

Model Predictive Task Sampling for Efficient and Robust Adaptation

Qi (Cheems) Wang^{1*} Zehao Xiao^{2*} Yixiu Mao^{1*} Yun Qu^{1*} Jiayi Shen² Yiqin Lv¹ Xiangyang Ji^{1†}

¹Department of Automation, Tsinghua University; ²Informatics Institute, University of Amsterdam

[†]Correspondence Author: xyji@tsinghua.edu.cn

Abstract

Foundation models have revolutionized general-purpose problem-solving, offering rapid task adaptation through pretraining, meta-training, and finetuning. Recent crucial advances in these paradigms reveal the importance of challenging task prioritized sampling to enhance adaptation robustness under distribution shifts. However, ranking task difficulties over iteration as a preliminary step typically requires exhaustive task evaluation, which is practically unaffordable in computation and data-annotation. This study provides a novel perspective to illuminate the possibility of leveraging the dual importance of adaptation robustness and learning efficiency, particularly in scenarios where task evaluation is risky or costly, such as iterative agent-environment interactions for robotic policy evaluation or computationally intensive inference steps for finetuning foundation models. Firstly, we introduce **Model Predictive Task Sampling (MPTS)**, a framework that bridges the task space and adaptation risk distributions, providing a theoretical foundation for robust active task sampling. MPTS employs a generative model to characterize the episodic optimization process and predicts task-specific adaptation risk via posterior inference. The resulting risk predictive model amortizes the costly evaluation of task adaptation performance and provably approximates task difficulty rankings. MPTS seamlessly integrates into *zero-shot*, *few-shot*, and *supervised finetuning* settings. Empirically, we conduct extensive experiments in pattern recognition using foundation models and sequential decision-making. Our results demonstrate that MPTS significantly enhances adaptation robustness for tail risk or out-of-distribution (OOD) tasks and improves learning efficiency compared to state-of-the-art (SoTA) methods. The code is available at the project site <https://github.com/thu-rllab/MPTS>.

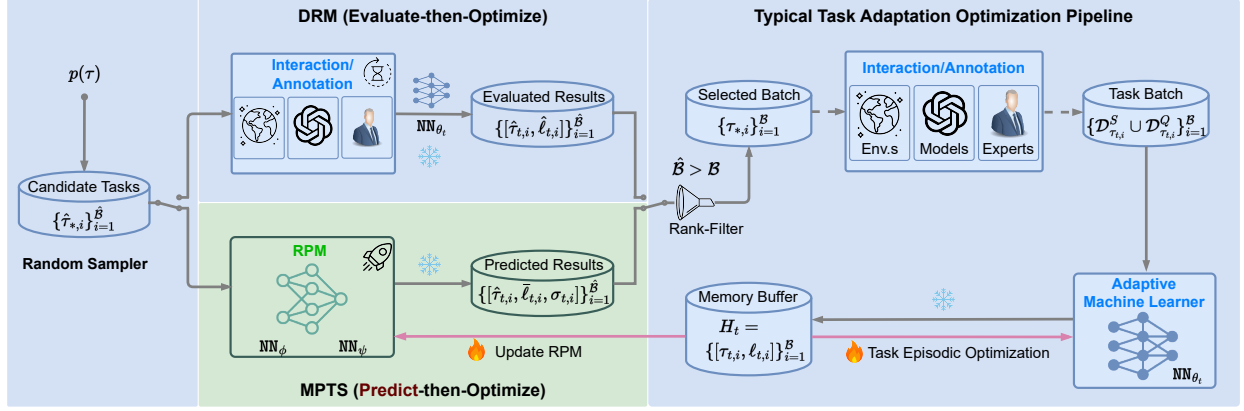
1 Introduction

Generalization across diverse scenarios remains a central challenge in artificial general intelligence. The rise of generative AI offers a promising solution, driving the development of foundation models¹⁻³. Unlike traditional task-specific models, which might fail in new tasks, foundation models enable fast deployment across diverse scenarios without learning from scratch. Their rapid problem-solving stems from widely adopted adaptation learning paradigms, including pretraining, meta-learning, and supervised finetuning (SFT).

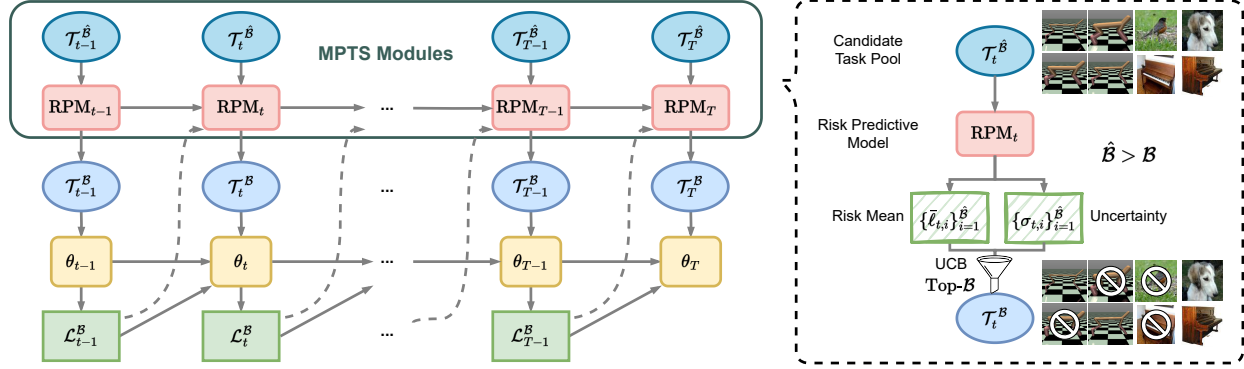
These paradigms train machine learners over a task distribution, consolidating past experience into problem-solving priors to handle unseen but related tasks in zero-shot or few-shot settings^{1,4}. Each iteration samples a task batch, e.g., from a uniform distribution, and then executes the learning-to-adapt step (see Fig. 1). Large language models (LLMs), for instance, treat episodic corpus datasets as tasks and perform in-context learning for adaptation⁵. Similarly, in obtaining generalist robotic policies, decision-making environments, such as Markov decision processes (MDPs), are randomized for robots to perform policy optimization. These task distributions are typically determined by identifiers; e.g., in Fig. 2a, varying physics parameters configure different MDPs as tasks for domain randomization (DR)⁶ and meta reinforcement learning (Meta-RL)⁷.

Research Motivations: Distribution shifts^{8,9} are prevalent in real-world scenarios, making task adaptation robustness at test time increasingly critical^{10,11}. In this context, task sampling strategies play a pivotal role, yet uniform sampling often causes catastrophic failures in risk-sensitive scenarios due to the undersampling of critical tasks. Two real-world applications highlight this case: (i) *Tail tasks*. In developing autonomous-driving systems, traffic accidents are rare in training datasets but disproportionately important for testing robustness¹². (ii) *OOD tasks*. Robots trained in controlled environments struggle in unstructured real-world settings, e.g., leading to errors in navigation and object manipulation. To improve robustness, challenging task prioritized sampling¹²⁻¹⁴ has gained traction, where assessing task difficulty is central to robust optimization.

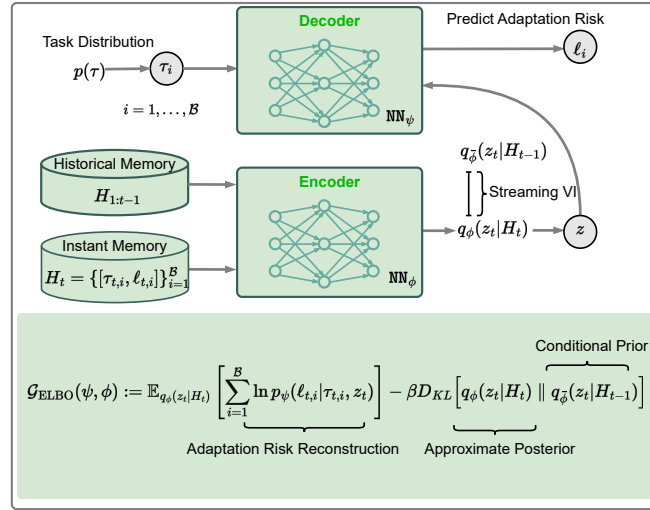
a. Incorporating Model Predictive Task Sampling into Adaptation Learning Pipelines (Pretrain/Meta-train/Fine-tune)



b. Model Predictive Task Sampling Pipelines



c. Amortized Inference and Risk Learner Architectures (Lightweight Risk Predictive Model in Optimization)



d. Amortized Evaluation and Task Prioritization (without Querying Task Dataset)

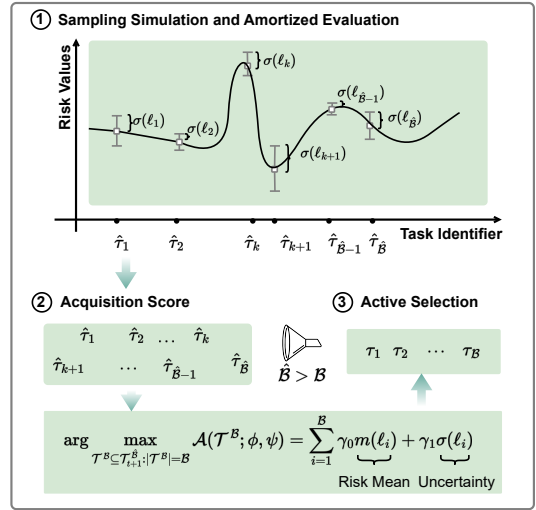


Figure 1: **Framework of MPTS in Adaptation Learning.** **a.** The left shows standard random sampling for generating candidate tasks. The middle in blue denotes costly evaluation of \hat{B} tasks (e.g., agent–environment interaction or foundation model forward pass) in DRM to select Top- B worst ones. The middle in green depicts MPTS, which predicts task difficulty via a lightweight generative model, avoiding expensive evaluation. The right illustrates the standard optimization pipeline in Meta-Learning, DR, or SFT (Snow: frozen models; Fire: updated models). **b.** MPTS samples candidate tasks, ranks their difficulty via the risk predictive model’s predictions for subset selection, and updates the learner, approximating CVAR $_{\alpha}$ ’s Monte Carlo optimization in a predictive manner. The gathered risk signals further update RPMs online. **c.** The RPM utilizes the risk history $H_{1:t}$ to train under a streaming VI framework. **d.** The RPM simulates adaptation outcomes $p(\ell|\tau, H_{1:t}; \theta_t)$ for \hat{B} candidate identifiers, computes acquisition scores, and selects the Top- B identifiers for the $(t + 1)$ -th iteration.

These methods^{13–18} evaluate, rank, and prioritize difficult tasks for iterative optimization (see Fig. 1a). However, precisely evaluating tasks via losses, human annotations, or gradients incurs high computational costs. For instance, in LLM alignment, task evaluation through SFT requires extensive forward passes, while preference optimization consumes millions of expert annotations¹⁹. Similarly, in DR and Meta-RL, agents must interact with numerous MDPs to collect post-adaptation episodes and compute returns. These challenges uncover the urgent need for more efficient learning strategies when enhancing robust adaptation, particularly when deploying foundation models or when environment interactions are costly.

Motivated by the above pressing demands, we dive into *robust active task sampling*, a paradigm that has the potential to eliminate unnecessary costs associated with task construction, intensive annotations, or computational overhead during the evaluation of a machine learner’s adaptation to specific tasks. In scenarios involving zero-shot learning, few-shot learning, or SFT, we aim to develop a task sampling strategy that requires *fewer* learning resources but retains *more* deployment benefits such as adaptation robustness in pattern recognition with foundation models and risk-averse sequential decision-making.

Developed Approach: Note that our brain is energy-efficient and simulates the outcome of decision-making in unencountered scenarios from accumulated experience, without actual trials. This capability arises from mechanisms like implicit information gating and active task selection^{20–22}. Inspired by this, we propose a model-based optimization approach for adaptive learning, dynamically adjusting task sampling strategies using predicted outcomes as feedback. This work explores the design of risk predictive models (RPMs) for robust task sampling based on two key insights: (i) Adaptation risk is probably predictable in episodic learning, providing a basis for task difficulty ranking and selection; (ii) Generative modeling of adaptation risk captures risk distributions throughout the task space with quantified uncertainty, aligning optimization with robustness principles.

To this end, we introduce Model Predictive Task Sampling (MPTS), a framework for risk-aware task selection. As shown in Fig. 1a-b, MPTS leverages historical risk to train a lightweight RPM, which forecasts adaptation risks across the task space to guide the task sampler and optimize the adaptive machine learner. This way amortizes expensive task evaluation for ranking their difficulty to select the worst subset (see the comparison in the middle block of Fig. 1a). The RPM in Fig. 1c adopts a variational autoencoder (VAE)²³ structure, generating adaptation risk estimates via posterior inference²⁴. Finally, the acquisition function in Fig. 1d integrates worst-case performance and predictive uncertainty into the rule of subset selection.

MPTS also draws inspiration from active inference²⁵, which operates through a loop of perception, action, and learning to minimize uncertainty about the planning environment. Here, subset selection from the task batch can be viewed as online planning to derive a robust machine learner. Technically, MPTS specifies or infers identifiers from the task distribution (see examples in Fig. 2a) to establish mappings between identifiers and adaptation risk. It employs streaming variational inference (VI)^{26,27} for the RPM’s training. Furthermore, by simulating adaptation outcomes in a larger identifier pool for subset selection, MPTS balances exploration (uncertainty minimization) and exploitation (worst-case robustness) across the task space. MPTS is also theoretically grounded, where the optimization proceeds in the direction of robustness enhancement while assuring convergence. As far as we know, this work is the first to examine the predictability of online task difficulty with generative models. In primary, our proposed MPTS enjoys several benefits in practice:

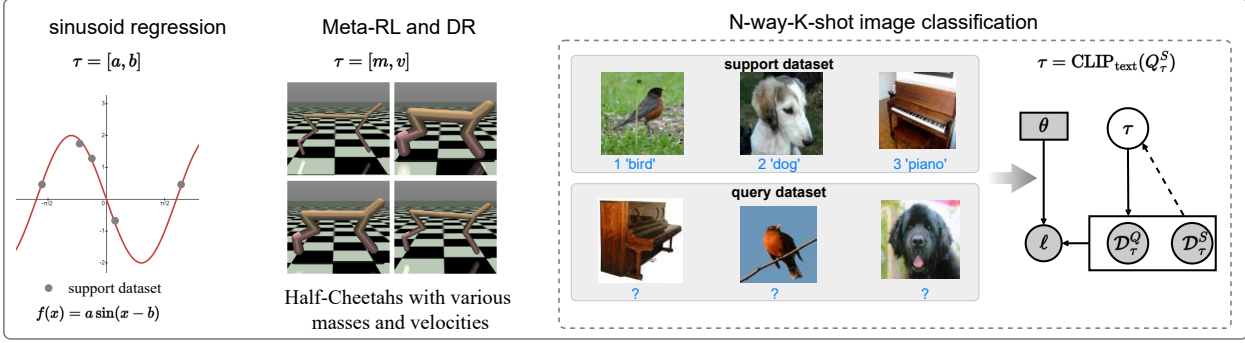
1. **Adaptation Robustness.** The optimization pipeline of MPTS can advance the machine learner’s adaptation robustness under severe task distribution shifts, such as tail risk or OOD task scenarios;
2. **Learning Efficiency.** Constructing the lightweight RPM to amortize expensive task evaluation, MPTS can diminish computational overhead, avoid unnecessary annotations, and promote efficient exploration in the task space;
3. **Framework Versatility.** Learning from risk histories, MPTS serves as a plug-and-play module to rank the task difficulties in optimization and allows seamlessly integration into robust zero-shot or few-shot learning and SFT.

This work evaluates MPTS across few-shot regression, image classification with foundation models, Meta-RL, robotic DR, and prompt-tuning foundation models. Empirical results demonstrate MPTS’s outstanding adaptation robustness across diverse scenarios. Compared to SoTA robust adaptation methods, MPTS significantly reduces computational overhead, memory usage, and environment interactions while, in some cases, accelerating learning.

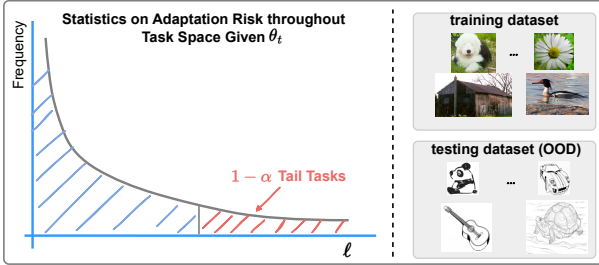
2 Adaptation and Robustness

Notations. We represent a task sample by $\tau \sim p(\tau)$, with \mathcal{T} denoting the task space. Each task τ within the distribution is specified by an identifier³¹, a real-valued vector $\boldsymbol{\tau}$, as illustrated in Fig. 2a. The task-specific risk function $\ell : \mathcal{D}_\tau^S \cup \mathcal{D}_\tau^Q \times \Theta \mapsto \mathbb{R}$ evaluates the adaptation performance of a machine learner θ on τ . For example, in regression, the support dataset $\mathcal{D}_\tau^S = \{[\mathbf{x}_i, \mathbf{y}_i]\}_{i=1}^K$ enables rapid adaptation to obtain the model $p_\theta(\mathbf{y}|\mathcal{D}_\tau^S, \mathbf{x})$; while the query dataset $\mathcal{D}_\tau^Q = \{[\mathbf{x}_i, \mathbf{y}_i]\}_{i=K+1}^{K+N}$ is used for post-adaptation evaluation as risk $\ell = -\frac{1}{N} \sum_{i=K+1}^{K+N} \ln p_\theta(\mathbf{y}_i|\mathcal{D}_\tau^S, \mathbf{x}_i)$.

a. Task Concept and Explicit or Implicit Task Identifiers



b. Tail Task Risk and OOD Scenario in Domain Generalization



c. Generative Model and Recognition Model for Episodic Learning

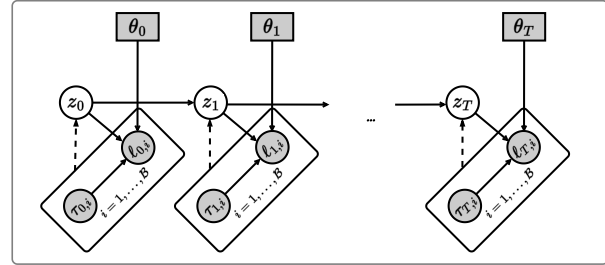


Figure 2: **Fundamental Concepts: Task Identifiers, Episodic Learning and Probabilistic Graphical Models.** **a.** The task distribution is uniform and defined over meaningful identifiers τ . For example, the amplitude and the phase $[a, b]$ specifies a sinusoid curve to complete with K-shot observed data points. Robots like Half-Cheetahs are trained to accomplish different locomotion tasks with varying masses and velocities. Some multimodal pattern recognition tasks' identifiers are implicit but can be described from a reference model, e.g., text encoders in CLIP¹. **b.** The tail task generalization corresponds to CVaR $_{\alpha}$, i.e., the integral of tail task risk values in red. In OOD generalization, this work prompt-tunes CLIP on ImageNet²⁸ to test on ImageNet-S²⁹. **c.** Here, the generative model includes grey units as observed variables and white ones as unobservable. The solid directed lines describe the *generative model*³⁰. We use the dash-directed lines to indicate the *recognition model* and approximate inference within autoencoding variational Bayes²³.

If $|\mathcal{D}_{\tau}^S| = \emptyset$, ℓ measures zero-shot adaptation; otherwise, it reflects few-shot adaptation risk. In SFT, each sample $(x, y) \in \mathcal{D}_{\text{SFT}}$ is treated as a task. The episodic task batch history is defined as $\hat{H}_t = \{\theta_t, \{(\tau_{t,i}, \mathcal{D}_{\tau_{t,i}}, \ell_{t,i})\}_{i=1}^B\}$, where B is the task batch size and θ_t represents the machine learner's parameter in t -th iteration. The tuple set $\{(\tau_{t,i}, \mathcal{D}_{\tau_{t,i}}, \ell_{t,i})\}_{i=1}^B$ includes the sampled task identifier batch $\{\tau_{t,i}\}_{i=1}^B$, the support and query dataset $\{\mathcal{D}_{\tau_{t,i}} := \mathcal{D}_{\tau_{t,i}}^S \cup \mathcal{D}_{\tau_{t,i}}^Q\}_{i=1}^B$, and the evaluated adaptation risk $\{\ell_{t,i}\}_{i=1}^B$. For simplicity, the risk history is expressed as $H_t = \{[\tau_{t,i}, \ell_{t,i}]\}_{i=1}^B$, which depends on θ_t .

Adaptation Risk Function and Robustness Concept. The learning setup optimizes the machine learner within $p(\tau)$. Our analysis is interested in the *adaptation risk* in the task space as illustrated in Fig. 1d. Such a perspective emphasizes the interplay between the task identifier τ , the task-specific dataset $\mathcal{D}_{\tau}^S \cup \mathcal{D}_{\tau}^Q$ and the adaptation risk function ℓ conditioned on θ . Regarding adaptation performance, we mainly examine zero-shot learning, few-shot learning, and SFT scenarios.

Zero-Shot Adaptation. During training, we evaluate ℓ on the query dataset \mathcal{D}_{τ}^Q conditioned on the machine learner θ , i.e., $\ell(\mathcal{D}_{\tau}^Q; \theta)$. With robotic DR⁶ as an example, $\ell(\mathcal{D}_{\tau}^Q; \theta)$ denotes the negative return of trajectories collected under the policy θ in MDP τ . This setup is without support information.

Few-Shot Adaptation. The form of ℓ is specific to meta-learning methods. For instance, MAML³² implements a bi-level optimization framework. In this case, $\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta)$ is written as $\ell(\mathcal{D}_{\tau}^Q; \theta_{\text{meta}} - \eta \nabla_{\theta} \ell(\mathcal{D}_{\tau}^S; \theta))$, where θ_{meta} denotes the meta initialization, and the inside-bracket term corresponds to finetuning θ_{meta} tailored to τ with η the learning rate.

Supervised Finetuning. The machine learner adapts to a specific downstream task using a labeled dataset. Formally, given a SFT dataset \mathcal{D}_{SFT} and a test dataset \mathcal{D}^Q , the objective is $\ell(\mathcal{D}^Q; \theta(\mathcal{D}_{\text{SFT}}, \theta_{\text{init}}))$, where θ_{init} is the pretrained initialization and θ denotes the finetuned parameter after adaptation. Unlike meta-learning, which optimizes over task distributions, SFT relies solely on the labeled samples from a specific task.

Definition 1 (Conditional Value-at-Risk, CVaR³³) Given the machine learner parameter θ , we denote the task specific random variable by $\ell_i := \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta)$. Throughout the task space \mathcal{T} , let the cumulative risk distribution and the quantile of

risk values respectively be $F(\ell)$ and $\ell^\alpha = \min_\ell \{\ell | F(\ell) \geq \alpha\}$. Then the CVaR at α -robustness level can be estimated as:

$$\text{CVaR}_\alpha[\ell(\mathcal{T}; \boldsymbol{\theta})] = \int \ell dF^\alpha(\ell; \boldsymbol{\theta}), \quad (1)$$

where we define the normalized cumulative distribution of task risk values by:

$$F^\alpha(\ell; \boldsymbol{\theta}) = \begin{cases} 0, & l < \ell^\alpha \\ \frac{F(\ell; \boldsymbol{\theta}) - \alpha}{1 - \alpha}, & l \geq \ell^\alpha. \end{cases} \quad (2)$$

And this induce the tail risk task density function denoted by $p_\alpha(\tau; \boldsymbol{\theta})$.

As illustrated in Fig. 2b, the robustness metric CVaR_α in Definition 1 is commonly used for measuring the expected risk in the worse scenarios, i.e., $1 - \alpha$ proportional tailed cases, with $\alpha \in (0, 1)$ a specific confidence level. The induced normalized distribution can be interpreted as a form of distributional shift. Beyond this notion, other robustness concepts exist, such as resilience to label noise^{34,35}, with representative approaches including ExcessMTL³⁵ that explicitly address such scenarios.

3 Results

This section reports primary findings in robust adaptation and analyzes the effect of MPTS. Prior to elaborating on the experimental setups, we outline the *predict-then-optimize* workflow underpinning MPTS. In other words, we propose a novel solution to robust adaptation: a computation- and annotation-efficient framework, as formalized in Definition 2.

Definition 2 (Model Predictive Task Sampling) *To amortize the expensive online task evaluation, MPTS reuses the risk history $H_{1:t} = \{(\tau_{i,j}, \ell_{i,j})\}_{j=1}^{\mathcal{B}}\}_{i=1}^t$ to construct the risk predictive model $p_\psi(\ell | \tau_i, H_{1:t}; \boldsymbol{\theta}_t)$ that guides the active selection of the task subset $\mathcal{T}^\mathcal{B}$ from a larger task pool $\hat{\mathcal{B}}$ for $(t+1)$ -th optimization step:*

$$\mathcal{T}_{t+1}^\mathcal{B} = \arg \max_{\mathcal{T}^\mathcal{B} \subseteq \mathcal{T}_{t+1}^\mathcal{B}, |\mathcal{T}^\mathcal{B}| = \mathcal{B}} \mathcal{A}(\mathcal{T}^\mathcal{B}, \bar{\mathcal{L}}^\mathcal{B}),$$

where $\bar{\mathcal{L}}^\mathcal{B} = \{\bar{\ell}_{t+1,i}\}_{i=1}^{\hat{\mathcal{B}}}$ denotes the predicted risk quantity of a larger task pool to roughly score task difficulties without exact evaluation at $(t+1)$ -th iteration. And $\mathcal{A}(\cdot)$ specifies certain subset acquisition rules, e.g., prioritizing tasks with worse adaptation losses for the purpose of robust optimization.

Optimization Outcome Prediction with Theoretical Guarantee and MPTS Guided Risk Minimization. First, we characterize the optimization pipeline for a typical family of robust adaptation methods, i.e., the Monte Carlo estimate for CVaR_α minimization³³:

$$\dots \xrightarrow{\text{update}} \boldsymbol{\theta}_{t-1} \xrightarrow{\text{evaluate}} \{[\hat{\tau}_{t-1,i}, \hat{\ell}_{t-1,i}]\}_{i=1}^{\hat{\mathcal{B}}} \xrightarrow{\text{Top-}\mathcal{B}} H_{t-1} := \{[\tau_{t-1,i}, \ell_{t-1,i}]\}_{i=1}^{\mathcal{B}} \xrightarrow{\text{update}} \boldsymbol{\theta}_t \xrightarrow{\text{evaluate}} \{[\hat{\tau}_{t,i}, \hat{\ell}_{t,i}]\}_{i=1}^{\hat{\mathcal{B}}} \xrightarrow{\text{Top-}\mathcal{B}} \dots, \quad (3)$$

which picks up the tail tasks to optimize in each iteration. Existing works to prioritize challenging tasks over iterations^{12–14,18} take the above steps yet suffer from: (i) learning efficiency issues, such as the need for extensive evaluation of the machine learner across numerous tasks for subset selection, and (ii) restricted batch sizes for evaluating or exploring the task space due to sample or memory constraints. Notably, nearly all of these approaches fail to leverage the optimization outcomes $H_{1:t}$.

Let us predict what to optimize from the cumulated risk episodes. MPTS differs from prior works and reuses $H_{1:t}$ to train the RPM. Coupling the identifier τ and adaptation risk $\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta})$ forms a streaming database to online learn. In **Methods**, Theorem 1 provides a provable basis for ranking tasks from predicted outcomes, suggesting stable ranking relation of task difficulties under perturbations in $\boldsymbol{\theta}$, e.g., a gradient update with a small learning rate. Thus, the candidate tasks $\mathcal{T}_{t+1}^\mathcal{B}$ at $\boldsymbol{\theta}_t$ probabilistically preserve their relative difficulty rank at $\boldsymbol{\theta}_{t+1}$. Theorem 2 further validates the convergence of the optimization pipeline. Moreover, learning adaptation risk provides a risk distribution over the task space.

MPTS surrogates CVaR_α optimization with efficiency and exploration benefits. Learning $p(\ell | \tau, H_{1:t}; \boldsymbol{\theta}_t)$ enables efficient evaluation across infinite tasks with minimal computation, expanding the pseudo batch size $\hat{\mathcal{B}}$ for subset selection and fostering exploration. As shown in **Lemma 3**, MPTS introduces external approximation bias into the gradient. For clarity, we treat MPTS as a risk minimization framework under specific acquisition criteria. As shown in Fig. 1 and Fig. 2c, its core workflow involves training the RPM $p(\ell | \tau, H_{1:t}; \boldsymbol{\theta}_t)$, evaluating task-specific adaptation risk via posterior inference, and screening task subsets using the upper confidence bound (UCB) principle³⁶ for $(t+1)$ -th optimization. These operations are formalized in Eq. (4), where the Monte Carlo estimate of the RPM yields the mean $m(\ell)$ and standard deviation $\sigma(\ell)$ of adaptation risk, while the acquisition function $\mathcal{A}(\cdot)$ quantifies total task subset risk.

$$\text{Approximate Post-Adaptation Results in Histories : } \max_{\psi \in \Psi} \mathcal{L}_{\text{ML}}(\psi) := \ln p_{\psi}(H_t | H_{1:t-1}) \quad (4a)$$

$$\text{Amortized Evaluation with RPMs : } p_{\psi}(\ell | \tau_i, H_{1:t}; \theta_t) \xrightarrow{\text{Monte Carlo Estimates}} \{m(\ell_i), \sigma(\ell_i)\}_{i=1}^{\hat{\mathcal{B}}} \quad (4b)$$

$$\text{Active Subset Selection under the UCB Principle : } \mathcal{T}_{t+1}^{\mathcal{B}} = \arg \max_{\mathcal{T}^{\mathcal{B}} \subseteq \mathcal{T}_{t+1}^{\mathcal{B}} : |\mathcal{T}^{\mathcal{B}}| = \mathcal{B}} \mathcal{A}(\mathcal{T}^{\mathcal{B}}; \phi, \psi) \quad (4c)$$

Approximating optimization outcome relies on streaming VI^{26,27}, with the RPM a lightweight model. Selecting a portion of challenging tasks to optimize, MPTS can be viewed as a biased surrogate of $\text{CVaR}_{1-\mathcal{B}/\hat{\mathcal{B}}}$ minimization while circumventing extra computations, annotations, or environment interactions. This design not only enhances learning efficiency but also aligns with the overarching goals of robust adaptation. Repeating the boxed steps of MPTS until convergence brings a robust adaptive machine learner, and the implementation details are attached in **Methods** and Algorithm 1.

Algorithm 1: Model Predictive Task Sampling

Input : Task distribution $p(\tau)$; Task batch size \mathcal{B} ; Candidate batch size $\hat{\mathcal{B}}$; Latest updated $\{\psi, \phi\}$; Latest history H_{t-1} ; Iteration number K ; Learning rate λ_2 .

Output : Selected identifier batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.

// **Train the Risk Predictive Model via Stochastic Gradient Variational Bayes**

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

 Perform gradient updates given H_{t-1} :

$\phi \leftarrow \phi + \lambda_2 \nabla_{\phi} \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (17b);

$\psi \leftarrow \psi + \lambda_2 \nabla_{\psi} \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (17b);

end

// **Simulate Zero-shot, Few-shot Adaptation and SFT Results with Trained RPMs**

Randomly sample $\{\hat{\tau}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ from $p(\tau)$;

Run amortized evaluation on candidate tasks $\{\delta_i := \gamma_0 m(\ell_i) + \gamma_1 \sigma(\ell_i)\}_{i=1}^{\hat{\mathcal{B}}}$ in Eq. (19);

// **Active Subset Selection from Predicted Results**

Rank $\{\delta_i\}_{i=1}^{\hat{\mathcal{B}}}$ and screen Top- \mathcal{B} values;

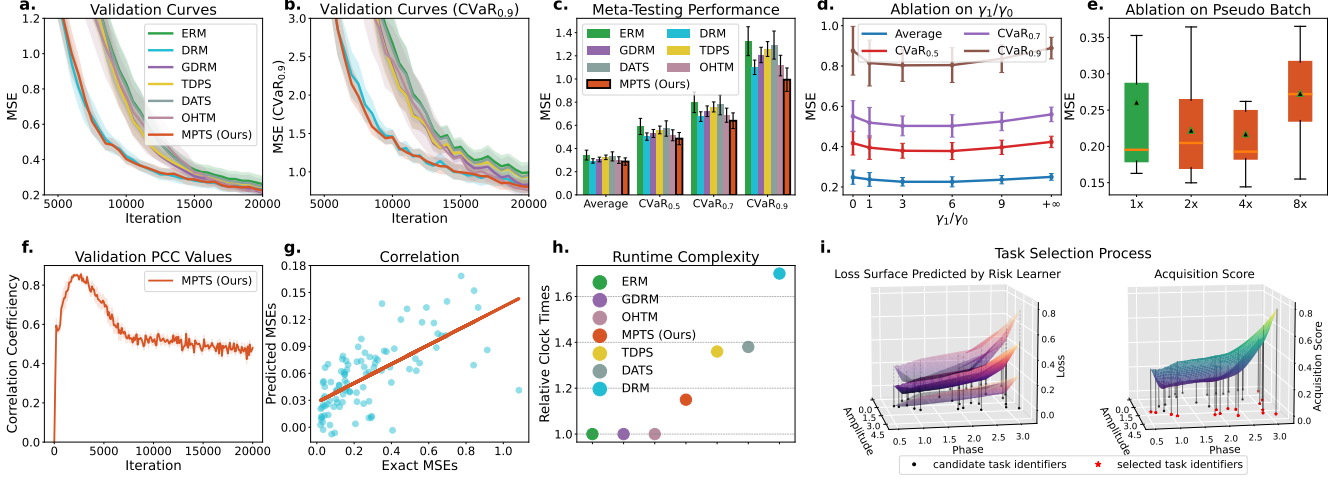
Return the screened identifier subset $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.

Adaptation Learning Benchmark. The experimental design considers the benchmark typicality and the practical challenges. Downstream tasks span pattern recognition and sequential decision-making, with certain experiments involving multimodal foundation models. These experiments mainly examine few-shot adaptation and include (1) K-shot sinusoid regression³², (2) N-way K-shot image classification^{37,38} with CLIP¹ models, and (3) Meta-RL³². Additionally, MPTS validates scenarios like (4) robotic DR³⁹ for zero-shot adaptation and (5) SFT CLIP models towards image classification.

Backbones and Task Robust Baselines. This study primarily compares MPTS with risk minimization principles and focuses on robustness improvement. While these methods, including MPTS, are agnostic to zero-shot, few-shot learning, or finetuning techniques, we adopt SoTA backbones for experiments. For sinusoid regression and Meta-RL, we use MAML³² as default. As CLIP has strong zero-shot performance, we extend it with MaPLe⁴⁰ for N-way K-shot image classification. For robotic DR in Ergo-Reacher and Lunar-Lander³⁹, we employ TD3⁴¹ due to its stability. In SFT, we again use MaPLe for prompt-tuning in image classification.

Baselines include Empirical Risk Minimization (ERM)⁴², Distributionally Robust Risk Minimization (DRM)^{14,16,17,43,44}, and Group Distributionally Robust Risk Minimization (GDRM)⁴⁵⁻⁴⁷. Meanwhile, for certain scenarios, we further include some robust optimization methods such as Difficulty-Aware Task Sampler (DATS)⁴⁸, Online Hard Task Mining Sampler (OHTM)⁴⁹, and Task Difficulty Prioritized Sampler (TDPS)⁵⁰ for comparison. Accordingly, adaptation robustness is evaluated via CVaR_{α} across validation/testing tasks with $\alpha = \{0.9, 0.7, 0.5\}$, including some OOD results. We also compare computational cost, memory usage, and sample efficiency. For fairness, all baselines share the same task batch \mathcal{B} in optimization, excluding pruned easier tasks. ERM, GDRM, DATS, TDPS, and OHTM use batch size \mathcal{B} , while DRM samples $\hat{\mathcal{B}} = 2\mathcal{B}$, filtering half for stable optimization. See Supplementary Notes F/G for details.

Meta Learning Backbone: MAML



Meta Learning Backbone: Reptile

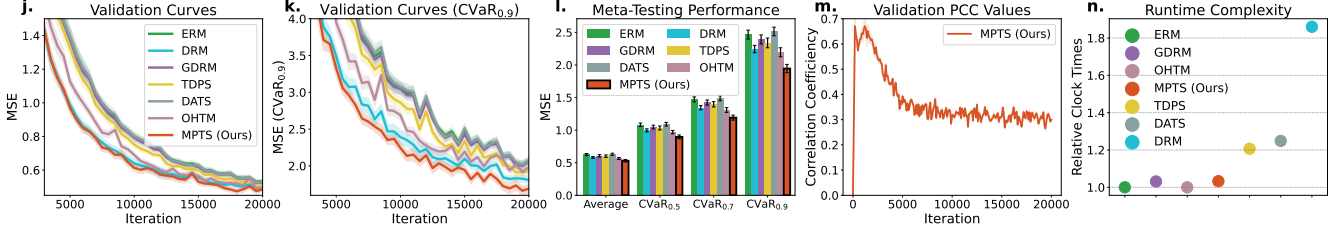


Figure 3: K-shot **Sinusoid Regression Results (10 Runs)**. Note that **a-h** are results with MAML as the backbone, while **i-l** reports the results with Reptile as the backbone. **a.** Shown are curves of averaged MSEs on the validation task set during meta-training for all methods. **b.** Curves illustrate the CVaR_{0.9} MSEs on the validation task set during the meta-training process. **c.** The meta-trained machine learners are tested on a fixed task set, reporting the average MSEs and CVaR values. **d.** Displayed are meta-testing results with MPTS machine learners trained by various γ_1/γ_0 ratios. **e.** Displayed are meta-testing results with MPTS machine learners trained in various pseudo batch sizes, i.e., $\hat{B} = \{1B, 2B, 4B, 8B\}$. **f.** The PCC values are tracked during meta-training. **g.** At a specific iteration, the statistical correlation between predicted and exact adaptation risk values of the task batch is visualized with overall $\rho_{\bar{\ell}, \ell} = 0.669$. **h.** The required relative run-time is computed for all methods during meta-training with ERM as the anchor. **i.** At some meta-training time-step, we visualize the subset selection from the pseudo batch under the RPM. **j.** We illustrate the curves of averaged MSEs on the validation task set during Reptile meta-training for all relevant baselines. **k.** We track the corresponding CVaR_{0.9} MSEs on the validation task set throughout the Reptile meta-training. **l.** Reported are the tested average MSEs and CVaR values of the meta-trained machine learners on a fixed task set. **m.** The PCC values are tracked for MPTS during Reptile meta-training. **n.** We compute the relative run-time for all methods during meta-training with ERM as the anchor.

3.1 Demonstration of the MPTS's role in K-shot sinusoid regression

In K-shot sinusoid regression³², the function family $\{f(x) = a_i \sin(x - b_i) | (a_i, b_i) \in [0.1, 5.0] \times [0.0, \pi]\}$ is specified by the identifier $\tau = [a, b]$. This serves as a toy case to illustrate MPTS and the role of the RPM.

The RPM allows for roughly scoring the task difficulty over iterations. In MPTS, for the screened subset at $(t + 1)$ -th iteration, we track the predicted Top- B risk values $\{\bar{\ell}_{t+1,i} \approx \mathbb{E}_{q_\phi(z_t|H_t)}[p_\psi(\ell|\tau_{t+1,i}, H_{1:t})]\}_{i=1}^B$ and corresponding exact evaluations $\{\ell_{t+1,i}\}_{i=1}^B$ from θ_{t+1} to compute the Pearson correlation coefficient (PCC) $\rho_{\bar{\ell}, \ell} := \frac{\sum_{i=1}^B (\bar{\ell}_{t+1,i} - \text{Mean}[\{\bar{\ell}_{t+1,\cdot}\}]) (\ell_{t+1,i} - \text{Mean}[\{\ell_{t+1,\cdot}\}])}{\sqrt{\sum_{i=1}^B (\bar{\ell}_{t+1,i} - \text{Mean}[\{\bar{\ell}_{t+1,\cdot}\})^2} \sqrt{\sum_{i=1}^B (\ell_{t+1,i} - \text{Mean}[\{\ell_{t+1,\cdot}\})^2}}$. For continuous risk values, PCC reasonably quantifies the correlation effect of ranking in a batch. The RPM amortizes the exact evaluation $\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta_t) \forall \tau \in \mathcal{T}$ and $\theta_t \in \Theta$ using risk histories, roughly scoring adaptation difficulty for next iteration. Only if the RPM approximately ranks tasks, MPTS can trust amortized evaluations for worst subset selection.

As shown in Fig. 3f, $\rho_{\bar{\ell}, \ell}$ remains between 0.4–0.8 across iterations, validating the reliability of the RPM in predicting adaptation difficulty. However, PCC declines over time—a trend also observed across experiments—likely due to model θ_t 's convergence.

This reduces task diversity, negatively affecting the RPM’s training after local task space overoptimization. Fig. 3g shows the statistical correlation between predicted and exact adaptation risk at a specific iteration. Scattered points demonstrate strong overall alignment, despite varying value scales between iterations. Notably, difficult tasks with high mean square errors (MSEs) are well identified and clustered around or above the correlation slope along the x -axis.

MPTS accelerates the learning process and improves comprehensive adaptation performance under active sampling. In Fig. 3a-b, MPTS converges faster in average and $\text{CVaR}_{0.9}$ MSEs and reaches performance comparable to other baselines with less iterations, e.g., 15K steps, due to its uncertainty-guided worst-case acquisition. DRM processes $2\mathcal{B}$ tasks to filter half per iteration, raising $0.7\times$ computational overhead over ERM (see Fig. 3h). In contrast, MPTS incurs only $0.14\times$ runtime increase, a marginal overhead. To illustrate active task sampling, Fig. 3i visualizes predicted risk values over the task space. Selected tasks favor regions with high deviations, clustering in high-risk areas.

In meta-testing, Fig. 3c shows that MPTS, OHTM, and DRM achieve the lowest average and CVaR_α MSEs, with their advantage over other baselines increasing at α . DATS and TDPS require more computational resources to achieve slight gains in robustness metrics. Prior work⁴⁴ confirms that DRM’s efficiency is sacrificed for robustness, relying on intensive task evaluation. Using MAML, gradient-based inner-loop adaptation further increases overhead, whereas MPTS bypasses it via probabilistic prediction, reducing computational cost.

The appropriate hyper-parameter setup secures performance and efficiency. We first analyze the acquisition function $\mathcal{A}(\mathcal{T}^\mathcal{B}; \phi, \psi)$ by varying trade-off parameters $\{\gamma_0, \gamma_1\}$ in Fig. 1d and Eq. (19). Meta-testing machine learners trained with $\frac{\gamma_1}{\gamma_0} = \{1.0, 3.0, 6.0, 9.0\}$, $\gamma_0 = 0.0$ and $\gamma_1 = 0.0$ shows that higher uncertainty weights lower average MSEs (Fig. 3d). However, removing worst-case considerations ($\gamma_0 = 0.0$) weakens performance. We further examine the impact of pseudo batch size $\hat{\mathcal{B}}$ in Fig. 3e. Increasing $\hat{\mathcal{B}}$ reduces average MSEs, but excessively large values (e.g., $\hat{\mathcal{B}} = 8 \times \mathcal{B}$) degrade performance. This occurs because an enlarged identifier set under worst-case selection might over-optimize local task regions, hindering global generalization. Thus, MPTS configuration follows two principles: (i) $\hat{\mathcal{B}}$ should be moderate to encourage exploration while preventing excessive local optimization. (ii) Since adaptation robustness is the priority, we consistently set $\gamma_0 \in \mathbb{R}_+$ as the default in all experiments.

MPTS is plug-and-play and agnostic to the backbone. We further conduct experiments using another SoTA meta-optimizer Reptile⁵¹ as the backbone. From Fig. 3j-k, the convergence trends of most algorithms are similar to those observed under MAML in Fig. 3a-b. OHTM lags slightly behind MPTS, and the overall MSEs are higher. This observation is corroborated by the meta-testing results in Fig. 3l, where the secondary method OHTM outperforms the remaining baselines with a lower computational cost. MPTS continues to exhibit a positive correlation between predicted and exactly evaluated adaptation results (Fig. 3m). The order of computational complexity across all methods remains the same as under MAML, though the relative scale in Fig. 3n varies slightly.

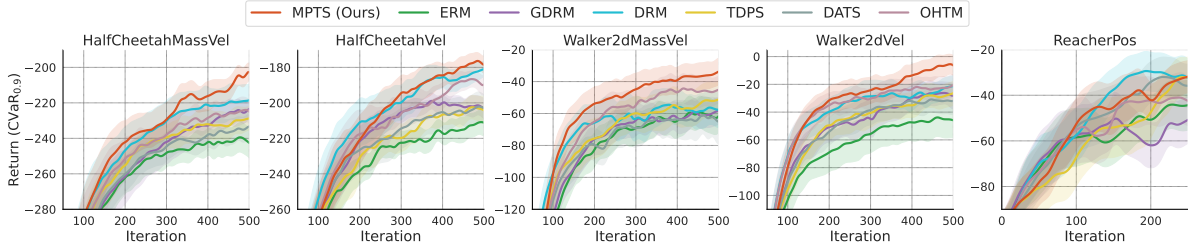
3.2 Few-Shot adaptation benefits from MPTS in robustness and learning efficiency

Result analysis in N-way K-shot image classification. We perform 5-way 1-shot image classification using MaPLE, with six meta-training datasets from ImageNet-CG⁵², ImageNet-CI⁵², ImageNet-CS⁵², ImageNet-A⁵³, ImageNet-S²⁹ and ImageNet-R⁵⁴. The efficiency row in Table 1 compares computational time and memory usage across methods during meta-training. The overhead from optimizing RPMs in MPTS is negligible, whereas DRM incurs $1.3\times$ computational time and $1.6\times$ memory usage relative to ERM. DATS and TDPS also significantly suffer from more computational complexity issues due to the additional gradient computations required to derive task-specific weights. As a memory retrieval strategy, OHTM occupies $2.54\times$ memory to prioritize difficult tasks during optimization.

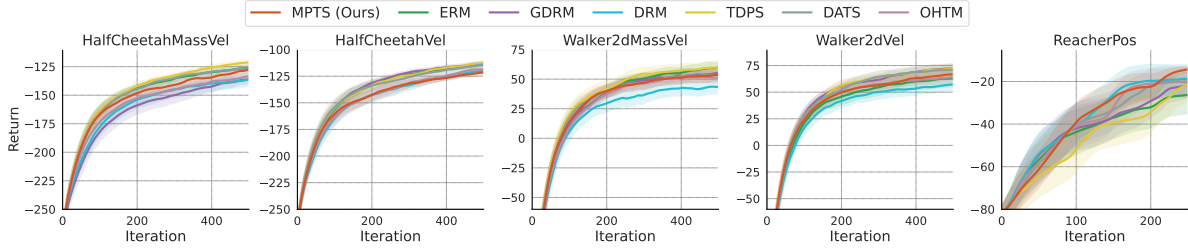
In Table 1, MPTS consistently surpasses risk-aware baselines and achieves significant performance gains, i.e., 0.36-2.32% accuracy increase, over ERM across all benchmark datasets. The improvement is particularly prominent under the stringent $\text{CVaR}_{0.9}$ metric, e.g., +2.32% on ImageNet-CS and +2.3% on ImageNet-S compared to ERM. In contrast, DRM provides moderate improvements after sacrificing more computations during meta-training. There is no evident runner-up. OHTM, DATS, and TDPS, suffer from diminishing average and CVaR accuracies on ImageNet-CI while mostly beating ERM in other cases. Overall, this benchmark result demonstrates the comprehensive merits of prioritizing challenging tasks, such as optimizing worst-case performance, which surprisingly improves the average results. This contradicts the following conclusion in Meta-RL yet is consistent with the consequence of prioritizing complex tasks during LLM pretraining¹³.

Result analysis in Meta-RL. We first analyze meta-training results in Fig. 4a-b. MPTS achieves the highest $\text{CVaR}_{0.9}$ validation returns on most benchmarks. DRM sacrifices average returns on HalfCheetahMassVel, HalfCheetahVel, and Walker2dVel, whereas MPTS maintains average performance comparable to ERM on HalfCheetahMassVel and Walker2dVel. GDRM, OHTM, DATS, and TDPS are entangled in learning curves and behave intermediate performance, while DRM balances

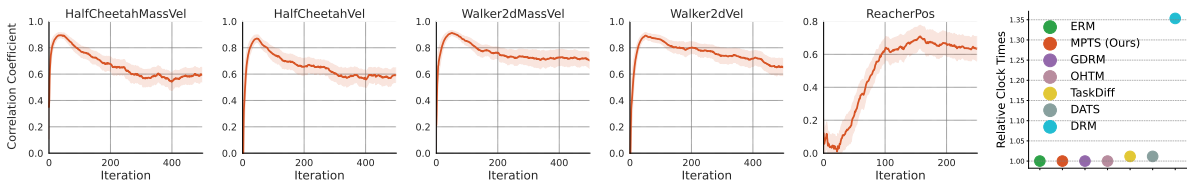
a. Validation $\text{CVaR}_{0.9}$ Task Returns during Meta-Training



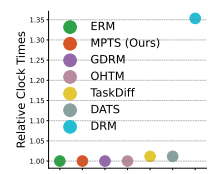
b. Validation Average Task Returns during Meta-Training



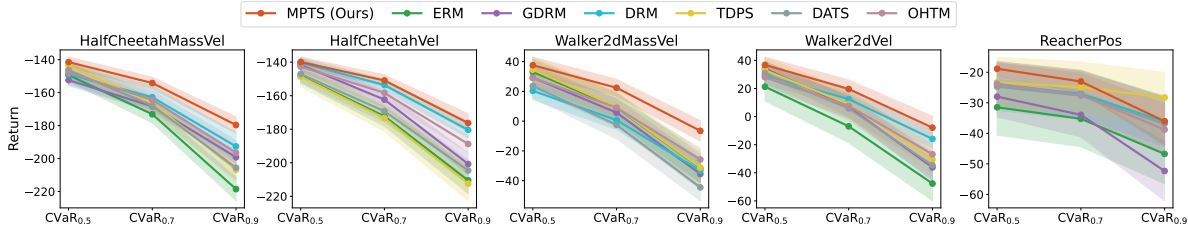
c. Tracked PCC Values in Task Batches during Meta-Training



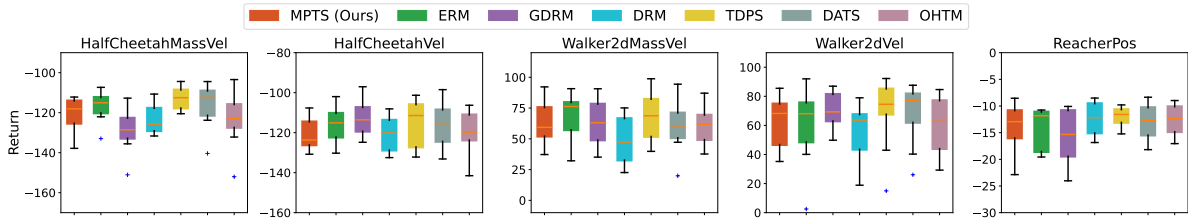
d. Runtime Complexity



e. CVaR_α Task Returns during Meta-Testing



f. Average Task Returns during Meta-Testing



g. Compatible with Other Meta Learning Backbones (e.g., PEARL)

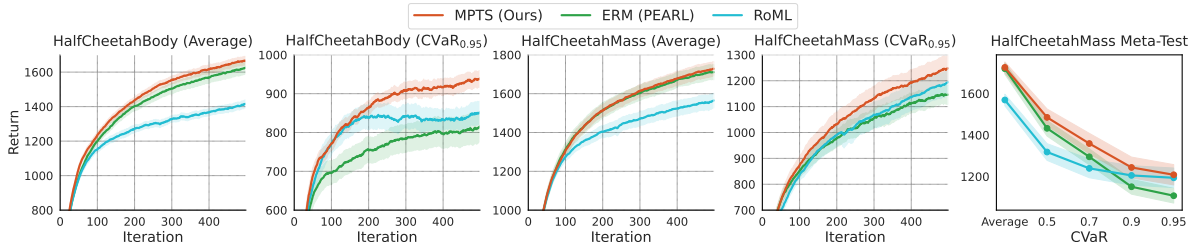


Figure 4: **Meta-RL Results on Mujoco Environments (10 Runs).** **a.** The cumulative returns with standard error of means (SEMs) belonging to $\text{CVaR}_{0.9}$ validation MDPs are displayed during meta-training. **b.** We compute the average cumulative returns with SEMs on validation MDPs during meta-training. **c.** Tracked are the RPM’s PCC values with SEMs over training iterations. **d.** The relative clock time quantifies the computational complexity for all methods on Walker2dVel, where ERM’s runtime works as the anchor. **e.** We report CVaR_α returns of meta-testing MDPs. **f.** The box-plot reports results averaged over meta-testing MDPs. **g.** With PEARL⁵⁵ as the Meta-RL backbone, we illustrate the learning curves and meta-testing results on HalfCheetahBody and HalfCheetahMass from RoML⁵⁶ baseline.

Table 1: 5-Way 1-Shot **Meta-Testing Classification Results on Various Datasets and Efficiency Comparison (10 Runs)**. We report testing $\text{CVaR}_{0.9}$, $\text{CVaR}_{0.7}$, $\text{CVaR}_{0.5}$ and average accuracies with corresponding SEMs evaluated by the meta-trained machine learner on ID and OOD datasets. The **best results** are in bold with the runner-up underlined, and MPTS’s performance gains over ERM $\Delta \uparrow$ are marked in blue. With experiments on ImageNet-A as an example, we report the memory cost and runtime relative to ERM during meta-training.

Dataset	Metrics	ERM	DRM	GDRM	OHTM	DATS	TDPS	MPTS (Ours)	$\Delta \uparrow$
ImageNet-CG ⁵²	$\text{CVaR}_{0.9}$	76.87 \pm 0.31	77.47 \pm 0.24	<u>77.56\pm0.19</u>	76.52 \pm 0.14	77.55 \pm 0.21	77.28 \pm 0.17	78.16\pm0.25	+1.29
	$\text{CVaR}_{0.7}$	82.04 \pm 0.21	82.47 \pm 0.20	<u>82.62\pm0.20</u>	81.34 \pm 0.17	82.31 \pm 0.18	82.33 \pm 0.24	82.95\pm0.26	+0.91
	$\text{CVaR}_{0.5}$	84.62 \pm 0.24	84.99 \pm 0.17	85.16 \pm 0.22	84.06 \pm 0.29	85.19 \pm 0.14	<u>85.29\pm0.23</u>	85.33\pm0.26	+0.71
	Avg	89.22 \pm 0.22	89.54 \pm 0.16	89.49 \pm 0.18	88.59 \pm 0.17	89.45 \pm 0.12	<u>89.69\pm0.11</u>	89.86\pm0.14	+0.64
ImageNet-CI ⁵²	$\text{CVaR}_{0.9}$	80.40 \pm 0.18	<u>80.49\pm0.22</u>	80.26 \pm 0.19	77.08 \pm 0.18	78.86 \pm 0.15	78.07 \pm 0.19	80.94\pm0.19	+0.54
	$\text{CVaR}_{0.7}$	84.98 \pm 0.17	<u>85.05\pm0.19</u>	84.90 \pm 0.22	81.60 \pm 0.11	83.66 \pm 0.22	83.16 \pm 0.11	85.70\pm0.24	+0.72
	$\text{CVaR}_{0.5}$	87.01 \pm 0.20	<u>87.47\pm0.26</u>	87.04 \pm 0.18	84.14 \pm 0.22	86.08 \pm 0.16	85.75 \pm 0.13	87.74\pm0.15	+0.73
	Avg	91.20 \pm 0.23	<u>91.40\pm0.18</u>	91.21 \pm 0.17	88.81 \pm 0.23	90.97 \pm 0.13	90.67 \pm 0.18	91.56\pm0.12	+0.36
ImageNet-CS ⁵²	$\text{CVaR}_{0.9}$	76.15 \pm 0.20	77.78 \pm 0.24	76.25 \pm 0.15	<u>77.84\pm0.23</u>	77.77 \pm 0.16	77.47 \pm 0.16	78.47\pm0.23	+2.32
	$\text{CVaR}_{0.7}$	81.55 \pm 0.19	82.64 \pm 0.16	81.42 \pm 0.17	82.59 \pm 0.20	<u>82.81\pm0.19</u>	82.31 \pm 0.17	83.13\pm0.16	+1.58
	$\text{CVaR}_{0.5}$	84.44 \pm 0.19	85.16 \pm 0.25	84.19 \pm 0.15	85.23 \pm 0.18	84.67 \pm 0.20	84.84 \pm 0.13	85.87\pm0.22	+1.43
	Avg	89.30 \pm 0.20	89.80 \pm 0.21	89.11 \pm 0.17	<u>89.87\pm0.25</u>	89.60 \pm 0.16	89.82 \pm 0.20	90.28\pm0.11	+0.98
ImageNet-A ⁵³	$\text{CVaR}_{0.9}$	76.71 \pm 0.22	77.36 \pm 0.26	77.45 \pm 0.22	76.81 \pm 0.18	<u>77.63\pm0.21</u>	77.12 \pm 0.20	77.96\pm0.11	+1.25
	$\text{CVaR}_{0.7}$	82.03 \pm 0.25	82.52 \pm 0.27	82.37 \pm 0.20	<u>82.72\pm0.10</u>	82.69 \pm 0.20	82.04 \pm 0.23	83.43\pm0.18	+1.40
	$\text{CVaR}_{0.5}$	84.50 \pm 0.22	85.15 \pm 0.19	85.34 \pm 0.25	85.20 \pm 0.23	<u>85.77\pm0.19</u>	84.92 \pm 0.20	86.74\pm0.24	+2.24
	Avg	89.38 \pm 0.21	90.31 \pm 0.17	90.29 \pm 0.17	<u>90.68\pm0.20</u>	90.58 \pm 0.18	90.03 \pm 0.23	91.10\pm0.21	+1.72
ImageNet-S ²⁹	$\text{CVaR}_{0.9}$	82.58 \pm 0.21	83.87 \pm 0.10	83.18 \pm 0.22	82.23 \pm 0.18	83.60 \pm 0.19	<u>83.93\pm0.15</u>	84.88\pm0.14	+2.30
	$\text{CVaR}_{0.7}$	87.23 \pm 0.23	<u>88.45\pm0.21</u>	87.32 \pm 0.14	87.08 \pm 0.13	88.32 \pm 0.19	88.06 \pm 0.16	89.27\pm0.12	+2.04
	$\text{CVaR}_{0.5}$	89.56 \pm 0.19	89.96 \pm 0.22	90.06 \pm 0.18	90.01 \pm 0.22	<u>91.16\pm0.22</u>	90.65 \pm 0.20	91.52\pm0.19	+1.96
	Avg	93.52 \pm 0.13	<u>94.28\pm0.10</u>	94.13 \pm 0.11	94.19 \pm 0.15	94.15 \pm 0.13	94.18 \pm 0.14	94.84\pm0.08	+1.32
ImageNet-R ⁵⁴	$\text{CVaR}_{0.9}$	88.28 \pm 0.13	88.89 \pm 0.15	87.33 \pm 0.27	88.73 \pm 0.16	<u>89.20\pm0.21</u>	88.52 \pm 0.20	89.78\pm0.16	+1.50
	$\text{CVaR}_{0.7}$	91.47 \pm 0.13	<u>91.95\pm0.16</u>	91.11 \pm 0.16	91.41 \pm 0.19	<u>91.95\pm0.23</u>	90.89 \pm 0.20	93.27\pm0.19	+1.80
	$\text{CVaR}_{0.5}$	93.42 \pm 0.12	<u>93.87\pm0.17</u>	93.20 \pm 0.20	92.75 \pm 0.15	93.27 \pm 0.21	92.48 \pm 0.18	94.69\pm0.18	+1.27
	Avg	96.02 \pm 0.14	<u>96.39\pm0.16</u>	95.90 \pm 0.18	95.32 \pm 0.17	95.85 \pm 0.16	95.64 \pm 0.17	96.91\pm0.12	+0.89
Efficiency	Runtime	1.00	1.33	1.00	1.09	1.68	1.68	1.01	
	Memory	1.00	1.61	1.00	2.54	1.99	1.99	1.01	

average and $\text{CVaR}_{0.9}$ returns, excelling on ReacherPos. Fig. 4c witnesses the RPM’s strong task difficulty discrimination capability, measured by $\rho_{\bar{\ell}, \ell}$. In Fig. 4d, DRM consumes 1.35 \times runtime on Walker2dVel due to extra environment interactions, while MPTS avoids this inefficiency.

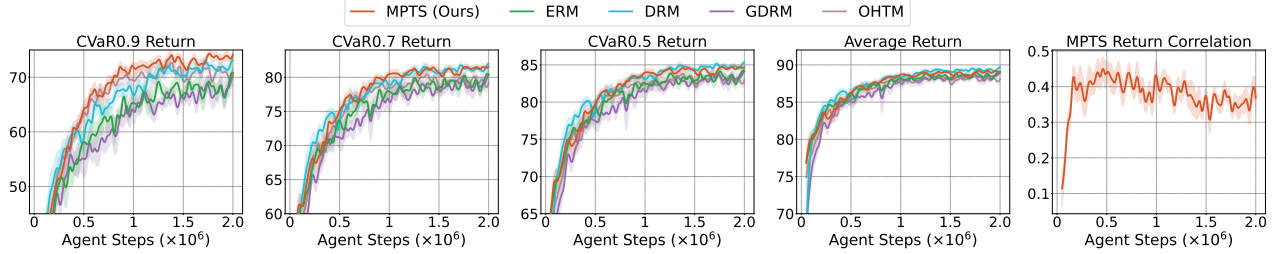
Meta-testing results in Fig. 4e-f highlight MPTS’s robustness, with more return gains at higher α values. With increase of α in CVaR_{α} , performance difference across methods grows significant. In extreme cases, i.e., $\text{CVaR}_{0.9}$, MPTS surpasses ERM by over 20% on all benchmarks, and DRM mostly dominates other baselines. Average performance varies: Walker2dMassVel and Walker2dVel show minor differences, while HalfCheetahMassVel favors MPTS with slightly higher variance. HalfCheetah marginally benefits GDRM and ERM, whereas ReacherPos favors MPTS and DRM with reduced variance. Overall, methods except DRM are mostly entangled in average performance, while DRM tends to sacrifice computational efficiency for adaptation robustness increase. The primary goal of robust optimization is to enhance performance in tail risk or OOD tasks. Hence, methods often trade off worst-case and average performance in Meta-RL, as implied in work¹⁷.

Moreover, Fig. 4g reports the performance of MPTS with PEARL⁵⁵ on environments and configurations adopted from RoML¹⁷, a representative robust Meta-RL algorithm. The learning curves verify the trade-off between average and worst-case performance in ERM and RoML. In contrast, MPTS maintains average returns while substantially outperforming RoML in $\text{CVaR}_{0.95}$, a behavior distinct from the MAML backbone results. The meta-testing results on HalfCheetahMass corroborate this trend, consistent with the meta-training phase, where MPTS outperforms ERM by over 9% in $\text{CVaR}_{0.95}$ and exceeds RoML in average return by approximately 10%.

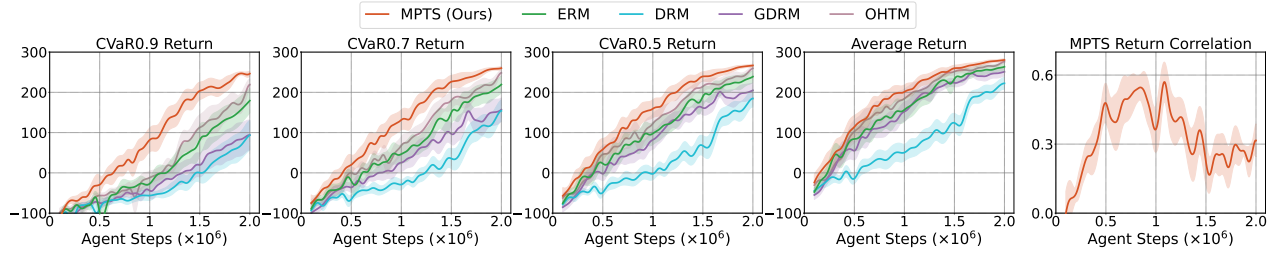
3.3 MPTS retains multi-faced advantages beyond robustness in zero-shot continuous control

MPTS dominates the overall performance in DR training. In Fig. 5a Ergo-Reacher, distinguished from Meta-RL finding, MPTS and DRM improve both average and CVaR_{α} performance. This likely stems from MPTS’s broader task exploration via larger \hat{B} . Meanwhile, $\rho_{\bar{\ell}, \ell}$ fluctuates near 0.4 throughout training. In Fig. 5b Lunar-Lander, MPTS maintains the leading trend in average and CVaR_{α} returns, followed by OHTM. In contrast, DRM and GDRM not only underperform in average returns but also achieve the lowest $\text{CVaR}_{0.9}$ values, failing in robust optimization. Naively worst-case selection or reweighting

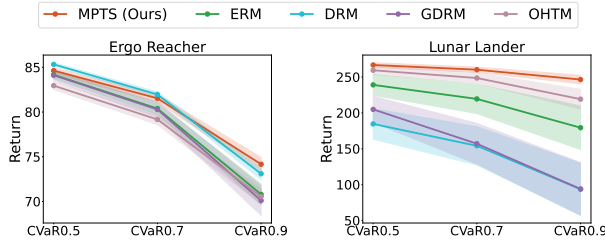
a. Validation Task Returns and Tracked PCC Values during DR-Training on Ergo-Reacher



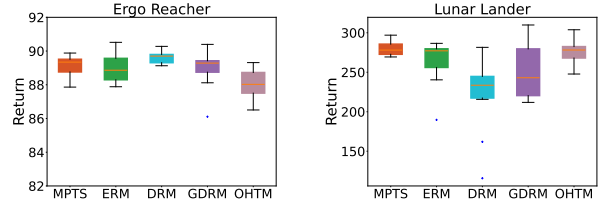
b. Validation Task Returns and Tracked PCC Values during DR-Training on Lunar-Lander



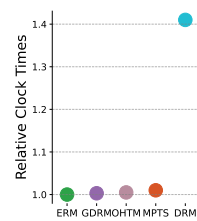
c. CVaR Task Returns during DR-Testing



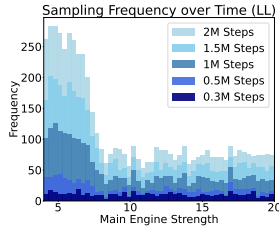
d. Average Task Returns during DR-Testing



e. Runtime Complexity



f. Task Statistics on Lunar-Lander



g. ID and OOD Task Returns during DR-Testing on Lunar-Lander

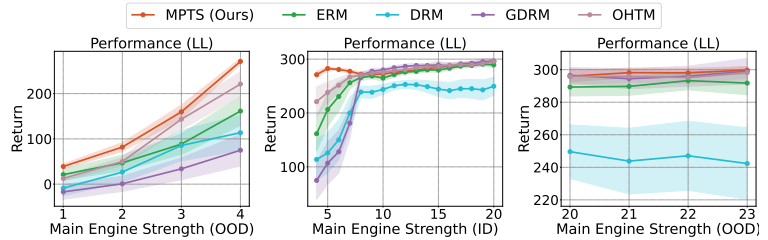


Figure 5: DR Results on Ergo-Reacher and Lunar-Lander (10 Runs). **a.** In Ergo-Reacher, the $\text{CVaR}_{0.9}$, $\text{CVaR}_{0.7}$, $\text{CVaR}_{0.5}$ and average cumulative returns on validation MDPs are reported together with the RPM’s PCC curve during DR training. **b.** In Lunar-Lander, the cumulative returns on validation MDPs are illustrated together with the RPM’s PCC curve during DR training. **c.** We test the DR-trained policies on the fixed MDP set and report the CVaR_α cumulative returns. **d.** The returns averaged over DR-testing MDPs are illustrated. **e.** The required runtime is computed for all methods on Lunar-Lander. **f.** In Lunar-Lander, shown are frequencies of sampled identifiers using MPTS during DR training. **g.** In Lunar-Lander, we test the trained policies in both in-distribution (ID) domains and out-of-distribution (OOD) domains to report each task’s average returns.

tends to degrade performance when unsolvable tasks are frequently sampled. MPTS mitigates this by balancing worst-case and uncertainty-guided selection, preventing over-optimization on a finite set of difficult MDPs. Here, $\rho_{\ell, \ell}$ peaks above 0.6 before stabilizing near 0.3, consistent with prior findings that task selection converges, reducing task difficulty discrimination. MPTS’s runtime in Lunar-Lander is comparable to ERM and GDRM in Fig. 5e. In Lunar-Lander, the identifier $\tau \in \mathbb{R}_+$ represents the main engine strength. Fig. 5f shows task sampling frequency, where MPTS favors lower-engine-strength MDPs while still exploring all engine-strengths early in training.

MPTS facilitates policy adaptation in the worst-case and OOD MDPs. For DR-trained policies, Fig. 5c-d confirm MPTS and DRM’s superior CVaR_α returns in Ergo-Reacher, while ERM exhibits a minor dip in average returns. In Lunar-Lander, MPTS attains the highest CVaR_α returns, remaining stable even as α increases—outperforming ERM by over 20%. Additionally, MPTS, OHTM and ERM yield top average returns with minimal variance. For OOD generalization, we shift τ ’s range from training interval $\tau \in [4.0, 20.0]$ to testing interval $\tau \in [1.0, 4.0] \cup (20.0, 23.0]$. All methods struggle in hard OOD tasks (Fig. 5g left), but MPTS dominates in difficult cases, demonstrating strong adaptation. DRM exhibits high variability and weak generalization, even for easier tasks (Fig. 5g right).

3.4 MPTS also reserves the potential of robust SFT

Table 2: **Testing Classification Results after Prompt-Tuning on ImageNet and Efficiency Comparison (10 Runs).** We report testing $\text{CVaR}_{0.9}$, $\text{CVaR}_{0.7}$, $\text{CVaR}_{0.5}$ and average accuracies together with corresponding SEMs evaluated by the prompt-tuned machine learner on ID and OOD datasets. Evaluation on OOD datasets corresponds to the domain generalization setting. The **best results** are in bold with the runner-up underlined, and MPTS’s performance gains over ERM $\Delta \uparrow$ are marked in blue. During prompt-tuning ImageNet, we report the memory cost and runtime relative to ERM for all methods.

Dataset	Metrics	ERM	DRM	GDRM	OHTM	MPTS (Ours)	$\Delta \uparrow$
ImageNet ²⁸ (ID)	$\text{CVaR}_{0.9}$	31.70 \pm 0.12	<u>32.23\pm0.15</u>	31.37 \pm 0.18	31.53 \pm 0.21	32.52\pm0.16	+0.82
	$\text{CVaR}_{0.7}$	42.87 \pm 0.14	<u>44.10\pm0.14</u>	42.97 \pm 0.18	43.28 \pm 0.18	44.28\pm0.14	+1.41
	$\text{CVaR}_{0.5}$	51.43 \pm 0.10	<u>52.61\pm0.11</u>	51.75 \pm 0.19	51.96 \pm 0.16	52.72\pm0.15	+1.29
	Avg	70.80 \pm 0.08	<u>70.90\pm0.10</u>	<u>71.00\pm0.08</u>	70.80 \pm 0.10	71.20\pm0.09	+0.40
ImageNet-A ⁵³ (OOD)	$\text{CVaR}_{0.9}$	15.35 \pm 0.21	15.50 \pm 0.25	<u>15.62\pm0.22</u>	15.44 \pm 0.25	18.38\pm0.22	+3.03
	$\text{CVaR}_{0.7}$	22.82 \pm 0.23	23.05 \pm 0.21	<u>23.13\pm0.19</u>	23.00 \pm 0.23	24.05\pm0.20	+1.23
	$\text{CVaR}_{0.5}$	30.05 \pm 0.21	29.53 \pm 0.18	<u>30.25\pm0.20</u>	29.96 \pm 0.23	31.26\pm0.20	+1.21
	Avg	<u>49.88\pm0.18</u>	48.42 \pm 0.25	49.53 \pm 0.14	49.70 \pm 0.25	51.10\pm0.18	+1.22
ImageNet-R ⁵⁴ (OOD)	$\text{CVaR}_{0.9}$	26.20 \pm 0.11	28.01 \pm 0.13	26.01 \pm 0.10	27.87 \pm 0.14	28.21\pm0.12	+2.01
	$\text{CVaR}_{0.7}$	43.55 \pm 0.12	<u>45.25\pm0.11</u>	43.95 \pm 0.11	44.89 \pm 0.15	45.50\pm0.10	+1.95
	$\text{CVaR}_{0.5}$	56.70 \pm 0.11	<u>58.25\pm0.12</u>	57.36 \pm 0.11	58.03 \pm 0.12	58.77\pm0.12	+2.07
	Avg	76.91 \pm 0.08	<u>77.40\pm0.10</u>	77.30 \pm 0.10	77.12 \pm 0.11	77.60\pm0.09	+0.69
ImageNet-S ²⁹ (OOD)	$\text{CVaR}_{0.9}$	12.24 \pm 0.20	<u>13.18\pm0.20</u>	12.20 \pm 0.19	12.73 \pm 0.19	13.61\pm0.17	+1.37
	$\text{CVaR}_{0.7}$	20.01 \pm 0.18	<u>21.10\pm0.22</u>	20.44 \pm 0.20	20.65 \pm 0.25	21.46\pm0.21	+1.45
	$\text{CVaR}_{0.5}$	26.70 \pm 0.20	<u>27.44\pm0.20</u>	27.42 \pm 0.25	27.40 \pm 0.26	27.98\pm0.20	+1.28
	Avg	48.85 \pm 0.23	48.89 \pm 0.21	48.91 \pm 0.23	<u>49.12\pm0.20</u>	49.62\pm0.22	+0.77
ImageNet-V ⁵⁷ (OOD)	$\text{CVaR}_{0.9}$	24.65 \pm 0.27	<u>25.52\pm0.20</u>	24.17 \pm 0.18	24.82 \pm 0.17	25.95\pm0.23	+1.30
	$\text{CVaR}_{0.7}$	34.72 \pm 0.20	<u>35.70\pm0.17</u>	35.00 \pm 0.23	35.32 \pm 0.25	35.86\pm0.16	+1.14
	$\text{CVaR}_{0.5}$	43.49 \pm 0.24	<u>44.20\pm0.19</u>	43.60 \pm 0.22	43.95 \pm 0.18	44.60\pm0.20	+1.11
	Avg	64.02 \pm 0.22	63.87 \pm 0.19	<u>64.10\pm0.20</u>	<u>64.10\pm0.17</u>	64.55\pm0.16	+0.53
Efficiency	Runtime	1.00	1.10	1.00	1.09	1.03	
	Memory	1.00	1.50	1.00	2.08	1.01	

In SFT, each labeled example in the dataset can be viewed as a task. Following MaPLE, we execute prompt tuning on ImageNet²⁸ and conduct standard evaluation. To assess post-SFT robustness, we test on four OOD datasets—ImageNet-A⁵³, ImageNet-S²⁹, ImageNet-R⁵⁴, and ImageNet-V⁵⁷ for capturing diverse domain shifts.

Table 2 shows MPTS consistently outperforms baselines in average and CVaR accuracies on ID and OOD datasets. MPTS achieves 0.82–3.03% higher $\text{CVaR}_{0.9}$, $\text{CVaR}_{0.7}$ and $\text{CVaR}_{0.5}$ scores over ERM, with greater OOD advantages than on ImageNet. On 4/5 datasets, DRM ranks second to MPTS in CVaR_α but matches ERM in average accuracy. GDRM’s performance varies with α , showing only marginal gains over ERM. OHTM slightly boosts ERM’s CVaR_α and mean accuracy in most settings, while consuming more memory during tuning processes. Still, DRM sacrifices both memory and computational efficiency for robustness in the bottom of Table 2. While MPTS shares DRM’s optimization goal, its risk predictive module and larger-batch simulation enable better task exploration at minimal computational cost, yielding a more robust machine learner.

4 Discussion

Rapid adaptation to novel scenarios is a cornerstone of artificial general intelligence. However, challenges such as safety, limited annotations, and computational constraints necessitate robust and efficient adaptation mechanisms. This study explores learn-to-adapt optimization via generative modeling and introduces MPTS, a versatile framework for robust active task sampling.

Experiments demonstrate the feasibility of predicting optimization outcomes for active task selection. Meanwhile, MPTS enhances adaptation robustness across diverse scenarios in an efficient manner. These results highlight MPTS’s potential to scale CVaR $_{\alpha}$ principles for foundation model development and large-scale decision-making, without additional learning resources.

5 Methods

In alignment with the realistic necessities, this work focuses on robust adaptation while securing learning efficiency, such as circumventing partial expensive evaluation. Such a purpose facilitates the birth of MPTS. As previously mentioned, the framework is agnostic to adaptation learning methods; hence, we leave out zero-shot learning, few-shot learning, and SFT details.

In Fig. 1a, several roles are involved in the optimization: (1) the **adaptive machine learner**, e.g., foundation models or generalist policies, learns to adapt given some optimizers; (2) the **risk predictive model** as a critic evaluates and forecasts the task-specific adaptation risk; (3) the **task sampler** as an actor works for screening the task subset for next iteration. These components participate in episodic learning until convergence.

Technically, this work recasts task episodic learning to sequence generation and presents MPTS as the task sampling strategy to balance exploration and exploitation. At first, we introduce the foundation of RPMs for ranking task difficulty. To reconcile theory and practice, we introduce a tractable optimization approach to enable functional posterior inference towards adaptation risk. Then, we devise the acquisition function informed by the captured risk landscapes. Finally, an understanding concerning the optimization pipeline is attached to conclude the **Methods** part.

5.1 Theoretical Feasibility of Constructing RPMs

Predictive Foundation. MPTS in Definition 2 leverages the risk history and generalization results under θ_t to approximately score the difficulty of task samples for the $(t + 1)$ -th iteration. An unbiased Monte Carlo estimate of $\text{CVaR}_{1-\hat{\mathcal{B}}/\mathcal{B}}$ would require fully evaluation of the $(t + 1)$ -th adaptation risks across $\hat{\mathcal{B}}$ tasks in order to identify the Top- \mathcal{B} subset, which is computationally prohibitive. Instead, MPTS circumvents this costly evaluation by utilizing the risk batch $\{\ell_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ from iteration t , together with historical information $H_{1:t-1}$, to train a generative model to rank difficulty of $\hat{\mathcal{B}}$ tasks for the $(t + 1)$ -th subset selection. This design inevitably introduces the look-ahead bias (see **Theorem 1** and **Lemma 1-2**) while avoiding direct access to the expensive evaluations of $\{\ell_{t+1,i}\}_{i=1}^{\hat{\mathcal{B}}}$.

We begin by introducing Assumptions 1/2/3/4, which characterize the smoothness, boundedness and margin conditions essential to the optimization framework. Specifically, under a fixed machine learner θ , it is reasonable to expect that similar tasks, represented by τ , will exhibit sufficiently close adaptation risk values.

Assumption 1 (Lipschitz Continuity) We assume the adaptation risk function $\ell(\cdot; \theta)$ reserves the Lipschitz continuity w.r.t. θ and τ , i.e.,

$$|\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) - \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta')| \leq \beta_1 \|\theta - \theta'\| \quad \text{and} \quad |\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) - \ell(\mathcal{D}_{\tau'}^Q, \mathcal{D}_{\tau'}^S; \theta)| \leq \beta_2 \|\tau - \tau'\|, \quad (5)$$

where $\forall \{\theta, \theta'\} \in \Theta$ and $\forall \{\tau, \tau'\} \in \mathcal{T}$ with Lipschitz constants β_1 and β_2 .

Assumption 2 (Bounded Sample Gradient) We assume the norm of the adaptation risk function’s gradient $\nabla \ell(\cdot; \theta_t)$ is bounded:

$$\sup_{\tau \in \mathcal{T}} \|\nabla_{\theta} \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta_t)\|_2 < G_t \quad \text{and} \quad \sup_{\tau \in \mathcal{T}, t \in \mathbb{N}_+} \|\nabla_{\theta} \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta_t)\|_2 < G, \quad (6)$$

where G_t is a positive constant and G is a overall bound.

Assumption 3 (Sub-Gaussian Stochastic Gradient) The stochastic gradient $\tilde{\mathbf{g}} := \mathbf{g} + \epsilon$ for the machine learner’s adaptation at t -th iteration is σ -sub-Gaussian, which means:

$$\mathbb{E}[\exp(\eta \mathbf{v}^T \epsilon)] \leq \exp\left(\frac{\eta^2 \sigma^2 \|\mathbf{v}\|_2^2}{2}\right) \quad \forall \eta \in \mathbb{R} \text{ and } \mathbf{v} \in \mathbb{R}^d, \quad (7)$$

where $\mathbb{E}[\tilde{\mathbf{g}}] = \mathbf{g}$, $\mathbb{E}[\|\tilde{\mathbf{g}} - \mathbf{g}\|_2^2] \leq \sigma^2$ and $\sigma \in \mathbb{R}_+$.

Assumption 4 (Margin Anti-Concentration) We assume $\Delta_{ij}(\boldsymbol{\theta}_t) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_t)$ has a density uniformly bounded by ρ in a neighborhood of zero. That is, $\forall \epsilon > 0$, the probability inequality holds:

$$\mathbb{P}(|\Delta_{ij}(\boldsymbol{\theta}_t)| \leq \epsilon) \leq \rho\epsilon. \quad (8)$$

Under the aforementioned assumptions, we derive **Theorem 1**. Specifically, we define a random variable $\Delta_{ij}(\boldsymbol{\theta}_t)$ as the sign of the adaptation risk difference and analyze its evolution following gradient updates across a population. Our theoretical analysis demonstrates that, under a sufficiently small learning rate for the machine learner update, a significant proportion of these sign variables remain largely unchanged in a probabilistic sense. This result establishes a rigorous foundation for evaluating relative task difficulty on $\boldsymbol{\theta}_{t+1}$ based on posterior inference outcomes derived from $\boldsymbol{\theta}_t$ and further guides amortizing the sample average Monte Carlo of CVaR $_{\alpha}$ optimization objective (see Fig. 1a-b).

Theorem 1 (Provably Approximately Invariant Task Difficulties) Given arbitrary K data points $\{(\tau_i, \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t))\}_{i=1}^K$, the adaptation gradient $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)$ as a σ -sub-Gaussian random variable and $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)$, we denote the relative difficulty via the difference $\Delta_{ij}(\boldsymbol{\theta}_{t+1}) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_{t+1}) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_{t+1})$ and $\Delta_{ij}(\boldsymbol{\theta}_t) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_t)$ between t -th and $(t+1)$ -th iterations, and the gradient difference as $\mathbf{v}_{ij} := \nabla_{\boldsymbol{\theta}} \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t) - \nabla_{\boldsymbol{\theta}} \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_t)$.

Under Assumption 1/2/3, the set of rank-preserving variable $E_{ij} := \mathbb{1}[\text{sign}(\Delta_{ij}(\boldsymbol{\theta}_{t+1})) = \text{sign}(\Delta_{ij}(\boldsymbol{\theta}_t))]$ satisfies the probability inequality:

$$\mathbb{P}(\bigcap_{i < j} E_{ij}) \geq 1 - \xi,$$

when $\eta_t \leq \frac{\delta_t}{2G_t M_t + \sqrt{8\sigma^2 G_t^2 \ln\left(\frac{K(K-1)}{2\xi}\right)}}$ with G_t in Assumption 2, $\delta_t := \min_{i \neq j} |\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_t)| \in \mathbb{R}_+$, the stochastic gradient norm $M_t := \|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)\|_2$.

Lemma 1 (Misranked Subset Quantity) In the presence of adaptation risk value, let N_f be the number of order flipped cross-pairs between the ground-truth Top- \mathcal{B} subset $\bar{\mathcal{T}}_t^{\mathcal{B}}$ and the remainder of the candidate tasks $\mathcal{T}_t^{\mathcal{C}} := \mathcal{T}_{t+1}^{\mathcal{B}} \setminus \bar{\mathcal{T}}_t^{\mathcal{B}}$ when the machine learner's parameter changes from $\boldsymbol{\theta}_t$ to $\boldsymbol{\theta}_{t+1}$. We denote the number of tasks that change Top- \mathcal{B} membership by $m_{t+1} = |\bar{\mathcal{T}}_t^{\mathcal{B}} \Delta \mathcal{T}_{t+1}^{\mathcal{B}}|$, where $\bar{\mathcal{T}}_{t+1}^{\mathcal{B}}$ is the ground-truth Top- \mathcal{B} task subset under $\boldsymbol{\theta}_{t+1}$. Then the inequality holds: $m_{t+1} \leq 2N_f$.

Lemma 2 (Rank-Preserving Bound in Expectation) With the risk difference notation $\Delta_{ij}(\boldsymbol{\theta}_t) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_t)$ and $\Delta_{ij}(\boldsymbol{\theta}_{t+1}) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_{t+1}) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_{t+1})$, when the (i, j) cross-pair flips its order from $\boldsymbol{\theta}_t$ to $\boldsymbol{\theta}_{t+1}$, we conclude that $|\Delta_{ij}(\boldsymbol{\theta}_t)| \leq |\Delta_{ij}(\boldsymbol{\theta}_{t+1}) - \Delta_{ij}(\boldsymbol{\theta}_t)|$.

Note that the rank of adaptation difficulty in a fixed task set might flip due to the model update. **Lemma 1** and **2** provides a bound for the theoretically misranked tasks in the next-iteration selected subset and the adaptation risk difference over iteration. Starting from the updated model parameter $\boldsymbol{\theta}_{t+1}$, the subsequent one-step gradient update can be written as

$$\boldsymbol{\theta}_{t+2} = \boldsymbol{\theta}_{t+1} - \eta_{t+1} \underbrace{\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_{t+1})}_{\text{Perturbed Gradient After Rank-Flipping}} \quad (9a)$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_{t+1}) = \underbrace{\bar{\mathbf{g}}_{t+1}}_{\text{Unbiased Average Gradient}} + \underbrace{\Delta \mathbf{g}_{t+1}}_{\text{Gradient Difference}}, \quad (9b)$$

where the perturbed gradient reflects the average task gradient after rank-flipping. We estimate the unbiased average gradient $\bar{\mathbf{g}}_{t+1}$ from the ground-truth Top- \mathcal{B} subset under $\boldsymbol{\theta}_{t+1}$, i.e., $\bar{\mathcal{T}}_{t+1}^{\mathcal{B}}$. The result of rank-flipping tasks after model update contributes to the gradient difference as $\Delta \mathbf{g}_{t+1} = \frac{2}{m_{t+1}} \left(\sum_{\tau \in \bar{\mathcal{T}}_t^{\mathcal{B}} \setminus \mathcal{T}_{t+1}^{\mathcal{B}}} \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}_{t+1}) - \sum_{\tau \in \mathcal{T}_{t+1}^{\mathcal{B}} \setminus \bar{\mathcal{T}}_t^{\mathcal{B}}} \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}_{t+1}) \right)$. These terms serve the convergence analysis $\forall t \in \mathbb{N}_+$ in the following contents.

Lemma 3 (Misranking Acute-Angle) Let $\bar{\mathbf{g}}_t$ and $\Delta \mathbf{g}_t$ respectively denote the unbiased Monte Carlo estimate of CVaR and the gradient difference between it and the biased gradient $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)$. For a given $c \in (0, 1)$ and arbitrary small $\epsilon_* > 0$, there exists T such that $\forall t \geq T$, the following dichotomy holds: either (i) $\|\bar{\mathbf{g}}_t\|_2 \leq \epsilon_*$, or (ii) $\|\Delta \mathbf{g}_t\|_2 \leq c\|\bar{\mathbf{g}}_t\|_2$.

Theorem 2 (Convergence with Diminishing Rank Flipping) Suppose there exists $c \in (0, 1)$ and $T \geq 0$ such that $\forall t \geq T$, the inequality holds: $\|\Delta \mathbf{g}_t\|_2 \leq c\|\bar{\mathbf{g}}_t\|_2$. Given the Assumptions 1/2/3/4, and the appropriate construction of the step sizes $\{\eta_t\}$ from Theorem 1, we conclude that: $\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)\|_2] = 0$. Hence, the iteration converges to first-order stationary points in expectation.

Theorem 2 provides a convergence guarantee of using noisy Top- \mathcal{B} task subset for robust optimization. It is worth noting that, provided the RPM possesses sufficient expressiveness, the accumulated optimization history can theoretically drive down the generalization error arising from function approximation, as supported by statistical learning theory. We leave more detailed discussions and MPTS relevant proof in Supplementary Notes E.

5.2 Generative Modeling Risk Functions and Posterior Inference

Here, we design the sampling strategy through the lens of risk landscapes and pay more attention to datasets of learning optimization outcome $\{H_t\}_{t=1}^T$. To characterize the adaptation risk during batch optimization, we introduce the latent variable z_t to summarize episodic information and present a versatile deep generative model as:

$$p(H_{0:T}, z_{0:T} | \theta_{0:T}) = p(z_0) \prod_{t=0}^T p_\psi(H_t | z_t; \theta_t) \prod_{t=0}^{T-1} p(z_{t+1} | z_t). \quad (10)$$

Within a Bayesian framework, we approximate the underlying function distribution with the latent variable, and the posterior $p(z_t | H_t)$ summarizes the historical risk information and accounts for uncertainty in distributions. The following writes the form of $p(z_t | H_t)$ according to the Bayes rule⁵⁸:

$$p(z_t | H_t) = \frac{p(H_t | z_t) p(z_t | H_{1:t-1})}{\int p(H_t | z_t) p(z_t | H_{1:t-1}) dz_t}, \quad (11)$$

where $p(z_t | H_{1:t-1})$ encodes the past evaluation results as the conditional prior. Moreover, $p(H_t | z_t)$ conveys the likelihood of producing observations of the task batch risk values in the t -th iteration. Notably, the exact computation *w.r.t.* the posterior is intractable due to the complicated integral in the denominator.

Generative Process. As illustrated in Fig. 2c, risk values of the task batch ℓ are correlated with the machine learner’s parameters θ . In specific, the factorization of the sequential optimization relevant variables arrives at:

$$p_\psi(H_t | H_{1:t-1}) = \int p_\psi(H_t | z_t) p(z_t | H_{1:t-1}) dz_t = \int \left[\prod_{i=1}^{\mathcal{B}} p_\psi(\ell_{t,i} | \tau_{t,i}, z_t; \theta_t) \right] p(z_t | H_{1:t-1}) dz_t, \quad (12)$$

where z_t in the probabilistic graphical model constitutes the distribution over risk functions (For the sake of simplicity, we skip over other variables less relevant to our learning purposes). Here, we assume the conditional independence between task-specific risk values given z and the machine learner’s parameter θ in Eq. (12). And the *primary optimization objective* is to $\max_{\psi \in \Psi} \ln p_\psi(H_t | H_{1:t-1})$ for the optimization outcome prediction.

Inference Process. The manner of episodic training, where the task batch and its evaluation arrive sequentially, inspires us to predict adaptation risk values online to actively sample tasks in a batch. However, the exact inference *w.r.t.* $p(z_t | H_t)$ is infeasible as there is no structural information regarding posteriors. In each iteration, the risk function distribution relies on the updated machine learner θ ; hence, such *non-stationarity* in the risk function distributions prompts us to involve the streaming VI^{26,27} to derive the approximate posterior.

To do so, we handle the streaming task batches and update the posterior in a recursive way:

$$\underbrace{p(z_t | H_t)}_{\text{Updated Posterior}} \propto \underbrace{p(H_t | z_t)}_{\text{Likelihood}} \underbrace{p(z_t | H_{1:t-1})}_{\text{Functional Prior}} \quad (13)$$

where $p(z_t | H_{1:t-1})$ represents the conditional prior using the last time updated posterior as the proxy. The role of the estimated functional posterior is to provide uncertainty-aware prediction and serves the task sampling strategy design, which will be detailed in Section 5.3.

As a result, we can formulate the evidence lower bound (ELBO) as a tractable optimization objective in Eq. (14) from approximate inference.

$$\max_{\psi \in \Psi, \phi \in \Phi} \hat{\mathcal{J}}_{\text{ELBO}}(\psi, \phi) := \mathbb{E}_{q_\phi(z_t | H_t)} \left[\sum_{i=1}^{\mathcal{B}} \ln p_\psi(\ell_{t,i} | \tau_{t,i}, z_t) \right] - D_{KL}[q_\phi(z_t | H_t) \parallel p(z_t | H_{1:t-1})] \quad (14)$$

For implementation convenience, we adopt the parameterized Gaussian distribution with diagonal covariance matrices as variational distributions similar to vanilla VAEs^{23,59} and neural processes (NPs)⁶⁰. In other words, these distribution parameters are approximated with neural networks, e.g., $q_\phi(z_t | H_t) = \mathcal{N}(z_t; \mu_\phi(H_t), \Sigma_\phi(H_t))$, and the reparameterization trick²³ is used for stochastic gradient estimate.

Definition 3 (Permutation Invariant Function) With an n -element permutation group \mathcal{S}_n , the operator $g \in \mathcal{S}_n$ maps the order set to itself:

$$g : [1, 2, \dots, n] \mapsto [g_1, g_2, \dots, g_n]. \quad (15)$$

Then the function Φ is called permutation invariant if for any set of data points $\mathbf{x}_1, \dots, \mathbf{x}_n$, the following condition holds:

$$\Phi(g \circ [\mathbf{x}_1, \dots, \mathbf{x}_n]) = \Phi([\mathbf{x}_{g_1}, \dots, \mathbf{x}_{g_n}]) = \Phi([\mathbf{x}_1, \dots, \mathbf{x}_n]) \quad \forall g \in \mathcal{S}_n. \quad (16)$$

As for the neural architecture, we employ the DeepSet encoding module⁶¹ to process the set dataset H_t , which corresponds to the permutation invariant function family in Definition 3. Also, in the context of streaming VI, $q_\phi(\mathbf{z}_t|H_{t-1})$ mostly works as the proxy for the conditional prior as default. Consequently, we can modify the exact ELBO in Eq. (14) and further translate the practical optimization process with the Lagrange multiplier β into:

$$\max_{\psi \in \Psi, \phi \in \Phi} \mathbb{E}_{q_\phi(\mathbf{z}_t|H_t)} \left[\sum_{i=1}^{\mathcal{B}} \ln p_\psi(\ell_{t,i}|\tau_{t,i}, \mathbf{z}_t) \right] \quad \text{s.t.} \quad D_{KL}[q_\phi(\mathbf{z}_t|H_t) \parallel q_{\bar{\phi}}(\mathbf{z}_t|H_{t-1})] \leq \epsilon \Leftrightarrow \quad (17a)$$

$$\max_{\psi \in \Psi, \phi \in \Phi} \mathcal{G}_{\text{ELBO}}(\psi, \phi) := \mathbb{E}_{q_\phi(\mathbf{z}_t|H_t)} \left[\sum_{i=1}^{\mathcal{B}} \ln p_\psi(\ell_{t,i}|\tau_{t,i}, \mathbf{z}_t) \right] - \beta D_{KL}[q_\phi(\mathbf{z}_t|H_t) \parallel q_{\bar{\phi}}(\mathbf{z}_t|H_{t-1})], \quad (17b)$$

where $\bar{\phi}$ indicates no gradients computed through ϕ in the term, and $\{\beta \in \mathbb{R}_+, \epsilon \in \mathbb{R}_+\}$ constrains the machine learner's parameter search in next iteration.

5.3 Task Sampling Strategy Design

In robust active task sampling, existing strategies evaluate task batches to rank their difficulties in adaptation and then prioritize challenging subsets for optimization^{14,16,17,43,44}. Besides the expensive evaluation cost, these strategies are weak in the efficient exploration of the task space.

As Theorem 1 has established the theoretical foundation of approximately rank task difficulty, this necessitates the development of the RPM from cumulated risk histories. With the model predictive results as amortized evaluation, specific rules can be incorporated into the acquisition function for active sampling. Meanwhile, it is fascinating for the RPM to evaluate the machine learner's adaptation to arbitrarily many tasks with minimal computational cost. Hence, we can easily enlarge the pseudo batch size $\hat{\mathcal{B}}$ for more selection candidates and exploit the epistemic uncertainty from the RPM, encouraging more exploration in the task space.

Evaluating Adaptation Performance through Stochastic Forward Passes. The RPM and estimated functional posteriors in Eq. (12)/(14) work as tools for the active selection of the task batch. Specifically, the predictive distribution can be depicted as:

$$\begin{aligned} p_\psi(\ell|\tau, H_{1:t}) &= \int p_\psi(\ell|\tau, \mathbf{z}_t) p(\mathbf{z}_t|H_{1:t}) d\mathbf{z}_t \triangleq \int p_\psi(\ell|\tau, \mathbf{z}_t) q_\phi(\mathbf{z}_t|H_t) d\mathbf{z}_t \\ &\approx \frac{1}{K} \sum_{k=1}^K p_\psi(\ell|\tau, \mathbf{z}_t^{(k)}), \quad \text{with } \mathbf{z}_t^{(k)} \sim q_\phi(\mathbf{z}_t|H_t) \quad \forall \tau \sim p(\tau). \end{aligned} \quad (18)$$

The above predictive distribution $p_\psi(\ell|\tau, H_{1:t})$ benefits from the Bayesian modeling and provides a tractable way to roughly assess difficulties of tasks throughout the whole task space.

Rank-Filtering the Next Task Batch to Episodically Train. After obtaining $p_\psi(\ell|\tau, H_{1:t})$, we draw up a batch sampling strategy on the basis of its quantified statistics. The criteria resembles the acquisition function in classical Bayesian optimization (BO), which includes a collection of available evaluation principles, such as expected improvement⁶², output information theoretical index⁶³ or UCB³⁶.

However, it is also necessary to clarify that the search space is on the sequential task batch instead of machine learners' parameters, which differs from the ultimate purpose in BO. Central to our approach is the principle of optimism in the face of uncertainty⁶⁴. We consider the difficult task's prioritization for robustness and the epistemic uncertainty as pivotal elements in developing acquisition functions. The grounds behind this idea are that (i) the subset with the worst performance deserves extra attention in optimization for adaptation robustness, and (ii) task regions with high predictive uncertainty tend to be underexplored in the last few iterations.

As a result, we present the acquisition function built on the UCB principle³⁶:

$$\begin{aligned} \mathcal{A}(\mathcal{T}^{\mathcal{B}}; \phi, \psi) &= \sum_{i=1}^{\mathcal{B}} a(\tau_i) = \sum_{i=1}^{\mathcal{B}} \overbrace{\gamma_0}^{\text{Risk Mean}} m(\ell_i) + \gamma_1 \overbrace{\sigma(\ell_i)}^{\text{Epistemic Uncertainty}}, \quad \text{where } \tau_i \sim p(\tau) \\ \text{with } m(\ell_i) &= \mathbb{E}_{q_\phi(\mathbf{z}_t|H_t)} [p_\psi(\ell|\tau_i, \mathbf{z}_t)] \quad \text{and } \sigma(\ell_i) = \mathbb{V}_{q_\phi(\mathbf{z}_t|H_t)}^{\frac{1}{2}} [p_\psi(\ell|\tau_i, \mathbf{z}_t)], \end{aligned} \quad (19)$$

where $m(\ell_i)$ and $\sigma(\ell_i)$ are, respectively, the adaptation risk mean and standard deviations, which can be estimated from multiple stochastic forward passes $z_t \sim p(z_t|H_{1:t})$ and $\ell \sim p_\psi(\ell|\tau_i, z_t)$ using the risk generative model. And $\{\gamma_0, \gamma_1\}$ are hyperparameters to balance considerations.

Then, the Simulate-Rank-Filter operation in Eq. (4)c arrives at the task batch for $(t + 1)$ -th iteration, i.e., $\mathcal{T}_{t+1}^{\mathcal{B}} = \arg \max_{\mathcal{T}^{\mathcal{B}} \subseteq \mathcal{T}_{t+1}^{\hat{\mathcal{B}}}, |\mathcal{T}^{\mathcal{B}}| = \mathcal{B}} \mathcal{A}(\mathcal{T}^{\mathcal{B}}; \phi, \psi)$. This characterizes the step of the active subset selection from $\mathcal{T}_{t+1}^{\hat{\mathcal{B}}}$, the randomly sampled identifier candidate set with $|\mathcal{T}_{t+1}^{\hat{\mathcal{B}}}| = \hat{\mathcal{B}}$. In an implementation, we still perform random sampling from $p(\tau)$ and forecast the task-wise acquisition score $a(\cdot)$ from the RPM. Candidates in Top- \mathcal{B} acquisition scores are screened to formulate the task batch $\mathcal{T}_{t+1}^{\mathcal{B}}$ for episodic optimization, as illustrated in Fig. 1d. These steps approximately solve Eq. (4)c and obtain $\mathcal{T}_{t+1}^{\mathcal{B}}$ in a heuristic way.

5.4 Sequentially Optimize the Adaptive Machine Learner

Given the screened \mathcal{T}_{t+1} , we execute optimization to update the machine learner’s parameters. The task-specific adaptation risk in $(t + 1)$ -th iteration is written as $\ell_{t+1,i}(\theta)$ for the selected task τ_i . The developed MPTS is agnostic to any-shot learning methods, and the following includes the standard update rule for zero-shot, few-shot, and SFT scenarios.

Machine Learner Updates in Zero-Shot Adaptation: The zero-shot setup does not require the support dataset to identify the task. Hence, taking the vanilla DR⁶⁵ as an instantiation, we can obtain the update rule as:

$$\theta_{t+1} = \theta_t - \frac{\lambda}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \nabla_{\theta} \ell(\mathcal{D}_{\tau_{t+1,i}}^Q; \theta_t), \quad (20)$$

where θ denotes the zero-shot learning model parameter with λ the learning rate.

Machine Learner Updates in Few-Shot Adaptation: Still, we take the typical optimization-based method MAML³² as an instantiation and provide the update rule as follows:

$$\ell_{t+1,i}(\theta) = \ell(\mathcal{D}_{\tau_{t+1,i}}^Q; \theta_t^{\text{meta}} - \lambda_{1,1} \nabla_{\theta} \ell(\mathcal{D}_{\tau_{t+1,i}}^S)) \quad (21a)$$

$$\theta_{t+1}^{\text{meta}} = \theta_t^{\text{meta}} - \frac{\lambda_{1,2}}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \nabla_{\theta} \ell_{t+1,i}(\theta), \quad \forall i \in \{1, \dots, \mathcal{B}\} \quad (21b)$$

where θ^{meta} denotes the meta initialization, and $\lambda_{1,1}$ and $\lambda_{1,2}$ are, respectively, learning rates in the inner and outer loops.

Machine Learner Updates in SFT: Here, we take finetuning pretrained models to downstream tasks⁶⁶ as an instantiation. In this case, each data point $[x, y]$ can be viewed as a task with either its embedding τ or x as the task identifier. Then the model update rule can be:

$$\theta_{t+1} = \theta_t - \frac{\lambda}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \nabla_{\theta} \ell([x_{t+1,i}, y_{t+1,i}]; \theta_t), \quad (22)$$

where $\{[x_{t+1,i}, y_{t+1,i}]\}_{i=1}^{\mathcal{B}}$ denote the sampled task batch for the $(t + 1)$ -th iteration.

5.5 Overall Algorithm and Interpretation

Implementation Pipelines. Here, we write the general form of MPTS in **Algorithm 1**, where the past risk episodes are reused to train the RPM and serve the active subset selection. We also provide some implementation examples by putting all the ingredients and optimization recipes together in the zero-shot, few-shot, and SFT scenarios. See Supplementary Notes in **Algorithm 2-7** for details. Since the first iteration in **Algorithm 2/4/6** does not involve active sampling, due to no latest history, and the task batch follows the standard random sampling setup.

Connection with Sequential Decision-making and Control. Intuitively, MPTS resembles model predictive control (MPC)⁶⁷ when treating task sampling under some criteria as an optimal planning problem. In this case, the episodic learning process specifies an underlying dynamical system for MPTS to predict with only one future time step in the simulation to assess the influence of selecting the task batch, and the feedback as exact adaptation risk information further helps improve the episodic risk prediction system. In addition, through the lens of sequential decision-making, we can interpret the optimization pipeline of MPTS from the actor-critic framework in RL⁶⁸. In detail, the RPM works as the critic that predicts adaptation performance

in the task τ given a fixed machine learner. Accordingly, the actor plays the role of selecting the task batch from the acquisition function and then executing the machine learner’s optimization. These two roles are entangled in the MPTS pipeline to achieve robust yet efficient adaptation.

Data Availability

Here, we declare that the source of data used in experiments is from open-source dataset repositories or widely adopted simulators. The sinusoid simulator is from work³². The six few-shot image classification datasets are from ImageNet-CG⁵², ImageNet-CI⁵², ImageNet-CS⁵², ImageNet-A⁵³, ImageNet-S²⁹ and ImageNet-R⁵⁴. The image dataset for SFT is from ImageNet²⁸ and four OOD datasets-ImageNet-A⁵³, ImageNet-S²⁹, ImageNet-R⁵⁴, and ImageNet-V⁵⁷. The Meta RL simulators are based on the Mujoco robotic software with configurations reported in our code. The DR environments Ergo-Reacher and Lunar-Lander are from work³⁹.

As for the data of experimental results, it is reflected in the Tables and Figures, which can be found in the attached manuscript file.

Code Availability

The code of Model Predictive Task Sampling and demonstration is accessible at <https://github.com/thu-rllab/MPTS>.

Acknowledgments and Disclosure of Funding

This work is funded by National Natural Science Foundation of China (NSFC) projects with Numbers # 62306326 and # 62495091.

References

- [1] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.
- [4] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [6] Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- [7] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RL2: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- [8] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International conference on machine learning*, pages 5637–5664. PMLR, 2021.
- [9] Anurag Ajay, Abhishek Gupta, Dibya Ghosh, Sergey Levine, and Pulkit Agrawal. Distributionally adaptive meta reinforcement learning. *Advances in Neural Information Processing Systems*, 35:25856–25869, 2022.
- [10] Jiuding Sun, Chantal Shaib, and Byron C Wallace. Evaluating the zero-shot robustness of instruction-tuned language models. In *International Conference on Learning Representations*. ICLR, 2024.
- [11] Lexin Zhou, Wout Schellaert, Fernando Martínez-Plumed, Yael Moros-Daval, Cèsar Ferri, and José Hernández-Orallo. Larger and more instructable language models become less reliable. *Nature*, pages 1–8, 2024.
- [12] Davis Rempe, Jonah Philion, Leonidas J Guibas, Sanja Fidler, and Or Litany. Generating useful accident-prone driving scenarios via a learned traffic prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17305–17315, 2022.
- [13] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [14] Qi Wang, Yiqin Lv, Zheng Xie, Jincui Huang, et al. A simple yet effective strategy to robustify the meta learning paradigm. *Advances in Neural Information Processing Systems*, 36, 2024.
- [15] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. *Advances in Neural Information Processing Systems*, 35:19523–19536, 2022.
- [16] Talfan Evans, Shreya Pathak, Hamza Merzic, Jonathan Schwarz, Ryutaro Tanno, and Olivier J Henaff. Bad students make great teachers: Active learning accelerates large-scale visual understanding. *arXiv preprint arXiv:2312.05328*, 2023.
- [17] Ido Greenberg, Shie Mannor, Gal Chechik, and Eli Meir. Train hard, fight easy: Robust meta reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [18] Talfan Evans, Nikhil Parthasarathy, Hamza Merzic, and Olivier J Henaff. Data curation via joint example selection further accelerates multimodal learning. *arXiv preprint arXiv:2406.17711*, 2024.
- [19] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [20] Jane X Wang, Zeb Kurth-Nelson, Dhharshan Kumaran, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Demis Hassabis, and Matthew Botvinick. Prefrontal cortex as a meta-reinforcement learning system. *Nature neuroscience*, 21(6):860–868, 2018.

- [21] Wei-Long Zheng, Zhongxuan Wu, Ali Hummos, Guangyu Robert Yang, and Michael M Halassa. Rapid context inference in a thalamocortical model using recurrent neural networks. *Nature Communications*, 15(1):8275, 2024.
- [22] Naomi P Friedman and Trevor W Robbins. The role of prefrontal cortex in cognitive control and executive function. *Neuropsychopharmacology*, 47(1):72–89, 2022.
- [23] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [24] Mandt Stephan, Matthew D Hoffman, David M Blei, et al. Stochastic gradient descent as approximate bayesian inference. *Journal of Machine Learning Research*, 18(134):1–35, 2017.
- [25] Karl Friston, Thomas FitzGerald, Francesco Rigoli, Philipp Schwartenbeck, Giovanni Pezzulo, et al. Active inference and learning. *Neuroscience & Biobehavioral Reviews*, 68:862–879, 2016.
- [26] Tamara Broderick, Nicholas Boyd, Andre Wibisono, Ashia C Wilson, and Michael I Jordan. Streaming variational bayes. *Advances in neural information processing systems*, 26, 2013.
- [27] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. *arXiv preprint arXiv:1710.10628*, 2017.
- [28] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015.
- [29] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. *Advances in Neural Information Processing Systems*, 32, 2019.
- [30] Jakub M Tomczak. *Deep Generative Modeling*. Springer Cham, 2024.
- [31] Jean Kaddour, Steindór Sæmundsson, et al. Probabilistic active meta-learning. *Advances in Neural Information Processing Systems*, 33:20813–20822, 2020.
- [32] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- [33] R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2:21–42, 2000.
- [34] Huaxiu Yao, Yu Wang, Ying Wei, Peilin Zhao, Mehrdad Mahdavi, Defu Lian, and Chelsea Finn. Meta-learning with an adaptive task scheduler. *Advances in Neural Information Processing Systems*, 34:7497–7509, 2021.
- [35] Yifei He, Shiji Zhou, Guojun Zhang, Hyokun Yun, Yi Xu, Belinda Zeng, Trishul Chilimbi, and Han Zhao. Robust multi-task learning with excess risks. In *International Conference on Machine Learning*, pages 18094–18114. PMLR, 2024.
- [36] P Auer. Finite-time analysis of the multiarmed bandit problem, 2002.
- [37] Chenghao Liu, Zhihao Wang, Doyen Sahoo, Yuan Fang, Kun Zhang, and Steven CH Hoi. Adaptive task sampling for meta-learning. In *European Conference on Computer Vision*, pages 752–769. Springer, 2020.
- [38] Muhammad Waleed Gondal, Jochen Gast, Inigo Alonso Ruiz, Richard Droste, Tommaso Macri, Suren Kumar, and Luitpold Staudigl. Domain aligned clip for few-shot classification. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5721–5730, 2024.
- [39] Bhairav Mehta, Manfred Diaz, Florian Golemo, Christopher J Pal, and Liam Paull. Active domain randomization. In *Conference on Robot Learning*, pages 1162–1176. PMLR, 2020.
- [40] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19113–19122, 2023.
- [41] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pages 1587–1596. PMLR, 2018.
- [42] Vladimir Naumovich Vapnik, Vlamimir Vapnik, et al. Statistical learning theory. 1998.
- [43] Aravind Rajeswaran, Sarvjeet Ghotra, Balaraman Ravindran, and Sergey Levine. Epopt: Learning robust neural network policies using model ensembles. In *International Conference on Learning Representations*, 2022.
- [44] Yiqin Lv, Cheems Wang, Dong Liang, and Zheng Xie. Theoretical investigations and practical enhancements on tail task risk minimization in meta learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=Mcrz0o0hwr>.
- [45] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2019.

- [46] Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model pretraining. *Advances in Neural Information Processing Systems*, 36, 2024.
- [47] Joey Hejna, Chethan Anand Bhateja, Yichen Jiang, Karl Pertsch, and Dorsa Sadigh. Remix: Optimizing data mixtures for large scale imitation learning. In *8th Annual Conference on Robot Learning*, 2024.
- [48] Maryam Toloubidokhti, Yubo Ye, Ryan Missel, Xiajun Jiang, Nilesh Kumar, Ruby Shrestha, and Linwei Wang. Dats: Difficulty-aware task sampler for meta-learning physics-informed neural networks. In *The Twelfth International Conference on Learning Representations*, 2023.
- [49] Ramnath Kumar, Tristan Deleu, and Yoshua Bengio. The effect of diversity in meta-learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 8396–8404, 2023.
- [50] Jingyao Wang, Wenwen Qiang, Xingzhe Su, Changwen Zheng, Fuchun Sun, and Hui Xiong. Towards task sampler learning for meta-learning. *International Journal of Computer Vision*, 132(12):5534–5564, 2024.
- [51] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- [52] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *ICLR*, 2019.
- [53] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15262–15271, 2021.
- [54] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. *ICCV*, 2021.
- [55] Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient off-policy meta-reinforcement learning via probabilistic context variables. In *International conference on machine learning*, pages 5331–5340. PMLR, 2019.
- [56] Ido Greenberg, Shie Mannor, Gal Chechik, and Eli Meirom. Train hard, fight easy: Robust meta reinforcement learning. *Advances in Neural Information Processing Systems*, 36:68276–68299, 2023.
- [57] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *ICML*, 2019.
- [58] Stephen M Stigler. Thomas bayes’s bayesian inference. *Journal of the Royal Statistical Society: Series A (General)*, 145(2):250–258, 1982.
- [59] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.
- [60] Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye Teh. Neural processes. *arXiv preprint arXiv:1807.01622*, 2018.
- [61] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30, 2017.
- [62] J Mockus, V Tiesis, and A Zilinskas. The application of bayesian methods for seeking the extremum, vol. 2. *L Dixon and G Szego. Toward Global Optimization*, 2, 1978.
- [63] Binxin Ru, Michael A Osborne, Mark McLeod, and Diego Granziol. Fast information-theoretic bayesian optimisation. In *International Conference on Machine Learning*, pages 4384–4392. PMLR, 2018.
- [64] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3 (Nov):397–422, 2002.
- [65] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017.
- [66] Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature Machine Intelligence*, 5(3):220–235, 2023.
- [67] Manfred Morari and Jay H Lee. Model predictive control: past, present and future. *Computers & chemical engineering*, 23(4-5):667–682, 1999.
- [68] Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.

- [69] Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet: Finding important examples early in training. *Advances in neural information processing systems*, 34:20596–20607, 2021.
- [70] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- [71] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of machine learning models. In *International Conference on Machine Learning*, pages 6950–6960. PMLR, 2020.
- [72] Yu Yang, Hao Kang, and Baharan Mirzasoleiman. Towards sustainable learning: Coresets for data-efficient deep learning. In *International Conference on Machine Learning*, pages 39314–39330. PMLR, 2023.
- [73] Siddharth Joshi and Baharan Mirzasoleiman. Data-efficient contrastive self-supervised learning: Most beneficial examples for supervised learning contribute the least. In *International conference on machine learning*, pages 15356–15370. PMLR, 2023.
- [74] Andreas Damianou and Neil D Lawrence. Deep gaussian processes. In *Artificial intelligence and statistics*, pages 207–215. PMLR, 2013.
- [75] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. *Advances in neural information processing systems*, 25, 2012.
- [76] James Wilson, Frank Hutter, and Marc Deisenroth. Maximizing acquisition functions for bayesian optimization. *Advances in neural information processing systems*, 31, 2018.
- [77] Roman Garnett. *Bayesian optimization*. Cambridge University Press, 2023.
- [78] Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM computing surveys (csur)*, 53(3):1–34, 2020.
- [79] Qi Wang, Yanghe Feng, Jincai Huang, Yiqin Lv, Zheng Xie, and Xiaoshan Gao. Large-scale generative simulation artificial intelligence: The next hotspot. *The Innovation*, 4(6), 2023.
- [80] Wenjia Xu, Yongqin Xian, Jiuniu Wang, Bernt Schiele, and Zeynep Akata. Attribute prototype network for any-shot learning. *International Journal of Computer Vision*, 130(7):1735–1753, 2022.
- [81] Shichen Liu, Mingsheng Long, Jianmin Wang, and Michael I Jordan. Generalized zero-shot learning with deep calibration network. *Advances in neural information processing systems*, 31, 2018.
- [82] Aoxue Li, Zhiwu Lu, Jiechao Guan, Tao Xiang, Liwei Wang, and Ji-Rong Wen. Transferrable feature and projection learning with class hierarchy for zero-shot learning. *International Journal of Computer Vision*, 128:2810–2827, 2020.
- [83] Rohit Keshari, Richa Singh, and Mayank Vatsa. Generalized zero-shot learning via over-complete distribution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13300–13308, 2020.
- [84] Yongqin Xian, Tobias Lorenz, Bernt Schiele, and Zeynep Akata. Feature generating networks for zero-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5542–5551, 2018.
- [85] Edgar Schonfeld, Sayna Ebrahimi, Samarth Sinha, Trevor Darrell, and Zeynep Akata. Generalized zero-and few-shot learning via aligned variational autoencoders. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8247–8255, 2019.
- [86] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):5149–5169, 2021.
- [87] Qi Wang, Marco Federici, and Herke van Hoof. Bridge the inference gaps of neural processes via expectation maximization. In *The Eleventh International Conference on Learning Representations*, 2022.
- [88] Muhammad Waleed Gondal, Shruti Joshi, Nasim Rahaman, Stefan Bauer, Manuel Wuthrich, and Bernhard Schölkopf. Function contrastive learning of transferable meta-representations. In *International Conference on Machine Learning*, pages 3755–3765. PMLR, 2021.
- [89] Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. Conditional neural processes. In *International conference on machine learning*, pages 1704–1713. PMLR, 2018.
- [90] Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. *Advances in neural information processing systems*, 31, 2018.
- [91] Momin Abbas, Quan Xiao, Lisha Chen, Pin-Yu Chen, and Tianyi Chen. Sharp-maml: Sharpness-aware model-agnostic meta learning. In *International conference on machine learning*, pages 10–32. PMLR, 2022.
- [92] Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. *Advances in neural information processing systems*, 32, 2019.

- [93] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017.
- [94] Kelsey Allen, Evan Shelhamer, Hanul Shin, and Joshua Tenenbaum. Infinite mixture prototypes for few-shot learning. In *International conference on machine learning*, pages 232–241. PMLR, 2019.
- [95] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- [96] Marcin Sendera, Marcin Przewięźlikowski, Konrad Karanowski, Maciej Zięba, Jacek Tabor, and Przemysław Spurek. Hypershot: Few-shot learning by kernel hypernetworks. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2469–2478, 2023.
- [97] Yonatan Oren, Shiori Sagawa, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust language modeling. *arXiv preprint arXiv:1909.02060*, 2019.
- [98] Liam Collins, Aryan Mokhtari, and Sanjay Shakkottai. Task-robust model-agnostic meta-learning. *Advances in Neural Information Processing Systems*, 33:18860–18871, 2020.
- [99] Yun Qu, Qi Wang, Yixiu Mao, Vincent Tao Hu, Björn Ommer, and Xiangyang Ji. Can prompt difficulty be online predicted for accelerating rl finetuning of reasoning models? *arXiv preprint arXiv:2507.04632*, 2025.
- [100] Jonathan Gordon, Wessel P Bruinsma, Andrew YK Foong, James Requeima, Yann Dubois, and Richard E Turner. Convolutional conditional neural processes. *arXiv preprint arXiv:1910.13556*, 2019.
- [101] Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. Attentive neural processes. *arXiv preprint arXiv:1901.05761*, 2019.
- [102] Andrew Foong, Wessel Bruinsma, Jonathan Gordon, Yann Dubois, James Requeima, and Richard Turner. Meta-learning stationary stochastic process prediction with convolutional neural processes. *Advances in Neural Information Processing Systems*, 33:8284–8295, 2020.
- [103] Qi Wang and Herke Van Hoof. Learning expressive meta-representations with mixture of expert neural processes. *Advances in neural information processing systems*, 35:26242–26255, 2022.
- [104] Philippe Rigollet and Jan-Christian Hütter. High-dimensional statistics. *arXiv preprint arXiv:2310.19244*, 2023.
- [105] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.
- [106] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16816–16825, 2022.
- [107] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages 5026–5033. IEEE, 2012.
- [108] Erin Catto. Box2d: A 2d physics engine for games, 2007. URL <http://box2d.org>.
- [109] Florian Golemo, Adrien Ali Taiga, Aaron Courville, and Pierre-Yves Oudeyer. Sim-to-real transfer with neural-augmented robot simulation. In *Conference on Robot Learning*, pages 817–828. PMLR, 2018.
- [110] Erwin Coumans. Bullet physics simulation. In *ACM SIGGRAPH 2015 Courses*, page 1. 2015.

Contents

1	Introduction	1
2	Adaptation and Robustness	3
3	Results	5
3.1	Demonstration of the MPTS’s role in K -shot sinusoid regression	7
3.2	Few-Shot adaptation benefits from MPTS in robustness and learning efficiency	8
3.3	MPTS retains multi-faced advantages beyond robustness in zero-shot continuous control	10
3.4	MPTS also reserves the potential of robust SFT	12
4	Discussion	13
5	Methods	13
5.1	Theoretical Feasibility of Constructing RPMs	13
5.2	Generative Modeling Risk Functions and Posterior Inference	15
5.3	Task Sampling Strategy Design	16
5.4	Sequentially Optimize the Adaptive Machine Learner	17
5.5	Overall Algorithm and Interpretation	17
A	Quick Guideline to MPTS	27
B	Research Background	29
B.1	Adaptation Learning for Cross-Task Generalization	29
B.2	Dataset Curation and Task-Level Robustness	29
B.3	Risk Minimization Principles and Prioritized Sampling	30
C	Task Construction and Identifiers	32
C.1	Tasks with Explicit Identifiers	32
C.2	Tasks with Implicit Identifiers	32
C.3	Scalability with Large Reasoning Models	33
D	Auto-Encoding Adaptation Risk through Streaming VI	33
D.1	Neural Modules to Parameterize Distributions	34
D.2	Formulation of ELBO & Stochastic Gradient Estimates	34
D.3	Theoretical Guarantee for Task Difficulties’ Scoring with Posterior Inference	35
E	Convergence Analysis	37
F	Backbone Methods & Experimental Details in Any-Shot Learning	40
F.1	MAML, Reptile and PEARL	40
F.2	Domain Randomization	41
F.3	Multi-Modal Prompt Learning	41

G	Experimental Setups & Implementation Details	42
G.1	Sinusoid Regression	42
G.2	N-way K-shot Image Classification	43
G.3	Prompt-based Few-shot Image Classification	43
G.4	Meta-RL	44
G.5	Robotic DR	44
G.6	Prompt-Tuning Multimodal Foundation Models	45
H	Computational Tools & Platforms & Model Design	45
I	Competing Interests & Author Contributions	45

Supplementary Notes: Model Predictive Task Sampling for Efficient and Robust Adaptation

Algorithm 2: MPTS for DR (Zero-Shot Scenarios)

Input : Task distribution $p(\tau)$; Task batch size \mathcal{B} ;
 Learning rate λ_1 .
Output : Adapted machine learner θ .
 Set the initial iteration number $t = 1$;
 Randomly initialize machine learner θ ;
 Randomly initialize RPM $\{\psi, \phi\}$;
while not converged do
 Execute **Algorithm 3** to access the batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$ and
 induced $\{\mathcal{D}_{\tau_{t,i}}^Q\}_{i=1}^{\mathcal{B}}$;
 // **Eval Adaptation Performance**
 Compute the task specific adaptation risk
 $\{\ell_{t,i} := \ell(\mathcal{D}_{\tau_{t,i}}^Q; \theta_t)\}_{i=1}^{\mathcal{B}}$;
 Return $H_t = \{[\tau_{t,i}, \ell_{t,i}]\}_{i=1}^{\mathcal{B}}$ as the Input to
 Algorithm 3;
 // **Update Machine Learner**
 Perform batch gradient updates:
 $\theta_{t+1} \leftarrow \theta_t - \frac{\lambda_1}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \nabla_{\theta} \ell_{t,i}$;
 Update the iteration number: $t \leftarrow t + 1$;
end

Algorithm 4: MPTS for Model Agnostic Meta Learning (Few-Shot Scenarios)

Input : Task distribution $p(\tau)$; Task batch size \mathcal{B} ;
 Learning rates: $\{\lambda_{1,1}, \lambda_{1,2}\}$.
Output : Meta-trained initialization θ^{meta} .
 Set the initial iteration number $t = 1$;
 Randomly initialize meta learner θ^{meta} ;
 Randomly initialize RPM $\{\psi, \phi\}$;
while not converged do
 Execute **Algorithm 5** to access the batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$ and
 $\{\mathcal{D}_{\tau_{t,i}}^S \cup \mathcal{D}_{\tau_{t,i}}^Q\}_{i=1}^{\mathcal{B}}$;
 // **Inner Loop to Fast Adapt**
 for $i = 1$ **to** K **do**
 Compute the task-specific gradient:
 $\nabla_{\theta} \ell(\mathcal{D}_{\tau_{t,i}}^S; \theta)$;
 Perform gradient updates as fast adaptation:
 $\theta_t^i \leftarrow \theta_t^{\text{meta}} - \lambda_{1,1} \nabla_{\theta} \ell(\mathcal{D}_{\tau_{t,i}}^S; \theta)$;
 end
 // **Outer Loop to Meta-train**
 Evaluate fast adaptation performance
 $\{\ell_{t,i} := \ell(\mathcal{D}_{\tau_{t,i}}^Q; \theta_t^i)\}_{i=1}^{\mathcal{B}}$;
 Return $H_t = \{[\tau_{t,i}, \ell_{t,i}]\}_{i=1}^{\mathcal{B}}$ as the Input to
 Algorithm 5;
 Perform meta initialization updates:
 $\theta_{t+1}^{\text{meta}} \leftarrow \theta_t^{\text{meta}} - \frac{\lambda_{1,2}}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \nabla_{\theta} \ell_{t,i}$;
 Update the iteration number: $t \leftarrow t + 1$;
end

Algorithm 3: Model Predictive Task Sampling

Input : Task distribution $p(\tau)$; Task batch size \mathcal{B} ;
 Candidate batch size $\hat{\mathcal{B}}$; Latest updated $\{\psi, \phi\}$;
 Latest history H_{t-1} ; Iteration number K ;
 Learning rate λ_2 .
Output : Task identifier batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.
 // **Functional Posterior Inference**
for $i = 1$ **to** K **do**
 Perform gradient updates given H_{t-1} :
 $\phi \leftarrow \phi + \lambda_2 \nabla_{\phi} \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (17b);
 $\psi \leftarrow \psi + \lambda_2 \nabla_{\psi} \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (17b);
end
 // **Simulating Adaptation Results**
 Randomly sample $\{\hat{\tau}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ from $p(\tau)$;
 Run amortized evaluation on candidate tasks
 $\{\delta_i := \gamma_0 m(\ell_i) + \gamma_1 \sigma(\ell_i)\}_{i=1}^{\hat{\mathcal{B}}}$ in Eq. (19);
 // **Active Subset Selection from Predicted Results**
 Rank $\{\delta_i\}_{i=1}^{\hat{\mathcal{B}}}$ and screen Top- \mathcal{B} values;
 Return the screened task batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.

Algorithm 5: Model Predictive Task Sampling

Input : Task distribution $p(\tau)$; Task batch size \mathcal{B} ;
 Candidate batch size $\hat{\mathcal{B}}$; Latest updated $\{\psi, \phi\}$;
 Latest history H_{t-1} ; Iteration number K ;
 Learning rate λ_2 .
Output : Task identifier batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.
 // **Functional Posterior Inference**
for $i = 1$ **to** K **do**
 Perform gradient updates given H_{t-1} :
 $\phi \leftarrow \phi + \lambda_2 \nabla_{\phi} \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (17b);
 $\psi \leftarrow \psi + \lambda_2 \nabla_{\psi} \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (17b);
end
 // **Simulating Adaptation Results**
 Randomly sample $\{\hat{\tau}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ from $p(\tau)$;
 Run amortized evaluation on candidate tasks
 $\{\delta_i := \gamma_0 m(\ell_i) + \gamma_1 \sigma(\ell_i)\}_{i=1}^{\hat{\mathcal{B}}}$ in Eq. (19);
 // **Active Subset Selection from Predicted Results**
 Rank $\{\delta_i\}_{i=1}^{\hat{\mathcal{B}}}$ and screen Top- \mathcal{B} values;
 Return the screened task batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.

Algorithm 6: MPTS for Pretrained Model Finetuning

Input : Task distribution $p(\mathbf{x})$; Task batch size \mathcal{B} ;
Learning rate λ_1 .
Output : Fine-tuned machine learner θ .
Set the initial iteration number $t = 1$;
Randomly initialize machine learner θ ;
Randomly initialize RPM $\{\psi, \phi\}$;
while not converged do
 Execute **Algorithm 7** to access the batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$ and $\{\mathbf{x}_{t,i}, \mathbf{y}_{t,i}\}_{i=1}^{\mathcal{B}}$;
 // **Eval Adaptation Performance**
 Compute the instance-specific adaptation risk $\{\ell_{t,i} := \ell(\mathbf{x}_{t,i}, \mathbf{y}_{t,i}; \theta_t)\}_{i=1}^{\mathcal{B}}$;
 Return $H_t = \{\tau_{t,i}, \ell_{t,i}\}_{i=1}^{\mathcal{B}}$ as the Input to **Algorithm 7**;
 // **Update Machine Learner**
 Perform batch gradient updates:
 $\theta_{t+1} \leftarrow \theta_t - \frac{\lambda_1}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \nabla_{\theta} \ell_{t,i}$;
 Update the iteration number: $t \leftarrow t + 1$;
end

Algorithm 7: Model Predictive Task Sampling

Input : Offline processed τ dataset; Task batch size \mathcal{B} ;
Candidate batch size $\hat{\mathcal{B}}$; Latest updated $\{\psi, \phi\}$;
Latest history H_{t-1} ; Iteration number K ;
Learning rate λ_2 .
Output : Task identifier batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.
// **Functional Posterior Inference**
for $i = 1$ **to** K **do**
 Perform gradient updates given H_{t-1} :
 $\phi \leftarrow \phi + \lambda_2 \nabla_{\phi} \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (17b);
 $\psi \leftarrow \psi + \lambda_2 \nabla_{\psi} \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (17b);
end
// **Simulating Adaptation Results**
Randomly sample $\{\hat{\tau}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ from $p(\tau)$;
Run amortized evaluation on candidate tasks $\{\delta_i := \gamma_0 m(\ell_i) + \gamma_1 \sigma(\ell_i)\}_{i=1}^{\hat{\mathcal{B}}}$ in Eq. (19);
Rank $\{\delta_i\}_{i=1}^{\hat{\mathcal{B}}}$ and screen Top- \mathcal{B} values;
// **Exact Evaluation or Active Annotations**
Return the screened batch $\{\mathbf{x}_{t,i}, \mathbf{y}_{t,i}\}_{i=1}^{\mathcal{B}}$.

A Quick Guideline to MPTS

Task episodic learning serves as a cornerstone in developing adaptive models by structuring diverse, context-rich learning experiences. One of the pivotal insights underpinning this process is the neural scaling law, which establishes a relationship between task volume, model complexity, and computational resources, offering a principled insight into training foundation models at a certain budget. Recent viewpoints have also shed light on the importance of task quality^{15,16,18,69–73}, prompting innovative data curation strategies to refine datasets for pretraining, meta-training, and post-training. Evidence suggests that carefully curated data can significantly reduce task sampling complexity, decrease computational demands, and enhance robustness against distributional shifts—sometimes achieving these goals simultaneously. Despite these advancements, a practical operation such as Evaluate-Rank-Filter still faces challenges associated with costly evaluations from intensive task queries, computational overhead, and massive annotations. Addressing these bottlenecks remains essential to fully realize the potential of task episodic learning in robust efficient foundation model training.

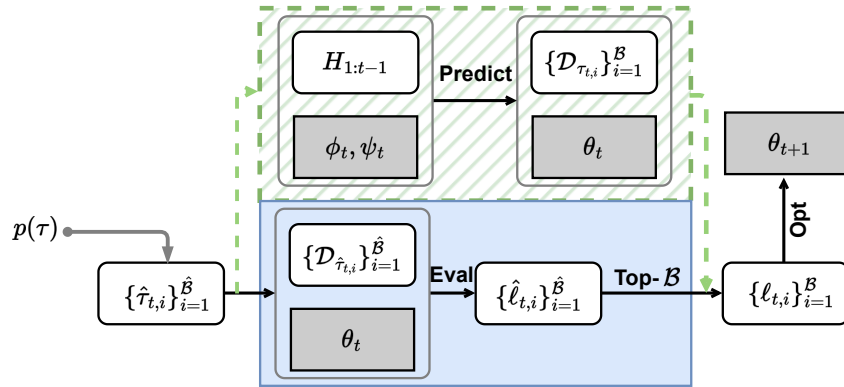


Figure 6: **Risk Predictive Module in MPTS for Active Subset Selection.** MPTS adopts a *predict-then-optimize* strategy and uses a predictive module in green to approximately score the task subset difficulty and obtain the preferred task subset. While the traditional method in blue exhausts $\hat{\mathcal{B}}$ tasks in construction and evaluation to filter preferred subset.

Computational Complexity Analysis. The involvement of the RPM inevitably brings extra computational overhead in optimization. However, the RPM used in this work is lightweight with the model complexity $\mathcal{O}(|\phi| + |\psi|) \ll \mathcal{O}(|\theta|)$. We can

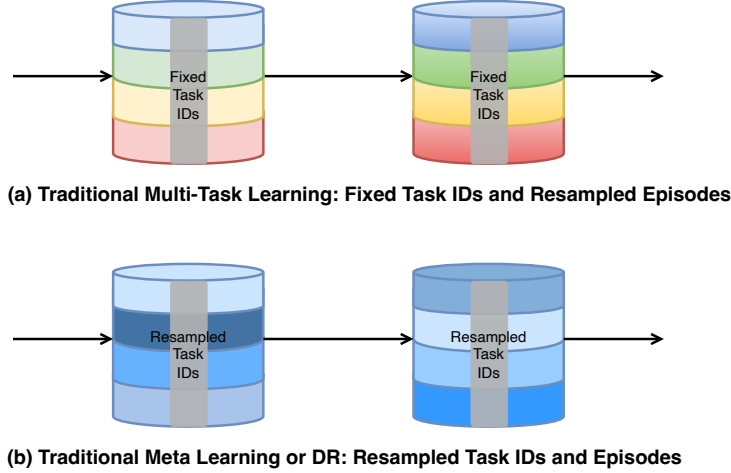


Figure 7: **Comparison between Traditional Multi-Task Learning and Meta Learning and Domain Randomization.** Traditional multi-task learning aims to exploit the correlation between heterogeneous tasks, and task IDs in a batch do not change³⁵. In comparison, meta learning and domain randomization typically consider the homogeneous task family and resample task IDs in a batch to optimize. The latter is our studied scenario, referred to as task episodic optimization. [Same types of colors mean the homogeneous tasks with depth to denote the various episode.]

roughly estimate these extra computations that arise from the predictive model as $\mathcal{O}((|\phi| + |\psi|)T_{\text{MPTS}})$ throughout the training phase. Moreover, the computational and task evaluation complexities of different methods are estimated in Table 3. Compared with DRM, MPTS retains more computational and task efficiency when the filtering ratio $\hat{\alpha}$ is high, and the machine learner θ is largely given similar convergence iteration steps. OHTM⁴⁹ uses the risk buffer to store historical difficult tasks and mixes it into the instant evaluated ones, hence, there is no additional computational cost. As for DATS⁴⁸ and TDPS⁵⁰, they rely on additional gradient updates to compute the weights and raises more computational cost.

Table 3: **Computational Complexities using Different Methods.** Here, we drop out the ranking or reweighting computational complexity as the model complexity of the machine learner considered in this analysis is major, such as the multimodal foundation models. T refers to the required iteration steps until the convergence for separate methods.

	ERM	DRM	GDRM	OHTM	MPTS (Ours)
computation	$\mathcal{O}(\theta T_{\text{ERM}})$	$\mathcal{O}(\frac{1}{1-\hat{\alpha}} \theta T_{\text{DRM}})$	$\mathcal{O}(\theta T_{\text{GDRM}})$	$\mathcal{O}(\theta T_{\text{OHTM}})$	$\mathcal{O}((\phi + \psi + \theta)T_{\text{MPTS}})$
task eval	$\mathcal{O}(\mathcal{B}T_{\text{ERM}})$	$\mathcal{O}(\frac{\mathcal{B}}{1-\hat{\alpha}}T_{\text{DRM}})$	$\mathcal{O}(\mathcal{B}T_{\text{GDRM}})$	$\mathcal{O}(\mathcal{B}T_{\text{OHTM}})$	$\mathcal{O}(\mathcal{B}T_{\text{MPTS}})$

Choice of Surrogate Models. Among MPTS’s core components, the RPM works to predict the adaptation risk values based on historical information and further serves the calculation of acquisition functions. Importantly, this work investigates the feasibility and effectiveness of risk predictive strategies and does not impose rigid constraints on the form of the RPM $p(\ell|\tau, H_{1:t}; \theta)$ too much in modeling. The design of this RPM $p(\ell|\tau, H_{1:t}; \theta)$ must meet several criteria: it is tractable in optimization, can process historical risk information, and offers uncertainty in prediction.

A series of candidate probabilistic models exist that probably apply to adaptation risk modeling. One alternative choice can be the Gaussian process⁷⁴, which provides an analytical form of the predictive distribution. However, its implementation (i) is less scalable in the case of relatively higher dimensional task identifiers, (ii) holds the cubic runtime complexity in obtaining the predictive covariance matrix, (iii) is sensitive to kernel selection, coupled with limited expressiveness of the Gaussian distribution in learned risk functions. Hence, for simplicity and computational efficiency, we adopt the basic VAE-like model and execute a handful of gradient updates to train the RPM. We leave more advanced RPM modeling for future exploration.

Bayesian Optimization for Black-box Functions. This work relates to active sampling and Bayesian optimization. The purpose of BO⁷⁵ is to sequentially find a global optimum of a black-box function $f(x)$ expensive to evaluate in \mathcal{S} , namely $x_* = \arg \max_{x \in \mathcal{S} \subset \mathbb{R}^d} f(x)$.

In each iteration $t = 1, \dots, T$, the BO method actively queries \mathbf{x}_t to evaluate $f(\mathbf{x}_t)$, yielding an output $\ell_t = f(\mathbf{x}_t) + \epsilon$ with a white noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Due to the high cost of function evaluation, the key to BO is constructing a surrogate model to guide the data point to query. The resulting acquisition function⁷⁶ works as an active sampling objective to maximize and obtain the candidate \mathbf{x}_t based on the previous sequence. BO requires limited function evaluations as observations and exploits the correlations in queried data points. These properties make it more theoretically data efficient than random or grid search in seeking the optimal solution⁷⁷. This work differs from standard BO as task episodic learning is not the optimal parameter search problem.

Specific Pseudo Algorithms in Considered Scenarios. The main paper provides the workflow of MPTS in Algorithm 1. For separate scenarios, we attach detailed pseudo algorithms as follows. These illustrated Algorithms are in the context of supervised learning. Regarding RL scenarios, such as meta RL and DR, there is a slight modification for MPTS. As simply picking up worst-case MDPs restricts the task subspace in optimization¹⁷, we adopt the mixture of the identifier subset from the random sampler and the identifier subset from the MPTS sampler. For example, in meta RL, with the pseudo batch size $\hat{\mathcal{B}} = 1.5\mathcal{B}$, there $1.5\mathcal{B}$ identifier candidates from the random sampler. We retain $0.5\mathcal{B}$ random ones and execute standard MPTS amortized evaluation and acquisition rule to obtain another $0.5\mathcal{B}$ ones from the rest random \mathcal{B} identifiers, formulating the mixed \mathcal{B} task batch for RL training. Such an operation makes RL over the MDP distribution stable in optimization. See the open-source code for more RL details.

Application Scope. MPTS considers the task distribution to conduct robust optimization, where tasks are homogeneous in implementations. The entire optimization pipeline adheres to task-episodic training, which means that a batch of tasks is resampled to train in each iteration. However, for classical multi-task learning scenarios, when tasks are fixed and heterogeneous, the task identity of the batch is not changed over iterations. Such a difference can be illustrated in Fig. ?? . Hence, vanilla MPTS might encounter the applicability issue. We leave developing model predictive strategies to balance the weight of tasks for robust optimization as future work.

B Research Background

B.1 Adaptation Learning for Cross-Task Generalization

Learning from zero-shot or few-shot examples has been identified as a crucial adaptation capability of the machine learner nowadays^{78,79}. In SFT, this work treats the individual example as each task to meet MPTS setup. As SFT techniques have been widely discussed in the field⁶⁶, we skip this part in the background introduction.

Zero-Shot Adaptation. This assesses the machine learner’s generalization capability when directly deploying in unseen scenarios without the help of a support dataset. Such a cross-task generalization is commonly studied in computer vision⁸⁰, and the core of the relevant methods is effective semantic representation either from embedding-based methods^{81–83} or generative-based methods^{84,85}. In the era of the foundation models, the pretraining mechanism between multimodality also sometimes empowers the machine learner, such as CLIP¹, with zero-shot capability. When it comes to sequential decision-making, a commonly seen method is DR^{39,65}, which places a distribution over environments for the agent to interact.

Few-Shot Adaptation. This examines the machine learner’s capability of resolving unseen tasks from some annotated examples as hints. Meta-learning, as the typical learning paradigm, has gained popularity over the past decade. It achieves few-shot adaptation by leveraging past experience and distilling knowledge to unseen but similar scenarios in a few-shot way⁸⁶. In brief, we categorize commonly seen methods into context-based, optimization-based, geometric-based, and others. (i) Formulated in an encoder-decoder structure, the context-based method resembles variational autoencoders and encodes the few-shot information into latent variables or embeddings. Typical ones are neural process families^{60,87–89}, which aim to constitute exchangeable deep stochastic processes with neural networks. (ii) The optimization-based methods, with their versatile nature and ability to enable cross-task skill transfer, have piqued the interest and engagement of researchers in the field. For example, MAML^{32,90–92} reduces meta-learning to a bi-level optimization in the parameter space, and its extensions have been widely investigated in the field. (iii) The deep metric-based methods^{93,94} attempt to embed tasks into the latent space and are more suitable for few-shot image classification tasks. Besides, there are other families, such as hyper-networks^{95,96}, recurrent meta-learning⁷, etc.

B.2 Dataset Curation and Task-Level Robustness

Task Curation in Robust Adaptation Learning Pipelines. Recent works^{15,16} demonstrate the effectiveness of challenging task prioritization over uniform sampling in improving cross-task generalization and adaptation robustness, particularly when the learning dataset is sufficiently large. Many methods^{13–18,43} adopt an Evaluate-Rank-Filter step for iterative model updates,

introducing a batch filtering ratio $\hat{\alpha} = 1 - \frac{\mathcal{B}}{\mathcal{B}}$ $\in [0, 1)$ to quantify the fraction of discarded tasks in a sample batch. This prioritization of "difficult" tasks aligns with minimizing CVaR $_{\alpha}$ ³³, a robustness metric for tail-case performance. Alternatively, other methods^{45–47} focus on constructing uncertainty sets and reweighting tasks within the batch to achieve robust adaptation. Additionally, coreset methods^{71–73} aim to select a small subset of tasks that effectively represent the utility of the full dataset, often through gradient approximation in optimization. These approaches address a subproblem of data efficiency, with the acquisition strategy in MPTS serving as an episodic coreset selection mechanism tailored for robustness.

Task Distributional Robustness. The CVaR $_{\alpha}$ or expected shortfall³³ is a statistical measure to assess the proportional worst-case performance of some models at certain levels. This is widely adopted in risk-averse applications. As implied in Definition 1, CVaR $_{\alpha}$ describes the expected risk under the normalized $(1 - \alpha)$ proportional tail risk task distribution, and this work specifies the distribution in the task space. Meanwhile, the normalized tail task distribution $p_{\alpha}(\tau; \theta)$ can be viewed as a shifted result from the initial task distribution $p(\tau)$; hence, such a measure provides robustness quantification in the presence of task distribution shifts^{14,17,44}.

Another indicator to evaluate the machine learner's robustness is the performance in OOD tasks. This refers to the case when the training and the testing task distributions are different. Particularly, in DR and prompt-tuning scenarios, we also use the OOD tasks that never appear in the training task distribution to test the trained policy, and this setup corresponds to domain generalization, a type of substantial distribution shift⁸.

Similar to the setup¹⁴, let $(\Omega_{\tau}, \mathcal{F}_{\tau}, \mathbb{P}_{\tau})$ be a probability space over tasks, where \mathcal{F}_{τ} is a σ -algebra on subsets of Ω_{τ} . Consider (\mathbb{R}, \mathbb{B}) with \mathbb{B} the Borel σ -algebra, defining a probability measure for the adaptation risk function $\ell(\mathcal{D}_{\tau}; \theta)$. For $\theta \in \Theta$, the adaptive optimization operator is $\mathcal{M}_{\theta} : \tau \mapsto \ell(\mathcal{D}_{\tau}; \theta)$.

Thus $\ell(\cdot)$ acts as a random variable to induce the risk distribution $p(\ell)$. The corresponding cumulative distribution is $F_{\ell}(\ell; \theta) = \mathbb{P}(\{\ell(\mathcal{D}_{\tau}; \theta) \leq \ell; \tau \in \Omega_{\tau}, \ell \in \mathbb{R}\})$. Note that $F(\ell; \theta)$ depends on θ and generally lacks a closed form. The resulting probability density function and the following explanation can be found in the main paper Definition 1.

Sample Average Approximation of CVaR. Note that the vanilla optimization objective of CVaR can be expressed in the form of dual representation, which corresponds to:

$$\min_{\theta \in \Theta, \zeta \in \mathbb{R}} \text{CVaR}_{\alpha}(\theta) := \zeta + \frac{1}{1 - \alpha} \mathbb{E}_{p(\tau)} [\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) - \zeta]^+, \quad (23)$$

where the signed function means $[\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) - \zeta]^+ = \max\{\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) - \zeta, 0\}$. With the help of sample average approximation in Monte Carlo, this can be further written as:

$$\min_{\theta \in \Theta, \zeta \in \mathbb{R}} \text{CVaR}_{\alpha}(\theta) := \zeta + \frac{1}{(1 - \alpha)\hat{\mathcal{B}}} \sum_{i=1}^{\hat{\mathcal{B}}} [\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta) - \zeta]^+. \quad (24)$$

Since the optimality for the auxiliary variable holds at the condition $\zeta = \text{VaR}_{\alpha}(\theta)$, selecting the Top- \mathcal{B} element in the set $\{\ell_i | \ell_i = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta)\}_{i=1}^{\hat{\mathcal{B}}}$ is an unbiased Monte Carlo estimate of CVaR. The gradient of Eq. (24) is the steepest direction in gradient optimization, and executing gradient descent decreases the CVaR value. In addition, the adaptation risk in this work $\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta)$, e.g., classification accuracies or MSEs, is typically bounded. With the monotonic improvement and bounded risk function values, optimizing the Monte Carlo CVaR (MC-CVaR) leads to convergence. Such a robust optimization method corresponds to the baseline DRM¹⁴ in this work.

B.3 Risk Minimization Principles and Prioritized Sampling

The risk minimization principles are entangled with task sampling and robust optimization.

Expected/Empirical Risk Minimization (ERM). With the fixed $p(\tau)$, the principle follows the statistical learning theory⁴² and minimizes the expectation of adaptation risk over the task space. As a result, we can have:

$$\min_{\theta \in \Theta} \mathbb{E}_{p(\tau)} [\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta)]. \quad (25)$$

It draws batches with a random task sampler to optimize iteratively.

Distributionally Robust Risk Minimization (DRM)^{14,16,17,43,44}. We retain the notation of task robust work¹⁴, which terms the tail task risk minimization as DRM. It aims to improve the robustness of adaptation to the tail tasks over iteration. No explicit form exists as the tail task distribution is θ -dependent. The optimization objective is derived as the CVaR $_{\alpha}(\theta)$ ³³:

$$\min_{\theta \in \Theta} \text{CVaR}_{\alpha}(\theta) := \mathbb{E}_{p_{\alpha}(\tau; \theta)} [\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta)], \quad (26)$$

where we write $p_\alpha(\tau; \theta)$ to express the $(1 - \alpha)$ proportional worst case for easier formulation. In other words, $\mathbb{E}_{p_\alpha(\tau; \theta)}[\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)]$ also relates to the task distribution with constraints. Also note that when α approaches 1, the problem degenerates to the worst-case risk minimization.

This work retains the setup in work¹⁴ and picks up the Top- \mathcal{B} in optimization, which corresponds to sample average Monte Carlo of CVaR $_\alpha$. This implies that the actual task batch to evaluate is $\frac{\mathcal{B}}{1-\alpha}$. And for fair comparison with MPTS and light computations, we retain the Monte Carlo estimator for the risk quantile in implementation. To ensure stable training, in all benchmarks, we keep the actual task batch $\hat{\mathcal{B}} = 2\mathcal{B}$ to evaluate and discard the easiest half before the machine learner’s optimization.

Group Distributionally Robust Risk Minimization (GDRM)⁴⁵. This can be interpreted as a min-max optimization problem. Such a principle⁴⁵ effectively improves robustness in distribution shifts and has shown positive effects on training foundation models^{46,47}. It constructs a collection of uncertainty sets over tasks and results in the optimization objective as follows:

$$\min_{\theta \in \Theta} \sup_{g \in \mathcal{G}} \mathbb{E}_{p_g(\tau)} [\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)], \quad (27)$$

where \mathcal{G} are groups of uncertainty sets, and $p_g(\tau)$ indicates the probability measure over the task group. The operation inside Eq. (27) prioritizes the worst group to optimize in a soft way.

GDRM increases the machine learner’s robustness by assigning more probability mass to worst cases in a reweighted manner. That means in each iteration with the best selected $p_{\hat{g}}(\tau)$, the optimization problem is reduced to

$$\min_{\theta \in \Theta} \mathbb{E}_{p_{\hat{g}}(\tau)} [\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)] = \mathbb{E}_{p(\tau)} \left[\frac{p_{\hat{g}}(\tau)}{p(\tau)} \ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \right], \quad (28)$$

where we use $\omega(\tau) = \frac{p_{\hat{g}}(\tau)}{p(\tau)}$ to denote the weight.

Given a fixed number of tasks, GDRM will heuristically or dynamically group them into clusters and then perform a reweighting mechanism according to the evaluated risk. In task episodic learning, there is no task grouping operation as the task batch is reset after each iteration. And the default computation of task-specific weights is $\omega(\tau_i) = \frac{\exp(\eta \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta))}{\sum_{b=1}^{\mathcal{B}} \exp(\eta \ell(\mathcal{D}_{\tau_b}^Q, \mathcal{D}_{\tau_b}^S; \theta))}$, where η is the temperature parameter and $\{\tau_b\}_{b=1}^{\mathcal{B}}$ is the identifier of the task batch. The implementation detail can be found in https://github.com/kohpangwei/group_DR0.

As revealed in works^{14,97,98}, the heuristic operation as the Evaluate-Rank-Filter or reweighting mechanism in GDRM is widely adopted for approximate optimization. For example, in task robust meta-learning scenarios, the prerequisite step in DR-MAML¹⁴ is to execute gradient updates in the inner loop for all candidate tasks and then screen the tail task subset to meta-optimize according to the evaluation results.

Difficulty-Aware Task Sampler (DATS)⁴⁸. Such a baseline is originally developed for few-shot learning problems of physics-informed neural networks (PINNs). Here, we modify it to satisfy the studied benchmarks in few-shot learning scenarios. The basic idea of scoring task difficulty is to compute the inner product of the task-specific gradient on the support dataset $\mathbf{g}_{t,i} := \nabla_{\theta} \ell(\mathcal{D}_{\tau_{t,i}}^S; \theta_t)$ and the task-average gradient on the query dataset of all tasks $\bar{\mathbf{g}}_t := \frac{1}{\mathcal{B}} \sum_{j=1}^{\mathcal{B}} \nabla_{\theta} \ell(\mathcal{D}_{\tau_{t,j}}^Q; \theta_t)$ in the batch.

In implementation, it computes the gradient inner product $s_{t,i} = \langle \mathbf{g}_{t,i}, \bar{\mathbf{g}}_t \rangle$ and then perform normalization as $\omega(\tau_i) = \frac{\exp(\eta s_{t,i})}{\sum_{i=1}^{\mathcal{B}} \exp(\eta s_{t,i})}$. The explanation for such a task weighting mechanism lies in the fact that higher $s_{t,i}$ reveals a more consistent gradient update direction to reduce the validation loss across all tasks, while lower $s_{t,i}$ probably encounter the conflicting gradient issue. Hence, DATS places higher weight on higher $s_{t,i}$. And the vanilla implementation involves the coefficient over the exponential term, and the default is $\frac{1}{\mathcal{B}}$ for all tasks, i.e., the uniform-KL regularization.

Online Hard Task Mining Sampler (OHTM)⁴⁹. The OHTM strategy selects the most challenging tasks from the set of tasks already encountered. To adapt OHTM for meta-learning, Kumar et al.⁴⁹ implements a hybrid scheme: half of each meta-batch is drawn using the OHTM sampler, while the other half is chosen uniformly at random. This design ensures diversity by including a broad range of tasks rather than limiting training to previously observed ones.

Task difficulty Prioritized Sampler (TDPS)⁵⁰. The vanilla Adaptive Sampler⁵⁰ is primarily designed to improve the overall generalization of meta-learning methods, identifying task difficulty, task diversity, and task entropy as key factors for acquiring high-quality tasks. It integrates multiple prioritized sampling criteria into a combined metric. To align with the objective of robust optimization and achieve fair comparison, we adopt only the difficulty-prioritized module of TDPS as our baseline and

use the difficulty-aware module. The core idea is to compute gradients on both the support and query sets and quantify their discrepancy, where the gradient inconsistency within \mathcal{D}_τ serves as a proxy for task difficulty in optimization.

Other Representative Sampling Strategies. Beyond the baselines adopted in this study, other domain-specific task sampling strategies have been proposed with distinct optimization objectives. For instance, Liu et al.³⁷ introduce greedy class-pair (GCP) sampling, which emphasizes class-based adaptive sampling for few-shot image classification by actively constructing challenging tasks for meta-learning. In contrast, our setting operates at the instance level without altering the task construction pipeline. Another representative method, probabilistic active meta-learning (PAML)³¹, aims to improve data efficiency by inferring task embeddings to quantify informativeness. However, PAML is not designed for robust optimization and is more suitable for robotic system identification problems. Adaptive task scheduler³⁴ also adopts the gradient inner product between the query and support dataset to reweigh task losses, and DATS and TDPS share a similar motivation. Given this, we include the previously mentioned baselines in comparison.

C Task Construction and Identifiers

Here we refer to the variables that sufficiently configure a task as the task identifier τ . In other literature work, these task identifiers can be viewed as the task representations in a lower dimensional space. To clarify these concepts, we provide more explanations in specific scenarios.

C.1 Tasks with Explicit Identifiers

K-shot Sinusoid Regression. In this setup³², meta learners aim at quickly adapting the model to an unseen function $f(x) = a \sin(x - b)$ with the help of K data points randomly sampled from the function. This case treats the amplitude and phase variables (a, b) as the task identifier to configure the task. And the task distribution is induced by the uniform distribution over the task identifier.

Meta Reinforcement Learning. Here, we take the ReacherPos task as an example. The goal of the robot arm is to reach an unobserved target location $[x_1, x_2]$. The end-effector position of the robot arm is initialized randomly, and the step-wise reward corresponds to the feedback to the agent after each move based on its distance to the target location. As the task distribution is specified by a uniform distribution over the target location, $\tau = [x_1, x_2]$ can be viewed as the task identifier. Similarly, we vary physics parameters in simulators to generate diverse MDPs. This constitutes different meta RL benchmarks.

Domain Randomization. DR is a promising paradigm to achieve zero-shot adaptation in unseen scenarios, which is widely adopted in robotics³⁹ and computer vision⁶⁵. The basic idea is to train the machine learner in a distribution over environments and then directly apply the learned model to new ones.

Table 4: **Benchmarks with Explicit Task Identifiers.** Here, we list the detail information about the task identifier to induce the task distribution.

Benchmarks	Identifier Meaning	Identifier Range
K-shot sinusoid regression	amplitude and phase (a, b)	$[0.1, 5.0] \times [0, \pi]$
Meta-RL: HalfCheetahMassVel	mass and velocity (m, v)	$[0.75, 1.25] \times [0, 2.0]$
Meta-RL: HalfCheetahVel	velocity v	$[0, 2.0]$
Meta-RL: ReacherPos	goal location (x_1, x_2)	$[-0.2, 0.2] \times [-0.2, 0.2]$
Meta-RL: Walker2dMassVel	mass and velocity (m, v)	$[0.75, 1.25] \times [0, 2.0]$
Meta-RL: Walker2dVel	velocity v	$[0, 2.0]$
DR: LunarLander	main engine strength s	$[4, 20]$
DR: ErgoReacher	joint damping d and max torque t ($\times 4$ joints)	$[0.1, 2.0] \times [2, 20]$

As noted in the main paper, we suppose that the task identifier contains semantics that reflects the difficulty of tasks to resolve and the adaptation risk function is smooth with respect to the identifier. In total, we summarize these benchmarks with explicit task identifiers in Table 4.

C.2 Tasks with Implicit Identifiers

As previously mentioned, we assume the existence of a statistical correlation between task identifiers and adaptation risk values given a specific adaptive machine learner. This implies that the task identifier preserves precise semantics about the task information. These provide the basis for establishing the RPM from the coupled dataset $\{[\tau_i, \ell_i]\}_{i=1}^B$.

Nevertheless, in several scenarios, it is intractable to access the explicit task identifier. For example, in few-shot image classification, the task information is just the coupled support and query dataset $\mathcal{D}_\tau = \mathcal{D}_\tau^S \cup \mathcal{D}_\tau^Q$. Similarly, in SFT for LLMs, the task can be in the form of the QA pair $\mathcal{D}_\tau = \mathcal{D}_\tau^Q$. There is no explicit representation method, such as τ , for these tasks, which brings difficulty in building up the RPM. Retaining the prior notation, the episodic task batch can be written as $\hat{H}_t = \{\theta_t, (\tau_{t,i}, \mathcal{D}_{\tau_{t,i}}, \ell_{t,i})\}_{i=1}^B$, where τ of our interest is unobservable. Some experiments in this work, such as few-shot image classification and SFT, encounter such circumstance.

Task Representation through Identifier Inference. To scale our approach under these circumstances, we propose an alternative candidate schema as the complementary. The probabilistic relationship between variables is depicted in Fig. 2. We consider obtaining the implicit identifier through inference from the task dataset. To do so, we include additional module f_ξ with $\xi \in \Xi$ to embed \mathcal{D}_τ^S and \mathcal{D}_τ^Q and further induce a vector $\tau = f_\xi(\mathcal{D}_\tau^S, \mathcal{D}_\tau^Q)$ as the approximate task identifier. These operations imply seeking the appropriate inference module directly influences the RPM’s performance.

Fortunately, there exist pretrained models that enable the task representation to be generalizable to downstream tasks. For example, in the N-way K-shot image classification, the task is in the form of support image-label pairs and the query images and the goal is to assign labels to the query images from the support dataset. With the help of CLIP models¹, for a fixed task in the form of \mathcal{D}_τ , we can access a N vectors $\{z_i\}_{i=1}^N$ by inputting the set of text-form classes $\{\mathcal{C}_i\}_{i=1}^N$ extracted from the support dataset \mathcal{D}_τ^S , i.e., $\text{CLIP}(\{\mathcal{C}_i\}_{i=1}^N) = [\text{CLIP}_{\text{text}}(\mathcal{C}_1), \dots, \text{CLIP}_{\text{text}}(\mathcal{C}_K)] := \tau$. As a result, we can obtain $H_t = \{[\tau_{t,i}, \ell_{t,i}]\}_{i=1}^B$ conditioned on θ_t for feasible task risk functional posterior inference. This helps our approach to circumvent the unavailability of exact task identifiers. And it is plausible for the RPM to optimize in learning $p(\ell|\tau, H_{1:t})$. It is worth noting that this case still prefers lightweight models for identifier inference, and the text encoder of CLIP well satisfies this requirement and can be used in the N-way K-shot image classification. Details on specific task identifier inference modules can be found in Section G.3 and F.

C.3 Scalability with Large Reasoning Models

The rise of large reasoning models (LRMs) this year makes MPTS even more critical in cutting off policy evaluation cost, i.e., expensive agent-environment interactions.

Finite tasks and unobservable variables. The optimization of large reasoning models (LRMs) often requires massive rollouts to validate outcomes, such as in mathematical problem-solving, particularly under the reinforcement learning from verified reward (RLVR) paradigm. In this setting, the prompt in RLVR corresponds to the task in our framework, yet the task dataset is typically finite and lacks explicit identifiers during RL-based finetuning. A common workflow involves sampling a batch of prompts, generating multiple rollouts per prompt, and extracting informative signals for optimization. For instance, in verifiable mathematical problem-solving, the average success rate estimated from multiple rollouts constitutes an unobservable variable for MPTS to predict, as it reflects task difficulty. The model predictive prompt selection (MoPPS)⁹⁹, a variant of MPTS tailored for LRMs, casts each prompt as a bandit problem. It leverages the 0–1 signals from rollout histories to perform online Bayesian inference over the success rate of each prompt. The inferred distribution then serves as a predictive prior for active prompt selection, avoiding direct interactions with the underlying LRM.

Compatibility with curriculum task selection. In the context of RL finetuning for LRMs, MoPPS and related approaches demonstrate that curriculum task selection can be more effective than naïvely prioritizing task difficulty. Incorporating curriculum criteria accelerates finetuning and highlights the versatility of MPTS, which can be seamlessly integrated with loss-oriented task samplers.

D Auto-Encoding Adaptation Risk through Streaming VI

Note that the basis of MPTS is to establish the bridge between the task identifier and the adaptation risk value over the course of the machine learner’s optimization. In other words, we are seeking a lightweight stochastic risk function in Definition 4 to approximate the posterior $p(\ell|\tau, H_{1:t})$ in the task space.

Definition 4 (Stochastic Risk Function) Let \mathfrak{X} denote the index set’s Cartesian product with the task identifier’s dimension $\tau \in \mathbb{N}^d$. For any $k \in \mathbb{N}$ and finite index sequence $\tau_1, \dots, \tau_k \in \mathfrak{X}$, we write some probability measure over \mathbb{R}^k as $\nu_{(\tau_1, \dots, \tau_k)}$. By introducing the probability space $(\Omega_\tau, \mathcal{F}_\theta, \mathcal{P})$ and $\forall \theta \in \Theta$, we can induce a stochastic function $\mathcal{F}_\theta : \mathcal{T} \times \Omega_\tau \mapsto \mathbb{R}^k$, so that $\nu_{(\tau_1, \dots, \tau_k)}(C_1 \times \dots \times C_k) = \mathcal{P}(\mathcal{F}_\theta(\tau_1) \in C_1, \dots, \mathcal{F}_\theta(\tau_k) \in C_k) \forall \tau_i \in \mathfrak{X}$ and $C_i \in \mathbb{R}$.

This section details steps in auto-encoding historical task risk information, parameterizing variational distributions, deriving the approximate optimization objective, and estimating the stochastic gradients of parameters.

D.1 Neural Modules to Parameterize Distributions

Here, we detail the neural modules to parameterize the distributions of interest. For the approximate posterior $q_\phi(z_t|H_t)$ and conditional prior $p(z_t|H_{1:t-1})$, the inputs of the module are a set of task risk pairs. The neural module requires the permutation invariance *w.r.t.* the order of the data points in the set H_t or $H_{1:t-1}$ in Definition 3. Hence, we adopt the DeepSet style neural network⁶¹ to process the collected H_t or $H_{1:t-1}$.

For example, we denote the neural network parameters by $\phi = \{\phi_1, \phi_{2,1}, \phi_{2,2}\}$ together with a mean pooling operator \oplus , we can have:

$$\mathbf{r}_i = h_{\phi_1}(\tau_{k,i}, \ell_{k,i}) \quad \forall i \in \{1, \dots, \mathcal{B}\}, \quad \bar{\mathbf{r}} = \oplus_{i=1}^{\mathcal{B}} \mathbf{r}_i, \quad \boldsymbol{\mu}_\phi = h_{\phi_{2,1}}(\bar{\mathbf{r}}) \text{ and } \boldsymbol{\Sigma}_\phi = h_{\phi_{2,2}}(\bar{\mathbf{r}}), \quad (29)$$

where the output corresponds to $q_\phi(z_t|H_t) = \mathcal{N}(\boldsymbol{\mu}_\phi, \boldsymbol{\Sigma}_\phi)$ (see Fig. 8 for details).

Regarding the task risk functional posterior inference module, this work has a close connection with the NP family^{60,87,89,100–103}. Both handle the set data points in probabilistic inference.

D.2 Formulation of ELBO & Stochastic Gradient Estimates

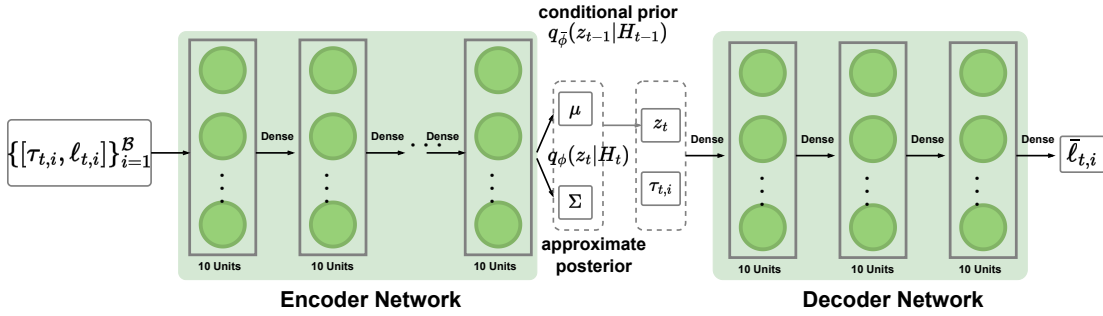


Figure 8: The Encoder-Decoder Neural Network to Parameterize the RPM.

Unlike previous risk minimization principles in task episodic learning, ours include an additional risk predictive module, which guides the task batch sampling. Importantly, we use the latent variable to summarize the historical information information and quantify uncertainty in predicting task-specific adaptation risk. The following details the steps.

$$\mathcal{L}_{\text{ML}}(\psi) := \ln p_\psi(H_t|H_{1:t-1}) = \ln \left[\int p_\psi(H_t|z_t) p(z_t|H_{1:t-1}) dz_t \right] \quad (30a)$$

$$= \ln \left[\int q_\phi(z_t|H_t) \frac{p(z_t|H_{1:t-1})}{q_\phi(z_t|H_t)} p_\psi(H_t|z_t) dz_t \right] \quad (30b)$$

$$\geq \mathbb{E}_{q_\phi(z_t|H_t)} \left[\ln p_\psi(H_t|z_t) \right] - D_{\text{KL}}[q_\phi(z_t|H_t) \parallel p(z_t|H_{1:t-1})] := \mathcal{G}_{\text{ELBO}}(\psi, \phi) \quad (30c)$$

Then, we can rewrite the ELBO with the help of reparameterization trick²³ in Eq. (31).

$$\hat{\mathcal{G}}_{\text{ELBO}}(\psi, \phi) = \mathbb{E}_{q_\phi(z_t|H_t)} \left[\ln p_\psi(H_t|z_t) \right] - D_{\text{KL}}[q_\phi(z_t|H_t) \parallel p(z_t|H_{1:t-1})] \quad (31a)$$

$$= \mathbb{E}_{p(\epsilon)} \left[\ln p_\psi(H_t|g_\phi(\epsilon, H_t)) \right] - D_{\text{KL}}[q_\phi(z_t|H_t) \parallel p(z_t|H_{1:t-1})] \quad (31b)$$

$$\approx \ln p_\psi(H_t|g_\phi(\epsilon, H_t)) - D_{\text{KL}}[q_\phi(z_t|H_t) \parallel p(z_t|H_{1:t-1})], \quad \text{with } \epsilon \sim \mathcal{N}(0, \mathbf{I}_d) \quad (31c)$$

Moreover, we estimate the stochastic gradients *w.r.t.* all model parameters based on the reparameterized latent variable distribution.

$$\nabla_\phi \mathcal{G}_{\text{ELBO}}(\psi, \phi) \approx \nabla_\phi \ln p_\psi(H_t|g_\phi(\epsilon, H_t)) - \frac{1}{2} \nabla_\phi \left(\text{Tr}(\hat{\Sigma}^{-1} \boldsymbol{\Sigma}_\phi) + (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_\phi)^T \hat{\Sigma} (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_\phi) - \ln(\det \boldsymbol{\Sigma}_\phi) \right) \quad (32a)$$

$$\text{with } q_\phi(z_t|H_t) = \mathcal{N}(\boldsymbol{\mu}_\phi, \boldsymbol{\Sigma}_\phi) \text{ and } p(z_t|H_{1:t-1}) = \mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\Sigma}) \quad (32b)$$

$$\nabla_\psi \mathcal{G}_{\text{ELBO}}(\psi, \phi) \approx \nabla_\psi \ln p_\psi(H_t|g_\phi(\epsilon, H_t)) \quad (32c)$$

As illustrated in Eq. (32), one stochastic forward pass is required for gradient estimates in the training process. For flexible implementation, we adopt a β -VAE strategy to turn Eq. (31) into

$$\max_{\psi \in \Psi, \phi \in \Phi} \mathcal{G}_{\text{ELBO}}(\psi, \phi) := \mathbb{E}_{q_\phi(\mathbf{z}_t | H_t)} \left[\sum_{i=1}^{\mathcal{B}} \ln p_\psi(\ell_{t,i} | \tau_{t,i}, \mathbf{z}_t) \right] - \beta D_{KL} \left[q_\phi(\mathbf{z}_t | H_t) \parallel q_{\bar{\phi}}(\mathbf{z}_t | H_{t-1}) \right] \quad (33)$$

D.3 Theoretical Guarantee for Task Difficulties' Scoring with Posterior Inference

Assumption 1 (Lipschitz Continuity) We assume the adaptation risk function $\ell(\cdot; \boldsymbol{\theta})$ reserves the Lipschitz continuity w.r.t. $\boldsymbol{\theta}$ and $\boldsymbol{\tau}$, i.e.,

$$|\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}) - \ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}')| \leq \beta_1 \|\boldsymbol{\theta} - \boldsymbol{\theta}'\| \quad \text{and} \quad |\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}) - \ell(\mathcal{D}_{\tau'}^Q, \mathcal{D}_{\tau'}^S; \boldsymbol{\theta})| \leq \beta_2 \|\boldsymbol{\tau} - \boldsymbol{\tau}'\|, \quad (34)$$

where $\forall \{\boldsymbol{\theta}, \boldsymbol{\theta}'\} \in \Theta$ and $\forall \{\boldsymbol{\tau}, \boldsymbol{\tau}'\} \in \mathcal{T}$ with Lipschitz constants β_1 and β_2 .

Assumption 2 (Bounded Sample Gradient) We assume the norm of the adaptation risk function's gradient $\nabla \ell(\cdot; \boldsymbol{\theta}_t)$ is bounded:

$$\sup_{\tau \in \mathcal{T}} \|\nabla_{\boldsymbol{\theta}} \ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}_t)\|_2 < G_t \quad \text{and} \quad \sup_{\tau \in \mathcal{T}, t \in \mathbb{N}_+} \|\nabla_{\boldsymbol{\theta}} \ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \boldsymbol{\theta}_t)\|_2 < G, \quad (35)$$

where G_t is a positive constant and G is a overall bound.

Assumption 3 (Sub-Gaussian Stochastic Gradient) The stochastic gradient $\tilde{\mathbf{g}} := \mathbf{g} + \boldsymbol{\epsilon}$ for the machine learner's adaptation at t -th iteration is σ -sub-Gaussian, which means:

$$\mathbb{E} \left[\exp(\eta \mathbf{v}^T \boldsymbol{\epsilon}) \right] \leq \exp \left(\frac{\eta^2 \sigma^2 \|\mathbf{v}\|_2^2}{2} \right) \quad \forall \eta \in \mathbb{R} \text{ and } \mathbf{v} \in \mathbb{R}^d, \quad (36)$$

where $\mathbb{E}[\tilde{\mathbf{g}}] = \mathbf{g}$, $\mathbb{E}[\|\tilde{\mathbf{g}} - \mathbf{g}\|_2^2] \leq \sigma^2$ and $\sigma \in \mathbb{R}_+$.

Given the Assumption 3 and the Chernoff bound¹⁰⁴, we can have the concentration inequality as:

$$\mathbb{P}(\|\tilde{\mathbf{g}} - \mathbf{g}\|_2 \geq t) \leq 2 \exp \left(-\frac{t^2}{2\sigma^2} \right) \quad \forall t \in \mathbb{R}. \quad (37)$$

Theorem 1 (Provably Approximately Invariant Task Difficulties) Given arbitrary K data points $\{(\tau_i, \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t))\}_{i=1}^K$, the adaptation gradient $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)$ as a σ -sub-Gaussian random variable and $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)$, we denote the relative difficulty via the difference $\Delta_{ij}(\boldsymbol{\theta}_{t+1}) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_{t+1}) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_{t+1})$ and $\Delta_{ij}(\boldsymbol{\theta}_t) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_t)$ between t -th and $(t+1)$ -th iterations, and the gradient difference as $\mathbf{v}_{ij} := \nabla_{\boldsymbol{\theta}} \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t) - \nabla_{\boldsymbol{\theta}} \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_t)$.

Under Assumption 1/2/3, the set of rank-preserving variable $E_{ij} := \mathbb{1}[\text{sign}(\Delta_{ij}(\boldsymbol{\theta}_{t+1})) = \text{sign}(\Delta_{ij}(\boldsymbol{\theta}_t))]$ satisfies the probability inequality:

$$\mathbb{P}(\bigcap_{i < j} E_{ij}) \geq 1 - \xi,$$

when $\eta_t \leq \frac{\delta_t}{2G_t M_t + \sqrt{8\sigma^2 G_t^2 \ln \left(\frac{K(K-1)}{2\xi} \right)}}$ with G_t in Assumption 2, $\delta_t := \min_{i \neq j} |\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \boldsymbol{\theta}_t)| \in \mathbb{R}_+$, the stochastic gradient norm $M_t := \|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)\|_2$.

The purpose of this part is to uncover the mechanism of the RPM in amortized evaluation of adaptation risk values and scoring the difficulty of tasks. The function of the RPM relies on Assumptions 1/2/3 and the posterior inference $p(\ell | \tau, H_{1:t}; \boldsymbol{\theta}_t)$ from the historical risk information $H_{1:t}$. The foundation of predicting the outcome of optimization in a rough granularity lies in the Theorem 1, and we detail the proof of such a theorem as below.

①. Any-Shot Adaptation After One-step Gradient Descent.

Here, we consider a set of data points for the RPM $\{(\tau_i, \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}_t))\}_{i=1}^K$ under an arbitrary fixed machine learner $\boldsymbol{\theta}_t$, where tasks in the set $\{\tau_i\}_{i=1}^K$ are randomly sampled from $p(\tau)$. Without loss of generality, we can assume that the adaptation risk values satisfy a rank ordering:

$$\ell(\mathcal{D}_{\tau_1}^Q, \mathcal{D}_{\tau_1}^S; \boldsymbol{\theta}_t) \geq \ell(\mathcal{D}_{\tau_2}^Q, \mathcal{D}_{\tau_2}^S; \boldsymbol{\theta}_t) \geq \dots \geq \ell(\mathcal{D}_{\tau_K}^Q, \mathcal{D}_{\tau_K}^S; \boldsymbol{\theta}_t). \quad (38)$$

The gradient descent as fast adaptation is denoted by:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t). \quad (39)$$

After the above operator, we can obtain another set of data points for the updated RPM $\{(\tau_i, \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_{t+1}))\}_{i=1}^K$.

②. Changes of Adaptation Risk Values and Pairwise Ranks.

Based on the Assumption 1, we can perform local approximation over $\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta)$ with the help of first-order Talor expansion w.r.t. the θ_t :

$$\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_{t+1}) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \eta_t \nabla_{\theta} \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t)^T \nabla_{\theta} \mathcal{L}(\theta_t) + \mathcal{O}(\eta_t^2 \|\nabla_{\theta} \mathcal{L}(\theta_t)\|_2^2) \quad (40a)$$

$$\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_{t+1}) \approx \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \nabla_{\theta} \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t)^T \nabla_{\theta} \mathcal{L}(\theta_t) \quad \forall i \in \{1, 2, \dots, K\}. \quad (40b)$$

One straightforward way to assess the task difficulty is to compare arbitrary paired tasks $\{\tau_i, \tau_j\}$'s adaptation risk values $\{\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t), \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t)\}$ with $i < j$. Then, we can estimate the relative difficulty via the difference as:

$$\Delta_{ij}(\theta_{t+1}) \approx \Delta_{ij}(\theta_t) - \eta_t (\nabla_{\theta} \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \nabla_{\theta} \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t))^T \nabla_{\theta} \mathcal{L}(\theta_t), \quad (41)$$

where we denote the relative difficulty via the difference as $\Delta_{ij}(\theta_{t+1}) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_{t+1}) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_{t+1})$ and $\Delta_{ij}(\theta_t) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t)$ between t -th and $(t+1)$ -th iterations. As $\Delta_{ij}(\theta_t)$ is positive, one feasible condition for $\Delta_{ij}(\theta_{t+1}) \in \mathbb{R}_+$ is:

$$\Delta_{ij}(\theta_{t+1}) \approx \Delta_{ij}(\theta_t) - \eta_t (\nabla_{\theta} \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \nabla_{\theta} \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t))^T \nabla_{\theta} \mathcal{L}(\theta_t) \quad (42a)$$

$$\geq \Delta_{ij}(\theta_t) - 2\eta_t G_t M > 0, \quad \Rightarrow \eta_t < \frac{\Delta_{ij}(\theta_t)}{2G_t M_t} \quad \text{with} \quad M_t := \|\nabla_{\theta} \mathcal{L}(\theta_t)\|_2. \quad (42b)$$

The above implies that when the learning rate η_t in gradient step is smaller enough, the relative difficulty between the task i and j can be preserved after the machine learner's update with the Assumption 2.

③. Probabilistic Inequality with a Nearly Invariant Ranking Guarantee.

In practice, the stochastic gradient descent is performed, which means the gradient is a random variable $\nabla_{\theta} \mathcal{L}(\theta_t) = \mathbf{g}_t + \epsilon$ with $\mathbb{E}[\epsilon] = 0$, $\mathbb{E}[\|\epsilon\|_2^2] < \sigma^2$ and $\mathbf{g}_t = \mathbb{E}[\nabla_{\theta} \mathcal{L}(\theta_t)]$. Meanwhile, we denote the gradient difference by $\mathbf{v}_{ij} := \nabla_{\theta} \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \nabla_{\theta} \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t)$, which leads to:

$$\|\mathbf{v}_{ij}\|_2 \leq 2 \sup_{\tau \in \mathcal{T}} \|\nabla_{\theta} \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta)\|_2 < 2G_t, \quad (43)$$

according to the Assumption 2. Another variable is introduced as the minimum separation between arbitrary paired adaptation risk values:

$$\delta_t := \min_{i \neq j} |\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t)| \in \mathbb{R}_+. \quad (44)$$

Still, to make sure the invariant rank, one necessary condition can be:

$$\Delta_{ij}(\theta_t) - \eta_t \mathbf{v}_{ij}^T \mathbf{g}_t \geq 0. \quad (45)$$

And the above inequality reasonably holds when $\eta_t \mathbf{v}_{ij}^T \mathbf{g}_t < \delta_t$. Here, we define the random event $E_{ij} := \mathbb{1}[\text{sign}(\Delta_{ij}(\theta_{t+1})) = \text{sign}(\Delta_{ij}(\theta_t))]$ from the task pair together with $E_{ij}^c := \mathbb{1}[\text{sign}(\Delta_{ij}(\theta_{t+1})) \neq \text{sign}(\Delta_{ij}(\theta_t))]$. With the help of σ -sub-Gaussian property in Assumption 3, we can bound the case of the rank flipping as (note some critical conditions that $\mathbf{v}_{ij}^T \mathbf{g}_t \in \mathbb{R}_+$ and $\eta_t \mathbf{v}_{ij}^T \mathbf{g}_t < \delta_t$ as the learning rate η_t can be typically smaller enough):

$$\mathbb{P}(E_{ij}^c) = \mathbb{P}(\eta_t \mathbf{v}_{ij}^T \nabla_{\theta} \mathcal{L}(\theta_t) \geq \delta_t) \leq \exp\left(-\frac{(\delta_t - \eta_t \mathbf{v}_{ij}^T \mathbf{g}_t)^2}{2\eta_t^2 \sigma^2 \|\mathbf{v}_{ij}\|_2^2}\right) < \exp\left(-\frac{(\delta_t - 2\eta_t G_t M_t)^2}{8\eta_t^2 \sigma^2 G_t^2}\right) \quad (46a)$$

$$\mathbb{P}\left(\bigcup_{i < j} E_{ij}^c\right) \leq \sum_{i < j} \mathbb{P}(E_{ij}^c) \leq \frac{K(K-1)}{2} \exp\left(-\frac{(\delta_t - 2\eta_t G_t M_t)^2}{8\eta_t^2 \sigma^2 G_t^2}\right) \quad (46b)$$

$$\mathbb{P}\left(\bigcap_{i < j} E_{ij}\right) = 1 - \mathbb{P}\left(\bigcup_{i < j} E_{ij}^c\right) \geq 1 - \frac{K(K-1)}{2} \exp\left(-\frac{(\delta_t - 2\eta_t G_t M_t)^2}{8\eta_t^2 \sigma^2 G_t^2}\right) \geq 1 - \xi. \quad (46c)$$

The condition for the above inequality holds is $\eta_t \leq \frac{\delta_t}{2G_t M_t + \sqrt{8\sigma^2 G_t^2 \ln\left(\frac{K(K-1)}{2\xi}\right)}}$. With the above **steps ①-③** and corresponding conditions, we complete the proof.

E Convergence Analysis

Predictive Foundation of MPTS. In brief, the foundation of MPTS is to use the risk history and generalization results under θ_t to approximately rank the difficulty of task samples in $(t+1)$ -th iteration. Note that precise sample average Monte Carlo for CVaR relies on the exact evaluation and access to the $(t+1)$ -th adaptation risk values in $\hat{\mathcal{B}}$ tasks to pick up the Top- \mathcal{B} ones. MPTS tries to circumvent the complete evaluation and uses the t -th adaptation risk batch $\{\ell_{t,i}\}_{i=1}^{\mathcal{B}}$ to train a generative model with $H_{1:t-1}$ to extract the prior to forecast the $\hat{\mathcal{B}}$ tasks' difficulty order for $(t+1)$ -th subset selection. (Importantly, such a setup necessitates MPTS when accessing expensive evaluation of $\{\ell_{t+1,i}\}_{i=1}^{\hat{\mathcal{B}}}$ is prohibited.)

The theoretical analysis focuses on the rank-flipping that occurs during the model's update $\theta_t \rightarrow \theta_{t+1}$, and we will demonstrate that, under certain conditions, this does not affect the convergence of the adopted strategy to surrogate MC-CVaR in a biased manner. The candidate task pool is $\bar{\mathcal{T}}_{t+1}^{\hat{\mathcal{B}}}$ used for the subset selection and robust adaptation at the $(t+1)$ -th iteration. The Top- \mathcal{B} task batch under θ_t 's risk evaluation is denoted by $\bar{\mathcal{T}}_t^{\mathcal{B}} \subset \bar{\mathcal{T}}_{t+1}^{\hat{\mathcal{B}}}$, and its complementary set is $\mathcal{T}_t^C := \bar{\mathcal{T}}_{t+1}^{\hat{\mathcal{B}}} \setminus \bar{\mathcal{T}}_t^{\mathcal{B}}$. The ground truth Top- \mathcal{B} task batch under θ_{t+1} 's risk evaluation is $\mathcal{T}_{t+1}^{\mathcal{B}}$.

Assumption 4 (Margin Anti-Concentration) We assume $\Delta_{ij}(\theta_t) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t)$ has a density uniformly bounded by ρ in a neighborhood of zero. That is, $\forall \epsilon > 0$, the probability inequality holds:

$$\mathbb{P}(|\Delta_{ij}(\theta_t)| \leq \epsilon) \leq \rho\epsilon. \quad (47)$$

The above assumption is reasonable in online robust optimization scenarios. If a large mass of task pairs has nearly identical loss, this suggests the Top- \mathcal{B} set is unstable, making it difficult for any difficulty prioritized sampling algorithm to approximate the CVaR subset stably. Anti-concentration ensures that near-ties are probabilistically rare, hence rank changes are infrequent over iterations.

Lemma 1 (Misranked Subset Quantity) In the presence of adaptation risk value, let N_f be the number of order flipped cross-pairs between the ground-truth Top- \mathcal{B} subset $\bar{\mathcal{T}}_t^{\mathcal{B}}$ under θ_t and the remainder of the candidate tasks $\mathcal{T}_t^C := \bar{\mathcal{T}}_{t+1}^{\hat{\mathcal{B}}} \setminus \bar{\mathcal{T}}_t^{\mathcal{B}}$ when the machine learner's parameter changes from θ_t to θ_{t+1} . We denote the number of tasks that change Top- \mathcal{B} membership by $m_{t+1} = |\bar{\mathcal{T}}_t^{\mathcal{B}} \triangle \mathcal{T}_{t+1}^{\mathcal{B}}|$, where $\mathcal{T}_{t+1}^{\mathcal{B}}$ is the ground-truth Top- \mathcal{B} subset under θ_{t+1} . Then the inequality holds: $m_{t+1} \leq 2N_f$.

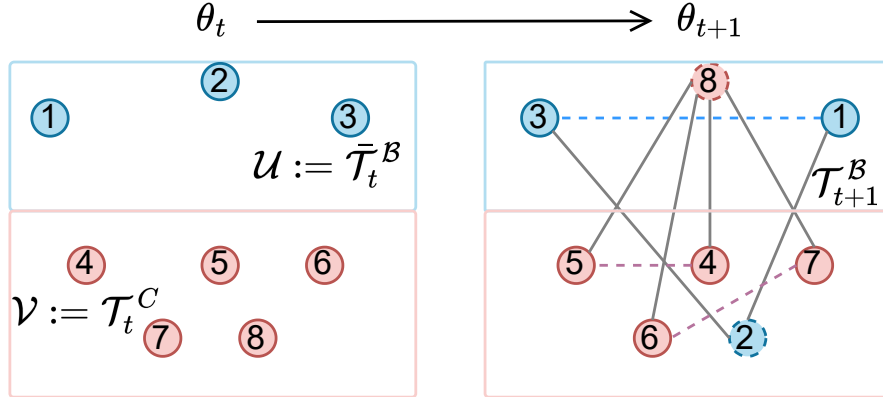


Figure 9: **Rank Flipping after One-Step Update.** We number tasks from the sampled pseudo task batch $\bar{\mathcal{T}}_{t+1}^{\hat{\mathcal{B}}}$ for the $(t+1)$ -th subset selection. The left figure partitions the set into $\mathcal{U} \cup \mathcal{V}$ as illustrated in the Proof of Lemma 1, with the blue ones, i.e., 1, 2, 3 as Top- \mathcal{B} tasks evaluated under θ_t and the red ones as the complementary set. Note that tasks 2 and 8 are swapped in relative Top- \mathcal{B} ranking after one-step model update; hence, the right figure depicts the Top- \mathcal{B} tasks as 1, 3, and 8. The gray lines indicate possible order-flipped cross-pairs between \mathcal{U} and \mathcal{V} . The blue and the red lines are probably flipped cross-pairs in separate partitioned sets' induced graphs.

Proof. Now we can construct a bipartite graph $\mathcal{G} = (\mathcal{U}, \mathcal{V}, \mathcal{E})$, where \mathcal{U} denotes the ground-truth Top- \mathcal{B} subset $\bar{\mathcal{T}}_t^{\mathcal{B}}$ while \mathcal{V} is the remainder of the candidate tasks $\mathcal{T}_t^C := \bar{\mathcal{T}}_{t+1}^{\hat{\mathcal{B}}} \setminus \bar{\mathcal{T}}_t^{\mathcal{B}}$. Each task in the set corresponds to a vertex in the graph. Also note that there are exactly $\mathcal{B}(\hat{\mathcal{B}} - \mathcal{B})$ ordered task cross-pairs.

For each ordered cross pair $(i \in \mathcal{U}; j \in \mathcal{V})$, which flips, i.e. the sign of $\Delta_{ij}(\theta_t)$ changes between θ_t and θ_{t+1} , we place an undirected edge e_{ij} between vertices i and j , and this constitutes the edge set \mathcal{E} . Further, the subset of vertices, where the order

flipping occurs, can be defined as $\mathcal{V}_E = \{\tau \in \mathcal{T}_{t+1}^{\hat{\mathcal{B}}} | \exists e \in \mathcal{E} \text{ with } \tau \in e\}$, which describes the set of tasks that participate in at least one flipped cross-pair. As a result, we can have $\mathcal{T}_t^{\mathcal{B}} \Delta \mathcal{T}_{t+1}^{\mathcal{B}} \subseteq \mathcal{V}_E$. Obviously, $m_{t+1} \leq |\mathcal{V}_E|$.

The total number of flipping pairs is N_f . As illustrated in Fig. 9, the flipped ordered cross-pairs comprises three disjoint cases or edge partition in \mathcal{G} : (i) order flipping between $\mathcal{T}_t^{\mathcal{B}}$, (ii) order flipping between $\mathcal{T}_t^{\mathcal{C}}$, and (iii) order flipping between $\mathcal{T}_t^{\mathcal{B}}$ and $\mathcal{T}_t^{\mathcal{C}}$. This indicates that $|\mathcal{E}| \leq N_f$.

Putting the above together, we can have $m_{t+1} \leq |\mathcal{V}_E| \leq 2|\mathcal{E}| \leq 2N_f$. And this establishes the connection between the relative task difficulty ranking and the direct task difficulty ranking. And this completes the proof of **Lemma 1**.

Lemma 2 (Rank-Preserving Bound in Expectation) With the risk difference notation $\Delta_{ij}(\theta_t) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t)$ and $\Delta_{ij}(\theta_{t+1}) = \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_{t+1}) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_{t+1})$, when the (i, j) cross-pair flips its order from θ_t to θ_{t+1} , we conclude that $|\Delta_{ij}(\theta_t)| \leq |\Delta_{ij}(\theta_{t+1}) - \Delta_{ij}(\theta_t)|$.

Proof. Now let us consider limited scenarios for the occurrence of the order flipping event $\bar{E}_{ij} := \mathbb{1}[\text{sign}(\Delta_{ij}(\theta_{t+1})) \neq \text{sign}(\Delta_{ij}(\theta_t))]$.

In scenario (1), where $\Delta_{ij}(\theta_t) < 0$ and $\Delta_{ij}(\theta_{t+1}) \geq 0$ hold, this implies that:

$$\underbrace{\Delta_{ij}(\theta_{t+1}) - \Delta_{ij}(\theta_t)}_{\text{Positive Value}} \geq \underbrace{-\Delta_{ij}(\theta_t)}_{\text{Positive Value}}. \quad (48)$$

In scenario (2), where $\Delta_{ij}(\theta_t) \geq 0$ and $\Delta_{ij}(\theta_{t+1}) < 0$ hold, this implies that:

$$\underbrace{\Delta_{ij}(\theta_{t+1}) - \Delta_{ij}(\theta_t)}_{\text{Non-Positive Value}} \leq \underbrace{0 - \Delta_{ij}(\theta_t)}_{\text{Non-Positive Value}}. \quad (49)$$

Putting the above together, we can derive that $|\Delta_{ij}(\theta_t)| \leq |\Delta_{ij}(\theta_{t+1}) - \Delta_{ij}(\theta_t)|$. And this completes the proof of **Lemma 2**.

Next, we analyze the term $|\Delta_{ij}(\theta_{t+1}) - \Delta_{ij}(\theta_t)|$ and associate it with $|\Delta_{ij}(\theta_t)|$:

$$|\Delta_{ij}(\theta_{t+1}) - \Delta_{ij}(\theta_t)| = |\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_{t+1}) - \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_{t+1}) + \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t)| \quad (50a)$$

$$\leq |\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_{t+1}) - \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta_t)| + |\ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_{t+1}) - \ell(\mathcal{D}_{\tau_j}^Q, \mathcal{D}_{\tau_j}^S; \theta_t)| \quad (50b)$$

$$\leq 2G_t \|\theta_{t+1} - \theta_t\|_2 \quad (50c)$$

$$\Rightarrow |\Delta_{ij}(\theta_t)| \leq 2G_t \|\theta_{t+1} - \theta_t\|_2. \quad (50d)$$

Connection between Rank Flipping and Gradient Difference. From **Lemma 2** together with the anti-concentration Assumption 4, we can bound the probability of the order flipping event \bar{E}_{ij} with the one-step iteration:

$$\mathbb{P}(\bar{E}_{ij} | \theta_t, \theta_{t+1}) \leq \mathbb{P}(|\Delta_{ij}(\theta_t)| \leq 2G_t \|\theta_{t+1} - \theta_t\|_2) \leq 2\rho G_t \|\theta_{t+1} - \theta_t\|_2. \quad (51)$$

Note that the graph only cares about pairwise relation between the Top- \mathcal{B} tasks and other $\hat{\mathcal{B}} - \mathcal{B}$ tasks. Now in a population level, the scale of plausible pairs in the task set $\mathcal{T}_{t+1}^{\hat{\mathcal{B}}}$ is $\mathcal{B}(\hat{\mathcal{B}} - \mathcal{B})$ and $|\mathcal{V}_E| = \sum_{i,j} \bar{E}_{ij}$. Also, with the linearity of expectation over \bar{E}_{ij} , **Lemma 1** and Eq. (50)/(51), we can bound the flipping count expectation:

$$\mathbb{E}[m_{t+1}] \leq \mathbb{E}[|\mathcal{V}_E|] \leq 2\rho \mathcal{B}(\hat{\mathcal{B}} - \mathcal{B}) G_t \mathbb{E}[\|\theta_{t+1} - \theta_t\|_2] \quad (52a)$$

$$r_{t+1} := \frac{\mathbb{E}[m_{t+1}]}{\mathcal{B}} \leq 2\rho(\hat{\mathcal{B}} - \mathcal{B}) G_t \mathbb{E}[\|\theta_{t+1} - \theta_t\|_2] \leq 2\rho(\hat{\mathcal{B}} - \mathcal{B}) G_t^2 \mathcal{O}(\eta_t). \quad (52b)$$

Next, we need to show that the selection bias from rank flipping vanishes quickly and the gradient norm declines to ensure the convergence. Recall that the expectation form of CVaR can be written as $\mathbb{E}_{p_\alpha(\tau; \theta_{t+1})}[\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)]$, and MC-CVaR is its unbiased estimate from the subset $\mathcal{T}_{t+1}^{\hat{\mathcal{B}}} \subset \mathcal{T}_{t+1}^{\hat{\mathcal{B}}}$. Accordingly, the expectation form of perturbed CVaR can be written as $\mathcal{L}(\theta) = \mathbb{E}_{q(\tau; \theta_{t+1})}[\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)]$ and $q(\tau; \theta_{t+1})$ denotes the distribution of biased subset due to miscalibrated cases, which produced the biased estimate of expected CVaR from $\tilde{\mathcal{T}}_t^{\mathcal{B}} \subset \mathcal{T}_{t+1}^{\hat{\mathcal{B}}}$.

Given the latest machine learner parameter θ_{t+1} , we can derive the equation after one-step gradient update as:

$$\text{One-Step Gradient Update: } \theta_{t+2} = \theta_{t+1} - \eta_{t+1} \underbrace{\nabla_{\theta} \mathcal{L}(\theta_{t+1})}_{\text{Perturbed Gradient After Rank-Flipping}} \quad (53a)$$

$$\text{with Gradient Decomposition } \nabla_{\theta} \mathcal{L}(\theta_{t+1}) = \underbrace{\bar{\mathbf{g}}_{t+1}}_{\text{Unbiased Average Gradient}} + \underbrace{\Delta \mathbf{g}_{t+1}}_{\text{Gradient Difference}}, \quad (53b)$$

where we refer to the average task gradient after rank-flipping as the *perturbed gradient*. The unbiased average gradient are estimated from $\mathcal{T}_{t+1}^{\mathcal{B}}$, the ground-truth Top- \mathcal{B} subset under θ_{t+1} . And term $\Delta \mathbf{g}_{t+1} = \frac{2}{m_{t+1}} \left(\sum_{\tau \in \mathcal{T}_t^{\mathcal{B}} \setminus \mathcal{T}_{t+1}^{\mathcal{B}}} \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta_{t+1}) - \sum_{\tau \in \mathcal{T}_{t+1}^{\mathcal{B}} \setminus \mathcal{T}_t^{\mathcal{B}}} \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta_{t+1}) \right)$ denotes the gradient difference between rank-flipping tasks after model update. Eq. (53) will work for the convergence analysis $\forall t \in \mathbb{N}_+$ in the following contents.

Lemma 3 (Misranking Acute-Angle) Let $\bar{\mathbf{g}}_t$ and $\Delta \mathbf{g}_t$ respectively denote the unbiased Monte Carlo estimate of CVaR and the gradient difference between it and the biased gradient $\nabla_{\theta} \mathcal{L}(\theta_t)$. For a given $c \in (0, 1)$ and arbitrary small $\epsilon_* > 0$, there exists T such that $\forall t \geq T$, the following dichotomy holds: either (i) $\|\bar{\mathbf{g}}_t\|_2 \leq \epsilon_*$, or (ii) $\|\Delta \mathbf{g}_t\|_2 \leq c\|\bar{\mathbf{g}}_t\|_2$.

Proof. In case (i) if $\|\bar{\mathbf{g}}_t\|_2 \leq \epsilon_*$, with arbitrary small ϵ_* , we can see that $\lim_{t \rightarrow \infty} \|\bar{\mathbf{g}}_t\|_2 = 0$. Also, recall the truth in Eq. (52), we can find that $\lim_{t \rightarrow \infty} r_{t+1} = 0$ when $\lim_{t \rightarrow \infty} \eta_t = 0$. Then we can have $\lim_{t \rightarrow \infty} \|\nabla_{\theta} \mathcal{L}(\theta_t)\|_2 \leq \lim_{t \rightarrow \infty} \|\bar{\mathbf{g}}_t\|_2 + \|\Delta \mathbf{g}_t\|_2 \leq \lim_{t \rightarrow \infty} \|\bar{\mathbf{g}}_t\|_2 + r_{t+1}G$, which implies the stationarity of $\|\nabla_{\theta} \mathcal{L}(\theta_t)\|_2$ and already guarantees the convergence of algorithms with diminishing rank flipping.

In case (ii) if $\|\bar{\mathbf{g}}_t\|_2 > \epsilon_*$, we perform analysis as: recall the previous finding about the misranking rate $r_{t+1} := \frac{\mathbb{E}[m_{t+1}]}{\mathcal{B}} \leq 2\rho(\hat{\mathcal{B}} - \mathcal{B})G_t^2\mathcal{O}(\eta_t)$. When $\lim_{t \rightarrow \infty} \eta_t = 0$, it is also trivial to see that $\lim_{t \rightarrow \infty} r_{t+1} = 0$. Hence, choosing a large enough T such that for all $t > T$, we can have $r_{t+1} \leq \frac{c}{G} \epsilon_*$. Then the inequality holds:

$$\|\Delta \mathbf{g}_t\|_2 \leq r_{t+1}G_t \leq r_{t+1}G \leq c\epsilon_*. \quad (54)$$

Then, together with the case condition $\|\bar{\mathbf{g}}_t\|_2 > \epsilon_*$, we can have:

$$\|\Delta \mathbf{g}_t\|_2 \leq c\epsilon_* \leq c\|\bar{\mathbf{g}}_t\|_2 \Rightarrow \|\Delta \mathbf{g}_t\|_2 \leq c\|\bar{\mathbf{g}}_t\|_2. \quad (55)$$

Finally, the dichotomy holds for the studied two cases. And this completes the proof of **Lemma 3**.

Theorem 2 (Convergence with Diminishing Rank Flipping) Suppose there exists $c \in (0, 1)$ and $T \geq 0$ such that $\forall t \geq T$, the inequality holds: $\|\Delta \mathbf{g}_t\|_2 \leq c\|\bar{\mathbf{g}}_t\|_2$. Given the Assumptions 1/2/3/4, and the appropriate construction of the step sizes $\{\eta_t\}$ from Theorem 1, we conclude that: $\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\theta} \mathcal{L}(\theta_t)\|_2] = 0$. Hence, the iteration converges to first-order stationary points in expectation.

Proof. The **Theorem 1** indicates that the probability inequality holds when the learning rate follows that: $\xi_t \geq$

$$\frac{K(K-1)}{2} \exp\left(-\frac{(\frac{\delta_t}{\eta_t} - 2G_t M_t)^2}{8\sigma^2 G_t^2}\right). \text{ Then for a given } \eta_t, \text{ we can guarantee the probability of rank-preserving for all pairs as}$$

$$1 - \xi_t^* := 1 - \frac{K(K-1)}{2} \exp\left(-\frac{(\frac{\delta_t}{\eta_t} - 2G_t M_t)^2}{8\sigma^2 G_t^2}\right).$$

Since G_t and M_t are bounded, the behavior of the exponent is determined by the coefficient $\frac{\delta_t}{\eta_t}$. According to the Assumption 4, for a sufficiently small constant $\delta_{\min} > 0$, there holds the probability inequality $\mathbb{P}(\delta_t > \delta_{\min}) > 1 - \rho\delta_{\min}$, which means $\delta_t > 0$ stands with a sufficiently large probability with the sampled tasks well scattered in random sampling. Hence, the limit of $\frac{\delta_t}{\eta_t} > \frac{\delta_{\min}}{\eta_t}$ goes to the infinity when $\lim_{t \rightarrow \infty} \eta_t = 0$ in a probabilistic sense. This suggests the tightening bound of rank preserving with the diminishing ξ_t^* .

In practice, we can construct the example of the plausible learning rate scheduling strategy by setting $\eta_t = t^{-b}$ with $b \in (\frac{1}{2}, 1]$, which naturally makes Robbins-Monro conditions hold: $\sum_t \eta_t = \infty$ and $\sum_t \eta_t^2 < \infty$. This decays the learning rate $\lim_{t \rightarrow \infty} \eta_t = 0$ over iteration, which probabilistically reduces the chance of rank flipping cases and tightens the inequality bound in **Theorem 1**. Note that in this theorem, we leave out the case of approaching stationarity, i.e., case (i) when the gradient norm limit approaches zero from **Lemma 3**.

The next step is to perform the Taylor expansion of the MC-CVaR risk function $\mathcal{L}(\theta)$ at the parameter θ_t , which obtains $\mathcal{L}(\theta) = \mathcal{L}(\theta_t) + \bar{\mathbf{g}}_t^T(\theta - \theta_t) + \mathcal{O}(\|\theta - \theta_t\|)$. Then, we input the biased gradient $\nabla_{\theta} \mathcal{L}(\theta_t)$ and the updated parameter to derive

the equation as:

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) = \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \bar{\mathbf{g}}_t^T \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t) + \mathcal{O}(\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|) \quad (56a)$$

$$\Rightarrow \mathcal{L}(\boldsymbol{\theta}_{t+1}) \approx \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \|\bar{\mathbf{g}}_t\|_2^2 - \eta_t \bar{\mathbf{g}}_t^T \Delta \mathbf{g}_t \quad (56b)$$

$$\leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \|\bar{\mathbf{g}}_t\|_2^2 + \eta_t \|\bar{\mathbf{g}}_t\|_2 \|\Delta \mathbf{g}_t\|_2 \quad (56c)$$

$$\leq \mathcal{L}(\boldsymbol{\theta}_t) - \eta_t \|\bar{\mathbf{g}}_t\|_2^2 + c\eta_t \|\bar{\mathbf{g}}_t\|_2^2 = \mathcal{L}(\boldsymbol{\theta}_t) - (1-c)\eta_t \|\bar{\mathbf{g}}_t\|_2^2 \quad (56d)$$

$$\Rightarrow (1-c)\eta_t \|\bar{\mathbf{g}}_t\|_2^2 \leq \mathcal{L}(\boldsymbol{\theta}_t) - \mathcal{L}(\boldsymbol{\theta}_{t+1}) \quad (56e)$$

$$\Rightarrow \sum_{t=T}^{T+N} \eta_t (1-c) \mathbb{E} [\|\bar{\mathbf{g}}_t\|_2^2] \leq \mathbb{E} [\mathcal{L}(\boldsymbol{\theta}_T) - \mathcal{L}(\boldsymbol{\theta}_{t+N+1})] < \infty \quad (56f)$$

$$\lim_{N \rightarrow \infty} \sum_{t=T}^{T+N} \eta_t (1-c) \mathbb{E} [\|\bar{\mathbf{g}}_t\|_2^2] < \infty. \quad (56g)$$

Together with the condition $\sum_t \eta_t = \infty$ and $r_{t+1} = \mathcal{O}(\eta_t)$, we can validate the convergence of the optimization process, i.e., $\lim_{t \rightarrow \infty} \mathbb{E} [\|\bar{\mathbf{g}}_t\|_2] = 0$. This further indicates the $\lim_{t \rightarrow \infty} \mathbb{E} [\|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)\|_2] \leq \lim_{t \rightarrow \infty} \mathbb{E} [\|\bar{\mathbf{g}}_t\|_2] + \mathbb{E} [\Delta \mathbf{g}_t] \leq \lim_{t \rightarrow \infty} \mathbb{E} [\|\bar{\mathbf{g}}_t\|_2] + r_{t+1}G = 0$. And this completes the proof of **Theorem 2**.

Diminishing Temporal Bias in the RPM. The above few sections have clarified the predictive foundation of MPTS and illustrated the theoretical convergence in stochastic optimization. Intuitively, the diminishing rank flipping can be viewed as the noise injected to the gradient of the MC-CVaR. In implementation, MPTS uses the function approximation capability of neural networks and introduces the RPM to utilize the optimization history $H_{1:t}$, generalizing from previously sampled and evaluated tasks to a new candidate pool $\mathcal{T}_{t+1}^{\mathcal{B}}$. That is the rank flipping can also be attributed to the error of the function approximation in the RPM. However, this can be well alleviated if seeking an optimal RPM. For better understanding, we analyze the temporal bias on arbitrary task τ through the adaptation risk approximation lens:

$$\text{Temporal Bias: } |\mathbb{E}_{p(\ell|\tau, H_{1:t})}[\ell] - \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}_{t+1})| \quad (57a)$$

$$\leq |\mathbb{E}_{p(\ell|\tau, H_{1:t})}[\ell] - \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}_t)| + |\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}_t) - \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}_{t+1})| \quad (57b)$$

$$\leq \epsilon_t + \beta_1 \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|_2 = \epsilon_t^{\text{approx error}} + \beta_1 \eta_t \|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)\|_2. \quad (57c)$$

Here $\epsilon_t^{\text{approx error}}$ denotes the RPM's approximation error due to the function approximation of the RPM, while β_1 is the Lipschitz from the Assumption 1. Note that if the RPM is expressive enough, the generalization error from function approximation can be theoretically well decreased with the accumulation of the optimization history based on the statistical learning theory. Hence, the limit of the right side term shrinks to extremely tiny values over iteration.

As specific analysis from multiple approximation bias depends on the design of the RPM. This work provides the VAE-like model in implementation, and we believe there exist more optimal candidates for the choice of RPMs. And this can be a promising research direction in the future.

To conclude, MPTS is computationally efficient and requires minimal annotation, such as agent-environment interactions; however, it also introduces approximation bias from the RPM and irreducible bias due to rank flipping when the model parameters change from $\boldsymbol{\theta}_t$ to $\boldsymbol{\theta}_{t+1}$. (Please see this point in Appendix Fig. 9)

F Backbone Methods & Experimental Details in Any-Shot Learning

F.1 MAML, Reptile and PEARL

In sinusoid regression and Meta-RL, MAML is used as the backbone algorithm. As previously discussed, MAML is widely applied in solving few-shot learning tasks. In mathematics, its optimization objective can be characterized as:

$$\min_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{p(\tau)} [\ell(\mathcal{D}_{\tau}^Q; \boldsymbol{\theta} - \lambda \nabla_{\boldsymbol{\theta}} \ell(\mathcal{D}_{\tau}^S; \boldsymbol{\theta}))], \quad (58)$$

where the term inside the bracket specifies the adaptation risk $\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta})$, and $\boldsymbol{\theta} - \lambda \nabla_{\boldsymbol{\theta}} \ell(\mathcal{D}_{\tau}^S; \boldsymbol{\theta})$ denotes the gradient update with the learning rate λ as fast adaptation to the task τ . After meta-training, we can access the meta initialization $\boldsymbol{\theta}$ that generalizes across the task space.

When it comes to reinforcement learning scenarios, \mathcal{D}_{τ} corresponds to episodic returns collected from MDPs with either the meta policy or the fast adapted policy. To ensure enough coverage of task space, we adopt a mixture strategy of MPTS and random sampling as an empirical regularizer in all RL scenarios, which is similar to work¹⁷.

Similar to MAML, Reptile⁵¹ is a simple yet effective meta-learning algorithm to learn an initialization of model parameters for rapid adaptation to new tasks with just a few gradient updates. Unlike MAML, which explicitly differentiates through the adaptation process, Reptile approximates the meta-objective by repeatedly sampling tasks, performing a few steps of stochastic gradient descent on each, and then moving the initialization toward the adapted weights.

PEARL⁵⁵ is an off-policy SoTA Meta-RL algorithm that enables efficient adaptation to new tasks by learning a latent probabilistic context variable representing task-specific information. Instead of relying on on-policy trajectories, PEARL uses off-policy data and amortized inference to infer this latent variable $q_\phi(z|\mathcal{C})$, where \mathcal{C} is a small context set of experience from the target task. The policy $\pi_\theta(a|s, z)$ and value functions are conditioned on z , allowing fast adaptation through posterior inference rather than gradient updates. This design benefits from sample efficiency of off-policy learning and rapid task adaptation.

F.2 Domain Randomization

Robotic DR refers to the setup that trains the agent in a collection of environments to obtain a generalizable policy. The diversity of environments tends to increase the robustness of policies in deployment. Such a setup does not require few-shot episodes in unseen but similar environments. In mathematics, we can express the optimization objective as:

$$\max_{\theta \in \Theta} \mathcal{J}(\theta) := \mathbb{E}_{\pi_\theta} \mathbb{E}_{p(\tau)} \left[\sum_{t=0}^H \gamma^t r_t \right] \quad (59)$$

where $p(\tau)$ defines the distribution over MDPs, and $\{r_t\}_{t=0}^H$ is the episodic stepwise reward after interacting with a specific MDP with H as the horizon. Once finishing the optimization of Eq. (59), we can access the policy π_θ as the zero-shot decision-maker in new environments. In this case, the adaptation risk can be in the form $\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) = -\sum_{t=0}^H \gamma^t r_t$.

Remember that MDP distribution $p(\tau)$ is mostly induced by physical parameters, e.g., mass, gravity, friction, etc., or the reward functions. In each training iteration, the machine learner resamples a batch of MDPs and gets the shared policy to interact with them to collect episodes. Consequently, the query dataset contains the episodes collected with no support dataset. Overall, policy optimization follows the standard TD3 algorithm⁴¹ due to its sample efficiency and stability.

F.3 Multi-Modal Prompt Learning

Multi-modal prompt learning is based on the backbone of the prompt tuning method MaPLe⁴⁰, which we use on both few-shot and SFT for image classifications.

Few-shot classification. The few-shot prompt learning refers to the common few-shot classification setting^{32,93}. We integrate the MaPLe backbone with the ProtoNet⁹³ to fully utilize the support sets in few-shot learning. As illustrated in Section G.3, we generate the model prediction using both the textual classifiers from the CLIP textual encoder and the visual classifiers from the support set. By freezing the CLIP model parameters, only the prompts are optimized during meta-training. In mathematics, the optimization objective can be formulated as:

$$\max_{\mathbf{u}} \mathbb{E}_{p(\tau)} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_\tau^Q} [\log p_\theta(y|\mathbf{x}, \mathbf{u})], \quad (60)$$

where \mathbf{u} denotes the learnable multi-modal prompts. $\theta = (\theta_t, \theta_i)$ contains the parameters of the CLIP textual and visual encoders, which are frozen during training. The prediction of each query image \mathbf{x} from task τ is calculated by Eq. (64). Loglikelihood maximization is implemented by minimizing the classification cross-entropy loss.

SFT for image classification. The prompt learning setting refers to the 16-shot classification task proposed in work^{105,106}. Based on the MaPLe backbone⁴⁰, we again freeze the CLIP model parameters and tune multi-modal prompts. The prompts are optimized on a selected ImageNet subset, with 16 samples from each category.

Note that in the SFT setting, we do not have pre-defined N-way K-shot tasks, either the splits of support and query sets in each task. Therefore, we replace the “tasks” in the meta-learning setting with training samples. Model predictive task sampling is then achieved through data sampling. In mathematics, the objective can be formulated as:

$$\max_{\mathbf{u}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\log p_\theta(y|\mathbf{x}, \mathbf{u})], \quad (61)$$

where \mathbf{u} and θ denote the prompts and frozen CLIP parameters as in the few-shot prompt learning setting, respectively. \mathbf{x} are training samples from the entire training set \mathcal{D} . The prediction of each image $p_\theta(y|\mathbf{x}, \mathbf{u})$ is calculated by $\frac{\exp(-d(f_{\theta_i}(\mathbf{x}, \mathbf{u}), \mathbf{t}_k))}{\sum_{k'} \exp(-d(f_{\theta_i}(\mathbf{x}, \mathbf{u}), \mathbf{t}_{k'}))}$ with the textual classifiers \mathbf{t} , obtained similarly to Eq. (62).

G Experimental Setups & Implementation Details

Practical Learning Efficiency and Robustness. Widely recognized in reinforcement learning is the high sample complexity in policy evaluation, which demands massive interactions with environments, while policy optimization over the MDP distribution makes this even more severe. In N -way K -shot image classification, we can create K -shot classification task from an arbitrary combination of N classes; then the task space complexity $\mathcal{O}(C_M^N)$ grows with the number of categories M in image datasets. Meanwhile, challenges arise when gradient updates of foundation models consume substantial computational power and memory with a large batch size. Similar circumstance also occurs in robust finetuning foundation models.

Table 5: **A Summary of the Considered Benchmarks.** Here, we list the primary expensive part in task episodic learning for each scenario together with backbone methods. Also note that N -way K -shot image classification and SFT requires implicit task identifiers while others can directly access explicit task identifiers as the lower dimensional task representation.

Benchmarks	Adaptation	Backbone	Expensive Part
K-shot sinusoid regression	few-shot	MAML/Reptile	computations
N -way K -shot image classification	few-shot	MaPLe	computation/memory
Meta-RL	few-shot	MAML/PEARL	interactions
Robotic DR	zero-shot	TD3	interactions
SFT	many-shot	MaPLe	computation/memory

Neural architecture of the RPM. As mentioned in the main paper, the RPM is in an encoder-decoder structure. For generality sake, we keep the neural architecture same for all benchmarks, including regression, classification and reinforcement learning. The encoder includes an embedding network with 4 hidden layers of size 128 (for Image Classification and Prompt-Tuning) or 10 (for Sinusoid Regression, Meta-RL, and Robotic DR) with the Rectified Linear Unit (ReLU) nonlinear activation units to encode $\{[\tau_{t,i}, \ell_{t,i}]\}$ batch for mean pooling and then maps to $[\mu, \Sigma]$ with an output layer. The decoder is a network with 3 hidden layers with nonlinear activation units to map $[z, \tau]$ to $\ell \in \mathbb{R}$. For further details, please refer to our code.

Visualized Results during Training Phases. Note that the active selection and the random sampler with different batches affect the reflection of the machine learner’s performance. Hence, learning curves in sinusoid regression (Fig. 3.a), Meta-RL (Fig. 4.a-b), and DR (Fig. 5.a-b) are actually evaluated in a uniformly sampled validation task dataset for fair comparison. Details on these validation task dataset are attached in the opensourced code.

Considerations in Experimental Design. To comprehensively evaluate MPTS, we consider benchmarks that span typicality and various levels of complexity and computational demands. For trait illustration, we adopt sinusoid regression, a classical and widely used benchmark in meta-learning research, as it provides an accessible and interpretable testbed for verifying the core traits of MPTS. To examine computational and memory efficiency in pattern recognition with foundation models, we employ few-shot learning and multimodal prompt-tuning tasks in image classification, which are standard benchmarks in the community. These tasks illuminate the advantage of MPTS in alleviating the intensive computational and memory costs typically associated with large backbone models under CVaR-based strategies. Finally, to investigate annotation and computational efficiency in sequential decision-making, we include Meta-RL and DR benchmarks. As standard environments for adaptive agent training, they are notoriously expensive due to extensive agent–environment interactions and the heavy task evaluations required by MC-CVaR strategies. In contrast, MPTS is capable of reducing such costs by partially surrogating agent–environment interactions and directly identifying a subset of MDPs from which to collect trajectories.

G.1 Sinusoid Regression

Task setup. For sinusoid regression, we retain the setup in MAML³², where the few-shot machine learner tries to complete a wave function with the support dataset. In specific, sampling the amplitude a and the phase b configures the wave function, and 10 data points are uniformly sampled from the interval $[-5.0, 5.0]$ coupled with $y = a \sin(x - b)$ to obtain the support dataset. This formulates the 10 – shot sinusoid regression task.

Meta training process and neural architectures for MAML. The machine learner is a neural network with 2 hidden layers of size 40 with two nonlinear activation units. The task batch for ERM, GDRM, OHTM, DATS and TDPS is 16, and that for DRM is 32 as default. DRM selects the 16 hardest tasks for optimization based on their real evaluation losses. OHTM selects the 8 hardest tasks from the memory and randomly samples the remaining ones from the current batch, following Kumar et al.⁴⁹. The temperature parameter is set to $\eta = 0.001$ in GDRM, $\eta = 0.02$ in DATS, and $\eta = 0.005$ in PDTS. The learning

rates for the inner loop and the outer loop are 0.001. For MPTS, the batch size of the candidate task identifier in training is 32, the Lagrange multiplier is set as 1, and the balance parameters $\{\gamma_0, \gamma_1\}$ is set to $\{1, 3\}$. We use the Adam optimizer with the learning rate $3e-4$ to update the RPM for 20000 step. In sinusoid regression and Meta-RL, we use the standard repository in MAML³².

Meta training process and neural architectures for Reptile. Reptile is implemented following Nichol et al.⁵¹, with 8 inner iterations, an inner learning rate of 0.01, and an outer learning rate of 0.5. The neural network architecture for the machine learner and the hyperparameters for all baselines are kept consistent with MAML, except that the temperature parameter in DATS is set to $\eta = 0.2$.

G.2 N-way K-shot Image Classification

Task Setup. This is a commonly seen benchmark in few-shot learning. It learns a model that can classify images from N distinct classes with support of K labeled examples for each class. The support dataset as reference is in the form $\mathcal{D}_\tau^S = \{[\mathbf{x}_{i,k}, y_{i,k} = i]\}_{k=1}^K\}_{i=1}^N$. And the query dataset corresponds to the image information for the model to classify. Hence, for a large image dataset with M classes, the complexity of the task space is $\mathcal{O}(C_M^N)$. Here, we include ImageNet-CG⁵², ImageNet-CI⁵², ImageNet-CS⁵², ImageNet-A⁵³, ImageNet-S²⁹ and ImageNet-R⁵⁴ as the dataset in evaluation.

Meta training process and neural architectures. Explicit τ are unavailable to specify the task; however, it can be approximately resolved by describing the identifier through a small reference model. Specifically, we leverage CLIP’s text encoder to obtain $\tau \approx [\text{CLIP}_{\text{text}}(\mathcal{C}_1), \dots, \text{CLIP}_{\text{text}}(\mathcal{C}_K)]$ with the tokenizations of K class texts $\mathcal{C}_{1:K}$ from \mathcal{D}_τ^S . The machine learner utilizes a prompt learning backbone following MaPLe, with the frozen CLIP model. The task batch for ERM and GDRM is 4, and that for DRM is 8 as default. The temperature parameter in GDRM is 0.001. The learning rate for the outer loop is 0.01. The learning rate for the inner loop follows that in MaPLe⁴⁰. The following is about extra optimization details or setups in MPTS. The task identifier is generated by the frozen CLIP text encoder using the input class names, with a dimensionality of 512. The batch size of the identifier in training is 8, the Lagrange multiplier is set as β , and we use the Adam optimizer with the learning rate 0.01 to update the RPM.

G.3 Prompt-based Few-shot Image Classification

We adopt the standard few-shot image classification setting^{32,93}, where tasks are constructed using the N-way K-shot paradigm for both meta-training and meta-testing. Each task comprises support and query sets. The support set contains K examples for each of the N classes, while the query set includes 15 examples per class. During meta-training, labels for both support and query data are accessible to the adaptive machine learner. During meta-testing, the query dataset’s labels are to be predicted given the labeled support dataset. The class categories of task datasets in the meta-training and meta-testing do not overlap. In experiments, we specifically consider a 5-way 1-shot image classification configuration. During meta-training, we set the task batch for ERM and DRM as $\mathcal{B} = 4$ (For implementation simplicity, the data loader samples 8 tasks and then randomly keeps half without ranking to optimize). The task batch for DRM is $\mathcal{B} = 8$ before the filtering operation; DRM filters half to optimization. Similarly, that for MPTS is $\mathcal{B} = 8$ and only 4 tasks are screened to optimize. For DATS and TDPS, we set the task batch as $\mathcal{B} = 4$ and compute the gradients of each task. The optimization of the 4 tasks are then weighted by gradient-derived difficulty scores (TDPS) or similarity scores (DATS). In OHTM, we introduce a hard task buffer as size ten to store the most difficult tasks during meta-training. We set $\mathcal{B} = 8$ and randomly sample 4 tasks from the buffer in each iteration.

To enable few-shot learning by prompt-tuning, we integrate the multimodal prompt learning methods MaPLe⁴⁰ and prototypical network (ProtoNet)⁹³. MaPLe operates on the CLIP model¹, capturing multimodal prompts to refine both visual and textual feature representations with frozen CLIP’s parameters. These refined textual features serve as classifiers for predicting refined image features. In parallel, ProtoNet is utilized to derive class-specific visual embeddings from the support set, which assist in distinguishing query samples.

To utilize both MaPLe and ProtoNet, we construct classifiers based on both textual features and visual embeddings. Predictions for query samples are generated by combining the classifiers through a weighted sum. This combination strategy is employed during meta-training to optimize the multimodal prompts. Meanwhile, the CLIP model’s parameters remain frozen throughout optimization. During meta-testing, these trained prompts are adopted to create textual and visual classifiers and process query image features. Final predictions for each query sample are made using the same weighted combination approach as in meta-training.

In mathematics, we can characterize the mentioned pipeline as:

$$\text{Textual Classifier from the Textual Features: } t_k = f_{\theta_t}(\mathbf{l}_k, \mathbf{u}), \quad (62)$$

$$\text{Prototypical Classifier from the Support Image: } \mathbf{c}_k = \frac{1}{|S_k|} \sum_{(\mathbf{x}_i, y_i) \in S_k} f_{\theta_i}(\mathbf{x}_i, \mathbf{u}), \quad (63)$$

Classification Likelihood from the Query Dataset:

$$p_{\theta}(y = k | \mathbf{x}, \mathbf{u}) = \lambda_1 \frac{\exp(-d(f_{\theta_i}(\mathbf{x}, \mathbf{u}), \mathbf{c}_k))}{\sum_{k'} \exp(-d(f_{\theta_i}(\mathbf{x}, \mathbf{u}), \mathbf{c}_{k'}))} + \lambda_2 \frac{\exp(-d(f_{\theta_i}(\mathbf{x}, \mathbf{u}), \mathbf{t}_k))}{\sum_{k'} \exp(-d(f_{\theta_i}(\mathbf{x}, \mathbf{u}), \mathbf{t}_{k'}))}, \quad (64)$$

where θ_t and θ_i denote the textual and visual encoders of the CLIP model, respectively. \mathbf{l}_k and \mathbf{u} respectively denote the textual descriptions of category k and the multimodal prompts. $(\mathbf{x}_i, y_i) \in S_k$ denotes the images and labels of category k in the support dataset \mathcal{D}_{τ}^S . Once the textual classifier \mathbf{t} and support visual classifier \mathbf{c} are obtained, we predict the query sample \mathbf{x} by the classifiers with hyperparameters $\lambda_1 = 0.25$ and $\lambda_2 = 1.0$.

G.4 Meta-RL

Task Setup. We construct MDP distributions based on Mujoco physics engines¹⁰⁷. These include HalfCheetahVel, HalfCheetahMassVel, Walker2dVel, Walker2dMassVel, and ReacherPos. The HalfCheetahVel and Walker2dVel tasks involve training the cheetah or walker robot to achieve a target velocity. These tasks define the reward function as the negative absolute difference between the robot’s current velocity and the target velocity, supplemented by a control penalty and an alive bonus to facilitate the learning process. The goals and rewards of HalfCheetahMassVel and Walker2dMassVel are the same as those of the corresponding velocity-related tasks, with the additional identifier of varying mass for the cheetah or walker robot. The ReacherPos task tries to move a two-jointed robot arm’s end effector close to a target position, and its reward function is defined as the negative $L-1$ distance between the robot arm’s position and the target position, supplemented by a control cost to ensure robustness.

Meta training process and neural architectures for MAML. The machine learner is a neural network with 2 hidden layers of size 64 with the Rectified Linear Unit (ReLU) nonlinear activation units. The task batch for ERM, GDRM, OHTM, DATS and TDPS is 20, and that for DRM is 40 as default. DRM selects the 20 hardest tasks for optimization based on their real evaluation losses. OHTM selects the 10 hardest tasks from the memory and randomly samples the remaining ones from the current batch, following Kumar et al.⁴⁹. The temperature parameter in GDRM is 0.001. For DATS, it is set to 0.05 for HalfCheetahMassVel, 0.1 for ReacherPos, HalfCheetahVel, and Walker2dVel, and 0.2 for Walker2dMassVel. For TDPS, it is set to 0.05 for HalfCheetahMassVel, Walker2dMassVel, and HalfCheetahVel, 0.1 for Walker2dVel, and 0.2 for ReacherPos. The learning rates for the inner loop and the outer loop are 0.1. For MPTS, the batch size of the candidate task identifier in training is 30, the Lagrange multiplier is set as $\beta = 0.0001$, and the balance parameters $\{\gamma_0, \gamma_1\}$ is set to $\{1, 5\}$. We use the Adam optimizer with the learning rate 0.005 to update the RPM.

Meta training process and neural architectures for PEARL. PEARL⁵⁵ and RoML¹⁷ are implemented based on the official RoML codebase available at <https://github.com/ido90/RoML-pearl>. Two tasks, HalfCheetahBody and HalfCheetahMass, are adopted from RoML. For MPTS, the batch size of the candidate task identifier during training is set to 5× that of the real training batch size (30) for HalfCheetahBody, and 2× for HalfCheetahMass. All other hyperparameters follow those used in MAML.

G.5 Robotic DR

Task Setup. We conduct experiments on LunarLander-v2 and ErgoReacher-v0 environments³⁹. LunarLander is a 2 degrees of freedom (DoF) environment in which the agent has to land a spacecraft on a designated landing pad without crashing, implemented using Box2D¹⁰⁸. The reward function of LunarLander awards positive rewards for successful landings, negative rewards for crashes, and additional penalties for fuel consumption and deviation from the landing pad, encouraging efficient and controlled landings. ErgoReacher is a 4 DoF arm environment from¹⁰⁹ in which the arm has to touch a goal with its end effector, implemented in the Bullet Physics Engine¹¹⁰. The reward function of ErgoReacher includes the negative distance between the end effector’s position and the target, along with other control costs to promote efficient and safe movements. In LunarLander, we randomize the engine strength, while in ErgoReacher, we randomize the joint damping and maximum torque for each of the 4 joints, resulting in a total of 8 parameters. The detailed ranges of the randomized parameters for each environment are provided in Table 4.

DR training process and neural architectures. Both the actor and critic are neural networks with 2 hidden layers of size 256 with the Rectified Linear Unit (ReLU) nonlinear activation units. The task batch for ERM, GDRM, and OHTM is 10, and that for DRM is 20 as default. DRM selects the 10 hardest tasks for optimization based on their real evaluation losses. OHTM selects the 5 hardest tasks from the memory and randomly samples the remaining ones from the current batch, following Kumar et al.⁴⁹. The temperature parameter in GDRM is 0.01. The learning rates for actor and critic are $3e-4$. The following is about

extra optimization details or setups in MPTS. The batch size of the candidate task identifier in training is 25 for LunarLander and 250 for ErgoReacher. The Lagrange multiplier is set as $\beta = 1.0$, and the balance parameters $\{\gamma_0, \gamma_1\}$ is set to $\{1, 5\}$. We use the Adam optimizer with the learning rate 0.005 to update the RPM.

G.6 Prompt-Tuning Multimodal Foundation Models

Task Setup of Prompt-tuning. We refer the reader to MaPLe’s implementation in <https://github.com/muzairkhattak/multimodal-prompt-learning>. For all baselines, we retain the MaPLe’s task construction in prompt-tuning.

Prompt-tuning process and neural architectures. The machine learner follows the prompt learning method MaPLe⁴⁰ based on the frozen CLIP model (ViT/B-16). The task batch for ERM and GDRM is 4, and that for DRM is 8 as default. The temperature parameter in GDRM is 0.001. The learning rate for the outer loop is 0.005. The learning rate for the inner loop follows that in MaPLe⁴⁰. The following is about extra optimization details or setups in MPTS. As for the neural architecture of the RPM, the encoder is a neural network with 5 hidden layers with 4 ReLU nonlinear activation units, and the decoder is a neural network with 4 hidden layers with 3 nonlinear activation units. The task identifier’s dimension is 512 with the latent embedding from CLIP encoders. The batch size of the identifier in training is 8, the Lagrange multiplier is set as β , and we use the Adam optimizer with the learning rate 0.005 to update the RPM for 8000 steps. During prompt-tuning, we set the task batch for ERM and DRM as $\mathcal{B} = 4$ (For implementation simplicity, the data loader samples 8 tasks and then randomly keeps half without ranking to optimize). The task batch for DRM is $\mathcal{B} = 8$ before the filtering operation; DRM filters half to optimization. Similarly, that for MPTS is $\mathcal{B} = 8$ and only 4 tasks are screened to optimize.

H Computational Tools & Platforms & Model Design

In this research project, we use the Pytorch as the package to implement all methods to run all deep learning experiments.

From Definition 2, we can see that the goal is to learn some low-dimensional signals online when the task identifier is low-dimensional and the loss is one-dimensional. According to the Occam’s Razor theorem, the best practice is to adopt the generative model that captures historical priors while exhibiting less model complexity. In our implementation, we typically use the same architecture of shallow networks from the Open-source code. However, we also believe it is a promising direction to develop more advanced RPMs to online tell task difficulty and achieve the purpose of Definition 2.

I Competing Interests & Author Contributions

The author list is Qi Cheems Wang (Q.C.W.), Zehao Xiao (Z.X.), Yixiu Mao (Y.M.), Yun Qu (Y.Q.), Jiayi Shen (J.S.), Yiqin Lv (Y.L.), and Xiangyang Ji (X.J.). The authors declare no competing interests in this work. X.J. launched and sponsored this research on the reliable and efficient adaptation learning project. Z.X. and J.S. from the University of Amsterdam attended this project, and this work was done during their remote visiting the Tsinghua University Intelligent Decision-Making Lab from April 2024 to August 2024. J.S. is now working at Facebook AI Research. The authors confirm their contributions to this work as follows:

Under the supervision of professor X.J., Q.C.W. conceptualized the idea of MPTS, designed the computational framework, formulated the mathematical part, and wrote the draft. Z.X. and J.S. implemented MPTS in sinusoid regression, prompt-based few-shot image classification, and SFT image classification, collected experimental results to visualize, and added implementation details in Supplementary Material. Y.M., Y.Q., and Y.L. implemented MPTS in sinusoid regression, Meta-RL and robotic DR, collected experimental results to visualize, and added implementation details in Supplementary Material. X.J. supervised the progress of MPTS, organized technical discussions, reviewed and revised the original draft. All authors have read the manuscript and approved the public version.

First Author Biography: Qi Wang received Ph.D. degree under supervision of Professor Max Welling and Associate Professor Herke van Hoof in 2022. He is now under supervision of Prof. Xiangyang Ji and works as a research assistant at Tsinghua University. His research focus is on generative modeling and intelligent decision-making. He has published several papers on top-tier conferences such as ICML/NeurIPS/ICLR and was awarded 2023 China Multi-Agent System Outstanding Doctoral Thesis Award.

Correspondence Author Biography: Xiangyang Ji received the B.S. degree in materials science and the M.S. degree in computer science from the Harbin Institute of Technology, Harbin, China, in 1999 and 2001, respectively, and the Ph.D. degree in computer science from the Institute of Computing Technology, Chinese Academy of Sciences, in 2008. He joined Tsinghua University, Beijing, in 2008, where he is currently a Professor with the Department of Automation, School of Information

Science and Technology. He has authored more than 100 refereed conference and journal papers. His current research interests include signal processing, computer vision, computational photography, and intelligent decision-making.