
POMRL: No-Regret Learning-to-Plan with Increasing Horizons

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the problem of planning under model uncertainty in an online meta-
2 reinforcement learning (RL) setting where an agent is presented with a sequence of
3 related tasks with limited interactions per task. The agent can use its experience in
4 each task *and* across tasks to estimate both the transition model and the distribution
5 over tasks. We propose an algorithm to meta-learn the underlying relatedness across
6 tasks, utilize it to plan in each task, and upper-bound the regret of the planning
7 loss. Our bound suggests that the average regret over tasks decreases as the number
8 of tasks increases and as the tasks are more similar. In the classical single-task
9 setting, it is known that the planning horizon should depend on the estimated
10 model’s accuracy, that is, on the number of samples within task. We generalize this
11 finding to meta-RL and study this dependence of planning horizon on the number
12 of tasks. Based on our theoretical findings, we derive heuristics for selecting slowly
13 increasing discount factors, and validate its significance empirically.

14 1 Introduction

15 *Meta-learning* (Caruana, 1997; Baxter, 2000; Thrun and Pratt, 1998; Finn et al., 2017; Denevi et al.,
16 2018) offers a powerful paradigm to leverage past experience to reduce the sample complexity of
17 learning future related tasks. *Online meta-learning* considers a sequential setting, where the agent
18 progressively accumulates knowledge and uses past experience to learn good priors and to quickly
19 adapt within each task Finn et al. (2019); Denevi et al. (2019). Robots acting in real world for instance
20 need to be responsive to and robust against perturbation inherent in the environment dynamics and
21 their decision making. When the tasks share a structure i.e. have similar transition dynamics and
22 are related, such approaches enable progressively faster convergence, or equivalently better model
23 accuracy with better sample complexity (Schmidhuber and Huber, 1991; Thrun and Pratt, 1998;
24 Baxter, 2000; Finn et al., 2017; Balcan et al., 2019).

25 In model-based reinforcement learning (RL), the agent uses an estimated model of the environment
26 to plan actions ahead towards the goal of maximizing rewards. A key component in the agent’s
27 decision making is the horizon used during planning. In general, an *evaluation horizon* is imposed by
28 the task itself, but the learner may want to use a different and potentially shorter *guidance horizon*.
29 In the discounted setting, the size of the evaluation horizon is of order $(1 - \gamma_{\text{eval}})^{-1}$, for some
30 discount factor $\gamma_{\text{eval}} \in (0, 1)$, and the agent may use $\gamma \neq \gamma_{\text{eval}}$ for planning. For instance, a classic
31 result known as Blackwell Optimality (Blackwell, 1962) states there exists a discount factor γ^*
32 and a corresponding optimal policy such that the policy is also optimal for any greater discount
33 factor $\gamma \geq \gamma^*$. Thus, an agent that plans with $\gamma = \gamma^*$ will be optimal for any $\gamma_{\text{eval}} > \gamma^*$. In the
34 Arcade Learning Environment (Bellemare et al., 2013) a discount factor of $\gamma_{\text{eval}} = 1$ is used for
35 evaluation, but typically a smaller γ is used for training (Mnih et al., 2015). Using a smaller discount
36 factor acts as a regularizer (Amit et al., 2020; Petrik and Scherrer, 2008; Van Seijen et al., 2009;

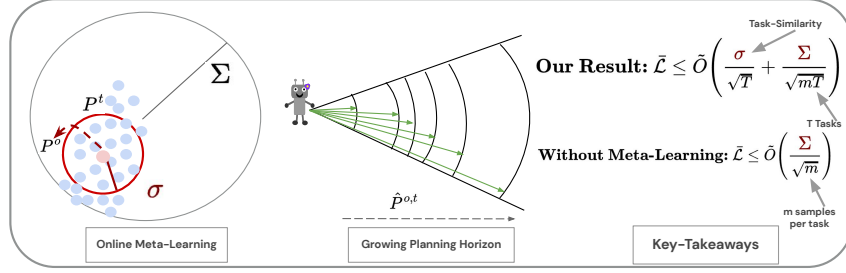


Figure 1: **Effective Planning Horizons in Meta-Reinforcement Learning.** The agent faces a sequence of tasks with transition vector $(P^t)_{t \in [T]}$ (probability vectors represented by blue dots) all close to each other ($\sigma < \Sigma = 1$). The agent builds a transition model for each task and plans with these inaccurate models. By using data from previous tasks, the agent meta-learns an initialization of the model ($\hat{P}^{o,t}$), which leads to better planning in new related but unseen tasks. We show an improved average regret upper bound that scales with task-similarity parameter σ and inversely with the number of tasks T : as knowledge accumulates, uncertainty diminishes, and the agent can plan with longer horizons. All tasks $P^t \sim \mathcal{P}$ are centered at some fixed but unknown P^o , depicted here by the shaded red dot and pointed by the arrow.

37 François-Lavet et al., 2019; Arumugam et al., 2018) and reduces planner over-fitting in random
 38 MDPs (Arumugam et al., 2018). Indeed, the choice of planning horizon plays a significant role in
 39 computation (Kearns et al., 2002), optimality (Kocsis and Szepesvári, 2006), and on the complexity of
 40 the policy class (Jiang et al., 2015). In addition, meta-learning discount factors has led to significant
 41 improvements in performance (Xu et al., 2018; Zahavy et al., 2020; Flennerhag et al., 2021, 2022;
 42 Luketina et al., 2022).

43 When doing model-based RL with a learned model, the optimal guidance planning horizon, called
 44 *effective horizon* by Jiang et al. (2015), depends on the accuracy of the model, and so on the amount
 45 of data used to estimate it. Jiang et al. (2015) show that when data is scarce, a guidance discount
 46 factor $\gamma < \gamma_{\text{eval}}$ should be preferred for planning. The reason for this is straightforward; if the model
 47 used for planning is inaccurate, then errors will tend to accumulate along the planned trajectory. A
 48 shorter effective planning horizon will accumulate less error and may lead to better performance,
 49 even when judged using the true γ_{eval} . While that work treated only the batch, single-task setting,
 50 the question of effective planning horizon remains open in the online meta-learning setting where the
 51 agent accumulates knowledge from many tasks, with limited interactions within each task.

52 In this work, we consider a *meta-reinforcement-learning* problem made of a sequence of **related**
 53 **tasks**. We leverage this structural task similarity to obtain model estimators with faster convergence
 54 as more tasks are seen. The central question of our work is: *Can we meta-learn the model across*
 55 *tasks and adapt the effective planning horizon accordingly?*

56 We take inspiration from the *Average Regret-Upper-Bound Analysis* [ARUBA] framework (Khodak
 57 et al., 2019) to generalize planning loss bounds to the meta-RL setting. A high-level, intuitive outline
 58 of our approach is presented in Fig. 1. **Our main contributions** are as follows: 1) We formalize
 59 planning in a model-based meta-RL setting as an *average planning loss* minimization problem, and
 60 we propose an algorithm to solve it, 2) Under a structural *task-similarity* assumption, we prove a novel
 61 high-probability task-averaged regret upper-bound on the planning loss of our algorithm, inspired by
 62 ARUBA. We also demonstrate a way to learn the task-similarity parameter σ on-the-fly. To the best
 63 of our knowledge, this is a first formal (ARUBA-style) analysis to show that meta-RL can be more
 64 efficient than RL, and 3) Our theoretical result highlights a new dependence of the planning horizon
 65 on the size of the within-task data m and on the number of tasks T . This observation allows us to
 66 propose two heuristics to adapt the planning horizon given the overall sample-size.

67 2 Background & Illustration

68 In practice, the true model of the world is unknown and must be estimated from data¹. One approach
 69 to approximately solve the optimization problem above is to construct a model, $\langle \hat{R}, \hat{P} \rangle$ from data, then

¹We defer the reader to Appendix Sec. B for detailed background.

70 find $\pi_{\hat{M},\gamma}^*$ for the corresponding MDP $\hat{M} = \langle \mathcal{S}, \mathcal{A}, \hat{R}, \hat{P}, \gamma \rangle$. This approach is called *model-based*
 71 *RL* or *certainty-equivalence (CE) control*. In this setting, [Jiang et al. \(2015\)](#) define the planning
 72 loss as the gap in expected return in MDP M when using $\gamma \leq \gamma_{\text{eval}}$ and the optimal policy for an
 73 approximate model \hat{M} :

$$\mathcal{L}(\hat{M}, \gamma | M, \gamma_{\text{eval}}) = \|V_{M, \gamma_{\text{eval}}}^{\pi_{\hat{M}, \gamma_{\text{eval}}}^*} - V_{M, \gamma_{\text{eval}}}^{\pi_{\hat{M}, \gamma}^*}\|_{\infty}.$$

74 Thus, the **optimal effective planning horizon** $(1 - \gamma^*)^{-1}$ is defined using the discount factor that
 75 minimizes the planning loss, *i.e.*, $\gamma^* := \min_{0 \leq \gamma \leq \gamma_{\text{eval}}} \mathcal{L}(\hat{M}, \gamma | M, \gamma_{\text{eval}})$.

76 **Theorem 1.** ([Jiang et al. \(2015\)](#)) *Let M be an MDP with non-negative bounded rewards and*
 77 *evaluation discount factor γ_{eval} . Let \hat{M} be the approximate MDP comprising the true reward*
 78 *function of M and the approximate transition model \hat{P} , estimated from $m > 0$ samples for each*
 79 *state-action pair. Then, with probability at least $1 - \delta$,*

$$\|V_{M, \gamma_{\text{eval}}}^{\pi_{\hat{M}, \gamma_{\text{eval}}}^*} - V_{M, \gamma_{\text{eval}}}^{\pi_{\hat{M}, \gamma}^*}\|_{\infty} \leq \frac{\gamma_{\text{eval}} - \gamma}{(1 - \gamma_{\text{eval}})(1 - \gamma)} + \frac{2\gamma R_{\max}}{(1 - \gamma)^2} \left(\sqrt{\frac{\Sigma}{2m} \log \frac{2SA|\Pi_{\gamma}|}{\delta}} \right) \quad (1)$$

80 where Σ is upper-bounded by 1 as $P, \hat{P} \in \Delta_{\mathcal{S}}$.

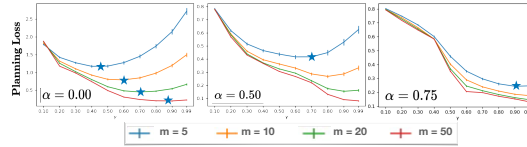


Figure 2: **On the role of incorporating a ground truth prior of transition model on planning horizon.** The planning loss is a function of the discount factor γ and is impacted by incorporating prior knowledge. The learner has m samples per task to estimate the model, corresponding to the curves in each sub figure. Inspecting any sub figure, we observe that larger values of m lead to lower planning loss and a larger effective discount factor. Besides, inspecting one value of m across tasks (e.g., $m = 5$), the same effect (lower planning loss and larger effective discount) occurs when the learner puts more weight on the ground truth prior through α .

81 These effects are illustrated in [Fig. B2](#) on a 10-state, 2-action random MDP. The leftmost plot uses
 82 the simple count-based model estimator and reproduces the results from [Jiang et al. \(2015\)](#). We then
 83 incorporate the true prior (mean model P^o as in [Fig 1](#) and defined above [Eq. 3](#) in [Assumption 1](#)) in
 84 the estimator with a growing mixing factor $\alpha \in (0, 1)$: $\hat{P}(m) = \alpha P^o + (1 - \alpha) \frac{\sum^i X^i}{m}$. We observe
 85 that increasing the weight $\alpha \in (0, 1)$ on good prior knowledge enables longer planning horizons and
 86 lower planning loss.

87 We consider an **online meta-RL problem** where an agent is presented with a sequence of tasks
 88 M_1, \dots, M_T , where for each $t \in [T]$, $M_t = \langle \mathcal{S}, \mathcal{A}, P^t, R, \gamma_{\text{eval}} \rangle$, that is, the MDPs only differ from
 89 each other by the transition matrix (dynamics model) P^t . The learner must sequentially estimate the
 90 model \hat{P}^t for each task t from a batch of m transitions simulated for each state-action pair. Its goal is
 91 to minimize the average planning loss also expressed in the form of task averaged regret suffered in
 92 planning and defined as

$$\bar{\mathcal{L}}(\hat{M}_{1:T}, \gamma | M_{1:T}, \gamma_{\text{eval}}) = \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\hat{M}_t, \gamma | M_t, \gamma_{\text{eval}}) = \frac{1}{T} \sum_{t=1}^T \|V_{M_t, \gamma_{\text{eval}}}^{\pi_{\hat{M}_t, \gamma_{\text{eval}}}^*} - V_{M_t, \gamma_{\text{eval}}}^{\pi_{\hat{M}_t, \gamma}^*}\|_{\infty} \quad (2)$$

93 3 Planning with Online Meta-Reinforcement Learning

94 We here formalize planning in a model-based meta-RL setting. We start by specifying all our
 95 assumptions in [Sec 3.1](#) including our main assumption about task relatedness in [Sec. 1](#), present our
 96 approach and explain the proposed algorithms POMRL and ada-POMRL in [Sec. 3.2](#). Our main result
 97 is a high-probability upper bound on the average planning loss under the assumed task relatedness,
 98 presented as [Theorem 2](#).

99 **3.1 Assumptions**

100 In many real world scenarios such as robotics, it is required to be responsive to changes in the
 101 environment and, at the same time, to be robust against perturbation inherent in the environment and
 102 their decision making. In such practical scenarios, the key reason to employ meta-learning is for the
 103 learner to leverage **task-similarity** (or task variance) across tasks.

104 **Assumption 1** (Structural Assumption Across Tasks: Task Relatedness). *In this work, we exploit the*
 105 *structural assumption that for all $t \in [T]$, $P^t \sim \mathcal{P}$ centered at some fixed but unknown $P^o \in \Delta_S^{S \times A}$*
 106 *and such that for any (s, a) ,*

$$\|P_{s,a}^t - P_{s,a}^o\|_\infty \leq \sigma = \max_{(s,a)} \sigma(s, a) \quad a.s. \quad (3)$$

107 This also implies that $\max_{t,t'} \|P_{s,a}^t - P_{s,a}^{t'}\|_\infty \leq 2\sigma$, and that the meta-distribution \mathcal{P} is bounded
 108 within a small subset of the simplex. It is immediate to extend our results under a high-probability
 109 assumption instead of the almost sure statement above. In our experiments, we will use Gaussian
 110 or Dirichlet priors over the simplex, whose moments are bounded with high-probability, not almost
 111 surely. Importantly, we will say that a multi-task environment is *strongly structured* when $\sigma < \Sigma$,
 112 *i.e.* when the effective diameter of the models is smaller than that of the entire feasible space.

113 **Assumption 2** (Access to a Simulator). *We assume that for each task $t \in [T]$ we have access to*
 114 *a simulator of transitions (Kearns et al., 2002) providing m i.i.d. samples $(X_{s,a}^{t,i})_{i=1..m} \in \mathcal{S}^m \sim$*
 115 *$P^t(\cdot|s, a)$ (categorical distribution).*

116 **Assumption 3** (Known Rewards). *Given a distribution of tasks, the rewards are known.*

117 **3.2 Our Approach**

118 With access to a simulator (Assumption 2); for each (s, a) , we can compute an empirical estimator
 119 for each $s' \in [S]$: $\hat{P}_{s,a}^t(s') = \sum_{i=1}^m \mathbb{1}\{X_{s,a}^{t,i} = s'\}/m$, with naturally $\sum_{s'} \hat{P}_{s,a}^t(s') = 1$. We perform
 120 meta-RL via alternating minimizing a batch *within-task* regularized least-squares loss, and an outer-
 121 loop step where we optimize the regularization to optimally balance bias and variance of the next
 122 estimator.

123 **Estimating dynamics model via regularized least squares.** We adapt the standard technique of
 124 meta-learned regularizer (see e.g. Baxter (2000); Cella et al. (2020) for supervised learning and
 125 bandit respectively) to this model estimation problem. At each round t , the **current model** $\hat{P}_{(s,a)}^t$ is
 126 estimated by minimizing a **regularized least square loss**: for a given **regularizer** h_t (to be specified
 127 below) and parameter $\lambda_t > 0$ for each $(s, a) \in \mathcal{S} \times \mathcal{A}$ we solve

$$\hat{P}_{(s,a)}^t = \arg \min_{P_{(s,a)} \in \Delta_S} \left\| \underbrace{\frac{1}{m} \sum_{i=1}^m \mathbb{1}\{X_{s,a}^{t,i}\}}_{\text{empirical transition prob.}} - P_{(s,a)} \right\|_2^2 + \lambda_t \|P_{(s,a)} - h_t\|_2^2, \quad (4)$$

128 where we use $\mathbb{1}\{X_{s,a}^{t,i}\}$ to denote the one-hot encoding of the state into a vector in \mathbb{R}^S . Importantly, h_t
 129 and λ_t are meta-learned in the outer-loop (see below) and affect the bias and variance of the resulting
 130 estimator. The solution of equation 4 can be computed in closed form as a convex combination of
 131 the empirical average (count-based) and the prior: $\hat{P}^t = \alpha_t h_t + (1 - \alpha_t) \bar{P}^t$ where $\alpha_t = \frac{\lambda_t}{1 + \lambda_t}$ is the
 132 current mixing parameter.

133 **Outer-loop: Meta-learning the regularization.** At the beginning of task $1 < t \leq T$, the learner
 134 has already observed $t - 1$ *related but different* tasks. We define h_t as an **average of Means (AoM)**:

$$h_{(s,a)}^t \leftarrow \hat{P}_{(s,a)}^{o,t} = \frac{1}{t-1} \sum_{j=1}^{t-1} \frac{\sum_{i=1}^m \mathbb{1}\{X_{(s,a)}^{j,i}\}}{m} := \frac{1}{t-1} \sum_{j=1}^{t-1} \bar{P}_{(s,a)}^j. \quad (5)$$

135 **Deriving the mixing rate.** To set α_t , we compute the Mean Squared Error (MSE) of $\hat{P}_{(s,a)}^t$, and
 136 minimize an upper bound (see details in Appendix D): $\text{MSE}(\hat{P}_{(s,a)}^t) \leq \alpha_t^2 \sigma^2 (1 + \frac{1}{t}) + (1 - \alpha_t)^2 \frac{1}{m}$,
 137 which leads to $\alpha_t = \frac{1}{\sigma^2(1+1/t)m+1}$.

138 **Alg 1** depicts the complete pseudo code. We note here that POMRL (σ) assumes, for now, that the under-
 139 lying task-similarity parameter σ is known, and we discuss a fully empirical extension further below
 140 (See Sec. 4). The learner does not know the number of tasks a priori and tasks are faced sequentially
 141 online. The learner performs meta-RL alternating between within-task estimation of the dynamics
 142 model \hat{P}^t via a batch of m samples for that task, and an outer loop step to meta-update the regularizer
 143 $\hat{P}^{o,t+1}$ alongside the mixing rate α_{t+1} . For each task, we use a γ -Selection-Procedure to choose
 144 planning horizon $\gamma^* \leq \gamma_{\text{eval}}$. We defer the details of this step to Sec. 6 as it is non-trivial and only a
 145 partial consequence of our theoretical analysis. Next, the learner performs planning with an imperfect
 146 model \hat{P}^t . For planning, we use dynamic programming, in particular policy iteration (a combination
 147 of policy evaluation, and improvement), and value iteration to obtain the optimal policy $\pi_{\hat{P}^t, \gamma^*}^*$ for
 148 the corresponding MDP \hat{M}_t .

Algorithm 1: POMRL (σ) – Planning with Online Meta-Reinforcement Learning

Input: Given task-similarity ($\sigma(s, a)$) a matrix of size $S \times A$. Initialize $\hat{P}^{o,1}$ to uniform, $\alpha_1 = 0$.

for task $t \in [T]$ **do**

for t^{th} batch of m samples **do**

$\hat{P}^t(m) = (1 - \alpha_t) \frac{1}{m} \sum_{i=1}^m X_i + \alpha_t \hat{P}^{o,t}$ // regularized least squares
 minimizer.

$\gamma^* \leftarrow \gamma\text{-Selection-Procedure}(m, \alpha_t, \sigma, T, S, A)$

$\pi_{\hat{P}^t, \gamma^*}^* \leftarrow \text{Planning}(\hat{P}^t(m))$ //

Output: $\pi_{\hat{P}^t, \gamma^*}^*$

 Update $\hat{P}^{o,t+1}, \alpha_{t+1} = \frac{1}{\sigma^2(1+1/t)m+1}$ // meta-update AoM (Eq. 5) & mixing rate

150 **3.3 Average Regret Bound for Planning with Online-meta-learning**

151 Our main theoretical result below controls the average regret of POMRL (σ), a version of Alg. 1 with
 152 additional knowledge of the underlying task relatedness, *i.e.*, the true $\sigma > 0$.

153 **Theorem 2.** *Using the notation of Theorem 1, we bound the average planning loss equation 10 for*
 154 *POMRL (σ):*

$$\bar{\mathcal{L}} \leq \frac{\gamma_{\text{eval}} - \gamma}{(1 - \gamma_{\text{eval}})(1 - \gamma)} + \frac{2\gamma S}{(1 - \gamma)^2} \tilde{O} \left(\frac{\sigma + \sqrt{\frac{1}{T} \left(\sigma + \sqrt{\sigma^2 + \frac{\Sigma}{m}} \right)}}{\sigma^2 m + 1} + \frac{\sigma^2 m \sqrt{\frac{\Sigma}{m}}}{\sigma^2 m + 1} \right) \quad (6)$$

155 *with probability at least $1 - \delta$, where $\sigma^2 < 1$ is the measure of the task-similarity and $\sigma =$*
 156 *$\max_{(s,a)} \sigma(s, a)$.*

157 The proof is provided in Appendix F and relies on a new concentration bound for the meta-learned
 158 model estimator. The last term on the r.h.s. corresponds to the uncertainty on the dynamics. First we
 159 verify that if $T = 1$ and m grows large, the second term dominates and is equivalent to $\tilde{O}(\sqrt{\frac{\Sigma}{m}})$ (as
 160 $\sigma^2/(\sigma^2 m + 1) \rightarrow 0$), which is similar to that of Jiang et al. (2015) as there is no meta-learning, with
 161 an additional $O(\frac{1}{m})$ but second order term due to the introduced bias. Then, if m is fixed and small,
 162 for small enough values of σ^2 (typically $\sigma < 1/\sqrt{m}$), the first term dominates and the r.h.s. boils
 163 down to $\tilde{O} \left((\sigma + \frac{1}{\sqrt{m}}) / \sqrt{T} \right)$. This highlights the interplay of our structural assumption parameter σ
 164 and the amount of data m available at each round. The regimes of the bound for various similarity
 165 levels are explored empirically in Sec. 5 (Q3). We also show the dependence of the regret upper
 166 bound on m and T for a fixed σ , in Appendix Fig. H5.

167 **4 Practical Considerations: Adaption On-The-Fly**

168 In this section we propose a variant of POMRL that meta learns the task similarity parameter, which
 169 we call ada-POMRL. We compare the two algorithms empirically in a 10 state, 2 action MDP with
 170 closely related tasks with a total of $T = 15$ tasks (details of the setup are deferred to Sec. 5).

171 **Performance of POMRL** . Recall that POMRL is primarily learning the regularizer and assumes the
 172 knowledge of the underlying task similarity (i.e. σ). We observe in Fig. 3 that with each round
 173 $t \in T$ POMRL is able to plan better as it learns and adapts the regularizer to the incoming tasks. The
 174 convergence rate and final performance corroborates with our theory.

175 **Can we also meta-learn the task-similarity parameter?** In practice, the parameter σ may
 176 not be known and must be estimated online and plugged in (see Appendix E for details).
 177

178 Alg. 2 ada-POMRL uses Welford’s algorithm to compute an online
 179 estimate of the variance after every task using the model estimators,
 180 and simply plugs-in this estimate wherever POMRL was using the true
 181 value. From the perspective of ada-POMRL, POMRL is an "oracle", i.e.
 182 the underlying task-similarity is known. However, in most practical
 183 scenarios, the learner does not have this information a priori.

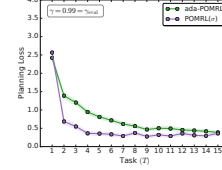


Figure 3: ada-POMRL enables meta-learning the task-similarity on-the-fly with a performance gap for the initial tasks compared to POMRL, but improves with more tasks

184 We compare empirically POMRL and ada-POMRL on a strongly structured
 185 problem ($\sigma \approx 0.01$) in Fig. 3 and observe that meta-learning the
 186 underlying task relatedness allows ada-POMRL to adapt to the incoming
 187 tasks accordingly. Adaptation on-the-fly with ada-POMRL comes at a
 188 cost i.e., the performance gap in comparison to POMRL but eventually
 189 converges albeit with a slower rate.

190 This online estimation of σ means that ada-POMRL now requires an initial value for $\hat{\sigma}_1$, which is
 191 a choice left to the practitioner, but will only affect the results of a finite number of tasks at the
 192 beginning. Using $\hat{\sigma}_1$ too small will give a slightly increased weight to the prior in initial tasks, which
 193 is not desirable as the latter is not yet learned and will result in an increased bias. On the other hand,
 194 setting $\hat{\sigma}_1$ too large (i.e close to 1/2) will decrease the weight of the prior and increase the variance
 195 of the returned solution; in particular, in cases where the true σ is small, a large initialization will
 196 slow down convergence and we observe empirical larger gaps between POMRL and ada-POMRL. In
 197 the extreme case where $\sigma \approx 0$, a large initialization will drastically slow down ada-POMRL as it will
 198 take many tasks before it *discovers* that the optimal behavior is essentially to aggregate the batches.

Algorithm 2: ada-POMRL – Planning with Online Meta-Reinforcement Learning

Input: Initialize $\hat{P}^{o,1}$ to uniform, $(\hat{\sigma})_1$ as a matrix of size $S \times A, \alpha_1 = 0$.

for task $t \in [T]$ **do**

for t^{th} batch of m samples **do**

$\hat{P}^t(m) = (1 - \alpha_t) \frac{1}{m} \sum_{i=1}^m X_i + \alpha_t \hat{P}^{o,t}$ // regularized least squares
 minimizer.

$\gamma^* \leftarrow \gamma$ -Selection-Procedure($m, \alpha_t, \sigma_t, T, S, A$)

$\pi_{\hat{P}^t, \gamma^*}^* \leftarrow$ Planning($\hat{P}^t(m)$)

Output: $\pi_{\hat{P}^t, \gamma^*}^*$

 Update $\hat{P}^{o,t+1}, \hat{\sigma}_{t+1} \leftarrow$ Welford’s online algorithm($(\hat{\sigma}_o)_t, \hat{P}^{o,t+1}, \hat{P}^{o,t}$) //

 meta-update AoM (Eq. 5) and task-similarity parameter.

 Update $\alpha_{t+1} = \frac{1}{\hat{\sigma}_{t+1}^2(1+1/t)m+1}$ // meta-update mixing rate, plug $\max(\sigma_{S \times A})$

200 **Tasks vary only in certain states and actions.** Thus far, we considered a *uniform* notion of task
 201 similarity as Eq. 3 holds for any (s, a) . However, in many practical settings the transition distribution
 202 might remains the same for most part of the state space but only vary on some states across different
 203 tasks. These scenarios are hard to analyse in general because local changes in the model parameters
 204 do not always imply changes in the optimal value function nor necessarily modify the optimal policy.
 205 Our Theorem 2 still remains valid, but it may not be tight when the meta-distribution has non-uniform
 206 noise levels. More precisely Theorem 2 in Appendix F remains locally valid for each (s, a) pair and
 207 one could easily replace the uniform σ with local $\sigma_{(s,a)}$, but this cannot directly imply a stronger
 208 bound on the average planning loss. Indeed, in our experiments, in both POMRL and ada-POMRL, the
 209 parameter σ and $\hat{\sigma}$ respectively, are $S \times A$ matrices of state-action dependent variances resulting in
 210 state-action dependent mixing rate α_t .

211 **5 Experiments**

212 We study the empirical behavior of planning with online meta-learning and affirmatively answer the
 213 following questions: **Q1.** Does meta-learning a good initialization of the dynamics model facilitate
 214 improved planning accuracy for the choice of $\gamma = \gamma_{\text{eval}}$? (Sec. 5.1) **Q2.** Does meta-learning a good
 215 initialization of the dynamics model enables longer planning horizons? (Sec. 5.2) **Q3.** How does
 216 performance depend on the amount of shared structure across tasks *i.e.*, σ ? (Sec. 5.3)

217 **Setting:** For each experiment, we fix a mean model $P^o \in \Delta_S^{S \times A}$ (see below how), and for each new
 218 task $t \in [T]$, we sample P^t from a Dirichlet distribution² centered at P^o . As prescribed by theory
 219 (see Sec. 3.2), we set³ $\sigma \approx 0.01 \lesssim 1/S\sqrt{m}$ unless otherwise specified (see Q3). Note that σ and
 220 $\hat{\sigma}$ respectively, are $S \times A$ matrices of state-action dependent variances that capture the directional
 221 variance as we used Dirichlet distributions as priors and these have non-uniform variance levels in
 222 the simplex, depending on how close to the simplex boundary the mean is located. Aligned with our
 223 theory, we use the max of the σ matrices resulting in the aforementioned single scalar value. As in
 224 Jiang et al. (2015), P^o (and each P^t) characterizes a random chain MDP with $S = 10$ states⁴ and
 225 $A = 2$ actions, which is drawn such that, for each state-action pair, the transition function $P(s, a, s')$
 226 is constructed by choosing randomly $k = 5$ states whose probability is set to 0. Then we draw the
 227 value of the $S - k$ remaining states uniformly in $[0, 1]$ and normalize the resulting vector.

228 **5.1 Meta-reinforcement learning leads to improved planning accuracy for $[\gamma_{\text{eval}}]$. [Q1.]**

229 We consider the aforementioned problem setting with a total of $T = 15$ closely related tasks and
 230 focus on the planning loss gains due to improved model accuracy. We fix $\gamma = \gamma_{\text{eval}}$, a rather
 231 naive γ -Selection-Procedure and show the planning loss of POMRL (Alg. 1) with the following
 232 **baselines:** 1) **Oracle Prior Knowledge** knows a priori the underlying task structure (P^o, σ) and
 233 uses an estimator (Eq. 4) with exact regularizer P^o and optimal mixing rate $\alpha_t = \frac{1}{\sigma^2(1+1/t)m+1}$.
 234 2) **Without Meta-Learning** simply uses $\hat{P}^t = \bar{P}^t$, the count-based estimated model using the m
 235 samples seen in each task, 3) POMRL (Alg. 1) meta-learns the regularizer but knows a priori the
 236 underlying task structure, and 4) ada-POMRL (Alg. 2) meta-learns not only the regularizer, but also
 237 the underlying task-similarity online. The oracle is a strong baseline that provides a minimally
 238 inaccurate model and should play the role of an "empirical lower bound". For all baselines, the
 239 number of samples per task $m = 5$. Results are averaged over 100 independent runs. Besides, we
 240 also propose and empirically validate competitive heuristics for γ -Selection-Procedure in Sec. 6.
 241 Besides, we also run another baseline called Aggregating($\alpha = 1$), that simply ignores the meta-RL
 242 structure and just plans assuming there is a single task (See Appendix H.2).

243 **Inspecting Fig. 4(a)**, we can see that our approach ada-POMRL (green) results in decreasing per-task
 244 planning loss as more tasks are seen, and decreasing variance as the estimated model gets more stable
 245 and approaches the optimal value returned by the oracle prior knowledge baseline (blue). On the
 246 contrary, without meta-learning (red), the agent struggles to cope as it faces new tasks every round,
 247 and its performance does not improve. ada-POMRL gradually improves as more tasks are seen whilst
 248 adaptation to learned task-similarity on-the-fly which is the primary cause of the performance gap
 249 in ada-POMRL and POMRL. Importantly, no prior knowledge about the underlying task relatedness
 250 enables a more practical algorithm with the same theoretical guarantees (See Sec. 4). Recall that
 251 oracle prior knowledge is a strong baseline as it corresponds to both known task relatedness and
 252 regularizer.

253 **5.2 Meta-learning the underlying task relatedness enables longer planning horizons. [Q2.]**

254 We run ada-POMRL for $T = 15$ (with $\sigma \approx 0.01$) as above and report planning losses for a range
 255 of values of guidance γ factors. Results are averaged over 100 independent runs and displayed on
 256 Fig. 4(b). We observe in Fig. 4(b) when the agent has seen fewer tasks T , an intermediate value

²The variance of this distribution is controlled by its coefficient parameters $\alpha_{1:S}$: the larger they are, the smaller is the variance. More details on our choices are given in Appendix H.1. Dirichlet distributions with small variance satisfy the high-probability version of our structural assumption 3 for $\sigma = \max_i \sigma_i$

³Our priors are multivariate Dirichlet distribution in dimension S so we divide the theoretical rate by S to ensure the max bounded by $1/\sqrt{m}$. See App. H for implementation details.

⁴We provide additional experiments with varying size of the state space in Appendix Fig. H7.

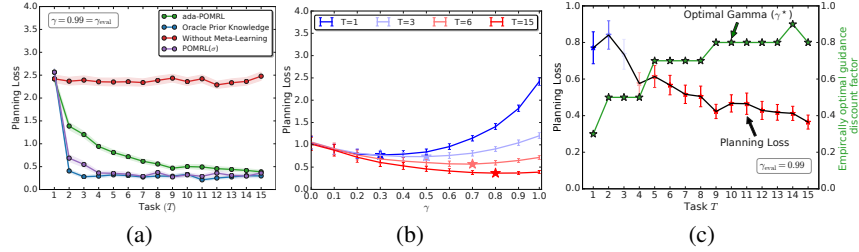


Figure 4: **Planning with Online Meta-Learning.** (a) **Per-task planning loss** of our algorithms POMRL and ada-POMRL compared to an Oracle, and Without Meta-learning baselines. All methods use a fixed $\gamma = \gamma_{eval} = 0.99$. (b) ada-POMRL’s **planning loss** decreases as more tasks are seen. Markers denote the γ that minimizes the planning loss in respective tasks. Error bars show standard error. (c) ada-POMRL’s **empirically optimal guidance discount factor** (right y axis) depicts the effective planning horizon, *i.e.*, one that minimizes the planning loss. Optimal γ aka the effective planning horizon is larger with online meta-learning. Planning loss (left y axis) shows the minimum planning loss achieved by the agent in that round T . Results are averaged over 100 independent runs and error bars represent 1-standard deviation.

257 of the discount is optimal, *i.e.*, one that minimizes the task-averaged planning loss ($\gamma^* < 0.5$). In
 258 the presence of strong underlying structure across tasks, **as the agent sees more tasks, the effective**
 259 **planning horizon** ($\gamma^* > 0.7$) **shifts to a larger value** - one that is closer to the gamma used for
 260 evaluation ($\gamma_{eval} = 0.99$).

261 As we incorporate the knowledge of the underlying task distribution, *i.e.*, meta-learned initialization
 262 of the dynamics model, we note that the adaptive mixing rate α_t puts increasing amounts of weight
 263 on the shared task-knowledge. Note that this conforms to the effect of increasing weight on the
 264 model initialization that we observed in Fig. B2. As predicted by theory, the per-task planning loss
 265 decreases as T grows and is minimized for progressively larger values of γ , meaning for longer
 266 planning horizons (See Fig. 4(c)). In addition, Appendix Fig. H6 depicts the effective planning
 267 horizon individually for ada-POMRL, Oracle and without meta learning baselines.

268 5.3 POMRL and ada-POMRL perform consistently well for varying task-similarity. [Q3.]

269 We have thus far studied scenarios where the learner can exploit strong task relatedness, *i.e.*, $\sigma \approx$
 270 $0.01 < 1/(S\sqrt{m})$ (for low data per task *i.e.*, $m = 5$) is small and we now illustrate the other regimes
 271 discussed in Section 3.2. We find that our algorithms remain consistently good for all amounts of
 272 task-similarity in Appendix Fig. A1.

273 6 Adaptation of planning horizon γ

274 We now propose and empirically validate two heuristics to design an adaptive schedule for γ based
 275 on existing work and on our average regret upper bound.

276 **Schedule adapted from Dong et al. (2021)** [$\gamma = f(m, \alpha_t, \sigma_t, T)$] Dong et al. (2021) study a
 277 continuous, never-ending RL setting. They divide the time into growing phases $(T_t)_{t \geq 0}$, and tune a
 278 discount factor $\gamma_t = 1 - 1/T_t^{1/5}$. We adapt their schedule to our problem, where the time is already
 279 naturally divided into tasks: for each $t \geq 0$, we define the phase size T_t and the corresponding γ_t as

$$T_0 = m, \quad T_t = \frac{SA}{L} \underbrace{\left((1 - \alpha_t)m + \alpha_t m(t - 1) \right)}_{\text{efficient sample size}}, \quad \gamma_t = 1 - \frac{1}{T_t^{1/5}},$$

280 where L is the maximum trajectory length. The size of each T_t , $t \geq 1$, is controlled by an "efficient
 281 sample size" which includes a combination of the current task’s samples and of the samples observed
 282 so far, as used to construct our estimator in POMRL.

283 **Using the upper bound to guide the schedule** [$\gamma = \min\{1, \gamma_0 + \tilde{\gamma}\}$] Having a second look at
 284 Theorem 2, we see that the r.h.s. is a function of γ of the form

$$U : \gamma \mapsto \frac{1}{1 - \gamma_{eval}} + \frac{1}{\gamma - 1} + C_{m,T,S,A,\sigma,\delta} \frac{\gamma}{(1 - \gamma)^2},$$

285 where the first term is positive and monotonically decreasing on $(0, \gamma_{\text{eval}})$ and the second term is
 286 positive and monotonically increasing on $(0, 1)$. We simplify and scale this constant, keeping only
 287 problem-related terms: $C_t = (\frac{1}{\sqrt{t}}(\sigma + \frac{1}{\sqrt{m}}))/(\sigma^2 m + 1) + \sigma^2 m \frac{1}{\sqrt{m}}/(\sigma^2 m + 1)$, which is of the
 288 order of the constant in equation 6. Optimizing γ by using the function U with constant C does not
 289 lead to a principled analytical value strictly speaking because U is derived from an upper bound that
 290 may be loose and may not reflect the true shape of the loss w.r.t. γ , but we may use the resulting
 291 growth schedule to guide our choices online. In general, the existence of a strict minimum for U in
 292 $(0, 1)$ is not always guaranteed: depending on the values of $C \approx C_{m,T,S,A,\sigma}$, the function may be
 293 monotonic and the minimum may be on the edges. We give explicit ranges in the proposition below,
 294 proved in Appendix G.

295 **Proposition 1.** *The existence of a strict minimum in $(0, 1)$ is determined by $C = C_{m,T,S,A,\sigma,\delta}$ (which
 296 can be computed) as follows:*

$$\tilde{\gamma} = \begin{cases} 0 & \text{if } C \geq 1 \\ 1 & \text{if } C < 1/2 \\ \frac{1-C}{1+C} & \text{otherwise, i.e if } 1/2 < C < 1 \end{cases}$$

297

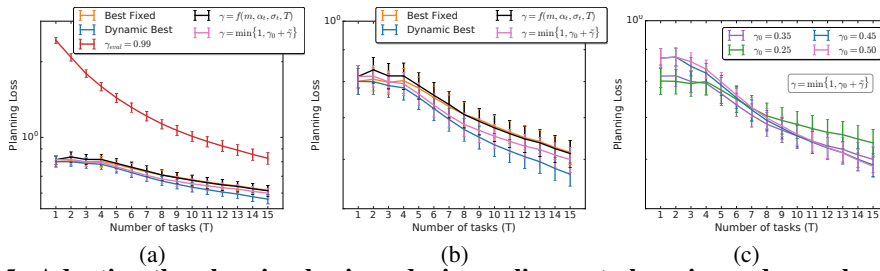


Figure 5: **Adapting the planning horizon during online meta-learning reduces planning loss.**

(a) Planning with online-meta learning shows that *all* baselines outperform using a constant discount factor. (b) Zoomed in plot of average planning loss over the progression of tasks T shows competitive performance with the proposed schedule of $\gamma = f(m, \alpha_t, \sigma_t, T)$ beating best-fixed as more tasks are seen. The γ schedule $\gamma = \min\{1, \gamma_0 + \tilde{\gamma}\}$ beats the best-fixed and is very competitive to the dynamic-best baseline. (c) Using the upper bound to guide the schedule significantly outperforms γ_{eval} and is shown for $\gamma_0 \in (0.25, 0.50)$. We plot 1-standard error for 600 independent runs.

298 **Empirical Validation** We consider the setup described in Sec. 5 with 15 tasks in a 10-state, 2-action
 299 random MDP distribution of tasks with $\sigma \approx 0.01$. In Fig. 5, we plot the planning loss obtained by
 300 POMRL with our schedules, a fixed γ_{eval} and two strong baselines: *best fixed* which considers the
 301 best fixed value of discount over all tasks estimated in hindsight and *dynamic best* which considers
 302 the best choice if we had used the optimal γ^* in each round as in Fig. 4(c). It is important to note
 303 that *dynamic best* is a lower bound that we cannot outperform. We observe in Fig. 5(a) that γ_{eval}
 304 results in a very high loss, potentially corresponding to trying to plan too far ahead despite model
 305 uncertainty. Upon inspecting Fig. 5(b), we observe that the proposed $\gamma = f(m, \alpha_t, \sigma_t, T)$ obtains
 306 similar performance to *best fixed* and is within the significance range of the lower bound. Our second
 307 heuristic, $\gamma = \min\{1, \gamma_0 + \tilde{\gamma}\}$ obtains similarly good performance, as seen in Fig. 5(b). Fig. 5(c)
 308 shows the effect of different values of γ_0 in the prescribed range. These results provide evidence that
 309 it is possible to adapt the planning horizon as a function of the problem’s structure (meta-learned
 310 task-similarity) and sample sizes. Adapting the planning horizon online is an open problem and
 311 beyond the scope of our work.

312 7 Conclusion

313 We presented connections between planning with inaccurate models and online meta-learning via a
 314 high-probability task-averaged regret upper-bound on the planning loss that primarily depends on
 315 task-similarity σ as opposed to the entire search space Σ . Algorithmically, we demonstrate that the
 316 agent can use its experience in each task *and* across tasks to estimate both the transition model and
 317 the distribution over tasks. Meta-learning the underlying task similarity and a good initialization of
 318 transition model across tasks enables longer planning horizons. See Appendix Sec. H.6 for extended
 319 discussion.

References

- 320
- 321 Amit, R., Meir, R., and Ciosek, K. (2020). Discount factor as a regularizer in reinforcement learning.
322 In *International conference on machine learning*, pages 269–278. PMLR.
- 323 Arumugam, D., Abel, D., Asadi, K., Gopalan, N., Grimm, C., Lee, J. K., Lehnert, L., and Littman,
324 M. L. (2018). Mitigating planner overfitting in model-based reinforcement learning. *arXiv preprint*
325 *arXiv:1812.01129*.
- 326 Balcan, M.-F., Khodak, M., and Talwalkar, A. (2019). Provable guarantees for gradient-based
327 meta-learning. In *International Conference on Machine Learning*, pages 424–433. PMLR.
- 328 Baxter, J. (2000). A model of inductive bias learning. *Journal of artificial intelligence research*,
329 12:149–198.
- 330 Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. (2013). The arcade learning environment:
331 An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279.
- 332 Blackwell, D. (1962). Discrete dynamic programming. *The Annals of Mathematical Statistics*, pages
333 719–726.
- 334 Caruana, R. (1997). Multitask learning. *Machine learning*, 28(1):41–75.
- 335 Cella, L., Lazaric, A., and Pontil, M. (2020). Meta-learning with stochastic linear bandits. In
336 *International Conference on Machine Learning*, pages 1360–1370. PMLR.
- 337 Denevi, G., Ciliberto, C., Stamos, D., and Pontil, M. (2018). Learning to learn around a common
338 mean. *Advances in Neural Information Processing Systems*, 31.
- 339 Denevi, G., Stamos, D., Ciliberto, C., and Pontil, M. (2019). Online-within-online meta-learning.
340 *Advances in Neural Information Processing Systems*, 32.
- 341 Dong, S., Van Roy, B., and Zhou, Z. (2021). Simple agent, complex environment: Efficient
342 reinforcement learning with agent state. *arXiv preprint arXiv:2102.05261*.
- 343 Fedus, W., Gelada, C., Bengio, Y., Bellemare, M. G., and Larochelle, H. (2019). Hyperbolic
344 discounting and learning over multiple horizons. *arXiv preprint arXiv:1902.06865*.
- 345 Finn, C., Abbeel, P., and Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep
346 networks. In *International Conference on Machine Learning*, pages 1126–1135. PMLR.
- 347 Finn, C., Rajeswaran, A., Kakade, S., and Levine, S. (2019). Online meta-learning. In *International*
348 *Conference on Machine Learning*, pages 1920–1930. PMLR.
- 349 Flennerhag, S., Schroecker, Y., Zahavy, T., van Hasselt, H., Silver, D., and Singh, S. (2021). Boot-
350 strapped meta-learning. *arXiv preprint arXiv:2109.04504*.
- 351 Flennerhag, S., Zahavy, T., O’Donoghue, B., van Hasselt, H., György, A., and Singh, S. (2022).
352 Optimistic meta-gradients. In *Sixth Workshop on Meta-Learning at the Conference on Neural*
353 *Information Processing Systems*.
- 354 François-Lavet, V., Rabusseau, G., Pineau, J., Ernst, D., and Fonteneau, R. (2019). On overfitting and
355 asymptotic bias in batch reinforcement learning with partial observability. *Journal of Artificial*
356 *Intelligence Research*, 65:1–30.
- 357 Jiang, N., Kulesza, A., Singh, S., and Lewis, R. (2015). The dependence of effective planning horizon
358 on model accuracy. In *Proceedings of the 2015 International Conference on Autonomous Agents*
359 *and Multiagent Systems*, pages 1181–1189. International Foundation for Autonomous Agents and
360 Multiagent Systems.
- 361 Kearns, M., Mansour, Y., and Ng, A. Y. (2002). A sparse sampling algorithm for near-optimal
362 planning in large markov decision processes. *Machine learning*, 49(2):193–208.
- 363 Khodak, M., Balcan, M.-F., and Talwalkar, A. (2019). Adaptive gradient-based meta-learning
364 methods. *arXiv preprint arXiv:1906.02717*.

- 365 Kocsis, L. and Szepesvári, C. (2006). Bandit based monte-carlo planning. In *European conference*
366 *on machine learning*, pages 282–293. Springer.
- 367 Luketina, J., Flennerhag, S., Schroecker, Y., Abel, D., Zahavy, T., and Singh, S. (2022). Meta-
368 gradients in non-stationary environments. In *ICLR Workshop on Agent Learning in Open-*
369 *Endedness*.
- 370 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A.,
371 Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. (2015). Human-level control through deep
372 reinforcement learning. *Nature*, 518(7540):529.
- 373 Müller, R. and Pacchiano, A. (2022). Meta learning mdps with linear transition models. In *International*
374 *Conference on Artificial Intelligence and Statistics*, pages 5928–5948. PMLR.
- 375 Oh, J., Hessel, M., Czarnecki, W. M., Xu, Z., van Hasselt, H., Singh, S., and Silver, D. (2020).
376 Discovering reinforcement learning algorithms. *arXiv preprint arXiv:2007.08794*.
- 377 Petrik, M. and Scherrer, B. (2008). Biasing approximate dynamic programming with a lower discount
378 factor. *Advances in neural information processing systems*, 21.
- 379 Pineau, J. (2019). The machine learning reproducibility checklist. *arxiv*.
- 380 Schmidhuber, J. and Huber, R. (1991). Learning to generate artificial fovea trajectories for target
381 detection. *International Journal of Neural Systems*, 2(01n02):125–134.
- 382 Tao, T. and Vu, V. (2015). Random matrices: universality of local spectral statistics of non-hermitian
383 matrices. *The Annals of Probability*, 43(2):782–874.
- 384 Thrun, S. and Pratt, L. (1998). Learning to learn: Introduction and overview. In *Learning to learn*,
385 pages 3–17. Springer.
- 386 Van Seijen, H., Van Hasselt, H., Whiteson, S., and Wiering, M. (2009). A theoretical and empirical
387 analysis of expected sarsa. In *2009 IEEE Symposium on Adaptive Dynamic Programming and*
388 *Reinforcement Learning*, pages 177–184. IEEE.
- 389 Xu, Z., van Hasselt, H., and Silver, D. (2018). Meta-gradient reinforcement learning. *arXiv preprint*
390 *arXiv:1805.09801*.
- 391 Zahavy, T., Xu, Z., Veeriah, V., Hessel, M., Oh, J., van Hasselt, H. P., Silver, D., and Singh, S. (2020).
392 A self-tuning actor-critic algorithm. *Advances in neural information processing systems*, 33.

393 **A How does performance depend on the amount of shared structure across**
 394 **tasks *i.e.*, σ ?**

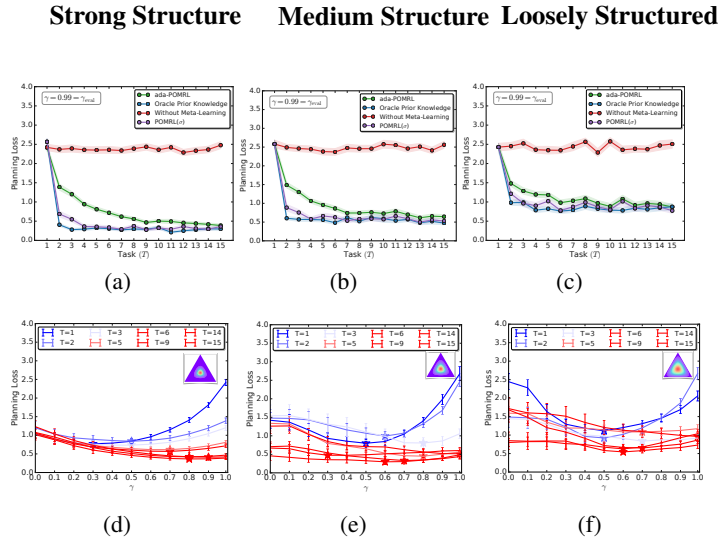


Figure A1: POMRL and ada-POMRL are robust to varying task-similarity σ for a small fixed amount of data $m = 5$ available at each round $t \in T$. A small value of σ reflects the fact that tasks are closely related to each other and share a good amount of structure whereas a much larger value indicates loosely related tasks (simplex plots illustrate the meta-distribution in dimension 2). In the former case, meta-learning the shared structure alongside a good model initialization leads to most gains. In the latter, the learner struggles to cope with new unseen tasks which differ significantly. Error bars represent 1-standard deviation of uncertainty across 100 independent runs.

395 We let σ vary to cover the **three regimes**: $\sigma \approx 0.01$ corresponding to fast convergence, $\sigma = 0.025$
 396 is in the intermediate regime (needs longer T), and $\sigma = 0.047$ is the loosely structured case where
 397 we don't expect much meta-learning to help improve model accuracy. The small inset figures in
 398 Fig. A1 represent the task distribution in the simplex. In all cases, ada-POMRL estimates σ online
 399 and we report the planning losses for a range of γ 's. Inspecting Fig. A1, we observe that while in the
 400 presence of closely related tasks (Fig. 1(a)) all methods perform well (except without meta-learning).
 401 As the underlying task relatedness decreases (for intermediate regime in Fig. 1(b)), both POMRL and
 402 ada-POMRL remain consistent in their performance as compared to the Oracle Prior Knowledge
 403 baseline. When the underlying tasks are loosely related (as in Fig. 1(c)), ada-POMRL and POMRL can
 404 still perform well in comparison to other baselines.

405 Next, we report and discuss the planning loss plot for ada-POMRL for the three cases are shown in
 406 Figures 1(d), 1(e), and 1(f) respectively. An intermediate value of task-similarity (Fig. 1(e)) still leads
 407 to gains, albeit at a lower speed of convergence. In contrast, a large value of $\sigma = 0.047$ indicates
 408 little relatedness across tasks resulting in minimal gains from meta-learning here as seen in Fig. 1(f).
 409 The learner struggles to learn a good initialization of the model dynamics as there is no natural one.
 410 All planning loss curves remain U-shaped and overall higher with an intermediate optimal guidance
 411 γ value (0.5). However, ada-POMRL does not do worse overall than the initial run $T = 1$, meaning
 412 that while there is not a significant improvement, our method does not hurt performance in loosely
 413 related tasks⁵. Recall that ada-POMRL has no apriori knowledge of the number of tasks (T), or the
 414 underlying task relatedness (σ) *i.e.*, adaptation is on-the-fly.

415 **Implications for degree of task-similarity *i.e.*, σ values:** Our bound suggests that the degree of
 416 improvement you can get from meta learning scales with the task similarity σ instead of the set size
 417 Σ . Thus, for $\sigma \leq \Sigma$, performing meta learning with Alg1 guarantees better learning measured via our
 418 improved regret bound when there is underlying structure in the problem space which we formalize

⁵The theoretical bound may lead to think that the average planning loss is higher due to the introduced bias, but in practice we do not observe that, which means our bound is pessimistic on the second order terms.

419 through Eq. 3. Should σ be large, the techniques will still hold and our bounds will simply scale
 420 accordingly.

421 **When $\sigma = 0$, all tasks are exactly the same.** Indeed, the mixing rate $\alpha_t \approx 1$ for all t , so our
 422 algorithm boils down to returning the average of means $\hat{P}^{o,t}$ for each task, which simply corresponds
 423 to solving the tasks as a continuous, uninterrupted stream of batches from the nearly same model
 424 that $\hat{P}^{o,t}$ aggregates. Unsurprisingly, our bound recovers that of (Jiang et al., 2015, Theorem 1):
 425 the bound below reflects that we have to estimate only one model in a space of “size” Σ with mT
 426 samples.

$$\bar{\mathcal{L}} \leq \frac{\gamma_{\text{eval}} - \gamma}{(1 - \gamma_{\text{eval}})(1 - \gamma)} + \frac{2\gamma S}{(1 - \gamma)^2} \tilde{O} \left(\sqrt{\frac{\Sigma}{mT}} \right) \quad (7)$$

427 **When $\sigma = 1$, then $\sigma = \Sigma = 1$, then the meta-learning assumption is not relevant but our**
 428 **bound remains valid and gracefully degrades to reflect it.** We need to estimate T models each
 429 with m samples. Then the second term $\frac{1}{\sqrt{m}}$ reflects the usual estimation error for each task while the
 430 first term is an added bias (second order in $\frac{1}{m}$) due to our regularization to our mean prior P^o that is
 431 not relevant here.

$$\bar{\mathcal{L}} \leq \frac{\gamma_{\text{eval}} - \gamma}{(1 - \gamma_{\text{eval}})(1 - \gamma)} + \frac{2\gamma S}{(1 - \gamma)^2} \tilde{O} \left(\frac{1}{m} \left(1 + \frac{1}{\sqrt{T}} \left(1 + \sqrt{1 + \frac{1}{m}} \right) \right) + \frac{1}{\sqrt{m}} \right) \quad (8)$$

432 **Connections to ARUBA.** As explained earlier, our metric is not directly comparable to that of
 433 ARUBA (Khodak et al., 2019) but it is interesting to make a parallel with the high-probability average
 434 regret bounds proved in their Theorem 5.1. They also obtain an upper bound in $\tilde{O}(1/\sqrt{m} + 1/\sqrt{mT})$
 435 if one upper bounds their average within-task regret $\bar{U} \leq B\sqrt{m}$.

436 **Remark 1** (Role of the task similarity σ in Eq. 2). *When $\sigma > 0$, POMRL naturally integrates each
 437 new data batch into the model estimation. The knowledge of σ is necessary to obtain this exact and
 438 intuitive update rule, and our theory only covers POMRL equipped with this prior knowledge, but we
 439 discuss how to learn and plug-in $\hat{\sigma}_t$ in practice. Note that it would be possible to extend our result
 440 to allow for using the empirical variance estimator with tools like the Bernstein inequality, but we
 441 believe this is out of the scope of this work as it would essentially give a similar bound as obtained in
 442 Theorem 2 with an additional lower order term in $O(1/T)$, and it would not provide much further
 443 intuition on the meta-planning problem we study.*

444 B Additional Detailed Background

445 **Reinforcement Learning.** We consider tabular Markov Decision Processes (MDPs) $\mathcal{M} =$
 446 $\langle \mathcal{S}, \mathcal{A}, R, P, \gamma_{\text{eval}} \rangle$, where \mathcal{S} is a finite set of states, \mathcal{A} is a finite set of actions and we denote
 447 the set cardinalities as $S = |\mathcal{S}|$ and $A = |\mathcal{A}|$. For each state $s \in \mathcal{S}$, and for each available action
 448 $a \in \mathcal{A}$, the probability vector $P(\cdot | s, a)$ defines a transition model over the state space and is a
 449 probability distribution in a set of feasible models $\mathcal{D}_P \subset \Delta_S$, where Δ_S the probability simplex of
 450 dimension $S - 1$. We denote $\Sigma \leq 1$ the diameter of \mathcal{D}_P . A policy is a function $\pi : \mathcal{S} \rightarrow \mathcal{A}$ and it
 451 characterizes the agent’s behavior.

452 We consider the bounded reward setting, *i.e.*, $R \in [0, R_{\text{max}}]$ and without loss of generality we set
 453 $R_{\text{max}} = 1$ (unless stated otherwise). Given an MDP, or task, M , for any policy π , let $V_{M,\gamma}^\pi \in \mathbb{R}^S$ be
 454 the value function when evaluated in MDP M with discount factor $\gamma \in (0, 1)$ (potentially different
 455 from γ_{eval}); defined as $V_{M,\gamma}^\pi(s) = \mathbb{E} \sum_{t=0}^{\infty} (\gamma^t R_{s_t} | s_0 = s)$. The goal of the agent is to find an
 456 optimal policy, $\pi_{M,\gamma}^* = \arg \max_{\pi} \mathbf{E}_{s \sim \rho} V_{M,\gamma}^\pi(s)$ where $\rho > 0$ is any positive measure, denoted π^*
 457 when there is no ambiguity. For given state and action spaces and reward function $(\mathcal{S}, \mathcal{A}, R)$, we
 458 denote Π_γ the set of *potentially* optimal policies for discount factor γ : $\Pi_\gamma = \{\pi | \exists P \text{ s.t. } \pi =$
 459 $\pi_{M,\gamma}^* \text{ where } M = \langle \mathcal{S}, \mathcal{A}, R, P, \gamma \rangle\}$. We use Big-O notation, $O(\cdot)$ and $\tilde{O}(\cdot)$, to hide respectively
 460 universal constants and poly-logarithmic terms in T, S, A and $\delta > 0$ (the confidence level).

461 **Model-based Reinforcement Learning.** In practice, the true model of the world is unknown and
 462 must be estimated from data. One approach to approximately solve the optimization problem
 463 above is to construct a model, (\hat{R}, \hat{P}) from data, then find $\pi_{\hat{M},\gamma}^*$ for the corresponding MDP $\hat{M} =$
 464 $\langle \mathcal{S}, \mathcal{A}, \hat{R}, \hat{P}, \gamma \rangle$. This approach is called *model-based RL* or *certainty-equivalence (CE) control*.

465 **Planning with inaccurate models.** In this setting, Jiang et al. (2015) define the planning loss as the
 466 gap in expected return in MDP M when using $\gamma \leq \gamma_{\text{eval}}$ and the optimal policy for an approximate
 467 model \hat{M} :

$$\mathcal{L}(\hat{M}, \gamma | M, \gamma_{\text{eval}}) = \|V_{M, \gamma_{\text{eval}}}^{\pi_{\hat{M}, \gamma_{\text{eval}}}^*} - V_{M, \gamma_{\text{eval}}}^{\pi_{\hat{M}, \gamma}^*}\|_{\infty}.$$

468 Thus, the **optimal effective planning horizon** $(1 - \gamma^*)^{-1}$ is defined using the discount factor that
 469 minimizes the planning loss, i.e., $\gamma^* := \min_{0 \leq \gamma \leq \gamma_{\text{eval}}} \mathcal{L}(\hat{M}, \gamma | M, \gamma_{\text{eval}})$.

470 **Theorem 1.** (Jiang et al. (2015)) Let M be an MDP with non-negative bounded rewards and
 471 evaluation discount factor γ_{eval} . Let \hat{M} be the approximate MDP comprising the true reward
 472 function of M and the approximate transition model \hat{P} , estimated from $m > 0$ samples for each
 473 state-action pair. Then, with probability at least $1 - \delta$,

$$\|V_{M, \gamma_{\text{eval}}}^{\pi_{\hat{M}, \gamma_{\text{eval}}}^*} - V_{M, \gamma_{\text{eval}}}^{\pi_{\hat{M}, \gamma}^*}\|_{\infty} \leq \frac{\gamma_{\text{eval}} - \gamma}{(1 - \gamma_{\text{eval}})(1 - \gamma)} + \frac{2\gamma R_{\text{max}}}{(1 - \gamma)^2} \left(\sqrt{\frac{\Sigma}{2m} \log \frac{2SA|\Pi_{\gamma}|}{\delta}} \right) \quad (9)$$

474 where Σ is upper-bounded by 1 as $P, \hat{P} \in \Delta_S$.

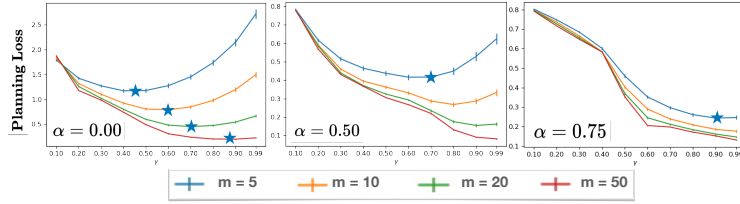


Figure B2: **On the role of incorporating a ground truth prior of transition model on planning horizon.** The planning loss is a function of the discount factor γ and is impacted by incorporating prior knowledge. The learner has $m = 5, 10, 20, 50$ samples per task to estimate the model, corresponding to the curves in each sub figure. Inspecting any of the sub figures, we observe that larger values of m lead to lower planning loss and a larger effective discount factor. Besides, inspecting one value of m across tasks (e.g., $m = 5$), we see that the same effect (lower planning loss and larger effective discount) occurs when the learner puts more weight on the ground truth through α .

475 This result holds for a count-based model estimator (i.e, empirical average of observed transitions)
 476 given by a generator model for each pair (s, a) . It gives an upper-bound on the planning loss as
 477 a function of the guidance discount factor $\gamma < 1$. The result decomposes the loss into two terms:
 478 the constant bias which decreases as γ tends to γ_{eval} , and the variance (or uncertainty) term which
 479 increases with γ but decreases as $1/\sqrt{m}$. As $m \rightarrow \infty$ that second factor vanishes, but in the
 480 low-sample regime the optimal effective planning horizon should trade-off both terms.

481 **Illustration.** These effects are illustrated in Fig. B2 on a simple 10-state, 2-action random MDP.
 482 The leftmost plot uses the simple count-based model estimator and reproduces the results from
 483 Jiang et al. (2015). We then incorporate the true prior (mean model P^o as in Fig 1 and defined
 484 above Eq. 3 in Assumption 1) in the estimator with a growing mixing factor $\alpha \in (0, 1)$: $\hat{P}(m) =$
 485 $\alpha P^o + (1 - \alpha) \frac{\sum^i X^i}{m}$. We observe that increasing the weight $\alpha \in (0, 1)$ on good prior knowledge
 486 enables longer planning horizons and lower planning loss.

487 **Online Meta-Learning and Regret.** We consider an online meta-RL problem where an agent is pre-
 488 sented with a sequence of tasks M_1, M_2, \dots, M_T , where for each $t \in [T]$, $M_t = \langle \mathcal{S}, \mathcal{A}, P^t, R, \gamma_{\text{eval}} \rangle$,
 489 that is, the MDPs only differ from each other by the transition matrix (dynamics model) P^t . The
 490 learner must sequentially estimate the model \hat{P}^t for each task t from a batch of m transitions simulated
 491 for each state-action pair⁶.

492 Its goal is to minimize the average planning loss also expressed in the form of task averaged regret
 493 suffered in planning and defined as

$$\bar{\mathcal{L}}(\hat{M}_{1:T}, \gamma | M_{1:T}, \gamma_{\text{eval}}) = \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\hat{M}_t, \gamma | M_t, \gamma_{\text{eval}}) = \frac{1}{T} \sum_{t=1}^T \|V_{M_t, \gamma_{\text{eval}}}^{\pi_{\hat{M}_t, \gamma_{\text{eval}}}^*} - V_{M_t, \gamma_{\text{eval}}}^{\pi_{\hat{M}_t, \gamma}^*}\|_{\infty} \quad (10)$$

⁶So a total of mSA samples.

494 Note that the reference MDP for each term is the true M_t , and the discount factor γ is the same in
 495 all tasks. One can see this objective as a stochastic dynamic regret: at each task $t \in [T]$, the learner
 496 competes against the optimal policy for the *current* true model, as opposed to competing against the
 497 best fixed policy in hindsight used in classical definitions of regret.

498 Note that **our dynamic regret is different from the one considered in ARUBA (Khodak et al.,**
 499 **2019)**. They consider the fully online setting where the data is observed as an arbitrary stream within
 500 each task, and each comparator is simply the minimum of the within-task loss in hindsight. In our
 501 model, however, given access to a simulator (See Sec. 2) allows us to get i.i.d transition samples as a
 502 batch at the beginning of each task, and consequently we define our regret with respect to the true
 503 generating parameter. One key consequence of this difference is that their regret bounds cannot be
 504 directly applied to our setting, and we prove new results further below.

505 C Additional Related Work

506 **Discount Factor Adaptation.** For almost all real-world applications, RL agents operate in a much
 507 larger environment than the agent capacity in the context of both the computational and memory
 508 complexity (e.g. the internet). Inevitably, it becomes crucial to adapt the planning horizon over time
 509 as opposed to using a relatively longer planning horizon from the start (which can be both expensive
 510 and sub-optimal). This has been extensively studied in the context of planning with inaccurate models
 511 in reinforcement learning (Jiang et al., 2015; Arumugam et al., 2018).

512 Dong et al. (2021) introduced a schedule for γ that we take inspiration from in Section 6. They
 513 consider a 'never-ending RL' problem in the infinite-horizon, average-regret setting in which the
 514 true horizon is 1, but show that adopting a different smaller discount value proportional to the time
 515 in the agent's life results in significant gains. Their focus and contributions are different from ours
 516 as they are interested in asymptotic rates, but we believe the connection between our findings is an
 517 interesting avenue for future research.

518 **Meta-Learning and Meta-RL**, or *learning-to-learn* has shown tremendous success in online dis-
 519 covery of different aspects of an RL algorithm, ranging from hyper-parameters (Xu et al., 2018) to
 520 complete objective functions (Oh et al., 2020). In recent years, many deep RL agents (Fedus et al.,
 521 2019; Zahavy et al., 2020) have gradually used higher discounts moving away from the traditional
 522 approach of using a fixed discount factor. However, to the best of our knowledge, existing works do
 523 not provide a formal understanding of why this is helping the agents in better performance, especially
 524 across varied tasks. Our analysis is motivated by the aforementioned empirical success of adapting the
 525 discount factor. While there has been significant progress in meta-learning-inspired meta-gradients
 526 techniques in RL (Xu et al., 2018; Zahavy et al., 2020; Flennerhag et al., 2021), they are largely
 527 focused on *empirical analysis* with lot or room for in-depth insights about the source of underlying
 528 gains.

529 D Closed-form solution of the regularized least squares

530 We note that each \hat{P} should be understood as $\hat{P}_{(s,a)}(s')$.

$$\nabla \ell(P|h) = -\frac{2}{m} \sum_{i=1}^m (X^i - P) + 2\lambda(P - h)$$

$$\nabla \ell(P|h) = 0 \iff P(1 + \lambda) = \frac{\sum_i X^i}{m} + \lambda h$$

$$\hat{P}_{(s,a)}(s'|h) = \frac{1}{1 + \lambda} \frac{\sum_i X^i}{m} + \frac{\lambda}{1 + \lambda} h \tag{11}$$

$$= \alpha h + (1 - \alpha) \frac{\sum_i X^i}{m} \quad \text{where } \alpha = \frac{\lambda}{1 + \lambda} \tag{12}$$

531 **Derivation of Mixing Rate α_t :** To choose α_t , we want to minimize the MSE of the final estimator.

$$\begin{aligned}\mathbb{E}_{X \sim P^t} \left(\hat{P}^t - P^t \right)^2 &= \mathbb{E}_{X \sim P^t} \left(\alpha_t h_t + (1 - \alpha_t) \bar{P}^t - P^t \right)^2 \\ &= \mathbb{E}_{X \sim P^t} \left(\alpha_t (h_t - P^t) + (1 - \alpha_t) (\bar{P}^t - P^t) \right)^2 \\ &= \alpha_t^2 (h_t - P^t)^2 + (1 - \alpha_t)^2 \mathbb{E}_{X \sim P^t} \left((\bar{P}^t - P^t) \right)^2\end{aligned}$$

532 where the cross term $2\alpha_t t(h_t - P^t)(1 - \alpha_t) \mathbb{E}_{X \sim P^t} E[\bar{P}^t - P^t] = 0$ since $\mathbb{E}[\bar{P}^t] = P^t$. This is the
533 classic bias-variance decomposition of an estimator and we see that the choice of h_t plays a role as
534 well as the variance of \bar{P}^t , which is upper bounded by $1/m$ (because each $X^{i,t}$ is bounded in $(0, 1)$).
535 For instance, for the choice $h_t = P^o$, by our structural assumption 3 we get:

$$\mathbb{E}_{X \sim P^t} \left(\hat{P}^t - P^t \right)^2 \leq \alpha^2 \sigma^2 + (1 - \alpha)^2 \frac{1}{m},$$

536 and we minimize this upper bound in α to obtain the mixing coefficient with smallest MSE: $\alpha^* =$
537 $\frac{1}{\sigma^2 m + 1}$, or equivalently $\lambda^* = \frac{1}{\sigma^2 m}$. Recall this is the within-task estimator's variance where we
538 consider the true P^o .

539 In practice, however, we meta-learn the prior, so for $t > 1$, $h_t = \hat{P}^{o,t} = \frac{1}{t-1} \sum_{j=1}^{t-1} \bar{P}^j$. Intuitively,
540 as m and t grow large, $\hat{P}^{o,t} \rightarrow P^o$ and we retrieve the result above (we show this formally to prove
541 Eq. 22 in the proof of our main theorem). To obtain a simple expression for α_t , we minimize the
542 "meta-MSE" of our estimator:

$$\begin{aligned}\mathbb{E}_{P^t \sim P^o} \left(\hat{P}^t - P^t \right)^2 &= \alpha_t^2 \mathbb{E}_{P^t \sim P^o} \mathbb{E}_{X \sim P^t} (h_t - P^t)^2 + (1 - \alpha_t)^2 \mathbb{E}_{P^t \sim P^o} \mathbb{E}_{X \sim P^t} \left((\bar{P}^t - P^t) \right)^2 \\ &\leq \alpha_t^2 \mathbb{E}_{P^t \sim P^o} \left(\frac{1}{t-1} \sum_{j=1}^{t-1} P^j - P^o + P^o - P^t \right)^2 + (1 - \alpha_t)^2 \frac{1}{m} \\ &\leq \alpha_t^2 \sigma^2 (1 + 1/t) + (1 - \alpha_t)^2 \frac{1}{m},\end{aligned}$$

543 where in the last inequality, we upper bounded the variance of $\frac{1}{t-1} \sum_{j=1}^{t-1} P^j$ (the "denoised" $\hat{P}^{o,t}$)
544 by σ^2/t since each P^t is bounded in $[P^o - \sigma, P^o + \sigma]$ by our structural assumption. Minimizing
545 that last upper bound in α_t leads to $\alpha_t = \frac{1}{(\sigma^2)(1+1/t)m+1} \xrightarrow{t \rightarrow \infty} \alpha^*$, when $t \rightarrow \infty$. This means that the
546 uncertainty on the prior implies that its weight in the estimator is smaller, but eventually converges at
547 a fast rate to the optimal value (when the exact optimal prior is known). This inequality holds with
548 probability $1 - \delta$ because we use the concentration of $\hat{P}^{o,t}$ (see proof of Theorem 19 below)

549 E Online Estimation

550 **Online Estimation of Prior.** At each task, the learner gets m interactions per state-action pair. At
551 task $t = 1$, learner can compute the prior based on the samples seen so far, i.e.:

$$\hat{P}_o^{t=1}(s'|s, a) = \frac{\{\sum_{i=1}^m X_i\}_{t=1}}{m}$$

552 At subsequent tasks,

$$\begin{aligned}\hat{P}_o^{t=2}(s'|s, a) &= \frac{\{\sum_{i=1}^{2m} X_i\}_{t=1:2}}{2m} = \frac{1}{2} \left(\frac{\{\sum_{i=1}^m X_i\}}{m} + \frac{\{\sum_{i=m+1}^{2m} X_i\}}{m} \right) \\ &= \frac{1}{2} \left(\hat{P}_o^{t=1}(s'|s, a) + \frac{\{\sum_{i=m+1}^{2m} X_i\}}{m} \right)\end{aligned}$$

553 Similarly,

$$\begin{aligned}\hat{P}_o^{t=3}(s'|s, a) &= \frac{\{\sum_{i=1}^{3m} X_i\}_{t=1:3}}{3m} = \frac{\{\sum_{i=1}^{2m} X_i\}_{t=1:2} + \{\sum_{i=2m+1}^{3m} X_i\}_{t=3}}{3m} \\ &= \frac{1}{3} \left(\frac{2\{\sum_{i=1}^{2m} X_i\}_{t=1:2}}{2m} + \frac{\{\sum_{i=2m+1}^{3m} X_i\}_{t=3}}{m} \right) \\ \implies \hat{P}_o^t(s'|s, a) &= \frac{1}{t} \left((t-1)\hat{P}_o^{t-1}(s'|s, a) + \frac{\sum_{i=(t-1)m+1}^{tm} X_i}{m} \right)\end{aligned}$$

554 Therefore,

$$\hat{P}_o^t(s'|s, a) = \left(1 - \frac{1}{t}\right)\hat{P}_o^{t-1}(s'|s, a) + \left(\frac{1}{t}\right)\frac{\sum_{i=(t-1)m+1}^{tm} X_i}{m} \quad (13)$$

555 **Online Estimation of Variance.** Similarly, we can derive the online estimate of the variance:

$$(\hat{\sigma}_o^2)^t = (\hat{\sigma}_o^2)^{t-1} + \frac{(X_{mt} - \hat{P}_o^{t-1})(X_{mt} - \hat{P}_o^t) - (\hat{\sigma}_o^2)^{t-1}}{t} \quad (14)$$

556 Since the above method is numerically unstable, we will employ Welford's online algorithm for
557 variance estimate.

558 F Concentration bounds and Proof of Theorem 2

559 F.1 Proof of Theorem 2

560 We begin the proof by decomposing each term of the loss:

561 **Lemma 1.** For a task t denoted by M , and its estimate denoted by $\hat{M}, \forall s \in S$,

$$V_{P^t, \gamma_{eval}}^{\pi_{P^t, \gamma_{eval}}} (s) - V_{P^t, \gamma}^{\pi_{\hat{P}^t, \gamma}} (s) = \underbrace{\left(V_{P^t, \gamma_{eval}}^{\pi_{P^t, \gamma_{eval}}} (s) - V_{P^t, \gamma}^{\pi_{P^t, \gamma_{eval}}} (s) \right)}_{A_t} + \underbrace{\left(V_{P^t, \gamma}^{\pi_{P^t, \gamma_{eval}}} (s) - V_{P^t, \gamma}^{\pi_{\hat{P}^t, \gamma}} (s) \right)}_{B_t}$$

562 We are going to bound each term separately. The term (A_t) corresponds to the bias constant due to
563 using γ instead of γ_{eval} and was already bounded by [Jiang et al. \(2015\)](#):

564 **Lemma 2.** [Jiang et al. \(2015\)](#) For any MDP \hat{M} with rewards in $[0, R_{max}]$, $\forall \pi : S \rightarrow A$ and
565 $\gamma \leq \gamma_{eval}$,

$$V_{P^t, \gamma}^{\pi} \leq V_{P^t, \gamma_{eval}}^{\pi} \leq V_{P^t, \gamma}^{\pi} + \frac{\gamma_{eval} - \gamma}{(1 - \gamma_{eval})(1 - \gamma)} R_{max} \quad (15)$$

566 We denote $C(\gamma) = \frac{\gamma_{eval} - \gamma}{(1 - \gamma_{eval})(1 - \gamma)} R_{max}$ and notice that $\sum_t A_t / T = C(\gamma)$ so that bounds the first
567 part of the average loss.

568 To bound the second term B_t , we first use Lemma 3 (Equation 18) in [Jiang et al. \(2015\)](#) to upper
569 bound

$$V_{P^t, \gamma}^{\pi_{P^t, \gamma_{eval}}} (s) - V_{P^t, \gamma_{eval}}^{\pi_{\hat{P}^t, \gamma}} (s) \leq 2 \max_{s \in S, \pi \in \Pi_{R, \gamma}} |V_{P^t, \gamma}^{\pi_{P^t, \gamma_{eval}}} (s) - V_{\hat{P}^t, \gamma_{eval}}^{\pi_{\hat{P}^t, \gamma}} (s)| \quad (16)$$

$$\leq 2 \max_{\substack{s \in S, a \in \mathcal{A}, \\ \pi \in \Pi_{R, \gamma}}} |Q_{P^t, \gamma}^{\pi_{P^t, \gamma_{eval}}} (s, a) - Q_{\hat{P}^t, \gamma_{eval}}^{\pi_{\hat{P}^t, \gamma}} (s, a)| \quad (17)$$

570 Using Lemma 4 from [Jiang et al. \(2015\)](#) and noticing that in our setting we do not estimate R so $\hat{R} =$
571 $R, Q_{P^t, \gamma}^{\pi} (s, a) = R(s, a) + \gamma \langle P^t(s, a, ;), V_{P^t, \gamma}^{\pi} \rangle$ and $Q_{\hat{P}^t, \gamma}^{\pi} (s, a) = R(s, a) + \gamma \langle \hat{P}^t(s, a, ;), V_{\hat{P}^t, \gamma}^{\pi} \rangle$,
572 we have

$$\max_{\substack{s \in S, a \in \mathcal{A}, \\ \pi \in \Pi_{R, \gamma}}} |Q_{P^t, \gamma}^{\pi_{P^t, \gamma_{eval}}} (s, a) - Q_{\hat{P}^t, \gamma_{eval}}^{\pi_{\hat{P}^t, \gamma}} (s, a)| \leq \frac{1}{(1 - \gamma)} \max_{\substack{s \in S, a \in \mathcal{A}, \\ \pi \in \Pi_{R, \gamma}}} \left| \gamma \langle \hat{P}^t(s, a, ;), V_{P^t, \gamma}^{\pi} \rangle - \gamma \langle P^t(s, a, ;), V_{P^t, \gamma}^{\pi} \rangle \right| \quad (18)$$

573 Notice that we are comparing the value functions of two different MDPs which is non-trivial and we
 574 leverage the result of Jiang et al. (2015). We refer the reader to the proof of Lemma 4 therein for
 575 intermediate steps.

576 Now summing over tasks, we have

$$\begin{aligned}
 \frac{\sum_t (B)_t}{T} &\leq \frac{1}{T} \sum_{t=1}^T \frac{2}{(1-\gamma)} \max_{\substack{s \in \mathcal{S}, a \in \mathcal{A}, \\ \pi \in \Pi_{R,\gamma}}} \left| \gamma \langle \hat{P}^t(s, a, ;), V_{\hat{P}^t, \gamma}^\pi \rangle - \gamma \langle P^t(s, a, ;), V_{P^t, \gamma}^\pi \rangle \right| \\
 &\leq \frac{2\gamma}{(1-\gamma)} \frac{1}{T} \sum_{t=1}^T \max_{\substack{s \in \mathcal{S}, a \in \mathcal{A}, \\ \pi \in \Pi_{R,\gamma}}} \left| \langle \hat{P}^t(s, a, ;), V_{\hat{P}^t, \gamma}^\pi \rangle - \langle P^t(s, a, ;), V_{P^t, \gamma}^\pi \rangle \right| \\
 &\leq \frac{2R_{\max}}{(1-\gamma)} \frac{1}{T} \sum_{t=1}^T \sum_{s' \in [S]} \max_{\substack{s \in \mathcal{S}, a \in \mathcal{A}, \\ \pi \in \Pi_{R,\gamma}}} \left| \hat{P}^t(s, a, s') - P^t(s, a, s') \right| |V_{\hat{P}^t, \gamma}^\pi| \\
 &\leq \frac{2R_{\max}\gamma}{(1-\gamma)^2} \frac{S}{T} \sum_{t=1}^T \max_{s, s' \in \mathcal{S}, a \in \mathcal{A}} \left| \hat{P}^t(s, a, s') - P^t(s, a, s') \right|
 \end{aligned}$$

577 where we upper-bounded the value function by $R_{\max}/(1-\gamma)$ and one sum over \mathcal{S} by $S \times \max_{s' \in \mathcal{S}} \dots$
 578 Note that this step differs from Jiang et al. (2015) and allows us to boil down to an average (worst-
 579 case) estimation error of the transition model. We finally upper bound the r.h.s using Theorem 2
 580 stated and proved below.

581 **Remark 2.** In Jiang et al. (2015), the argument is slightly more direct and involves directly controlling
 582 the deviations of the scalar random variables $R(s, a) + \gamma \langle \hat{P}^t(s, a, ;), V_{\hat{P}^t, \gamma}^\pi \rangle$, arguing that it is
 583 bounded and centered at $Q_{\hat{P}^t, \gamma}^\pi(s, a)$. This approach is followed by taking a union bound over the
 584 policy space $\Pi_{R,\gamma}$ and results in a factor $\log(\Pi_{R,\gamma})$ under the square root. We could have followed
 585 this approach and obtained a similar result but we made the alternative choice above as we believe it
 586 is informative. In our case, this factor is replaced (and upper bounded) by the extra S term. As a result,
 587 we lose the direct dependence on the size of the policy class, which is a function of γ and should play a
 588 role in the bound. In turn, and at the price of this extra looseness, we get a slightly more "exploitable"
 589 bound (see our heuristic for a gamma schedule in Section 6). It is easy and straightforward to adapt
 590 our concentration bound below to directly bound $R(s, a) + \gamma \langle \hat{P}^t(s, a, ;), V_{\hat{P}^t, \gamma}^\pi \rangle - Q_{\hat{P}^t, \gamma}^\pi(s, a)$ as
 591 in Jiang et al. (2015), and one would obtain a similar bound as Eq. equation 6 without the factor S ,
 592 but with an extra $\log(\Pi_{R,\gamma})$.

593 F.2 Concentration of the model estimator

594 To avoid clutter in the notation of this section, we drop the (s, a, s') everywhere, as we did in
 595 Appendix D above. All definitions of \hat{P} and \hat{P}_0 are as stated in the latter section.

596 **Theorem 2.** with probability $1 - \delta$:

$$\begin{aligned}
 \max_{s, a, s'} |\hat{P}^t - P^t| &\leq \frac{1}{\sigma^2 m + 1} \left(\sqrt{\frac{\log(\frac{6T}{\delta}) \log(\frac{TS^2A}{\delta}) (\sigma^2 + \frac{\Sigma \log^2(\frac{6T^2}{\delta})}{m})}{T}} + \sigma \sqrt{\frac{\log(\frac{3TS^2A}{\delta})}{T}} + 2\sigma \right) \\
 &\quad + \frac{\sigma^2 m}{2\sigma^2 m + 1} \sqrt{\frac{\Sigma \log(\frac{3TS^2A}{\delta})}{2m}} \quad (19)
 \end{aligned}$$

597 For any $t \in [T]$, s, a, s' and $\pi \in \Pi_{R,\gamma}$, define $\hat{P}^{t,*} = \alpha_t P^o + (1 - \alpha_t) \bar{P}_m^t$ the optimally regularized
 598 estimator (using the true unknown P^o for each t). We have

$$\begin{aligned}
 \left| \hat{P}^t - P^t \right| &\leq \left| \hat{P}^t - \hat{P}^{t,*} \right| + \left| \hat{P}^{t,*} - P^t \right| \\
 &\leq \underbrace{\alpha_t |\hat{P}^{o,t} - P^o|}_{(A)} + \underbrace{(1 - \alpha_t) |\bar{P}_m^t - P^t|}_{(B)} + \underbrace{\alpha_t |P^o - P^t|}_{\leq 2 \cdot \sigma \text{ by assum.}} \quad (20)
 \end{aligned}$$

599 **Bounding Term A**

600 Substituting the estimator $\hat{P}_t = \alpha_t \frac{1}{m} \sum_i^m X_i + (1 - \alpha_t) \hat{P}_o^t$,

$$\begin{aligned} A &\leq \alpha_t \left(\left| \hat{P}_o^t - \frac{1}{t-1} \sum_{j=1}^{t-1} P^j \right| + \left| \frac{1}{t-1} \sum_{j=1}^{t-1} P_j - P_o \right| \right) \\ &\leq \frac{1}{\sigma^2 m + 1} \left(\underbrace{\left| \hat{P}_o^t - \frac{1}{t-1} \sum_{j=1}^{t-1} P^j \right|}_{A_1} + \underbrace{\left| \frac{1}{t-1} \sum_{j=1}^{t-1} P_j - P_o \right|}_{A_2} \right) \end{aligned}$$

601 where α_t is simply upper bounded by its initial value $\frac{1}{\sigma^2 m + 1}$ and we introduced the *denoised*
602 (expected) average $\frac{1}{t-1} \sum_{j=1}^{t-1} P^j = \mathbb{E}_{P^1, \dots, P^{t-1}} \hat{P}^{o,t}$. Indeed, by assumption, $\mathbb{E}_{P \sim \mathcal{P}} \frac{1}{t-1} \sum_{j=1}^{t-1} P^j =$
603 P^o and the variance of this estimator is bounded by $\sigma^2/(t-1)$ by our structure assumption. This
604 allows to naturally bound A_2 using Hoeffding's inequality for bounded random variables: with
605 probability at least $1 - \delta/3$,

$$\max_{s, a, s'} A_2 \leq \sigma \sqrt{\frac{\log(6S^2 AT/\delta)}{T}} \quad (21)$$

606 We now bound A_1

$$A_1 = \left| \frac{1}{t-1} \sum_j (\bar{P}_m^j - P^j) \right|$$

607 We note here that the first term in A_1 is indeed a martingale $M_t = \sum_{j=1}^{t-1} Z_j$, where $Z_j = \bar{P}_m^j - P^j$,
608 such that each increment is bounded with high probability: for each j , $|Z_j| \leq c_j$ w.p $1 - \frac{\delta}{6}$, where
609 $c_j = \sqrt{\frac{\Sigma}{m} \log(\frac{6T^2}{\delta})}$. Moreover, the differences $|Z_j - Z_{j+1}|$ are also bounded with high probability:

$$|Z_j - Z_{j+1}| \leq |P^j - P^{j+1}| + |\bar{P}^j - \bar{P}^{j+1}| < 2\sigma + 2c_j = D_j = 2 \left(\sigma + \frac{\sqrt{\Sigma} \log(\frac{6T^2}{\delta})}{\sqrt{m}} \right)$$

610 Then by (Tao and Vu, 2015, Prop. 34), for any $\epsilon > 0$,

$$P\left(\left| \frac{M_t}{t-1} \right| \geq \frac{\epsilon}{t-1} \sqrt{\sum_{j=1}^{t-1} D_j^2} \right) \leq 2 \exp(-2\epsilon^2) + \sum_{j=1}^{t-1} \frac{\delta}{6T^2}$$

611 Choosing $\epsilon = \sqrt{\frac{1}{2} \log(\frac{12T}{\delta})}$, we get

$$P\left(\left| \frac{1}{t-1} \sum_{j=1}^{t-1} \bar{P}_m^j - P_j \right| \geq \sqrt{\frac{(\sigma + \frac{\sqrt{\Sigma} \log(\frac{6T^2}{\delta})}{\sqrt{m}})^2 \log(\frac{6T}{\delta})}{T}} \right) \leq \frac{\delta}{6T} + \frac{\delta}{6T} = \frac{\delta}{3T}$$

612 With a union bound as before, we get that with probability at least $1 - \delta/3$,

$$A_1 \leq \sqrt{\frac{\log(\frac{6T}{\delta}) \log(\frac{TS^2 A}{\delta}) (\sigma^2 + \frac{\Sigma \log^2(\frac{6T^2}{\delta})}{m})}{T}} \quad (22)$$

613 because $(\sigma + \sqrt{\frac{\Sigma}{m}})^2 \geq \sigma^2 + \frac{\Sigma}{m}$.

614 By combining equation 22 and equation 21, we get:

$$\max_{s,a,s'} \alpha_t |P^{o,t} - P^o| \leq \frac{1}{\sigma^2 m + 1} \left(\sqrt{\frac{\log(\frac{6T}{\delta}) \log(\frac{TS^2 A}{\delta}) (\sigma^2 + \frac{\log^2(\frac{6T^2}{\delta})}{m})}{T}} + \sigma \sqrt{\frac{\log(3S^2 AT/\delta)}{T}} \right) \quad (23)$$

615 Bounding Term B

616 Term B is simply the concentration of the average of bounded variables $\bar{P}_m^t = \frac{1}{m} \sum_i X_i$, whose
 617 variance is bounded by 1. So by Hoeffding's inequality, and a union bound, with probability at least
 618 $1 - \delta/4$

$$\max_{s,a,s'} |\bar{P}_m^t - P^t| \leq \sqrt{\frac{\Sigma \log(4TS^2 A/\delta)}{2m}}$$

619 To bound term B, it remains to upper bound $1 - \alpha_t$ for all $t \in [T]$:

$$1 - \alpha_t = \frac{\sigma^2(1 + \frac{1}{t})m}{\sigma^2(1 + \frac{1}{t})m + 1} \leq \frac{\sigma^2 m}{2\sigma^2 m + 1}$$

620 We get that with probability $1 - \delta/3$

$$\max_{s,a,s'} (B) \leq \frac{\sigma^2 m}{2\sigma^2 m + 1} \sqrt{\frac{\Sigma \log(3TS^2 A/\delta)}{2m}} \quad (24)$$

621 Combining all bounds

622 To conclude, we combine the bounds on the terms in equation 20, replacing with equation 23, equation
 623 24, and with a union bound, we get that with probability $1 - \delta$,

$$\max_{s,a,s'} |\hat{P}^t - P^t| \leq \frac{1}{\sigma^2 m + 1} \left(\sqrt{\frac{\log(\frac{6T}{\delta}) \log(\frac{TS^2 A}{\delta}) (\sigma^2 + \frac{\Sigma \log^2(\frac{6T^2}{\delta})}{m})}{T}} + \sigma \sqrt{\frac{\log(3S^2 AT/\delta)}{T}} + 2\sigma \right) + \frac{\sigma^2 m}{2\sigma^2 m + 1} \sqrt{\frac{\Sigma \log(3TS^2 A/\delta)}{2m}} \quad (25)$$

624 Discussion

625 The bound has 4 main terms respectively in $\tilde{O}(\sqrt{\frac{1}{mT}})$, $\tilde{O}(\sqrt{\frac{1}{T}})$, $\tilde{O}(\frac{1}{m})$ and $\tilde{O}(\sqrt{\frac{1}{m}})$, all scaled by
 626 some factor depending on σ^2 and m . A first remark is that when m is large and $T = 1$, the last part
 627 in $\tilde{O}(\sqrt{\frac{1}{m}})$ dominates due to the factor $\frac{\sigma^2 m}{\sigma^2 m + 1} \rightarrow 1$, while the coefficient of the first two terms goes
 628 to 0 fast (in $1/(\sigma^2 m)$).

629 G Proof of Proposition 1

630 We study the function U defined by

$$U : \gamma \mapsto \frac{1}{1 - \gamma_{\text{eval}}} + \frac{1}{\gamma - 1} + C_{m,T,S,A,\sigma,\delta} \frac{\gamma}{(1 - \gamma)^2},$$

631 where γ_{eval} is a fixed constant and $C := C_{m,T,S,A,\sigma,\delta}$ is seen as a parameter whose value controls
 632 the general "shape" of the function. We differentiate with respect to γ :

$$\frac{dU}{d\gamma} = -\frac{-C(\gamma + 1) + (1 - \gamma)}{(1 - \gamma)^3}.$$

633 We see that the sign of the derivative is affected by the value of the parameter C :

634 • If $\forall \gamma \in (0, 1)$, $-C(\gamma + 1) + (1 - \gamma) > 0$ then U is monotonically decreasing on $(0, 1)$ and
 635 the minimum is reached for $\gamma = 1$,

$$\forall \gamma \in (0, 1), -C(\gamma + 1) + (1 - \gamma) > 0 \iff -2C + 1 > 0 \iff C < 1/2.$$

636 • Similarly, if C is really large, U may be monotonically increasing on $(0, 1)$:

$$\forall \gamma \in (0, 1), -C(\gamma + 1) + (1 - \gamma) < 0 \iff C \geq 1;$$

637 • Finally, if $C \in (1/2, 2)$, the minimum exists inside $(0, 1)$ and is reached for

$$-C\gamma - C + 1 - \gamma = 0 \iff \gamma = \gamma^* = \frac{1 - C}{1 + C}$$

638 We note that we use these values as a guide. Typically, when $T = 1$ and m is small, the multiplicative
 639 term C is large and the bound is not really informative (concentration has not happened yet), and γ
 640 should be small, potentially close to but not equal to zero. As a heuristic, we propose to simply offset
 641 $\tilde{\gamma}$ by an additional γ_0 such that the guidance discount factor is $\gamma = \min\{1, \gamma_0 + \tilde{\gamma}\}$, where γ_0 should
 642 be reasonably chosen by the practitioner to allow for some short-horizon planning at the beginning of
 643 the interaction. Empirically, $\gamma_0 \in (0.25, 0.50)$ seems reasonable for our random MDP setting as it
 644 corresponds to the empirical minima on Fig 4(b).

645 H Experiments: Implementation Details, Ablations & Additional Results

646 H.1 Implementation details

647 We consider a Dirichlet distribution of tasks such that all tasks $t \in [T]$, $P^t \sim \mathcal{P}$ are centered
 648 at some fixed mean $P^o \in \Delta_S^{S \times A}$ as shown in Figure H3. The mean of the task distribution P^o
 649 is chosen as a sampled random MDP and variance of this distribution is determined such that
 650 $\|P_{s,a}^t - P_{s,a}^o\|_\infty \leq \sigma < 1$. Next, we compute the variance of this distribution $\sigma_i = \frac{\tilde{\alpha}_i(1-\tilde{\alpha}_i)}{\alpha_0+1}$, where
 $\tilde{\alpha}_i = \frac{\alpha_i}{\alpha_0}$ and $\alpha_0 = \sum_i^S \alpha_i$.

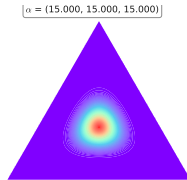


Figure H3: **Dirichlet Task Distribution** for $S = 3$ states, with $Dir(\alpha)$ where $\alpha = [15, 15, 15]$, resulting in our task-similarity measure approximately to be $\sigma = 0.0129$.

651

652 H.2 Ablations

653 We also run ablations with **Aggregating**($\alpha = 1$), a naive baseline that simply ignores the meta-RL
 654 structure and just plans assuming there is a single task. We observe in Fig. H4 the aggregating baseline
 655 works at-par with our method POMRL which is intuitive when the tasks are strongly related to each other
 656 in this case. However, as the underlying task structure decreases, we note that **Aggregating**($\alpha = 1$)
 657 as though it is one single task is problematic and suffers from a non-vanishing bias due to which
 658 for each new task there is on average an error which does not go to zero. More importantly, the
 659 **Aggregating**($\alpha = 1$) baseline cannot have the same guarantees as POMRL and ada-POMRL .

660 H.3 Additional Experiments

661 We examine more properties of ada-POMRL , namely **Effect of m , and T on Planning Loss** in Fig.
 662 **H5**, **Individual Baseline's Performance** in Fig. **H6**, and **Varying State Space $|S|$, m , and T** in Fig.
 663 **H7**.

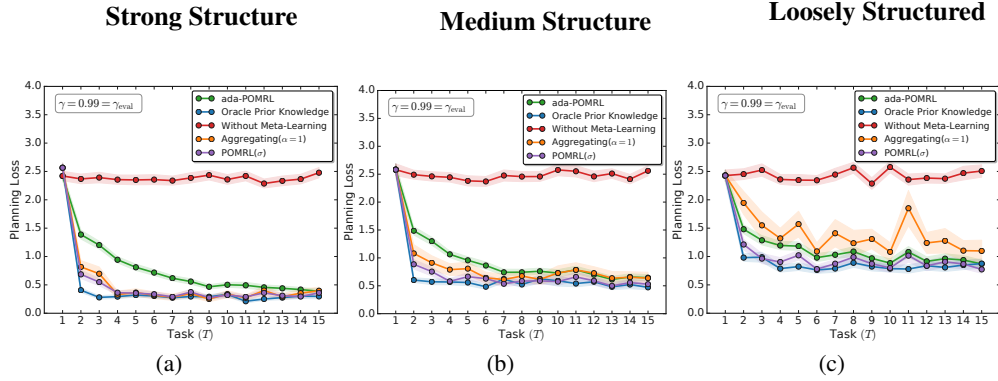


Figure H4: **Ablations for Efficacy of POMRL and ada-POMRL for varying task-similarity.** depicts the effect of the task-similarity parameter σ for a small fixed amount of data $m = 5$ available at each round. We run another baseline called Aggregating (orange) that simply ignores the meta-RL structure and acts as if it is all one single task. In the presence of strong structure, meta-learning the shared structure alongside a good model initialization leads to most gains and even naively aggregating the tasks transitions might seem to work well. However, such a naive method is not reliable as the underlying task similarity decreases - the learner struggles to cope with new unseen tasks which differ significantly and the planning loss doesn't improve. Error bars represent 1-standard deviation of uncertainty across 100 independent runs.

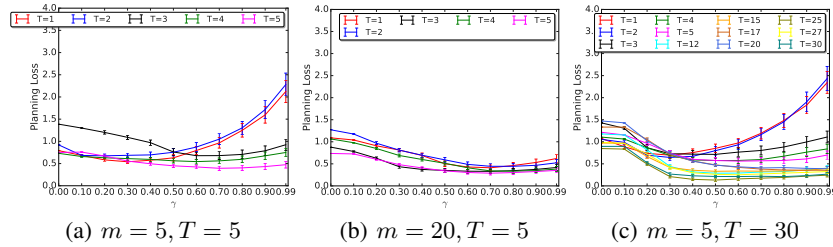


Figure H5: **Effect of m and T on Average Regret Upper Bound on Planning:** for a fixed value of task similarity σ , depends on the number of samples per task m and the number of tasks T . (a) For $m = T$, smaller loss is obtained with very small discount factor. This implies that with a lot of uncertainty it is not interesting to plan far too ahead, (b) For $m \gg T$, each task has enough samples to inform itself resulting in slightly larger effective discount factors. Not a lot is gained in this scenario from meta-learning, (c) $m \ll T$ is the most interesting case as samples seen in each individual task are very limited due to small m . However, the number of tasks are much more resulting in huge gains from leveraging shared structure across tasks.

664 H.4 Reproducibility

665 We follow the reproducibility checklist by Pineau (2019) to ensure this research is reproducible. For
 666 all algorithms presented, we include a clear description of the algorithm and source code is included
 667 with these supplementary materials. For any theoretical claims, we include: a statement of the result,
 668 a clear explanation of any assumptions, and complete proofs of any claims. For all figures that
 669 present empirical results, we include: the empirical details of how the experiments were run, a clear
 670 definition of the specific measure or statistics used to report results, and a description of results with
 671 the standard error in all cases.

672 H.5 Computing and Open source libraries.

673 All experiments were conducted using Google Colab instances⁷.

⁷<https://colab.research.google.com/>

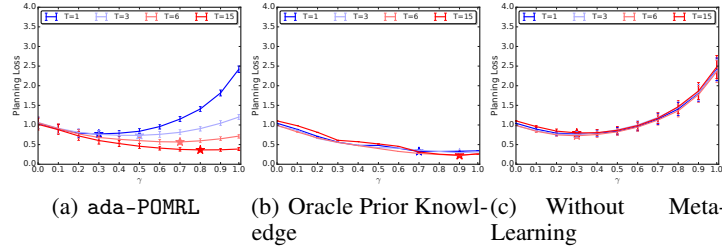


Figure H6: **Planning with Online Meta Learning - Baselines.** (a) **ada-POMRL**. Meta updates include learning P_o , σ , α as a function of tasks. (b) **Oracle Prior Knowledge** considers the optimal α , true mean of the task distribution P_o and actual underlying task similarity σ as known apriori, (c) **Without Meta-Learning** estimates the transition kernel in each round T without any meta-learning. All baselines are obtained with $T = 15$ tasks and $m = 5$ samples per task.

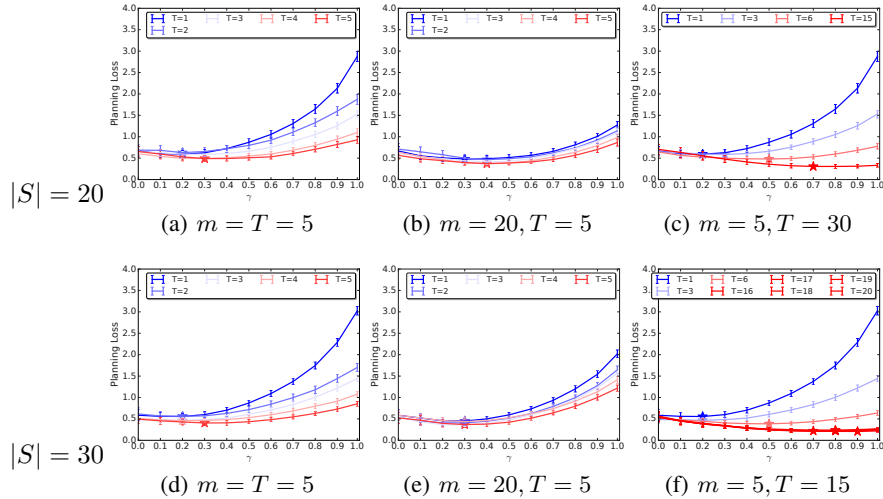


Figure H7: **Varying the size of state-space S , number of samples per task m , and number of tasks T , on Task-averaged Regret Upper Bound on Planning:** for a fixed value of task similarity σ , We note that despite larger state-space we observe the same effect i.e. (a,d,g) For $m = T$, smaller loss is obtained with very small discount factor i.e. a lot of uncertainty and inability to plan far too ahead, (b,e,h) For $m \gg T$, each task has enough samples to inform itself resulting in slightly larger effective discount factors. Not a lot is gained in this scenario from meta-learning. (c,f,i) $m \ll T$ is the most interesting case as samples seen in each individual task are very limited due to small m . Meta-learning has most significant gains in this case by leveraging the structure across tasks. Results are averaged over 20 independent runs and error bars represent 1-standard deviation.

674 **H.6 Extended Discussion**

675 **Beyond the tabular case:** Function approximation is at the heart of practical RL so a natural question
676 is how to extend our work to parametrized models. For linear MDPs, Müller and Pacchiano (2022)
677 recently derived regret bounds in the fixed-horizon setting for an algorithm using meta-regularizers
678 similar to ours. One question is whether this idea could be extended to infinite horizons and further
679 to non-linear, richer representations. Another, and perhaps deeper question, is around designing
680 and evaluating better planning strategies. Should we revisit such line of work under the light of
681 the planning loss rather than the regret? **On- or Off- Policy Meta-Learning without a simulator:**
682 Realistic problem settings in RL involve using sequentially learnt policies to collect data instead
683 of the simulator. One direction could be to extend our approach to model-based RL algorithms
684 via meta-gradient updates as in ARUBA or MAML, and seek regret guarantees induced by our
685 concentration results. **Non-stationary meta-distribution:** Many real-world scenarios have (slow or
686 sudden) drifts in the underlying distribution itself, e.g. weather. A promising future direction is to
687 consider non-stationary environments where the optimal initialization varies over time.