the negative squared distance to the goal. The agents are trained for one gradient update, employing policy gradient with the generalized advantage estimation [67] in the inner loop and trust-region policy optimization (TRPO) [66] in the outer loop update. The learning rate for the one-step gradient update is set to 0.1.

I.4 Distribution Adversary Implementations

The distribution adversary is implemented with the help of normalizing flows. For AR-MAML, we adopt the distribution adversary with the neural network as follows. We employ a 2-layer Planar flow [61] to transform the initial distribution. In each layer, the dimension of the latent variable is 2, and the activation function is leaky ReLU. In the last layer, the min-max normalization is used. We use Adam with a cosine learning rate scheduler for the distribution adversary optimizer.

J ADDITIONAL EXPERIMENTAL RESULTS

Tail Risk Robustness. Due to the page limit in the main paper, we present experimental results on tail risk robustness across various confidence levels α for sinusoid regression and acrobot system identification. As illustrated in Figures 14/15, the advantage of AR-MAML is similar to that in pendulum system identification. AR-MAML exhibits a more significant performance gain in the adversarial distribution than in the initial distribution.

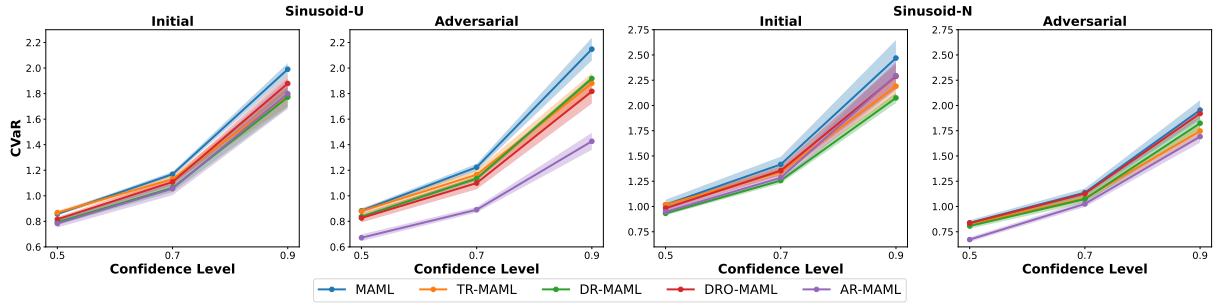


Figure 14: $CVaR_{\alpha}$ MSEs with Various Confidence Level α . Sinusoid-U/N denotes Uniform/Normal as the initial distribution type. The plots report meta testing $CVaR_{\alpha}$ MSEs in initial and adversarial distributions with standard error bars in shadow regions.

Impacts of Shift Distribution Constraints. Still, by varying the Lagrange multipliers λ , we include the learned task structures and the meta-testing results in Figures 16/17/18. For sinusoid and acrobot cases, empirical findings are similar to those in the main paper. In point robot navigation, we observe high probability density regions significantly change, and even the meta-testing results are improved in the initial task distribution with increasing λ values.

K PLATFORMS & COMPUTATIONAL TOOLS

This project uses NVIDIA A100 GPUs in numeric computation. And we employ Pytorch [55] as the deep learning toolkit in implementing experiments.

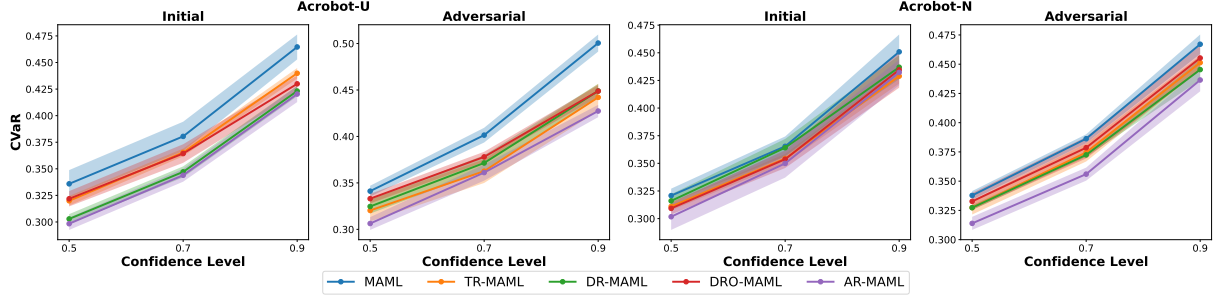


Figure 15: CVaR_α MSEs with Various Confidence Level α . Acrobot-U/N denotes Uniform/Normal as the initial distribution type. The plots report meta testing CVaR_α MSEs in initial and adversarial distributions with standard error bars in shadow regions.

Table 5: Entropy of initial distribution and adversarial distribution under different λ values.

Benchmark	Meta-Test Distribution	$\lambda = 0.0$	$\lambda = 0.1$	$\lambda = 0.2$
Sinusoid-U	Initial	2.734	2.734	2.734
	Adversarial	2.15 ± 0.01	2.44 ± 0.01	2.46 ± 0.01
Sinusoid-N	Initial	1.922	1.922	1.922
	Adversarial	-2.92 ± 0.00	1.11 ± 0.02	1.72 ± 0.00

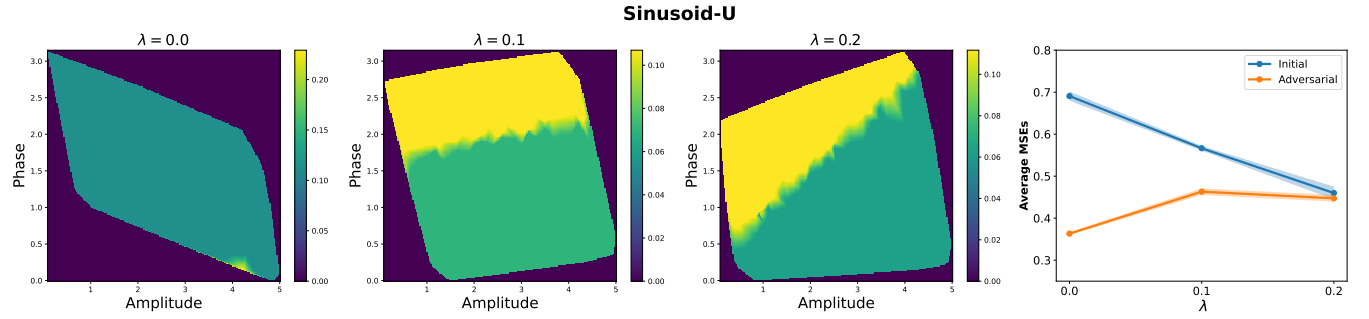


Figure 16: The first three plots show adversarial task probability distributions with varying Lagrange multipliers λ in the sinusoid-U benchmark. The last plot depicts meta testing MSEs across different values of λ .

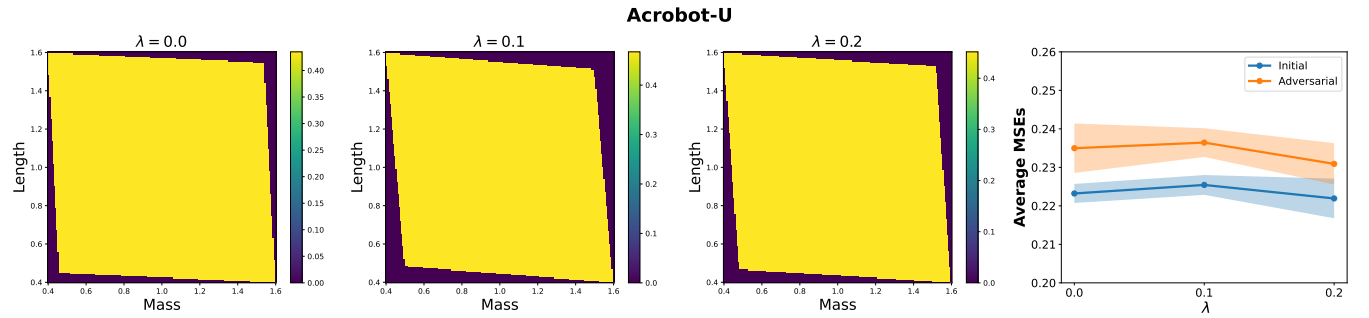


Figure 17: The first three plots show adversarial task probability distributions with varying Lagrange multipliers λ in the acrobot-U benchmark. The last plot depicts meta testing MSEs across different values of λ .

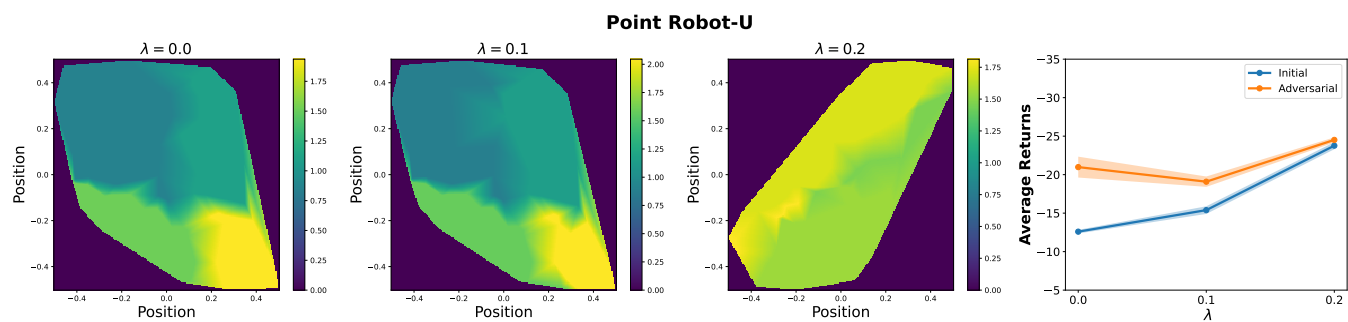


Figure 18: The first three plots show adversarial task probability distributions with varying Lagrange multipliers λ in the point robot-U benchmark. The last plot depicts meta testing MSEs across different values of λ .