
Supervised Pretraining Can Learn In-Context Reinforcement Learning

Jonathan N. Lee^{*1} Annie Xie^{*1} Aldo Pacchiano² Yash Chandak¹
Chelsea Finn¹ Ofir Nachum³ Emma Brunskill¹

Abstract

Large transformer models trained on diverse datasets have shown a remarkable ability to *learn in-context*, achieving high few-shot performance on tasks they were not explicitly trained to solve. In this paper, we study the in-context learning capabilities of transformers in decision-making problems, i.e., reinforcement learning (RL) for bandits and Markov decision processes. To do so, we introduce and study *Decision-Pretrained Transformer (DPT)*, a supervised pretraining method where the transformer predicts an optimal action given a query state and an in-context dataset of interactions, across a diverse set of tasks. This procedure, while simple, produces a model with several surprising capabilities. We find that the pretrained transformer can be used to solve a range of RL problems in-context, exhibiting both exploration online and conservatism offline, despite not being explicitly trained to do so. The model also generalizes beyond the pretraining distribution to new tasks and automatically adapts its decision-making strategies to unknown structure. Theoretically, we show DPT can be viewed as an efficient implementation of Bayesian posterior sampling, a provably sample-efficient RL algorithm. We further leverage this connection to provide guarantees on the regret of the in-context algorithm yielded by DPT, and prove that it can learn faster than algorithms used to generate the pretraining data. These results suggest a promising yet simple path towards instilling strong in-context decision-making abilities in transformers.

1. Introduction

For supervised learning, transformer-based models trained at scale have shown impressive abilities to perform tasks given an input context, often referred to as few-shot prompting or in-context learning (Brown et al., 2020). In this setting, a pretrained model is presented with a small number of supervised input-output examples in its context, and is then asked to predict the most likely completion (i.e. output) of an unpaired input, without parameter updates. Over the last few years in-context learning has been applied to solve a range of tasks (Min et al., 2022) and a growing number works are beginning to understand and analyze in-context learning for supervised learning (Chan et al., 2022; Garg et al., 2022; Razeghi et al., 2022; Akyürek et al., 2022). In this work, our focus is to study and understand in-context learning applied to sequential decision-making, specifically in the context of reinforcement learning (RL) settings. Decision-making (e.g. RL) is considerably more dynamic and complex than supervised learning. Understanding and leveraging in-context learning here could potentially unlock significant improvements in an agent’s ability to adapt and make few-shot decisions in response to observations from the world. Such capabilities are instrumental for practical applications ranging from robotics to recommendation systems.

For in-context decision-making (Laskin et al., 2022; Xu et al., 2022; 2023), rather than input-output tuples, the context takes the form of state-action-reward tuples representing a dataset of interactions with an unknown environments. The agent must leverage these interactions to understand the dynamics of the world and what actions lead to good outcomes. A hallmark of good decision-making in online RL algorithms is a judicious balance of selecting exploratory actions to gather information and selecting increasingly optimal actions by exploiting that information (Sutton & Barto, 2018). In contrast, an RL agent with access to only a suboptimal offline dataset should produce a policy that conservatively selects actions (Levine et al., 2020). An ideal in-context decision-maker should exhibit similar behaviors.

To study in-context decision-making formally, we propose a new simple supervised pretraining objective, namely, to train (via supervised learning) a transformer to predict an

^{*}Equal contribution ¹Stanford University ²Microsoft Research ³Google DeepMind. Correspondence to: JNL <jnl@stanford.edu>, AX <annixie@stanford.edu>.

optimal action label¹ given a query state and an in-context dataset of interactions, across a diverse set of tasks. We refer to the pretrained model as a Decision-Pretrained Transformer (DPT). Once trained, DPT can be deployed as either an online or offline RL algorithm in a new task by passing it an in-context dataset of interactions and querying it for predictions of the optimal action in different states. For example, online, the in-context dataset is initially empty and DPT’s predictions are uncertain because the new task is unknown, but it fills the dataset with its interactions as it learns and becomes more confident about the optimal action. We show empirically and theoretically that DPT yields a surprisingly effective in-context decision-maker with regret guarantees. As it turns out, DPT effectively performs posterior sampling — a provably sample-efficient Bayesian RL algorithm that has historically been limited by its computational burden (Osband et al., 2013). We summarize our main findings below.

- **Predicting optimal actions alone gives rise to near-optimal decision-making algorithms.** The DPT objective is solely based on predicting optimal actions from in-context interactions. At the outset, it is not immediately apparent that these predictions at test-time would yield good decision-making when the task is unknown and behaviors such as online exploration are necessary to solve it. Intriguingly, DPT as an algorithm is capable of dealing with this uncertainty in-context. For example, despite not being explicitly trained to explore, DPT exhibits an exploration strategy on par with hand-designed algorithms, as a means to discover the optimal actions.
- **DPT generalizes to new decision-making problems, offline and online.** We show DPT can handle reward distributions unseen in its pretraining data on bandit problems as well as unseen goals, dynamics, and datasets in simple MDPs. This suggests that the in-context strategies learned during pretraining are robust and generalizable without any parameter updates at test time.
- **DPT improves over the data used to pretrain it by exploiting latent structure.** As an example, in parametric bandit problems, specialized algorithms can leverage structure (such as linear rewards) and offer provably better regret, but a representation must be known in advance. Perhaps surprisingly, we find that pretraining on linear bandit problems, even with unknown representations, leads DPT to select actions and explore in a way that matches an efficient linear bandit algorithm. This holds even when the source pretraining data comes from a suboptimal algorithm (i.e., one that does not exploit any structure), demonstrating the ability to learn improved in-context strategies beyond what it was trained on.
- **Posterior sampling can be implemented via in-context**

¹If not explicitly known, the optimal action can be determined by running any (potentially inefficient) minimax-optimal regret algorithm for each pretraining task.

learning. Posterior sampling (PS), a generalization of Thompson Sampling, can provably sample-efficiently solve online RL problems (Osband et al., 2013), but a common criticism is the lack of computationally efficient ways to update and sample from a posterior distribution. DPT can be viewed as learning a posterior distribution over optimal actions, shortcutting the PS procedure. Under some conditions, we show theoretically that DPT in-context is equivalent to PS. Furthermore, DPT’s prior and posterior updates are grounded in data rather than needing to be specified *a priori*. This suggests that in-context learning could help unlock practical and efficient RL via posterior sampling.

2. Related Work

Meta-learning. Algorithmically, in-context learning falls under the meta-learning framework (Schaul & Schmidhuber, 2010; Bengio et al., 1990). At a high-level, these methods attempt to learn some underlying shared structure of the training distribution of tasks to accelerate learning of new tasks. For decision-making and RL, there is a often choice in what shared ‘structure’ is specifically learned such as the dynamics of the task (Fu et al., 2016; Nagabandi et al., 2018; Landolfi et al., 2019), a task context identifier (Rakelly et al., 2019; Humplik et al., 2019; Zintgraf et al., 2019; Liu et al., 2021), temporally extended skills and options (Perkins et al., 1999; Gupta et al., 2018; Jiang et al., 2022), or initialization of a neural network policy (Finn et al., 2017; Rothfuss et al., 2018)). In-context learning can be viewed as taking a more agnostic approach by learning the learning algorithm itself, more similar to (Duan et al., 2016; Wang et al., 2016; Mishra et al., 2017). Algorithm Distillation (AD) (Laskin et al., 2022; Lu et al., 2023) also falls under this category, applying autoregressive supervised learning to distill (sub-sampled) traces of a single-task RL algorithm into a task-agnostic model. While DPT also leverages autoregressive SL, it does not distill an existing RL algorithm in order to imitate how to learn. Instead, we pretrain DPT to predict optimal actions, yielding potentially emergent online and offline strategies at test time that automatically leverage the task structure to behave similarly to posterior sampling.

Autoregressive transformers for decision-making. In decision-making fields such as RL and imitation learning, transformer models trained using autoregressive supervised action prediction have proliferated (Yang et al., 2023), inspired by the successes of these techniques for large language models (Vaswani et al., 2017; Raffel et al., 2020; Brown et al., 2020). For example, Decision Transformer (DT) (Chen et al., 2021; Janner et al., 2021) uses a transformer to autoregressively model sequences of actions from offline experience data, conditioned on the achieved return. During inference, one can then query the model conditioned on a desired return value. This approach has been shown to

scale favorably to large models and multi-task settings (Lee et al., 2022), at times exceeding the performance of large-scale multi-task imitation learning with transformers (Reed et al., 2022; Brohan et al., 2022; Shafiq et al., 2022). However, DT is known to be provably (and unboundedly) sub-optimal in common scenarios (Brandfonbrener et al., 2022; Yang et al., 2022). A common criticism of DT, and supervised learned transformers in general, is their inability to improve upon the dataset. For example, there is little reason for DT to output meaningful behavior if conditioned on return higher than any observed in training, without strong extrapolation assumptions (Brandfonbrener et al., 2022). In contrast, a major contribution of our work is theoretical and empirical evidence for the ability of DPT to improve over behaviors seen in the dataset in terms of regret.

Value and policy-based offline RL. Offline RL algorithms offer the opportunity to learn from existing datasets. To address distributional shift, many prior algorithms incorporate the principle of value pessimism (Kumar et al., 2020; Yu et al., 2021; Liu et al., 2020; Ghasemipour et al., 2022), or policy regularization (Fujimoto et al., 2019; Kumar et al., 2019; Wu et al., 2019; Siegel et al., 2020; Liu et al., 2019). To reduce the amount of offline data required in a new task, methods for offline meta-RL can reuse interactions collected in a set of related tasks (Li et al., 2020; Mitchell et al., 2021; Dorfman et al., 2021). However, they still must address distribution shift, requiring solutions such as policy regularization (Li et al., 2020) or additional online interactions (Pong et al., 2022). DPT follows the success of autoregressive models like DT and AD, avoiding these issues. With our pretraining objective, DPT also leverages offline datasets for new tasks more effectively than AD.

3. In-Context Learning Model

Basic decision models. The basic decision model of our study is the finite-horizon Markov decision process (MDP). An MDP is specified by the tuple $\tau = \langle \mathcal{S}, \mathcal{A}, T, R, H, \rho \rangle$ to be solved, where \mathcal{S} is the state space, \mathcal{A} is the action space, $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the transition function, $R : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathbb{R})$ is the reward function, $H \in \mathbb{N}$ is the horizon, and $\rho \in \Delta(\mathcal{S})$ is the initial state distribution. A learner interacts with the environment through the following protocol: (1) an initial state s_1 is sampled from ρ ; (2) at time step h , the learner chooses an action a_h and transitions to state $s_{h+1} \sim T(\cdot | s_h, a_h)$, and receives a reward $r_h \sim R(\cdot | s_h, a_h)$. The episode ends after H steps. A policy π maps states to distributions over actions and can be used to interact with the MDP. We denote the optimal policy as π^* , which maximizes the value function $V(\pi^*) = \max_{\pi} V(\pi) := \max_{\pi} \mathbb{E}_{\pi} \sum_h r_h$. When necessary, we use the subscript τ to distinguish V_{τ} and π_{τ}^* for the specific MDP τ . We assume the state space is partitioned by

$h \in [H]$ so that π^* is notationally independent of h . Note this framework encompasses multi-armed bandit settings where the state space is a single point, e.g. $\mathcal{S} = \{1\}$, $H = 1$, and the optimal policy is $a^* \in \operatorname{argmax}_{a \in \mathcal{A}} \mathbb{E} [r_1 | a_1 = a]$.

Algorithm 1 Decision-Pretrained Transformer (DPT): Training and Deployment

```

1: // Collecting pretraining dataset
2: Initialize empty pretraining dataset  $\mathcal{B}$ 
3: for  $i$  in  $[N]$  do
4:   Sample task  $\tau \sim \mathcal{T}_{\text{pre}}$ , in-context dataset  $D \sim \mathcal{D}_{\text{pre}}(\cdot; \tau)$ , query state  $s_{\text{query}} \sim \mathcal{D}_{\text{query}}$ 
5:   Sample label  $a^* \sim \pi_{\tau}^*(\cdot | s_{\text{query}})$  and add  $(s_{\text{query}}, D, a^*)$  to  $\mathcal{B}$ 
6: end for
7: // Pretraining model on dataset
8: Initialize model  $M_{\theta}$  with parameters  $\theta$ 
9: while not converged do
10:  Sample  $(s_{\text{query}}, D, a^*)$  from  $\mathcal{B}$  and predict  $\hat{p}_j(\cdot) = M_{\theta}(\cdot | s_{\text{query}}, D_j)$  for all  $j \in [n]$ 
11:  Compute loss in (1) with respect to  $a^*$  and backpropagate to update  $\theta$ .
12: end while
13: // Offline test-time deployment
14: Sample unknown task  $\tau \sim \mathcal{T}_{\text{test}}$ , sample dataset  $D \sim \mathcal{D}_{\text{test}}(\cdot; \tau)$ 
15: Deploy  $M_{\theta}$  in  $\tau$  by choosing  $a_h \in \operatorname{argmax}_{a \in \mathcal{A}} M_{\theta}(a | s_h, D)$  at step  $h$ 
16: // Online test-time deployment
17: Sample unknown task  $\tau \sim \mathcal{T}_{\text{test}}$  and initialize empty  $D = \{\}$ 
18: for  $\text{ep}$  in  $\text{max\_eps}$  do
19:   Deploy  $M_{\theta}$  by sampling  $a_h \sim M_{\theta}(\cdot | s_h, D)$  at step  $h$ 
20:   Add  $(s_1, a_1, r_1, \dots)$  to  $D$ 
21: end for

```

Pretraining. We give pseudocode in Algorithm 1 and a visualization in Figure 1. Let \mathcal{T}_{pre} be a distribution over tasks at the time of pretraining. A task $\tau \sim \mathcal{T}_{\text{pre}}$ can be viewed as a specification of an MDP, $\tau = \langle \mathcal{S}, \mathcal{A}, T, R, H, \rho \rangle$. The distribution \mathcal{T}_{pre} can span different reward and transition functions and even different state and action spaces. We then sample a context (or a prompt) which consists of a dataset $D \sim \mathcal{D}_{\text{pre}}(\cdot; \tau)$ of interactions between the learner and the MDP specified by τ . $D = \{s_j, a_j, s'_j, r_j\}_{j \in [n]}$ is a collection of transition tuples taken in τ . We refer to D as the *in-context dataset* because it provides the contextual information about τ . D could be generated through variety of means, such as: (1) random interactions within τ , (2) demonstrations from an expert, and (3) rollouts of an algorithm. Additionally, we independently sample a query state s_{query} from the distribution $\mathcal{D}_{\text{query}}$ over states \mathcal{S} and a label a^* is sampled from the optimal policy $\pi_{\tau}^*(\cdot | s_{\text{query}})$ for task τ (see Section 5.3 for how to implement this in com-

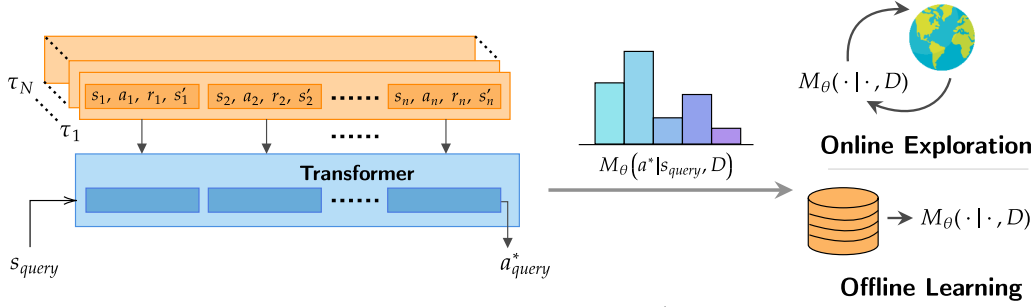


Figure 1. A transformer model M_θ is pretrained to predict an optimal action a_{query}^* from a state s_{query} in a task, given a dataset of interactions from that task. The resulting Decision-Pretrained Transformer (DPT) learns a distribution over the optimal action conditioned on an in-context dataset. M_θ can be deployed in *new* tasks online by collecting data on the fly, or offline by immediately conditioning on a static dataset.

mon practical scenarios). We denote the joint pretraining distribution over tasks, in-context datasets, query states, and action labels as P_{pre} :

$$\begin{aligned} P_{\text{pre}}(\tau, D, s_{\text{query}}, a^*) \\ = \mathcal{T}_{\text{pre}}(\tau) \mathcal{D}_{\text{pre}}(D; \tau) \mathcal{D}_{\text{query}}(s_{\text{query}}) \pi_\tau^*(a^* | s_{\text{query}}) \end{aligned}$$

Given the in-context dataset D and a query state s_{query} , we can train a model to predict the optimal action a^* in response simply via supervised learning. Let $D_j = \{(s_1, a_1, s'_1, r_1), \dots, (s_j, a_j, s'_j, r_j)\}$ denote the partial dataset up to j samples. Formally, we aim to train a causal GPT-2 transformer model M parameterized by θ , which outputs a distribution over actions \mathcal{A} , to minimize the expected loss over samples from the pretraining distribution:

$$\min_{\theta} \mathbb{E}_{P_{\text{pre}}} \sum_{j \in [n]} \ell(M_\theta(\cdot | s_{\text{query}}, D_j), a^*) \quad (1)$$

Generally, we set the loss to be the negative log-likelihood with $\ell(M_\theta(\cdot | s_{\text{query}}, D_j), a^*) := -\log M_\theta(a^* | s_{\text{query}}, D_j)$. This framework can work for both discrete and continuous \mathcal{A} . For our experiments with discrete \mathcal{A} , we use a softmax parameterization for the distribution of M_θ , essentially treating this as a classification problem. The resulting output model M_θ can be viewed as an algorithm that takes in a dataset of interactions D and can be queried with a forward pass for predictions of the optimal action via inputting a query state s_{query} . We refer to the trained model M_θ as a Decision-Pretrained Transformer (DPT).

Testing. After pretraining, a new task (MDP) τ is sampled from a test-task distribution $\mathcal{T}_{\text{test}}$. If the DPT is to be tested *offline*, then a dataset (prompt) is a sampled $D \sim \mathcal{D}_{\text{test}}(\cdot; \tau)$ and the policy that the model in-context learns is given conditionally as $M_\theta(\cdot | \cdot, D)$. Namely, we evaluate the policy by selecting action $a_h \in \arg\max_a M_\theta(a | s_h, D)$ when the learner visits state s_h . If the model is to be tested *online* through multiple episodes of interaction, then the dataset is initialized as empty $D = \{\}$. At each episode, $M_\theta(\cdot | \cdot, D)$ is deployed where the model samples $a_h \sim M_\theta(\cdot | s_h, D)$ upon observing state s_h . Throughout a full episode, it collects interactions $\{s_1, a_1, r_1, \dots, s_H, a_H, r_H\}$ which are

subsequently appended to D . The model then repeats the process with another episode, and so on until a specified number of episodes has been reached. A key distinction of the testing phase is that there are no updates to the parameters of M_θ . This is in contrast to hand-designed RL algorithms that would perform parameter updates or maintain statistics using D to learn from scratch. Instead, the model M_θ performs a computation through its forward pass to generate a distribution over actions conditioned on the in-context D and query state s_h .

4. Learning in Bandits

We begin with an empirical investigation of DPT in a multi-armed bandit, a well-studied special case of the MDP where the state space \mathcal{S} is a singleton and the horizon $H = 1$ is a single step. We will examine the performance of DPT both when aiming to select a good action from offline historical data and for online learning where the goal is to maximize cumulative reward from scratch. Offline, it is critical to account for uncertainty due to noise as certain actions may not be sampled well enough. Online, it is critical to judiciously balance exploration and exploitation to minimize overall regret. For detailed descriptions of the experiment setups, see Appendix A.

Pretraining distribution. For the pretraining task distribution \mathcal{T}_{pre} , we sample 5-armed bandits ($|\mathcal{A}| = 5$). The reward function for arm a is a normal distribution $R(\cdot | s, a) = \mathcal{N}(\mu_a, \sigma^2)$ where $\mu_a \sim \text{Unif}[0, 1]$ independently and $\sigma = 0.3$. To generate in-context datasets \mathcal{D}_{pre} , we randomly generate action frequencies by sampling probabilities from a Dirichlet distribution and mixing them with a point-mass distribution on one random arm (see details in Appendix A.3). Then we sample the actions accordingly from this distribution. This encourages diversity of the in-context datasets. The optimal policy π_τ^* for bandit τ is $\arg\max_a \mu_a$, which we can easily compute during pretraining. We pretrain the model M_θ to predict a^* from D as described in Section 3 for datasets up to size $n = 500$.

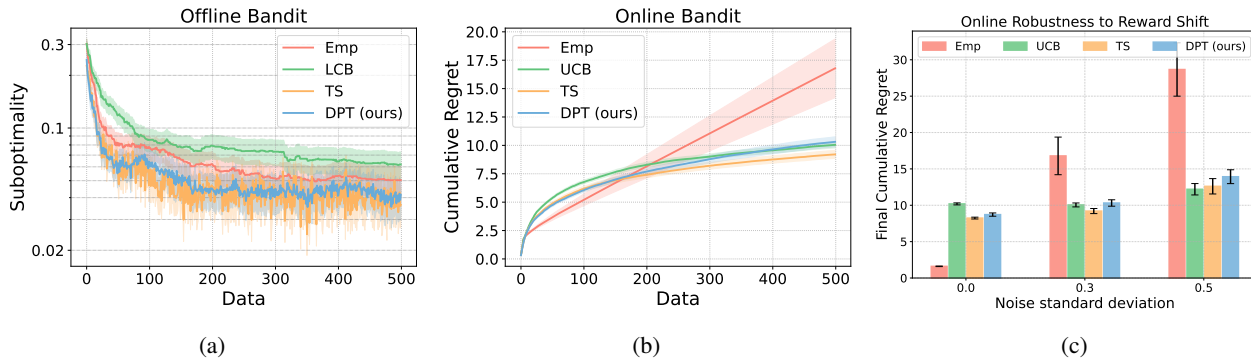


Figure 2. (a) Offline performance on in-distribution bandits, given random in-context datasets. (b) Online cumulative regret on bandits. (c) Final (after 500 steps) cumulative regret on out-of-distribution bandits with different Gaussian noise standard deviations. The mean and standard error are computed over 200 test tasks.

Comparisons. We compare to several well-known algorithms for bandits². All of the algorithms are designed to reason in a particular way about uncertainty based on their observations.

- Empirical mean algorithm (Emp) selects the action with the highest empirical mean reward naively.
- Upper Confidence Bound (UCB) selects the action with the highest upper confidence bound.
- Lower Confidence Bound (LCB) selects the action with the highest lower confidence bound.
- Thompson Sampling (TS) selects the action with the highest sampled mean from a posterior distribution over reward models. The prior and likelihood functions are Gaussian.

Emp and TS (Russo et al., 2018; Thompson, 1933) can both be used for offline or online learning; UCB (Auer et al., 2002) is known to be provably optimal online by ensuring exploration through optimism under uncertainty; and LCB (Xiao et al., 2021; Jin et al., 2021) is used to minimize suboptimality given an offline dataset by selecting actions pessimistically. It is the opposite of UCB. We evaluate algorithms with standard bandit metrics. Offline, we use the suboptimality $\mu_{a^*} - \mu_{\hat{a}}$ where \hat{a} is the chosen action. Online, we use cumulative regret: $\sum_k \mu_{a^*} - \mu_{\hat{a}_k}$ where \hat{a}_k is the k th action chosen.

DPT learns to reason through uncertainty. As shown in Figure 2a, in the offline setting, DPT significantly exceeds the performance of Emp and LCB while matching the performance of TS, when the in-context datasets are sampled from the same distribution as during pretraining. The results suggest that the transformer is capable of reasoning through uncertainty caused by the noisy rewards in the dataset. Unlike Emp which can be fooled by noisy, undersampled actions, the transformer has learned to *hedge* to a degree. However, it also suggests that this hedging is fundamentally different from what LCB does, at least on this

²See Appendix A.2 for additional details such as hyperparameters.

specific distribution³. Interestingly, the same transformer produces an extremely effective online bandit algorithm when sampling actions instead of taking an argmax. As shown in Figure 2b, DPT matches the performance of classical optimal algorithms, UCB and TS, which are specifically designed for exploration. This is notable because DPT was not explicitly trained to explore, but its emergent strategy is on par with some of the best. In Figure 2c, we show this property is robust to noise in the rewards not seen during pre-training by varying the standard deviation. In Appendix B, we show this generalization happens offline too and even with unseen Bernoulli rewards.

Leveraging structure from suboptimal data. We now investigate whether DPT can learn to leverage the inherent structure of a problem class, even without prior knowledge of this structure and even when learning from in-context datasets that do not explicitly utilize it. More precisely, we consider \mathcal{T}_{pre} to be a distribution over *linear* bandits, where the reward function is given by $\mathbb{E}[r | a, \tau] = \langle \theta_\tau, \phi(a) \rangle$ and $\theta_\tau \in \mathbb{R}^d$ is a task-specific parameter vector and $\phi : \mathcal{A} \rightarrow \mathbb{R}^d$ is fixed feature vector that is the same for all tasks. Given the feature representation ϕ , LinUCB (Abbasi-Yadkori et al., 2011), a UCB-style algorithm that leverages ϕ , should achieve regret $\tilde{O}(d\sqrt{K})$ over K steps, a substantial gain over UCB and TS when $d \ll |\mathcal{A}|$. Here, we pretrain a DPT model with in-context datasets gathered by TS, which does not leverage the linear structure. Figures 3a and 3b show that DPT can exploit the unknown linear structure, essentially learning a surrogate for ϕ , allowing it to do more informed exploration online and decision-making offline. It is nearly on par with LinUCB (which is given ϕ) and significantly outperforms the source, TS, which does not know or use the structure. These results suggest that (1) DPT can automatically leverage structure, and (2) supervised learning-based approaches to RL *can* learn novel explorations that tran-

³Note our randomly generated environments are equally likely to have expert-biased datasets and adversarial datasets, so LCB is not expected to outperform here (Xiao et al., 2021).

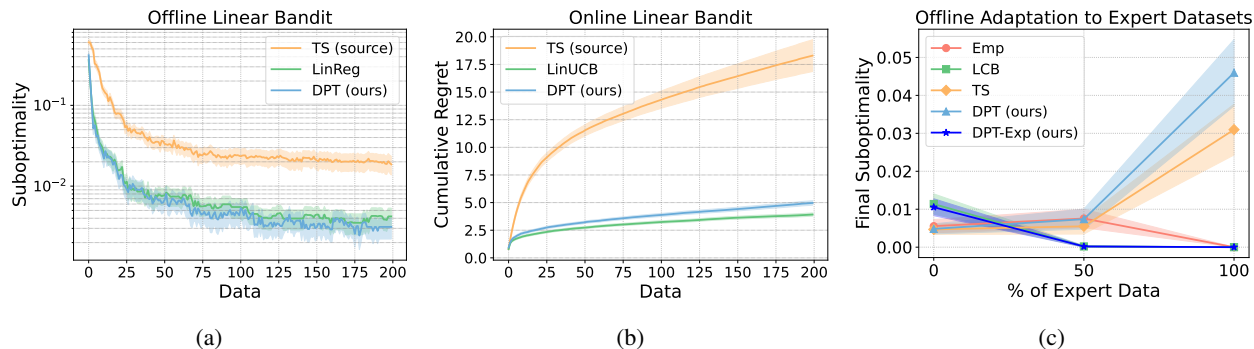


Figure 3. (a) Offline performance of DPT trained on linear bandits from TS source data. LinReg does linear regression and outputs the greedy action. (b) Online cumulative regret of the same model. The mean and standard error are computed over 200 test tasks. (c) Offline performance on expert-biased datasets. DPT pretrained on a different prior continues to match TS, but DPT-Exp trained from a more representative prior excels.

scend the quality of their pretraining data.

Adapting to expert-biased datasets. A common assumption in offline RL is that datasets tend to be a mixture between optimal data (e.g. expert demonstrations) and suboptimal data (e.g. random interactions) (Rashidinejad et al., 2021). Hence, LCB is generally effective in practice and the pretraining and testing distributions should be biased towards this setting. Motivated by this, we pretrain a second DPT model where \mathcal{D}_{pre} is generated by mixing the in-context datasets with varying fractions of expert data, biasing \mathcal{D}_{pre} towards datasets that contain more examples of the optimal action. We denote this model by DPT-Exp. In Figure 3c, we plot the test-time performance of both pretrained models when evaluated on new offline datasets with varying percentages of expert data⁴. Our results suggest that when the pretraining distribution is also biased towards expert-suboptimal data, DPT-Exp behaves similarly to LCB, while DPT continues to resemble TS. This is quite interesting as for other methods, such as TS, it is less clear how to automatically incorporate the right amount of expert bias to get the same effect, but DPT can learn this from pretraining.

5. Learning in Markov Decision Processes

We next study how DPT can tackle Markov decision processes by testing its ability to perform exploration and credit assignment. In the following experiments, the DPT demonstrates generalization to new tasks, scalability to image-based observations, and capability to stitch in-context behaviors (Section 5.2). This section also examines whether DPT can be pretrained with datasets and action labels generated by a different RL algorithm, rather than the exact optimal policy (Section 5.3).

⁴That is, 0% is fully random while 100% has only optimal actions in the in-context dataset.

5.1. Experimental Setup

Environments. We consider environments that require targeted exploration to solve the task. The first is Dark Room (Zintgraf et al., 2019; Laskin et al., 2022), a 2D discrete environment where the agent must locate the unknown goal location in a 10×10 room, and only receives a reward of 1 when at the goal. We hold out a set of goals for generalization evaluation. Our second environment is Miniworld (Chevalier-Boisvert, 2018), a 3D visual navigation problem to test the scalability of DPT to image observations. The agent is in a room with four boxes of different colors, and must find the target box, the color of which is unknown to the agent initially. It receives a reward of 1 only when near the correct box. Details on these environments and the pre-training datasets are in App. A.4 and A.5.

Comparisons. Our experiments aim to understand the effectiveness of DPT in comparison to that of other context-based meta-RL algorithms. To that end, we compare to meta-RL algorithms based on supervised and RL objectives.

- Proximal Policy Optimization (PPO) (Schulman et al., 2017): We compare to this single-task RL algorithm, which trains from scratch without any pretraining data, to contextualize the performance of DPT and other meta-RL algorithms.
- Algorithm Distillation (AD) (Laskin et al., 2022): AD first generates a dataset of learning histories by running an RL algorithm in each training task. Then, given a sampled subsequence $h_j = (s_j, a_j, r_j, \dots, s_{j+c})$ from a learning history, a transformer is trained to predict the next action a_{j+c} from the learning history.
- RL^2 (Duan et al., 2016): This online meta-RL comparison uses a recurrent neural network to adapt the agent’s policy from the given context. Unlike AD and DPT, which are trained with a supervised objective, the RL^2 agent is trained to maximize the expected return with PPO.

PPO and RL^2 are online algorithms, while AD is capable of

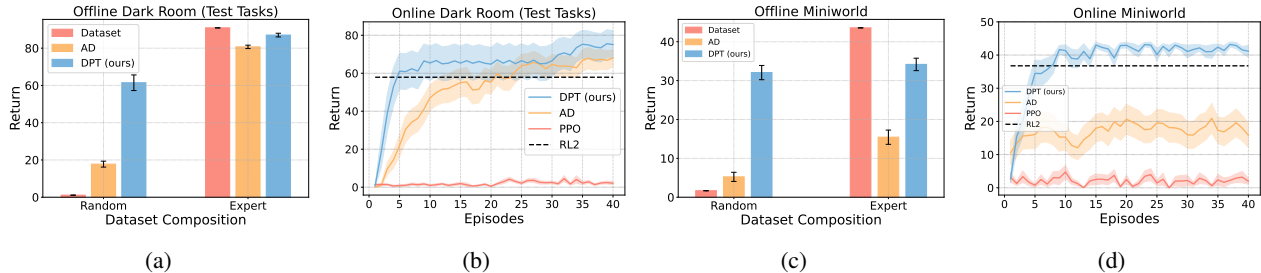


Figure 4. (a) Offline performance on held-out Dark Room goals, given random and expert datasets. (b) Online performance on held-out Dark Room goals. (c) Offline performance on Miniworld (images), given random and expert datasets. (d) Online performance on Miniworld (images) after 40 episodes. We report the average and standard error of the mean over 100 different offline datasets in (a) and (c) and 20 online trials in (b) and (d).

learning both offline and online. Details on the implementation of these algorithms can be found in Appendix A.2.

5.2. Main Results

Generalizing to new offline datasets and tasks. To study the generalization capabilities of DPT, we evaluate the model in Dark Room on a set of 20 held-out goals not in the pretraining dataset. When given an expert dataset, DPT achieves near-optimal performance. Even when given a random dataset, which has an average total reward of 1.1, DPT obtains a much higher average return of 61.5 (see Fig. 4a). Qualitatively, we observe that when the in-context dataset contains a transition to the goal, DPT immediately exploits this and takes a direct path to the goal. In contrast, while AD demonstrates strong offline performance with expert data, it performs worse in-context learning with random data compared to DPT. The difference arises because AD is trained to infer a better policy than the in-context data, but not necessarily the optimal one.

We next evaluate DPT, AD, RL², and PPO online without any prior data from the 20 test-time Dark Room tasks, shown in Fig. 4b. After 40 episodes, PPO does not make significant progress towards the goal, highlighting the difficulty of learning from such few interactions alone. RL² is trained to perform adaptation within four episodes each of length 100, and we report the performance after the four adaptation episodes. Notably, DPT on average solves each task faster than AD and reaches a higher final return than RL², demonstrating its capability to explore effectively online even in MDPs. In Appendix B, we also present results on generalization to new dynamics.

Learning from image-based observations. In Miniworld, the agent receives RGB image observations of 25×25 pixels. As shown in Fig. 4d, DPT can solve this high-dimensional task offline from both random and expert datasets. Compared to AD and RL², DPT also learns online more efficiently.

Stitching novel trajectories from in-context subse-

quences. A desirable ability of some offline RL algorithms stitching suboptimal subsequences from the offline dataset into new trajectories with higher return. To test whether DPT exhibits stitching, we design the *Dark Room (Three Tasks)* environment in which there are three possible tasks. The pretraining data consists only of expert demonstrations of two of them. At test-time DPT is evaluated on third unseen task, but its offline dataset is only expert demonstrations of the original two. Despite this, it leverages the data to infer a path solving the third task (see Fig. 5a).

5.3. Learning from Algorithm-Generated Policies and Rollouts

So far, we have only considered action labels provided by an optimal policy. However, in some tasks, an optimal policy is not readily available even in pretraining. In this experiment, we use actions labeled by a policy learned via PPO and in-context datasets sampled from PPO replay buffers. We train PPO agents in each of the 80 train tasks for 1K episodes to generate 80K total rollouts, from which we sample the in-context datasets. This variant, DPT (PPO, PPO), performs on par with DPT and still better than AD, as shown in Figures 5b and 5c. DPT (PPO, PPO) can be viewed as a direct comparison between our pretraining objective and that of AD, given the same pretraining data but just used differently. We also evaluated a variant, DPT (Rand, PPO), which pretrains on random in-context datasets (like DPT), but still using PPO action labels. The performance is worse than the other DPT variants in some settings, but only marginally so. In Appendix B, we analyze the sensitivity of DPT to other hyperparameters.

6. Theory

We now shed light on the observations of the previous empirical results through a theoretical analysis. Our main result shows that DPT (under a slight modification to pretraining) essentially performs in-context posterior sampling (PS). PS is a generalization of Thompson Sampling for RL in

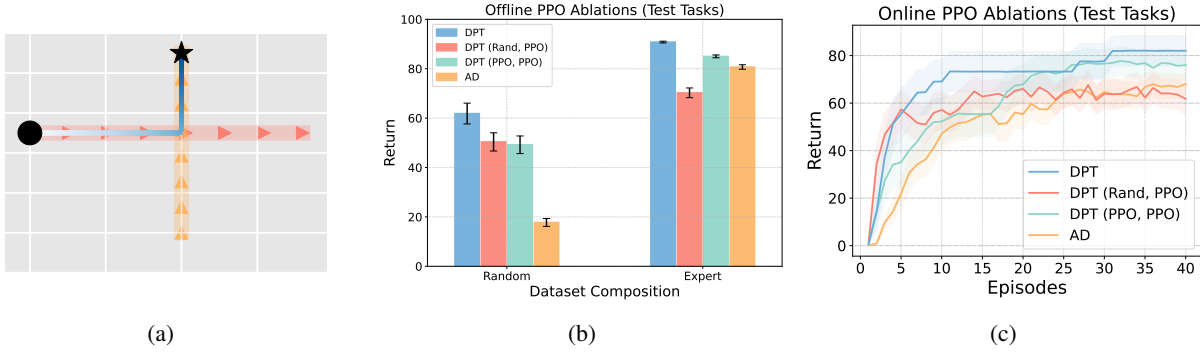


Figure 5. (a) In *Dark Room (Three Tasks)*, DPT stitches a new, optimal trajectory to the goal (blue) given two in-context demonstrations of other tasks (pink and orange). (b) Offline *Dark Room* performance of DPT trained on PPO data. (c) Online *Dark Room* performance of DPT trained on PPO data.

MDPs. It maintains and samples from a posterior over tasks τ given historical data D and executes optimal policies π_τ^* (see Appendix C for a formal outline). It is provably sample-efficient with online Bayesian regret guarantees (Osband et al., 2013), but maintaining posteriors is generally computationally intractable. The ability for DPT to perform PS in-context suggests a path towards computation- and provably sample-efficient RL with priors learned from the data.

6.1. History-Dependent Pretraining and Assumptions

We start with a modification to the pretraining of DPT. Rather than conditioning only on s_{query} and D to predict $a^* \sim \pi_\tau^*(\cdot | s_{\text{query}})$, we propose also conditioning on a sequence $\xi_h = (s_{1:h}, a_{1:h}^*)$ where $s_{1:h} \sim \mathfrak{S}_h \in \Delta(\mathcal{S}^h)$ is a distribution over sets of states, independent of τ , and $a_{h'}^* \sim \pi_\tau^*(\cdot | s_{h'})$ for $h' \in [h]$. Thus, we use π_τ^* to label both the query state (which is the prediction label) and the sequence of states sampled from \mathfrak{S}_h . Note that this does not require any environment interactions and hence no sampling from either T_τ or R_τ . At test-time at step h , this will allow us to condition on the history ξ_{h-1} of states that M_θ visits and the actions that it takes in those states. Formally, the learned M_θ is deployed as follows, given D . (1) At $h = 0$, initialize $\xi_0 = ()$ to be empty. (2) At step h , visit s_h and find a_h by sampling from $M_\theta(\cdot | s_{\text{query}}, D, \xi_{h-1})$. (3) Append (s_h, a_h) to ξ_{h-1} to get ξ_h . Note for bandits and contextual bandits ($H = 1$), there is no difference between this and the original pretraining procedure of prior sections because ξ_0 is empty. For MDPs, the original DPT can be viewed as a convenient approximation.

We now make several assumptions to simplify the analysis. First, assume $\mathcal{D}_{\text{query}}$, \mathcal{D}_{pre} , and \mathfrak{S} have sufficient support such that all conditional probabilities of P_{pre} are well defined. Similar to other studies of in-context learning (Xie et al., 2021), we assume M_θ fits the pretraining distribution exactly with enough coverage and data, so that the focus of

the analysis is just the in-context learning abilities.

Assumption 6.1. (Learned model is consistent). Let M_θ denote the pretrained model. For all $(s_{\text{query}}, D, \xi_h)$, we have $P_{\text{pre}}(a | s_{\text{query}}, D, \xi_h) = M_\theta(a | s_{\text{query}}, D, \xi_h)$ for all $a \in \mathcal{A}$.

To provide some cursory justification, if M_θ is the global minimizer of (1), then $\mathbb{E}_{P_{\text{pre}}} \|P_{\text{pre}}(\cdot | s_{\text{query}}, D, \xi_h) - M_\theta(\cdot | s_{\text{query}}, D, \xi_h)\|_1^2 \rightarrow 0$ as the number of pretraining samples $N \rightarrow \infty$ with high probability for transformer model classes of bounded complexity (see Proposition C.1). Approximate versions of the above assumptions are easily possible but obfuscate the key elements of the analysis. We also assume that the in-context dataset $D \sim \mathcal{D}_{\text{pre}}$ is compliant (Jin et al., 2021), meaning that the actions from D can depend only on the observed history and not additional confounders. Note that this still allows \mathcal{D}_{pre} to be very general — it could be generated randomly or from adaptive algorithms like PPO or TS.

Definition 6.2 (Compliance). The in-context dataset distribution $\mathcal{D}_{\text{pre}}(\cdot; \tau)$ is *compliant* if, for all $i \in [n]$, the i th action of the dataset, a_i , is conditionally independent of τ given the i th state s_i and partial dataset, D_{i-1} , so far. In other words, $\mathcal{D}_{\text{pre}}(a_i | s_i, D_{i-1}; \tau)$ is invariant to τ .

Generally, \mathcal{D}_{pre} can influence M_θ . In Proposition 6.6, we show that all compliant \mathcal{D}_{pre} form a sort of equivalence class that generate the same M_θ . For the remainder, we assume all \mathcal{D}_{pre} are compliant.

6.2. Main Results

Equivalence of DPT and PS. We now state our main result which shows that the trajectories generated by a pretrained M_θ will follow the same distribution as those from a well-specified PS algorithm. In particular, let PS use the well-specified prior \mathcal{T}_{pre} . Let τ_c be an arbitrary task. Let $P_{\text{ps}}(\cdot | D, \tau_c)$ and $P_{M_\theta}(\cdot | D, \tau_c)$ denote the distributions over trajectories $\xi_H \in (\mathcal{S} \times \mathcal{A})^H$ generated from running PS and $M_\theta(\cdot | \cdot, D, \cdot)$, respectively, in task τ_c given historical

data D .

Theorem 6.3 (DPT \iff PS). *Let the above assumptions hold. Then, $P_{ps}(\xi_H \mid D, \tau_c) = P_{M_\theta}(\xi_H \mid D, \tau_c)$ for all trajectories ξ_H .*

Regret implications. To see this result in action, let us specialize to the finite MDP setting (Osband et al., 2013). Suppose we pretrain M_θ on a distribution \mathcal{T}_{pre} over MDPs with $S := |\mathcal{S}|$ and $A := |\mathcal{A}|$. Let \mathcal{D}_{pre} be constructed by uniform sampling (s_i, a_i) and observing (r_i, s'_i) for $i \in [KH]$. Let $\mathbb{E}[r_h | s_h, a_h] \in [0, 1]$. And let \mathcal{D}_{query} and \mathcal{S}_h be uniform over \mathcal{S} and \mathcal{S}^h (for all h) respectively. Finally, let \mathcal{T}_{test} be the distribution over test tasks with the same cardinalities. For a task τ , define the online cumulative regret of DPT over K episodes as $\text{Reg}_\tau(M_\theta) := \sum_{k \in [K]} V_\tau(\pi_\tau^*) - V_\tau(\hat{\pi}_k)$ where $\hat{\pi}_k(\cdot | s_h) = M_\theta(\cdot | s_h, D_{(k-1)}, \xi_{h-1})$ and $D_{(k)}$ contains the first k episodes collected from $\hat{\pi}_{1:k}$.

Corollary 6.4 (Finite MDPs). *Suppose that $\sup_\tau \mathcal{T}_{test}(\tau) / \mathcal{T}_{pre}(\tau) \leq C$ for some $C > 0$. For the above MDP setting, the pretrained model M_θ satisfies $\mathbb{E}_{\mathcal{T}_{test}}[\text{Reg}_\tau(M_\theta)] \leq \tilde{O}(CH^{3/2}S\sqrt{AK})$.*

A similar analysis due to (Russo & Van Roy, 2014) allows us to prove why pretraining on (latently) linear bandits can lead to substantial empirical gains, even when the in-context datasets are generated by algorithms unaware of this structure. We observed this empirically in Section 4. Consider a similar setup as there where \mathcal{S} is a singleton, \mathcal{A} is finite but large, $\theta_\tau \in \mathbb{R}^d$ is sampled as $\theta_\tau \sim \mathcal{N}(0, I/d)$, $\phi : \mathcal{A} \rightarrow \mathbb{R}^d$ is a fixed feature map with $\sup_{a \in \mathcal{A}} \|\phi(a)\|_2 \leq 1$, and the reward of $a \in \mathcal{A}$ in task τ is distributed as $\mathcal{N}(\langle \theta_\tau, \phi(a) \rangle, 1)$. This time, we let $\mathcal{D}_{pre}(\cdot; \tau)$ be given by running Thompson Sampling with Gaussian priors and likelihood functions on τ .

Corollary 6.5 (Latent representation learning in linear bandits). *For $\mathcal{T}_{test} = \mathcal{T}_{pre}$ in the above linear bandit setting, M_θ satisfies $\mathbb{E}_{\mathcal{T}_{test}}[\text{Reg}_\tau(M_\theta)] \leq \tilde{O}(d\sqrt{K})$.*

This significantly improves over the $\tilde{O}(\sqrt{|\mathcal{A}|K})$ upper regret bound for TS that does not leverage the linear structure. This highlights how DPT can have provably tighter upper bounds on future bandit problems than the algorithms used to generate its (pretraining) data.

Invariance of M_θ to compliant \mathcal{D}_{pre} . Our final result sheds light on how \mathcal{D}_{pre} impacts the final DPT behavior M_θ . Combined with Assumption 6.1, M_θ is invariant to \mathcal{D}_{pre} satisfying Definition 6.2.

Proposition 6.6. *Let P_{pre}^1 and P_{pre}^2 be pretraining distributions that differ only by their in-context dataset distributions, denoted by \mathcal{D}_{pre}^1 and \mathcal{D}_{pre}^2 . If \mathcal{D}_{pre}^1 and \mathcal{D}_{pre}^2 are compliant with the same support, then $P_{pre}^1(a^* | s_{query}, D, \xi_h) = P_{pre}^2(a^* | s_{query}, D, \xi_h)$ for all a^*, s_{query}, D, ξ_h .*

That is, if we generate in-context datasets D by running various algorithms that depend only on the observed data in the current task, we will end up with the same M_θ . For example, TS could be used for \mathcal{D}_{pre}^1 and PPO for \mathcal{D}_{pre}^2 . Expert-biased datasets discussed in Section 4 violate Definition 6.2, since knowledge of τ is being used. This helps explain our empirical results that pretraining on expert-biased datasets leads to a qualitatively different learned model at test-time.

7. Discussion

In this paper, we studied the problem of in-context decision-making. We introduced a new pretraining method and transformer model, DPT, which is trained via supervised learning to predict optimal actions given an in-context dataset of interactions. Through in-depth evaluations in classic decision problems in bandits and MDPs, we showed that this simple objective naturally gives rise to an in-context RL algorithm that is capable of online exploration and offline decision-making, unlike other algorithms that are explicitly trained or designed to do these. Our empirical and theoretical results provide first steps towards understanding these capabilities that arise from DPT and what factors are important for it to succeed. The inherent strength of pretraining lies in its simplicity—we can sidestep the complexities of hand-designing exploration or conservatism in RL algorithms and while simultaneously allowing the transformer to derive novel strategies that best leverage problem structure. These findings underscore the potential of supervised pretraining in equipping transformer models with in-context decision-making abilities.

Limitations and future work. One limitation of DPT is the requirement of optimal actions at pretraining. Empirically, we find that this requirement can be relaxed by using actions generated by another RL-trained agent during pretraining, which only leads to a slight loss in performance. However, fully understanding this problem and how best to leverage multi-task decision-making datasets remains a key open problem. We also discussed that the practical implementation for MDPs differs from true posterior sampling. It would be interesting to further understand and bridge this empirical-theoretical gap in the future. Our preliminary analysis shows promise for DPT to generalize to new tasks beyond its pretraining distribution. This suggests that diversifying the task distributions during pretraining could significantly enhance the model’s ability to generalize to new tasks. This possibility holds an exciting avenue for future work. Finally, further investigation is required to understand the implications of these findings for existing foundation models, such as instruction-finetuned models, that are increasingly being deployed in decision-making settings (Wang et al., 2023).

References

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.
- Abernethy, J., Agarwal, A., Marinov, T. V., and Warmuth, M. K. A mechanism for sample-efficient in-context learning for sparse retrieval tasks. *arXiv preprint arXiv:2305.17040*, 2023.
- Agarwal, A., Kakade, S., Krishnamurthy, A., and Sun, W. Flambe: Structural complexity and representation learning of low rank mdps. *Advances in neural information processing systems*, 33:20095–20107, 2020.
- Agrawal, S. and Goyal, N. Near-optimal regret bounds for thompson sampling. *Journal of the ACM (JACM)*, 64(5): 1–24, 2017.
- Agrawal, S. and Jia, R. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. *Advances in Neural Information Processing Systems*, 30, 2017.
- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47:235–256, 2002.
- Bengio, Y., Bengio, S., and Cloutier, J. *Learning a synaptic learning rule*. Citeseer, 1990.
- Brandfonbrener, D., Bietti, A., Buckman, J., Laroche, R., and Bruna, J. When does return-conditioned supervised learning work for offline reinforcement learning? *arXiv preprint arXiv:2206.01079*, 2022.
- Brohan, A., Brown, N., Carbajal, J., Chebotar, Y., Dabis, J., Finn, C., Gopalakrishnan, K., Hausman, K., Herzog, A., Hsu, J., et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Chan, S., Santoro, A., Lampinen, A., Wang, J., Singh, A., Richemond, P., McClelland, J., and Hill, F. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891, 2022.
- Chen, L., Lu, K., Rajeswaran, A., Lee, K., Grover, A., Laskin, M., Abbeel, P., Srinivas, A., and Mordatch, I. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- Chevalier-Boisvert, M. Miniworld: Minimalistic 3d environment for rl and robotics research, 2018.
- Dorfman, R., Shenfeld, I., and Tamar, A. Offline meta reinforcement learning—identifiability challenges and effective data collection strategies. *Advances in Neural Information Processing Systems*, 34:4607–4618, 2021.
- Duan, Y., Schulman, J., Chen, X., Bartlett, P. L., Sutskever, I., and Abbeel, P. RL2: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Fu, J., Levine, S., and Abbeel, P. One-shot learning of manipulation skills with online dynamics adaptation and neural network priors. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4019–4026. IEEE, 2016.
- Fujimoto, S., Meger, D., and Precup, D. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*, pp. 2052–2062. PMLR, 2019.
- Garg, S., Tsipras, D., Liang, P. S., and Valiant, G. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Ghasemipour, K., Gu, S. S., and Nachum, O. Why so pessimistic? estimating uncertainties for offline rl through ensembles, and why their independence matters. *Advances in Neural Information Processing Systems*, 35: 18267–18281, 2022.
- Goyal, A., Friesen, A., Banino, A., Weber, T., Ke, N. R., Badia, A. P., Guez, A., Mirza, M., Humphreys, P. C., Konyushova, K., et al. Retrieval-augmented reinforcement learning. In *International Conference on Machine Learning*, pp. 7740–7765. PMLR, 2022.
- Gupta, A., Mendonca, R., Liu, Y., Abbeel, P., and Levine, S. Meta-reinforcement learning of structured exploration strategies. *Advances in neural information processing systems*, 31, 2018.

- Humplik, J., Galashov, A., Hasenclever, L., Ortega, P. A., Teh, Y. W., and Heess, N. Meta reinforcement learning as task inference. *arXiv preprint arXiv:1905.06424*, 2019.
- Janner, M., Li, Q., and Levine, S. Offline reinforcement learning as one big sequence modeling problem. *Advances in neural information processing systems*, 34: 1273–1286, 2021.
- Jiang, Y., Liu, E., Eysenbach, B., Kolter, J. Z., and Finn, C. Learning options via compression. *Advances in Neural Information Processing Systems*, 35:21184–21199, 2022.
- Jin, Y., Yang, Z., and Wang, Z. Is pessimism provably efficient for offline rl? In *International Conference on Machine Learning*, pp. 5084–5096. PMLR, 2021.
- Ke, N. R., Chiappa, S., Wang, J., Goyal, A., Bornschein, J., Rey, M., Weber, T., Botvinic, M., Mozer, M., and Rezende, D. J. Learning to induce causal structure. *arXiv preprint arXiv:2204.04875*, 2022.
- Kirsch, L., Harrison, J., Sohl-Dickstein, J., and Metz, L. General-purpose in-context learning by meta-learning transformers. *arXiv preprint arXiv:2212.04458*, 2022.
- Kumar, A., Fu, J., Soh, M., Tucker, G., and Levine, S. Stabilizing off-policy q-learning via bootstrapping error reduction. *Advances in Neural Information Processing Systems*, 32, 2019.
- Kumar, A., Zhou, A., Tucker, G., and Levine, S. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33: 1179–1191, 2020.
- Landolfi, N. C., Thomas, G., and Ma, T. A model-based approach for sample-efficient multi-task reinforcement learning. *arXiv preprint arXiv:1907.04964*, 2019.
- Laskin, M., Wang, L., Oh, J., Parisotto, E., Spencer, S., Steigerwald, R., Strouse, D., Hansen, S., Filos, A., Brooks, E., et al. In-context reinforcement learning with algorithm distillation. *arXiv preprint arXiv:2210.14215*, 2022.
- Lee, K.-H., Nachum, O., Yang, M. S., Lee, L., Freeman, D., Guadarrama, S., Fischer, I., Xu, W., Jang, E., Michalewski, H., et al. Multi-game decision transformers. *Advances in Neural Information Processing Systems*, 35: 27921–27936, 2022.
- Levine, S., Kumar, A., Tucker, G., and Fu, J. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- Li, L., Yang, R., and Luo, D. Focal: Efficient fully-offline meta-reinforcement learning via distance metric learning and behavior regularization. *arXiv preprint arXiv:2010.01112*, 2020.
- Li, Y., Ildiz, M. E., Papailiopoulos, D., and Oymak, S. Transformers as algorithms: Generalization and implicit model selection in in-context learning. *arXiv preprint arXiv:2301.07067*, 2023.
- Liu, E. Z., Raghunathan, A., Liang, P., and Finn, C. Decoupling exploration and exploitation for meta-reinforcement learning without sacrifices. In *International conference on machine learning*, pp. 6925–6935. PMLR, 2021.
- Liu, Y., Swaminathan, A., Agarwal, A., and Brunskill, E. Off-policy policy gradient with state distribution correction. *UAI*, 2019.
- Liu, Y., Swaminathan, A., Agarwal, A., and Brunskill, E. Provably good batch off-policy reinforcement learning without great exploration. *Advances in neural information processing systems*, 33:1264–1274, 2020.
- Lu, C., Schroecker, Y., Gu, A., Parisotto, E., Foerster, J., Singh, S., and Behbahani, F. Structured state space models for in-context reinforcement learning. *arXiv preprint arXiv:2303.03982*, 2023.
- Lu, X. and Van Roy, B. Ensemble sampling. *Advances in neural information processing systems*, 30, 2017.
- Min, S., Lyu, X., Holtzman, A., Artetxe, M., Lewis, M., Hajishirzi, H., and Zettlemoyer, L. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- Mishra, N., Rohaninejad, M., Chen, X., and Abbeel, P. A simple neural attentive meta-learner. *arXiv preprint arXiv:1707.03141*, 2017.
- Mitchell, E., Rafailov, R., Peng, X. B., Levine, S., and Finn, C. Offline meta-reinforcement learning with advantage weighting. In *International Conference on Machine Learning*, pp. 7780–7791. PMLR, 2021.
- Nagabandi, A., Clavera, I., Liu, S., Fearing, R. S., Abbeel, P., Levine, S., and Finn, C. Learning to adapt in dynamic, real-world environments through meta-reinforcement learning. *arXiv preprint arXiv:1803.11347*, 2018.
- Nguyen, T. and Grover, A. Transformer neural processes: Uncertainty-aware meta learning via sequence modeling. *arXiv preprint arXiv:2207.04179*, 2022.
- Olsson, C., Elhage, N., Nanda, N., Joseph, N., DasSarma, N., Henighan, T., Mann, B., Askell, A., Bai, Y., Chen, A., et al. In-context learning and induction heads. *arXiv preprint arXiv:2209.11895*, 2022.
- Osband, I., Russo, D., and Van Roy, B. (more) efficient reinforcement learning via posterior sampling. *Advances in Neural Information Processing Systems*, 26, 2013.

- Osband, I., Blundell, C., Pritzel, A., and Van Roy, B. Deep exploration via bootstrapped dqn. *Advances in neural information processing systems*, 29, 2016.
- Osband, I., Aslanides, J., and Cassirer, A. Randomized prior functions for deep reinforcement learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Perkins, T. J., Precup, D., et al. Using options for knowledge transfer in reinforcement learning. Technical report, Citeseer, 1999.
- Pong, V. H., Nair, A. V., Smith, L. M., Huang, C., and Levine, S. Offline meta-reinforcement learning with on-line self-supervision. In *International Conference on Machine Learning*, pp. 17811–17829. PMLR, 2022.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- Raffin, A., Hill, A., Gleave, A., Kanervisto, A., Ernestus, M., and Dormann, N. Stable-baselines3: Reliable reinforcement learning implementations. *The Journal of Machine Learning Research*, 22(1):12348–12355, 2021.
- Rakelly, K., Zhou, A., Finn, C., Levine, S., and Quillen, D. Efficient off-policy meta-reinforcement learning via probabilistic context variables. In *International conference on machine learning*, pp. 5331–5340. PMLR, 2019.
- Rashidinejad, P., Zhu, B., Ma, C., Jiao, J., and Russell, S. Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information Processing Systems*, 34:11702–11716, 2021.
- Razeghi, Y., Logan IV, R. L., Gardner, M., and Singh, S. Impact of pretraining term frequencies on few-shot reasoning. *arXiv preprint arXiv:2202.07206*, 2022.
- Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S. G., Novikov, A., Barth-Maron, G., Gimenez, M., Sulsky, Y., Kay, J., Springenberg, J. T., et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- Rothfuss, J., Lee, D., Clavera, I., Asfour, T., and Abbeel, P. Prompt: Proximal meta-policy search. *arXiv preprint arXiv:1810.06784*, 2018.
- Russo, D. and Van Roy, B. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- Russo, D. J., Van Roy, B., Kazerouni, A., Osband, I., Wen, Z., et al. A tutorial on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.
- Schaul, T. and Schmidhuber, J. Metalearning. *Scholarpedia*, 5(6):4650, 2010.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Shafiq, N. M., Cui, Z., Altanzaya, A. A., and Pinto, L. Behavior transformers: Cloning k modes with one stone. *Advances in neural information processing systems*, 35:22955–22968, 2022.
- Shin, S., Lee, S.-W., Ahn, H., Kim, S., Kim, H., Kim, B., Cho, K., Lee, G., Park, W., Ha, J.-W., et al. On the effect of pretraining corpora on in-context learning by a large-scale language model. *arXiv preprint arXiv:2204.13509*, 2022.
- Siegel, N. Y., Springenberg, J. T., Berkenkamp, F., Abdolmaleki, A., Neunert, M., Lampe, T., Hafner, R., Heess, N., and Riedmiller, M. Keep doing what worked: Behavioral modelling priors for offline reinforcement learning. *arXiv preprint arXiv:2002.08396*, 2020.
- Strens, M. A bayesian framework for reinforcement learning. In *ICML*, volume 2000, pp. 943–950, 2000.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Thompson, W. R. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A., and Vladymyrov, M. Transformers learn in-context by gradient descent. *arXiv preprint arXiv:2212.07677*, 2022.
- Wang, G., Xie, Y., Jiang, Y., Mandlkar, A., Xiao, C., Zhu, Y., Fan, L., and Anandkumar, A. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.

- Wang, J. X., Kurth-Nelson, Z., Tirumala, D., Soyer, H., Leibo, J. Z., Munos, R., Blundell, C., Kumaran, D., and Botvinick, M. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.
- Wies, N., Levine, Y., and Shashua, A. The learnability of in-context learning. *arXiv preprint arXiv:2303.07895*, 2023.
- Wu, Y., Tucker, G., and Nachum, O. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- Xiao, C., Wu, Y., Mei, J., Dai, B., Lattimore, T., Li, L., Szepesvari, C., and Schuurmans, D. On the optimality of batch policy optimization algorithms. In *International Conference on Machine Learning*, pp. 11362–11371. PMLR, 2021.
- Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An explanation of in-context learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*, 2021.
- Xu, M., Shen, Y., Zhang, S., Lu, Y., Zhao, D., Tenenbaum, J., and Gan, C. Prompting decision transformer for few-shot policy generalization. In *International Conference on Machine Learning*, pp. 24631–24645. PMLR, 2022.
- Xu, M., Lu, Y., Shen, Y., Zhang, S., Zhao, D., and Gan, C. Hyper-decision transformer for efficient online policy adaptation. *arXiv preprint arXiv:2304.08487*, 2023.
- Yang, M., Schuurmans, D., Abbeel, P., and Nachum, O. Dichotomy of control: Separating what you can control from what you cannot. *arXiv preprint arXiv:2210.13435*, 2022.
- Yang, S., Nachum, O., Du, Y., Wei, J., Abbeel, P., and Schuurmans, D. Foundation models for decision making: Problems, methods, and opportunities. *arXiv preprint arXiv:2303.04129*, 2023.
- Yu, T., Kumar, A., Rafailov, R., Rajeswaran, A., Levine, S., and Finn, C. Combo: Conservative offline model-based policy optimization. *Advances in neural information processing systems*, 34:28954–28967, 2021.
- Zintgraf, L., Shiarlis, K., Igl, M., Schulze, S., Gal, Y., Hofmann, K., and Whiteson, S. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. *arXiv preprint arXiv:1910.08348*, 2019.
- Zintgraf, L., Shiarlis, K., Igl, M., Schulze, S., Gal, Y., Hofmann, K., and Whiteson, S. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. In *International Conference on Learning Representation (ICLR)*, 2020.

Acknowledgments

We thank Evan Liu, Sherry Yang, and Lucy Shi for helpful discussions and feedback. This work was supported in part by NSF grant 2112926 and ONR grant N00014-21-1-2685. JNL acknowledges support from the NSF GRFP.

Additional Related Work

In-context learning. Beyond decision-making and reinforcement learning, our approach takes inspiration from general in-context learning, a phenomenon observed most prominently in large language models in which large-scale autoregressive modelling can surprisingly lead to a model that exhibits meta-learning capabilities (Brown et al., 2020). Recently, there has been great interest in understanding the capabilities and properties of in-context learning in various settings (Garg et al., 2022; von Oswald et al., 2022; Razeghi et al., 2022; Olsson et al., 2022; Kirsch et al., 2022; Shin et al., 2022; Li et al., 2023; Wies et al., 2023; Akyürek et al., 2022; Abernethy et al., 2023; Ke et al., 2022; Goyal et al., 2022). While a common hypothesis suggests that this phenomenon is due to properties of the data used to train large language models (Chan et al., 2022), our work suggests that this phenomenon can also be encouraged in general settings via adjustments to the pre-training objective. In fact, DPT could be interpreted as explicitly encouraging the ability to perform Bayesian inference, which is a popular explanation for the mechanism behind in-context learning for large language models (Xie et al., 2021). For meta-learning in-context, work by (Nguyen & Grover, 2022) proposed Transformer Neural Processes to address problems traditionally addressed with neural processes. Their work focuses on supervised learning-style input-output pairs for meta-learning, but they also applied this to contextual bandits to learn rewards.

Posterior Sampling. Posterior sampling originates from the seminal work of (Thompson, 1933), and has been popularized and thoroughly investigated in recent years by a number of authors (Russo et al., 2018; Agrawal & Goyal, 2017; Strens, 2000; Osband et al., 2013; Agrawal & Jia, 2017; Russo & Van Roy, 2014). For bandits, it is often referred to as Thompson Sampling, but the framework is easily generalizable to RL. The principle is as follows: begin with a prior over possible models (i.e. reward and transition functions), and maintain a posterior distribution over models by updating as new interactions are made. At decision-time, sample a model from the posterior and execute its optimal policy. The aforementioned prior works have developed strong theoretical guarantees on Bayesian and frequentist regret for posterior sampling. Despite its desirable theoretical characteristics, a major limitation is that computing the posterior is often computationally intractable, leading practitioners to rely on approximation-based solutions (Lu & Van Roy, 2017; Osband et al., 2016; 2018). In Section 6, we show that a version of the DPT model learned from pretraining can be viewed as implementing posterior sampling as it should be without resorting to approximations or deriving complicated posterior updates. Instead, the posterior update is implicitly learned through pretraining to predict the optimal action. This suggests that in-context learning (or meta-learning more generally) could be a key in unlocking practically applicable posterior sampling for RL.

A. Implementation and Experiment Details

Algorithm 2 Decision-Pretrained Transformer (detailed)

```

1: // Collecting pretraining dataset
2: Initialize empty dataset  $\mathcal{B}$ 
3: for  $i$  in  $[N]$  do
4:   Sample task  $\tau \sim \mathcal{T}_{\text{pre}}$ 
5:   Sample interaction dataset  $D \sim \mathcal{D}_{\text{pre}}(\cdot; \tau)$  of length  $n$ 
6:   Sample  $s_{\text{query}} \sim \mathcal{D}_{\text{query}}$  and  $a^* \sim \pi_{\tau}^*(\cdot | s_{\text{query}})$ 
7:   Add  $(s_{\text{query}}, D, a^*)$  to  $\mathcal{B}$ 
8: end for
9: // Training model on dataset
10: Initialize model  $M_{\theta}$  with parameters  $\theta$ 
11: while not converged do
12:   Sample  $(s_{\text{query}}, D, a^*)$  from  $\mathcal{B}$ 
13:   Predict  $\hat{p}_j(\cdot) = M_{\theta}(\cdot | s_{\text{query}}, D_j)$  for all  $j \in [n]$ .
14:   Compute loss in (4) with respect to  $a^*$  and backpropagate to update  $\theta$ .
15: end while

```

Algorithm 3 Offline test-time deployment (detailed)

```

1: // Task and offline dataset are generated without learner's control
2: Sample unknown task  $\tau \sim \mathcal{T}_{\text{test}}$ 
3: Sample dataset  $D \sim \mathcal{D}_{\text{test}}(\cdot; \tau)$ 
4: // Deploying offline policy  $M_\theta(\cdot, D)$ 
5:  $s_1 = \text{reset}(\tau)$ 
6: for  $h$  in  $[H]$  do
7:    $a_h = \text{argmax}_{a \in \mathcal{A}} M_\theta(\cdot | s_h, D)$  // Most likely action
8:    $s_{h+1}, r_h = \text{step}(\tau, a_h)$ 
9: end for

```

Algorithm 4 Online test-time deployment (detailed)

```

1: // Online, dataset is empty as learning is from scratch
2: Initialize  $D = \{\}$ 
3: Sample unknown task  $\tau \sim \mathcal{T}_{\text{test}}$ 
4: for ep in max_eps do
5:    $s_1 = \text{reset}(\tau)$ 
6:   for  $h$  in  $[H]$  do
7:      $a_h \sim M_\theta(\cdot | s_h, D)$  // Sample action from predicted distribution
8:      $s_{h+1}, r_h = \text{step}(\tau, a_h)$ 
9:   end for
10: // Experience from previous episode added to dataset
11: Add  $(s_1, a_1, r_1, \dots)$  to  $D$ 
12: end for

```

A.1. DPT Architecture: Formal Description

In this section, we provide a detailed description of the architecture alluded to in Section 3 and Figure 1. See hyperparameter details for models in their respective sections. The model is implemented in Python with PyTorch (Paszke et al., 2019). The backbone of the transformer architecture we use is an autoregressive GPT-2 model from the HuggingFace `transformers` library.

For the sake of exposition, we suppose that \mathcal{S} and \mathcal{A} are subsets of \mathbb{R}^{d_S} and \mathbb{R}^{d_A} respectively. We handle discrete state and action spaces with one-hot encoding. Consider a single training datapoint derived from an (potentially unknown) task τ : we have a dataset D of interactions within τ , a query state s_{query} , and its corresponding optimal action $a^* = \pi_\tau^*(s_{\text{query}})$. We construct the embeddings to be passed to the GPT-2 backbone in the following way. From the dataset $D = \{(s_j, a_j, s'_j, r_j)\}_{j \in [n]}$, we construct vectors $\xi_j = (s_j, a_j, s'_j, r_j)$ by stacking the elements of the transition tuple into dimension $d_\xi := 2d_S + d_A + 1$ for each j in the sequence. This sequence of n elements is concatenated with another vector $v := (s_{\text{query}}, \mathbf{0})$ where the $\mathbf{0}$ vector is a vector of zeros of sufficient length to make the entire element dimension d_ξ . The $(n+1)$ -length sequence is given by $X = (v, \xi_1, \dots, \xi_n)$. As order does not often matter for the dataset D^5 , we do not use positional encoding in order to take advantage of this invariance. We first apply a linear layer `Linear(X)` and pass the result to the transformer, which outputs the sequence $Y = (\hat{y}_0, \hat{y}_1, \dots, \hat{y}_n)$. In the continuous action case, these can be used as is for predictions of a^* . For the discrete action case, we use them as logits to be converted to either a distribution over actions in \mathcal{A} or one-hot vector predictions of a^* . Here, we compute action probabilities

$$\hat{p}_j = \text{softmax}(\hat{y}_j) \in \Delta(\mathcal{A}) \tag{2}$$

Because of the GPT-2 causal architecture (we defer details to the original papers (Radford et al., 2019; Brown et al., 2020)), we note that \hat{p}_j depends only on s_{query} and the partial dataset $D_j = \{(s_k, a_k, s'_k, r_k)\}_{k \in [j]}$, which is why we write the model notation,

$$M_\theta(\cdot | s_{\text{query}}, D_j) = \hat{p}_j(\cdot), \tag{3}$$

⁵This is not always true such as when data comes from an algorithm such as PPO or Thompson Sampling.

to denote that the predicted probabilities of the j th element only depend on D_j and not the entire D for the model M with parameters $\theta \in \Theta$. For example, with $j = 0$, the prediction of a^* is made without any contextual information about the task τ except for s_{query} , which can be interpreted as the prior over a^* . We measure loss of this training example via the cross entropy for each $j \in [n]$:

$$-\sum_{j \in [n]} \log \hat{p}_j(a^*) \quad (4)$$

Intuition. Elements of the inputs sequence X represent transitions in the environment. When passed through the GPT-2 transformer, the model learns to associate elements of the sequence via the standard query-key-value mechanism of the attention model. The query state s_{query} is demarcated by its zeros vector (which also acts as padding). Unlike other examples of transformers used for decision-making such as the Decision Transformer (Chen et al., 2021) and Algorithm Distillation (Laskin et al., 2022), DPT does not separate the individual (s, a, s', r) into their own embeddings to be made into one long sequence. This is because we view the transition tuples in the dataset as their own singletons, to be related with other singletons in the dataset through the attention mechanism. We note that there are various other implementation variations one could take, but we found success and robustness with this one.

A.2. Implementation Details

A.2.1. BANDIT ALGORITHMS

First, we describe the comparisons from the bandit experiments with hyperparameters.

Empirical Mean (Emp). Emp has no hyperparameters, but we give it some mechanism to avoid degenerate scenarios. In the offline setting, Emp will only choose from actions that have at least one example in the dataset. This gives Emp and LCB-style effect when actions are missing. Similarly, online, Emp will sample each action at least once before defaulting to its real strategy. These changes only improve Emp.

Upper Confidence Bound (UCB). According to the Hoeffding bound, we choose actions as $\hat{a} \in \operatorname{argmax}_{a \in \mathcal{A}} \left\{ \hat{\mu}_a + \sqrt{1/n_a} \right\}$ where $\hat{\mu}_a$ is the empirical mean so far for action a and n_a is the number of times a has been chosen so far. To arrive at this constant for the bonus, we coarsely tried a set of plausible values given the noise and found this to perform the best.

Lower Confidence Bound (LCB). We choose actions as $\hat{a} \in \operatorname{argmax}_{a \in \mathcal{A}} \left\{ \hat{\mu}_a - \sqrt{1/n_a} \right\}$ where $\hat{\mu}_a$ is the empirical mean so far for action a and n_a is the number of times a has been chosen so far.

Thompson Sampling (TS). Since the means are sampled uniformly from $[0, 1]$, Gaussian TS is partially misspecified; however, we set prior mean and variance to $\frac{1}{2}$ and $\frac{1}{12}$ to match the true ones. The noise model was well-specified with the correct variance. In the linear experiments of Figure 3a and Figure 3b, we set the prior mean and variance to 0 and 1 to fit the true ones better.

LinUCB. We choose $\hat{a}_t \in \operatorname{argmax}_{a \in \mathcal{A}} \langle \hat{\theta}_t, \phi(a) \rangle + \beta \|\phi(a)\|_{\hat{\Sigma}_t^{-1}}$ where $\beta = 1$ and $\hat{\Sigma}_t = I + \sum_{s \in [t-1]} \phi(a_s) \phi(a_s)^\top$ and $\hat{\theta}_t = \hat{\Sigma}_t^{-1} \sum_{s \in [t-1]} r_s \phi(a_s)$. Here, r_s and a_s are the reward and action observed at time s .

LinReg. LinReg (offline) is the same as LinUCB except we set $\beta = 0$ to greedily choose actions.

DPT. The transformer for DPT has an embedding size of 32, context length of 500 for basic bandits and 200 for linear bandits, 4 hidden layers, and 4 attention heads per attention layer for all bandits. We use the AdamW optimizer with weight decay $1e-4$, learning rate $1e-4$, and batch-size 64. For all experiments, we shuffle the in-context dataset D since order does not matter except in the linear bandit.

A.2.2. RL ALGORITHMS

Below, we describe the comparisons from the MDP experiments and their hyperparameters.

Proximal Policy Optimization (PPO). The reported results for PPO use the Stable Baselines3 implementation (Raffin et al., 2021) with the default hyperparameters, which successfully learns each task given 100K environment steps in Dark Room and 125K environment steps in Miniworld. In Dark Room, the policy is implemented as a multi-layer perceptron with two hidden layers of 64 units each. In Miniworld, the policy is a convolutional neural network with two convolutional layers with $16\ 3 \times 3$ kernels each, followed by a linear layer with output dimension of 8.

Algorithm Distillation (AD). We first collect learning histories with PPO for each of the training tasks. Then, given a cross-episodic context of length H , where H is the task horizon, the model is trained to predict the actions taken K episodes later (given the states visited in that episode). This was shown to lead to faster algorithms in (Laskin et al., 2022). We evaluated AD across different values of K . Between $K = 10, 50, 100$, we found $K = 100$ to be most performant in the Dark Room environment. In Miniworld, we also subsampled with $K = 100$. In Dark Room, the transformer has similar hyperparameters as DPT: an embedding size of 32, context length of 100 steps, 4 hidden layers, and 4 attention heads per attention layer. In Miniworld, as with DPT, we first encode the image with a convolutional network with two convolutional layers with $16\ 3 \times 3$ kernels each, followed by a linear layer with output dimension of 8.

RL². The reported results for RL² use an open-sourced implementation from (Zintgraf et al., 2020). The implementation uses PPO as the RL algorithm and defines a single trial as four consecutive episodes. The policy is implemented with one hidden layer of 32 units in Dark Room. In Miniworld, the policy is parameterized with a convolutional neural network with two convolutional layers with $16\ 3 \times 3$ kernels each, followed by a linear layer with output dimension of 8.

DPT. The transformer for DPT has an embedding size of 32, context length of 100 steps, 4 hidden layers, and 4 attention heads per attention layer in Dark Room. In Miniworld, the image is first passed through a convolutional network with two convolutional layers $16\ 3 \times 3$ kernels each, followed by a linear layer with output dimension of 8. The transformer model that processes these image embeddings otherwise has the same hyperparameters as in Dark Room. We use the AdamW optimizer with weight decay $1e-4$, learning rate $1e-3$, and batch-size 128.

A.3. Bandit Pretraining and Testing

Basic Bandit. Offline, to generate the in-context datasets for pretraining, we used a Dirichlet distribution to sample action frequencies in order to generate datasets with diverse compositions (i.e. some more uniform, some that only choose a few actions, etc.): $p_1 \sim \text{Dir}(\mathbb{1})$ where $p_1 \in \Delta(\mathcal{A})$ and $\mathbb{1} \in \mathbb{R}^{|\mathcal{A}|}$. We also mixed this with a distribution that has all mass on one action: $\hat{a} \sim \text{Unif}(\mathcal{A})$ and $p_2(\hat{a}) = 1$ and $p_2(a) = 0$ for all $a \neq \hat{a}$. The final action distribution is $p = (1 - \omega)p_1 + \omega p_2$ where $\omega \sim \text{Unif}(0.1[10])$. We train on 100,000 pretraining samples for 300 epochs with an 80/20 train/validation split. In Figure 2a, $\mathcal{D}_{\text{test}}$ is generated in the same way.

Expert-Biased Bandit. To generate expert-biased datasets for pretraining, we compute the action frequencies to bias the dataset towards the optimal action. Let a^* be the optimal one. As before, we take $p_1 \sim \text{Dir}(\mathbb{1})$. Then, $p_2(a^*) = 1$ and $p_2(a) = 0$ for all $a \neq a^*$. For of bias of ω , we take $p = (1 - \omega)p_1 + \omega p_2$ with $\omega \sim \text{Unif}(0.1[10])$. We use the same pretraining sample size and epochs as before. For testing, $\mathcal{D}_{\text{test}}$ is generated the same way except we fix a particular $\omega \in \{0, 0.5, 1\}$ to test on.

Linear Bandit. We consider the case where $|\mathcal{A}| = 10$ and $d = 2$. To generate environments from \mathcal{T}_{pre} , we first sampled a fixed set of actions from $\mathcal{N}(\mathbf{0}, I_d/d)$ in \mathbb{R}^d to represent the features. Then, for each τ , we sampled $\theta_\tau \sim \mathcal{N}(\mathbf{0}, I_d/d)$ to produce the means $\mu_a = \langle \theta_\tau, \phi(a) \rangle$ for $a \in \mathcal{A}$. To generate the in-context dataset, we ran Gaussian TS (which does not leverage ϕ) over $n = 200$ steps (see hyperparameters in previous section). Because order matters, we did not shuffle and used 1,000,000 pretraining samples over 200 epochs with an 80/20 train/validation split. At test time, we set $\mathcal{T}_{\text{test}} = \mathcal{T}_{\text{pre}}$ and $\mathcal{D}_{\text{test}} = \mathcal{D}_{\text{pre}}$. Note that ϕ is fixed over all τ , as is standard for a linear bandit.

A.4. MDP Environment Details

Dark Room. The agent must navigate a 10×10 grid to find the goal within $H = 100$ steps. The agent’s observation is its xy -position, the allowed actions are left, right, up, down, and stay, and the reward is only $r = 1$ when the agent is at the goal, and $r = 0$ otherwise. At test time, the agent begins at the $(0, 0)$ position. We randomly designate 80 of the 100 grid squares to be goals for the training tasks, and hold out the remaining 20 for evaluation.

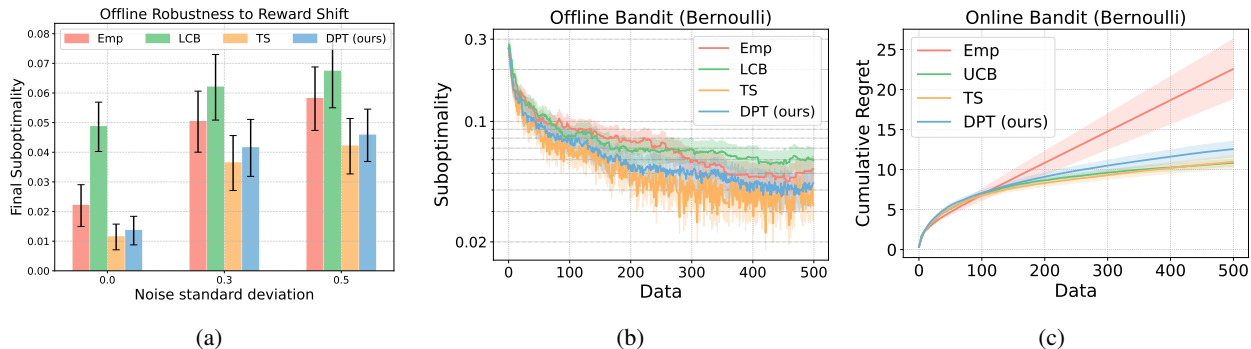


Figure 6. (a) Final (after 500 steps) offline suboptimality on out-of-distribution bandits with different Gaussian noise standard deviations. (b) Offline performance on out-of-distribution Bernoulli bandits, given random in-context datasets. (c) Online cumulative regret on Bernoulli bandits. The mean and standard error are computed over 200 test tasks.

Miniworld. The agent must navigate to the correct box, which is initially unknown, from 25×25 RGB image observations. The agent is additionally conditioned on its own direction vector. In each episode, the environment is initialized with four boxes of different colors, one in each corner of the square room. The agent can turn left, turn right, or move forward. The reward is only $r = 1$ when the agent is near the correct box and $r = 0$ otherwise, and each episode is 50 time-steps long. At test time, the agent begins in the middle of the room.

A.5. MDP Pretraining Datasets

Dark Room. In Dark Room, we collect 100K in-context datasets, each of length $H = 100$ steps, with a uniform-random policy. The 100K datasets are evenly collected across the 100 goals. The query states are uniformly sampled from the state space, and the optimal actions are computed as follows: move up/down until the agent is on the same y -position as the goal, then move left/right until the agent is on the x -position as the goal. Of the 100K collections of datasets, query states, and optimal actions, we use the first 80K (corresponding to the first 80 goals) for training and the remaining 20K for validation.

Miniworld. While this task is solved from image-based observations, we also note that there are only four distinct tasks (one for each colored box), and the agent does not need to handle new tasks at test time. Hence, the number of in-context datasets required in pretraining is fewer – we use 40K datasets each of length $H = 50$ steps. So as to reduce computation, the in-context datasets only have only (s, a, r) tuples. The query states, which consist of image and direction are sampled uniformly from the entire state space, i.e., the agent is placed uniformly at random in the environment, pointing in a random direction. The optimal actions are computed as follows: turn towards the correct box if the agent is not yet facing it (within ± 15 degrees), otherwise move forward. Of the 40K collections of datasets, query states, and optimal actions, we use 32K for training and the remaining 8K for validation.

B. Additional Experimental Results

B.1. Bandits

This section reports additional experimental results in bandit environments.

Out-of-distribution reward variances. In Figures 2c and 6a, we demonstrate the robustness of the basic pretrained model under shifts in the reward distribution at test time by varying the amount of noise observed in the rewards. DPT maintains robustness to these shifts similar to TS.

Bernoulli rewards. We test the out-of-distribution ability of DPT further by completely changing the reward distribution from Gaussian to Bernoulli bandits. Despite being trained only on Gaussian tasks during pretraining, DPT maintains strong performance both offline and online in Figures 6b and 6c.

B.2. Markov Decision Processes

Supervised Pretraining Can Learn In-Context Reinforcement Learning

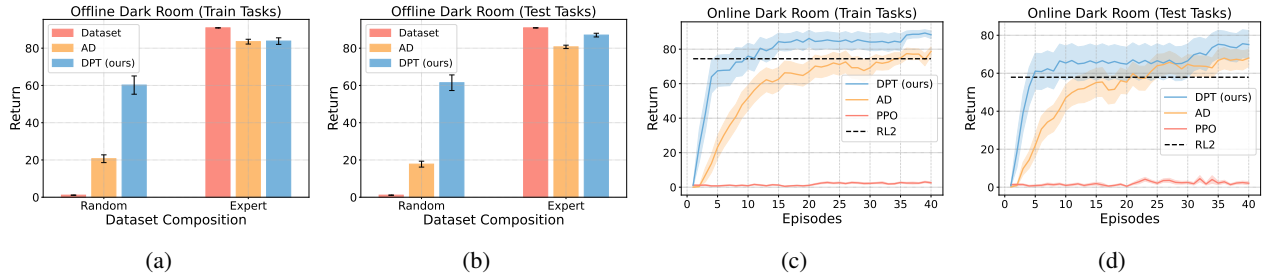


Figure 8. All comparisons in Dark Room evaluated on the tasks that were seen during pretraining, displayed next to their evaluations on test task counterparts from the main text.

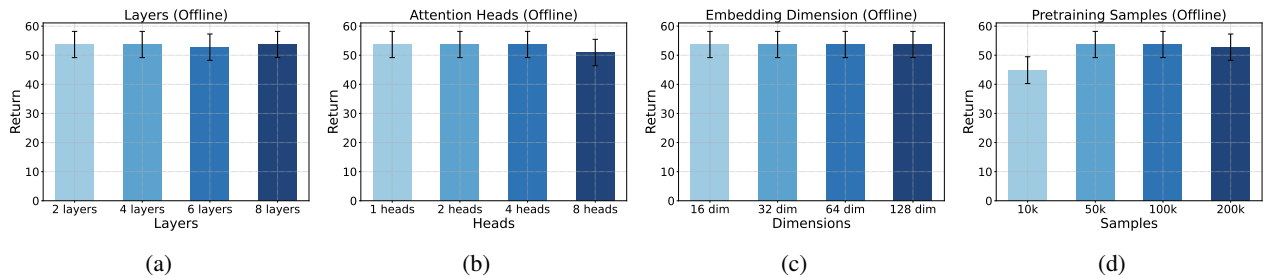


Figure 9. Sensitivity analysis of the offline Dark Room task over the GPT-2 transformer’s hyperparameters: (a) layers (b) attention heads (c) embedding dimensions (d) pretraining samples.

This section reports additional experimental results in the Dark Room and Miniworld environments.

Performance on training tasks. In Fig. 8, we show the performance of each method on the training tasks in Dark Room. Offline, DPT and AD demonstrate comparable performance as on the training tasks, indicating a minimal generalization gap to new goals. Online, DPT, AD, and RL² also achieve performance on the training tasks similar to that on the test tasks.

Generalization to new dynamics. In this experiment, we study generalization to variations in a different aspect of the MDP, namely the dynamics. We design *Dark Room (Permuted)*, a variant of Dark Room in which the goal is fixed to a corner but the action space is randomly permuted. Hence, the agent must leverage its historical context to infer the effect of each action. On a held-out set of 20 permutations, DPT infers the optimal policy correctly every time offline, given only 100 offline samples, matching the optimal policy at 83 return. Similarly, the online performance immediately snaps to a near optimal policy in one episode once it identifies the novel permutation in Figure 7.

B.3. Sensitivity Analysis

We next seek to understand the sensitivity of DPT to different hyperparameters, including the model size and size of the pretraining dataset. These experiments are performed in the Dark Room environment. As shown in Fig. 9, the performance of DPT is robust to the model size; it is the same across different embedding sizes, number of layers, and number of attention heads. Notably, the performance is slightly worse with 8 attention heads, which may be attributed to slight overfitting. We do see that when the pretraining dataset is reduced to 10% of its original size (10000 samples) the performance degrades, but otherwise has similar performance with larger pretraining datasets.

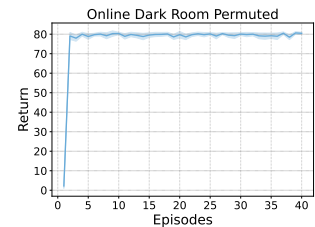


Figure 7. Online evaluation of DPT on Dark Room when tested on novel actions set permutations.

C. Additional Theory and Omitted Proofs

We start with a well-known concentration inequality for the maximum-likelihood estimate (MLE) to provide some more justification for the approximation made in Assumption 6.1. We state a version from (Agarwal et al., 2020). Let \mathcal{F} be a finite function class used to model a conditional distribution $p_{Y|X}(y|x)$ for $x \in \mathcal{X}$ and $y \in \mathcal{Y}$. Assume there is $f^* \in \mathcal{F}$ such that $p(y|x) = f^*(y|x)$ (realizable), and $f(\cdot|x) \in \Delta(\mathcal{Y})$ for all $x \in \mathcal{X}$ and $f \in \mathcal{F}$ (proper). Let $D = \{x_i, y_i\}_{i \in [N]}$ denote a dataset of i.i.d samples where $x_i \sim p_X$ and $y_i \sim p_{Y|X}(\cdot|x_i)$. Let

$$\hat{f} = \operatorname{argmax}_{f \in \mathcal{F}} \sum_{i \in [N]} \log f(y_i|x_i) \quad (5)$$

Proposition C.1 (Theorem 21 of (Agarwal et al., 2020)). *Let D and \hat{f} be given as above under the aforementioned conditions. Then, with probability at least $1 - \delta$,*

$$\mathbb{E}_{x \sim p_X} \|\hat{f}(\cdot|x) - p_{Y|X}(\cdot|x)\|_1^2 \leq \frac{8 \log(|\mathcal{F}|/\delta)}{N} \quad (6)$$

The finiteness of \mathcal{F} is done for simplicity, but we can see that this yields dependence on the log-cardinality, a common measure of complexity. Extensions to infinite \mathcal{F} of bounded statistical complexity can be readily made to replace this. For our setting, the bound suggests that $\mathbb{E}_{P_{pre}} \|P_{pre}(\cdot|s_{query}, D, \xi_h) - M_\theta(\cdot|s_{query}, D, \xi_h)\|_1^2 \rightarrow 0$ as $N \rightarrow \infty$ with high probability, provided the function class of M_θ has bounded statistical complexity.

C.1. Posterior Sampling

Posterior sampling is most generally described with the following procedure (Osband et al., 2013). Initialize a prior distribution $\mathcal{T}_1 = \mathcal{T}_{pre}$ and dataset $D = \{\}$. For $k \in [K]$

1. Sample $\tau_k \sim \mathcal{T}_k$ and compute $\hat{\pi}_{\tau_k}$
2. Execute $\pi_{\tau_k}^*$ and add interactions to D
3. Update posterior distribution $\mathcal{T}_{k+1}(\tau) = P(\tau|D)$.

The prior and posteriors are typically over models such as reward functions in bandits or transition dynamics in MDPs.

C.2. Proof of Theorem 6.3

Theorem 6.3 (DPT \iff PS). *Let the above assumptions hold. Then, $P_{ps}(\xi_H | D, \tau_c) = P_{M_\theta}(\xi_H | D, \tau_c)$ for all trajectories ξ_H .*

Proof. Without loss of generality, for a task τ , we take $\pi_\tau^*(\cdot|s)$ to be deterministic and denote the optimal action in state s as $\pi_\tau^*(s)$. Recall that we consider a fixed current task τ_c and a fixed in-context dataset D . Define $\xi_h = (s_1, a_1, \dots, s_h, a_h)$.

We now formally state the variant of the full joint distribution from which we sample during pretraining. Let τ and D' be an arbitrary task and dataset and let $a^* \in \mathcal{A}$, $s_{query} \in \mathcal{S}$, $\xi_{H-1} \in (\mathcal{S} \times \mathcal{A})^{H-1}$, and $h \in [0, H-1]$ be arbitrary.

$$P_{pre}(\tau, a^*, s_{query}, D', \xi_{H-1}, h) = \mathcal{T}_{pre}(\tau) \mathcal{D}_{pre}(D'; \tau) \mathcal{D}_{query}(s_{query}) \mathfrak{S}_H(s_{1:H}) \pi_\tau^*(a^*|s_{query}) \quad (7)$$

$$\times \operatorname{Unif}[0, H-1] \prod_{i \in [H]} \pi_\tau^*(a_i|s_i) \quad (8)$$

The $\operatorname{Unif}[0, H-1]$ is due to the fact that we sample $h \sim \operatorname{Unif}[0, H-1]$ and then truncate ξ_h from ξ_{H-1} (or, equivalently, sample $\xi_h \sim \mathfrak{S}_h$ directly), marginalizing out the other variables. For $h' \leq h-1$, recall that we also use the notation $\mathfrak{S}_{h'}(s_{1:h'})$ to denote the marginalization of the full joint \mathfrak{S}_H . We will eventually work with the posterior of this distribution given the data D and history ξ_h :

$$P_{pre}(\tau|D, \xi_h) \propto \mathcal{T}_{pre}(\tau) \mathcal{D}_{pre}(D; \tau) \prod_{i \in [h]} \pi_\tau^*(a_i|s_i) \quad (9)$$

$$\propto P_{pre}(\tau|D) \prod_{i \in [h]} \pi_\tau^*(a_i|s_i) \quad (10)$$

We define the following random sequences and subsequences:

$$\Xi_{ps}(h; D) = (S_1^{ps}, A_1^{ps}, \dots, S_h^{ps}, A_h^{ps}) \quad (11)$$

where the variables are generated according to the following conditional process: $\tau_{ps} \sim P(\cdot|D)$, $S_1^{ps} \sim \rho_{\tau_c}$, $A_h^{ps} \sim \pi_{\tau_{ps}}^*(\cdot|S_h^{ps})$, and $S_{h+1}^{ps} \sim T_{\tau_c}(\cdot|S_h^{ps}, A_h^{ps})$. We also define $\Xi_{ps}(h' : h; D)$ to be the last $h - h'$ elements of $\Xi_{ps}(h; D)$. Analogously, we define

$$\Xi_{pre}(h; D) = (S_1^{pre}, A_1^{pre}, \dots, S_h^{pre}, A_h^{pre}) \quad (12)$$

where the variables are from the process: $S_1^{pre} \sim \rho_{\tau_c}$, $A_h^{pre} \sim P_{pre}(\cdot|S_h^{pre}, D, \Xi_{pre}(h-1; D))$, and $S_{h+1}^{pre} \sim T_{\tau_c}(\cdot|S_h^{pre}, A_h^{pre})$. Note that A_h^{pre} is sampled conditioned on the sequence $\Xi_{pre}(h; D)$ so far.

We will show that $\Xi_{ps}(h; D)$ and $\Xi_{pre}(h; D)$ follow the same distribution for all $h \in [H]$. For convenience, we will drop notational dependence on D , except where it resolves ambiguity. Also, because of Assumption 6.1, we have that $P_{pre}(\cdot|S_h^{pre}, D, \Xi_{pre}(h-1)) = M_\theta(\cdot|S_h^{pre}, D, \Xi_{pre}(h-1))$, so we will just work with P_{pre} for the remainder of the proof. We will also make use of the following lemma.

Lemma C.2. *If \mathcal{D}_{pre} is compliant, then $P_{pre}(\tau|D) = P(\tau_{ps} = \tau|D)$.*

Proof. From the definition of posterior sampling (using the same prior, \mathcal{T}_{pre}), we have that

$$P(\tau_{ps} = \tau|D) \propto P(D|\tau)\mathcal{T}_{pre}(\tau) \quad (13)$$

$$\propto \mathcal{T}_{pre}(\tau) \prod_{j \in [n]} T_\tau(s'_j|s_j, a_j)R_\tau(r_j|s_j, a_j) \quad (14)$$

$$\propto \mathcal{T}_{pre}(\tau) \prod_{j \in [n]} T_\tau(s'_j|s_j, a_j)R_\tau(r_j|s_j, a_j)\mathcal{D}_{pre}(a_j|s_j, D_{j-1}) \quad (15)$$

$$= \mathcal{T}_{pre}(\tau)\mathcal{D}_{pre}(D; \tau) \quad (16)$$

$$= P_{pre}(\tau|D) \quad (17)$$

where the second line crucially uses the fact that posterior sampling chooses actions based only on the prior and history so far. Similarly, the third line uses the fact that \mathcal{D}_{pre} is compliant. Since the two sides are proportional in τ , they are equivalent. \square

We will prove Theorem 6.3 via induction for each $h \in [H]$. First, consider the base case for a sequence of length $h = 1$. Recall that ρ_{τ_c} denotes the initial state distribution of τ_c . We have that the densities can be written as

$$P(\Xi_{ps}(1) = \xi_1) = P(S_1^{ps} = s_1, A_1^{ps} = a_1) \quad (18)$$

$$= \rho_{\tau_c}(s_1)P(A_1^{ps} = a_1|S_1^{ps} = s_1) \quad (19)$$

$$= \rho_{\tau_c}(s_1) \int_{\tau} P(A_1^{ps} = a_1, \tau_{ps} = \tau|S_1^{ps} = s_1)d\tau \quad (20)$$

$$= \rho_{\tau_c}(s_1) \int_{\tau} \pi_\tau^*(a_1|s_1)P_{ps}(\tau_{ps} = \tau|D, S_1^{ps} = s_1)d\tau \quad (21)$$

$$= \rho_{\tau_c}(s_1) \int_{\tau} \pi_\tau^*(a_1|s_1)P_{ps}(\tau_{ps} = \tau|D)d\tau \quad (22)$$

$$= \rho_{\tau_c}(s_1)P_{pre}(A_1^{pre} = a_1|s_1, D) \quad (23)$$

$$= P(\Xi_{pre}(1) = \xi_1) \quad (24)$$

where the second line uses the sampling process of S_1^{pre} ; the third marginalizes over τ_{ps} , which is the task that posterior sampling samples to find the optimal policy; the fourth decomposes this into the optimal policy and the posterior over τ_{ps} given D and S_1^{ps} . Since S_1^{ps} is independent of sampling of τ_{ps} this dependence goes away in the next line. The sixth line applies Lemma C.2 and then, for $h = 1$, there is no history to condition on.

Now, we leverage the inductive hypothesis to prove the full statement. Suppose that the hypothesis holds for $h - 1$. Then,

$$P(\Xi_{ps}(h) = \xi_h) = P(\Xi_{ps}(h-1) = \xi_{h-1})P(S_h^{ps} = s_h, A_h^{ps} = a_h | \Xi_{ps}(h-1) = \xi_{h-1}) \quad (25)$$

$$(26)$$

By the hypothesis, we have that $P(\Xi_{ps}(h-1) = \xi_{h-1}) = P(\Xi_{pre}(h-1) = \xi_{h-1})$. For the second factor,

$$P(S_h^{ps} = s_h, A_h^{ps} = a_h | \Xi_{ps}(h-1) = \xi_{h-1}) \quad (27)$$

$$= T_{\tau_c}(s_h | s_{h-1}, a_{h-1}) \cdot P(A_h^{ps} = a_h | S_h^{ps} = s_h, \Xi_{ps}(h-1) = \xi_{h-1}) \quad (28)$$

$$= T_{\tau_c}(s_h | s_{h-1}, a_{h-1}) \cdot \int_{\tau} P(A_h^{ps} = a_h, \tau_{ps} = \tau | S_h^{ps} = s_h, \Xi_{ps}(h-1) = \xi_{h-1}) d\tau \quad (29)$$

As before, we can further rewrite the last factor as

$$P(A_h^{ps} = a_h, \tau_{ps} = \tau | S_h^{ps} = s_h, \Xi_{ps}(h-1) = \xi_{h-1}) \quad (30)$$

$$= \pi_{\tau}^*(a_h | s_h) \cdot P(\tau_{ps} = \tau | S_h^{ps} = s_h, \Xi_{ps}(h-1) = \xi_{h-1}) \quad (31)$$

where

$$P(\tau_{ps} = \tau | S_h^{ps} = s_h, \Xi_{ps}(h-1) = \xi_{h-1}) = \frac{P(S_h^{ps} = s_h, \Xi_{ps}(h-1) = \xi_{h-1} | \tau_{ps} = \tau) P(\tau_{ps} = \tau | D)}{P(S_h^{ps} = s_h, \Xi_{ps}(h-1) = \xi_{h-1})} \quad (32)$$

$$\propto P_{pre}(\tau | D) \prod_{i \in [h-1]} T_{\tau_c}(s_{i+1} | s_i, a_i) \pi_{\tau}^*(a_i | s_i) \quad (33)$$

$$\propto P_{pre}(\tau | D) \prod_{i \in [h-1]} \pi_{\tau}^*(a_i | s_i) \quad (34)$$

$$\propto P_{pre}(\tau | D) \mathcal{D}_{\text{query}}(s_h) \mathfrak{S}_{h-1}(s_{1:h-1}) \prod_{i \in [h-1]} \pi_{\tau}^*(a_i | s_i) \quad (35)$$

$$\propto P_{pre}(\tau | s_h, D, \xi_{h-1}) \quad (36)$$

$$(37)$$

where \propto denotes that the two sides are equal up to multiplicative factors independent of τ . In the first line, we used Bayes rule. In the second line, given that $\tau_{ps} = \tau$ (i.e. posterior sampling selected τ to deploy), we decompose the probability of observing that sequence of states of actions. We also used Lemma C.2. The denominator does not depend on τ . Similarly, for the third and fourth lines, T_{τ_c} and \mathfrak{S} do not depend on τ . The final line follows from the definition of the joint pretraining distribution in this regime.

Therefore, we conclude that the posterior over the value of τ_{ps} is the same as the posterior over the task in the pretraining distribution, given s_h, D, ξ_{h-1} . Substituting back through all the previous equations, we have

$$P(\Xi_{ps}(h) = \xi_h) \quad (38)$$

$$= P(\Xi_{pre}(h-1) = \xi_{h-1}) \cdot T_{\tau_c}(s_h | s_{h-1}, a_{h-1}) \int_{\tau} \pi_{\tau}^*(a_h | s_h) P_{pre}(\tau | s_h, D, \xi_{h-1}) d\tau \quad (39)$$

$$= P(\Xi_{pre}(h-1) = \xi_{h-1}) \cdot T_{\tau_c}(s_h | s_{h-1}, a_{h-1}) P_{pre}(a_h | s_h, D, \xi_{h-1}) \quad (40)$$

$$= P(\Xi_{pre}(h) = \xi_h) \quad (41)$$

This concludes the proof. \square

C.3. Proof of Corollary 6.4

Corollary C.3 (Finite MDPs). *Suppose that $\sup_{\tau} \mathcal{T}_{\text{test}}(\tau) / \mathcal{T}_{\text{pre}}(\tau) \leq \mathcal{C}$ for some $\mathcal{C} > 0$. For the above MDP setting, the pretrained model M_{θ} satisfies $\mathbb{E}_{\mathcal{T}_{\text{test}}} [\text{Reg}_{\tau}(M_{\theta})] \leq \tilde{\mathcal{O}}(CH^{3/2}S\sqrt{AK})$.*

Proof. Note that \mathcal{D}_{pre} is clearly compliant since it is generated by random sampling. We use the equivalence between M_θ and posterior sampling established in Theorem 6.3. The proof then follows immediately from Theorem 1 of (Osband et al., 2013) to guarantee that

$$\mathbb{E}_{\mathcal{T}_{\text{pre}}} [\text{Reg}_\tau(M_\theta)] \leq \tilde{\mathcal{O}} \left(H^{3/2} S \sqrt{AK} \right) \quad (42)$$

where the notation $\tilde{\mathcal{O}}$ omits polylogarithmic dependence. The bound on the test task distribution follows from the assumed bound on the likelihood ratio under the priors:

$$\int \mathcal{T}_{\text{test}}(\tau) \text{Reg}_\tau(M_\theta) d\tau \leq \mathcal{C} \int \mathcal{T}_{\text{pre}}(\tau) \text{Reg}_\tau(M_\theta) d\tau. \quad (43)$$

□

C.4. Proof of Corollary 6.5

Corollary C.4 (Latent representation learning in linear bandits). *For $\mathcal{T}_{\text{test}} = \mathcal{T}_{\text{pre}}$ in the above linear bandit setting, M_θ satisfies $\mathbb{E}_{\mathcal{T}_{\text{test}}} [\text{Reg}_\tau(M_\theta)] \leq \tilde{\mathcal{O}}(d\sqrt{K})$.*

Proof. The distribution \mathcal{D}_{pre} satisfies compliance by definition because it is generated by an adaptive algorithm TS. The proof once again follows by immediately deferring to the established result of (Russo & Van Roy, 2014) (Proposition 3) for linear bandits by the posterior sampling equivalence of Theorem 6.3. This ensures that posterior sampling achieves regret $\tilde{\mathcal{O}}(d\sqrt{K})$. It remains, however, to justify that $P_{\text{pre}}(\cdot|D_k)$ will be covered by Gaussian Thompson Sampling for all D_k with $k \in [K]$. This is verified by noting that $P_{\text{ps}}(a|D_k) > 0$ for non-degenerate Gaussian Thompson Sampling (positive variances of the prior and likelihood functions) and finite K . This guarantees that any D_k will have support. □

C.5. Proof of Proposition 6.6

Proposition C.5. *Let P_{pre}^1 and P_{pre}^2 be pretraining distributions that differ only by their in-context dataset distributions, denoted by $\mathcal{D}_{\text{pre}}^1$ and $\mathcal{D}_{\text{pre}}^2$. If $\mathcal{D}_{\text{pre}}^1$ and $\mathcal{D}_{\text{pre}}^2$ are compliant with the same support, then $P_{\text{pre}}^1(a^*|s_{\text{query}}, D, \xi_h) = P_{\text{pre}}^2(a^*|s_{\text{query}}, D, \xi_h)$ for all $a^*, s_{\text{query}}, D, \xi_h$.*

Proof. The proof follows by direct inspection of the pretraining distributions. For P_{pre}^1 , we have

$$P_{\text{pre}}^1(a^*|s_{\text{query}}, D, \xi) = \int_{\tau} \pi_{\tau}^*(a^*|s_{\text{query}}) P_{\text{pre}}^1(\tau|s_{\text{query}}, D, \xi) d\tau \quad (44)$$

The posterior distribution over tasks is simply

$$P_{\text{pre}}^1(\tau|s_{\text{query}}, D, \xi) = \frac{P_{\text{pre}}^1(\tau, s_{\text{query}}, D, \xi)}{P_{\text{pre}}^1(s_{\text{query}}, D, \xi)} \quad (45)$$

$$\propto P_{\text{pre}}^1(\tau) P_{\text{pre}}^1(\xi|\tau) \mathcal{D}_{\text{query}}(s_{\text{query}}) \mathcal{D}_{\text{pre}}^1(D; \tau) \quad (46)$$

$$= P_{\text{pre}}^2(\tau) P_{\text{pre}}^2(\xi|\tau) \mathcal{D}_{\text{query}}(s_{\text{query}}) \mathcal{D}_{\text{pre}}^1(D; \tau) \quad (47)$$

Then, the distribution over the in-context dataset can be decomposed as

$$\mathcal{D}_{\text{pre}}^1(D; \tau) = \prod_{i \in [n]} R_\tau(r_i|s_i, a_i) T_\tau(s'_i|s_i, a_i) \mathcal{D}_{\text{pre}}^1(a_i|s_i, D_{i-1}; \tau) \quad (48)$$

$$= \prod_{i \in [n]} R_\tau(r_i|s_i, a_i) T_\tau(s'_i|s_i, a_i) \mathcal{D}_{\text{pre}}^1(a_i|s_i, D_{i-1}) \quad (49)$$

$$\propto \prod_{i \in [n]} R_\tau(r_i|s_i, a_i) T_\tau(s'_i|s_i, a_i) \mathcal{D}_{\text{pre}}^2(a_i|s_i, D_{i-1}) \quad (50)$$

$$= \mathcal{D}_{\text{pre}}^2(D; \tau), \quad (51)$$

where the second equality holds because $\mathcal{D}_{\text{pre}}^1(a_j|s_j, D_j; \tau)$ is assumed to be invariant to τ by compliance, and the fifth equality holds because $\mathcal{D}_{\text{pre}}^2(a_j|s_j, D_j; \tau)$ is assumed to be invariant to τ .

Therefore, we conclude that, for any s, D, ξ ,

$$P_{\text{pre}}^1(\tau|s, D, \xi) \propto P_{\text{pre}}^2(\tau)P_{\text{pre}}^2(\xi|\tau)\mathcal{D}_{\text{query}}(s)\mathcal{D}_{\text{pre}}^2(D; \tau) \quad (52)$$

$$\propto P_{\text{pre}}^2(\tau|s, D, \xi). \quad (53)$$

Since also $\int_{\tau} P_{\text{pre}}^1(\tau|s, D, \xi) = 1 = \int_{\tau} P_{\text{pre}}^2(\tau|s, D, \xi)$, then

$$P_{\text{pre}}^1(\tau|s, D, \xi) = P_{\text{pre}}^2(\tau|s, D, \xi). \quad (54)$$

Substituting this back into Equation 44 yields $P_{\text{pre}}^1(a^*|s, D, \xi) = P_{\text{pre}}^1(a^*|s, D, \xi)$. □