
Fast and Robust: Task Sampling with Posterior and Diversity Synergies for Adaptive Decision-Makers in Randomized Environments

Anonymous Authors¹

Abstract

Task robust adaptation is a long-standing pursuit in sequential decision-making. Some risk-averse strategies, e.g., the conditional value-at-risk principle, are incorporated in domain randomization or meta reinforcement learning to prioritize difficult tasks in optimization, which demand costly intensive evaluations. The efficiency issue prompts the development of robust active task sampling to train adaptive policies, where risk-predictive models can surrogate policy evaluation. This work characterizes robust active task sampling as a secret Markov decision process, posits theoretical and practical insights, and constitutes robustness concepts in risk-averse scenarios. Importantly, we propose an easy-to-implement method, referred to as Posterior and Diversity Synergized Task Sampling (PDTs), to accommodate fast and robust sequential decision-making. Extensive experiments show that PDTs unlocks the potential of robust active task sampling, significantly improves the zero-shot and few-shot adaptation robustness in challenging tasks, and even accelerates the learning process under certain scenarios.

1. Introduction

Deep reinforcement learning (RL) has garnered remarkable progress in solving complicated sequential decision-making problems in the past few years (Sutton & Barto, 2018). However, an existing challenge is effectively transferring the RL policy to unseen but similar scenarios without learning from scratch. A commonly used strategy is to randomize the environment, e.g., placing a distribution over Markov decision processes (MDPs), for policy search in a zero-shot or few-shot manner. This facilitates the rise of domain ran-

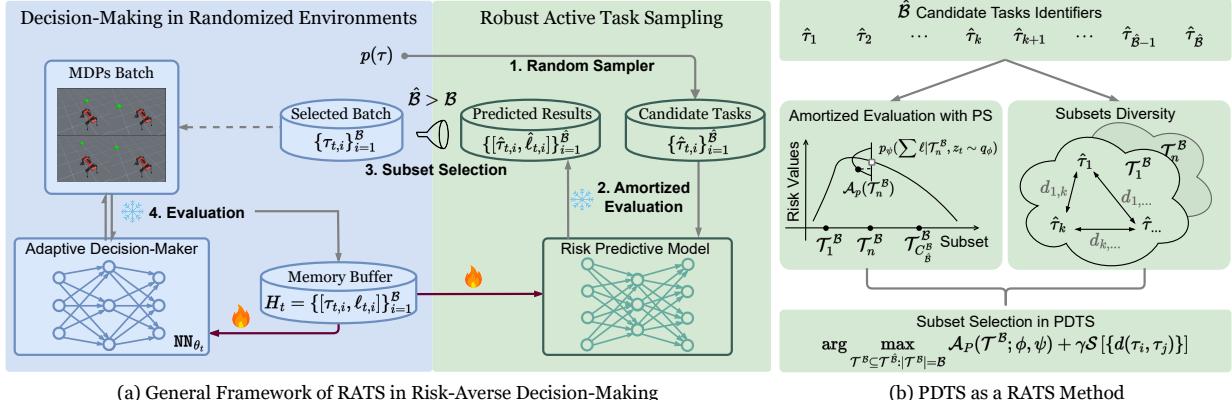
domization (DR) (Muratore et al., 2018) and Meta-RL (Finn et al., 2017) paradigms, which train adaptive policies in task episodic learning. Simultaneously, adaptation robustness to worst-case scenarios is catching increasing attention as most real-world decision-making scenarios are inherently risk-sensitive, where failures in adaptation can cause catastrophic outcomes, e.g., damage to robots (Carpin et al., 2016) or accidents in autonomous driving (Rempe et al., 2022).

Active Inference’s Promise for Adaptive Robust Decision-Maker: When risk-averse principles are incorporated in DR and Meta-RL to enhance adaptation robustness (Wang et al., 2024; Lv et al., 2024; Greenberg et al., 2024); prioritizing challenging tasks to optimize demands intensive and expensive policy evaluation in massive environments over iteration. To overcome the efficiency bottleneck, Wang et al. (2025) constructs a risk predictive model to actively infer MDP difficulties for worst subset selection in policy search, which we identify as a method of robust active task sampling (RATS) paradigm in Fig. 1. As policy evaluation in arbitrary MDPs can be amortized by executing this risk predictive model, the resulting model predictive task sampling (MPTS) (Wang et al., 2025) enables expanding the scope of task subset selection, e.g., screening \mathcal{B} from \mathcal{B} MDPs, to learn adaptive policies without extra environment interaction cost. This reflects RTAS’s huge potential for efficient, robust decision-making when exhaustive policy evaluation is prohibitive in the vast MDP space.

Challenges in Theoretical Analysis and Implementations: Despite RATS’s promise in decision-making, we can still perceive several issues from its latest SOTA method MPTS. (i) No versatile tool is developed for theoretical analysis, e.g., the robustness concept in optimization, which is an indispensable consideration in risk-averse cases. (ii) It requires a dedicated pseudo batch size \hat{B} and other configurations. As implied in Meta-RL results (Wang et al., 2025), appropriately mixing up the random and predictive samplers is decisive in task subset selection; otherwise, it degrades generalization and robustness (See Sec. 3.2). (iii) There lack comprehensive discussions about acquisition principles in RATS. The adopted upper confidence bound (UCB) principle must strike a balance between worst-case and uncertainty, with hyper-parameters hard to adjust in practice.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.



(a) General Framework of RATS in Risk-Averse Decision-Making

(b) PDTS as a RATS Method

Figure 1. (a) General RATS in risk-averse decision-making. The pipeline involves amortized evaluation of task difficulties, robust subset selection, policy optimization in the MDP batch, and risk predictive models’ update. [fire: updates; snow: evaluation] (b) PDTS as a RATS method. PDTS treats task subsets as bandit arms, evaluates values through posterior sampling, and solves a regularized problem.

Making Sense of RATS in Decision-Making: Regarding theoretical understanding, we first abstract the general task episodic learning process as a secret MDP \mathcal{M} , construct an infinitely many-armed bandit (i-MAB) (Carpentier & Valko, 2015) and demonstrate MPTS (Wang et al., 2025) as a special solution to the i-MAB. Aim at exploring more optimal subsets with the risk predictive model; we make a diagnosis of the concentration issue under a larger $\hat{\mathcal{B}}$ and enhance the acquisition function with the diversity regularization (Wang et al., 2023; Borodin et al., 2017). To simplify the amortized evaluation and exploit the stochastic optimism in i-MAB, we adopt the posterior sampling strategy (Russo & Van Roy, 2014) to search for the optimal task subset with fewer configurations. Under these modifications, we present the Posterior and Diversity Synergized Task Sampling (PDTs) as a competitive RATS method.

Contributions and Fascinating Discoveries in PDTs. This work is built on empirical findings and the risk predictive model in MPTS, while the research focus is orthogonal (See Table 1). Surrounding risk-averse decision-making, the primary contributions are:

1. Our constructed i-MAB provides a versatile model to achieve RATS under various principles, including but not limited to MPTS and PDTs. The separate robustness concepts can be refined accordingly.
2. The designed diversity regularized acquisition function fixes the concentration issue, allows for exploration in a wider range of task sets (e.g., $\hat{\mathcal{B}} = 64\mathcal{B}$), and secures nearly worst-case MDP robustness.
3. The resulting PDTs is easy-to-implement and benefits from stochastic optimism in posterior sampling for decision-making.

Empirically, the most thrilling finding is that PDTs exhibits adaptation robustness superior to existing SOTA baselines in typical DR and Meta-RL benchmarks without complicated configurations. Even in more realistic and challenging sce-

narios, such as vision-based decision-making, PDTs retains a remarkable performance over others. Feel free to access our code at anonymous.4open.science/r/PDTs/.

2. Research Background

Notations. For conciseness and coherence, we retain most notations in (Wang et al., 2025) and leave the MDP distribution $p(\tau)$ details and RL preliminaries in Appendix A. Both DR and Meta-RL perform policy search in $p(\tau)$, where MDPs as tasks are specified by physics identifiers $\tau \in \mathbb{R}^d$. The primary goal of DR (Tobin et al., 2017; Muratore et al., 2018) and Meta-RL (Beck et al., 2023) is to seek a policy $\theta \in \Theta$ that adapts well to a new MDP with zero or a few rollouts. We denote the task-specific dataset by $\mathcal{D}_\tau = \mathcal{D}_\tau^S \cup \mathcal{D}_\tau^Q$, where \mathcal{D}_τ^S are the K -rollouts for fast adaptation in MDP τ and \mathcal{D}_τ^Q are rollouts for after-adaptation policy evaluation. DR differs from Meta-RL and learns to adapt in a zero-shot manner, indicating $\mathcal{D}_\tau^S = \emptyset$. The risk function is $\ell : \mathcal{D}_\tau \times \Theta \mapsto \mathbb{R}$, mapping to the adaptation risk value, e.g., negative average returns of query rollouts.

In task episodic learning for DR and Meta-RL, the optimization history is written as $\hat{H}_t = \{\theta_t, (\tau_{t,i}, \mathcal{D}_{\tau_{t,i}}, \ell_{t,i})\}_{i=1}^{\mathcal{B}}$, with \mathcal{B} the number of tasks to optimize in t -th iteration. MPTS (Wang et al., 2025) further prepares it as $H_t = \{(\tau_{t,i}, \mathcal{D}_{\tau_{t,i}}, \ell_{t,i})\}_{i=1}^{\mathcal{B}}$ to feed into the risk learner to predict adaptation risk on arbitrary task. Importantly, it introduces the pseudo task set $\mathcal{T}_t^{\hat{\mathcal{B}}}$ with $\hat{\mathcal{B}} > \mathcal{B}$ and employ the acquisition function $\mathcal{A}(\cdot)$ to actively select the subset for next iteration. Throughout this work, we leave the exact optimization task batch size \mathcal{B} fixed and same for all methods in both theoretical analysis and evaluation.

2.1. Task Robust Optimization Methods

This part recaps risk minimization principles, which are incorporated in DR and Meta-RL for robust decision-making.

110 **Expected/Empirical Risk Minimization (ERM).** ERM
 111 originates from the statistical learning theory (Vapnik et al.,
 112 1998). Such a principle minimizes the expectation of risk
 113 values under $p(\tau)$:

$$115 \min_{\theta \in \Theta} \mathbb{E}_{p(\tau)} [\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)], \quad (1)$$

117 where $p(\tau)$ is mostly a uniform distribution as default.

118 **Group Distributionally Robust Risk Minimization**
 119 (**GDRM**). Such a principle aims to boost the adaptation ro-
 120 bustness under certain proportional scenarios, e.g., a group
 121 of some under-sampled yet challenging tasks (Sagawa et al.,
 122 2019). In mathematics, we can write the expression as:

$$125 \min_{\theta \in \Theta} \max_{g \in \mathcal{G}} \mathbb{E}_{p_g(\tau)} [\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)], \quad (2)$$

126 with \mathcal{G} to denote a collection of groups over a task dataset.
 127 GDRM handles extremely worst subpopulation shifts (Koh
 128 et al., 2021) through a risk-reweighting mechanism.

129 **Distributionally Robust Risk Minimization (DRM).** As a
 130 typical strategy in DRM, CVaR $_\alpha$ selects $(1-\alpha)$ proportional
 131 worst tasks after exact evaluation to optimize:

$$135 \min_{\theta \in \Theta} \text{CVaR}_\alpha(\theta) := \mathbb{E}_{p_\alpha(\tau; \theta)} [\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)]. \quad (3)$$

136 In Meta-RL, this corresponds to the hard MDP priori-
 137 zation (Greenberg et al., 2024; Lv et al., 2024). With
 138 $\min_{\theta \in \Theta} \lim_{\alpha \rightarrow 1} \text{CVaR}_\alpha(\theta)$, it degenerates to the worst-
 139 case optimization in (Collins et al., 2020).

2.2. RATS Preliminaries & MPTS Modules

140 Here, we specify RATS paradigm, which slightly differs
 141 from traditional active learning purposes (Cohn et al., 1996).
 142 RATS is mainly incorporated into risk-averse learning with
 143 traits: (i) active inference towards task difficulties with lim-
 144 ited cost, and (ii) acquisition rules to select the subset from
 145 the pseudo task set to optimize for robustness.

146 For example, MPTS designs a risk predictive model
 147 $p(\ell|\tau, H_{1:t}; \theta_t)$ to surrogate expensive evaluation, e.g., neg-
 148 ative average rollout returns ℓ of the policy θ_t in a MDP τ .
 149 Hence, we identify it as a method of RATS. In particular,
 150 the empirical evidence in (Wang et al., 2025) Fig. 5 vali-
 151 dates its risk predictive model's feasibility of approximately
 152 scoring MDPs' difficulties with high Pearson correlation
 153 coefficients between the model predictive ones and exact
 154 evaluation. Next, we overview its construction together with
 155 the acquisition function.

156 **Generative Modeling Adaptation Optimization.** MPTS
 157 treats the optimization history as sequence generation and

158 involves latent variables z_t to summarize batches of adapta-
 159 tion risk over iterations, which leads to:

$$160 p(\mathcal{L}_{0:T}^B, z_{0:T} | \theta_{0:T}) = p(z_0) \prod_{t=0}^T p(\mathcal{L}_t^B | z_t, \theta_t) \prod_{t=0}^{T-1} p(z_{t+1} | z_t), \quad (4)$$

161 with the evaluation risk batch $\mathcal{L}_t^B = \{(\tau_{t,i}, \ell_{t,i})\}_{i=1}^B$.

162 With the Bayes rule and the streaming variational inference
 163 (Broderick et al., 2013; Nguyen et al., 2017) w.r.t. Eq. (4),
 164 it obtains the approximate evidence lower bound of the risk
 165 learner to maximize in each batch:

$$166 \max_{\psi \in \Psi, \phi \in \Phi} \mathcal{G}_{\text{ELBO}}(\psi, \phi) := \mathbb{E}_{q_\phi(z_t | H_t)} \left[\sum_{i=1}^B \ln p_\psi(\ell_{t,i} | \tau_{t,i}, z_t) \right] \\ - \beta D_{KL} [q_\phi(z_t | H_t) \| q_{\bar{\phi}}(z_t | H_{t-1})], \quad (5)$$

167 with $\bar{\phi}$ the fixed conditioned prior from the last update
 168 and $\beta \in \mathbb{R}^+$ the penalty weight. Then the amortized
 169 evaluation of adaptation performance is approximated
 170 as $p(\ell | \tau, H_{1:t}; \theta_t) \approx \mathbb{E}_{q_\phi(z_t | H_t)} [p_\psi(\ell | \tau, z_t; \theta_t)]$ through
 171 Monte Carlo estimates.

$$\max_{\psi \in \Psi} \mathcal{L}_{\text{ML}}(\psi) := \ln p_\psi(H_t | H_{1:t-1}) \quad (6a)$$

$$p(\ell | \hat{\tau}_i, H_{1:t}; \theta_t) \xrightarrow{\text{MC}} \{m(\ell_i), \sigma(\ell_i)\}_{i=1}^{\hat{B}} \quad (6b)$$

$$\mathcal{T}_{t+1}^{B*} = \arg \max_{\mathcal{T}_{t+1}^B \subseteq \mathcal{T}_{t+1}^{\hat{B}} : |\mathcal{T}_{t+1}^B| = B} \mathcal{A}_{\text{U}}(\mathcal{T}_{t+1}^B) \quad (6c)$$

172 **Optimization Pipeline.** MPTS contains three critical steps
 173 in accordance with Eq. (6). ① In *approximate posterior*
 174 *inference* for Eq. (6)a, MPTS optimizes the risk predictive
 175 model through maximizing Eq. (5) with H_t ; ② In *amortized*
 176 *evaluation*, it samples \mathcal{T}_{t+1}^B from $p(\tau)$ and runs stochastic
 177 forward passes to estimate $m(\ell_i)$ and $\sigma(\ell_i) \forall \hat{\tau} \in \mathcal{T}_t^{\hat{B}}$ in
 178 Eq. (6)b; ③ In *subset selection*, it picks up the subset \mathcal{T}_{t+1}^B
 179 among Top-B acquisition scores from $\mathcal{T}_{t+1}^{\hat{B}}$ for next iteration
 180 optimization. By repeating ①-③ steps for T -rounds, MPTS
 181 derives the robust adaptive decision-maker.

182 **Acquisition Function in MPTS.** The subset selection rule
 183 is built on the principle of optimism in the face of uncer-
 184 tainty (OFU) (Auer, 2002b), and tasks with worse adap-
 185 tation performance and higher epistemic uncertainty are
 186 prioritized. This leads to the UCB principle, i.e., $\mathcal{A}_{\text{U}}(\mathcal{T}_{t+1}^B)$
 187 (Auer, 2002a; Garivier & Moulines, 2008):

$$188 \mathcal{A}_{\text{U}}(\mathcal{T}_{t+1}^B) = \sum_{i=1}^B \gamma_0 m(\ell_i) + \gamma_1 \sigma(\ell_i), \text{ with } \tau_i \in \mathcal{T}_{t+1}^B, \quad (7)$$

189 where the mean $m(\ell_i) = \mathbb{E}_{q_\phi(z_t | H_t)} [p_\psi(\ell | \tau_i, z_t)]$ and
 190 the standard deviation $\sigma(\ell_i) = \mathbb{V}_{q_\phi(z_t | H_t)}^{\frac{1}{2}} [p_\psi(\ell | \tau_i, z_t)]$
 191 of task-specific adaptation risk and are estimated through
 192 multiple stochastic forward passes, i.e., $z_t \sim q_\phi(z_t | H_t)$
 193 and $\ell_i \sim p_\psi(\ell_i | \tau_i, z_t)$. $\{\gamma_0, \gamma_1\}$ are trade-off parameters.

Table 1. Comparison between MPTS and PDTS in Contributions.

	Research Lens	Robustness Concept	Subset Selection	Acquisition Rule
MPTS (Wang et al., 2025)	Generative Modeling	Nearly CVaR $_{\alpha}$	Max-Sum (Top- \mathcal{B})	UCB
PDTs (Ours)	MDP & i-MABs	Nearly Worst-Case	Max-Sum Diversity	Posterior Sampling

3. Theoretical Investigations and Practical Enhancements

This section first studies RATS with a secret MDP, proposes i-MABs as a versatile tool for RATS methods and establishes its connections with MPTS. To promote RATS's use in risk-averse decision-making, we analyze acquisition rules' influence on robustness concepts and present PDTs as an easy-to-implement yet powerful scheme.

3.1. Enable Robust Active Task Sampling with i-MABs

When RATS meets decision-making, it involves scoring MDPs' difficulty from amortized evaluation, selecting a subset, and performing adaptive policy search in either a zero-shot or few-shot manner. The following introduces a theoretical tool to analyze these steps in RATS.

Task Robust Episodic Learning as a Secret MDP. The primary insight lies in a finite-horizon Markov decision process (Puterman, 2014), denoted by $\mathcal{M} = \langle \mathbf{S}, \mathbf{A}, \mathbf{P}, \mathbf{R} \rangle$. Here, we specify the essential components of \mathcal{M} as:

- **State Space.** \mathcal{M} treats the feasible machine learner's parameter, such as Meta-RL policies, as the reachable state, i.e., $\mathbf{S} = \{\boldsymbol{\theta} \in \Theta\}$;
- **Action Space.** The action space constitutes a collection of task subsets with cardinality constraints $\mathbf{A}_t = \{\mathcal{T}_t^{\mathcal{B}} \subseteq \mathcal{T}_t^{\hat{\mathcal{B}}} \text{ with } |\mathcal{T}_t^{\mathcal{B}}| = \mathcal{B}\}$ in RATS or CVaR $_{\alpha}$ methods;
- **Transition Dynamics.** We describe the dynamical system as $p(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}}) \in \mathbf{P}$ conditioned on $\boldsymbol{\theta}_t$ and the action $\mathcal{T}_{t+1}^{\mathcal{B}}$ with the transited state (after-adaptation) $\boldsymbol{\theta}_{t+1}$;
- **Reward Function.** The step-wise reward quantifies adaptation robustness improvement after state transitions. With CVaR $_{\alpha}$ as a risk-averse measure, we define it as $R(\boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}}) := \text{CVaR}_{\alpha}(\boldsymbol{\theta}_t) - \text{CVaR}_{\alpha}(\boldsymbol{\theta}_{t+1})$.

Given \mathcal{M} , the maximum iteration step T also corresponds to the total interaction rounds. The agent actively selects $\mathcal{T}_{t+1}^{\mathcal{B}}$ from the time-varying $\mathcal{T}_{t+1}^{\hat{\mathcal{B}}}$ and optimize the machine learner to increase adaptation robustness in T -rounds.

Accordingly, the sequential actions are the outcome of a series of deterministic functions as the policy set $\Pi_{0:T-1} = \{\pi_t\}_{t=0}^{T-1}$. Here, the policy maps $H_{0:t}$ to the next subset from $\mathcal{T}_{t+1}^{\hat{\mathcal{B}}}$ for optimization, i.e., $\pi_t : H_{0:t} \mapsto \mathcal{T}_{t+1}^{\mathcal{B}}$. These ingredients can depict \mathcal{M} in a probabilistic graph as Fig. 8.

Formally, we can express the agent's ultimate goal as maximizing the cumulative reward in a sequential manner,

$$\Pi_{0:T-1}^* = \arg \max_{\Pi_{0:T-1}} \sum_{t=0}^{T-1} R(\boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}}). \quad (8)$$

Note that solving Eq. (8) is equivalent to reaching the optimal state $\boldsymbol{\theta}_T^*$ with lowest CVaR $_{\alpha}(\boldsymbol{\theta})$ after repeating T -round optimization steps.

We also denote the policy subset of an arbitrary intermediate decision-making sequence by $\Pi_{k:j} := \{\pi_i\}_{i=k}^j$ with its optimal solution marked in $*$. The associated cumulated return is $\mathbf{R}(\Pi_{k:j}^*; \boldsymbol{\theta}_k) = \sum_{i=k}^j R(\boldsymbol{\theta}_i, \mathcal{T}_{i+1}^{\mathcal{B}*})$ conditioned on the starting state $\boldsymbol{\theta}_k$.

Remark 3.1 (Bellman Optimality). For the studied T -horizon \mathcal{M} , we write its Bellman optimality as:

$$\mathbf{R}(\Pi_{0:T-1}^*; \boldsymbol{\theta}_0) = \max_{\Pi_{0:T-1}} \sum_{i=0}^{T-1} R(\boldsymbol{\theta}_i, \mathcal{T}_{i+1}^{\mathcal{B}}) + \mathbf{R}(\Pi_{i:T-1}^*; \boldsymbol{\theta}_i), \quad (9)$$

revealing that the optimal solution to the sub-problem also reserves its global optimality. Hence, we can break the problem-solving into optimal subset selection in each round.

From i-MABs to MPTS's Robustness Concept. Finding plausible strategies to solve \mathcal{M} is non-trivial as policy search is considered in a discrete space with the varying action set \mathbf{A}_t . To this end, we further simplify \mathcal{M} into an infinite MAB (Mahajan & Teneketzis, 2008) to enable an online search of the optimal subset and interpret the optimization pipeline of MPTS as a special case.

Distinguished from the vanilla multi-armed bandit, \mathcal{M} involves the state $\boldsymbol{\theta}_t$ and \mathbf{A}_{t+1} induced by resampled $\mathcal{T}_{t+1}^{\hat{\mathcal{B}}}$. The arm corresponds to a feasible subset $\mathcal{T}_{t+1}^{\mathcal{B}} \in \mathbf{A}_{t+1}$. Hence, we can associate the reward distribution $p(R | \boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}})$ with the chosen action $\mathcal{T}_{t+1}^{\mathcal{B}}$. As a result, one optimistic strategy for i-MABs is to execute greedy search $\mathcal{T}_{t+1}^{\mathcal{B}*} = \arg \max_{\mathcal{T}_{t+1}^{\mathcal{B}} \in \mathbf{A}_{t+1}} \mathbb{E}[R | \boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}}]$ in each round.

The regret in the i-MAB measures the performance difference between the cumulated rewards of the optimal policy $\Pi_{0:T-1}^*$ and the actually executed policy $\Pi_{0:T-1}$ over T -rounds:

$$\text{Regret}(T, \Pi_{0:T-1}) = \sum_{t=0}^{T-1} [R(\boldsymbol{\theta}_t^*, \mathcal{T}_{t+1}^{\mathcal{B}*}) - R(\boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}})]. \quad (10)$$

We can further simplify the definition expression as $\text{Regret}(T, \Pi_{0:T-1}) = \text{CVaR}_{\alpha}(\boldsymbol{\theta}_T) - \text{CVaR}_{\alpha}(\boldsymbol{\theta}_T^*)$, which quantifies the performance gap between the optimal state and the practical state under $(1 - \alpha)$ -tail robustness.

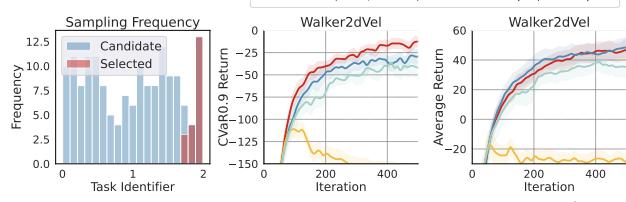


Figure 2. MPTS’s Performance Collapse with Greater $\hat{\mathcal{B}}$. We report the performance collapses of MPTS on Walker2dVel in the case $\hat{\mathcal{B}} = 8\mathcal{B}$. The task sampling frequency reveals the presence of the concentration issue.

Proposition 3.2 (MPTS as a UCB-guided Solution to i-MABs). *Executing MPTS pipeline in Eq. (6) is equivalent to approximately solving \mathcal{M} with the i-MAB under the UCB principle.*

Consequently, our i-MAB is a theoretical model for inducing RATS methods, and MPTS can be viewed as a special case in **Proposition 3.2**.

3.2. Acquisition Functions Matter in Improving Coverage & Boosting Robustness

Several decision-making scenarios, e.g., robotics, are risk-averse, which makes worst-case optimization more advantageous (Greenberg et al., 2024). However, (i) the search scope to worst cases is restricted by the batch size, and (ii) minimax optimization requires carefully designed relaxation in the field (Collins et al., 2020; Sagawa et al., 2019).

Benefits of Enlarging $\hat{\mathcal{B}}$ in Subset Selection. Note that the feasible subset is the arm in the i-MAB, enlarging $\hat{\mathcal{B}}$ increases its number to $|\mathbf{A}_t| = C_{\hat{\mathcal{B}}}^{\mathcal{B}}$. Also, one promising trait of the risk predictive model in MPTS is to amortize the policy evaluation in arbitrary MDP without exact interactions. Hence, greater $\hat{\mathcal{B}}$ in implementation reserves at least two bonus: (i) it encourages exploration in the task space with more candidate subsets at no actual interaction cost; (ii) under high-risk prioritization rule, the optimization pipeline in RATS approximately executes $\text{CVaR}_{1-\frac{\mathcal{B}}{\hat{\mathcal{B}}}}$ in each iteration, i.e., worst-case optimization with $\hat{\mathcal{B}} \rightarrow \infty$.

Unfortunately, MPTS might encounter severe performance collapses with greater $\hat{\mathcal{B}}$, as indicated in Fig. 2 when $\hat{\mathcal{B}} = 8\mathcal{B}$ without a remedy. This empirical result can be attributed to the selected subset’s concentration in a narrow range, which motivates the proposal of heuristic tricks in Meta-RL to adopt a mixed sampling strategy to ensure sufficient visitations to the whole task space (Wang et al., 2025; Greenberg et al., 2024). Although these heuristics mitigate the issue, they do not provide a complete solution, as performance degradation persists with increasing $\hat{\mathcal{B}}$.

Theoretical Diagnosis of Sealed Exploration Potential in MPTS. In practice, we propose to circumvent complicated heuristics and retain the simplicity to enable sufficient explo-

ration in i-MABs. To this end, we analyze the concentration issue and report the theoretical analysis as **Proposition 3.3**.

Proposition 3.3 (Concentration Issue in Average Top- \mathcal{B} Selection). *Let $f(\tau) : \mathbb{R}^d \rightarrow \mathbb{R}$ be a unimodal and continuous function, where $d \in \mathbb{N}^+$ and $\tau \in \mathbb{R}^d$, with a maximum value $f(\tau^*)$ at τ^* . We uniformly sample a set of points $\mathcal{T}^{\hat{\mathcal{B}}} = \{\tau_i\}_{i=1}^{\hat{\mathcal{B}}}$, where τ_i are i.i.d. with a probability p_ϵ of falling within a ϵ -neighborhood of τ^* as $|f(\tau) - f(\tau^*)| \leq \epsilon$. Following MPTS, we select the Top- \mathcal{B} samples with the largest function values, i.e.,*

$$\mathcal{T}^{\mathcal{B}} = \text{Top-}\mathcal{B}(\mathcal{T}^{\hat{\mathcal{B}}}, f), \quad \hat{\mathcal{B}}, \mathcal{B} \in \mathbb{N}^+, \quad \mathcal{B} \leq \hat{\mathcal{B}},$$

For any $\epsilon > 0$ such that $p_\epsilon < \frac{\hat{\mathcal{B}} - \mathcal{B} + 2}{\hat{\mathcal{B}} + 1}$, the concentration probability

$$\mathbb{P}(|f(\tau) - f(\tau^*)| \leq \epsilon | \forall \tau \in \mathcal{T}^{\mathcal{B}})$$

increases with $\hat{\mathcal{B}}$ and converges to 1 with $\hat{\mathcal{B}} \rightarrow \infty$.

As implied, with the increase of $\hat{\mathcal{B}}$ and fixed \mathcal{B} , the Top- \mathcal{B} operator tends to select the subset in a small neighborhood of $\arg \max_{\tau} f(\tau)$, which corresponds to the subset concentration issue in MPTS, over-optimizes a local region and hampers task space exploration. Due to the use of Top- \mathcal{B} operator in batch optimization, we also hypothesize a similar phenomenon in DRM (Wang et al., 2024; Lv et al., 2024). Next, we will propose a natural plausible mechanism to enlarge $\hat{\mathcal{B}}$ in RATS without suffering concentration pains.

Diversity Regularization & Robustness Concept. Our strategy is to encourage the coverage of the task space during subset selection. Specifically, rather than selecting individual candidate tasks based on their acquisition score, we evaluate task batches based on both adaptation risk and task diversities in the subset. This formulation constructs a diversity maximization problem (Wang et al., 2023):

$$\max_{\mathcal{T}^{\mathcal{B}} \subseteq \mathcal{T}^{\hat{\mathcal{B}}} : |\mathcal{T}^{\mathcal{B}}| = \mathcal{B}} \mathcal{A}(\mathcal{T}^{\mathcal{B}}) + \gamma \mathcal{S}[\{d(\tau_i, \tau_j)\}] \quad (11)$$

where \mathcal{S} measures the diversity of the subset $\mathcal{T}^{\mathcal{B}}$ from identifiers’ pairwise distances, e.g., $\sum_{i,j} \| \tau_i - \tau_j \|_2^2$.

Though the involvement of the regularization term makes the subset selection problem NP-hard, it can be solved using simple approximate algorithms (Borodin et al., 2017; Wang et al., 2023).

Proposition 3.4 (Nearly Worst-Case Optimization with PDTS). *When $\hat{\mathcal{B}}$ grows large enough, optimizing the subset from Eq. (11) achieves nearly worst-case optimization.*

Proposition 3.4 indicates that the regularized acquisition rule in Eq. (11) enables exploration of the worst subset from numerous candidate arms while retaining the subset’s coverage of the task space to a certain level. Compared with (Collins et al., 2020; Sagawa et al., 2019), such a regularization will show its effectiveness in experiments.

275 3.3. Practical Sampling with Stochastic Optimism

276 So far, we have presented the i-MAB as a theoretical tool
 277 for RATS and interpreted MPTS as the UCB-guided solution.
 278 However, the UCB acquisition rule requires multiple
 279 stochastic forward passes for each task’s evaluation, and
 280 computations grow with $\hat{\mathcal{B}}$. It also demands calibration of
 281 exploration and exploitation weights in subset search.
 282

283 To cut off unnecessary computations and retain the uncer-
 284 tainty optimism, we adopt the posterior sampling strategy
 285 as the acquisition principle for RATS. The posterior sam-
 286 pling (Osband et al., 2013; Asmuth et al., 2012) is an exten-
 287 sion of Thompson sampling (Thompson, 1933) in solving
 288 MDPs with infinite actions. The reward or action value
 289 is treated as a randomized function, and each arm’s value,
 290 sampled from the posterior once, serves action selection,
 291 i.e., $\mathcal{A}_P(\mathcal{T}^{\mathcal{B}}) = \sum_{i=1}^{\hat{\mathcal{B}}} \hat{\ell}_{t+1,i}$ in Eq. (12).

$$293 \text{One Forward Pass: } z_t \sim q_\phi(z_t | H_t) \quad (12a)$$

$$294 \hat{\ell}_{t+1,i} \sim p_\psi(\ell | \hat{\tau}_i, z_t) \quad \forall i \in \{1, \dots, \hat{\mathcal{B}}\} \quad (12b)$$

$$295 \mathcal{T}_{t+1}^{\mathcal{B}^*} = \arg \max_{\substack{\mathcal{T}^{\mathcal{B}} \subseteq \mathcal{T}^{\hat{\mathcal{B}}} \\ |\mathcal{T}^{\mathcal{B}}| = \mathcal{B}}} \mathcal{A}_P(\mathcal{T}^{\mathcal{B}}) + \gamma \delta [\{d(\tau_i, \tau_j)\}] \quad (12c)$$

299 As implied in (Russo & Van Roy, 2014), posterior sampling
 300 avoids over-exploitation of inaccurate estimated uncertainty,
 301 and benefits more from stochasticity. Together with diver-
 302 sity regularization, we obtain PDTS as a new RATS method
 303 in Eq. (12). Intuitively, diversity penalty and posterior sam-
 304 pling suppress inaccurate Top- \mathcal{B} operations and synergize
 305 exploration in the task space. PDTS enjoys implementation
 306 simplicity, computational efficiency and stochastic optimism
 307 in decision-making.

309 3.4. Overall Implementation

310 Putting the above modules together, we present PDTS in
 311 Algorithm 1, which is easily wrapped into DR and Meta-RL
 312 (See Algorithm 2 and 4).

314 4. Experiments

317 This section presents experimental results on Meta-RL and
 318 robotics DR involving randomized physics properties.

319 **Risk-Averse Baselines.** For fair comparison, we retain
 320 the standard setup in (Wang et al., 2025) and use SOTA
 321 baselines as described in Sec. 2.1: ERM (Vapnik et al.,
 322 1998), DRM (Wang et al., 2024; Greenberg et al., 2024),
 323 GDRM (Sagawa et al., 2019), and MPTS. Following im-
 324 plementations in (Wang et al., 2025; Tao et al., 2024), we
 325 use MAML (Finn et al., 2017) as the backbone algorithm
 326 in Meta-RL and employ TD3 (Fujimoto et al., 2018) and
 327 PPO (Schulman et al., 2017) for domain randomization,
 328 respectively.

Algorithm 1 Posterior-Diversity Synergized 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 task identifier batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.

// **Optimize Risk Predictive Module**

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. (5);

$\psi \leftarrow \psi + \lambda_2 \nabla_\psi \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (5);

end

// **Simulating After-Adaptation Results**

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

// **Posterior Sampling Outcome**

Amortized evaluation $\{\hat{\ell}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ with one stochastic forward pass for all candidate tasks through executing Eq. (12)a-b;

// **Diversity-Guided Subset Search**

Run approximate algorithms to solve Eq. (12)c;

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

We run each algorithm on seven random seeds and report the average performance with standard error of means. For adaptation robustness, we compute CVaR $_\alpha$ values across the validation tasks, with $\alpha = \{0.9, 0.7, 0.5\}$, and also report some out-of-distribution (OOD) results. Importantly, the pseudo batch size is set to $\hat{\mathcal{B}} = 64\mathcal{B}$ as default for PDTs.

4.1. Meta Reinforcement Learning

We consider Meta-RL continuous control scenarios based on MuJoCo (Todorov et al., 2012): Reacher, Walker2d, and HalfCheetah. These include identifiers: target position, velocity, and mass (Wang et al., 2025).

Fig. 3 reveals PDTs consistent superiority over others in CVaR $_{0.9}$ validation return, demonstrating excellent adaptation robustness. Owing to intrinsic sampling randomness, PDTs and MPTS achieve average performance comparable to ERM, while DRM struggles to improve robustness and sacrifices more average returns. Surprisingly, benefiting from exploration, PDTs boosts both adaptation robustness and task efficiency on ReacherPos. We will further discuss this in Sec. 4.2. Regarding meta-testing, PDTs excels across CVaR values, consistent with the learning curves. Notably, PDTs’s performance advantage increases with higher α values. In the extreme scenario CVaR $_{0.9}$, PDTs outperforms others by more than 15% on all benchmarks except HalfCheetahMassVel.

PDTs is agnostic to the meta-learning backbone. PDTs is a plug-and-play module agnostic to the meta-learning backbone. Here, we integrate PDTs with PEARL (Rakelly et al., 2019) and compare it with MPTS and RoML (Greenberg et al., 2024), built on the same backbone. As shown in

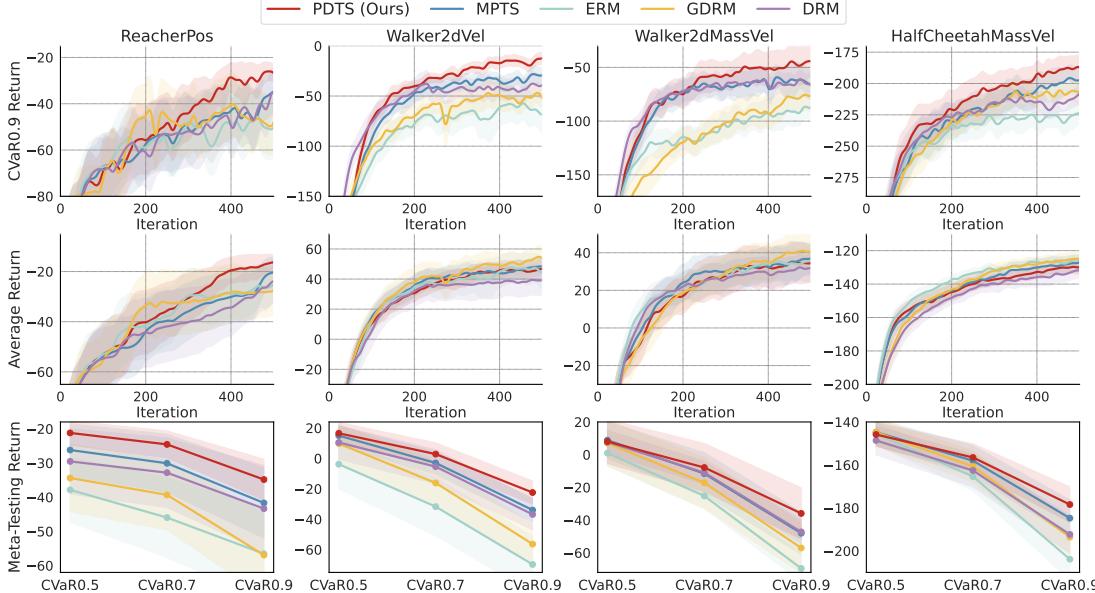


Figure 3. **Meta-RL results.** The top depicts the cumulative return curves for $\text{CVaR}_{0.9}$ validation MDPs during meta-training; the middle shows the average cumulative returns curves during meta-training; and the bottom presents the meta-testing results with various α .

Table 2, we evaluate these methods on two Meta-RL scenarios from RoML. PDTs outperforms others in terms of CVaR return, further examining its effectiveness and scalability in risk-averse decision-making.

Table 2. PDTs’s compatibility with PEARL backbone. Baselines are MPTS and RoML (Greenberg et al., 2024).

$\text{CVaR}_{0.95}$	Return	PDTs	MPTS	RoML	PEARL
HalfCheetahBody	993±26	945±26	855±35	847±42	
HalfCheetahMass	1296±41	1209±45	1197±59	1118±51	

4.2. Physical Robotics Domain Randomization

We evaluate PDTs on three DR scenarios introduced by Mehta et al. (2020): a game scenario, LunarLander, and two robotic arm control scenarios, Pusher and ErgoReacher.

Empirical findings in Fig. 4(a) are consistent with those in Meta-RL, where PDTs dominates the robustness performance. Notably, PDTs’s advantage is more pronounced on Pusher and LunarLander, where the identifiers’ dimension is lower than ErgoReacher’s. This is likely because low-dimensional task identifiers exacerbate the concentration issue, limiting the performance of MPTS. GDRM and DRM perform poorly in both average performance and robustness. Given a fixed batch size, this weakness may stem from their limited exploration capacity. Regarding final testing performance, PDTs excels across CVaR values. Specifically, PDTs outperforms ERM by more than 8% on all benchmarks in $\text{CVaR}_{0.9}$, and by as much as 73% on LunarLander.

PDTs improves task efficiency in specific scenarios. Beyond adaptation robustness, PDTs shows the potential to accelerate training in specific scenarios. As shown in Fig. 4(a), PDTs achieves average returns comparable to ERM’s best performance but with fewer training steps, achieving $2.4\times$

acceleration in Pusher and $1.3\times$ in LunarLander. We attribute this acceleration to the efficient exploration of acquisition criteria. Hence, PDTs holds promise for achieving robustness with reduced training steps in broader scenarios.

PDTs achieves superior zero-shot policy adaptation in OOD MDPs. On Lundarlander, Fig. 4(b) witnesses PDTs’s highest average test returns across sampled tasks. To evaluate zero-shot adaptation in OOD MDPs, we shift the identifier interval $\tau \in [4.0, 20.0]$ to the OOD range $\tau \in [1.0, 4.0] \cup (20.0, 23.0]$. Notably, PDTs exhibits the smallest performance degradation, particularly in the most challenging OOD tasks ($\tau \in [1.0, 4.0)$).

4.3. Visual Robotics Domain Randomization

Vision-based robotics control is more challenging and underexplored in MPTS. With the latest robotics simulator, ManiSkill3 (Tao et al., 2024), we design two visual DR scenarios: one randomizing lighting with a table-top two-finger gripper arm robot called LiftPegUpright_Light, and another randomizing goal locations with a quadruped robot called AnymalCReach_Goal, as illustrated in Fig. 5(a).

PDTs achieves best robustness performance in Fig. 5(b), same as symbolic scenarios. Moreover, PDTs also exhibits a trend of accelerated training in terms of average performance, highlighting its potential for realistic scenarios requiring both fast and robust adaptation. Due to insufficient task exploration, DRM performs poorly and probably gets trapped in challenging tasks. ERM obtains the worst final performance on AnymalCReach_Goal. These findings underscore the importance of robust optimization with sufficient exploration, which may explain PDTs’s success.

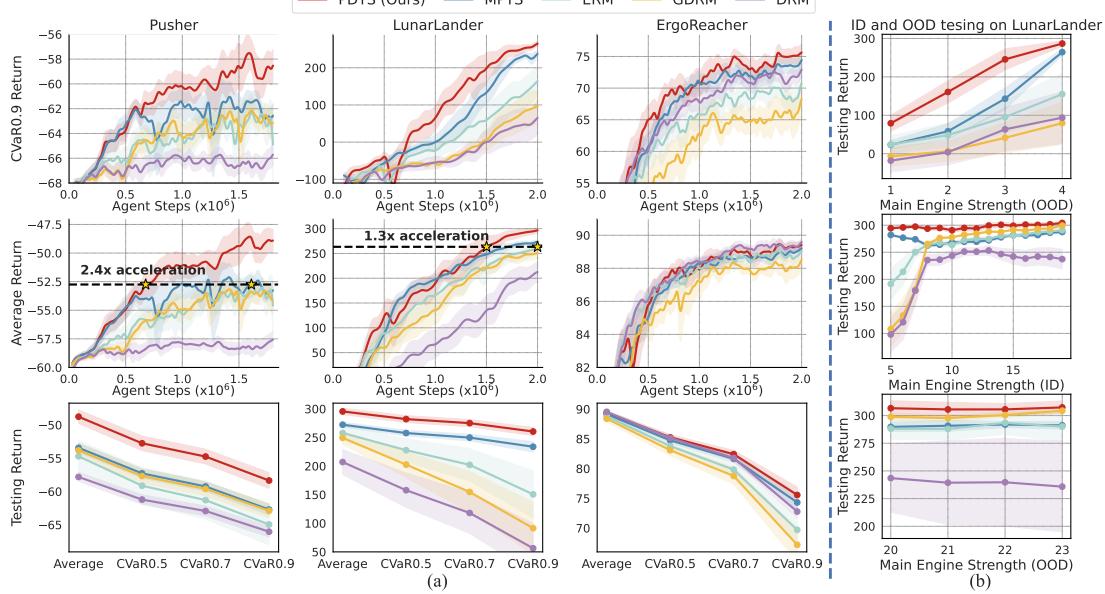


Figure 4. **Physical Robotics DR results.** (a) The top shows the cumulative return curves for CVaR_{0.9} validation MDPs during training; the middle displays the average cumulative return curves across all validation MDPs during training; and the bottom presents the test results at various CVaR _{α} . (b) We evaluate the trained policies in both in-distribution (ID) and out-of-distribution (OOD) domains on LunarLander, reporting the average returns for each sampled task.

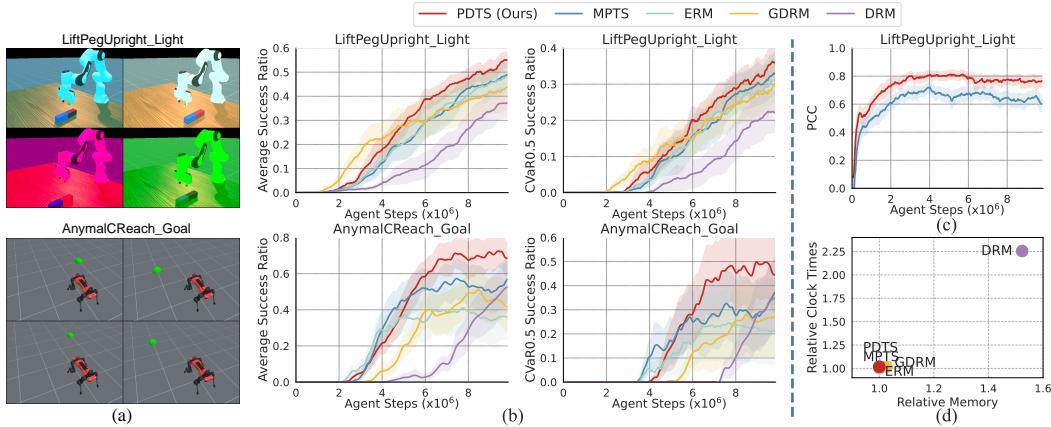


Figure 5. **Visual Robotics DR results.** (a) Illustrations of two scenarios. (b) Curves of the average success ratio and the CVaR_{0.5} success ratio on validation tasks during training. (c) Training curves of PCC values between predicted and true episode returns. (d) Memory cost and clock time relative to ERM during meta-training.

PDTs Can Discriminate Task Difficulties. We evaluate the Pearson Correlation Coefficient (PCC) between the predicted episode returns and the exact values during training. As shown in Fig. 5(c), PDTs inherits MPTS’s difficulty scoring capability of MPTS, with both PCC values greater than 0.5. In particular, PDTs exceeds MPTS in predicting accuracies, which may be attributed to the extended exploration and the stochastic optimism in posterior sampling.

PDTs Offers Advantages at Minimal Cost. Like MPTS, PDTs avoids additional evaluation costs associated with interaction and rendering, making it more beneficial in vision-based RL tasks. From relative clock time and memory usage in Fig. 5(d), we find PDTs retains computational efficiency with others except DRM. As a result, PDTs is a robust and scalable approach for complex, interaction-intensive tasks.

5. Conclusion

This work studies RATS’s use for learning adaptive decision-makers in risk-averse decision-making. We present the i-MAB as a theoretical tool to enable RATS and establish theoretical connections with its latest method (Wang et al., 2025). Built on these insights, we propose PDTs as a competitive RATS method to achieve nearly worst-case task optimization. Extensive DR and Meta-RL experiments validate the effectiveness of RATS and show its efficiency potential in risk-averse scenarios.

Future explorations can include discovering more risk-predictive models to enhance amortized evaluation reliability and integrating more robust optimization rules into i-MABs to develop practical RATS methods.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, including enhancing real-world applications of generalized robotics and other decision-making systems. Although this study significantly improves adaptation robustness, adaptation failures remain possible, potentially leading to serious consequences in real-world scenarios. Thus, human supervision remains essential.

References

- Abbas, M., Xiao, Q., Chen, L., Chen, P.-Y., and Chen, T. Sharp-maml: Sharpness-aware model-agnostic meta learning. In *International conference on machine learning*, pp. 10–32. PMLR, 2022.
- Asmuth, J., Li, L., Littman, M. L., Nouri, A., and Wingate, D. A bayesian sampling approach to exploration in reinforcement learning. *arXiv preprint arXiv:1205.2664*, 2012.
- Auer, P. Finite-time analysis of the multiarmed bandit problem, 2002a.
- Auer, P. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002b.
- Beck, J., Vuorio, R., Liu, E. Z., Xiong, Z., Zintgraf, L., Finn, C., and Whiteson, S. A survey of meta-reinforcement learning. *arXiv preprint arXiv:2301.08028*, 2023.
- Borodin, A., Jain, A., Lee, H. C., and Ye, Y. Max-sum diversification, monotone submodular functions, and dynamic updates. *ACM Transactions on Algorithms (TALG)*, 13(3):1–25, 2017.
- Broderick, T., Boyd, N., Wibisono, A., Wilson, A. C., and Jordan, M. I. Streaming variational bayes. *Advances in neural information processing systems*, 26, 2013.
- Carpentier, A. and Valko, M. Simple regret for infinitely many armed bandits. In *International Conference on Machine Learning*, pp. 1133–1141. PMLR, 2015.
- Carpin, S., Chow, Y.-L., and Pavone, M. Risk aversion in finite markov decision processes using total cost criteria and average value at risk. In *2016 ieee international conference on robotics and automation (icra)*, pp. 335–342. IEEE, 2016.
- Catto, E. Box2d: A 2d physics engine for games, 2007. URL <http://box2d.org>.
- Chow, Y. *Risk-sensitive and data-driven sequential decision making*. PhD thesis, Stanford University, 2017.
- Chow, Y., Tamar, A., Mannor, S., and Pavone, M. Risk-sensitive and robust decision-making: a cvar optimization approach. *Advances in neural information processing systems*, 28, 2015.
- Chow, Y., Ghavamzadeh, M., Janson, L., and Pavone, M. Risk-constrained reinforcement learning with percentile risk criteria. *Journal of Machine Learning Research*, 18(167):1–51, 2018.
- Chow, Y., Nachum, O., Faust, A., Dueñez-Guzman, E., and Ghavamzadeh, M. Safe policy learning for continuous control. In *Conference on Robot Learning*, pp. 801–821. PMLR, 2021.
- Cohn, D. A., Ghahramani, Z., and Jordan, M. I. Active learning with statistical models. *Journal of artificial intelligence research*, 4:129–145, 1996.
- Collins, L., Mokhtari, A., and Shakkottai, S. Task-robust model-agnostic meta-learning. *Advances in Neural Information Processing Systems*, 33:18860–18871, 2020.
- Coumans, E. Bullet physics simulation. In *ACM SIGGRAPH 2015 Courses*, pp. 1. 2015.
- Duan, Y., Schulman, J., Chen, X., Bartlett, P. L., Sutskever, I., and Abbeel, P. RL2: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Finn, C., Xu, K., and Levine, S. Probabilistic model-agnostic meta-learning. *Advances in neural information processing systems*, 31, 2018.
- Fujimoto, S., Hoof, H., and Meger, D. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- Gal, Y., Islam, R., and Ghahramani, Z. Deep bayesian active learning with image data. In *International conference on machine learning*, pp. 1183–1192. PMLR, 2017.
- Garivier, A. and Moulines, E. On upper-confidence bound policies for non-stationary bandit problems. *arXiv preprint arXiv:0805.3415*, 2008.
- Greenberg, I., Chow, Y., Ghavamzadeh, M., and Mannor, S. Efficient risk-averse reinforcement learning. *Advances in Neural Information Processing Systems*, 35:32639–32652, 2022.

- 495 Greenberg, I., Mannor, S., Chechik, G., and Meirom, E.
 496 Train hard, fight easy: Robust meta reinforcement learn-
 497 ing. *Advances in Neural Information Processing Systems*,
 498 36, 2024.
- 499 Muratore, F., Treede, F., Gienger, M., and Peters, J. Domain
 500 randomization for simulation-based policy optimization
 501 with transferability assessment. In *Conference on Robot
 502 Learning*, pp. 700–713. PMLR, 2018.
- 503 Gupta, A., Mendonca, R., Liu, Y., Abbeel, P., and Levine,
 504 S. Meta-reinforcement learning of structured exploration
 505 strategies. *Advances in neural information processing
 506 systems*, 31, 2018.
- 507 Muratore, F., Eilers, C., Gienger, M., and Peters, J. Data-
 508 efficient domain randomization with bayesian optimiza-
 509 tion. *IEEE Robotics and Automation Letters*, 6(2):911–
 510 918, 2021.
- 511 Nguyen, C. V., Li, Y., Bui, T. D., and Turner,
 512 R. E. Variational continual learning. *arXiv preprint
 513 arXiv:1710.10628*, 2017.
- 514 Osband, I., Russo, D., and Van Roy, B. (more) efficient
 515 reinforcement learning via posterior sampling. *Advances
 516 in Neural Information Processing Systems*, 26, 2013.
- 517 Pan, X., Seita, D., Gao, Y., and Canny, J. Risk averse
 518 robust adversarial reinforcement learning. In *2019 Inter-
 519 national Conference on Robotics and Automation (ICRA)*,
 520 pp. 8522–8528. IEEE, 2019.
- 521 Puterman, M. L. *Markov decision processes: discrete
 522 stochastic dynamic programming*. John Wiley & Sons,
 523 2014.
- 524 Qi, Y., Ban, Y., Wei, T., Zou, J., Yao, H., and He, J. Meta-
 525 learning with neural bandit scheduler. *Advances in Neural
 526 Information Processing Systems*, 36, 2024.
- 527 Rajeswaran, A., Finn, C., Kakade, S. M., and Levine, S.
 528 Meta-learning with implicit gradients. *Advances in neural
 529 information processing systems*, 32, 2019.
- 530 Rakelly, K., Zhou, A., Finn, C., Levine, S., and Quillen, D.
 531 Efficient off-policy meta-reinforcement learning via prob-
 532 abilistic context variables. In *International conference on
 533 machine learning*, pp. 5331–5340. PMLR, 2019.
- 534 Rempe, D., Phlion, J., Guibas, L. J., Fidler, S., and Litany,
 535 O. Generating useful accident-prone driving scenarios via
 536 a learned traffic prior. In *Proceedings of the IEEE/CVF
 537 Conference on Computer Vision and Pattern Recognition*,
 538 pp. 17305–17315, 2022.
- 539 Mahajan, A. and Teneketzis, D. Multi-armed bandit prob-
 540 lems. In *Foundations and applications of sensor manage-
 541 ment*, pp. 121–151. Springer, 2008.
- 542 Mehta, B., Diaz, M., Golemo, F., Pal, C. J., and Paull, L.
 543 Active domain randomization. In *Conference on Robot
 544 Learning*, pp. 1162–1176. PMLR, 2020.
- 545 Mukhoti, J., Kirsch, A., van Amersfoort, J., Torr, P. H., and
 546 Gal, Y. Deep deterministic uncertainty: A new simple
 547 baseline. In *Proceedings of the IEEE/CVF Conference
 548 on Computer Vision and Pattern Recognition*, pp. 24384–
 549 24394, 2023.
- 550 Ren, P., Xiao, Y., Chang, X., Huang, P.-Y., Li, Z., Gupta,
 551 B. B., Chen, X., and Wang, X. A survey of deep active
 552 learning. *ACM computing surveys (CSUR)*, 54(9):1–40,
 553 2021.
- 554 Rigter, M., Lacerda, B., and Hawes, N. Risk-averse bayes-
 555 adaptive reinforcement learning. *Advances in Neural
 556 Information Processing Systems*, 34:1142–1154, 2021.
- 557 Ritter, S., Wang, J., Kurth-Nelson, Z., Jayakumar, S., Blun-
 558 dell, C., Pascanu, R., and Botvinick, M. Been there, done
 559 that: Meta-learning with episodic recall. In *International
 560 Conference on Machine Learning*, pp. 7353–7362. PMLR,
 561 2021.

- 550 conference on machine learning, pp. 4354–4363. PMLR,
 551 2018.
- 552 Rockafellar, R. T., Uryasev, S., et al. Optimization of conditional
 553 value-at-risk. *Journal of risk*, 2:21–42, 2000.
- 554
- 555 Russo, D. and Van Roy, B. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39
 556 (4):1221–1243, 2014.
- 557
- 558 Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P.
 559 Distributionally robust neural networks. In *International Conference on Learning Representations*, 2019.
- 560
- 561 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and
 562 Klimov, O. Proximal policy optimization algorithms.
 563 *arXiv preprint arXiv:1707.06347*, 2017.
- 564
- 565 Setlur, A., Garg, S., Smith, V., and Levine, S. Prompting
 566 is a double-edged sword: Improving worst-group robustness of foundation models. In *Forty-first International Conference on Machine Learning*, 2024.
- 567
- 568 Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- 569
- 570 Tamar, A., Chow, Y., Ghavamzadeh, M., and Mannor, S.
 571 Policy gradient for coherent risk measures. *Advances in neural information processing systems*, 28, 2015.
- 572
- 573 Tao, S., Xiang, F., Shukla, A., Qin, Y., Hinrichsen, X.,
 574 Yuan, X., Bao, C., Lin, X., Liu, Y., Chan, T.-k., et al.
 575 Maniskill3: Gpu parallelized robotics simulation and
 576 rendering for generalizable embodied ai. *arXiv preprint arXiv:2410.00425*, 2024.
- 577
- 578 Thompson, W. R. On the likelihood that one unknown
 579 probability exceeds another in view of the evidence of
 580 two samples. *Biometrika*, 25(3-4):285–294, 1933.
- 581
- 582 Tiboni, G., Klink, P., Peters, J., Tommasi, T., D’Eramo, C.,
 583 and Chalvatzaki, G. Domain randomization via entropy
 584 maximization. *arXiv preprint arXiv:2311.01885*, 2023.
- 585
- 586 Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W.,
 587 and Abbeel, P. Domain randomization for transferring
 588 deep neural networks from simulation to the real world.
 589 In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pp. 23–30. IEEE, 2017.
- 590
- 591 Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics
 592 engine for model-based control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp.
 593 5026–5033. IEEE, 2012.
- 594
- 595 Vapnik, V. N., Vapnik, V., et al. Statistical learning theory.
 596 1998.
- 597
- 598 Wang, Q. and Van Hoof, H. Learning expressive meta-
 599 representations with mixture of expert neural processes.
 600 *Advances in neural information processing systems*, 35:
 601 26242–26255, 2022a.
- 602
- 603 Wang, Q. and Van Hoof, H. Model-based meta reinforcement learning using graph structured surrogate models
 604 and amortized policy search. In *International Conference on Machine Learning*, pp. 23055–23077. PMLR, 2022b.
- 605
- 606 Wang, Q., Lv, Y., Xie, Z., Huang, J., et al. A simple yet
 607 effective strategy to robustify the meta learning paradigm.
 608 *Advances in Neural Information Processing Systems*, 36,
 609 2024.
- 610
- 611 Wang, Q. C., Xiao, Z., Mao, Y., Qu, Y., Shen, J., Lv, Y., and
 612 Ji, X. Beyond any-shot adaptation: Predicting optimization
 613 outcome for robustness gains without extra pay, 2025.
 614 URL <https://arxiv.org/abs/2501.11039>.
- 615
- 616 Wang, Y., Mathioudakis, M., Li, J., and Fabbri, F. Max-
 617 min diversification with fairness constraints: Exact and
 618 approximation algorithms. In *Proceedings of the 2023 SIAM International Conference on Data Mining (SDM)*,
 619 pp. 91–99. SIAM, 2023.
- 620
- 621 Wu, J., Chen, J., and Huang, D. Entropy-based active
 622 learning for object detection with progressive diversity
 623 constraint. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9397–
 624 9406, 2022.
- 625
- 626 Xu, F., Jiang, S., Yin, H., Zhang, Z., Yu, Y., Li, M., Li, D.,
 627 and Liu, W. Enhancing context-based meta-reinforcement
 628 learning algorithms via an efficient task encoder (student
 629 abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 15937–15938, 2021.
- 630
- 631 Yoon, J., Kim, T., Dia, O., Kim, S., Bengio, Y., and Ahn,
 632 S. Bayesian model-agnostic meta-learning. *Advances in neural information processing systems*, 31, 2018.
- 633
- 634 Zhu, X., Lafferty, J., and Ghahramani, Z. Combining active
 635 learning and semi-supervised learning using gaussian
 636 fields and harmonic functions. In *ICML 2003 workshop on the continuum from labeled to unlabeled data in machine learning and data mining*, volume 3, pp. 58–65, 2003.
- 637
- 638 Zintgraf, L., Shiarlis, K., Igl, M., Schulze, S., Gal, Y., Hof-
 639 mann, K., and Whiteson, S. Varibad: A very good method
 640 for bayes-adaptive deep rl via meta-learning. *arXiv preprint arXiv:1910.08348*, 2019.

605 A. Adaptive Decision-Making in Randomized Environments

606 A.1. Zero-Shot and Few-Shot Sequential Decision-Making

608 In this part, we include typical MDPs and distributions, which align with randomized environments for sequential decision-making.
 609 Meanwhile, we show examples of DR and Meta-RL in this context, which can well specify the form of the
 610 task-specific adaptation risk ℓ in zero-shot and a few-shot setup.

611 **Definition A.1** (Identifier-Induced MDP Distribution). *We consider a distribution over MDPs $p(\mathcal{M}_\tau)$ induced by the
 612 identifier distribution $p(\tau)$. Here, the MDP as the sequential decision-making environment can be characterized as
 613 $\mathcal{M}_\tau = \langle \mathbf{S}, \mathbf{A}, \mathbf{P}_\tau, \mathbf{R}_\tau, \gamma \rangle$, where \mathbf{S} and \mathbf{A} are the state and the action space shared across all MDPs.*

615 *MDPs vary in either the transition dynamics $\mathbf{P}_\tau := p_\tau(s_{t+1}|s_t, a_t)$ or the reward function $\mathbf{R}_\tau :=$
 616 $R_\tau(s_t, a_t)$ or the both, specified by corresponding identifiers τ . The H -horizon rollout under a policy is
 617 $(s_0, a_0, R(s_0, a_0), \dots, s_{H-1}, a_{H-1}, R(s_{H-1}, a_{H-1}), s_H)$ with its cumulative reward $\sum_{t=0}^{H-1} \gamma^t R(s_t, a_t)$ and $\gamma \in \mathbb{R}^+$.*

618 With Definition A.1, we can provide some instantiations of DR and Meta-RL under specific policy optimization backbones.

619 **Meta-RL with Model Agnostic Meta Learning.** Meta-RL seeks to train agents capable of rapidly adapting to new tasks or
 620 environments by utilizing knowledge from related tasks. We adopt MAML (Finn et al., 2017) as the standard backbone
 621 algorithm due to its widespread application in addressing few-shot sequential decision-making problems. This also retains
 622 the same configuration in MPTS (Wang et al., 2025). The optimization objective of MAML is mathematically defined as:
 623

$$625 \min_{\theta \in \Theta} \mathbb{E}_{p(\tau)} [\ell(\mathcal{D}_\tau^Q; \theta - \lambda \nabla_\theta \ell(\mathcal{D}_\tau^S; \theta))] \quad (13a)$$

$$627 \min_{\theta \in \Theta} \mathbb{E}_{p(\tau)} \left[-\mathbb{E}_{\pi_{\theta'}, \mathcal{D}_\tau^Q} \left[\sum_{t=0}^H \gamma^t r_t \right]; \theta' = \theta - \lambda \nabla_\theta \left(-\mathbb{E}_{\pi_\theta, \mathcal{D}_\tau^S} \left[\sum_{t=0}^H \gamma^t r_t \right] \right) \right], \quad (13b)$$

630 where \mathcal{D}_τ denotes episodes collected from MDPs specified by τ under either the meta-policy or the fast-adapted policy. The
 631 term inside the brackets specifies the adaptation risk $\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)$, which represents either the negative cumulative returns
 632 or the negative final rewards. The expression $\theta - \lambda \nabla_\theta \ell(\mathcal{D}_\tau^S; \theta)$ indicates the gradient update for fast task adaptation, where
 633 λ is the learning rate. Upon completion of meta-training, the resulting meta-policy θ can be generalized across the task
 634 space.

635 **Robotics Domain Randomization.** Robotics domain randomization refers to a training paradigm in which an agent is
 636 trained across a collection of diverse environments to develop a generalizable policy. The diversity of these environments
 637 enhances the robustness of the resulting policy during deployment. Notably, this approach eliminates the need for few-shot
 638 episodes in unseen but similar environments. Mathematically, the optimization objective can be expressed as:

$$640 \max_{\theta \in \Theta} \mathcal{J}(\theta) := \mathbb{E}_{\pi_\theta} \mathbb{E}_{p(\tau)} \left[\sum_{t=0}^H \gamma^t r_t \right], \quad (14)$$

644 where $p(\tau)$ denotes the distribution over MDPs specified by τ , and $\{r_t\}_{t=0}^H$ represents the stepwise rewards obtained from
 645 interacting with a specific MDP, with H as the horizon.

646 After solving the optimization problem in Eq. (14), the resulting policy π_θ serves as a zero-shot decision-maker in new
 647 environments. In this context, the adaptation risk can be expressed as $\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) = -\sum_{t=0}^H \gamma^t r_t$. For physical domain
 648 randomization, policy optimization utilizes the TD3 algorithm (Fujimoto et al., 2018), an off-policy method known for its
 649 sample efficiency and stability. For visual domain randomization, we employ the PPO algorithm (Schulman et al., 2017), an
 650 on-policy method recommended by ManiSkill3 (Tao et al., 2024).

652 A.2. Related Work

654 **Risk-Averse Methods and Robustness Evaluation.** Real-world decision-making scenarios are risk-sensitive (Chow et al.,
 655 2021), and this makes robust optimization an indispensable component of policy optimization (Tamar et al., 2015; Chow
 656 et al., 2015; Chow, 2017; Chow et al., 2018; Pan et al., 2019; Rigter et al., 2021; Greenberg et al., 2022). Recent advances
 657 turn to some risk-averse measures to improve robustness in the task distribution. These are detailed in Section 2.1, which
 658 includes the DRM (Lv et al., 2024; Wang et al., 2024), GDRM (Sagawa et al., 2019), and MPTS (Wang et al., 2025). The
 659

660 existing worst-case optimization in meta-learning is (Collins et al., 2020), which relies on special relaxation tricks; otherwise,
 661 the optimization can be unstable. Besides, other inspiring robust methods in task sampling (Qi et al., 2024) have not been
 662 applied in DR or Meta-RL. As MPTS is the latest SOTA RATS method, we mainly compared PDTS with it in experiments.
 663 RoML (Greenberg et al., 2024) also considers the CVaR optimization in Meta-RL; hence, we include some comparisons
 664 given the PEARL (Rakelly et al., 2019) backbone. In terms of evaluation, the subpopulation shift and out-of-distribution
 665 scenarios are commonly used in the field (Koh et al., 2021).

666 **Cross-Task Adaptation in Sequential Decision-Making.** This work focuses on the zero-shot and few-shot decision-making
 667 scenarios. In zero-shot decision-making, the primary technique is to randomize the environments and train the agent in the
 668 distribution over environments, which is called domain randomization (Tobin et al., 2017; Muratore et al., 2018; Mehta et al.,
 669 2020; Muratore et al., 2021; Tiboni et al., 2023). Meta-learning is a promising paradigm for achieving few-shot adaptation to
 670 unseen tasks, avoiding learning from scratch (Hospedales et al., 2021). The secret behind its few-shot adaptation capability
 671 is to leverage past experience and consolidate these as the prior for fast problem-solving. There are three primary types of
 672 policy adaptation methods that can be seamlessly incorporated into deep RL scenarios. (i) The optimization-based methods
 673 aim to seek a robust meta-initialization of the model that can be adapted to unseen scenarios through fine-tuning (Finn
 674 et al., 2017; 2018; Abbas et al., 2022; Rajeswaran et al., 2019; Yoon et al., 2018; Gupta et al., 2018; Qi et al., 2024). (ii)
 675 The context-based methods mostly adopt an encoder and decoder structure in policy networks (Xu et al., 2021; Wang
 676 & Van Hoof, 2022a;b; Rakelly et al., 2019; Zintgraf et al., 2019; Li et al., 2020). The encoder summarizes the support
 677 trajectories into a latent variable to guide the policy adaptation. (iii) The recurrent methods mainly employ a recurrent neural
 678 network to encode the sequential decision-making episodes as the adaptation prior (Duan et al., 2016; Ritter et al., 2018).
 679 All of these zero-shot or few-shot learning is performed in a task episodic way, which means that a batch of MDPs are
 680 resampled in each iteration to train the adaptive policy during DR and Meta-RL.

A.3. CVaR $_{\alpha}$ as Risk-Averse Metrics

684 **Definition A.2** (CVaR $_{\alpha}$). *With θ -parameterized model, e.g., a deep RL policy, we can induce a random variable $\ell_i :=$
 685 $\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta)$ from $p(\tau)$. Then, we can define the cumulative adaptation risk distribution and its quantile by $F(\ell)$ and
 686 $\ell^{\alpha} = \min_{\ell} \{\ell | F(\ell) \geq \alpha\}$. Following (Rockafellar et al., 2000), CVaR at α -level robustness can be expressed as:*

$$\text{CVaR}_{\alpha}[\ell(\mathcal{T}; \theta)] = \int \ell dF^{\alpha}(\ell; \theta), \quad (15)$$

687 Accordingly, the normalized tail risk task distribution is

$$F^{\alpha}(\ell; \theta) = \begin{cases} 0, & \ell < \ell^{\alpha} \\ \frac{F(\ell; \theta) - \alpha}{1 - \alpha}, & \ell \geq \ell^{\alpha}, \end{cases} \quad (16)$$

688 with $p_{\alpha}(\tau; \theta)$ as its probability density function.

689 **Assumption 1** (Lipschitz Continuity). The adaptation risk function $\ell(\cdot; \theta)$ is β_{τ} -Lipschitz continuous w.r.t. θ and β_{θ} -
 690 Lipschitz w.r.t. τ , i.e.,

$$\begin{aligned} |\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) - \ell(\mathcal{D}_{\tau'}^Q, \mathcal{D}_{\tau'}^S; \theta')| &\leq \beta_{\tau} \|\theta - \theta'\| \\ |\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) - \ell(\mathcal{D}_{\tau'}^Q, \mathcal{D}_{\tau'}^S; \theta)| &\leq \beta_{\theta} \|\tau - \tau'\|, \end{aligned} \quad (17)$$

701 where $\forall \{\theta, \theta'\} \in \Theta$ and $\forall \{\tau, \tau'\} \in \mathcal{T}$.

702 **Assumption 2** (Bounded Functions). Given arbitrary $\theta \in \Theta$ and the task $\tau \in \mathcal{T}$, the adaptation risk function $\ell(\cdot; \theta)$ satisfies:

$$\max_{\tau \in \mathcal{T}} \ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \theta) \leq \ell_{\max}. \quad (18)$$

703 CVaR $_{\alpha}$ is a risk-averse measure in computational finance and management science (Linsmeier & Pearson, 2000; Rockafellar
 704 et al., 2000), and we write it in the Definition A.2. As for Assumptions 1 and 2, these are commonly seen in analysis
 705 (Lv et al., 2024; Wang et al., 2024) and are also necessary to MPTS and PDTS. These constitute the predictability of the
 706 adaptation risk value over iterations as the relative task difficulty remains invariant after the model's parameter is perturbed a
 707 bit.

A.4. Pseudo Algorithm in Meta-RL and DR

713 In this part, we show how PDTS is incorporated in DR and Meta-RL. These are Algorithm 2/3 and Algorithm 4/5.

715
 716 **Algorithm 2** PDTS for Domain Randomization (Zero-Shot
 717 Scenarios)
 718 **Input** : Task distribution $p(\tau)$; Task batch size \mathcal{B} ; Learn-
 719 ing rate λ_1 .
 720 **Output** : Adapted policy θ .
 721 Set the initial iteration number $t = 1$;
 722 Randomly initialize policy θ ;
 723 Randomly initialize risk learner $\{\psi, \phi\}$;
 724 **while** not converged **do**
 725 Execute **Algorithm 3** to access the task batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$;
 726 Sample trajectories for each task with policy θ to induce
 727 $\{\mathcal{D}_{\tau_{t,i}}^Q\}_{i=1}^{\mathcal{B}}$;
 728 // **Evaluation Performance**
 729 Compute the task specific adaptation risk $\{\ell_{t,i} :=$
 730 $-\mathbb{E}_{\pi_\theta, \mathcal{D}_{\tau_{t,i}}^Q} \left[\sum_{t=0}^H \gamma^t r_t \right]\}_{i=1}^{\mathcal{B}}$;
 731 Return $H_t = \{[\tau_{t,i}, \ell_{t,i}]\}_{i=1}^{\mathcal{B}}$ as the Input to **Algorithm**
 732 3;
 733 // **Update Policy**
 734 Perform batch gradient updates:
 735 $\theta_{t+1} \leftarrow \theta_t - \frac{\lambda_1}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \nabla_\theta \ell_{t,i}$;
 736 Update the iteration number: $t \leftarrow t + 1$;
 737 **end**

738
 739 **Algorithm 3** Posterior-Diversity Synergized Task Sampling
 740 **Input** : Task distribution $p(\tau)$; Task batch size \mathcal{B} ; Candi-
 741 date batch size $\hat{\mathcal{B}}$; Latest updated $\{\psi, \phi\}$; Latest
 742 history H_{t-1} ; Iteration number K ; Learning rate
 743 λ_2 .
 744 **Output**: Selected task identifier batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.
 745 // **Optimize Risk Predictive Module**
 746 **for** $i = 1$ **to** K **do**
 747 Perform gradient updates given H_{t-1} :
 748 $\phi \leftarrow \phi + \lambda_2 \nabla_\phi \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (5);
 749 $\psi \leftarrow \psi + \lambda_2 \nabla_\psi \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (5);
 750 **end**
 751 // **Simulating After-Adaptation Results**
 752 Randomly sample $\{\hat{\tau}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ from $p(\tau)$;
 753 // **Posterior Sampling Outcome**
 754 Amortized evaluation $\{\hat{\ell}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ with one stochastic forward
 755 pass for all candidate tasks through executing Eq. (12)a-b;
 756 // **Diversity-Guided Subset Search**
 757 Run approximate algorithms to solve Eq. (12)c;
 758 Return the screened subset $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$ for next iteration.

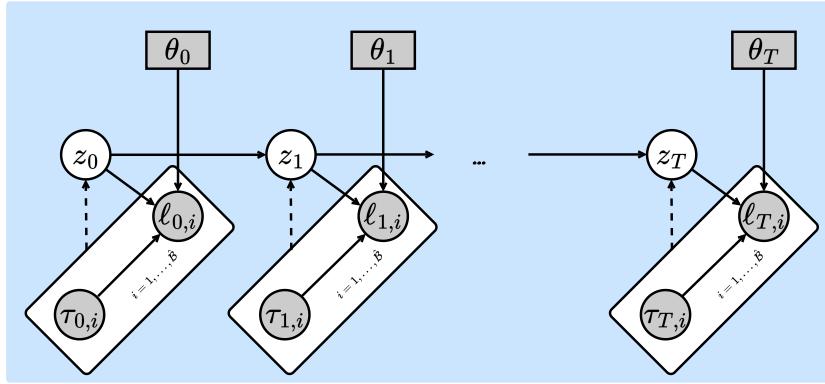
759
 760 **Algorithm 4** PDTS for Model Agnostic Meta Learning
 761 (Few-Shot Scenarios: Meta-RL, Sinusoid Regression)
 762 **Input** : Task distribution $p(\tau)$; Task batch size \mathcal{B} ; Learn-
 763 ing rates: $\{\lambda_{1,1}, \lambda_{1,2}\}$.
 764 **Output**: Meta-trained initialization θ^{meta} .
 765 Set the initial iteration number $t = 1$;
 766 Randomly initialize meta policy θ^{meta} ;
 767 Randomly initialize risk learner $\{\psi, \phi\}$;
 768 **while** not converged **do**
 769 Execute **Algorithm 5** to access the batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$;
 770 Sample trajectories for each task with policy θ^{meta} to
 771 induce $\{\mathcal{D}_{\tau_{t,i}}^S\}_{i=1}^{\mathcal{B}}$;
 772 // **Inner Loop to Fast Adapt**
 773 **for** $i = 1$ **to** K **do**
 774 Compute the task-specific adaptation risk:
 775 $\ell(\mathcal{D}_{\tau_{t,i}}^S; \theta) = -\mathbb{E}_{\pi_\theta, \mathcal{D}_{\tau_{t,i}}^S} \left[\sum_{t=0}^H \gamma^t r_t \right]$;
 776 Perform gradient updates as fast adaptation:
 777 $\theta_t^i \leftarrow \theta_t^{\text{meta}} - \lambda_{1,1} \nabla_\theta \ell(\mathcal{D}_{\tau_{t,i}}^S; \theta)$;
 778 Sample trajectories $\mathcal{D}_{\tau_{t,i}}^Q$ with policy θ_t^i for task τ_i ;
 779 **end**
 780 // **Outer Loop to Meta-train**
 781 Evaluate fast adaptation performance $\{\ell_{t,i} :=$
 782 $-\mathbb{E}_{\pi_{\theta_t^i}, \mathcal{D}_{\tau_{t,i}}^Q} \left[\sum_{t=0}^H \gamma^t r_t \right]\}_{i=1}^{\mathcal{B}}$;
 783 Return $H_t = \{[\tau_{t,i}, \ell_{t,i}]\}_{i=1}^{\mathcal{B}}$ as the Input to **Algorithm**
 784 5;
 785 Perform meta initialization updates:
 786 $\theta_{t+1}^{\text{meta}} \leftarrow \theta_t^{\text{meta}} - \frac{\lambda_{1,2}}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \nabla_\theta \ell_{t,i}$;
 787 Update the iteration number: $t \leftarrow t + 1$;
 788 **end**

789
 790 **Algorithm 5** Posterior-Diversity Synergized Task Sampling
 791 **Input** : Task distribution $p(\tau)$; Task batch size \mathcal{B} ; Candi-
 792 date batch size $\hat{\mathcal{B}}$; Latest updated $\{\psi, \phi\}$; Latest
 793 history H_{t-1} ; Iteration number K ; Learning rate
 794 λ_2 .
 795 **Output**: Selected task identifier batch $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$.
 796 // **Optimize Risk Predictive Module**
 797 **for** $i = 1$ **to** K **do**
 798 Perform gradient updates given H_{t-1} :
 799 $\phi \leftarrow \phi + \lambda_2 \nabla_\phi \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (5);
 800 $\psi \leftarrow \psi + \lambda_2 \nabla_\psi \mathcal{G}_{\text{ELBO}}(\psi, \phi)$ in Eq. (5);
 801 **end**
 802 // **Simulating After-Adaptation Results**
 803 Randomly sample $\{\hat{\tau}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ from $p(\tau)$;
 804 // **Posterior Sampling Outcome**
 805 Amortized evaluation $\{\hat{\ell}_{t,i}\}_{i=1}^{\hat{\mathcal{B}}}$ with one stochastic forward
 806 pass for all candidate tasks through executing Eq. (12)a-b;
 807 // **Diversity-Guided Subset Search**
 808 Run approximate algorithms to solve Eq. (12)c;
 809 Return the screened subset $\{\tau_{t,i}\}_{i=1}^{\mathcal{B}}$ for next iteration.

770 A.5. Details in Risk Predictive Modules

771 As introduced in Section 2.2, a key feature of RATS is its use of risk predictive models, e.g., $p(\ell|\tau, H_{1:t}; \theta_t)$, as surrogates
 772 for expensive evaluations. To the best of our knowledge, MPTS (Wang et al., 2025) is the first to develop such a model
 773 for robust optimization. As introduced in Section 2.2, MPTS leverages generative modeling for adaptation optimization
 774 and employs approximate posterior inference to forecast adaptation risk values after one-step optimization. Since small
 775 deviations in the machine learner’s parameters do not alter the relative difficulty scores within the task batch, this provides a
 776 tractable approach to coarse-grained task difficulty evaluation. The probabilistic graphical model is provided in Fig. 6.
 777

778 Empirically, a series of evaluation on DR and Meta-RL in (Wang et al., 2025) has validated the plausibility of such a schema.
 779 Particularly, it achieves high Pearson correlation coefficient (PCC) values between the risk learner’s predicted adaptation
 780 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 the corresponding exact adaptation risk values $\{\ell_{t+1,i}\}_{i=1}^B$,
 781 indicating the risk predictive model’s ability to score difficulty. The Pearson correlation coefficient value is computed as
 782 $\rho_{\bar{\ell}, \ell} := \frac{\sum_{i=1}^B (\bar{\ell}_{t+1,i} - \text{Mean}[\{\bar{\ell}_{t+1,i}\}])(\ell_{t+1,i} - \text{Mean}[\{\ell_{t+1,i}\}])}{\sqrt{\sum_{i=1}^B (\bar{\ell}_{t+1,i} - \text{Mean}[\{\bar{\ell}_{t+1,i}\}])^2} \sqrt{\sum_{i=1}^B (\ell_{t+1,i} - \text{Mean}[\{\ell_{t+1,i}\}])^2}}$.
 783



797 Figure 6. The probabilistic graphical model of the risk predictive model in (Wang et al., 2025), where gray units denote observed variables
 798 with the white as unobservable ones. The solid directed lines depict the generative model, and the dashed directed lines indicate the
 799 recognition model and approximate inference (Kingma & Welling, 2013).
 800

801 **Formulation of ELBO & Stochastic Gradient Estimates** For the sake of easy and fast understanding of PDTS for
 802 the reviewer, we include the risk predictive model part as follows. These are modified from the MPTS team’s open-
 803 sourced manuscript in (Wang et al., 2025). Here, the optimization objective of the risk predictive model is in Eq. (5). The
 804 risk predictive model uses a latent variable to summarize historical information and quantify uncertainty in predicting
 805 task-specific adaptation risk. The following outlines the steps for deriving the evidence lower bound in optimization.
 806

$$\mathcal{L}_{\text{ML}}(\psi) := \ln p_\psi(H_t|H_{1:t-1}) = \ln \left[\int p_\psi(H_t|z_t)p(z_t|H_{1:t-1}) dz \right] \quad (19a)$$

$$= \ln \left[\int q_\phi(z_t|H_t) \frac{p(z_t|H_{1:t-1})}{q_\phi(z_t|H_t)} p_\psi(H_t|z_t) dz_t \right] \quad (19b)$$

$$\geq \mathbb{E}_{q_\phi(z_t|H_t)} \left[\ln p_\psi(H_t|z_t) \right] - D_{KL} \left[q_\phi(z_t|H_t) \parallel p(z_t|H_{1:t-1}) \right] := \mathcal{G}_{\text{ELBO}}(\psi, \phi) \quad (19c)$$

811 Then, the ELBO is derived with the help of the reparameterization trick (Kingma & Welling, 2013). And Wang et al. (2025)
 812 further relax the ELBO to derive the β -VAE version, which corresponds to Eq. (5) in this work.
 813

814 **Neural Architecture of the Risk Predictive Model.** For practicality, stability and fair comparison, we adopt the same
 815 risk predictive model design proposed in MPTS (Wang et al., 2025). As shown in Fig. 7, the risk learner follows an
 816 encoder-decoder structure. For consistency across benchmarks, we employ the same neural architecture for all experiments.
 817 The encoder consists of an embedding network with four hidden layers, each containing 10 units and using Rectified
 818 Linear Unit (ReLU) activations. It encodes the batch $\{[\tau_{t,i}, \ell_{t,i}]\}_{i=1}^B$ into r through mean pooling, subsequently mapping
 819 it to $[\mu, \sigma]$. The latent variable z is sampled from a normal distribution defined by $[\mu, \sigma]$. The decoder is a three-layer
 820 neural network with nonlinear activation functions, mapping $\{[\tau_{t,i}, z]\}_{i=1}^B$ to the predicted risk $\{\hat{\ell}_{t,i}\}_{i=1}^B$. The optimization
 821 objective is given in Eq. (5). For further details on the implementation, please refer to our code repository.
 822
 823

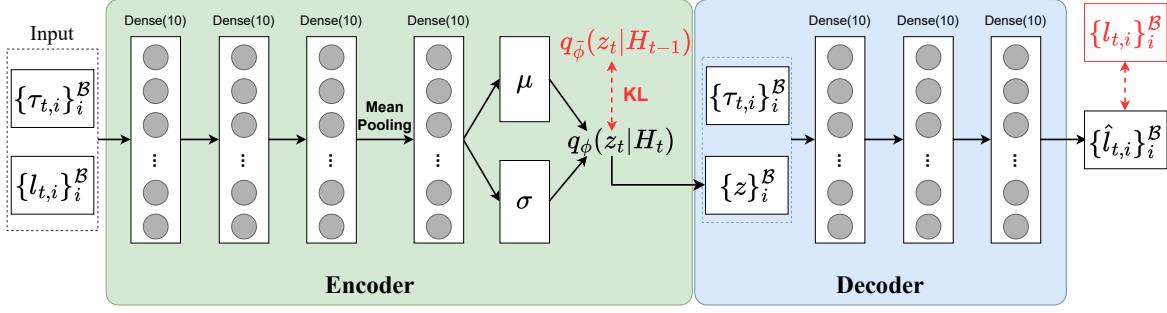


Figure 7. Illustration of the neural architecture of the risk predictive model in (Wang et al., 2025). The risk predictive model follows an encoder-decoder structure which encodes the batch $[\tau_{t,i}, \ell_{t,i}]$ into a latent variable z and then decodes it into predicted adaptation risks $\hat{\ell}_{t,i}$. We reuse it in PDTS as a standard risk predictive model backbone.

B. Theoretical Analysis and Proofs

B.1. Preliminaries of MPTS

In MPTS (Wang et al., 2025), the streaming adaptation outcome in task episodic learning is characterized as:

$$\theta_0 \xrightarrow{\text{eval}} \dots \xrightarrow{\text{eval}} \{(\tau_{t-1,i}, \ell_{t-1,i})\}_{i=1}^B \xrightarrow{\text{opt}} \theta_t \xrightarrow{\text{eval}} \{(\tau_{t,i}, \ell_{t,i})\}_{i=1}^B \xrightarrow{\text{opt}} \dots \xrightarrow{\text{opt}} \theta_T. \quad (20)$$

The most encouraging finding is that these cumulated task identifiers and adaptation risk values can be coupled to learn helpful adaptation prior and serve the evaluation of the task difficulty in a rough granularity. The resulting dataset is used to train the risk predictive model associated with the modified optimization objective as Eq. (5).

B.2. One Secret MDP Steers Adaptation to Many MDPs

The Operator for the State Transition in the Secret MDP. Here, given a smooth and fixed machine learning optimizer, we can define the operator in zero-shot or few-shot adaptation optimization as:

$$\text{State Transition Operation after Adaptation} \quad \mathcal{F} : \theta \times \mathcal{T}^B \mapsto \theta', \quad (21)$$

this corresponds to the deterministic transition function in the secret MDP \mathcal{M} . Here, take the DR case as an example, the above operator \mathcal{F} results in the following probably transited states:

$$\begin{aligned} \theta_{t+1}^* &= \arg \min_{\theta \in \Theta_{t+1}} \mathbb{E}_{p_\alpha(\tau; \theta)} [\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)], \\ \Theta_{t+1} &= \left\{ \theta_t - \eta \frac{1}{B} \nabla_\theta \sum_{b=1}^B \ell(\mathcal{D}_{\tau_b}^Q, \mathcal{D}_{\tau_b}^S; \theta) \mid \tau_b \in \mathcal{T}_{t+1}^B, |\mathcal{T}_{t+1}^B| = B, \mathcal{T}_{t+1}^B \subseteq \mathcal{T}_{t+1}^B \right\}, \end{aligned} \quad (22)$$

where $p_\alpha(\tau; \theta)$ corresponds to the probability density function of the $(1 - \alpha)$ -tailed tasks.

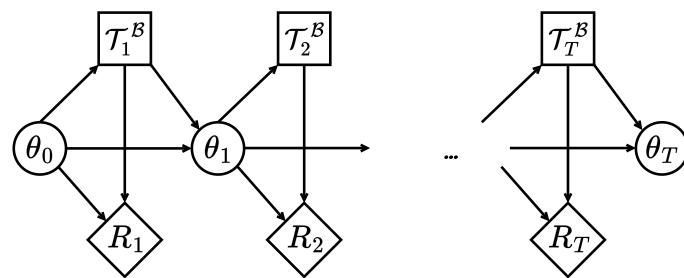


Figure 8. Active Task Episodic Learning as a Secret MDP.

Probabilistic Graphical Model of the Secret MDP. Here, we can write the probabilistic form of the constructed MDP as:

$$p(\mathcal{T}_{1:T}^{\mathcal{B}}, R_{1:T}, \boldsymbol{\theta}_{0:T} | \Pi_{1:T}, H_{0:T-1}) = \underbrace{p(\boldsymbol{\theta}_0)}_{\text{Initial State}} \prod_{t=1}^T \underbrace{p(R_t | \boldsymbol{\theta}_{t-1}, \mathcal{T}_t^{\mathcal{B}})}_{\text{Step-Wise Reward}} \prod_{t=0}^{T-1} \underbrace{\pi_t(\mathcal{T}_{t+1}^{\mathcal{B}} | H_{0:t})}_{\text{Task Sampler}} \prod_{t=0}^{T-1} \underbrace{p(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}})}_{\text{State Transition}}. \quad (23)$$

where $\{\mathcal{T}_{0:T}^{\mathcal{B}}, r_{0:T}, \boldsymbol{\theta}_{0:T}\}$ records the trajectory information during policy search. The corresponding probabilistic graphical model is in Fig. 23.

B.3. Proof of Proposition 3.2

Proposition 3.2 (MPTS as a UCB-guided Solution to i-MABs) *Executing MPTS pipeline in Eq. (6) is equivalent to approximately solving \mathcal{M} with the i-MAB under the UCB principle.*

To demonstrate the claim in **Proposition 3.2**, we revisit fundamental concepts in CVaR $_{\alpha}$ optimization and break the problem into several steps. These include # **Step1** the dual form of CVaR $_{\alpha}$ with its Monte Carlo estimates, **Lemma B.1** to determine the optimal step-wise action, and # **Step2** nearly CVaR $_{\alpha}$ approximation/UCB-Guided Solution to i-MABs in MPTS.

Step1. The Greedy Action as the Subset with Top- \mathcal{B} Adaptation Risk Values.

Unbiased Monte Carlo Estimate of CVaR $_{\alpha}$. Here, we can rewrite the optimization objective of CVaR $_{\alpha}$ in the dual form:

$$\min_{\boldsymbol{\theta} \in \Theta, \zeta \in \mathbb{R}} \text{CVaR}_{\alpha}(\boldsymbol{\theta}) := \zeta + \frac{1}{1-\alpha} \mathbb{E}_{p(\tau)} [\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}) - \zeta]^+, \quad (24)$$

where the conditional function means $[\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}) - \zeta]^+ = \max\{\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta}) - \zeta, 0\}$. And its Monte Carlo sample average approximation corresponds to

$$\min_{\boldsymbol{\theta} \in \Theta, \zeta \in \mathbb{R}} \text{CVaR}_{\alpha}(\boldsymbol{\theta}) := \zeta + \frac{1}{(1-\alpha)\mathcal{B}} \sum_{i=1}^{\mathcal{B}} [\ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \boldsymbol{\theta}) - \zeta]^+. \quad (25)$$

As the optimal auxiliary variable satisfies $\zeta = \text{VaR}_{\alpha}(\boldsymbol{\theta})$, the Top- \mathcal{B} element in the set $\{\ell_i | \ell_i = \ell(\mathcal{D}_{\hat{\tau}_i}^Q, \mathcal{D}_{\hat{\tau}_i}^S; \boldsymbol{\theta})\}_{i=1}^{\hat{\mathcal{B}}}$ can be viewed as the unbiased Monte Carlo estimate w.r.t. Eq. (24) and the equivalent form of Eq. (25). In implementation, (Wang et al., 2024; Lv et al., 2024) typically evaluate the machine learner's adaptation performance on candidate tasks, rank their values, and filter out $(1-\alpha)\hat{\mathcal{B}}$ to optimize.

Lemma B.1 (Subset with Top- \mathcal{B} Risk Values as the Optimal Arm). *Given the secret MDP \mathcal{M} in Section 3.1 and its step-wise reward $R(\boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}}) := \text{CVaR}_{\alpha}(\boldsymbol{\theta}_t) - \text{CVaR}_{\alpha}(\boldsymbol{\theta}_{t+1})$, selecting the subset with average Top- \mathcal{B} risk values inherently maximizes $R(\boldsymbol{\theta}_t, \mathcal{T}_{t+1}^{\mathcal{B}})$.*

Proof. Given the state $\boldsymbol{\theta}_t$ and the available action set \mathbf{A}_t in \mathcal{M} , we need to show that selecting the $\mathcal{T}_{t+1}^{\mathcal{B}*} \in \mathbf{A}_t$ brings largest CVaR decrease as $\text{CVaR}_{\alpha}(\boldsymbol{\theta}_t) - \text{CVaR}_{\alpha}(\boldsymbol{\theta}_{t+1}^*)$.

Let $\mathcal{T}_{t+1}^{\mathcal{B}*}$ denote the subset with highest average Top- \mathcal{B} risk values in the set $\{\ell_i | \ell_i = \ell(\mathcal{D}_{\hat{\tau}_i}^Q, \mathcal{D}_{\hat{\tau}_i}^S; \boldsymbol{\theta})\}_{i=1}^{\hat{\mathcal{B}}}$, which can be viewed as the Monte Carlo sample, i.e., unbiased estimate of $p_{\alpha}(\tau; \boldsymbol{\theta}_t)$. This is also in accordance with the dual form of CVaR $_{\alpha}$ in Eq. (25). Meanwhile, we can define its transited state as $\boldsymbol{\theta}_{t+1}^*$. In a similar way, let some arbitrary feasible action, i.e., the subset be $\mathcal{T}_{t+1}^{\mathcal{B}'}$, we can define the corresponding transited state as $\boldsymbol{\theta}_{t+1}^*$. Inspired by the notation in Eq. (A.2), we can re-express the subset $\mathcal{T}_{t+1}^{\mathcal{B}'}$ as the Monte Carlo sample from another task distribution $q(\tau)$, which differs from $p_{\alpha}(\tau; \boldsymbol{\theta}_t)$.

Next, we can estimate the robustness improvement under different actions from the one-step Taylor expansion trick.

$$\text{Unbiased Objective } \mathcal{L}(\boldsymbol{\theta}) := \mathbb{E}_{p_{\alpha}(\tau; \boldsymbol{\theta}_t)} [\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta})] \quad \text{Biased Objective } \hat{\mathcal{L}}(\boldsymbol{\theta}) := \mathbb{E}_{q(\tau)} [\ell(\mathcal{D}_{\tau}^Q, \mathcal{D}_{\tau}^S; \boldsymbol{\theta})] \quad (26a)$$

$$\boldsymbol{\theta}'_{t+1} = \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \quad \hat{\boldsymbol{\theta}}'_{t+1} = \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}} \hat{\mathcal{L}}(\boldsymbol{\theta}) \quad (26b)$$

$$\mathcal{L}(\boldsymbol{\theta}'_{t+1}) = \mathcal{L}(\boldsymbol{\theta}_t) - \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})^T \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) + \mathcal{O}(\|\boldsymbol{\theta}'_{t+1} - \boldsymbol{\theta}_t\|_2^2) \Rightarrow \mathcal{L}(\boldsymbol{\theta}_t) - \mathcal{L}(\boldsymbol{\theta}'_{t+1}) \approx \eta \|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})\|_2^2 \quad (26c)$$

$$\mathcal{L}(\boldsymbol{\theta}_t) - \mathcal{L}(\hat{\boldsymbol{\theta}}'_{t+1}) \approx \eta \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})^T \nabla_{\boldsymbol{\theta}} \hat{\mathcal{L}}(\boldsymbol{\theta}) = \eta \|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})\|_2^2 \cos \alpha_q \leq \eta \|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})\|_2^2 = \mathcal{L}(\boldsymbol{\theta}_t) - \mathcal{L}(\boldsymbol{\theta}'_{t+1}), \quad (26d)$$

where we typically assume the norms of stochastic gradients for $\nabla_{\theta}\mathcal{L}(\theta)^T$ and $\nabla_{\theta}\hat{\mathcal{L}}(\theta)$ are the same and their angle is α_q . As a result, we can see the optimal stochastic gradient should be the unbiased estimate of the Eq. (24), which corresponds to the Top- \mathcal{B} tasks in the pseudo batch $\mathcal{T}_{t+1}^{\hat{\mathcal{B}}}$ with $\frac{\mathcal{B}}{\hat{\mathcal{B}}} = 1 - \alpha$. This completes the proof of Lemma B.1. \square

Meanwhile, remember that in the secret MDP \mathcal{M} , the agent will never revisit the previous state θ_t due to the nature of the stochastic gradient descent in the operator \mathcal{F} in Eq. (21). Finally, we can claim that maximizing the state action value $Q(\theta_t, \mathcal{T}_{t+1}^{\mathcal{B}})$ in the i-MAB actually corresponds to maximizing the step-wise reward due to the Bellman optimality in the main paper, and picking up the worst subset secretly maximizes $R(\theta_t, \mathcal{T}_{t+1}^{\mathcal{B}})$ in a greedy way. This implies that $\pi_t^* = \arg \max_{\mathcal{T}_{t+1}^{\mathcal{B}} \subseteq \mathcal{T}_{t+1}^{\hat{\mathcal{B}}}} Q(\theta_t, \mathcal{T}_{t+1}^{\mathcal{B}})$ and $Q(\theta_t, \mathcal{T}_{t+1}^{\mathcal{B}}) \propto \frac{1}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \ell_{t+1,i}$, where $\{\ell_{t+1,i}\}_{i=1}^{\mathcal{B}}$ denotes the evaluated adaptation performance of a feasible subset $\mathcal{T}_{t+1}^{\mathcal{B}}$ conditioned on θ_t . In the presence of RATS, the adaptation performance is evaluated by the risk predictive model $p(\ell|\tau, H_{1:t}; \theta_t)$ in an amortized way, which suggests $Q(\theta_t, \mathcal{T}_{t+1}^{\mathcal{B}}) \propto \frac{1}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \hat{\ell}_{t+1,i}$ with $\hat{\ell}$ the predicted value.

Step2. UCB-Guided Solution to i-MABs in MPTS.

Assumption 3 (Randomized Adaptation Risk Value Function). *Given the secret MDP \mathcal{M} in Section 3, we assume the distribution of the adaptation risk value follows an implicit Gaussian distribution, i.e., $p(\ell_{t+1,i}|\tau_i, H_{1:t}; \theta_t) = \mathcal{N}(\mu_{t+1,i}, \sigma_{t+1,i}^2)$.*

Note that tasks with their identifiers are sampled in an *i.i.d.* way; this induces the conditional independence and the distribution of their summation as Eq. (27) with the help of Assumption 3.

$$\begin{aligned} p(\mathcal{L}_{t+1}^{\mathcal{B}} | \mathcal{T}_{t+1}^{\mathcal{B}}, H_{1:t}; \theta_t) &= \prod_{i=1}^{\mathcal{B}} p(\ell_{t+1,i} | \tau_i, H_{1:t}; \theta_t) = \prod_{i=1}^{\mathcal{B}} \mathcal{N}(\mu_{t+1,i}, \sigma_{t+1,i}^2) \\ &\implies p\left(\sum_{i=1}^{\mathcal{B}} \ell_{t+1,i} | \mathcal{T}_{t+1}^{\mathcal{B}}; \theta_t\right) = \mathcal{N}\left(\sum_{i=1}^{\mathcal{B}} \mu_{t+1,i}, \sum_{i=1}^{\mathcal{B}} \sigma_{t+1,i}^2\right) := \mathcal{N}(\mu_{t+1}^{\mathcal{B}}, \sigma_{t+1}^{\mathcal{B}^2}) \end{aligned} \quad (27)$$

At the same time, we find in MPTS, the multiple stochastic forward passes in Eq. (7) are performed to obtain the MC estimated distribution parameters $\{m(\ell_i) := \hat{\mu}_{t+1,i}, \sigma_i(\ell_i) := \hat{\sigma}_{t+1,i}\}$ for each task identifier τ_i in the batch $\mathcal{T}_{t+1}^{\mathcal{B}}$. This can be further associated with the factorization in the risk predictive model as Eq. (28):

$$p(\mathcal{L}_{t+1}^{\mathcal{B}} | \mathcal{T}_{t+1}^{\mathcal{B}}, H_{1:t}; \theta_t) \approx \int p_{\psi}(\mathcal{L}_{t+1}^{\mathcal{B}} | \mathcal{T}_{t+1}^{\mathcal{B}}, z_t) q_{\phi}(z_t | H_{1:t}) dz_t \quad (28a)$$

$$= \int \prod_{i=1}^{\mathcal{B}} p(\ell_{t+1,i} | \tau_{t+1,i}, z_t) q(z_t | H_{1:t}) dz_t. \quad (28b)$$

The last step shows that the worst subset corresponds to the unbiased Monte Carlo estimate of CVaR $_{\alpha}$, and its sample average adaptation risk value can be treated as the proxy of the reward. Hence, picking up the subset with each element in the Top- \mathcal{B} risk values is doing exploitation in robust fast adaptation. Rethinking MPTS's acquisition function in Eq. (7), we can find the following inequality:

$$\sqrt{\sum_{i=1}^{\mathcal{B}} \sigma_{t+1,i}^2} \leq \sum_{i=1}^{\mathcal{B}} \sigma_{t+1,i} \quad \text{with } \forall \sigma_{t+1,i} \in \mathbb{R}^+ \quad (29a)$$

$$\implies \gamma_1 \sqrt{\sum_{i=1}^{\mathcal{B}} \sigma_{t+1,i}^2 + \sum_{i=1}^{\mathcal{B}} \gamma_0 \mu_{t+1,i}} \leq \sum_{i=1}^{\mathcal{B}} \gamma_1 \sigma_{t+1,i} + \gamma_0 \mu_{t+1,i} \quad (29b)$$

$$\implies \underbrace{\gamma_1 \sqrt{\sum_{i=1}^{\mathcal{B}} \sigma(\ell_i)^2} + \gamma_0 \sum_{i=1}^{\mathcal{B}} m(\ell_i)}_{\text{UCB}} \leq \underbrace{\sum_{i=1}^{\mathcal{B}} \gamma_1 \sigma(\ell_i) + \gamma_0 m(\ell_i)}_{\text{Approximate UCB}} := \mathcal{A}_{\text{U}}(\mathcal{T}^{\mathcal{B}}), \quad (29c)$$

which means MPTS actually executes the approximate UCB to balance the exploitation (picking up the subset with the estimated worst performance) and the exploration of the task space (picking up the arm with the nearly highest epistemic uncertainty captured by the risk predictive model).

990 The above two steps complete the proof of **Proposition 3.2**.

992 B.4. Proof of Proposition 3.3

993 **Proposition 3.3** (Concentration Issue in Average Top- \mathcal{B} Selection) *Let $f(\tau) : \mathbb{R}^d \rightarrow \mathbb{R}$ be a unimodal and continuous*
 994 *function, where $d \in \mathbb{N}^+$ and $\tau \in \mathbb{R}^d$, with a maximum value $f(\tau^*)$ at τ^* . We uniformly sample a set of points $\mathcal{T}^{\hat{\mathcal{B}}} = \{\tau_i\}_{i=1}^{\hat{\mathcal{B}}}$,*
 995 *where τ_i are i.i.d. with a probability p_ϵ of falling within a ϵ -neighborhood of τ^* as $|f(\tau) - f(\tau^*)| \leq \epsilon$. Following MPTS,*
 996 *we select the Top- \mathcal{B} samples with the largest function values, i.e.,*

$$997 \quad \mathcal{T}^{\mathcal{B}} = \text{Top-}\mathcal{B}(\mathcal{T}^{\hat{\mathcal{B}}}, f), \quad \hat{\mathcal{B}}, \mathcal{B} \in \mathbb{N}^+, \quad \mathcal{B} \leq \hat{\mathcal{B}},$$

1000 For any $\epsilon > 0$ such that $p_\epsilon < \frac{\hat{\mathcal{B}}-\mathcal{B}+2}{\hat{\mathcal{B}}+1}$, the concentration probability

$$1002 \quad \mathbb{P}(|f(\tau) - f(\tau^*)| \leq \epsilon \mid \forall \tau \in \mathcal{T}^{\mathcal{B}})$$

1003 increases with $\hat{\mathcal{B}}$ and converges to 1 with $\hat{\mathcal{B}} \rightarrow \infty$.

1004 *Proof.* We define p_ϵ as the probability that a random variable τ uniformly sampled from the domain of definition falls
 1005 within the neighborhood of the maximum value τ^* , i.e.,

$$1006 \quad p_\epsilon = \mathbb{P}(|f(\tau) - f(\tau^*)| \leq \epsilon \mid \tau \sim \text{Unif}(\cdot)).$$

1007 Next, consider the probability that at least \mathcal{B} random variables from the set $\mathcal{T}^{\hat{\mathcal{B}}} = \{\tau_i\}_{i=1}^{\hat{\mathcal{B}}}$, where $\hat{\mathcal{B}}, \mathcal{B} \in \mathbb{N}^+$ and $\mathcal{B} \leq \hat{\mathcal{B}}$,
 1008 are within the neighborhood of τ^* . Since the τ_i 's are i.i.d. and $\tau_i \sim \text{Unif}(\cdot)$, the probability is given by

$$1009 \quad P^{\hat{\mathcal{B}}, \mathcal{B}} = 1 - \left[\sum_{i=1}^{\mathcal{B}} p_\epsilon^{\hat{\mathcal{B}}-i+1} (1-p_\epsilon)^{i-1} \binom{\hat{\mathcal{B}}}{i-1} \right].$$

1010 Since f is a unimodal and continuous function, we can directly relate the concentration probability as

$$1011 \quad \mathbb{P}(|f(\tau) - f(\tau^*)| \leq \epsilon \mid \forall \tau \in \mathcal{T}^{\mathcal{B}}) = P^{\hat{\mathcal{B}}, \mathcal{B}}.$$

1012 To establish the monotonicity of $P^{\hat{\mathcal{B}}, \mathcal{B}}$ with respect to $\hat{\mathcal{B}}$, observe that the term $p_\epsilon^{\hat{\mathcal{B}}-i+1} (1-p_\epsilon)^{i-1} \binom{\hat{\mathcal{B}}}{i-1}$ is monotonically
 1013 decreasing in $\hat{\mathcal{B}}$ for fixed i , given that $p_\epsilon < \frac{\hat{\mathcal{B}}-i+2}{\hat{\mathcal{B}}+1}$. To see this, we compute the ratio of consecutive terms:

$$1014 \quad \frac{p_\epsilon^{\hat{\mathcal{B}}-i+1} (1-p_\epsilon)^{i-1} \binom{\hat{\mathcal{B}}}{i-1}}{p_\epsilon^{\hat{\mathcal{B}}-i+2} (1-p_\epsilon)^{i-1} \binom{\hat{\mathcal{B}}+1}{i-1}} = \frac{\hat{\mathcal{B}}-i+2}{p_\epsilon (\hat{\mathcal{B}}+1)} > 1.$$

1015 Thus, the sequence $p_\epsilon^{\hat{\mathcal{B}}-i+1} (1-p_\epsilon)^{i-1} \binom{\hat{\mathcal{B}}}{i-1}$ decreases with $\hat{\mathcal{B}}$, and consequently, $P^{\hat{\mathcal{B}}, \mathcal{B}}$ increases monotonically in $\hat{\mathcal{B}}$ when
 1016 \mathcal{B} is fixed, provided that $p_\epsilon < \frac{\hat{\mathcal{B}}-\mathcal{B}+2}{n+1}$.

1017 Therefore, we conclude that for any $\epsilon > 0$ such that $p_\epsilon < \frac{\hat{\mathcal{B}}-\mathcal{B}+2}{\hat{\mathcal{B}}+1}$, the probability $\mathbb{P}(|f(\tau) - f(\tau^*)| \leq \epsilon \mid \forall \tau \in \mathcal{T}^{\mathcal{B}})$
 1018 increases monotonically with $\hat{\mathcal{B}}$ for fixed \mathcal{B} . □

1019
1020

1021 B.5. Proof of Proposition 3.4

1022 **Proposition 3.4** (Nearly Worst-Case Optimization with PDTS) *When $\hat{\mathcal{B}}$ grows large enough, optimizing the subset from Eq.*
 1023 *(11) achieves nearly worst-case optimization.*

1024 *Proof.* As the exact size of the subset \mathcal{B} is fixed in the optimization, the ratio $\frac{\mathcal{B}}{\hat{\mathcal{B}}}$ goes to nearly 0 with the increase of $\hat{\mathcal{B}}$ to a
 1025 certain scale. As the consequence, the number of arms grows to $C_{\hat{\mathcal{B}}}^{\mathcal{B}}$ and the robustness concept is $\text{CVaR}_{1-\frac{\mathcal{B}}{\hat{\mathcal{B}}}}$. Since the
 1026 involvement of the diversity regularization perturbs the worst arm selection, this induces the nearly worst-case optimization
 1027 in PDTS. □

1045 C. Experimental Setups & Implementation Details

1046 C.1. Risk-Averse Baseline Details

1048 These baselines are SOTA methods published in NeurIPS/ICLR conferences and the latest open-sourced version (Wang
 1049 et al., 2024; Lv et al., 2024; Sagawa et al., 2019; Wang et al., 2025; Greenberg et al., 2024).

1051 **GDRM (Sagawa et al., 2019; Setlur et al., 2024).** As briefly introduced in the main paper, GDRM can be viewed as a
 1052 min-max optimization problem. The core concept of enhancing the machine learner's robustness involves reallocating more
 1053 probability mass to the worst-case scenarios in a weighted manner. In each iteration with the optimal $p_{\hat{g}}(\tau)$, the optimization
 1054 problem simplifies to:

$$1056 \min_{\theta \in \Theta} \mathbb{E}_{p_{\hat{g}}(\tau)} [\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)] = \mathbb{E}_{p(\tau)} \left[\frac{p_{\hat{g}}(\tau)}{p(\tau)} \ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta) \right], \quad (30)$$

1059 where we use $\omega(\tau) = \frac{p_{\hat{g}}(\tau)}{p(\tau)}$ to denote the weight.

1060 In general, for a fixed number of tasks, GDRM organizes tasks heuristically or dynamically into clusters, followed by a
 1061 reweighting mechanism based on assessed risks. However, in task-episodic learning, no task grouping is performed because
 1062 the task batch is reset after each iteration (Wang et al., 2024; Lv et al., 2024). Task-specific weights are calculated as
 1063 $\omega(\tau_i) = \frac{\exp(\eta \ell(\mathcal{D}_{\tau_i}^Q, \mathcal{D}_{\tau_i}^S; \theta))}{\sum_{b=1}^B \exp(\eta \ell(\mathcal{D}_{\tau_b}^Q, \mathcal{D}_{\tau_b}^S; \theta))}$, where η denotes the temperature parameter, and $\{\tau_b\}_{b=1}^B$ represents the identifiers of the
 1064 task batch. Additional implementation details are available at https://github.com/kohpangwei/group_DRO.

1065 **DRM (Wang et al., 2024; Lv et al., 2024).** As outlined in the main paper, we adopt the standard approach in DRM, where
 1066 the optimization objective is expressed as $\mathbb{E}_{p_\alpha(\tau; \theta)} [\ell(\mathcal{D}_\tau^Q, \mathcal{D}_\tau^S; \theta)]$. A widely used practical strategy involves evaluating the
 1067 performance of the machine learner, ranking task-specific adaptation risks, and optimizing over the worst ($1 - \alpha$) proportion
 1068 of tasks.

1069 Following the setup in Wang et al. (2024), we select the Top- B tasks during optimization, which corresponds to evaluating a
 1070 task batch of size $\frac{B}{1-\alpha}$. To ensure a fair comparison with PDTs while preserving computational efficiency, we employ the
 1071 same Monte Carlo estimator for the risk quantile as in Wang et al. (2024). For all benchmarks, we set the actual task batch
 1072 size to $\hat{B} = 2B$, discarding the easiest half before optimizing the machine learner.

1073 **MPTS (Wang et al., 2025).** We have thoroughly introduced the MPTS pipeline in Section 2.2 and provided details
 1074 of the risk learner in Section A.5. For additional configurations, we adopt those from the official repository at <https://github.com/thu-rllab/MPTS>, including the heuristic random mixture strategy and the specific values of \hat{B} .

1082 **Table 3. Details of Task Identifiers and Algorithm Backbones Across Benchmarks.** Here, we provide detailed information about the
 1083 task identifiers used to induce task distributions and the algorithm backbones, including MAML (Finn et al., 2017), TD3 (Fujimoto et al.,
 1084 2018), and PPO (Schulman et al., 2017), employed in various benchmarks.

Benchmarks	Identifier Meaning	Identifier Range	Backbone
K-shot sinusoid regression	amplitude and phase (a, b)	$[0.1, 5.0] \times [0, \pi]$	MAML
Meta-RL: ReacherPos	goal location (x_1, x_2)	$[-0.2, 0.2] \times [-0.2, 0.2]$	
Meta-RL: Walker2dVel	velocity v	$[0, 2.0]$	
Meta-RL: Walker2dMassVel	mass and velocity (m, v)	$[0.75, 1.25] \times [0, 2.0]$	MAML
Meta-RL: HalfCheetahMassVel	mass and velocity (m, v)	$[0.75, 1.25] \times [0, 2.0]$	
DR: Pusher	puck friction loss f and puck joint damping d	$[0.004, 0.01] \times [0.01, 0.025]$	
DR: LunarLander	main engine strength s	$[4, 20]$	TD3
DR: ErgoReacher	joint damping d and max torque t ($\times 4$ joints)	$[0.1, 2.0] \times [2, 20]$	
VisualDR: LiftPegUpright.Light	ambient light l (x3 dimensions)	$[-1.0, 2.0]$	
VisualDR: AnymalCReach_Goal	goal location (x_1, x_2)	$[0.0, 1.0] \times [0.0, 1.0]$	PPO

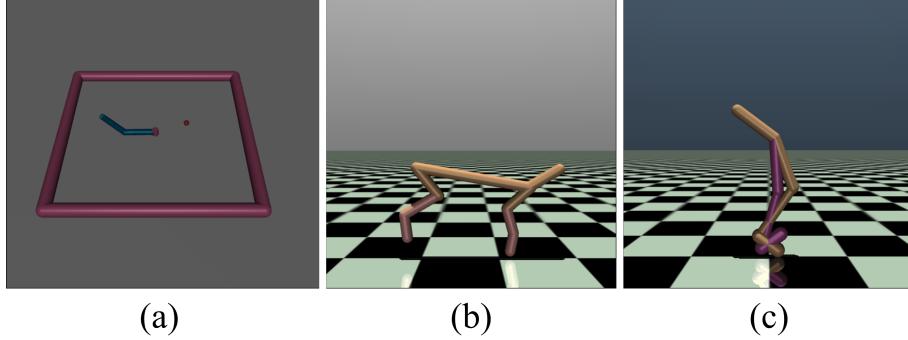


Figure 9. Illustrations of three types of robots in Meta-RL based on the Mujoco. (a) Reacher, (b) HalfCheetah, and (c) Walker2d.

C.2. Sinusoid Regression

Following standard setups in prior works (Finn et al., 2017; Wang et al., 2025), we define the sinusoid regression problem as a toy example of supervised meta-learning. The goal of this problem is to predict the wave function $y = a \sin(x - b)$, where the amplitude a and phase b are sampled as task-specific parameters, i.e., task identifier $\tau = (a, b)$, using a few-shot support dataset. We adopt the 10-shot sinusoid regression setting, where 10 data points are uniformly sampled from the interval $[-5.0, 5.0]$ to form the support dataset for each task.

Implementation Details. The machine learner is a neural network with 2 hidden layers, each of size 40, using the Rectified Linear Unit (ReLU) as the nonlinear activation function. The task batch size is set to 16 for ERM and GDRM, while a batch size of 32 is used as the default for DRM. The temperature parameter in GDRM is $\eta = 0.001$, and the learning rates for both the inner and outer loops are fixed at 0.001. The task identifier has a dimensionality of 2. For MPTS and PDTs, the identifier batch size during training is set to 32 ($2\times$) and 512 ($32\times$), respectively. We use the Adam optimizer with a learning rate of 5×10^{-4} to update the risk learner over 15,000 steps. The label for the risk learner is the average MSE loss value for each task. For sinusoid regression and meta-reinforcement learning, we use the standard repository provided by MAML (Finn et al., 2017). Regarding validation during training, we use a separate uniform task sampler with a fixed random seed to select 1,000 tasks for validating the training checkpoints of various methods.

C.3. Meta Reinforcement Learning

We adopt four Meta-RL scenarios—ReacherPos, Walker2dVel, Walker2dMassVel, and HalfCheetahMassVel—from MPTS (Wang et al., 2025), which represent distinct MDP distributions based on the Mujoco physics engine (Todorov et al., 2012). These scenarios involve three types of robots (HalfCheetah, Walker2d, and Reacher) and three randomized meta-learning objectives: velocity, goal location, and body mass.

- The objective in velocity-based scenarios, e.g., Walker2dVel, is to train the robot to achieve a target velocity, with the reward function defined as the negative absolute difference between the robot’s current velocity and the target velocity. This reward is augmented by a control penalty and an alive bonus to facilitate learning. As the task distribution is specified by a uniform distribution over the target velocity, $\tau = v$ can be viewed as the task identifier.
- The body mass and velocity scenarios, e.g., Walker2dMassVel and HalfCheetahMassVel, share the same objective as the velocity scenarios but additionally feature varying robot masses. As the task distribution is specified by a uniform distribution over the body mass and target velocity, $\tau = (m, v)$ can be viewed as the task identifier.
- The goal location scenario, e.g., ReacherPos, requires moving a two-jointed robot arm’s end effector close to a target position. Its reward function is defined as the negative L_1 distance between the end effector’s position and the target, supplemented by a control cost to encourage robustness. As the task distribution is specified by a uniform distribution over the goal location, $\tau = (x_1, x_2)$ can be viewed as the task identifier.

Implementation Details. The machine learner is implemented as a neural network with 2 hidden layers, each consisting of 64 units, and utilizes ReLU activations for nonlinearity. For ERM and GDRM, the default task batch size is set to 20, while

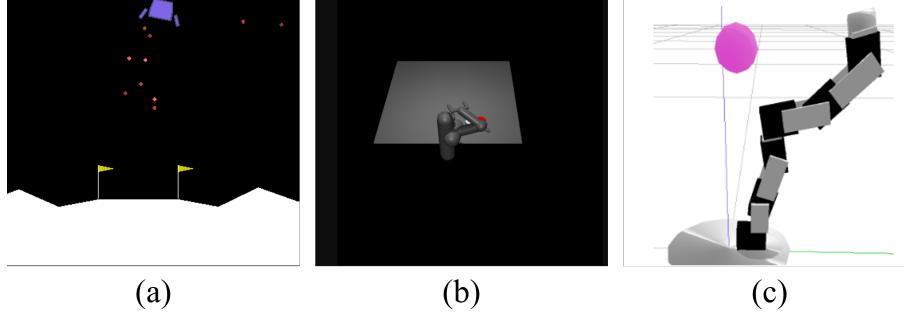


Figure 10. Illustrations of three physical robotics domain randomization scenarios: (a) LunarLander, (b) Pusher, and (c) ErgoReacher. The illustration of ErgoReacher is adapted from (Mehta et al., 2020).

DRM uses a batch size of 40. The temperature parameter for GDRM is configured as 0.001. Both the inner and outer loop learning rates are fixed at 0.1. Besides, the identifier batch size during training is 30 ($1.5\times$) for MPTS and 1280 ($64\times$) for PDTs. The risk learner is updated using the Adam optimizer with a learning rate of 5×10^{-3} . The label for the risk learner is the negative average reward value at the final step for each task. For validation during training, we uniformly sample 40 tasks from the task space at fixed intervals to validate the training checkpoints of different methods. For meta-testing after training, we uniformly sample 100 tasks from the task space to test the trained models of different methods.

C.4. Physical Robotics Domain Randomization

As shown in Fig. 10, we adopt three scenarios for robotics domain randomization from Mehta et al. (2020): Pusher, LunarLander, and ErgoReacher. The task distribution is defined as in (Wang et al., 2025). Specifically:

- LunarLander is a 2-degree-of-freedom (DoF) environment where the agent must softly land a spacecraft, implemented in Box2D (Catto, 2007). Its reward function provides positive rewards for successful landings, negative rewards for crashes, and penalties for fuel consumption and deviations from the landing pad, thereby promoting efficient and controlled landings. The task distribution is specified by a uniform distribution over the main engine strength s , which can be viewed as the task identifier $\tau = s$.
- Pusher is a 3-DoF robotic arm control environment based on MuJoCo (Todorov et al., 2012), where the agent pushes a puck to a target. The reward function penalizes the L_2 distance between the puck and the target, augmented by a control penalty. The task distribution is specified by a uniform distribution over the puck friction loss f and puck joint damping d , which can be viewed as the task identifier $\tau = (f, d)$.
- ErgoReacher involves a 4-DoF robotic arm implemented in the Bullet Physics Engine (Coumans, 2015), tasked with reaching a goal using its end effector. Its reward function penalizes the distance between the end effector and the target, combined with control penalties. The task distribution is specified by a uniform distribution over the joint damping d and max torque t across 4 joints, which can be viewed as the task identifier $\tau = (d_1, d_2, d_3, d_4, t_1, t_2, t_3, t_4)$.

Details of the randomized task identifier range are presented in Table 3.

Implementation Details. The machine learner is a neural network with two hidden layers, each consisting of 10 units, and uses ReLU activation functions. For ERM and GDRM, the task batch size is set to 10, while DRM uses a batch size of 20. GDRM employs a temperature parameter of 0.01. We adopt TD3 algorithm (Fujimoto et al., 2018) as the algorithm backbone. The actor and critic learning rates are both set to 3×10^{-4} . For MPTS, the identifier batch size during training is 25 ($2.5\times$) for LunarLander, 50 ($5\times$), and 250 ($25\times$) for ErgoReacher. In contrast, for PDTs, the identifier batch size during training is 640 ($64\times$) for all environments, with no additional requirements for fine-tuning. The risk learner is updated using the Adam optimizer with a learning rate of 0.005. The label for the risk learner is the negative average return for each task. For validation during training, we uniformly sample 100 tasks from the task space to validate the training checkpoints of different methods.

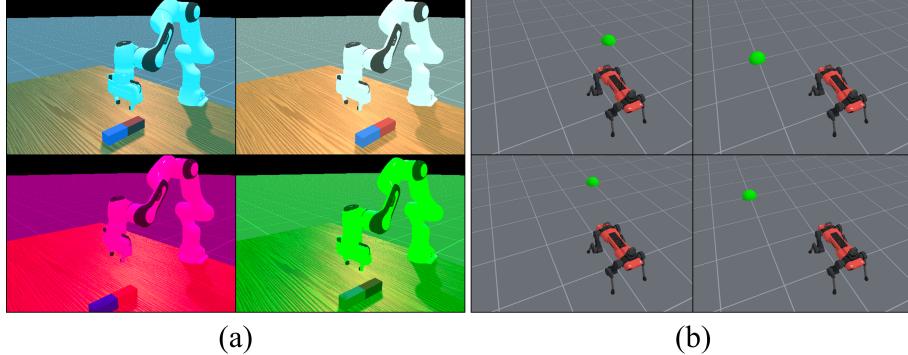


Figure 11. Illustrations of two visual robotics domain randomization scenarios: controlling (a) a table-top robotic arm and (b) a quadruped robot, operating under randomized lighting conditions and varying goal locations, respectively.

C.5. Visual Robotics Domain Randomization

Visual-based robotics control is common and crucial in real-world scenarios. As illustrated in Fig. 11, based on the latest robotics simulator, ManiSkill3 (Tao et al., 2024), we design two scenarios for visual robotics domain randomization: LiftPegUpright.Light and AnymalCReach.Goal. These scenarios involve controlling a tabletop two-finger gripper arm robot and a quadruped robot, respectively, under randomized lighting conditions and goal locations.

Specifically:

- LiftPegUpright.Light is derived from the LiftPegUpright-v1 scenario in ManiSkill3, where the objective is to move a peg lying on the table to an upright position. To emulate the complex lighting conditions found in real-world environments, we modify this scenario to make the ambient lighting controllable. An illustration of this setup is shown in Fig. 5. The task distribution is specified by a uniform distribution over the configurations of ambient light, which can be viewed as the task identifier $\tau = (l_1, l_2, l_3)$.
- AnymalCReach.Goal is adapted from the AnymalC-Reach-v1 scenario in ManiSkill3. The task is to control the AnymalC robot to reach a target location in front of it. Inspired by point-robot navigation scenarios in prior works (Finn et al., 2017), we randomize the goal location, which remains visible to the robot. The task distribution is specified by a uniform distribution over the goal location, which can be viewed as the task identifier $\tau = (x_1, x_2)$.

The detailed randomization configurations are provided in Table 3.

Implementation Details. We use the PPO algorithm (Schulman et al., 2017), along with its hyperparameters and network architecture, as provided in the official ManiSkill3 codebase. To accommodate GPU memory limitations, we set the task batch size to 16 for ERM and GDRM and 32 for DRM, with 500 steps per iteration on LiftPegUpright.Light. For AnymalCReach.Goal, the task batch size is 128 for ERM and GDRM and 256 for DRM, with 200 steps per iteration. For MPTS, the identifier batch size during training is 32 for LiftPegUpright.Light and 256 for AnymalCReach.Goal ($2.5\times$). For PDTTS, the identifier batch size during training is $64\times$ the default for all environments. The risk learner is updated using the Adam optimizer with a learning rate of 0.005. The label for the risk learner is the negative average reward for each task. For validation during training, we uniformly sample 64 tasks from the task space to validate the training checkpoints of different methods.

D. Additional Experiment Results

D.1. Beyond Decision-Making: Robust Supervised Meta-Learning

This work primarily focuses on risk-averse sequential decision-making. However, PDTTS is readily extendable to other risk-averse scenarios, such as supervised meta-learning. As shown in Fig. 12, we evaluate PDTTS on sinusoid regression, a commonly-used toy example introduced in Finn et al. (2017), which involves adapting quickly to new functions using only 10 samples. Consistent with the results observed in decision-making, PDTTS achieves superior performance compared to all baselines, demonstrating faster average performance and more robust adaptation.

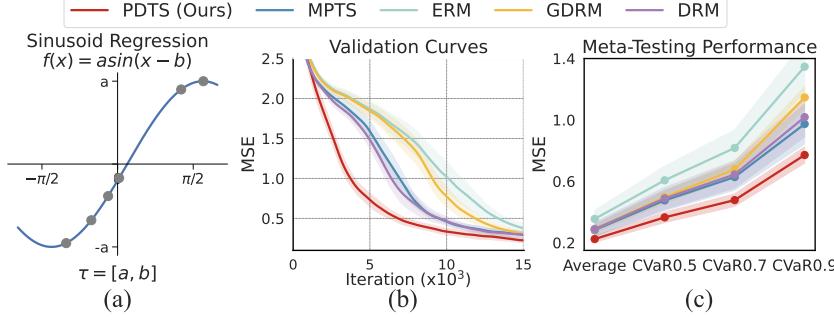


Figure 12. Few-Shot Sinusoid Regression Results. (a) Illustration of the sinusoid regression problem, where the task identifier τ consists of the amplitude and phase $[a, b]$. (b) Curves of averaged MSEs on the validation task set for all methods during meta-training. (c) MSE values for all methods at various α levels of CVaR $_{\alpha}$ during meta-testing.

D.2. Ablation Studies and Additional Analysis

In this section, we perform additional experiments to carry out ablation studies on the hyperparameters and components of PDTS, demonstrate the presence of the concentration issue in MPTS, and validate the effectiveness of PDTS in addressing it. For computational efficiency, we use sinusoid regression as the testbed.

Ablation Study on Diversity Regularization Weight γ . As shown in Fig. 13(a), we evaluate the effect of varying values of the diversity regularization weight γ . It is evident that diversity regularization plays a crucial role in PDTS, as its absence ($\gamma = 0$) leads to a dramatic performance drop. PDTS demonstrates robustness to the choice of γ within a certain range (e.g., [1, 2]). However, excessively large values of γ degrade performance, emphasizing the importance of balancing diversity and robust optimization.

Ablation Study on Key Components: Posterior Sampling and Diversity Regularization. In simple terms, PDTS can be viewed as the combination of MPTS and diversity regularization, with UCB replaced by posterior sampling. We conduct an ablation study to highlight the significance of each component. As shown in Fig. 13(b), we replace the UCB in MPTS with posterior sampling (MPTS+P) and incorporate diversity regularization into MPTS (MPTS+D). The results demonstrate that each component contributes significantly to the superiority of PDTS.

The Presence of the Concentration Issue in MPTS. As introduced in Sec. 3.2, we theoretically prove the presence of a concentration issue in MPTS. To empirically demonstrate the concentration issue and its impact, we evaluate MPTS on sinusoid regression. Fig. 13(c) shows that, as the candidate batch size increases, MPTS tends to select tasks concentrated within a small region—a phenomenon referred to as the concentration issue in the main paper. As illustrated in Fig. 13(d), this concentration issue leads to catastrophic performance degradation in MPTS as the candidate batch size increases.

PDTs Addresses the Concentration Issue and Benefits from Improved Coverage. In contrast to MPTS, Fig. 13(e) demonstrates that by incorporating diversity regularization, our method, PDTs, avoids the concentration issue and does not experience performance collapse as the candidate batch size increases. More impressively, the performance of PDTs improves with increasing candidate batch size, demonstrating the benefits of encouraging broader coverage of the task space during subset selection as proposed by PDTs.

E. Other Clarifications

Relation with Traditional Active Learning. Traditional active learning (Ren et al., 2021) aims at reducing the sampling redundancies during optimization and exploiting historical optimization information to improve learning efficiency, such as annotations and computations. The active query strategies (Zhu et al., 2003; Gal et al., 2017; Kirsch et al., 2019; Wu et al., 2022; Mukhoti et al., 2023) also rely on some predictive models and utilize principles like uncertainty, diversity, etc. RATS, such as MPTS (Wang et al., 2025) and PDTs in this work, stresses the importance of robustness during active sampling. Hence, the predictive model requires scoring the task difficulties without exact evaluation.

When MPTS Meets Diversity Regularization. The diagnosis of the concentration issue in Sec. 3.2 identifies the diversity

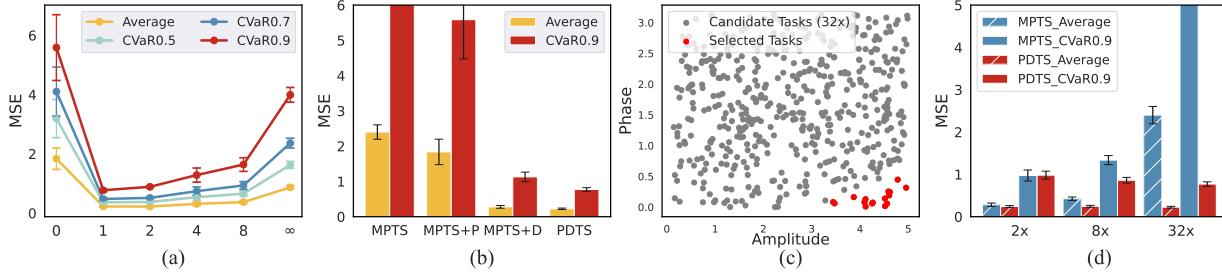


Figure 13. (a) Meta-testing results trained with different hyperparameter, γ . (b) Ablation studies on key components, including posterior sampling (P) and diversity regularization (D). (c) Visualization of the distribution of candidate tasks (gray points) and selected tasks (red points) in MPTS, highlighting the concentration issue. (d) Comparison of the average and CVaR_{0.9} meta-testing performance of PDTTS and MPTS with increasing pseudo-batch sizes.

regularization as a plausible solution to encourage the exploration of the task space and bring more worst-case robust solutions as more arms are constructed by increasing \mathcal{B} . Actually, we also examine this part and include the ablation studies on the sinusoid regression in Fig. 13(b). For other evaluations in the main paper, we assume that MPTS’s group (Wang et al., 2025) adopts the optimal hyper-parameter configurations. Hence, their setup is adopted to produce MPTS results. Meanwhile, the posterior sampling’s advantage lies in (i) no extra hyperparameter adjustment, unlike UCB used in MPTS, and (ii) stochastic optimism when the uncertainty is difficult to estimate.

PDTTS is Agnostic to the Risk Predictive Model. As RATS in this work is a rarely investigated concept in the field, limited methods have been developed to score task difficulties, particularly MDPs’ difficulties under a policy. The risk predictive model in MPTS has approximately achieved the purpose. Hence, we reuse their module as the backbone. In reality, our theoretical analysis and PDTTS will be compatible with other risk-predictive models in the future.

Though robust active learning is a little niche in sequential decision-making, some investigations to reduce computational and sample expense and boost robustness are necessary for establishing large decision models. The trial in this work presents a tractable strategy with impressive performance to achieve RATS purposes; however, such a setup relies on the identifier information and inherent smoothness inside the adaptation risk value function.