

PRETRAINING DECISION TRANSFORMERS WITH REWARD PREDICTION FOR IN-CONTEXT MULTI-TASK STRUCTURED BANDIT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we study the multi-task structured bandit problem where the goal is to learn a near-optimal algorithm that minimizes cumulative regret. The tasks share a common structure and the algorithm exploits the shared structure to minimize the cumulative regret for an unseen but related test task. We use a transformer as a decision-making algorithm to learn this shared structure so as to generalize to the test task. The prior work of pretrained decision transformers like **DPT** requires access to the optimal action during training which may be hard in several scenarios. Diverging from these works, our learning algorithm does not need the knowledge of optimal action per task during training but predicts a reward vector for each of the actions using only the observed offline data from the diverse training tasks. Finally, during inference time, it selects action using the reward predictions employing various exploration strategies in-context for an unseen test task. We show that our model outperforms other SOTA methods like **DPT**, and Algorithmic Distillation (**AD**) over a series of experiments on several structured bandit problems (linear, bilinear, latent, non-linear). Interestingly, we show that our algorithm, without the knowledge of the underlying problem structure, can learn a near-optimal policy in-context by leveraging the shared structure across diverse tasks. We further extend the field of pre-trained decision transformers by showing that they can leverage unseen tasks with new actions and still learn the underlying latent structure to derive a near-optimal policy. We validate this over several experiments to show that our proposed solution is very general and has wide applications to potentially emergent online and offline strategies at test time. Finally, we theoretically analyze the performance of our algorithm and obtain generalization bounds in the in-context multi-task learning setting.

1 INTRODUCTION

In this paper, we study multi-task bandit learning with the goal of learning an algorithm that discovers and exploits structure in a family of related tasks. In multi-task bandit learning, we have multiple distinct bandit tasks for which we want to learn a policy. Though distinct, the tasks share some structure, which we hope to leverage to speed up learning on new instances in this task family. Traditionally, the study of such structured bandit problems has relied on knowledge of the problem structure like linear bandits (Li et al., 2010; Abbasi-Yadkori et al., 2011; Degenne et al., 2020), bilinear bandits (Jun et al., 2019), hierarchical bandits (Hong et al., 2022a;b), Lipschitz bandits (Bubeck et al., 2008; 2011; Magureanu et al., 2014), other structured bandits settings (Riquelme et al., 2018; Lattimore & Szepesvári, 2019; Dong et al., 2021) and even linear and bilinear multi-task bandit settings (Yang et al., 2022a; Du et al., 2023; Mukherjee et al., 2023). When structure is unknown an alternative is to adopt sophisticated model classes, such as kernel machines or neural networks, exemplified by kernel or neural bandits (Valko et al., 2013; Chowdhury & Gopalan, 2017; Zhou et al., 2020; Dai et al., 2022). However, these approaches are also costly as they learn complex, nonlinear models from the ground up without any prior data (Justus et al., 2018; Zambaldi et al., 2018).

In this paper, we consider an alternative approach of synthesizing a bandit algorithm from historical data where the data comes from recorded bandit interactions with past instances of our target task family. Concretely, we are given a set of state-action-reward tuples obtained by running some bandit

algorithm in various instances from the task family. We then aim to train a transformer (Vaswani et al., 2017) from this data such that it can learn in-context to solve new task instances. Laskin et al. (2022) consider a similar goal and introduce the Algorithm Distillation (AD) method, however, AD aims to copy the algorithm used in the historical data and thus is limited by the ability of the data collection algorithm. Lee et al. (2023) develop an approach, DPT, that enables learning a transformer that obtains lower regret in-context bandit learning compared to the algorithm used to produce the historical data. However, this approach requires knowledge of the optimal action at each stage of the decision process. In real problems, this assumption is hard to satisfy and we will show that DPT performs poorly when the optimal action is only approximately known. With this past work in mind, the goal of this paper is to answer the question:

Can we learn an in-context bandit learning algorithm that obtains lower regret than the algorithm used to produce the training data without knowledge of the optimal action in each training task?

To answer this question, we introduce a new pre-training methodology, called **Pre-trained Decision Transformer with Reward Estimation (PreDeToR)** that obviates the need for knowledge of the optimal action in the in-context data — a piece of information that is often inaccessible. Our key observation is that while the mean rewards of each action change from task to task, certain probabilistic dependencies are persistent across all tasks with a given structure (Yang et al., 2020; 2022a; Mukherjee et al., 2023). These probabilistic dependencies can be learned from the pretraining data and exploited to better estimate mean rewards and improve performance in a new unknown test task. The nature of the probabilistic dependencies depends on the specific structure of the bandit and can be complex (i.e., higher-order dependencies beyond simple correlations). We propose to use transformer models as a general-purpose architecture to capture the unknown dependencies by training transformers to predict the mean rewards in each of the given trajectories (Mirchandani et al., 2023; Zhao et al., 2023). The key idea is that transformers have the capacity to discover and exploit complex dependencies in order to predict the rewards of all possible actions in each task from a *small* history of action-reward pairs in a new task. This paper demonstrates how such an approach can achieve lower regret by outperforming state-of-the-art baselines, relying solely on historical data, without the need for any supplementary information like the action features or knowledge of the complex reward models. We also show that the shared actions across the tasks are vital for PreDeToR to exploit the latent structure. We show that PreDeToR learns to adapt, in-context, to novel actions and new tasks as long as the number of new actions is small compared to shared actions across the tasks.

Contributions

1. We introduce a new pre-training procedure of learning the underlying reward structure and a decision algorithm. Moreover, PreDeToR by predicting the next reward for all arms circumvents the issue of requiring access to the optimal (or approximately optimal) action during training time.
2. We demonstrate empirically that this training procedure results in lower regret in a wide series of tasks (such as linear, nonlinear, bilinear, and latent bandits) compared to prior in-context learning algorithms and bandit algorithms with privileged knowledge of the common structure.
3. We also show that our training procedure leverages the shared latent structure and is robust to a small number of new actions introduced both during training and testing time.
4. Finally, we theoretically analyze the generalization ability of PreDeToR through the lens of algorithmic stability and new results for the transformer setting.

2 BACKGROUND

In this section, we first introduce our notation and the multi-task, structured bandit setting. We then formalize the in-context bandit learning model studied in Laskin et al. (2022); Lee et al. (2023); Sinii et al. (2023); Lin et al. (2023); Ma et al. (2023); Liu et al. (2023c;a).

2.1 PRELIMINARIES

In this paper, we consider the multi-task linear bandit setting (Du et al., 2023; Yang et al., 2020; 2022a). In the multi-task setting, we have a family of related bandit problems that share an action set \mathcal{A} and also a common action feature space \mathcal{X} . The actions in \mathcal{A} are indexed by $a = 1, 2, \dots, A$. The feature of each action is denoted by $\mathbf{x}(a) \in \mathbb{R}^d$ and $d \ll A$. A policy, π , is a probability distribution over the actions.

Define $[n] = \{1, 2, \dots, n\}$. In a multi-task structured bandit setting the expected reward for each action in each task is assumed to be an unknown function of the hidden parameter and action features (Lattimore & Szepesvári, 2020; Gupta et al., 2020). The interaction proceeds iteratively over n rounds for each task $m \in [M]$. At each round $t \in [n]$ for each task $m \in [M]$, the learner selects an action $I_{m,t} \in \mathcal{A}$ and observes the reward $r_{m,t} = f(\mathbf{x}(I_{m,t}), \boldsymbol{\theta}_{m,*}) + \eta_{m,t}$, where $\boldsymbol{\theta}_{m,*} \in \mathbb{R}^d$ is the hidden parameter specific to the task m to be learned by the learner. The function $f(\cdot, \cdot)$ is the unknown reward structure. This can be $f(\mathbf{x}(I_{m,t}), \boldsymbol{\theta}_{m,*}) = \mathbf{x}(I_{m,t})^\top \boldsymbol{\theta}_{m,*}$ for the linear setting or even more complex correlation between features and $\boldsymbol{\theta}_{m,*}$ (Filippi et al., 2010; Abbasi-Yadkori et al., 2011; Riquelme et al., 2018; Lattimore & Szepesvári, 2019; Dong et al., 2021).

In our paper, we assume that there exist weak demonstrators denoted by π^w . These weak demonstrators are stochastic A -armed bandit algorithms like Upper Confidence Bound (UCB) (Auer et al., 2002; Auer & Ortner, 2010) or Thompson Sampling (Thompson, 1933; Agrawal & Goyal, 2012; Russo et al., 2018; Zhu & Tan, 2020). We refer to these algorithms as weak demonstrators because they do not use knowledge of task structure or arm feature vectors to plan their sampling policy. In contrast to a weak demonstrator, a strong demonstrator, like LinUCB, uses feature vectors and knowledge of task structure to conduct informative exploration. Whereas weak demonstrators always exist, there are many real-world settings with no known strong demonstrator algorithm or where the feature vectors are unobserved and the learner can only use the history of rewards and actions.

2.2 IN-CONTEXT LEARNING MODEL

Similar to Lee et al. (2023); Sinii et al. (2023); Lin et al. (2023); Ma et al. (2023); Liu et al. (2023c;a) we assume the in-context learning model. We first discuss the pretraining procedure.

Pretraining: Let \mathcal{T}_{pre} denote the distribution over tasks m at the time of pretraining. Let \mathcal{D}_{pre} be the distribution over all possible interactions that the π^w can generate. We first sample a task $m \sim \mathcal{T}_{\text{pre}}$ and then a context \mathcal{H}_m which is a sequence of interactions for n rounds conditioned on the task m such that $\mathcal{H}_m \sim \mathcal{D}_{\text{pre}}(\cdot | m)$. So $\mathcal{H}_m = \{I_{m,t}, r_{m,t}\}_{t=1}^n$. We call this dataset \mathcal{H}_m an in-context dataset as it contains the contextual information about the task m . We denote the samples in \mathcal{H}_m till round t as $\mathcal{H}_m^t = \{I_{m,s}, r_{m,s}\}_{s=1}^{t-1}$. This dataset \mathcal{H}_m can be collected in several ways: (1) random interactions within m , (2) demonstrations from an expert, and (3) rollouts of an algorithm. Finally, we train a causal GPT-2 transformer model TF parameterized by Θ on this dataset \mathcal{D}_{pre} . Specifically, we define $\text{TF}_\Theta(\cdot | \mathcal{H}_m^t)$ as the transformer model that observes the dataset \mathcal{H}_m^t till round t and then produces a distribution over the actions. Our primary novelty lies in our training procedure which we explain in detail in Section 3.1.

Testing: We now discuss the testing procedure for our setting. Let $\mathcal{T}_{\text{test}}$ denote the distribution over test tasks $m \in [M_{\text{test}}]$ at the time of testing. Let $\mathcal{D}_{\text{test}}$ denote a distribution over all possible interactions that can be generated by π^w during test time. At deployment time, the dataset $\mathcal{H}_m^0 \leftarrow \{\emptyset\}$ is initialized empty. At each round t , an action is sampled from the trained transformer model $I_t \sim \text{TF}_\Theta(\cdot | \mathcal{H}_m^t)$. The sampled action and resulting reward, r_t , are then added to \mathcal{H}_m^t to form \mathcal{H}_m^{t+1} and the process repeats for n total rounds. Finally, note that in this testing phase, the model parameter Θ is not updated. Finally, the goal of the learner is to minimize cumulative regret for all task $m \in [M_{\text{test}}]$ defined as follows: $\mathbb{E}[R_n] = \frac{1}{M_{\text{test}}} \sum_{m=1}^{M_{\text{test}}} \sum_{t=1}^n \max_{a \in \mathcal{A}} f(\mathbf{x}(a), \boldsymbol{\theta}_{m,*}) - f(\mathbf{x}(I_t), \boldsymbol{\theta}_{m,*})$.

2.3 RELATED IN-CONTEXT LEARNING ALGORITHMS

In this section, we discuss related algorithms for in-context decision-making. For completeness, we describe the **DPT** and **AD** training procedure and algorithm now. During training, **DPT** first samples $m \sim \mathcal{T}_{\text{pre}}$ and then an in-context dataset $\mathcal{H}_m \sim \mathcal{D}_{\text{pre}}(\cdot | m)$. It adds this \mathcal{H}_m to the training dataset $\mathcal{H}_{\text{train}}$, and repeats to collect M_{pre} such training tasks. For each task m , **DPT** requires the optimal action $a_{m,*} = \arg \max_a f(\mathbf{x}(m, a), \boldsymbol{\theta}_{m,*})$ where $f(\mathbf{x}(m, a), \boldsymbol{\theta}_{m,*})$ is the expected reward for the action a in task m . Since the optimal action is usually not known in advance, in Section 4 we introduce a practical variant of **DPT** that approximates the optimal action with the best action

162 identified during task interaction. During training **DPT** minimizes the cross-entropy loss:

$$163 \mathcal{L}_t^{\text{DPT}} = \text{cross-entropy}(\text{TF}_{\Theta}(\cdot|\mathcal{H}_m^t), p(a_{m,*})) \quad (1)$$

164 where $p(a_{m,*}) \in \Delta^A$ is a one-hot vector such that $p(j) = 1$ when $j = a_{m,*}$ and 0 otherwise. This loss
165 is then back-propagated and used to update the model parameter Θ .

166 During test time evaluation for online setting the **DPT** selects $I_t \sim \text{softmax}_a^{\tau}(\text{TF}_{\Theta}(\cdot|\mathcal{H}_m^t))$
167 where we define the $\text{softmax}_a^{\tau}(\mathbf{v})$ over a A dimensional vector $\mathbf{v} \in \mathbb{R}^A$ as $\text{softmax}_a^{\tau}(\mathbf{v}(a)) =$
168 $\exp(\mathbf{v}(a)/\tau) / \sum_{a'=1}^A \exp(\mathbf{v}(a')/\tau)$ which produces a distribution over actions weighted by the
169 temperature parameter $\tau > 0$. Therefore this sampling procedure has a high probability of choos-
170 ing the predicted optimal action as well as induce sufficient exploration. In the online setting, the
171 **DPT** observes the reward $r_t(I_t)$ which is added to \mathcal{H}_m^t . So the \mathcal{H}_m during online testing consists
172 of $\{I_t, r_t\}_{t=1}^n$ collected during testing. This interaction procedure is conducted for each test task
173 $m \in [M_{\text{test}}]$. In the testing phase, the model parameter Θ is not updated.

174 An alternative to **DPT** that does *not* require knowledge of the optimal action is the **AD** approach
175 (**Laskin et al., 2022; Lu et al., 2023**). In **AD**, the learner aims to predict the next action of the
176 demonstrator. So it minimizes the cross-entropy loss as follows:

$$177 \mathcal{L}_t^{\text{AD}} = \text{cross-entropy}(\text{TF}_{\Theta}(\cdot|\mathcal{H}_m^t), p(I_{m,t})) \quad (2)$$

178 where $p(I_{m,t})$ is a one-hot vector such that $p(j) = 1$ when $j = I_{m,t}$ (the true action taken by the
179 demonstrator) and 0 otherwise. At deployment time, **AD** selects $I_t \sim \text{softmax}_a^{\tau}(\text{TF}_{\Theta}(\cdot|\mathcal{H}_m^t))$. Note
180 that the objective of **AD** is to match the performance of the demonstrator. In the next section, we
181 introduce a new method that can improve upon the demonstrator without knowledge of the optimal
182 action.

183 3 PROPOSED ALGORITHM **PREDETOR**

184 We now introduce our main algorithmic contribution, **PreDeToR** (which stands for **Pre-trained**
185 **Decision Transformer with Reward Estimation**).

186 3.1 PRE-TRAINING NEXT REWARD PREDICTION

187 The key idea behind **PreDeToR** is to leverage the in-context learning ability of transformers to infer
188 the reward of each arm in a given test task. By training this in-context ability on a set of training
189 tasks, the transformer can implicitly learn structure in the task family and exploit this structure
190 to infer rewards without trying every single arm. Thus, in contrast to **DPT** and **AD** that output
191 actions directly, **PreDeToR** outputs a scalar value reward prediction for each arm. To this effect, we
192 append a linear layer of dimension A on top of a causal GPT2 model, denoted by $\text{TF}^r_{\Theta}(\cdot|\mathcal{H}_m)$,
193 and use a least-squares loss to train the transformer to predict the reward for each action with these
194 outputs. Note that we use $\text{TF}^r_{\Theta}(\cdot|\mathcal{H}_m)$ to denote a reward prediction transformer and $\text{TF}_{\Theta}(\cdot|\mathcal{H}_m)$
195 as the transformer that predicts a distribution over actions (as in **DPT** and **AD**). At every round
196 t the transformer predicts the *next reward* for each of the actions $a \in \mathcal{A}$ for the task m based on
197 $\mathcal{H}_m^t = \{I_{m,s}, r_{m,s}\}_{s=1}^{t-1}$. This predicted reward is denoted by $\hat{r}_{m,t+1}(a)$ for each $a \in \mathcal{A}$.

198 **Loss calculation:** For each training task, m , we calculate the loss at each round, t , using the
199 transformer’s prediction $\hat{r}_{m,t}(I_{m,t})$ and the actual observed reward $r_{m,t}$ that followed action $I_{m,t}$.
200 We use a least-squares loss function:

$$201 \mathcal{L}_t = (\hat{r}_{m,t}(I_{m,t}) - r_{m,t})^2 \quad (3)$$

202 and hence minimizing this loss will minimize the mean squared-error of the transformer’s predictions.
203 The loss is calculated using equation 3 and is backpropagated to update the model parameter Θ .

204 **Exploratory Demonstrator:** Observe from the loss definition in equation 3 that it is calculated
205 from the observed true reward and action from the dataset \mathcal{H}_m . In order for the transformer to learn
206 accurate reward predictions during training, we require that the weak demonstrator is sufficiently
207 exploratory such that it collects \mathcal{H}_m such that \mathcal{H}_m contains some reward $r_{m,t}$ for each action a . We
208 discuss in detail the impact of the demonstrator on **PreDeToR** ($-\tau$) training in Appendix A.13.

3.2 DEPLOYING PREDETOR

At deployment time, **PreDeToR** learns in-context to predict the mean reward of each arm on an unseen task and acts greedily with respect to this prediction. That is, at deployment time, a new task is sampled, $m \sim \mathcal{T}_{\text{test}}$, and the dataset \mathcal{H}_m^0 is initialized empty. Then at every round t , **PreDeToR** chooses $I_t = \arg \max_{a \in \mathcal{A}} \text{TF}^{\text{r}}_{\Theta}(\hat{r}_{m,t}(a) \mid \mathcal{H}_m^t)$ which is the action with the highest predicted reward and $\hat{r}_{m,t}(a)$ is the predicted reward of action a . Note that **PreDeToR** is a greedy policy and thus may fail to conduct sufficient exploration. To remedy this potential limitation, we also introduce a soft variant, **PreDeToR- τ** that chooses $I_t \sim \text{softmax}_a^{\tau}(\text{TF}^{\text{r}}_{\Theta}(\hat{r}_{m,t}(a) \mid \mathcal{H}_m^t))$. For both **PreDeToR** and **PreDeToR- τ** , the observed reward $r_t(I_t)$ is added to the dataset \mathcal{H}_m and then used to predict the reward at the next round $t + 1$. The full pseudocode of using **PreDeToR** for online interaction is shown in Algorithm 1. In Appendix A.15, we discuss how **PreDeToR** ($-\tau$) can be deployed for offline learning.

Algorithm 1 Pre-trained Decision Transformer with Reward Estimation (**PreDeToR**)

- 1: **Collecting Pretraining Dataset**
 - 2: Initialize empty pretraining dataset $\mathcal{H}_{\text{train}}$
 - 3: **for** i in $[M_{\text{pre}}]$ **do**
 - 4: Sample task $m \sim \mathcal{T}_{\text{pre}}$, in-context dataset $\mathcal{H}_m \sim \mathcal{D}_{\text{pre}}(\cdot \mid m)$ and add this to $\mathcal{H}_{\text{train}}$.
 - 5: **end for**
 - 6: **Pretraining model on dataset**
 - 7: Initialize model $\text{TF}^{\text{r}}_{\Theta}$ with parameters Θ
 - 8: **while** not converged **do**
 - 9: Sample \mathcal{H}_m from $\mathcal{H}_{\text{train}}$ and predict $\hat{r}_{m,t}$ for action $(I_{m,t})$ for all $t \in [n]$
 - 10: Compute loss in equation 3 with respect to $r_{m,t}$ and backpropagate to update model parameter Θ .
 - 11: **end while**
 - 12: **Online test-time deployment**
 - 13: Sample unknown task $m \sim \mathcal{T}_{\text{test}}$ and initialize empty $\mathcal{H}_m^0 = \{\emptyset\}$
 - 14: **for** $t = 1, 2, \dots, n$ **do**
 - 15: Use $\text{TF}^{\text{r}}_{\Theta}$ on m at round t to choose

$$I_t \begin{cases} = \arg \max_{a \in \mathcal{A}} \text{TF}^{\text{r}}_{\Theta}(\hat{r}_{m,t}(a) \mid \mathcal{H}_m^t), & \text{PreDeToR} \\ \sim \text{softmax}_a^{\tau} \text{TF}^{\text{r}}_{\Theta}(\hat{r}_{m,t}(a) \mid \mathcal{H}_m^t), & \text{PreDeToR-}\tau \end{cases}$$
 - 16: Add $\{I_t, r_t\}$ to \mathcal{H}_m^t to form \mathcal{H}_m^{t+1} .
 - 17: **end for**
-

4 EMPIRICAL STUDY: NON-LINEAR STRUCTURE

Having introduced **PreDeToR**, we now investigate its performance in diverse bandit settings compared to other in-context learning algorithms. In our first set of experiments, we use a bandit setting with a common non-linear structure across tasks. Ideally, a good learner would leverage the structure, however, we choose the structure such that no existing algorithms are well-suited to the non-linear structure. This setting is thus a good testbed for establishing that in-context learning can discover and exploit common structure. Moreover, each task only consists of a few rounds of interactions. This setting is quite common in recommender settings where user interaction with the system lasts only for a few rounds and has an underlying non-linear structure (Kwon et al., 2022; Tomkins et al., 2020). We show that **PreDeToR** achieves lower regret than other in-context algorithms for the non-linear structured bandit setting. We study the performance of **PreDeToR** in the large horizon setting in Appendix A.7.

Baselines: We first discuss the baselines used in this setting.

- (1) **PreDeToR:** This is our proposed method shown in Algorithm 1.
- (2) **PreDeToR- τ :** This is the proposed exploratory method shown in Algorithm 1 and we fix $\tau = 0.05$.
- (3) **DPT-greedy:** This baseline is the greedy approximation of the **DPT** algorithm from Lee et al. (2023) which is discussed in Section 2.3. Note that we choose **DPT-greedy** as a representative

example of similar in-context decision-making algorithms studied in Lee et al. (2023); Sinii et al. (2023); Lin et al. (2023); Ma et al. (2023); Liu et al. (2023c;a) all of which require the optimal action (or its greedy approximation). **DPT-greedy** estimates the optimal arm using the reward estimates for each arm during each task.

(4) **AD**: This is the Algorithmic Distillation method (Laskin et al., 2022; Lu et al., 2021) discussed in Section 2.3.

(5) **Thomp**: This baseline is the celebrated stochastic A -action bandit Thompson Sampling algorithm from Thompson (1933); Agrawal & Goyal (2012); Russo et al. (2018); Zhu & Tan (2020). We choose **Thomp** as the weak demonstrator π^w as it does not make use of arm features. **Thomp** is also a stochastic algorithm that induces more exploration in the demonstrations.

(6) **LinUCB**: (Linear Upper Confidence Bound): This baseline is the Upper Confidence Bound algorithm for the linear bandit setting that leverages the linear structure and feature of the arms to select the most promising action as well as conducting exploration. We choose **LinUCB** as a baseline for each test task to show the limitations of algorithms that use linear feedback structure as an underlying assumption to select actions. Note that **LinUCB** requires oracle access to features to select actions per task.

(7) **MLinGreedy**: This is the multi-task linear regression bandit algorithm proposed by Yang et al. (2021). This algorithm assumes that there is a common low-dimensional feature extractor shared between the tasks and the reward of each task linearly depends on this feature extractor. We choose **MLinGreedy** as a baseline to show the limitations of algorithms that use linear feedback structure *across tasks* as an underlying assumption to select actions. Note that **MLinGreedy** requires oracle access to the action features to select actions as opposed to **DPT**, **AD**, and **PreDeToR**.

We describe in detail the baselines **Thomp**, **LinUCB**, and **MLinGreedy** for interested readers in Appendix A.2.2.

Outcomes: Before presenting the result we discuss the main outcomes from our experimental results in this section:

Finding 1: **PreDeToR** ($-\tau$) lowers regret compared to other baselines under unknown, non-linear structure. It learns to exploit the latent structure of the underlying tasks from in-context data even when it is trained without the optimal action $a_{m,*}$ (or its approximation) and without action features \mathcal{X} .

Experimental Result: These findings are reported in Figure 1. In Figure 1(a) we show the non-linear bandit setting for horizon $n = 50$, $M_{\text{pre}} = 100000$, $M_{\text{test}} = 200$, $A = 6$, and $d = 2$. The demonstrator π^w is the **Thomp** algorithm. We observe that **PreDeToR** ($-\tau$) has lower cumulative regret than **DPT-greedy**. Note that for this low data regime (short horizon) the **DPT-greedy** does not have a good estimation of $\hat{a}_{m,*}$ which results in a poor prediction of optimal action $\hat{a}_{m,t,*}$. This results in higher regret. The **PreDeToR** ($-\tau$) has lower regret than **LinUCB**, and **MLinGreedy**, which fail to perform well in this non-linear setting due to their algorithmic design and linear feedback assumption. Finally, **PreDeToR** ($-\tau$) performs slightly better than **PreDeToR** in both settings as it conducts more exploration.

In Figure 1(b) we show the non-linear bandit setting for horizon $n = 25$, $M_{\text{pre}} = 100000$, $M_{\text{test}} = 200$, $A = 6$, and $d = 2$ where the norm of the $\theta_{m,*}$ determines the reward of the actions which also is a non-linear function $\theta_{m,*}$ and action features. This setting is similar to the wheel bandit setting of Riquelme et al. (2018). Again, we observe that **PreDeToR** has lower cumulative regret than all the other baselines.

Finally in Figure 1(c) and Figure 1(d) we show the performance of **PreDeToR** against other baselines in real-world datasets MovieLens and Yelp. The MovieLens dataset consists of more than 32 million ratings of 200,000 users and 80,000 movies (Harper & Konstan, 2015) where each entry consists of user-id, movie-id, rating, and timestamp. The Yelp dataset (Asghar, 2016) consists of ratings of 1300 business categories by 150,000 users. Each entry is summarized as user-id, business-id, rating, and timestamp. Previously structured bandit works (Deshpande & Montanari, 2012; Hong et al., 2023) directly fit a linear structure or low-rank factorization to estimate the $\theta_{m,*}$ and simulate the

ratings. However, we directly use the user-ids and movie-ids (or business-ids) to build a histogram of ratings per user and calculate the mean rating per movie (or business-id) per task. Define this as the $\{\mu_{m,a}\}_{a=1}^A$. This is then used to simulate the rating for n horizon per movie per task where the data collection algorithm is uniform sampling. Note that this does not require estimation of user or movie features, and **PreDeToR** ($-\tau$) learns to exploit the latent structure of user-movie (or business) rating correlations directly from the data. From Figure 1(c) and Figure 1(d) we see that **PreDeToR**, and **PreDeToR**- τ outperform all the other baselines in these settings.

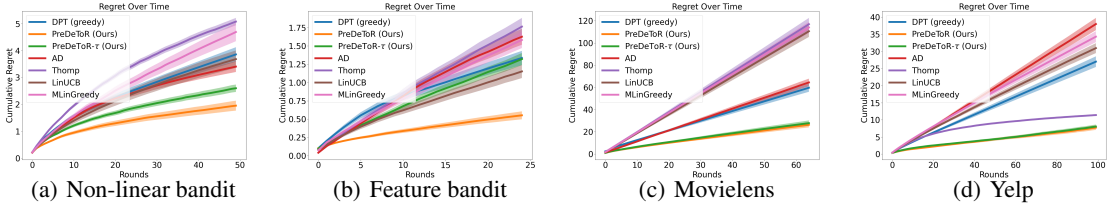


Figure 1: Non-linear regime. The horizontal axis is the number of rounds. Confidence bars show one standard error.

5 EMPIRICAL STUDY: LINEAR STRUCTURE AND UNDERSTANDING THE EXPLORATION OF **PREDETOR**

The previous experiments were conducted in a non-linear structured setting where we are unaware of a provably near-optimal algorithm. To assess how close **PreDeToR**'s regret is to optimal, in this section, we consider a *linear* setting for which there exist well-understood algorithms (Abbasi-Yadkori et al., 2011; Lattimore & Szepesvári, 2020). Such algorithms provide a strong upper bound for **PreDeToR**. We summarize the key finding below:

Finding 2: **PreDeToR** ($-\tau$) matches the performance of the optimal algorithm **LinUCB** in linear bandit setting as it learns to exploit the latent structure across tasks from in-context data and without access to features.

In Figure 2 we first show the linear bandit setting for horizon $n = 25$, $M_{\text{pre}} = 200000$, $M_{\text{test}} = 200$, $A = 10$, and $d = 2$. Note that the length of the context (the number of rounds) is an artifact of the transformer architecture and computational complexity. This is because the self-attention takes in as input a length- n sequence of tokens of size d , and requires $O(dn^2)$ time to compute the output (Keles et al., 2023). Further empirical setting details are stated in Appendix A.2.

We observe from Figure 2 that **PreDeToR** ($-\tau$) has lower cumulative regret than **DPT-greedy**, and **AD**. Note that for this low data (short horizon) regime, the **DPT-greedy** does not have a good estimation of $\hat{a}_{m,*}$ which results in a poor prediction of optimal action $\hat{a}_{m,t,*}$. This results in higher regret. Observe that **PreDeToR** ($-\tau$) performs quite similarly to **LinUCB** and lowers regret compared to **Thomp** which also shows that **PreDeToR** is able to exploit the latent linear structure and reward correlation of the underlying tasks. Note that **LinUCB** is close to the optimal algorithm for this linear bandit setting. **PreDeToR** outperforms **AD** as the main objective of **AD** is to match the performance of its demonstrator. In this short horizon, we see that **MLinGreedy** performs similarly to **LinUCB**.

We also show how the prediction error of the optimal action by **PreDeToR** is small compared to **LinUCB** in the linear bandit setting. In Figure 2(b) we first show how the 10 actions are distributed in the $M_{\text{test}} = 200$ test tasks. In Figure 2(b) for each bar, the frequency indicates the number of tasks where the action (shown in the x-axis) is the optimal action. Then, in Figure 2(c), we show the prediction error of **PreDeToR** ($-\tau$) for each task $m \in [M_{\text{test}}]$. The prediction error is calculated as $(\hat{\mu}_{m,n,*}(a) - \mu_{m,*}(a))^2$ where $\hat{\mu}_{m,n,*}(a) = \max_a \hat{\theta}_{m,n}^\top \mathbf{x}_m(a)$ is the empirical mean at the end of round n , and $\mu_{m,*}(a) = \max_a \theta_{m,*}^\top \mathbf{x}_m(a)$ is the true mean of the optimal action in task m . Then we average the prediction error for the action $a \in [A]$ by the number of times the action a is the optimal action in some task m . From the Figure 2(c), we see that for actions $\{2, 3, 5, 6, 7, 10\}$, the prediction error of **PreDeToR** is either close or smaller than **LinUCB**. Note that **LinUCB** estimates the empirical

mean directly from the test task, whereas **PreDeToR** has a strong prior based on the training data. So **PreDeToR** is able to estimate the reward of the optimal action quite well from the training dataset \mathcal{D}_{pre} . This shows the power of **PreDeToR** to go beyond the in-context decision-making setting studied in Lee et al. (2023); Lin et al. (2023); Ma et al. (2023); Sini et al. (2023); Liu et al. (2023c) which require long horizons/trajectories and optimal action during training to learn a near-optimal policy. We discuss how exploration of **PreDeToR** ($-\tau$) results in low cumulative regret in Appendix A.11.

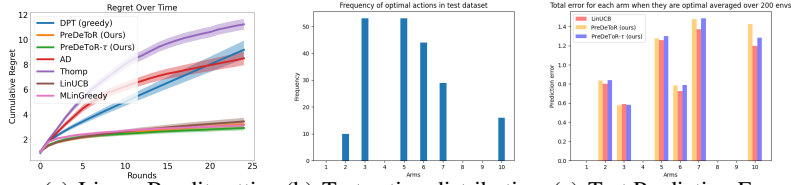


Figure 2: Linear Expt. The horizontal axis is the number of rounds. Confidence bars show one standard error.

6 EMPIRICAL STUDY: IMPORTANCE OF SHARED STRUCTURE AND INTRODUCING NEW ARMS

One of our central claims is that **PreDeToR** ($-\tau$) internally learns and leverages the shared structure across the training and testing tasks. To validate this claim, in this section, we consider the introduction of new actions at test time that do *not* follow the structure of training time. These experiments are particularly important as they show the extent to which **PreDeToR** ($-\tau$) is leveraging the latent structure and the shared correlation between the actions and rewards.

Invariant actions: We denote the set of actions fixed across the different tasks in the pretraining in-context dataset as \mathcal{A}^{inv} . Therefore these action features $\mathbf{x}(a) \in \mathbb{R}^d$ for $a \in \mathcal{A}^{\text{inv}}$ are fixed across the different tasks m . Note that these invariant actions help the transformer TF_w to learn the latent structure and the reward correlation across the different tasks. Therefore, as the structure breaks down, **PreDeToR** starts performing worse than other baselines.

New actions: However, we also want to test how robust is **PreDeToR** ($-\tau$) to new actions not seen during training time. To this effect, for each task $m \in [M_{\text{pre}}]$ and $m \in [M_{\text{test}}]$ we introduce $A - |\mathcal{A}^{\text{inv}}|$ new actions. *That is both for train and test tasks, we introduce new actions.* For each of these new actions $a \in [A - |\mathcal{A}^{\text{inv}}|]$ we choose the features $\mathbf{x}(m, a)$ randomly from $\mathcal{X} \subseteq \mathbb{R}^d$. Note the transformer now trains on a dataset $\mathcal{H}_m \subseteq \mathcal{D}_{\text{pre}} \neq \mathcal{D}_{\text{test}}$.

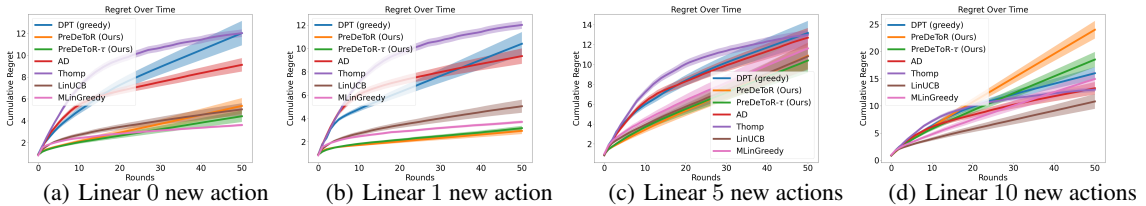


Figure 3: New action experiments. The horizontal axis is the number of rounds. Confidence bars show one standard error.

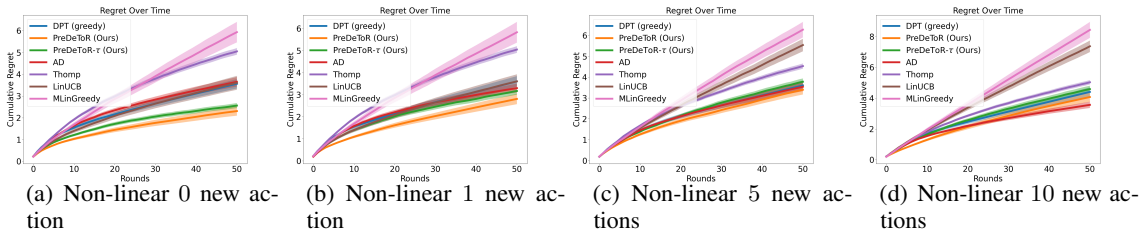


Figure 4: New action experiments with non-linear setting.

Baselines: We implement the same baselines discussed in Section 4.

Outcomes: Again before presenting the result we discuss the main outcomes from our experimental results of introducing new actions during data collection and evaluation:

Finding 3: PreDeToR ($-\tau$) performance degrades as the shared structure breaks down.

Experimental Result: We observe these outcomes in Figure 3 and Figure 4. We consider the linear and non-linear bandit setting of horizon $n = 50$, $M_{\text{pre}} = 100000$, $M_{\text{test}} = 200$, $A = 10$, and $d = 2$. Here during data collection and during collecting the test data, we randomly select between 0, 1, 5, and 10 new actions from \mathbb{R}^d for each task m . So the number of invariant actions is $|\mathcal{A}^{\text{inv}}| \in \{10, 5, 1, 0\}$. Again, the demonstrator π^w is the **Thomp** algorithm. From Figure 3(a), 3(b), 3(c), and 3(d), we observe that when the number of invariant actions is less than **PreDeToR** ($-\tau$) has lower cumulative regret than **DPT-greedy**, and **AD**. Observe that **PreDeToR** ($-\tau$) matches **LinUCB** and has lower regret than **DPT-greedy**, and **AD** when $|\mathcal{A}^{\text{inv}}| \in \{10, 5, 1\}$. This shows that **PreDeToR** ($-\tau$) is able to exploit the latent linear structure of the underlying tasks. However, as the number of invariant actions decreases we see that **PreDeToR** ($-\tau$) performance drops and becomes similar to the unstructured bandits **Thomp**.

Similarly in Figure 4(a), 4(b), 4(c), and 4(d) we show the performance of **PreDeToR** in the non-linear bandit setting. Observe that **LinUCB**, **MLinGreedy** fails to perform well in this non-linear setting due to their assumption of linear rewards. Again note that **PreDeToR** ($-\tau$) has lower regret than **DPT-greedy**, and **AD** when $|\mathcal{A}^{\text{inv}}| \in \{10, 1\}$. This shows that **PreDeToR** ($-\tau$) is able to exploit the latent linear structure of the underlying tasks. However, as the number of invariant actions decreases we see that **PreDeToR** ($-\tau$) performance drops and becomes similar to **AD**.

We also empirically study the test performance of **PreDeToR** ($-\tau$) in other *non-linear* bandit settings such as bilinear bandits (Appendix A.3), latent bandits (Appendix A.4), draw a connection between **PreDeToR** and Bayesian estimators (Appendix A.5), and perform sensitivity and ablation studies in Appendix A.6, A.8, A.9, A.10. We discuss data collection algorithms in Appendix A.13 and the offline setting in Appendix A.15. Due to space constraints, we refer the interested reader to the relevant section in the appendices.

7 THEORETICAL ANALYSIS OF GENERALIZATION

In this section, we present a theoretical analysis of how **PreDeToR**- τ generalizes to an unknown target task given a set of source tasks. We observe that **PreDeToR**- τ 's performance hinges on a low excess error on the predicted reward of the actions of the unknown target task based on the in-context data. Thus, in our analysis, we show that, in low-data regimes, **PreDeToR**- τ has a low expected excess risk for the unknown target task as the number of source tasks increases. This is summarized as follows:

Finding 4: **PreDeToR** ($-\tau$) has a low expected excess risk for the unknown target task as the number of source tasks increases. Moreover, the transfer learning risk of **PreDeToR**- τ (once trained on the M source tasks) scales with $\tilde{O}(1/\sqrt{M})$.

To show this, we proceed as follows: Suppose we have the training data set $\mathcal{H}_{\text{all}} = \{\mathcal{H}_m\}_{m=1}^{M_{\text{pre}}}$, where the task $m \sim \mathcal{T}$ with a distribution \mathcal{T} and the task data \mathcal{H}_m is generated from a distribution $\mathcal{D}_{\text{pre}}(\cdot|m)$. For illustration purposes, here we consider the training data distribution $\mathcal{D}_{\text{pre}}(\cdot|m)$ where the actions are sampled following soft-LinUCB (a stochastic variant of LinUCB) (Chu et al., 2011). Given the loss function in Equation (3), we can define the task m training loss of **PreDeToR**- τ as $\hat{\mathcal{L}}_m(\text{TF}^{\mathbf{r}}_{\Theta}) = \frac{1}{n} \sum_{t=1}^n \ell(r_{m,t}, \text{TF}^{\mathbf{r}}_{\Theta}(\hat{r}_{m,t}(I_{m,t})|\mathcal{H}_m^t)) = \frac{1}{n} \sum_{t=1}^n (\text{TF}^{\mathbf{r}}_{\Theta}(\hat{r}_{m,t}(I_{m,t})|\mathcal{H}_m^t) - r_{m,t})^2$. We drop the notation Θ, \mathbf{r} from $\text{TF}^{\mathbf{r}}_{\Theta}$ for simplicity and let $M = M_{\text{pre}}$. We define

$$\widehat{\text{TF}} = \arg \min_{\text{TF} \in \text{Alg}} \hat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\text{TF}) := \frac{1}{M} \sum_{m=1}^M \hat{\mathcal{L}}_m(\text{TF}), \quad (\text{ERM}) \quad (4)$$

where Alg denotes the space of algorithms induced by the TF. Let $\mathcal{L}_m(\text{TF}) = \mathbb{E}_{\mathcal{H}_m} [\hat{\mathcal{L}}_m(\text{TF})]$ and $\mathcal{L}_{\text{MTL}}(\text{TF}) = \mathbb{E} [\hat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\text{TF})] = \frac{1}{M} \sum_{m=1}^M \mathcal{L}_m(\text{TF})$ be the corresponding population risks. For

the ERM in equation 4, we want to bound the following excess Multi-Task Learning (MTL) risk of **PreDeToR- τ**

$$\mathcal{R}_{\text{MTL}}(\widehat{\text{TF}}) = \mathcal{L}_{\text{MTL}}(\widehat{\text{TF}}) - \min_{\text{TF} \in \text{Alg}} \mathcal{L}_{\text{MTL}}(\text{TF}). \quad (5)$$

Note that for in-context learning, a training sample (I_t, r_t) impacts all future decisions of the algorithm from time step $t + 1$ to n . Therefore, we need to control the stability of the input perturbation of the learning algorithm learned by the transformer. We introduce the following stability condition.

Assumption 7.1. (Error stability (Bousquet & Elisseeff, 2002; Li et al., 2023)). Let $\mathcal{H} = (I_t, r_t)_{t=1}^n$ be a sequence in $[A] \times [0, 1]$ with $n \geq 1$ and \mathcal{H}' be the sequence where the t' 'th sample of \mathcal{H} is replaced by $(I_{t'}, r_{t'})$. Error stability holds for a distribution $(I, r) \sim \mathcal{D}$ if there exists a $K > 0$ such that for any $\mathcal{H}, (I_t, r_t) \in ([A] \times [0, 1]), t \leq n$, and $\text{TF} \in \text{Alg}$, we have

$$|\mathbb{E}_{(I,r)} [\ell(r, \text{TF}(\widehat{r}(I)|\mathcal{H})) - \ell(r, \text{TF}(\widehat{r}(I)|\mathcal{H}'))]| \leq \frac{K}{n}.$$

Let ρ be a distance metric on Alg. Pairwise error stability holds if for all $\text{TF}, \text{TF}' \in \text{Alg}$ we have

$$|\mathbb{E}_{(x,y)} [\ell(r, \text{TF}(\widehat{r}(I)|\mathcal{H})) - \ell(r, \text{TF}'(\widehat{r}(I)|\mathcal{H})) - \ell(r, \text{TF}(\widehat{r}(I)|\mathcal{H}')) + \ell(r, \text{TF}'(\widehat{r}(I)|\mathcal{H}'))]| \leq \frac{K\rho(\text{TF}, \text{TF}')}{n}.$$

Now we present the Multi-task learning (MTL) risk of **PreDeToR- τ** .

Theorem 7.2. (**PreDeToR risk**) Suppose error stability Assumption 7.1 holds and assume loss function $\ell(\cdot, \cdot)$ is C -Lipschitz for all $r_t \in [0, B]$ and horizon $n \geq 1$. Let $\widehat{\text{TF}}$ be the empirical solution of (ERM) and $\mathcal{N}(\text{Alg}, \rho, \epsilon)$ be the covering number of the algorithm space Alg following Definition C.2 and C.3. Then with probability at least $1 - 2\delta$, the excess MTL risk of **PreDeToR- τ** is bounded by

$$\mathcal{R}_{\text{MTL}}(\widehat{\text{TF}}) \leq 4\frac{C}{\sqrt{nM}} + 2(B + K \log n)\sqrt{\frac{\log(\mathcal{N}(\text{Alg}, \rho, \epsilon)/\delta)}{cnM}},$$

where $\mathcal{N}(\text{Alg}, \rho, \epsilon)$ is the covering number of transformer $\widehat{\text{TF}}$ and $\epsilon = 1/\sqrt{nM}$.

The proof of Theorem 7.2 is provided in Appendix C.1. From Theorem 7.2 we see that in low-data regime with a small horizon n , as the number of tasks M increases the MTL risk decreases. We further discuss the stability factor K and covering number $\mathcal{N}(\text{Alg}, \rho, \epsilon)$ in Remark C.4, and C.5.

We now present the transfer learning risk of **PreDeToR- τ** for an unknown target task $g \sim \mathcal{T}$ with the test dataset $\mathcal{H}_g \sim \mathcal{D}_{\text{test}}(\cdot|g)$. Note that the test data distribution $\mathcal{D}_{\text{test}}(\cdot|g)$ is such that the actions are sampled following soft-LinUCB.

Theorem 7.3. (**Transfer risk**) Consider the setting of Theorem 7.2 and assume the training source tasks are independently drawn from task distribution \mathcal{T} . Let $\widehat{\text{TF}}$ be the empirical solution of (ERM) and $g \sim \mathcal{T}$. Define the expected excess transfer learning risk $\mathbb{E}_g[\mathcal{R}_g] = \mathbb{E}_g[\mathcal{L}_g(\widehat{\text{TF}})] - \arg \min_{\text{TF} \in \text{Alg}} \mathbb{E}_g[\mathcal{L}_g(\text{TF})]$. Then with probability at least $1 - 2\delta$, the $\mathbb{E}_g[\mathcal{R}_g] \leq 4\frac{C}{\sqrt{M}} + 2B\sqrt{\frac{\log(\mathcal{N}(\text{Alg}, \rho, \epsilon)/\delta)}{M}}$, where $\mathcal{N}(\text{Alg}, \rho, \epsilon)$ is the covering number of $\widehat{\text{TF}}$ and $\epsilon = \frac{1}{\sqrt{M}}$.

The proof is given in Appendix C.2. This shows that for the transfer learning risk of **PreDeToR- τ** (once trained on the M source tasks) scales with $\tilde{O}(1/\sqrt{M})$. This is because the unseen target task $g \sim \mathcal{T}$ induces a distribution shift, which, typically, cannot be mitigated with more samples n per task. A similar observation is provided in Lin et al. (2023). We further discuss this in Remark C.7. We also observe a similar phenomenon empirically; see the discussion in Appendix A.14.

8 CONCLUSIONS, LIMITATIONS AND FUTURE WORKS

In this paper, we studied the supervised pretraining of decision transformers in the multi-task structured bandit setting when the knowledge of the optimal action is unavailable. Moreover, our proposed methods **PreDeToR (- τ)** do not need to know the action representations or the reward structure and learn these in-context with the help of offline data. The **PreDeToR (- τ)** predict the reward for the next action of each action during pretraining and can generalize well in-context in several regimes spanning low-data, new actions, and structured bandit settings like linear, non-linear, bilinear, latent bandits. The **PreDeToR (- τ)** outperforms other in-context algorithms like **AD**, **DPT-greedy** in most of the experiments. Finally, we theoretically analyze **PreDeToR- τ** and show that pretraining it in M source tasks leads to a low expected excess error on a target task drawn from the same task distribution \mathcal{T} . In future, we want to extend our **PreDeToR (- τ)** to MDP setting (Sutton & Barto, 2018; Agarwal et al., 2019), and constraint MDP setting (Efroni et al., 2020; Gu et al., 2022).

REFERENCES

- 540
541
542 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic
543 bandits. *Advances in neural information processing systems*, 24, 2011.
- 544 Alekh Agarwal, Nan Jiang, Sham M Kakade, and Wen Sun. Reinforcement learning: Theory and
545 algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, 32, 2019.
- 546
547 Shipra Agrawal and Navin Goyal. Analysis of thompson sampling for the multi-armed bandit problem.
548 In *Conference on learning theory*, pp. 39–1. JMLR Workshop and Conference Proceedings, 2012.
- 549 Nabiha Asghar. Yelp dataset challenge: Review rating prediction. *arXiv preprint arXiv:1605.05362*,
550 2016.
- 551
552 Peter Auer and Ronald Ortner. Ucb revisited: Improved regret bounds for the stochastic multi-armed
553 bandit problem. *Periodica Mathematica Hungarica*, 61(1-2):55–65, 2010.
- 554 Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time Analysis of the Multiarmed Bandit
555 Problem. *Machine Learning*, 47(2):235–256, May 2002. ISSN 1573-0565. doi: 10.1023/A:
556 1013689704352. URL <https://doi.org/10.1023/A:1013689704352>.
- 557
558 Yoshua Bengio, Samy Bengio, and Jocelyn Cloutier. *Learning a synaptic learning rule*. Université
559 de Montréal, Département d’informatique et de recherche . . . , 1990.
- 560
561 C Bishop. Pattern recognition and machine learning. *Springer google schola*, 2:531–537, 2006.
- 562 Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine Learning*
563 *Research*, 2:499–526, 2002.
- 564
565 George EP Box and George C Tiao. *Bayesian inference in statistical analysis*. John Wiley & Sons,
566 2011.
- 567
568 David Brandfonbrener, Alberto Bietti, Jacob Buckman, Romain Laroche, and Joan Bruna. When
569 does return-conditioned supervised learning work for offline reinforcement learning? *Advances in*
Neural Information Processing Systems, 35:1542–1553, 2022.
- 570
571 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn,
572 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics
transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- 573
574 Sébastien Bubeck, Gilles Stoltz, Csaba Szepesvári, and Rémi Munos. Online optimization in x-armed
575 bandits. *Advances in Neural Information Processing Systems*, 21, 2008.
- 576
577 Sébastien Bubeck, Gilles Stoltz, and Jia Yuan Yu. Lipschitz bandits without the lipschitz constant. In
578 *Algorithmic Learning Theory: 22nd International Conference, ALT 2011, Espoo, Finland, October*
5-7, 2011. Proceedings 22, pp. 144–158. Springer, 2011.
- 579
580 Bradley P Carlin and Thomas A Louis. *Bayesian methods for data analysis*. CRC press, 2008.
- 581
582 Kamalika Chaudhuri, Prateek Jain, and Nagarajan Natarajan. Active heteroscedastic regression. In
International Conference on Machine Learning, pp. 694–702. PMLR, 2017.
- 583
584 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
585 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
586 modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- 587
588 Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International*
Conference on Machine Learning, pp. 844–853. PMLR, 2017.
- 589
590 Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff
591 functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and*
Statistics, pp. 208–214. JMLR Workshop and Conference Proceedings, 2011.
- 592
593 Zhongxiang Dai, Yao Shu, Arun Verma, Flint Xiaofeng Fan, Bryan Kian Hsiang Low, and Patrick
Jaillet. Federated neural bandit. *arXiv preprint arXiv:2205.14309*, 2022.

- 594 Rémy Degenne, Pierre Ménard, Xuedong Shang, and Michal Valko. Gamification of pure exploration
595 for linear bandits. In *International Conference on Machine Learning*, pp. 2432–2442. PMLR,
596 2020.
- 597 Yash Deshpande and Andrea Montanari. Linear bandits in high dimension and recommendation
598 systems. In *2012 50th Annual Allerton Conference on Communication, Control, and Computing*
599 (*Allerton*), pp. 1750–1754. IEEE, 2012.
- 600 Kefan Dong, Jiaqi Yang, and Tengyu Ma. Provable model-based nonlinear bandit and reinforcement
601 learning: Shelve optimism, embrace virtual curvature. *Advances in neural information processing*
602 *systems*, 34:26168–26182, 2021.
- 603 Yihan Du, Longbo Huang, and Wen Sun. Multi-task representation learning for pure exploration in
604 linear bandits. *arXiv preprint arXiv:2302.04441*, 2023.
- 605 Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RI
606 2: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*,
607 2016.
- 608 Yonathan Efroni, Shie Mannor, and Matteo Pirota. Exploration-exploitation in constrained mdps.
609 *arXiv preprint arXiv:2003.02189*, 2020.
- 610 Valerii Vadimovich Fedorov. *Theory of optimal experiments*. Elsevier, 2013.
- 611 Sarah Filippi, Olivier Cappé, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The
612 generalized linear case. *Advances in neural information processing systems*, 23, 2010.
- 613 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of
614 deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- 615 Vincent François-Lavet, Peter Henderson, Riashat Islam, Marc G Bellemare, Joelle Pineau, et al. An
616 introduction to deep reinforcement learning. *Foundations and Trends® in Machine Learning*, 11
617 (3-4):219–354, 2018.
- 618 Justin Fu, Sergey Levine, and Pieter Abbeel. One-shot learning of manipulation skills with online
619 dynamics adaptation and neural network priors. In *2016 IEEE/RSJ International Conference on*
620 *Intelligent Robots and Systems (IROS)*, pp. 4019–4026. IEEE, 2016.
- 621 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without
622 exploration. In *International conference on machine learning*, pp. 2052–2062. PMLR, 2019.
- 623 Yao Ge, Yuting Guo, Yuan-Chi Yang, Mohammed Ali Al-Garadi, and Abeer Sarker. Few-shot
624 learning for medical text: A systematic review. *arXiv preprint arXiv:2204.14081*, 2022.
- 625 Kamyar Ghasemipour, Shixiang Shane Gu, and Ofir Nachum. Why so pessimistic? estimating
626 uncertainties for offline rl through ensembles, and why their independence matters. *Advances in*
627 *Neural Information Processing Systems*, 35:18267–18281, 2022.
- 628 Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and
629 Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. *arXiv*
630 *preprint arXiv:2205.10330*, 2022.
- 631 Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Meta-
632 reinforcement learning of structured exploration strategies. *Advances in neural information*
633 *processing systems*, 31, 2018.
- 634 Samarth Gupta, Shreyas Chaudhari, Subhojyoti Mukherjee, Gauri Joshi, and Osman Yağan. A unified
635 approach to translate classical bandit algorithms to the structured bandit setting. *IEEE Journal on*
636 *Selected Areas in Information Theory*, 1(3):840–853, 2020.
- 637 F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm*
638 *transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- 639 Joey Hong, Branislav Kveton, Manzil Zaheer, Yinlam Chow, Amr Ahmed, and Craig Boutilier. Latent
640 bandits revisited. *Advances in Neural Information Processing Systems*, 33:13423–13433, 2020.

- 648 Joey Hong, Branislav Kveton, Sumeet Katariya, Manzil Zaheer, and Mohammad Ghavamzadeh.
649 Deep hierarchy in bandits. In *International Conference on Machine Learning*, pp. 8833–8851.
650 PMLR, 2022a.
- 651 Joey Hong, Branislav Kveton, Manzil Zaheer, and Mohammad Ghavamzadeh. Hierarchical bayesian
652 bandits. In *International Conference on Artificial Intelligence and Statistics*, pp. 7724–7741.
653 PMLR, 2022b.
- 654 Joey Hong, Branislav Kveton, Manzil Zaheer, Sumeet Katariya, and Mohammad Ghavamzadeh.
655 Multi-task off-policy learning from bandit feedback. In *International Conference on Machine
656 Learning*, pp. 13157–13173. PMLR, 2023.
- 657 Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence
658 modeling problem. *Advances in neural information processing systems*, 34:1273–1286, 2021.
- 659 Yiding Jiang, Evan Liu, Benjamin Eysenbach, J Zico Kolter, and Chelsea Finn. Learning options via
660 compression. *Advances in Neural Information Processing Systems*, 35:21184–21199, 2022.
- 661 Richard Arnold Johnson, Dean W Wichern, et al. Applied multivariate statistical analysis. 2002.
- 662 Kwang-Sung Jun, Rebecca Willett, Stephen Wright, and Robert Nowak. Bilinear bandits with
663 low-rank structure. In *International Conference on Machine Learning*, pp. 3163–3172. PMLR,
664 2019.
- 665 Daniel Justus, John Brennan, Stephen Bonner, and Andrew Stephen McGough. Predicting the
666 computational cost of deep learning models. In *2018 IEEE international conference on big data
667 (Big Data)*, pp. 3873–3882. IEEE, 2018.
- 668 Yue Kang, Cho-Jui Hsieh, and Thomas Chun Man Lee. Efficient frameworks for generalized low-rank
669 matrix bandit problems. *Advances in Neural Information Processing Systems*, 35:19971–19983,
670 2022.
- 671 Feyza Duman Keles, Pruthuvi Mahesakya Wijewardena, and Chinmay Hegde. On the computational
672 complexity of self-attention. In *International Conference on Algorithmic Learning Theory*, pp.
673 597–619. PMLR, 2023.
- 674 Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy
675 q-learning via bootstrapping error reduction. *Advances in Neural Information Processing Systems*,
676 32, 2019.
- 677 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
678 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- 679 Branislav Kveton, Csaba Szepesvári, Anup Rao, Zheng Wen, Yasin Abbasi-Yadkori, and S Muthukr-
680 ishnan. Stochastic low-rank bandits. *arXiv preprint arXiv:1712.04644*, 2017.
- 681 Jeongyeol Kwon, Yonathan Efroni, Constantine Caramanis, and Shie Mannor. Tractable optimality
682 in episodic latent mabs. *Advances in Neural Information Processing Systems*, 35:23634–23645,
683 2022.
- 684 Nicholas C Landolfi, Garrett Thomas, and Tengyu Ma. A model-based approach for sample-efficient
685 multi-task reinforcement learning. *arXiv preprint arXiv:1907.04964*, 2019.
- 686 Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald,
687 DJ Strouse, Steven Hansen, Angelos Filos, Ethan Brooks, et al. In-context reinforcement learning
688 with algorithm distillation. *arXiv preprint arXiv:2210.14215*, 2022.
- 689 Tor Lattimore and Csaba Szepesvári. An information-theoretic approach to minimax regret in partial
690 monitoring. In *Conference on Learning Theory*, pp. 2111–2139. PMLR, 2019.
- 691 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 692 Jonathan N Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and Emma
693 Brunskill. Supervised pretraining can learn in-context reinforcement learning. *arXiv preprint
694 arXiv:2306.14892*, 2023.

- 702 Kuang-Huei Lee, Ofir Nachum, Mengjiao Sherry Yang, Lisa Lee, Daniel Freeman, Sergio Guar-
703 rama, Ian Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, et al. Multi-game decision
704 transformers. *Advances in Neural Information Processing Systems*, 35:27921–27936, 2022.
- 705 Lanqing Li, Rui Yang, and Dijun Luo. Focal: Efficient fully-offline meta-reinforcement learning via
706 distance metric learning and behavior regularization. *arXiv preprint arXiv:2010.01112*, 2020.
- 707 Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to
708 personalized news article recommendation. In *Proceedings of the 19th international conference on*
709 *World wide web*, pp. 661–670, 2010.
- 710 Lihong Li, Yu Lu, and Dengyong Zhou. Provably optimal algorithms for generalized linear contextual
711 bandits. In *International Conference on Machine Learning*, pp. 2071–2080. PMLR, 2017.
- 712 Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers
713 as algorithms: Generalization and stability in in-context learning. In *International Conference on*
714 *Machine Learning*, pp. 19565–19594. PMLR, 2023.
- 715 Licong Lin, Yu Bai, and Song Mei. Transformers as decision makers: Provable in-context reinforce-
716 ment learning via supervised pretraining. *arXiv preprint arXiv:2310.08566*, 2023.
- 717 Evan Z Liu, Aditi Raghunathan, Percy Liang, and Chelsea Finn. Decoupling exploration and
718 exploitation for meta-reinforcement learning without sacrifices. In *International conference on*
719 *machine learning*, pp. 6925–6935. PMLR, 2021.
- 720 Xiaoqian Liu, Jianbin Jiao, and Junge Zhang. Self-supervised pretraining for decision foundation
721 model: Formulation, pipeline and challenges. *arXiv preprint arXiv:2401.00031*, 2023a.
- 722 Xin Liu, Daniel McDuff, Geza Kovacs, Isaac Galatzer-Levy, Jacob Sunshine, Jiening Zhan, Ming-
723 Zher Poh, Shun Liao, Paolo Di Achille, and Shwetak Patel. Large language models are few-shot
724 health learners. *arXiv preprint arXiv:2305.15525*, 2023b.
- 725 Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Off-policy policy gradient with
726 state distribution correction. *arXiv preprint arXiv:1904.08473*, 2019.
- 727 Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Provably good batch off-policy
728 reinforcement learning without great exploration. *Advances in neural information processing*
729 *systems*, 33:1264–1274, 2020.
- 730 Zhihan Liu, Hao Hu, Shenao Zhang, Hongyi Guo, Shuqi Ke, Boyi Liu, and Zhaoran Wang. Reason
731 for future, act for now: A principled framework for autonomous llm agents with provable sample
732 efficiency. *arXiv preprint arXiv:2309.17382*, 2023c.
- 733 Chris Lu, Yannick Schroecker, Albert Gu, Emilio Parisotto, Jakob Foerster, Satinder Singh, and
734 Feryal Behbahani. Structured state space models for in-context reinforcement learning. *arXiv*
735 *preprint arXiv:2303.03982*, 2023.
- 736 Yangyi Lu, Amirhossein Meisami, and Ambuj Tewari. Low-rank generalized linear bandit problems.
737 In *International Conference on Artificial Intelligence and Statistics*, pp. 460–468. PMLR, 2021.
- 738 Yi Ma, Chenjun Xiao, Hebin Liang, and Jianye Hao. Rethinking decision transformer via hierarchical
739 reinforcement learning. *arXiv preprint arXiv:2311.00267*, 2023.
- 740 Andrea Madotto, Zhaojiang Lin, Genta Indra Winata, and Pascale Fung. Few-shot bot: Prompt-based
741 learning for dialogue systems. *arXiv preprint arXiv:2110.08118*, 2021.
- 742 Stefan Magureanu, Richard Combes, and Alexandre Proutiere. Lipschitz bandits: Regret lower bound
743 and optimal algorithms. In *Conference on Learning Theory*, pp. 975–999. PMLR, 2014.
- 744 Odalric-Ambrym Maillard and Shie Mannor. Latent bandits. In *International Conference on Machine*
745 *Learning*, pp. 136–144. PMLR, 2014.
- 746 Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke
747 Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv*
748 *preprint arXiv:2202.12837*, 2022.

- 756 Suvir Mirchandani, Fei Xia, Pete Florence, Brian Ichter, Danny Driess, Montserrat Gonzalez Arenas,
757 Kanishka Rao, Dorsa Sadigh, and Andy Zeng. Large language models as general pattern machines.
758 *arXiv preprint arXiv:2307.04721*, 2023.
- 759
760 Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-
761 learner. *arXiv preprint arXiv:1707.03141*, 2017.
- 762 Eric Mitchell, Rafael Rafailov, Xue Bin Peng, Sergey Levine, and Chelsea Finn. Offline meta-
763 reinforcement learning with advantage weighting. In *International Conference on Machine*
764 *Learning*, pp. 7780–7791. PMLR, 2021.
- 765
766 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
767 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*
768 *arXiv:1312.5602*, 2013.
- 769 Subhojyoti Mukherjee, Qiaomin Xie, Josiah P Hanna, and Robert Nowak. Multi-task representation
770 learning for pure exploration in bilinear bandits. *arXiv preprint arXiv:2311.00327*, 2023.
- 771
772 Samuel Müller, Noah Hollmann, Sebastian Pineda Arango, Josif Grabocka, and Frank Hutter.
773 Transformers can do bayesian inference. *arXiv preprint arXiv:2112.10510*, 2021.
- 774 Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S Fearing, Pieter Abbeel, Sergey Levine, and
775 Chelsea Finn. Learning to adapt in dynamic, real-world environments through meta-reinforcement
776 learning. *arXiv preprint arXiv:1803.11347*, 2018.
- 777
778 Behnam Neyshabur, Ryota Tomioka, Ruslan Salakhutdinov, and Nathan Srebro. Geometry of
779 optimization and implicit regularization in deep learning. *arXiv preprint arXiv:1705.03071*, 2017.
- 780 Soumyabrata Pal, Arun Sai Suggala, Karthikeyan Shanmugam, and Prateek Jain. Optimal algorithms
781 for latent bandits with cluster structure. In *International Conference on Artificial Intelligence and*
782 *Statistics*, pp. 7540–7577. PMLR, 2023.
- 783
784 Theodore J Perkins and Doina Precup. Using options for knowledge transfer in reinforcement learning
785 title2, 1999.
- 786 Vitchyr H Pong, Ashvin V Nair, Laura M Smith, Catherine Huang, and Sergey Levine. Offline
787 meta-reinforcement learning with online self-supervision. In *International Conference on Machine*
788 *Learning*, pp. 17811–17829. PMLR, 2022.
- 789
790 Friedrich Pukelsheim. *Optimal design of experiments*. SIAM, 2006.
- 791 Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient off-policy
792 meta-reinforcement learning via probabilistic context variables. In *International conference on*
793 *machine learning*, pp. 5331–5340. PMLR, 2019.
- 794
795 Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel
796 Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist
797 agent. *arXiv preprint arXiv:2205.06175*, 2022.
- 798
799 Carlos Riquelme, George Tucker, and Jasper Snoek. Deep bayesian bandits showdown: An empirical
800 comparison of bayesian deep networks for thompson sampling. *arXiv preprint arXiv:1802.09127*,
2018.
- 801
802 Adam Roberts, Colin Raffel, Katherine Lee, Michael Matena, Noam Shazeer, Peter J Liu, Sharan
803 Narang, Wei Li, and Yanqi Zhou. Exploring the limits of transfer learning with a unified text-to-text
804 transformer. 2019.
- 805
806 Jonas Rothfuss, Dennis Lee, Ignasi Clavera, Tamim Asfour, and Pieter Abbeel. Promp: Proximal
meta-policy search. *arXiv preprint arXiv:1810.06784*, 2018.
- 807
808 Daniel J Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, Zheng Wen, et al. A tutorial on
thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.
- 809
Tom Schaul and Jürgen Schmidhuber. Metalearning. *Scholarpedia*, 5(6):4650, 2010.

- 810 Sina Semnani, Violet Yao, Heidi Zhang, and Monica Lam. Wikichat: Stopping the hallucination of
811 large language model chatbots by few-shot grounding on wikipedia. In *Findings of the Association*
812 *for Computational Linguistics: EMNLP 2023*, pp. 2387–2413, 2023.
- 813 Nur Muhammad Shafiullah, Zichen Cui, Ariuntuya Arty Altanzaya, and Lerrel Pinto. Behavior
814 transformers: Cloning k modes with one stone. *Advances in neural information processing systems*,
815 35:22955–22968, 2022.
- 816 Noah Y Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki, Michael Neunert,
817 Thomas Lampe, Roland Hafner, Nicolas Heess, and Martin Riedmiller. Keep doing what worked:
818 Behavioral modelling priors for offline reinforcement learning. *arXiv preprint arXiv:2002.08396*,
819 2020.
- 820 Viacheslav Sinii, Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, and Sergey Kolesnikov.
821 In-context reinforcement learning for variable action spaces. *arXiv preprint arXiv:2312.13327*,
822 2023.
- 823 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit
824 bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57,
825 2018.
- 826 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 827 William R Thompson. On the likelihood that one unknown probability exceeds another in view of
828 the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- 829 Sabina Tomkins, Peng Liao, Predrag Klasnja, Serena Yeung, and Susan Murphy. Rapidly person-
830 alizing mobile health treatment policies with limited data. *arXiv preprint arXiv:2002.09971*,
831 2020.
- 832 Michal Valko, Nathaniel Korda, Rémi Munos, Ilias Flaounas, and Nelo Cristianini. Finite-time
833 analysis of kernelised contextual bandits. *arXiv preprint arXiv:1309.6869*, 2013.
- 834 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz
835 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
836 *systems*, 30, 2017.
- 837 Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos,
838 Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *arXiv*
839 *preprint arXiv:1611.05763*, 2016.
- 840 Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning.
841 *arXiv preprint arXiv:1911.11361*, 2019.
- 842 Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context
843 learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- 844 Adam X Yang, Maxime Robeyns, Xi Wang, and Laurence Aitchison. Bayesian low-rank adaptation
845 for large language models. *arXiv preprint arXiv:2308.13111*, 2023.
- 846 Jiaqi Yang, Wei Hu, Jason D Lee, and Simon S Du. Impact of representation learning in linear bandits.
847 *arXiv preprint arXiv:2010.06531*, 2020.
- 848 Jiaqi Yang, Wei Hu, Jason D Lee, and Simon Shaolei Du. Impact of representation learning in linear
849 bandits. In *International Conference on Learning Representations*, 2021.
- 850 Jiaqi Yang, Qi Lei, Jason D Lee, and Simon S Du. Nearly minimax algorithms for linear bandits with
851 shared representation. *arXiv preprint arXiv:2203.15664*, 2022a.
- 852 Lin Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and
853 regret bound. In *International Conference on Machine Learning*, pp. 10746–10756. PMLR, 2020.
- 854 Mengjiao Yang, Dale Schuurmans, Pieter Abbeel, and Ofir Nachum. Dichotomy of control: Separat-
855 ing what you can control from what you cannot. *arXiv preprint arXiv:2210.13435*, 2022b.

864 Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn.
865 Combo: Conservative offline model-based policy optimization. *Advances in neural information*
866 *processing systems*, 34:28954–28967, 2021.
867

868 Vinicius Zambaldi, David Raposo, Adam Santoro, Victor Bapst, Yujia Li, Igor Babuschkin, Karl
869 Tuyls, David Reichert, Timothy Lillicrap, Edward Lockhart, et al. Relational deep reinforcement
870 learning. *arXiv preprint arXiv:1806.01830*, 2018.

871 Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: Llm
872 agents are experiential learners. *arXiv preprint arXiv:2308.10144*, 2023.
873

874 Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based exploration.
875 In *International Conference on Machine Learning*, pp. 11492–11502. PMLR, 2020.

876 Qiuyu Zhu and Vincent Tan. Thompson sampling algorithms for mean-variance bandits. In *Interna-*
877 *tional Conference on Machine Learning*, pp. 11599–11608. PMLR, 2020.
878

879 Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and
880 Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning.
881 *arXiv preprint arXiv:1910.08348*, 2019.
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918 A APPENDIX

919 A.1 RELATED WORKS

920 In this section, we briefly discuss related works.

921 In-context decision making (Laskin et al., 2022; Lee et al., 2023) has emerged as an attractive alterna-
 922 tive in Reinforcement Learning (RL) compared to updating the model parameters after collection of
 923 new data (Mnih et al., 2013; François-Lavet et al., 2018). In RL the contextual data takes the form of
 924 state-action-reward tuples representing a dataset of interactions with an unknown environment (task).
 925 In this paper, we will refer to this as the in-context data. Recall that in many real-world settings, the
 926 underlying task can be structured with correlated features, and the reward can be highly non-linear.
 927 So specialized bandit algorithms fail to learn in these tasks. To circumvent this issue, a learner can
 928 first collect in-context data consisting of just action indices I_t and rewards r_t . Then it can leverage
 929 the representation learning capability of deep neural networks to learn a pattern across the in-context
 930 data and subsequently derive a near-optimal policy (Lee et al., 2023; Mirchandani et al., 2023). We
 931 refer to this learning framework as an in-context decision-making setting.

932 The in-context decision-making setting of Sinii et al. (2023) also allows changing the action space
 933 by learning an embedding over the action space yet also requires the optimal action during training.
 934 In contrast we do not require the optimal action as well as show that we can generalize to new
 935 actions without learning an embedding over them. Similarly, Lin et al. (2023) study the in-context
 936 decision-making setting of Laskin et al. (2022); Lee et al. (2023), but they also require a greedy
 937 approximation of the optimal action. The Ma et al. (2023) also studies a similar setting for hierarchical
 938 RL where they stitch together sub-optimal trajectories and predict the next action during test time.
 939 Similarly, Liu et al. (2023c) studies the in-context decision-making setting to predict action instead
 940 of learning a reward correlation from a short horizon setting. In contrast we do not require a greedy
 941 approximation of the optimal action, deal with short horizon setting and changing action sets during
 942 training and testing, and predict the estimated means of the actions instead of predicting the optimal
 943 action. A survey of the in-context decision-making approaches can be found in Liu et al. (2023a).

944 In the in-context decision-making setting, the learning model is first trained on supervised input-
 945 output examples with the in-context data during training. Then during test time, the model is asked to
 946 complete a new input (related to the context provided) without any update to the model parameters
 947 (Xie et al., 2021; Min et al., 2022). Motivated by this, Lee et al. (2023) recently proposed the
 948 Decision Pretrained Transformers (DPT) that exhibit the following properties: (1) During supervised
 949 pretraining of DPT, predicting optimal actions alone gives rise to near-optimal decision-making
 950 algorithms for unforeseen task during test time. Note that DPT does not update model parameters
 951 during test time and, therefore, conducts in-context learning on the unforeseen task. (2) DPT improves
 952 over the in-context data used to pretrain it by exploiting latent structure. However, DPT either requires
 953 the optimal action during training or if it needs to approximate the optimal action. For approximating
 954 the optimal action, it requires a large amount of data from the underlying task.

955 At the same time, learning the underlying data pattern from a few examples during training is
 956 becoming more relevant in many domains like chatbot interaction (Madotto et al., 2021; Semnani
 957 et al., 2023), recommendation systems, healthcare (Ge et al., 2022; Liu et al., 2023b), etc. This is
 958 referred to as few-shot learning. However, most current RL decision-making systems (including
 959 in-context learners like DPT) require an enormous amount of data to learn a good policy.

960 The in-context learning framework is related to the meta-learning framework (Bengio et al., 1990;
 961 Schaul & Schmidhuber, 2010). Broadly, these techniques aim to learn the underlying latent shared
 962 structure within the training distribution of tasks, facilitating faster learning of novel tasks during
 963 test time. In the context of decision-making and reinforcement learning (RL), there exists a frequent
 964 choice regarding the specific 'structure' to be learned, be it the task dynamics (Fu et al., 2016;
 965 Nagabandi et al., 2018; Landolfi et al., 2019), a task context identifier (Rakelly et al., 2019; Zintgraf
 966 et al., 2019; Liu et al., 2021), or temporally extended skills and options (Perkins & Precup, 1999;
 967 Gupta et al., 2018; Jiang et al., 2022).

968 However, as we noted in the Section 1, one can do a greedy approximation of the optimal action
 969 from the historical data using a weak demonstrator and a neural network policy (Finn et al., 2017;
 970 Rothfuss et al., 2018). Moreover, the in-context framework generally is more agnostic where it learns

the policy of the demonstrator (Duan et al., 2016; Wang et al., 2016; Mishra et al., 2017). Note that both **DPT-greedy** and **PreDeToR** are different than algorithmic distillation (Laskin et al., 2022; Lu et al., 2023) as they do not distill an existing RL algorithm. moreover, in contrast to **DPT-greedy** which is trained to predict the optimal action, the **PreDeToR** is trained to predict the reward for each of the actions. This enables the **PreDeToR** (similar to **DPT-greedy**) to show to potentially emergent online and offline strategies at test time that automatically align with the task structure, resembling posterior sampling.

As we discussed in the Section 1, in decision-making, RL, and imitation learning the transformer models are trained using autoregressive action prediction (Yang et al., 2023). Similar methods have also been used in Large language models (Vaswani et al., 2017; Roberts et al., 2019). One of the more notable examples is the Decision Transformers (abbreviated as DT) which utilizes a transformer to autoregressively model sequences of actions from offline experience data, conditioned on the achieved return (Chen et al., 2021; Janner et al., 2021). This approach has also been shown to be effective for multi-task settings (Lee et al., 2022), and multi-task imitation learning with transformers (Reed et al., 2022; Brohan et al., 2022; Shafullah et al., 2022). However, the DT methods are not known to improve upon their in-context data, which is the main thrust of this paper (Brandfonbrener et al., 2022; Yang et al., 2022b).

Our work is also closely related to the offline RL setting. In offline RL, the algorithms can formulate a policy from existing data sets of state, action, reward, and next-state interactions. Recently, the idea of pessimism has also been introduced in an offline setting to address the challenge of distribution shift (Kumar et al., 2020; Yu et al., 2021; Liu et al., 2020; Ghasemipour et al., 2022). Another approach to solve this issue is policy regularization (Fujimoto et al., 2019; Kumar et al., 2019; Wu et al., 2019; Siegel et al., 2020; Liu et al., 2019), or reuse data for related task (Li et al., 2020; Mitchell et al., 2021), or additional collection of data along with offline data (Pong et al., 2022). However, all of these approaches still have to take into account the issue of distributional shifts. In contrast **PreDeToR** and **DPT-greedy** leverages the decision transformers to avoid these issues. Both of these methods can also be linked to posterior sampling. Such connections between sequence modeling with transformers and posterior sampling have also been made in Chen et al. (2021); Müller et al. (2021); Lee et al. (2023); Yang et al. (2023).

A.2 EXPERIMENTAL SETTING INFORMATION AND DETAILS OF BASELINES

In this section, we describe in detail the experimental settings and some baselines.

A.2.1 EXPERIMENTAL DETAILS

Linear Bandit: We consider the setting when $f(\mathbf{x}, \boldsymbol{\theta}_*) = \mathbf{x}^\top \boldsymbol{\theta}_*$. Here $\mathbf{x} \in \mathbb{R}^d$ is the action feature and $\boldsymbol{\theta}_* \in \mathbb{R}^d$ is the hidden parameter. For every experiment, we first generate tasks from \mathcal{T}_{pre} . Then we sample a fixed set of actions from $\mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)$ in \mathbb{R}^d and this constitutes the features. Then for each task $m \in [M]$ we sample $\boldsymbol{\theta}_{m,*} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)$ to produce the means $\mu(m, a) = \langle \boldsymbol{\theta}_{m,*}, \mathbf{x}(m, a) \rangle$ for $a \in \mathcal{A}$ and $m \in [M]$. Finally, note that we do not shuffle the data as the order matters. Also in this setting $\mathbf{x}(m, a)$ for each $a \in \mathcal{A}$ is fixed for all tasks m .

Non-Linear Bandit: We now consider the setting when $f(\mathbf{x}, \boldsymbol{\theta}_*) = 1/(1 + 0.5 \cdot \exp(2 \cdot \exp(-\mathbf{x}^\top \boldsymbol{\theta}_*)))$. Again, here $\mathbf{x} \in \mathbb{R}^d$ is the action feature, and $\boldsymbol{\theta}_* \in \mathbb{R}^d$ is the hidden parameter. Note that this is different than the generalized linear bandit setting (Filippi et al., 2010; Li et al., 2017). Again for every experiment, we first generate tasks from \mathcal{T}_{pre} . Then we sample a fixed set of actions from $\mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)$ in \mathbb{R}^d and this constitutes the features. Then for each task $m \in [M]$ we sample $\boldsymbol{\theta}_{m,*} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d/d)$ to produce the means $\mu(m, a) = 1/(1 + 0.5 \cdot \exp(2 \cdot \exp(-\mathbf{x}(m, a)^\top \boldsymbol{\theta}_{m,*})))$ for $a \in \mathcal{A}$ and $m \in [M]$. Again note that in this setting $\mathbf{x}(m, a)$ for each $a \in \mathcal{A}$ is fixed for all tasks m .

We use NVIDIA GeForce RTX 3090 GPU with 24GB RAM to load the GPT 2 Large Language Model. This requires less than 2GB RAM without data, and with large context may require as much as 20GB RAM.

A.2.2 DETAILS OF BASELINES

(1) **Thomp**: This baseline is the stochastic A -action bandit Thompson Sampling algorithm from Thompson (1933); Agrawal & Goyal (2012); Russo et al. (2018); Zhu & Tan (2020). We briefly describe the algorithm below: At every round t and each action a , **Thomp** samples $\gamma_{m,t}(a) \sim \mathcal{N}(\hat{\mu}_{m,t-1}(a), \sigma^2/N_{m,t-1}(a))$, where $N_{m,t-1}(a)$ is the number of times the action a has been selected till $t-1$, and $\hat{\mu}_{m,t-1}(a) = \frac{\sum_{s=1}^{t-1} \hat{r}_{m,s} \mathbf{1}(I_s=a)}{N_{m,t-1}(a)}$ is the empirical mean. Then the action selected at round t is $I_t = \arg \max_a \gamma_{m,t}(a)$. Observe that **Thomp** is not a deterministic algorithm like **UCB** (Auer et al., 2002). So we choose **Thomp** as the weak demonstrator π^w because it is more exploratory than **UCB** and also chooses the optimal action, $a_{m,*}$, a sufficiently large number of times. **Thomp** is a weak demonstrator as it does not have access to the feature set \mathcal{X} for any task m .

(2) **LinUCB**: (Linear Upper Confidence Bound): This baseline is the Upper Confidence Bound algorithm for the linear bandit setting that selects the action I_t at round t for task m that is most optimistic and reduces the uncertainty of the task unknown parameter $\theta_{m,*}$. To balance exploitation and exploration between choosing different items the **LinUCB** computes an upper confidence value to the estimated mean of each action $\mathbf{x}_{m,a} \in \mathcal{X}$. This is done as follows: At every round t for task m , it calculates the ucb value $B_{m,a,t}$ for each action $\mathbf{x}_{m,a} \in \mathcal{X}$ such that $B_{m,a,t} = \mathbf{x}_{m,a}^\top \hat{\theta}_{m,t-1} + \alpha \|\mathbf{x}_{m,a}\|_{\Sigma_{m,t-1}^{-1}}$ where $\alpha > 0$ is a constant and $\hat{\theta}_{m,t}$ is the estimate of the model parameter $\theta_{m,*}$ at round t . Here, $\Sigma_{m,t-1} = \sum_{s=1}^{t-1} \mathbf{x}_{m,s} \mathbf{x}_{m,s}^\top + \lambda \mathbf{I}_d$ is the data covariance matrix or the arms already tried. Then it chooses $I_t = \arg \max_a B_{m,a,t}$. Note that **LinUCB** is a *strong* demonstrator that we give oracle access to the features of each action; other algorithms do not observe the features. Hence, in linear bandits, **LinUCB** provides an approximate upper bound on the performance of all algorithms.

(3) **MLinGreedy**: This is the multi-task linear regression bandit algorithm proposed by Yang et al. (2021). This algorithm assumes that there is a common low dimensional feature extractor $\mathbf{B} \in \mathbb{R}^{k \times d}$, $k \leq d$ shared between the tasks and the rewards per task m are linearly dependent on a hidden parameter $\theta_{m,*}$. Under a diversity assumption (which may not be satisfied in real data) and $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_M]$ they assume $\Theta = [\theta_{1,*}, \dots, \theta_{M,*}] = \mathbf{B}\mathbf{W}$. During evaluation **MLinGreedy** estimates the $\hat{\mathbf{B}}$ and $\hat{\mathbf{W}}$ from training data and fit $\hat{\theta}_m = \hat{\mathbf{B}}\hat{\mathbf{w}}_m$ per task and selects action greedily based on $I_{m,t} = \arg \max_a \mathbf{x}_{m,a}^\top \hat{\theta}_{m,*}$. Finally, note that **MLinGreedy** requires access to the action features to estimate $\hat{\theta}_m$ and select actions as opposed to **DPT**, **AD**, and **PreDeToR**.

A.3 EMPIRICAL STUDY: BILINEAR BANDITS

In this section, we discuss the performance of **PreDeToR** against the other baselines in the bilinear setting. Again note that the number of tasks $M_{\text{pre}} \gg A \geq n$. Through this experiment, we want to evaluate the performance of **PreDeToR** to exploit the underlying latent structure and reward correlation when the horizon is small, the number of tasks is large, and understand its performance in the bilinear bandit setting (Jun et al., 2019; Lu et al., 2021; Kang et al., 2022; Mukherjee et al., 2023). Note that this setting also goes beyond the linear feedback model (Abbasi-Yadkori et al., 2011; Lattimore & Szepesvári, 2020) and is related to matrix bandits (Yang & Wang, 2020).

Bilinear bandit setting: In the bilinear bandits the learner is provided with two sets of action sets, $\mathcal{X} \subseteq \mathbb{R}^{d_1}$ and $\mathcal{Z} \subseteq \mathbb{R}^{d_2}$ which are referred to as the left and right action sets. At every round t the learner chooses a pair of actions $\mathbf{x}_t \in \mathcal{X}$ and $\mathbf{z}_t \in \mathcal{Z}$ and observes a reward

$$r_t = \mathbf{x}_t^\top \Theta_* \mathbf{z}_t + \eta_t$$

where $\Theta_* \in \mathbb{R}^{d_1 \times d_2}$ is the unknown hidden matrix which is also low-rank. The η_t is a σ^2 sub-Gaussian noise. In the multi-task bilinear bandit setting we now have a set of M tasks where the reward for the m -th task at round t is given by

$$r_{m,t} = \mathbf{x}_{m,t}^\top \Theta_{m,*} \mathbf{z}_{m,t} + \eta_{m,t}.$$

Here $\Theta_{m,*} \in \mathbb{R}^{d_1 \times d_2}$ is the unknown hidden matrix for each task m , which is also low-rank. The $\eta_{m,t}$ is a σ^2 sub-Gaussian noise. Let κ be the rank of each of these matrices $\Theta_{m,*}$.

A special case is the rank 1 structure where $\Theta_{m,*} = \theta_{m,*} \theta_{m,*}^\top$ where $\Theta_{m,*} \in \mathbb{R}^{d \times d}$ and $\theta_{m,*} \in \mathbb{R}^d$ for each task m . Let the left and right action sets be also same such that $\mathbf{x}_{m,t} \in \mathcal{X} \subseteq \mathbb{R}^d$. Observe then that the reward for the m -th task at round t is given by

$$r_{m,t} = \mathbf{x}_{m,t}^\top \Theta_{m,*} \mathbf{x}_{m,t} + \eta_{m,t} = (\mathbf{x}_{m,t}^\top \theta_{m,*})^2 + \eta_{m,t}.$$

This special case is studied in [Chaudhuri et al. \(2017\)](#).

Baselines: We again implement the same baselines discussed in Section 4. The baselines are [PreDeToR](#), [PreDeToR- \$\tau\$](#) , [DPT-greedy](#), and [Thomp](#). Note that we do not implement the [LinUCB](#) and [MLinGreedy](#) for the bilinear bandit setting. However, we now implement the [LowOFUL](#) ([Jun et al., 2019](#)) which is optimal in the bilinear bandit setting.

LowOFUL: The [LowOFUL](#) algorithm first estimates the unknown parameter $\Theta_{m,*}$ for each task m using E-optimal design ([Pukelsheim, 2006](#); [Fedorov, 2013](#); [Jun et al., 2019](#)) for n_1 rounds. Let $\hat{\Theta}_{m,n_1}$ be the estimate of $\Theta_{m,*}$ at the end of n_1 rounds. Let the SVD of $\hat{\Theta}_{m,n_1}$ be given by $\text{SVD}(\hat{\Theta}_{m,n_1}) = \hat{\mathbf{U}}_{m,n_1} \hat{\mathbf{S}}_{m,n_1} \hat{\mathbf{V}}_{m,n_1}^\top$. Then [LowOFUL](#) rotates the actions as follows:

$$\mathcal{X}'_m = \left\{ \left[\hat{\mathbf{U}}_{m,n_1} \hat{\mathbf{U}}_{m,n_1}^\perp \right]^\top \mathbf{x}_m : \mathbf{x}_m \in \mathcal{X} \right\} \text{ and } \mathcal{Z}'_m = \left\{ \left[\hat{\mathbf{V}}_{m,n_1} \hat{\mathbf{V}}_{m,n_1}^\perp \right]^\top \mathbf{z}_m : \mathbf{z}_m \in \mathcal{Z} \right\}.$$

Then defines a vectorized action set for each task m so that the last $(d_1 - \kappa) \cdot (d_2 - \kappa)$ components are from the complementary subspaces:

$$\tilde{\mathcal{A}}_m = \left\{ \left[\text{vec}(\mathbf{x}_{m,1:\kappa} \mathbf{z}_{m,1:\kappa}^\top); \text{vec}(\mathbf{x}_{m,\kappa+1:d_1} \mathbf{z}_{m,1:\kappa}^\top); \text{vec}(\mathbf{x}_{m,1:\kappa} \mathbf{z}_{m,\kappa+1:d_2}^\top); \text{vec}(\mathbf{x}_{m,\kappa+1:d_1} \mathbf{z}_{m,\kappa+1:d_2}^\top) \right] \in \mathbb{R}^{d_1 d_2} : \mathbf{x}_m \in \mathcal{X}'_m, \mathbf{z}_m \in \mathcal{Z}'_m \right\}.$$

Finally for $n_2 = n - n_1$ rounds, [LowOFUL](#) invokes the specialized OFUL algorithm ([Abbasi-Yadkori et al., 2011](#)) for the rotated action set $\tilde{\mathcal{A}}_m$ with the low dimension $k = (d_1 + d_2) \kappa - \kappa^2$. Note that the [LowOFUL](#) runs the per-task low dimensional OFUL algorithm rather than learning the underlying structure across the tasks ([Mukherjee et al., 2023](#)).

Outcomes: We first discuss the main outcomes of our experimental results for increasing the horizon:

Finding 5: [PreDeToR \(- \$\tau\$ \)](#) outperforms [DPT-greedy](#), [AD](#), and matches the performance of [LowOFUL](#) in bilinear bandit setting.

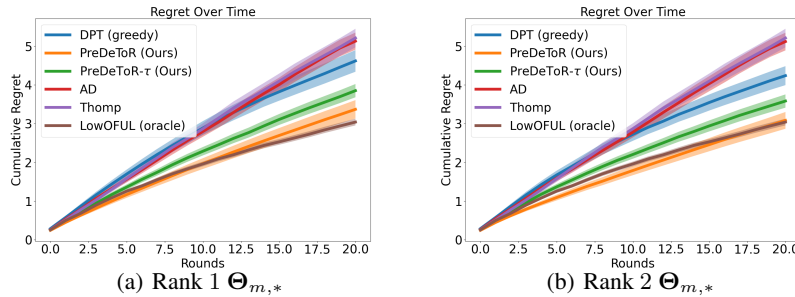


Figure 5: Experiment with bilinear bandits. The y-axis shows the cumulative regret.

Experimental Result: We observe these outcomes in Figure 5. In Figure 5(a) we experiment with rank 1 hidden parameter $\Theta_{m,*}$ and set horizon $n = 20$, $M_{\text{pre}} = 200000$, $M_{\text{test}} = 200$, $A = 30$, and $d = 5$. In Figure 5(b) we experiment with rank 2 hidden parameter $\Theta_{m,*}$ and set horizon $n = 20$, $M_{\text{pre}} = 250000$, $M_{\text{test}} = 200$, $A = 25$, and $d = 5$. Again, the demonstrator π^w is the [Thomp](#) algorithm. We observe that [PreDeToR](#) has lower cumulative regret than [DPT-greedy](#), [AD](#) and [Thomp](#). Note that for any task m for the horizon 20 the [Thomp](#) will be able to sample all the actions at most once. Note that for this small horizon setting the [DPT-greedy](#) does not have a good estimation of $\hat{a}_{m,t,*}$ which results in a poor prediction of optimal action $\hat{a}_{m,t,*}$. In contrast [PreDeToR](#) learns the

1134 correlation of rewards across tasks and can perform well. Observe from Figure 5(a), and 5(b) that
 1135 **PreDeToR** has lower regret than **Thomp** and matches **LowOFUL**. Also, in this low-data regime it
 1136 is not enough for **LowOFUL** to learn the underlying $\Theta_{m,*}$ with high precision. Hence, **PreDeToR**
 1137 also has slightly lower regret than **LowOFUL**. Note that the main objective of **AD** is to match the
 1138 performance of its demonstrator. Most importantly it shows that **PreDeToR** can exploit the underlying
 1139 latent structure and reward correlation better than **DPT-greedy**, and **AD**.

1141 A.4 EMPIRICAL STUDY: LATENT BANDITS

1142 In this section, we discuss the performance of **PreDeToR** ($-\tau$) against the other baselines in the
 1143 latent bandit setting and create a generalized bilinear bandit setting. Note that the number of tasks
 1144 $M_{\text{pre}} \gg A \geq n$. Using this experiment, we want to evaluate the ability of **PreDeToR** ($-\tau$) to exploit
 1145 the underlying reward correlation when the horizon is small, the number of tasks is large, and
 1146 understand its performance in the latent bandit setting (Hong et al., 2020; Maillard & Mannor, 2014;
 1147 Pal et al., 2023; Kveton et al., 2017). We create a latent bandit setting which generalizes the bilinear
 1148 bandit setting (Jun et al., 2019; Lu et al., 2021; Kang et al., 2022; Mukherjee et al., 2023). Again note
 1149 that this setting also goes beyond the linear feedback model (Abbasi-Yadkori et al., 2011; Lattimore
 1150 & Szepesvári, 2020) and is related to matrix bandits (Yang & Wang, 2020).

1151 **Latent bandit setting:** In this special multi-task latent bandits the learner is again provided with two
 1152 sets of action sets, $\mathcal{X} \subseteq \mathbb{R}^{d_1}$ and $\mathcal{Z} \subseteq \mathbb{R}^{d_2}$ which are referred to as the left and right action sets. The
 1153 reward for the m -th task at round t is given by

$$1154 r_{m,t} = \mathbf{x}_{m,t}^\top \underbrace{(\Theta_{m,*} + \mathbf{U}\mathbf{V}^\top)}_{\mathbf{Z}_{m,*}} \mathbf{z}_{m,t} + \eta_{m,t}.$$

1155 Here $\Theta_{m,*} \in \mathbb{R}^{d_1 \times d_2}$ is the unknown hidden matrix for each task m , which is also low-rank.
 1156 Additionally, all the tasks share a *common latent parameter matrix* $\mathbf{U}\mathbf{V}^\top \in \mathbb{R}^{d_1 \times d_2}$ which is also
 1157 low rank. Hence the learner needs to learn the latent parameter across the tasks hence the name latent
 1158 bandits. Finally, the $\eta_{m,t}$ is a σ^2 sub-Gaussian noise. Let κ be the rank of each of these matrices
 1159 $\Theta_{m,*}$ and $\mathbf{U}\mathbf{V}^\top$. Again special case is the rank 1 structure where the reward for the m -th task at
 1160 round t is given by

$$1161 r_{m,t} = \mathbf{x}_{m,t}^\top \underbrace{(\theta_{m,*} \theta_{m,*}^\top + \mathbf{u}\mathbf{v}^\top)}_{\mathbf{Z}_{m,*}} \mathbf{x}_{m,t} + \eta_{m,t}.$$

1162 where $\theta_{m,*} \in \mathbb{R}^d$ for each task m and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$. Note that the left and right action sets are the same
 1163 such that $\mathbf{x}_{m,t} \in \mathcal{X} \subseteq \mathbb{R}^d$.

1164 **Baselines:** We again implement the same baselines discussed in Section 4. The baselines are
 1165 **PreDeToR**, **PreDeToR- τ** , **DPT-greedy**, **AD**, **Thomp**, and **LowOFUL**. However, we now implement a
 1166 special **LowOFUL** (stated in Appendix A.3) which has knowledge of the shared latent parameters \mathbf{U} ,
 1167 and \mathbf{V} . We call this the **LowOFUL (oracle)** algorithm. Therefore **LowOFUL (oracle)** has knowledge
 1168 of the problem parameters in the latent bandit setting and hence the name. Again note that we do not
 1169 implement the **LinUCB** and **MLinGreedy** for the latent bandit setting.

1170 **Outcomes:** We first discuss the main outcomes of our experimental results for increasing the horizon:

1171 **Finding 6:** **PreDeToR** ($-\tau$) outperforms **DPT-greedy**, **AD**, and matches the performance of
 1172 **LowOFUL (oracle)** in latent bandit setting.

1173 **Experimental Result:** We observe these outcomes in Figure 6. In Figure 6(a) we experiment with
 1174 rank 1 hidden parameter $\theta_{m,*} \theta_{m,*}^\top$ and latent parameters $\mathbf{u}\mathbf{v}^\top$ shared across the tasks and set horizon
 1175 $n = 20$, $M_{\text{pre}} = 200000$, $M_{\text{test}} = 200$, $A = 30$, and $d = 5$. In Figure 6(b) we experiment with rank
 1176 2 hidden parameter $\Theta_{m,*}$, and latent parameters $\mathbf{U}\mathbf{V}^\top$ and set horizon $n = 20$, $M_{\text{pre}} = 250000$,
 1177 $M_{\text{test}} = 200$, $A = 25$, and $d = 5$. In Figure 6(c) we experiment with rank 3 hidden parameter $\Theta_{m,*}$,
 1178 and latent parameters $\mathbf{U}\mathbf{V}^\top$ and set horizon $n = 20$, $M_{\text{pre}} = 300000$, $M_{\text{test}} = 200$, $A = 25$, and
 1179 $d = 5$. Again, the demonstrator π^w is the **Thomp** algorithm. We observe that **PreDeToR** ($-\tau$) has

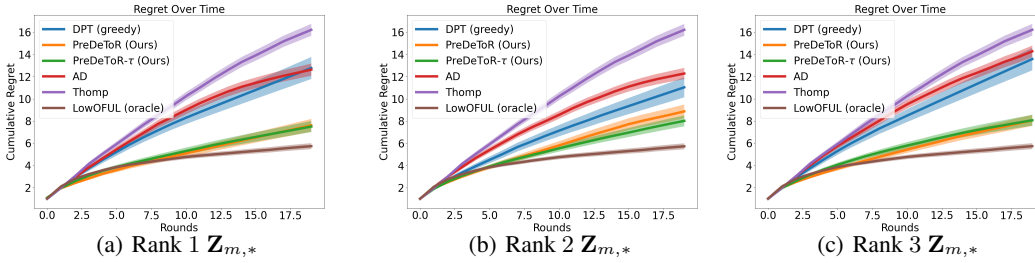


Figure 6: Experiment with latent bandits. The y-axis shows the cumulative regret.

lower cumulative regret than **DPT-greedy**, **AD** and **Thomp**. Note that for any task m for the horizon 20 the **Thomp** will be able to sample all the actions at most once. Note that for this small horizon setting the **DPT-greedy** does not have a good estimation of $\hat{a}_{m,*}$ which results in a poor prediction of optimal action $\hat{a}_{m,t,*}$. In contrast **PreDeToR** ($-\tau$) learns the correlation of rewards across tasks and is able to perform well. Observe from Figure 6(a), 6(b), and 6(c) that **PreDeToR** has lower regret than **Thomp** and has regret closer to **LowOFUL (oracle)** which has access to the problem-dependent parameters. Hence, **LowOFUL (oracle)** outperforms **PreDeToR** ($-\tau$) in this setting. This shows that **PreDeToR** is able to exploit the underlying latent structure and reward correlation better than **DPT-greedy**, and **AD**.

A.5 CONNECTION BETWEEN **PREDETOR** AND LINEAR MULTIVARIATE GAUSSIAN MODEL

In this section, we try to understand the behavior of **PreDeToR** and its ability to exploit the reward correlation across tasks under a *linear multivariate Gaussian model*. In this model, the hidden task parameter, θ_* , is a random variable drawn from a multi-variate Gaussian distribution (Bishop, 2006) and the feedback follows a linear model. We study this setting since we can estimate the Linear Minimum Mean Square Estimator (LMMSE) in this setting (Carlin & Louis, 2008; Box & Tiao, 2011). This yields a posterior prediction for the mean of each action over all tasks on average, by leveraging the linear structure when θ_* is drawn from a multi-variate Gaussian distribution. So we can compare the performance of **PreDeToR** against such an LMMSE and evaluate whether it is exploiting the underlying linear structure and the reward correlation across tasks. We summarize this as follows:

Finding 7: **PreDeToR** learns the reward correlation covariance matrix from the in-context training data $\mathcal{H}_{\text{train}}$ and acts greedily on it.

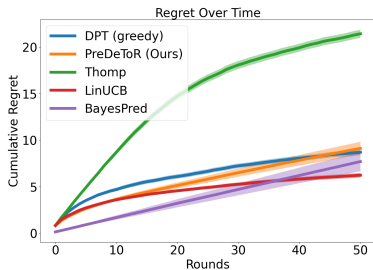


Figure 7: **BayesPred** Performance

PreDeToR for fair comparison. Then define the reward vector $\mathbf{Y}_n \in \mathbb{R}^n$ where the t -th component is the reward r_t observed for the action I_t for $t \in [n]$ in the pretraining data. Define the diagonal matrix

Consider the linear feedback setting consisting of A actions and the hidden task parameter $\theta_* \sim \mathcal{N}(0, \sigma_\theta^2 \mathbf{I}_A)$. The reward of the action \mathbf{x}_t at round t is given by $r_t = \mathbf{x}_t^\top \theta_* + \eta_t$, where η_t is σ^2 sub-Gaussian. Let π^w collect n rounds of pretraining in-context data and observe $\{I_t, r_t\}_{t=1}^n$. Let $N_n(a)$ denote the total number of times the action a is sampled for n rounds. Note that we drop the task index m in these notations as the random variable θ_* corresponds to the task. Define the matrix $\mathbf{H}_n \in \mathbb{R}^{n \times A}$ where the t -th row represents the action I_t for $t \in [n]$. The t -th row of \mathbf{H}_n is a one-hot vector with the I_t -th component being 1. We represent each action by one hot vector because we assume that this LMMSE does not have access to the feature vectors of the actions similar to the

$\mathbf{D}_A \in \mathbb{R}^{A \times A}$ estimated from pretraining data as follows

$$\mathbf{D}_A(i, i) = \begin{cases} \frac{\sigma^2}{N_n(a)}, & \text{if } N_n(a) > 0 \\ = 0, & \text{if } N_n(a) = 0 \end{cases} \quad (6)$$

where the reward noise being σ^2 sub-Gaussian is known. Finally define the estimated reward covariance matrix $\mathbf{S}_A \in \mathbb{R}^{A \times A}$ as $\mathbf{S}_A(a, a') = \hat{\mu}_n(a)\hat{\mu}_n(a')$, where $\hat{\mu}_n(a)$ is the empirical mean of action a estimated from the pretraining data. This matrix captures the reward correlation between the pairs of actions $a, a' \in [A]$. Then the posterior average mean estimator $\hat{\mu} \in \mathbb{R}^A$ over all tasks is given by the following lemma. The proof is given in Appendix B.1.

Lemma 1. *Let \mathbf{H}_n be the action matrix, \mathbf{Y}_n be the reward vector and \mathbf{S}_A be the estimated reward covariance matrix. Then the posterior prediction of the average mean reward vector $\hat{\mu}$ over all tasks is given by*

$$\hat{\mu} = \sigma_\theta^2 \mathbf{S}_A \mathbf{H}_n^\top (\sigma_\theta^2 \mathbf{H}_n (\mathbf{S}_A + \mathbf{D}_A) \mathbf{H}_n^\top)^{-1} \mathbf{Y}_n. \quad (7)$$

The $\hat{\mu}$ in equation 7 represents the posterior mean vector averaged on all tasks. So if some action $a \in [A]$ consistently yields high rewards in the pretraining data then $\hat{\mu}(a)$ has high value. Since the test distribution is the same as pretraining, this action on average will yield a high reward during test time.

We hypothesize that the **PreDeToR** is learning the reward correlation covariance matrix from the training data $\mathcal{H}_{\text{train}}$ and acting greedily on it. To test this hypothesis, we consider the greedy **BayesPred** algorithm that first estimates \mathbf{S}_A from the pretraining data. It then uses the LMMSE estimator in Lemma 1 to calculate the posterior mean vector $\hat{\mu}$, and then selects $I_t = \arg \max_a \hat{\mu}(a)$ at each round t . Note that **BayesPred** is a greedy algorithm that always selects the most rewarding action (exploitation) without any exploration of sub-optimal actions. Also the **BayesPred** is an LMMSE estimator that leverages the linear reward structure and estimates the reward covariance matrix, and therefore can be interpreted as a lower bound to the regret of **PreDeToR**. The hypothesis that **BayesPred** is a lower bound to **PreDeToR** is supported by Figure 7. In Figure 7 the reward covariance matrix for **BayesPred** is estimated from the $\mathcal{H}_{\text{train}}$ by first running the **Thomp** (π^w). Observe that the **BayesPred** has a lower cumulative regret than **PreDeToR** and almost matches the regret of **PreDeToR** towards the end of the horizon. Also note that **LinUCB** has lower cumulative regret towards the end of horizon as it leverages the linear structure and the feature of the actions in selecting the next action.

A.6 EMPIRICAL STUDY: INCREASING NUMBER OF ACTIONS

In this section, we discuss the performance of **PreDeToR** when the number of actions is very high so that the weak demonstrator π^w does not have sufficient samples for each action. However, the number of tasks $M_{\text{pre}} \gg A > n$.

Baselines: We again implement the same baselines discussed in Section 4. The baselines are **PreDeToR**, **PreDeToR- τ** , **DPT-greedy**, **AD**, **Thomp**, and **LinUCB**.

Outcomes: We first discuss the main outcomes from our experimental results of introducing more actions than the horizon (or more dimensions than actions) during data collection and evaluation:

Finding 8: **PreDeToR** ($-\tau$) outperforms **DPT-greedy**, and **AD**, even when $A > n$ but $M_{\text{pre}} \gg A$.

Experimental Result: We observe these outcomes in Figure 8. In Figure 8(a) we show the linear bandit setting for $M_{\text{pre}} = 250000$, $M_{\text{test}} = 200$, $A = 100$, $n = 50$ and $d = 5$. Again, the demonstrator π^w is the **Thomp** algorithm. We observe that **PreDeToR** ($-\tau$) has lower cumulative regret than **DPT-greedy** and **AD**. Note that for any task m the **Thomp** will not be able to sample all the actions even once. The weak performance of **DPT-greedy** can be attributed to both short horizons and the inability to estimate the optimal action for such a short horizon $n < A$. The **AD** performs similar to the demonstrator **Thomp** because of its training. Observe that **PreDeToR** ($-\tau$) has similar regret to **LinUCB** and lower regret than **Thomp** which also shows that **PreDeToR** is exploiting the

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

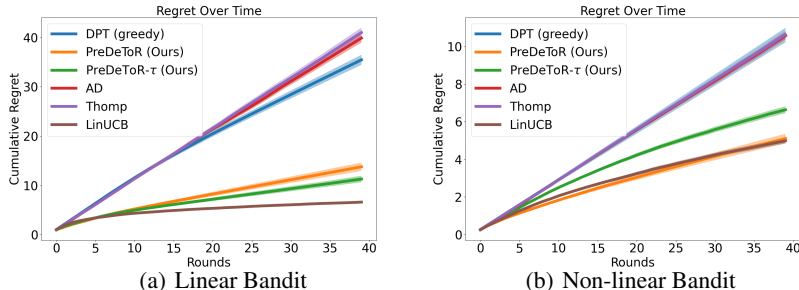


Figure 8: Testing the limit experiments. The horizontal axis is the number of rounds. Confidence bars show one standard error.

latent linear structure of the underlying tasks. In Figure 8(b) we show the non-linear bandit setting for horizon $n = 40$, $M_{\text{pre}} = 200000$, $A = 60$, $d = 2$, and $|\mathcal{A}^{\text{inv}}| = 5$. The demonstrator π^w is the Thomp algorithm. Again we observe that PreDeToR ($-\tau$) has lower cumulative regret than DPT-greedy, AD and LinUCB which fails to perform well in this non-linear setting due to its algorithmic design.

A.7 EMPIRICAL STUDY: INCREASING HORIZON

In this section, we discuss the performance of PreDeToR with respect to an increasing horizon for each task $m \in [M]$. However, note that the number of tasks $M_{\text{pre}} \geq n$. Note that Lee et al. (2023) studied linear bandit setting for $n = 200$. We study the setting up to a similar horizon scale.

Baselines: We again implement the same baselines discussed in Section 4. The baselines are PreDeToR, PreDeToR- τ , DPT-greedy, AD, Thomp, and LinUCB.

Outcomes: We first discuss the main outcomes of our experimental results for increasing the horizon:

Finding 9: PreDeToR ($-\tau$) outperforms DPT-greedy, and AD with increasing horizon.

Experimental Result: We observe these outcomes in Figure 9. In Figure 9 we show the linear bandit setting for $M_{\text{pre}} = 150000$, $M_{\text{test}} = 200$, $A = 20$, $n = \{20, 40, 60, 100, 120, 140, 200\}$ and $d = 5$. Again, the demonstrator π^w is the Thomp algorithm. We observe that PreDeToR ($-\tau$) has lower cumulative regret than DPT-greedy, and AD. Note that for any task m for the horizon 20 the Thomp will be able to sample all the actions at most once. Observe from Figure 9(a), 9(b), 9(c), Figure 9(d), 9(e), 9(f) and 9(g) that PreDeToR ($-\tau$) is closer to LinUCB and outperforms Thomp which also shows that PreDeToR ($-\tau$) is learning the latent linear structure of the underlying tasks. In Figure 9(h) we plot the regret of all the baselines with respect to the increasing horizon. Again we see that PreDeToR ($-\tau$) is closer to LinUCB and outperforms DPT-greedy, AD and Thomp. This shows that PreDeToR ($-\tau$) is able to exploit the latent structure and reward correlation across the tasks for varying horizon length.

A.8 EMPIRICAL STUDY: INCREASING DIMENSION

In this section, we discuss the performance of PreDeToR with respect to an increasing dimension for each task $m \in [M]$. Again note that the number of tasks $M_{\text{pre}} \gg A \geq n$. Through this experiment, we want to evaluate the performance of PreDeToR and see how it exploits the underlying reward correlation when the horizon is small as well as for increasing dimensions.

Baselines: We again implement the same baselines discussed in Section 4. The baselines are PreDeToR, PreDeToR- τ , DPT-greedy, AD, Thomp, and LinUCB.

Outcomes: We first discuss the main outcomes of our experimental results for increasing the horizon:

Finding 10: PreDeToR ($-\tau$) outperforms DPT-greedy, AD with increasing dimension and has lower regret than LinUCB for larger dimension.

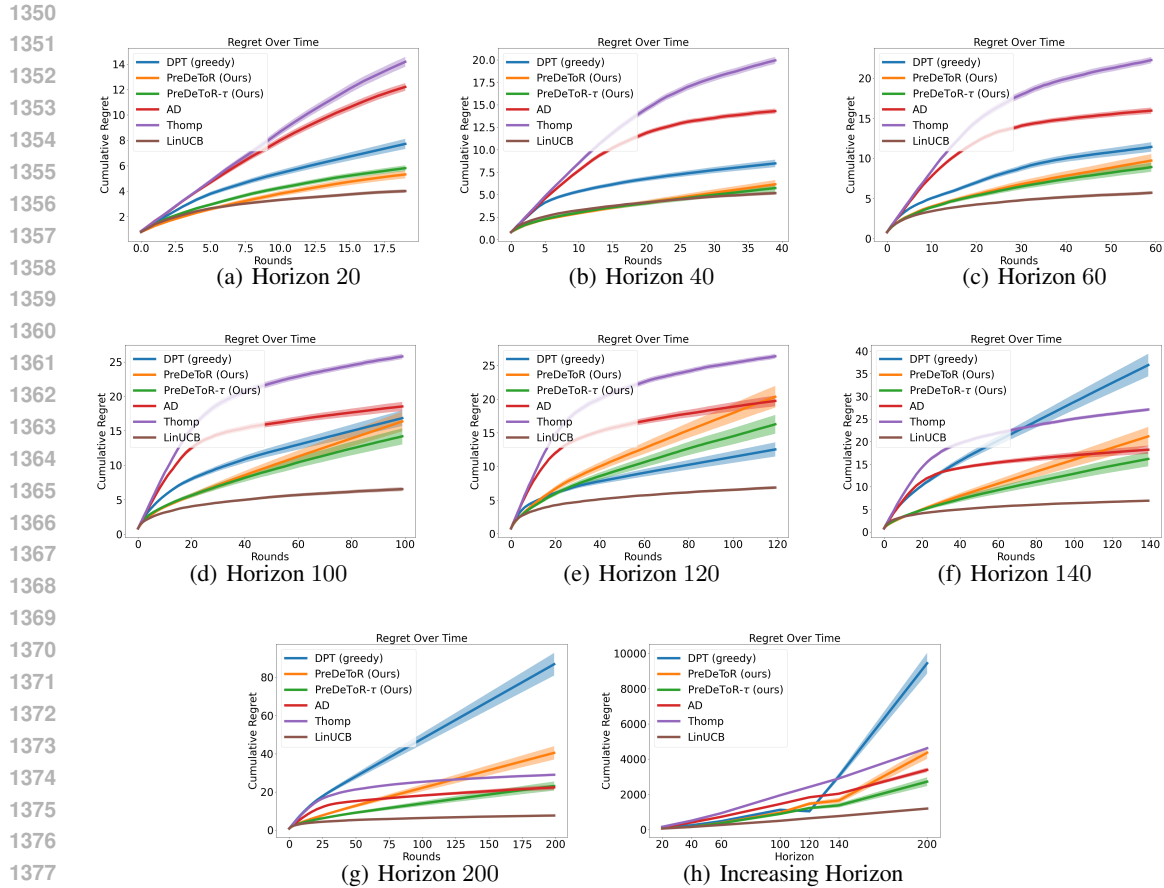


Figure 9: Experiment with increasing horizon. The y-axis shows the cumulative regret.

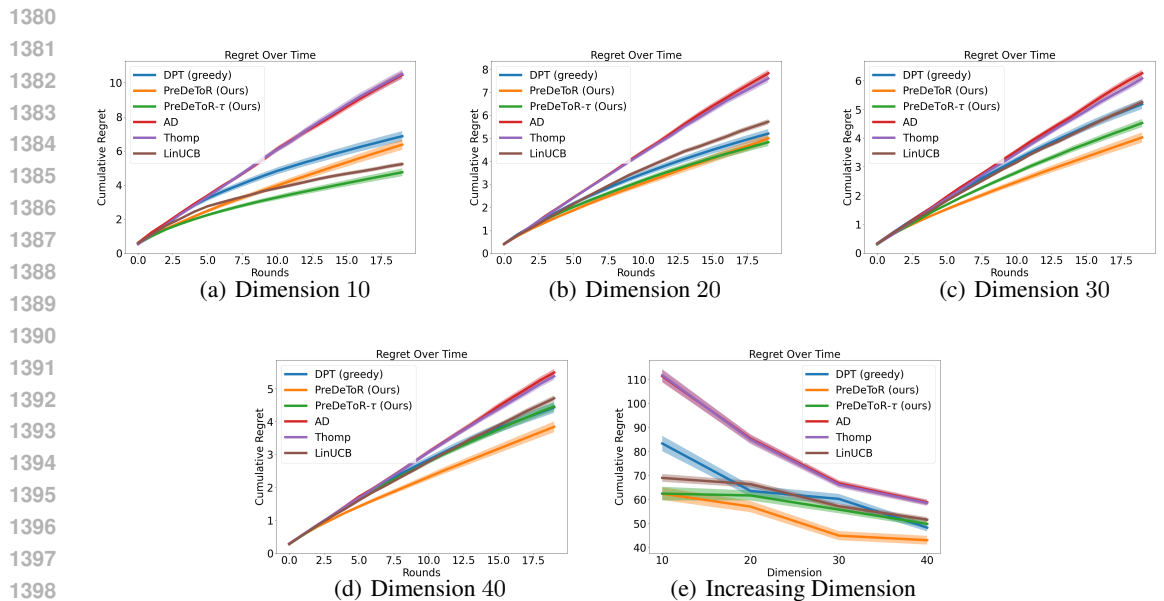


Figure 10: Experiment with increasing dimension. The y-axis shows the cumulative regret.

Experimental Result: We observe these outcomes in Figure 9. In Figure 9 we show the linear bandit setting for horizon $n = 20$, $M_{\text{pre}} = 160000$, $M_{\text{test}} = 200$, $A = 20$, and $d = \{10, 20, 30, 40\}$. Again, the demonstrator π^w is the Thomp algorithm. We observe that PreDeToR ($-\tau$) has lower cumulative regret than DPT-greedy, AD. Note that for any task m for the horizon 20 the Thomp

will be able to sample all the actions at most once. Observe from Figure 10(a), 10(b), 10(c), and 10(d) that **PreDeToR** ($-\tau$) is closer to **LinUCB** and has lower regret than **Thomp** which also shows that **PreDeToR** ($-\tau$) is exploiting the latent linear structure of the underlying tasks. In Figure 10(e) we plot the regret of all the baselines with respect to the increasing dimension. Again we see that **PreDeToR** ($-\tau$) has lower regret than **DPT-greedy**, **AD** and **Thomp**. Observe that with increasing dimension **PreDeToR** is able to outperform **LinUCB**. This shows that the **PreDeToR** ($-\tau$) is able to exploit reward correlation across tasks for varying dimensions.

A.9 EMPIRICAL STUDY: INCREASING ATTENTION HEADS

In this section, we discuss the performance of **PreDeToR** with respect to an increasing attention heads for the transformer model for the non-linear feedback model. Again note that the number of tasks $M_{\text{pre}} \gg A \geq n$. Through this experiment, we want to evaluate the performance of **PreDeToR** to exploit the underlying reward correlation when the horizon is small and understand the representative power of the transformer by increasing the attention heads. Note that we choose the non-linear feedback model and low data regime to leverage the representative power of the transformer.

Baselines: We again implement the same baselines discussed in Section 4. The baselines are **PreDeToR**, **PreDeToR**- τ , **DPT-greedy**, **AD**, **Thomp**, and **LinUCB**.

Outcomes: We first discuss the main outcomes of our experimental results for increasing the horizon:

Finding 11: **PreDeToR** ($-\tau$) outperforms **DPT-greedy**, and **AD** with increasing attention heads.

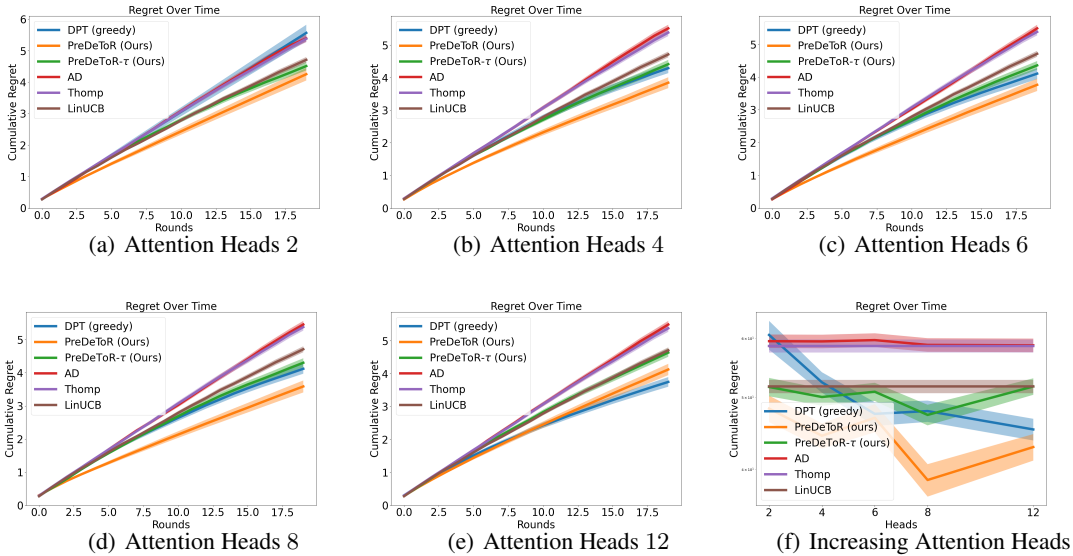


Figure 11: Experiment with increasing attention heads. The y-axis shows the cumulative regret.

Experimental Result: We observe these outcomes in Figure 11. In Figure 11 we show the non-linear bandit setting for horizon $n = 20$, $M_{\text{pre}} = 160000$, $M_{\text{test}} = 200$, $A = 20$, heads = $\{2, 4, 6, 8\}$ and $d = 5$. Again, the demonstrator π^w is the **Thomp** algorithm. We observe that **PreDeToR** ($-\tau$) has lower cumulative regret than **DPT-greedy**, **AD**. Note that for any task m for the horizon 20 the **Thomp** will be able to sample all the actions at most once. Observe from Figure 11(a), 11(b), 11(c), and 11(d) that **PreDeToR** ($-\tau$) has lower regret than **AD**, **Thomp** and **LinUCB** which also shows that **PreDeToR** ($-\tau$) is exploiting the latent linear structure of the underlying tasks for the non-linear setting. In Figure 11(f) we plot the regret of all the baselines with respect to the increasing attention heads. Again we see that **PreDeToR** ($-\tau$) regret decreases as we increase the attention heads.

A.10 EMPIRICAL STUDY: INCREASING NUMBER OF TASKS

In this section, we discuss the performance of **PreDeToR** with respect to the increasing number of tasks for the linear bandit setting. Again note that the number of tasks $M_{\text{pre}} \gg A \geq n$. Through this experiment, we want to evaluate the performance of **PreDeToR** to exploit the underlying reward correlation when the horizon is small and the number of tasks is changing. Finally, recall that when the horizon is small the weak demonstrator π^w does not have sufficient samples for each action. This leads to a poor approximation of the greedy action.

Baselines: We again implement the same baselines discussed in Section 4. The baselines are **PreDeToR**, **PreDeToR- τ** , **DPT-greedy**, **AD**, **Thomp**, and **LinUCB**.

Outcomes: We first discuss the main outcomes of our experimental results for increasing the horizon:

Finding 12: **PreDeToR (- τ)** fails to exploit the underlying latent structure and reward correlation from in-context data when the number of tasks is small.

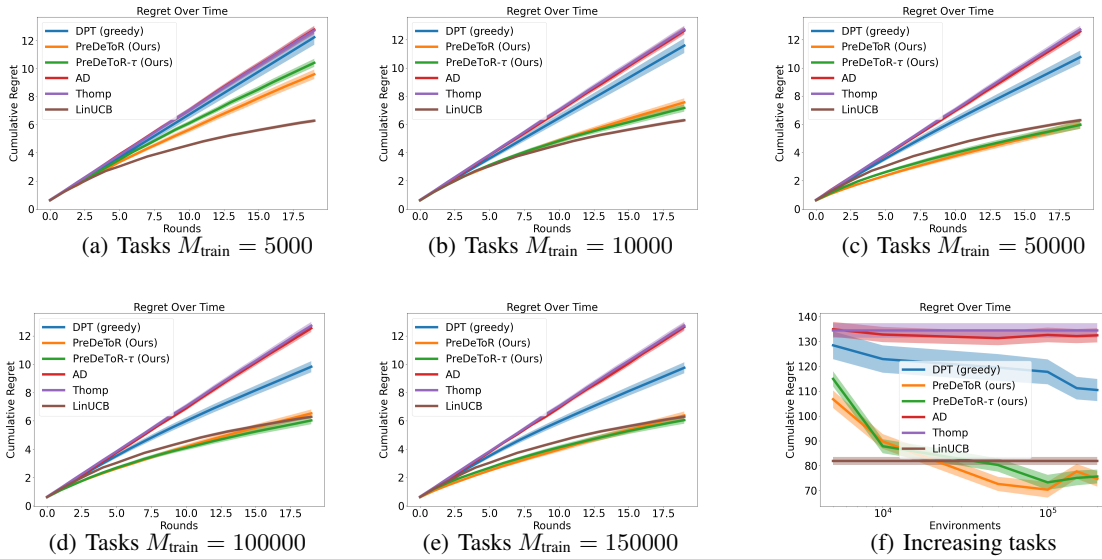


Figure 12: Experiment with an increasing number of tasks. The y-axis shows the cumulative regret.

Experimental Result: We observe these outcomes in Figure 12. In Figure 12 we show the linear bandit setting for horizon $n = 20$, $M_{\text{pre}} \in \{5000, 10000, 50000, 100000, 150000\}$, $M_{\text{test}} = 200$, $A = 20$, and $d = 40$. Again, the demonstrator π^w is the **Thomp** algorithm. We observe that **PreDeToR (- τ)**, **AD** and **DPT-greedy** suffer more regret than the **LinUCB** when the number of tasks is small ($M_{\text{train}} \in \{5000, 10000\}$) in Figure 12(a), and 12(b). However in Figure 12(c), 12(d), 12(e), and 12(f) we show that **PreDeToR** has lower regret than **Thomp** and matches **LinUCB**. This shows that **PreDeToR (- τ)** is exploiting the latent linear structure of the underlying tasks for the non-linear setting. Moreover, observe that as M_{train} increases the **PreDeToR** has lower cumulative regret than **DPT-greedy**, **AD**. Note that for any task m for the horizon 20 the **Thomp** will be able to sample all the actions at most once. Therefore **DPT-greedy** does not perform as well as **PreDeToR**. Finally, note that the result shows that **PreDeToR (- τ)** is able to exploit the reward correlation across the tasks better as the number of tasks increases.

A.11 EXPLORATION OF **PREDETOER(- τ)**

In this section, we discuss the exploration of **PreDeToR** in the linear bandit setting discussed in Section 4. Recall that the linear bandit setting consist of horizon $n = 25$, $M_{\text{pre}} = 200000$, $M_{\text{test}} = 200$, $A = 10$, and $d = 2$. The demonstrator π^w is the **Thomp** algorithm and we observe that **PreDeToR (- τ)** has lower cumulative regret than **DPT-greedy**, **AD** and matches the performance of

LinUCB. Therefore PreDeToR ($-\tau$) behaves almost optimally in this setting and so we analyze how PreDeToR conducts exploration for this setting.

Outcomes: We first discuss the main outcomes of our analysis of exploration in the low-data regime:

Finding 13: The PreDeToR ($-\tau$) has a two phase exploration. In the first phase, it explores with a strong prior over the in-context training data. In the second phase, once the task data has been observed for a few rounds (in-context) it switches to task-based exploration.

We first show in Figure 13(a) the training distribution of the optimal actions. For each bar, the frequency indicates the number of tasks where the action (shown in the x-axis) is the optimal action.

Then in Figure 13(b) we show how the sampling distribution of DPT-greedy, PreDeToR and PreDeToR- τ change in the first 10 and last 10 rounds for all the tasks where action 5 is optimal. To plot this graph we first sum over the individual pulls of the action taken by each algorithm over the first 10 and last 10 rounds. Then we average these counts over all test tasks where action 5 is optimal. From the figure Figure 13(b) we see that PreDeToR($-\tau$) consistently pulls the action 5 more than DPT-greedy. It also explores other optimal actions like {2, 3, 6, 7, 10} but discards them quickly in favor of the optimal action 5 in these tasks. This shows that PreDeToR ($-\tau$) only considers the optimal actions seen from the training data. Once sufficient observation have been observed for the task it switches to task-based exploration and samples the optimal action more than DPT-greedy.

Finally, we plot the feasible action set considered by DPT-greedy, PreDeToR, and PreDeToR- τ in Figure 13(c). To plot this graph again we consider the test tasks where the optimal action is 5. Then we count the number of distinct actions that are taken from round t up until horizon n . Finally we average this over all the considered tasks where the optimal action is 5. We call this the candidate action set considered by the algorithm. From the Figure 13(c) we see that DPT-greedy explores the least and gets stuck with few actions quickly (by round 10). Note that the actions DPT-greedy samples are sub-optimal and so it suffers a high cumulative regret (see Figure 2). PreDeToR explore slightly more than DPT-greedy, but PreDeToR- τ explores the most.

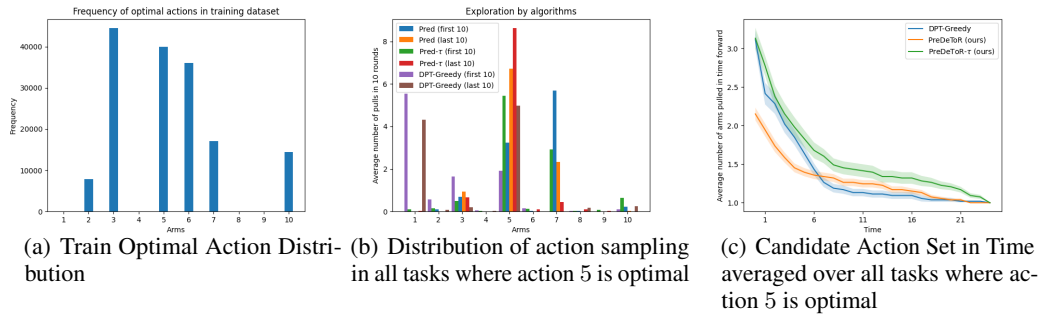


Figure 13: Exploration Analysis of PreDeToR($-\tau$)

A.12 EXPLORATION OF PREDETOR($-\tau$) IN NEW ARMS SETTING

In this section, we discuss the exploration of PreDeToR ($-\tau$) in the linear and non-linear new arms bandit setting discussed in Section 6. Recall that we consider the linear bandit setting of horizon $n = 50$, $M_{pre} = 200000$, $M_{test} = 200$, $A = 20$, and $d = 5$. Here during data collection and during collecting the test data, we randomly select one new action from \mathbb{R}^d for each task m . So the number of invariant actions is $|\mathcal{A}^{inv}| = 19$.

Outcomes: We first discuss the main outcomes of our analysis of exploration in the low-data regime:

Finding 14: The **PreDeToR** ($-\tau$) is robust to changes when the number of in-variant actions is large. **PreDeToR** ($-\tau$) performance drops as shared structure breaks down.

We first show in Figure 14(a) the training distribution of the optimal actions. For each bar, the frequency indicates the number of tasks where the action (shown in the x-axis) is the optimal action.

Then in Figure 14(b) we show how the sampling distribution of **DPT-greedy**, **PreDeToR** and **PreDeToR** ($-\tau$) change in the first 10 and last 10 rounds for all the tasks where action 17 is optimal. We plot this graph the same way as discussed in Appendix A.11. From the figure Figure 14(b) we see that **PreDeToR** ($-\tau$) consistently pulls the action 17 more than **DPT-greedy**. It also explores other optimal actions like $\{1, 2, 3, 8, 9, 15\}$ but discards them quickly in favor of the optimal action 17 in these tasks.

Finally, we plot the feasible action set considered by **DPT-greedy**, **PreDeToR**, and **PreDeToR** ($-\tau$) in Figure 14(c). To plot this graph again we consider the test tasks where the optimal action is 17. Then we count the number of distinct actions that are taken from round t up until horizon n . Finally we average this over all the considered tasks where the optimal action is 17. We call this the candidate action set considered by the algorithm. From the Figure 14(c) we see that **PreDeToR** ($-\tau$) explores more than **PreDeToR** in this setting.

We also show how the prediction error of the optimal action by **PreDeToR** compared to **LinUCB** in this 1 new arm linear bandit setting. In Figure 15(a) we first show how the 20 actions are distributed in the $M_{\text{test}} = 200$ test tasks. In Figure 15(a) for each bar, the frequency indicates the number of tasks where the action (shown in the x-axis) is the optimal action. Then in Figure 15(b) we show the prediction error of **PreDeToR** ($-\tau$) for each task $m \in [M_{\text{test}}]$. The prediction error is calculated the same way as stated in Section 6 From the Figure 15(b) we see that for most actions the prediction error of **PreDeToR** ($-\tau$) is closer to **LinUCB** showing that the introduction of 1 new action does not alter the prediction error much. Note that **LinUCB** estimates the empirical mean directly from the test task, whereas **PreDeToR** has a strong prior based on the training data. Therefore we see that **PreDeToR** is able to estimate the reward of the optimal action quite well from the training dataset \mathcal{D}_{pre} .

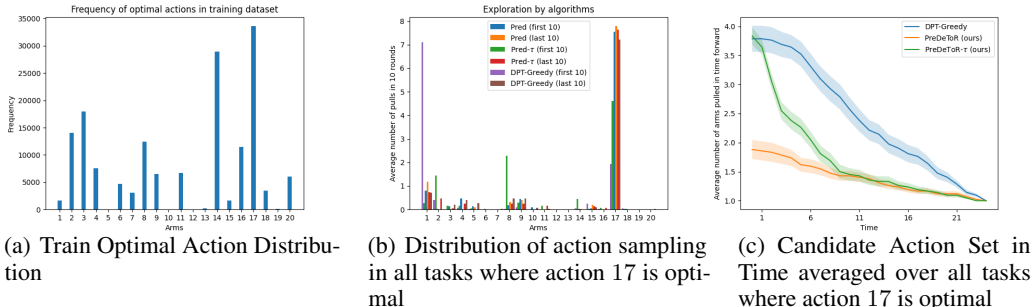


Figure 14: Exploration Analysis of **PreDeToR** ($-\tau$) in linear 1 new arm setting

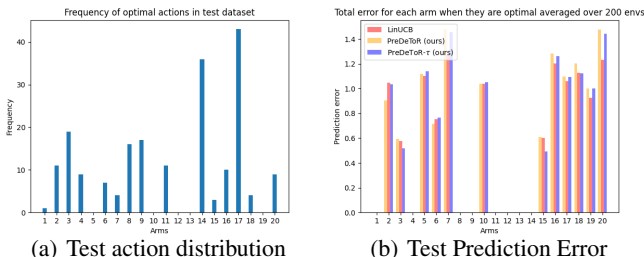
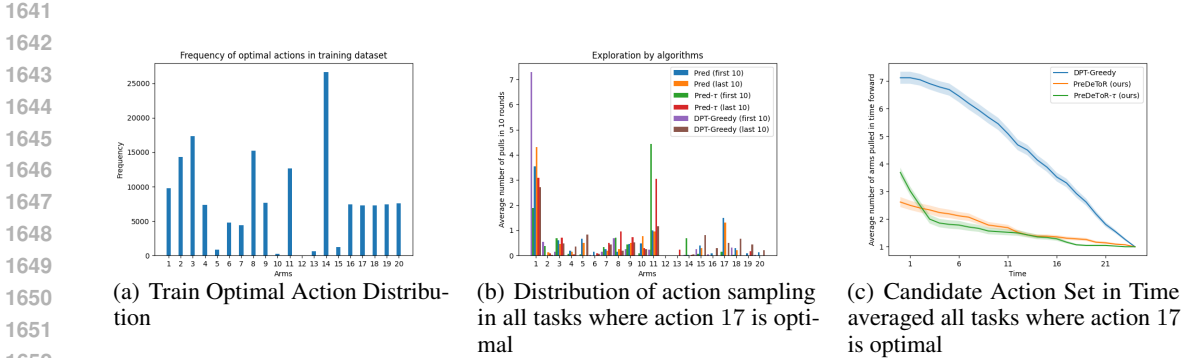


Figure 15: Prediction error of **PreDeToR** ($-\tau$) in linear 1 new arm setting

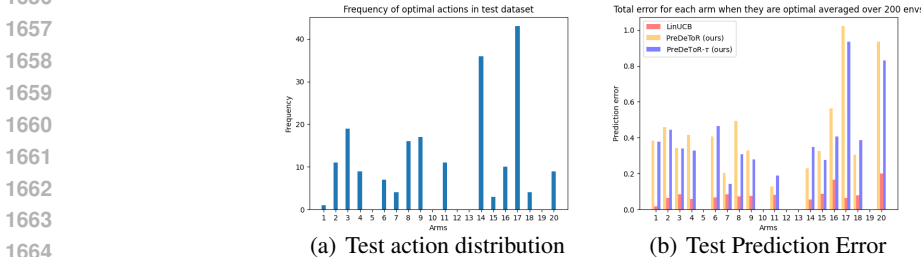
1620 We now consider the setting where the number of invariant actions is $|\mathcal{A}^{\text{inv}}| = 15$. We again show in
 1621 Figure 16(a) the training distribution of the optimal actions. For each bar, the frequency indicates the
 1622 number of tasks where the action (shown in the x-axis) is the optimal action. Then in Figure 16(b)
 1623 we show how the sampling distribution of **DPT-greedy**, **PreDeToR** and **PreDeToR- τ** change in the
 1624 first 10 and last 10 rounds for all the tasks where action 17 is optimal. We plot this graph the same
 1625 way as discussed in Appendix A.11. From the figure Figure 16(b) we see that none of the algorithms
 1626 **PreDeToR**, **PreDeToR- τ** , **DPT-greedy** consistently pulls the action 17 more than other actions. This
 1627 shows that the common underlying actions across the tasks matter for learning the exploration.

1628 Finally, we plot the feasible action set considered by **DPT-greedy**, **PreDeToR**, and **PreDeToR- τ** in
 1629 Figure 16(c). To plot this graph again we consider the test tasks where the optimal action is 17. We
 1630 build the candidate set the same way as before. From the Figure 16(c) we see that none of the three
 1631 algorithms **DPT-greedy**, **PreDeToR**, **PreDeToR- τ** , is able to sample the optimal action 17 sufficiently
 1632 high number of times.

1633 We also show how the prediction error of the optimal action by **PreDeToR** compared to **LinUCB** in
 1634 this 1 new arm linear bandit setting. In Figure 17(a) we first show how the 20 actions are distributed
 1635 in the $M_{\text{test}} = 200$ test tasks. In Figure 17(a) for each bar, the frequency indicates the number of
 1636 tasks where the action (shown in the x-axis) is the optimal action. Then in Figure 17(b) we show the
 1637 prediction error of **PreDeToR (- τ)** for each task $m \in [M_{\text{test}}]$. The prediction error is calculated the
 1638 same way as stated in Section 6. From the Figure 17(b) we see that for most actions the prediction
 1639 error is higher than **LinUCB** showing that the introduction of 5 new actions (and thereby decreasing
 1640 the invariant action set) significantly alters the prediction error.



1650 Figure 16: Exploration Analysis of **PreDeToR(- τ)** in linear 5 new arm setting



1654 Figure 17: Prediction error of **PreDeToR(- τ)** in linear 1 new arm setting

1655 A.13 DATA COLLECTION ANALYSIS

1656 In this section, we analyze the performance of **PreDeToR**, **PreDeToR- τ** , **DPT-greedy**, **AD**, **Thomp**,
 1657 and **LinUCB** when the weak demonstrator π^w is **Thomp**, **LinUCB**, or **Uniform**. We again consider the
 1658 linear bandit setting discussed in Section 4. Recall that the linear bandit setting consist of horizon
 1659 $n = 25$, $M_{\text{pre}} = 200000$, $M_{\text{test}} = 200$, $A = 10$, and $d = 2$. Finally, we show the cumulative regret by
 1660 the above baselines in Figure 18(a), 18(b), and 18(b) when data is collected through **Thomp**, **LinUCB**,
 1661 and **Uniform** respectively.

Outcomes: We first discuss the main outcomes of our experimental results for different data collection:

Finding 15: The **PreDeToR** ($-\tau$) excels in exploiting the underlying latent structure and reward correlation from in-context data when the data diversity is high.

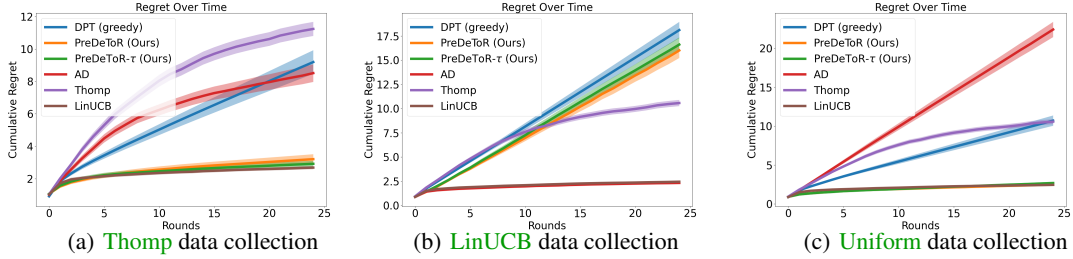


Figure 18: Data Collection with various algorithms and Performance analysis

Experimental Result: We observe these outcomes in Figure 18. In Figure 18(a) we see that the A -actioned **Thomp** is explorative enough as it does not explore with the knowledge of feature representation. So it pulls the sub-optimal actions sufficiently high number of times before discarding them in favor of the optimal action. Therefore the training data is diverse enough so that **PreDeToR** ($-\tau$) can predict the reward vectors for actions sufficiently well. Consequently, **PreDeToR** ($-\tau$) almost matches the **LinUCB** algorithm. Both **DPT-greedy** and **AD** perform poorly in this setting.

In Figure 18(b) we see that the **LinUCB** algorithm is not explorative enough as it explores with the knowledge of feature representation and quickly discards the sub-optimal actions in favor of the optimal action. Therefore the training data is not diverse enough so that **PreDeToR** ($-\tau$) is not able to correctly predict the reward vectors for actions. Note that **DPT-greedy** also performs poorly in this setting when it is not provided with the optimal action information during training. The **AD** matches the performance of its demonstrator **LinUCB** because of its training procedure of predicting the next action of the demonstrator.

Finally, in Figure 18(c) we see that the A -armed **Uniform** is fully explorative as it does not intend to minimize regret (as opposed to **Thomp**) and does not explore with the knowledge of feature representation. Therefore the training data is very diverse which results in **PreDeToR** ($-\tau$) being able to predict the reward vectors for actions very well. Consequently, **PreDeToR** ($-\tau$) perfectly matches the **LinUCB** algorithm. Note that **AD** performs the worst as it matches the performance of its demonstrator whereas the performance of **DPT-greedy** suffers due to the lack of information on the optimal action during training.

A.14 EMPIRICAL VALIDATION OF THEORETICAL RESULT

In this section, we empirically validate the theoretical result proved in Section 7. We again consider the linear bandit setting discussed in Section 4. Recall that the linear bandit setting consist of horizon $n = 25$, $M_{\text{pre}} = \{100000, 200000\}$, $M_{\text{test}} = 200$, $A = 10$, and $d = 2$. The demonstrator π^w is the **Thomp** algorithm and we observe that **PreDeToR** ($-\tau$) has lower cumulative regret than **DPT-greedy**, **AD** and matches the performance of **LinUCB**.

Baseline (LinUCB- τ): We define soft LinUCB (**LinUCB- τ**) as follows: At every round t for task m , it calculates the ucb value $B_{m,a,t}$ for each action $\mathbf{x}_{m,a} \in \mathcal{X}$ such that $B_{m,a,t} = \mathbf{x}_{m,a}^\top \hat{\boldsymbol{\theta}}_{m,t-1} + \alpha \|\mathbf{x}_{m,a}\|_{\boldsymbol{\Sigma}_{m,t-1}^{-1}}$ where $\alpha > 0$ is a constant and $\hat{\boldsymbol{\theta}}_{m,t}$ is the estimate of the model parameter $\boldsymbol{\theta}_{m,*}$ at round t . Here, $\boldsymbol{\Sigma}_{m,t-1} = \sum_{s=1}^{t-1} \mathbf{x}_{m,s} \mathbf{x}_{m,s}^\top + \lambda \mathbf{I}_d$ is the data covariance matrix or the arms already tried. Then it chooses $I_t \sim \text{softmax}_a^\tau(B_{m,a,t})$, where $\text{softmax}_a^\tau(\cdot) \in \Delta^A$ denotes a softmax

distribution over the actions and τ is a temperature parameter (See Section 4 for definition of $\text{softmax}_a^\tau(\cdot)$).

Outcomes: We first discuss the main outcomes of our experimental results:

Finding 16: **PreDeToR** ($-\tau$) excels in predicting the rewards for test tasks when the number of training (source) tasks is large.

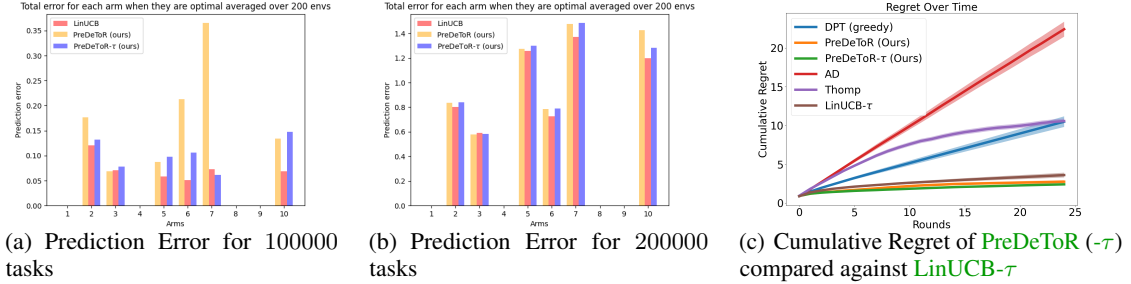


Figure 19: Empirical validation of theoretical analysis

Experimental Result: These findings are reported in Figure 19. In Figure 19(a) we show the prediction error of **PreDeToR** ($-\tau$) for each task $m \in [M_{\text{test}}]$. The prediction error is calculated as $(\hat{\mu}_{m,n,*}(a) - \mu_{m,*}(a))^2$ where $\hat{\mu}_{m,n,*}(a) = \max_a \hat{\theta}_{m,n}^\top \mathbf{x}_m(a)$ is the empirical mean at the end of round n , and $\mu_{*,m}(a) = \max_a \theta_{m,*}^\top \mathbf{x}_m(a)$ is the true mean of the optimal action in task m . Then we average the prediction error for the action $a \in [A]$ by the number of times the action a is the optimal action in some task m . We see that when the source tasks are 100000 the reward prediction falls short of **LinUCB** prediction for all actions except action 2.

In Figure 19(b) we again show the prediction error of **PreDeToR** ($-\tau$) for each task $m \in [M_{\text{test}}]$ when the source tasks are 200000. Note that in both these settings, we kept the horizon $n = 25$, and the same set of actions. We now observe that the reward prediction almost matches **LinUCB** prediction in almost all the optimal actions.

In Figure 19(c) we compare **PreDeToR** ($-\tau$) against **LinUCB**- τ and show that they almost match in the linear bandit setting discussed in Section 4 when the source tasks are 100000.

A.15 EMPIRICAL STUDY: OFFLINE PERFORMANCE

In this section, we discuss the offline performance of **PreDeToR** when the number of tasks $M_{\text{pre}} \gg A \geq n$.

We first discuss how **PreDeToR** ($-\tau$) is modified for the offline setting. In the offline setting, the **PreDeToR** first samples a task $m \sim \mathcal{T}_{\text{test}}$, then the test dataset $\mathcal{H}_m \sim \mathcal{D}_{\text{test}}(\cdot|m)$. Then **PreDeToR** and **PreDeToR**- τ act similarly to the online setting, but based on the entire offline dataset \mathcal{H}_m . The full pseudocode of **PreDeToR** is in Algorithm 2.

Recall that $\mathcal{D}_{\text{test}}$ denote a distribution over all possible interactions that can be generated by π^w during test time. For offline testing, first, a test task $m \sim \mathcal{T}_{\text{test}}$, and then an in-context test dataset \mathcal{H}_m is collected such that $\mathcal{H}_m \sim \mathcal{D}_{\text{test}}(\cdot|m)$. Observe from Algorithm 2 that in the offline setting, **PreDeToR** first samples a task $m \sim \mathcal{T}_{\text{test}}$, and then a test dataset $\mathcal{H}_m \sim \mathcal{D}_{\text{test}}(\cdot|m)$ and acts greedily. Crucially in the offline setting the **PreDeToR** does not add the observed reward r_t at round t to the dataset. Through this experiment, we want to evaluate the performance of **PreDeToR** to learn the underlying latent structure and reward correlation when the horizon is small. Finally, recall that when the horizon is small the weak demonstrator π^w does not have sufficient samples for each action. This leads to a poor approximation of the greedy action.

Algorithm 2 Pre-trained Decision Transformer with Reward Estimation (PreDeToR)

- 1: **Collecting Pretraining Dataset**
- 2: Initialize empty pretraining dataset $\mathcal{H}_{\text{train}}$
- 3: **for** i in $[M_{\text{pre}}]$ **do**
- 4: Sample task $m \sim \mathcal{T}_{\text{pre}}$, in-context dataset $\mathcal{H}_m \sim \mathcal{D}_{\text{pre}}(\cdot|m)$ and add this to $\mathcal{H}_{\text{train}}$.
- 5: **end for**
- 6: **Pretraining model on dataset**
- 7: Initialize model TF_{Θ} with parameters Θ
- 8: **while** not converged **do**
- 9: Sample \mathcal{H}_m from $\mathcal{H}_{\text{train}}$ and predict $\hat{r}_{m,t}$ for action $(I_{m,t})$ for all $t \in [n]$
- 10: Compute loss in equation 3 with respect to $r_{m,t}$ and backpropagate to update model parameter Θ .
- 11: **end while**
- 12: **Offline test-time deployment**
- 13: Sample unknown task $m \sim \mathcal{T}_{\text{test}}$, sample dataset $\mathcal{H}_m \sim \mathcal{D}_{\text{test}}(\cdot|m)$
- 14: Use TF_{Θ} on m at round t to choose

$$I_t \begin{cases} = \arg \max_{a \in \mathcal{A}} \text{TF}_{\Theta}(\hat{r}_{m,t}(a) | \mathcal{H}_m), & \text{PreDeToR} \\ \sim \text{softmax}_a^{\tau} \text{TF}_{\Theta}(\hat{r}_{m,t}(a) | \mathcal{H}_m), & \text{PreDeToR-}\tau \end{cases}$$

Baselines: We again implement the same baselines discussed in Section 4. The baselines are PreDeToR, PreDeToR- τ , DPT-greedy, AD, Thomp, and LinUCB. During test time evaluation for offline setting the DPT selects $I_t = \hat{a}_{m,t,*}$ where $\hat{a}_{m,t,*} = \arg \max_a \text{TF}_{\Theta}(a | \mathcal{H}_m^t)$ is the predicted optimal action.

Outcomes: We first discuss the main outcomes of our experimental results for increasing the horizon:

Finding 17: PreDeToR ($-\tau$) performs comparably to DPT-greedy and AD in the offline setting.

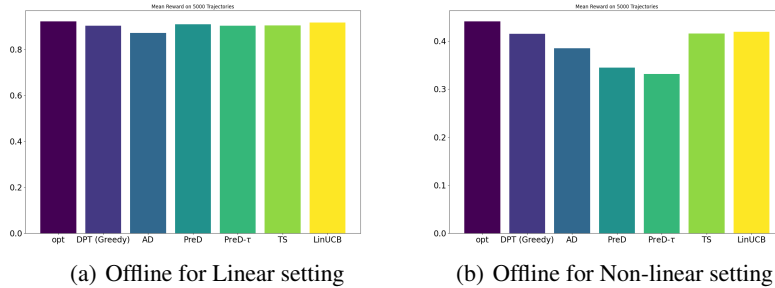


Figure 20: Offline experiment. The y-axis shows the cumulative reward.

Experimental Result: We observe these outcomes in Figure 20. In Figure 20 we show the linear bandit setting for horizon $n = 20$, $M_{\text{pre}} = 200000$, $M_{\text{test}} = 5000$, $A = 20$, and $d = 5$ for the low data regime. Again, the demonstrator π^w is the Thomp algorithm. We observe that PreDeToR ($-\tau$) has comparable cumulative regret to DPT-greedy. Note that for any task m for the horizon $n = 20$ the Thomp will be able to sample all the actions at most once. In the non-linear setting of Figure 20(b) the $n = 40$, $M_{\text{pre}} = 100000$, $A = 6$, $d = 2$. Observe that in all of these results, the performance of PreDeToR ($-\tau$) is comparable with respect to cumulative regret against DPT-greedy.

1836 B THEORETICAL ANALYSIS

1837 B.1 PROOF OF LEMMA 1

1840 *Proof.* The learner collects n rounds of data following π^w . The weak demonstrator π^w only observes
 1841 the $\{I_t, r_t\}_{t=1}^n$. Recall that $N_n(a)$ denotes the total number of times the action a is sampled for n
 1842 rounds. Define the matrix $\mathbf{H}_n \in \mathbb{R}^{n \times A}$ where the t -th row represents the action sampled at round
 1843 $t \in [n]$. The t -th row is a one-hot vector with 1 as the a -th component in the vector for $a \in [A]$. Then
 1844 define the reward vector $\mathbf{Y}_n \in \mathbb{R}^n$ as the reward vector where the t -th component is the observed
 1845 reward for the action I_t for $t \in [n]$. Finally define the diagonal matrix $\mathbf{D}_A \in \mathbb{R}^{A \times A}$ as in equation 6
 1846 and the estimated reward covariance matrix as $\mathbf{S}_A \in \mathbb{R}^{A \times A}$ such that $\mathbf{S}_A(a, a') = \hat{\mu}_n(a)\hat{\mu}_n(a')$.
 1847 This matrix captures the reward correlation between the pairs of actions $a, a' \in [A]$.

1848 Assume $\mu \sim \mathcal{N}(0, \mathbf{S}_*)$ where $\mathbf{S}_* \in \mathbb{R}^{A \times A}$. Then the observed mean vector \mathbf{Y}_n is

$$1849 \mathbf{Y}_n = \mathbf{H}_n \mu + \mathbf{H}_n \mathbf{D}_A^{1/2} \eta_n$$

1851 where, η_n is the noise vector over the $[n]$ training data. Then the posterior mean of $\hat{\mu}$ by Gauss
 1852 Markov Theorem (Johnson et al., 2002) is given by

$$1853 \hat{\mu} = \mathbf{S}_* \mathbf{H}_n^\top (\mathbf{H}_n (\mathbf{S}_* + \mathbf{D}_A) \mathbf{H}_n^\top)^{-1} \mathbf{Y}_n. \quad (8)$$

1855 However, the learner does not know the true reward co-variance matrix. Hence it needs to estimate
 1856 the \mathbf{S}_* from the observed data. Let the estimate be denoted by \mathbf{S}_A .

1857 **Assumption B.1.** We assume that π^w is sufficiently exploratory so that each action is sampled at
 1858 least once.

1860 The Assumption B.1 ensures that the matrix $(\sigma_\theta^2 \mathbf{H}_n (\mathbf{S}_A + \mathbf{D}_A) \mathbf{H}_n^\top)^{-1}$ is invertible. Under Assump-
 1861 tion B.1, plugging the estimate \mathbf{S}_A back in equation 8 shows that the average posterior mean over all
 1862 the tasks is

$$1863 \hat{\mu} = \mathbf{S}_A \mathbf{H}_n^\top (\mathbf{H}_n (\mathbf{S}_A + \mathbf{D}_A) \mathbf{H}_n^\top)^{-1} \mathbf{Y}_n. \quad (9)$$

1865 The claim of the lemma follows. \square

1866 C GENERALIZATION AND TRANSFER LEARNING PROOF FOR PREDETOR

1869 C.1 GENERALIZATION PROOF

1871 Alg is the space of algorithms induced by the transformer TF_Θ .

1872 **Theorem C.1. (PreDeToR risk)** Suppose error stability Assumption 7.1 holds and assume loss
 1873 function $\ell(\cdot, \cdot)$ is C -Lipschitz for all $r_t \in [0, B]$ and horizon $n \geq 1$. Let $\widehat{\text{TF}}$ be the empirical solution
 1874 of (ERM) and $\mathcal{N}(\mathcal{A}, \rho, \epsilon)$ be the covering number of the algorithm space Alg following Definition
 1875 C.2 and C.3. Then with probability at least $1 - 2\delta$, the excess Multi-task learning (MTL) risk of
 1876 PreDeToR- τ is bounded by

$$1877 \mathcal{R}_{\text{MTL}}(\widehat{\text{TF}}) \leq 4 \frac{C}{\sqrt{nM}} + 2(B + K \log n) \sqrt{\frac{\log(\mathcal{N}(\text{Alg}, \rho, \epsilon)/\delta)}{cnM}}$$

1880 where, $\mathcal{N}(\text{Alg}, \rho, \epsilon)$ is the covering number of transformer $\widehat{\text{TF}}$.

1882 *Proof.* We consider a meta-learning setting. Let M source tasks are i.i.d. sampled from a task
 1883 distribution \mathcal{T} , and let $\widehat{\text{TF}}$ be the empirical Multitask (MTL) solution. Define $\mathcal{H}_{\text{all}} = \bigcup_{m=1}^M \mathcal{H}_m$.
 1884 We drop the Θ, \mathbf{r} from transformer notation $\text{TF}^\mathbf{r}_\Theta$ as we keep the architecture fixed as in Lin et al.
 1885 (2023). Note that this transformer predicts a reward vector over the actions. To be more precise we
 1886 denote the reward predicted by the transformer at round t after observing history \mathcal{H}_m^{t-1} and then
 1887 sampling the action a_{mt} as $\text{TF}(\hat{r}_{mt}(a_{mt}) | \mathcal{H}_m^{t-1}, a_{mt})$. Define the training risk

$$1888 \hat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\text{TF}) = \frac{1}{nM} \sum_{m=1}^M \sum_{t=1}^n \ell(r_{mt}(a_{mt}), \text{TF}(\hat{r}_{mt}(a_{mt}) | \mathcal{H}_m^{t-1}, a_{mt}))$$

and the test risk

$$\mathcal{L}_{\text{MTL}}(\text{TF}) = \mathbb{E} \left[\widehat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\text{TF}) \right].$$

Define empirical risk minima $\widehat{\text{TF}} = \arg \min_{\text{TF} \in \text{Alg}} \widehat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\text{TF})$ and population minima

$$\text{TF}^* = \arg \min_{\text{TF} \in \text{Alg}} \mathcal{L}_{\text{MTL}}(\text{TF})$$

In the following discussion, we drop the subscripts MTL and \mathcal{H}_{all} . The excess MTL risk is decomposed as follows:

$$\begin{aligned} \mathcal{R}_{\text{MTL}}(\widehat{\text{TF}}) &= \mathcal{L}(\widehat{\text{TF}}) - \mathcal{L}(\text{TF}^*) \\ &= \underbrace{\mathcal{L}(\widehat{\text{TF}}) - \widehat{\mathcal{L}}(\widehat{\text{TF}})}_a + \underbrace{\widehat{\mathcal{L}}(\widehat{\text{TF}}) - \widehat{\mathcal{L}}(\text{TF}^*)}_b + \underbrace{\widehat{\mathcal{L}}(\text{TF}^*) - \mathcal{L}(\text{TF}^*)}_c. \end{aligned}$$

Since $\widehat{\text{TF}}$ is the minimizer of empirical risk, we have $b \leq 0$.

Step 1: (Concentration bound $|\mathcal{L}(\text{TF}) - \widehat{\mathcal{L}}(\text{TF})|$ for a fixed $\text{TF} \in \text{Alg}$) Define the test/train risks of each task as follows:

$$\widehat{\mathcal{L}}_m(\text{TF}) := \frac{1}{n} \sum_{t=1}^n \ell(r_{mt}(a_{mt}), \text{TF}(\widehat{r}_{mt}(a_{mt}) | \mathcal{H}_m^{t-1}, a_{mt})), \quad \text{and}$$

$$\mathcal{L}_m(\text{TF}) := \mathbb{E}_{\mathcal{H}_m} \left[\widehat{\mathcal{L}}_m(\text{TF}) \right] = \mathbb{E}_{\mathcal{H}_m} \left[\frac{1}{n} \sum_{t=1}^n \ell(r_{mt}(a_{mt}), \text{TF}(\widehat{r}_{mt}(a_{mt}) | \mathcal{H}_m^{t-1}, a_{mt})) \right], \quad \forall m \in [M].$$

Define the random variables $X_{m,t} = \mathbb{E} \left[\widehat{\mathcal{L}}_t(\text{TF}) | \mathcal{H}_m^t \right]$ for $t \in [n]$ and $m \in [M]$, that is, $X_{m,t}$ is the expectation over $\widehat{\mathcal{L}}_t(\text{TF})$ given training sequence $\mathcal{H}_m^t = \{(a_{mt'}, r_{mt'})\}_{t'=1}^t$ (which are the filtrations). With this, we have that $X_{m,n} = \mathbb{E} \left[\widehat{\mathcal{L}}_m(\text{TF}) | \mathcal{H}_m^n \right] = \widehat{\mathcal{L}}_m(\text{TF})$ and $X_{m,0} = \mathbb{E} \left[\widehat{\mathcal{L}}_m(\text{TF}) \right] = \mathcal{L}_m(\text{TF})$. More generally, $(X_{m,0}, X_{m,1}, \dots, X_{m,n})$ is a martingale sequence since, for every $m \in [M]$, we have that $\mathbb{E} [X_{m,t} | \mathcal{H}_m^{t-1}] = X_{m,t-1}$. For notational simplicity, in the following discussion, we omit the subscript m from a, r and \mathcal{H} as they will be clear from the left-hand-side variable $X_{m,t}$. We have that

$$\begin{aligned} X_{m,t} &= \mathbb{E} \left[\frac{1}{n} \sum_{t'=1}^n \ell(r_{t'}, \text{TF}(\widehat{r}_{t'} | \mathcal{H}^{t'-1}, a_{t'})) \middle| \mathcal{H}^t \right] \\ &= \frac{1}{n} \sum_{t'=1}^t \ell(r_{t'}, \text{TF}(\widehat{r}_{t'} | \mathcal{H}^{t'-1}, a_{t'})) + \frac{1}{n} \sum_{t'=t+1}^n \mathbb{E} \left[\ell(r_{t'}, \text{TF}(\widehat{r}_{t'} | \mathcal{H}^{t'-1}, a_{t'})) \middle| \mathcal{H}^t \right] \end{aligned}$$

Using the similar steps as in [Li et al. \(2023\)](#) we can show that

$$|X_{m,t} - X_{m,t-1}| \stackrel{(a)}{\leq} \frac{B}{n} + \sum_{t'=t+1}^n \frac{K}{t'n} \leq \frac{B + K \log n}{n}.$$

where, (a) follows by using the fact that loss function $\ell(\cdot, \cdot)$ is bounded by B , and error stability assumption.

Recall that $|\mathcal{L}_m(\text{TF}) - \widehat{\mathcal{L}}_m(\text{TF})| = |X_{m,0} - X_{m,n}|$ and for every $m \in [M]$, we have $\sum_{t=1}^n |X_{m,t} - X_{m,t-1}|^2 \leq \frac{(B+K \log n)^2}{n}$. As a result, applying Azuma-Hoeffding's inequality, we obtain

$$\mathbb{P} \left(\left| \mathcal{L}_m(\text{TF}) - \widehat{\mathcal{L}}_m(\text{TF}) \right| \geq \tau \right) \leq 2e^{-\frac{n\tau^2}{2(B+K \log n)^2}}, \quad \forall m \in [M] \quad (10)$$

Let us consider $Y_m := \mathcal{L}_m(\text{TF}) - \widehat{\mathcal{L}}_m(\text{TF})$ for $m \in [M]$. Then, $(Y_m)_{m=1}^M$ are i.i.d. zero mean sub-Gaussian random variables. There exists an absolute constant $c_1 > 0$ such that, the subgaussian

norm, denoted by $\|\cdot\|_{\psi_2}$, obeys $\|Y_m\|_{\psi_2}^2 < \frac{c_1(B+K \log n)^2}{n}$ via Proposition 2.5.2 of (Vershynin, 2018). Applying Hoeffding's inequality, we derive

$$\mathbb{P}\left(\left|\frac{1}{M}\sum_{m=1}^M Y_t\right| \geq \tau\right) \leq 2e^{-\frac{cnM\tau^2}{(B+K \log n)^2}} \implies \mathbb{P}(|\widehat{\mathcal{L}}(\text{TF}) - \mathcal{L}(\text{TF})| \geq \tau) \leq 2e^{-\frac{cnM\tau^2}{(B+K \log n)^2}}$$

where $c > 0$ is an absolute constant. Therefore, we have that for any $\text{TF} \in \text{Alg}$, with probability at least $1 - 2\delta$,

$$|\widehat{\mathcal{L}}(\text{TF}) - \mathcal{L}(\text{TF})| \leq (B + K \log n) \sqrt{\frac{\log(1/\delta)}{cnM}} \quad (11)$$

Step 2: (Bound $\sup_{\text{TF} \in \text{Alg}} |\mathcal{L}(\text{TF}) - \widehat{\mathcal{L}}(\text{TF})|$ where Alg is assumed to be a continuous search space). Let

$$h(\text{TF}) := \mathcal{L}(\text{TF}) - \widehat{\mathcal{L}}(\text{TF})$$

and we aim to bound $\sup_{\text{TF} \in \text{Alg}} |h(\text{TF})|$. Following Definition C.3, for $\varepsilon > 0$, let Alg_ε be a minimal ε -cover of Alg in terms of distance metric ρ . Therefore, Alg_ε is a discrete set with cardinality $|\text{Alg}_\varepsilon| := \mathcal{N}(\text{Alg}, \rho, \varepsilon)$. Then, we have

$$\sup_{\text{TF} \in \text{Alg}} |\mathcal{L}(\text{TF}) - \widehat{\mathcal{L}}(\text{TF})| \leq \sup_{\text{TF} \in \text{Alg}'} \min_{\text{TF}' \in \text{Alg}_\varepsilon} |h(\text{TF}) - h(\text{TF}')| + \max_{\text{TF} \in \text{Alg}_\varepsilon} |h(\text{TF})|.$$

We will first bound the quantity $\sup_{\text{TF} \in \text{Alg}'} \min_{\text{TF}' \in \text{Alg}_\varepsilon} |h(\text{TF}) - h(\text{TF}')|$. We will utilize that loss function $\ell(\cdot, \cdot)$ is C -Lipschitz. For any $\text{TF} \in \text{Alg}$, let $\text{TF}' \in \text{Alg}_\varepsilon$ be its neighbor following Definition C.3. Then we can show that

$$\begin{aligned} & \left| \widehat{\mathcal{L}}(\text{TF}) - \widehat{\mathcal{L}}(\text{TF}') \right| \\ &= \left| \frac{1}{nM} \sum_{m=1}^M \sum_{t=1}^n (\ell(r_{mt}(a_{mt}), \text{TF}(\widehat{r}_{mt}(a_{mt})|\mathcal{H}_m^{t-1}, a_{mt})) - \ell(r_{mt}(a_{mt}), \text{TF}'(\widehat{r}_{mt}(a_{mt})|\mathcal{H}_m^{t-1}, a_{mt}))) \right| \\ &\leq \frac{L}{nM} \sum_{m=1}^M \sum_{t=1}^n \|\text{TF}(\widehat{r}_{mt}(a_{mt})|\mathcal{H}_m^{t-1}, a_{mt}) - \text{TF}'(\widehat{r}_{mt}(a_{mt})|\mathcal{H}_m^{t-1}, a_{mt})\|_{\ell_2} \\ &\leq L\varepsilon. \end{aligned}$$

Note that the above bound applies to all data-sequences, we also obtain that for any $\text{TF} \in \text{Alg}$,

$$|\mathcal{L}(\text{TF}) - \mathcal{L}(\text{TF}')| \leq L\varepsilon.$$

Therefore we can show that,

$$\begin{aligned} & \sup_{\text{TF} \in \text{Alg}} \min_{\text{TF}' \in \text{Alg}_\varepsilon} |h(\text{TF}) - h(\text{TF}')| \\ &\leq \sup_{\text{TF} \in \text{Alg}} \min_{\text{TF}' \in \text{Alg}_\varepsilon} \left| \widehat{\mathcal{L}}(\text{TF}) - \widehat{\mathcal{L}}(\text{TF}') \right| + |\mathcal{L}(\text{TF}) - \mathcal{L}(\text{TF}')| \leq 2L\varepsilon. \quad (12) \end{aligned}$$

Next we bound the second term $\max_{\text{TF} \in \text{Alg}_\varepsilon} |h(\text{TF})|$. Applying union bound directly on Alg_ε and combining it with equation 11, then we will have that with probability at least $1 - 2\delta$,

$$\max_{\text{TF} \in \text{Alg}_\varepsilon} |h(\text{TF})| \leq (B + K \log n) \sqrt{\frac{\log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)/\delta)}{cnM}}$$

Combining the upper bound above with the perturbation bound equation 12, we obtain that

$$\max_{\text{TF} \in \text{Alg}} |h(\text{TF})| \leq 2C\varepsilon + (B + K \log n) \sqrt{\frac{\log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)/\delta)}{cnM}}.$$

It follows then that

$$\mathcal{R}_{\text{MTL}}(\widehat{\text{TF}}) \leq 2 \sup_{\text{TF} \in \text{Alg}} |\mathcal{L}(\text{TF}) - \widehat{\mathcal{L}}(\text{TF})| \leq 4C\varepsilon + 2(B + K \log n) \sqrt{\frac{\log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)/\delta)}{cnM}}$$

Again by setting $\varepsilon = 1/\sqrt{nM}$

$$\mathcal{L}(\widehat{\text{TF}}) - \mathcal{L}(\text{TF}^*) \leq \frac{4C}{\sqrt{nM}} + 2(B + K \log n) \sqrt{\frac{\log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)/\delta)}{cnM}}$$

The claim of the theorem follows. \square

Definition C.2. (Covering number) Let Q be any hypothesis set and $d(q, q') \geq 0$ be a distance metric over $q, q' \in Q$. Then, $\bar{Q} = \{q_1, \dots, q_N\}$ is an ε -cover of Q with respect to $d(\cdot, \cdot)$ if for any $q \in Q$, there exists $q_i \in \bar{Q}$ such that $d(q, q_i) \leq \varepsilon$. The ε -covering number $\mathcal{N}(Q, d, \varepsilon)$ is the cardinality of the minimal ε -cover.

Definition C.3. (Algorithm distance). Let Alg be an algorithm hypothesis set and $\mathcal{H} = (a_t, r_t)_{t=1}^n$ be a sequence that is admissible for some task $m \in [M]$. For any pair $\text{TF}, \text{TF}' \in \text{Alg}$, define the distance metric $\rho(\text{TF}, \text{TF}') := \sup_{\mathcal{H}} \frac{1}{n} \sum_{t=1}^n \|\text{TF}(\hat{r}_t | \mathcal{H}^{t-1}, a_t) - \text{TF}'(\hat{r}_t | \mathcal{H}^{t-1}, a_t)\|_{\ell_2}$.

Remark C.4. (Stability Factor) The work of [Li et al. \(2023\)](#) also characterizes the stability factor K in Assumption 7.1 with respect to the transformer architecture. Assuming loss $\ell(\cdot, \cdot)$ is C -Lipschitz, the algorithm induced by $\text{TF}(\cdot)$ obeys the stability assumption with $K = 2C((1 + \Gamma)e^\Gamma)^L$, where the norm of the transformer weights are upper bounded by $O(\Gamma)$ and there are L -layers of the transformer.

Remark C.5. (Covering Number) From Lemma 16 of [Lin et al. \(2023\)](#) we have the following upper bound on the covering number of the transformer class TF_Θ as

$$\log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)) \leq O(L^2 D^2 J)$$

where L is the total number of layers of the transformer and J and D denote the upper bound to the number of heads and hidden neurons in all the layers respectively. Note that this covering number holds for the specific class of transformer architecture discussed in section 2 of [\(Lin et al., 2023\)](#).

C.2 GENERALIZATION ERROR TO NEW TASK

Theorem C.6. (Transfer Risk) Consider the setting of Theorem 7.2 and assume the source tasks are independently drawn from task distribution \mathcal{T} . Let $\widehat{\text{TF}}$ be the empirical solution of (ERM) and $g \sim \mathcal{T}$. Then with probability at least $1 - 2\delta$, the expected excess transfer learning risk is bounded by

$$\mathbb{E}_g \left[\mathcal{R}_g(\widehat{\text{TF}}) \right] \leq 4 \frac{C}{\sqrt{M}} + B \sqrt{\frac{2 \log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)/\delta)}{M}}$$

where, $\mathcal{N}(\text{Alg}, \rho, \varepsilon)$ is the covering number of transformer $\widehat{\text{TF}}$.

Proof. Let the target task g be sampled from \mathcal{T} , and the test set $\mathcal{H}_g = \{a_t, r_t\}_{t=1}^n$. Define empirical and population risks on g as $\widehat{\mathcal{L}}_g(\text{TF}) = \frac{1}{n} \sum_{t=1}^n \ell(r_t(a_{mt}), \text{TF}(\hat{r}_t(a_{mt}) | \mathcal{H}_g^{t-1}, a_t))$ and $\mathcal{L}_g(\text{TF}) = \mathbb{E}_{\mathcal{H}_g} [\widehat{\mathcal{L}}_g(\text{TF})]$. Again we drop Θ from the transformer notation. Then the expected excess transfer risk following (ERM) is defined as

$$\mathbb{E}_g \left[\mathcal{R}_g(\widehat{\text{TF}}) \right] = \mathbb{E}_{\mathcal{H}_g} \left[\mathcal{L}_g(\widehat{\text{TF}}) \right] - \arg \min_{\text{TF} \in \text{Alg}} \mathbb{E}_{\mathcal{H}_g} [\mathcal{L}_g(\text{TF})]. \quad (13)$$

where \mathcal{A} is the set of all algorithms. The goal is to show a bound like this

$$\mathbb{E}_g \left[\mathcal{R}_g(\widehat{\text{TF}}) \right] \leq \min_{\varepsilon \geq 0} \left\{ 4C\varepsilon + B \sqrt{\frac{2 \log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)/\delta)}{T}} \right\}$$

where $\mathcal{N}(\text{Alg}, \rho, \varepsilon)$ is the covering number.

Step 1 ((Decomposition): Let $\text{TF}^* = \arg \min_{\text{TF} \in \text{Alg}} \mathbb{E}_g [\mathcal{L}_g(\text{TF})]$. The expected transfer learning excess test risk of given algorithm $\widehat{\text{TF}} \in \text{Alg}$ is formulated as

$$\widehat{\mathcal{L}}_m(\text{TF}) := \frac{1}{n} \sum_{t=1}^n \ell(r_{mt}(a_{mt}), \text{TF}(\hat{r}_{mt}(a_{mt}) | \mathcal{D}_m^{t-1}, a_{mt})), \quad \text{and}$$

$$\mathcal{L}_m(\text{TF}) := \mathbb{E}_{\mathcal{H}_m} \left[\widehat{\mathcal{L}}_m(\text{TF}) \right] = \mathbb{E}_{\mathcal{H}_m} \left[\frac{1}{n} \sum_{t=1}^n \ell(r_{mt}(a_{mt}), \text{TF}(\hat{r}_{mt}(a_{mt}) | \mathcal{D}_m^{t-1}, a_{mt})) \right], \quad \forall m \in [M].$$

Then we can decompose the risk as

$$\begin{aligned} \mathbb{E}_g \left[\mathcal{R}_g(\widehat{\text{TF}}) \right] &= \mathbb{E}_g \left[\mathcal{L}_g(\widehat{\text{TF}}) \right] - \mathbb{E}_g \left[\mathcal{L}_g(\text{TF}^*) \right] \\ &= \underbrace{\mathbb{E}_g \left[\mathcal{L}_g(\widehat{\text{TF}}) \right] - \widehat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\widehat{\text{TF}})}_a + \underbrace{\widehat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\widehat{\text{TF}}) - \widehat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\text{TF}^*)}_b + \underbrace{\widehat{\mathcal{L}}_{\mathcal{H}_{\text{all}}}(\text{TF}^*) - \mathbb{E}_g \left[\mathcal{L}_g(\text{TF}^*) \right]}_c. \end{aligned}$$

Here since $\widehat{\text{TF}}$ is the minimizer of training risk, $b < 0$. Then we obtain

$$\mathbb{E}_g \left[\mathcal{R}_g(\widehat{\text{TF}}) \right] \leq 2 \sup_{\text{TF} \in \text{Alg}} \left| \mathbb{E}_g \left[\mathcal{L}_g(\text{TF}) \right] - \frac{1}{M} \sum_{m=1}^M \widehat{\mathcal{L}}_m(\text{TF}) \right|. \quad (14)$$

Step 2 (Bounding equation 14) For any $\text{TF} \in \text{Alg}$, let $X_t = \widehat{\mathcal{L}}_t(\text{TF})$ and we observe that

$$\mathbb{E}_{m \sim \mathcal{T}} [X_t] = \mathbb{E}_{m \sim \mathcal{T}} \left[\widehat{\mathcal{L}}_m(\text{TF}) \right] = \mathbb{E}_{m \sim \mathcal{T}} \left[\mathcal{L}_m(\text{TF}) \right] = \mathbb{E}_g \left[\mathcal{L}_g(\text{TF}) \right]$$

Since $X_m, m \in [M]$ are independent, and $0 \leq X_m \leq B$, applying Hoeffding's inequality obeys

$$\mathbb{P} \left(\left| \mathbb{E}_g \left[\mathcal{L}_g(\text{TF}) \right] - \frac{1}{M} \sum_{m=1}^M \widehat{\mathcal{L}}_m(\text{TF}) \right| \geq \tau \right) \leq 2e^{-\frac{2M\tau^2}{B^2}}.$$

Then with probability at least $1 - 2\delta$, we have that for any $\text{TF} \in \text{Alg}$,

$$\left| \mathbb{E}_g \left[\mathcal{L}_g(\text{TF}) \right] - \frac{1}{M} \sum_{m=1}^M \widehat{\mathcal{L}}_m(\text{TF}) \right| \leq B \sqrt{\frac{\log(1/\delta)}{2M}}. \quad (15)$$

Next, let Alg_ε be the minimal ε -cover of Alg following Definition C.2, which implies that for any task $g \sim \mathcal{T}$, and any $\text{TF} \in \text{Alg}$, there exists $\text{TF}' \in \text{Alg}_\varepsilon$

$$\left| \mathcal{L}_g(\text{TF}) - \mathcal{L}_g(\text{TF}') \right|, \left| \widehat{\mathcal{L}}_g(\text{TF}) - \widehat{\mathcal{L}}_g(\text{TF}') \right| \leq C\varepsilon. \quad (16)$$

Since the distance metric following Definition 3.4 is defined by the worst-case datasets, then there exists $\text{TF}' \in \text{Alg}_\varepsilon$ such that

$$\left| \mathbb{E}_g \left[\mathcal{L}_g(\text{TF}) \right] - \frac{1}{M} \sum_{m=1}^M \widehat{\mathcal{L}}_m(\text{TF}) \right| \leq 2C\varepsilon.$$

Let $\mathcal{N}(\text{Alg}, \rho, \varepsilon) = |\text{Alg}_\varepsilon|$ be the ε -covering number. Combining the above inequalities (equation 14, equation 15, and equation 16), and applying union bound, we have that with probability at least $1 - 2\delta$,

$$\mathbb{E}_g \left[\mathcal{R}_g(\widehat{\text{TF}}) \right] \leq \min_{\varepsilon \geq 0} \left\{ 4C\varepsilon + B \sqrt{\frac{2 \log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)/\delta)}{M}} \right\}$$

Again by setting $\varepsilon = 1/\sqrt{M}$

$$\mathcal{L}(\widehat{\text{TF}}) - \mathcal{L}(\text{TF}^*) \leq \frac{4C}{\sqrt{M}} + 2B \sqrt{\frac{\log(\mathcal{N}(\text{Alg}, \rho, \varepsilon)/\delta)}{cM}}$$

The claim of the theorem follows. \square

Remark C.7. (Dependence on n) In this remark, we briefly discuss why the expected excess risk for target task \mathcal{T} does not depend on samples n . The work of Li et al. (2023) pointed out that the MTL pretraining process identifies a favorable algorithm that lies in the span of the M source tasks. This is termed as inductive bias (see section 4 of Li et al. (2023)) (Soudry et al., 2018; Neyshabur et al., 2017). Such bias would explain the lack of dependence of the expected excess transfer risk on n during transfer learning. This is because given a target task $g \sim \mathcal{T}$, the TF can use the learnt favorable algorithm to conduct a discrete search over span of the M source tasks and return the source task that best fits the new target task. Due to the discrete search space over the span of M source tasks, it is not hard to see that, we need $n \propto \log(M)$ samples (which is guaranteed by the M source tasks) rather than $n \propto d$ (for the linear setting).

C.3 TABLE OF NOTATIONS

Notations	Definition
M	Total number of tasks
d	Dimension of the feature
\mathcal{A}_m	Action set of the m -th task
\mathcal{X}_m	Feature space of m -th task
M_{test}	Tasks for testing
M_{pre}	Total Tasks for pretraining
$\mathbf{x}(m, a)$	Feature of action a in task m
$\theta_{m,*}$	Hidden parameter for the task m
\mathcal{T}_{pre}	Pretraining distribution on tasks
$\mathcal{T}_{\text{test}}$	Testing distribution on tasks
n	Total horizon for each task m
$\mathcal{H}_m = \{I_t, r_t\}_{t=1}^n$	Dataset sampled for the m -th task containing n samples
$\mathcal{H}_m^t = \{I_s, r_s\}_{s=1}^t$	Dataset sampled for the m -th task containing samples from round $s = 1$ to t
\mathbf{w}	Transformer model parameter
$\text{TF}_{\mathbf{w}}$	Transformer with model parameter \mathbf{w}
\mathcal{D}_{pre}	Pretraining in-context distribution
$\mathcal{H}_{\text{train}}$	Training in-context dataset
$\mathcal{D}_{\text{test}}$	Testing in-context distribution

Table 1: Table of Notations