

# DIFFUSION POSTERIOR SAMPLING FOR NONLINEAR CONTEXTUAL BANDITS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We study multi-task nonlinear contextual bandits, where different tasks share the same reward structure but are characterized by distinct model parameters drawn from a common unknown prior distribution. The goal is to leverage information from past tasks to minimize regret on a new task with limited online interactions. Thompson Sampling (TS) is a popular approach for solving contextual bandits, maintaining a posterior over the model parameter that is updated each round using a hand-specified conjugate prior (e.g., Gaussian) and the observed rewards. However, such priors cannot capture the rich cross-task structure in multi-task settings, leading to misspecified posteriors and suboptimal exploration. To address this, we train a diffusion model on data from past tasks to learn a flexible prior distribution over task parameters. In a new bandit task, parameters are estimated via a conditional reverse-diffusion process, where each step combines: (i) an unconditional drift from the diffusion prior, (ii) a likelihood-driven drift from the interaction history, and (iii) a noise term enabling randomized exploration. We instantiate this framework in two ways. **DLTS** integrates history into the diffusion prior at every reverse step to form a conditional posterior, from which approximate samples are drawn. **DPSG** first performs unconditional reverse sampling from the pretrained diffusion prior and then applies a single history-guided gradient correction. Both methods adhere to the same framework but differ in how they incorporate interaction history from the new task: DLTS explicitly constructs the conditional posterior, while DPSG provides a lightweight approximation by coupling unconditional sampling with one corrective step. In theory, we formalize oracle TS (OTS) and its diffusion counterpart (ODTS) and prove they are equivalent when the diffusion prior matches the true prior. We bound the per-round expected regret gap between ODTS and OTS by the cumulative score estimation error across diffusion levels. Our empirical evaluation demonstrates that our proposed methods are competitive with specialized baselines in linear settings and outperform baselines benefiting from the diffusion prior in challenging nonlinear bandit environments.

## 1 INTRODUCTION

Sequential decision-making under uncertainty hinges on balancing exploration and exploitation (Lattimore & Szepesvári, 2020). One prominent approach to address this trade-off is Thompson Sampling (TS) (Thompson, 1933). In contextual bandits, TS maintains a posterior over the model parameter that is updated each round using a hand-specified conjugate prior (e.g., Gaussian) and the observed rewards (Agrawal & Goyal, 2013a). In multi-task settings where different tasks share the same reward structure but are characterized by distinct model parameters drawn from a common unknown prior distribution, such simple priors fail to capture rich cross-task structure, leading to misspecified posteriors and suboptimal exploration (Chapelle & Li, 2011). To address this, Hong et al. (2022) proposed to initialize TS with a mixture prior for multi-task bandits. Nevertheless, fixed parametric mixtures still struggle with complex structure—e.g., multimodal task families where the parameter distribution has many distinct peaks (modes) that can be uneven in size or far apart (Vucelja et al., 2019; Finn et al., 2018), and heavy-tailed families with extreme outliers (Bubeck et al., 2013; Forbes & Wraith, 2014). Such patterns are hard to capture with a small, fixed set of components, leading to underfit posteriors. This motivates learning a flexible prior that fully leverages past tasks to enable fast adaptation on a new task with limited interaction.

Deep generative models have achieved remarkable success in producing high-quality synthetic data across modalities (Saharia et al., 2022; Rombach et al., 2022; Liu et al., 2023). These results highlight their ability to model complex, multi-modal distributions. In online decision-making, the policy must incorporate newly collected interaction data and update its strategy frequently. Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020) align well with this need: their reverse process is iterative and naturally supports conditioning on new observations at every round. In score-based diffusion (Song et al., 2021), this conditioning is implemented by augmenting the score with a likelihood term (Chung et al., 2023). Given this intuition, prior work (Hsieh et al., 2023) tackles multi-task multi-armed bandits by training a diffusion model to learn a flexible prior over task parameters and coupling it with Thompson Sampling, with practical steps such as variance calibration to obtain reliable uncertainty in the reverse process (Hsieh et al., 2023). In parallel, diffusion priors have been adapted to online posterior sampling for linear and generalized linear contextual bandits via closed-form updates at each reverse step (Kveton et al., 2024). However, these prior work focus on specialized updates for linear or generalized linear bandits, leaving multi-task nonlinear contextual bandits largely unexplored. The core difficulty is the absence of closed-form diffusion reverse updates for nonlinear reward models, which complicates algorithm design. Without closed forms, sampling error can be large; combined with limited interactions per task and potential inaccuracies in the diffusion learned prior, this can destabilize posterior sampling and misguide exploration. Consequently, a natural question arises:

*Can we design diffusion posterior sampling algorithms for general nonlinear contextual bandits?*

In this paper, we provide an affirmative answer to this question by first providing a diffusion-based posterior sampling framework for multi-task nonlinear contextual bandits. We learn a flexible prior over task parameters with a diffusion model trained on past tasks, and on a new task perform conditional reverse sampling that blends prior drift, likelihood guidance, and stochastic exploration at each step. We instantiate this with two methods: DLTS, which conditions at every reverse step and samples via Langevin Monte Carlo (LMC), and DPSG, which draws an unconditional sample and applies a single likelihood-guided correction. To stabilize learning, we also propose DPSG-MP, a practical DPSG variant that replaces the single correction with a short inner loop of likelihood-gradient updates with projection.

**Our contributions** are summarized as follows:

- **Unified Framework.** We present a unified diffusion-based posterior sampling framework for multi-task nonlinear contextual bandits. In our framework, the algorithm estimates the reward model parameters via a conditional reverse-diffusion process that, at each step, consists of (i) an unconditional drift from the diffusion prior, (ii) a likelihood-driven drift from the interaction history, and (iii) a noise term enabling randomized exploration.
- **Algorithm Design.** We instantiate the unified framework in two ways: (i) Diffusion Langevin Thompson Sampling (DLTS), which explicitly constructs the conditional posterior and then draws approximate samples; and (ii) Diffusion Posterior Sampling with Guidance (DPSG), which couples unconditional sampling with a single history-guided correction step for simplicity and speed. We also propose DPSG with Multi-step Projection (DPSG-MP), a practical DPSG variant to stabilize learning, which improves empirical performance.
- **Theoretical analysis.** We formalize oracle TS (OTS), oracle diffusion TS (ODTS) and prove they are equivalent when the diffusion prior exactly matches the true prior. We also bound the per-round expected regret gap between ODTS and OTS by the cumulative score-estimation error across diffusion levels.
- **Extensive experiments.** We evaluate our algorithms in simulations of both linear and non-linear contextual bandits. In the linear setting, our algorithms achieve performance comparable to baseline algorithms that use closed-form updates, such as LinTS and DiffTS. Furthermore, our methods demonstrate strong performance in the non-linear bandit setting, where the learned diffusion prior provides a significant advantage.

## 2 PRELIMINARIES

**Nonlinear Contextual Bandits** Contextual bandits form a broad class of sequential decision problems where the player chooses based on an observed action set represented by feature vectors (contexts). At each round  $t$ , the player observes arm set  $\mathcal{X}_t \subseteq \mathbb{R}^d$ , selects an arm  $\mathbf{x}_t \in \mathcal{X}_t$  and receives

reward  $y_t$  from the environment. Assume that the mean reward for a feature  $\mathbf{x} \in \mathbb{R}^d$  is generated by an underlying function  $f(\boldsymbol{\theta}^*; \mathbf{x})$ , and the observed reward satisfies  $y(\mathbf{x}) = f(\boldsymbol{\theta}^*; \mathbf{x}) + \varepsilon$ , where  $\boldsymbol{\theta}^* \in \mathbb{R}^d$  is an unknown parameter shared by all arms and  $\varepsilon$  is observation noise. In this work, we consider contextual bandits with nonlinear reward model  $f(\boldsymbol{\theta}^*; \mathbf{x})$ . The goal of a bandit algorithm is to maximize cumulative reward over a horizon  $T$ , equivalently to minimize the pseudo-regret (Lattimore & Szepesvári, 2020):  $R(T) = \sum_{t=1}^T (f(\boldsymbol{\theta}^*; \mathbf{x}_t^*) - f(\boldsymbol{\theta}^*; \mathbf{x}_t))$ , where  $\mathbf{x}_t^* = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}_t} f(\boldsymbol{\theta}^*; \mathbf{x})$  is the arm with the highest expected reward at round  $t$ .

**Multi-task Bandit Problem** We consider the multi-task nonlinear contextual bandit problem. All the bandit tasks share the same reward model structure  $y(\mathbf{x}) = f(\boldsymbol{\theta}^*; \mathbf{x}) + \varepsilon$  and the underlying parameter  $\boldsymbol{\theta}^*$  varies across tasks but is drawn independently from a common prior distribution  $p_0(\boldsymbol{\theta}^*)$ . In our problem setup, we have model parameter  $\{\boldsymbol{\theta}^{(i)}\}_{i=1}^N$  from previous  $N$  bandit tasks, which can be ground-truth (perfect data) or estimation (imperfect data). Our goal is to utilize these data to help solve for a new bandit task ( $N + 1$ ).

**Diffusion Models** Score-based diffusion models define a generative process as the reverse of a continuous-time noising process that gradually perturbs data toward a tractable reference distribution. Let  $\boldsymbol{\theta}_s \in \mathbb{R}^d$  denote the state at continuous time  $s \in [0, S]$ , with  $\boldsymbol{\theta}_0 \sim p_0$  and  $\boldsymbol{\theta}_S \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . In the variance-preserving (VP) formulation (Song et al., 2021), the forward SDE is:  $d\boldsymbol{\theta}_s = -\frac{1}{2}\beta_s \boldsymbol{\theta}_s ds + \sqrt{\beta_s} d\mathbf{w}_s$ , where  $\beta_s > 0$  is a noise schedule and  $\mathbf{w}_s$  is a standard  $d$ -dimensional Wiener process. The corresponding reverse-time diffusion that transports  $\boldsymbol{\theta}_S \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  back to  $p_0$ . The reverse diffusion SDE from  $s = S$  down to  $s = 0$  is as follows:

$$d\boldsymbol{\theta}_s = [-\frac{1}{2}\beta_s \boldsymbol{\theta}_s - \beta_s \nabla_{\boldsymbol{\theta}} \log p_s(\boldsymbol{\theta}_s)] ds + \sqrt{\beta_s} d\bar{\mathbf{w}}_s,$$

where  $\bar{\mathbf{w}}_s$  a standard Wiener process in reverse time and  $\nabla_{\boldsymbol{\theta}} \log p_s(\boldsymbol{\theta}_s)$  is the time-dependent score function where  $p_s(\boldsymbol{\theta}_s) = \int p_0(\boldsymbol{\theta}_0) p(\boldsymbol{\theta}_s | \boldsymbol{\theta}_0) d\boldsymbol{\theta}_0$ . Since score function  $\nabla_{\boldsymbol{\theta}} \log p_s(\boldsymbol{\theta}_s)$  is unknown, a neural network  $s_{\psi}(\boldsymbol{\theta}_s, s)$  is trained to approximate it via denoising score matching:

$$\psi^* = \operatorname{argmin}_{\psi} \mathbb{E}_{s \sim \mathcal{U}(\epsilon, 1), \boldsymbol{\theta}_0 \sim p_0, \boldsymbol{\theta}_s \sim p(\boldsymbol{\theta}_s | \boldsymbol{\theta}_0)} \|s_{\psi}(\boldsymbol{\theta}_s, s) - \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}_s | \boldsymbol{\theta}_0)\|_2^2,$$

where  $\epsilon > 0$  is a small constant and  $p(\boldsymbol{\theta}_s | \boldsymbol{\theta}_0)$  is Gaussian under the VP forward process. Plugging  $s_{\psi^*}$  into the reverse SDE and discretizing yields a sampler that approximately draws from  $p_0$ .

To obtain a practical training and sampling procedure, we discretize the time index into  $L$  noise levels indexed by  $\ell \in \{1, \dots, L\}$  and adopt discrete VP formulation. We define a discrete noise schedule  $\beta_{\ell} > 0$ , and define  $\alpha_{\ell} = 1 - \beta_{\ell}$ ,  $\bar{\alpha}_{\ell} = \prod_{j=1}^{\ell} \alpha_j$ . The forward process has the closed form:  $\boldsymbol{\theta}_{\ell} = \sqrt{\bar{\alpha}_{\ell}} \boldsymbol{\theta}_0 + \sqrt{1 - \bar{\alpha}_{\ell}} \mathbf{z}$  where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . The reverse process involves iteratively sampling  $\boldsymbol{\theta}_{\ell-1}$  from the posterior distribution  $p(\boldsymbol{\theta}_{\ell-1} | \boldsymbol{\theta}_{\ell})$ . Instead of directly approximating the score function  $\nabla_{\boldsymbol{\theta}_{\ell}} \log p_s(\boldsymbol{\theta}_{\ell})$ , it is common practice to reparameterize the model to predict the added noise  $\mathbf{z}$  at step  $\ell$  by denoiser network  $\varepsilon_{\phi^*}$ . The reverse update then becomes an ancestral Gaussian sampling step:  $\boldsymbol{\theta}_{\ell-1} \sim \mathcal{N}(\boldsymbol{\mu}_{\ell}, \boldsymbol{\Sigma}_{\ell})$ , where  $\boldsymbol{\mu}_{\ell} = \frac{1}{\sqrt{\alpha_{\ell}}} (\boldsymbol{\theta}_{\ell} - \frac{1 - \alpha_{\ell}}{\sqrt{1 - \bar{\alpha}_{\ell}}} \varepsilon_{\phi^*}(\boldsymbol{\theta}_{\ell}, \ell))$ ,  $\boldsymbol{\Sigma}_{\ell} = \frac{1 - \bar{\alpha}_{\ell-1}}{1 - \bar{\alpha}_{\ell}} \beta_{\ell} \mathbf{I}$ . This formulation is equivalent to DDPM (Ho et al., 2020) and serves as backbone for our diffusion prior.

### 3 DIFFUSION-BASED POSTERIOR SAMPLING IN MULTI-TASK NONLINEAR CONTEXTUAL BANDIT

In this section, we introduce the diffusion posterior sampling framework for multi-task nonlinear contextual bandits. We begin with the core components of a unified diffusion reverse update, and then instantiate it in two ways, following the prevailing classes of diffusion models in the literature.

#### 3.1 A UNIFIED DIFFUSION REVERSE UPDATE FOR DECISION-MAKING

We first recall posterior sampling in contextual bandits. Thompson Sampling (TS) maintains a posterior over reward model parameters, constructed from a hand-specified conjugate prior and the observed rewards. At round  $t$ , given history  $\mathcal{H}_t$  and prior  $p(\boldsymbol{\theta})$ , TS samples a parameter estimate  $\tilde{\boldsymbol{\theta}}_t \sim p(\boldsymbol{\theta} | \mathcal{H}_t) \propto p(\boldsymbol{\theta}) p(\mathcal{H}_t | \boldsymbol{\theta})$  and selects the arm  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}_t} f(\tilde{\boldsymbol{\theta}}_t; \mathbf{x})$ . However, such simple priors cannot capture the rich cross-task structure in multi-task settings, leading to misspecified posteriors and suboptimal exploration. To address this, we train a diffusion model on data from

past tasks to learn a flexible prior distribution over task parameters. For a new task, we run a conditional reverse-diffusion process that incorporates the interaction history. At each round, we initialize from noise and perform an  $L$ -step conditional reverse diffusion to produce a posterior sample (i.e., a parameter estimate) for arm selection.

We unify diffusion-based posterior sampling for decision-making with a single per-level reverse update (omitting coefficients):

$$\underbrace{\theta_{\ell-1}}_{\text{next state}} \leftarrow \underbrace{\theta_{\ell}}_{\text{current state}} + \underbrace{\text{unconditional term}}_{\text{diffusion prior drift}} + \underbrace{\text{likelihood term}}_{\text{data drift}} + \underbrace{\text{noise term}}_{\text{randomized exploration}}, \quad (3.1)$$

The **unconditional term** is the standard reverse step from a pretrained diffusion model, implemented via either a score network or a denoiser network. The **likelihood term** incorporates bandit interaction history up to round  $t - 1$  by defining a loss  $L_t$ , evaluating it at a diffusion-informed argument  $\Phi_{\ell}(\theta_{\ell})$  (the rescaled state or Tweedie estimate), and applying a gradient step  $-\eta_{\ell} \nabla_{\theta} L_t(\Phi_{\ell}(\theta_{\ell}))$ . The **noise term** enables randomized exploration, either as DDPM noise in the diffusion prior drift or Langevin noise from approximate sampling. This decomposition (3.1) separates offline knowledge (the diffusion prior drift) from online adaptation (the data drift).

In the next two subsections, we instantiate (3.1) in two ways. First, **DLTS** follows the TS workflow and performs Langevin updates at each level using the rescaled state, yielding a diffusion-prior approximate sampler for nonlinear reward models. Second, **DPSG** preserves the reverse-diffusion drift and adds a single likelihood drift per level via the Tweedie estimate, resulting in a simple and efficient conditional sampler. Both methods follows the unified diffusion posterior sampling framework in (3.1), while each admits its own interpretation and derivation.

### 3.2 DIFFUSION LANGEVIN THOMPSON SAMPLING

To handle contextual bandits with complex priors, Kveton et al. (2024) proposed DiffTS for linear and generalized linear models with specialized closed-form updates. For nonlinear bandits, such closed forms are unavailable, making algorithm design harder. Motivated by LMC-TS (Xu et al., 2022), an approximate sampling method for nonlinear bandits, we propose **Diffusion Langevin Thompson Sampling (DLTS)** (Algorithm 1). DLTS extends approximate TS to nonlinear rewards (including neural networks) by using a learned diffusion prior.

---

#### Algorithm 1 Diffusion Langevin Thompson Sampling (DLTS)

---

**Input:** Diffusion denoiser  $\varepsilon_{\phi^*}(\cdot, \cdot)$ , noise schedule  $\{\beta_{\ell}\}_{\ell=1}^L$ , learning rate  $\{\eta_{\ell}\}_{\ell=1}^L$ .

- 1:  $\alpha_{\ell} \leftarrow 1 - \beta_{\ell}$  and  $\bar{\alpha}_{\ell} \leftarrow \prod_{j=1}^{\ell} \alpha_j$  for all  $\ell$
  - 2: **for**  $t = 1, 2, \dots, T$  **do**
  - 3:   Receive contextual vector  $\{\mathbf{x}_t(a)\}_{a \in \mathcal{A}}$ .
  - 4:    $(\boldsymbol{\mu}_{L+1}, \boldsymbol{\Sigma}_{L+1}) \leftarrow (\mathbf{0}, \mathbf{I})$ .
  - 5:   **for**  $\ell = L + 1, L, \dots, 1$  **do**
  - 6:      $\theta_{\ell,0} \leftarrow \boldsymbol{\mu}_{\ell}$ .
  - 7:     **for**  $k = 0, 1, \dots, K_{\ell} - 1$  **do**
  - 8:       Sample  $\boldsymbol{\xi}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .
  - 9:        $\theta_{\ell,k+1} \leftarrow \left(1 - \frac{(1-\bar{\alpha}_{\ell})\eta_{\ell}}{(1-\bar{\alpha}_{\ell-1})\beta_{\ell}}\right)\theta_{\ell,k} + \frac{(1-\bar{\alpha}_{\ell})\eta_{\ell}}{(1-\bar{\alpha}_{\ell-1})\beta_{\ell}}\boldsymbol{\mu}_{\ell} - \eta_{\ell} \nabla_{\theta} L_t(\theta_{\ell,k}/\sqrt{\bar{\alpha}_{\ell-1}}) + \sqrt{2\eta_{\ell}\zeta_{\ell}^{-1}}\boldsymbol{\xi}_k$ .
  - 10:     **end for**
  - 11:      $\theta_{\ell-1} \leftarrow \theta_{\ell,K_{\ell}}$ .
  - 12:      $\boldsymbol{\mu}_{\ell-1} \leftarrow \frac{1}{\sqrt{\bar{\alpha}_{\ell-1}}}(\theta_{\ell-1} - \frac{1-\alpha_{\ell-1}}{\sqrt{1-\bar{\alpha}_{\ell-1}}}\varepsilon_{\phi^*}(\theta_{\ell-1}, \ell - 1))$ .
  - 13:   **end for**
  - 14:   Let  $\theta_t \leftarrow \theta_0$ , take action  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x}_t(a)} f(\tilde{\theta}_t; \mathbf{x}_t(a))$ .
  - 15:   Receive reward  $y_t$  and update history  $\mathcal{H}_{t+1} = \{(\mathbf{x}_i, y_i)\}_{i=1}^t$ .
  - 16: **end for**
- 

Our method follows the TS workflow, which builds a posterior over parameters and samples for arm selection, but it differs in three ways: (1) It replaces a hand-crafted prior with a diffusion prior trained on past tasks. (2) Instead of a single posterior update, it runs a reverse diffusion process and injects interaction history at each level, building the conditional posterior gradually. Formally, the reverse process is a Markov chain with the per-level transition in the form:

$$p(\theta_{\ell-1} | \theta_{\ell}, \mathcal{H}_t) \propto p(\mathcal{H}_t | \theta_{\ell-1}) p(\theta_{\ell-1} | \theta_{\ell}), \quad \ell = L, \dots, 1, \quad (3.2)$$

where  $\theta_\ell \in \mathbb{R}^d$  is the noisy parameter at level  $\ell$  and the terminal  $\theta_0$  is used as the parameter sample for action selection. (3) It applies iterative LMC updates to approximately sample from (3.2).

After training the diffusion denoiser  $\varepsilon_{\phi^*}(\cdot, \cdot)$  (see Appendix G for more details), we obtain the diffusion model parameter  $(\mu_\ell, \Sigma_\ell)_{\ell \in [L+1]}$ <sup>1</sup> for unconditional sampling in each reverse step:  $p(\theta_{\ell-1}|\theta_\ell) = \mathcal{N}(\theta_{\ell-1}; \mu_\ell, \Sigma_\ell)$ , where  $\mu_\ell = \frac{1}{\sqrt{\alpha_\ell}}(\theta_\ell - \frac{1-\alpha_\ell}{\sqrt{1-\alpha_\ell}}\varepsilon_{\phi^*}(\theta_\ell, \ell))$ ,  $\Sigma_\ell = \frac{1-\bar{\alpha}_{\ell-1}}{1-\alpha_\ell}\beta_\ell\mathbf{I}$ . Note that at step  $t$ , the nonlinear reward model is defined as  $y_t = f(\theta^*; \mathbf{x}_t) + \varepsilon_t$ . When we assume Gaussian noise  $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$ , then the likelihood becomes  $p(\mathcal{H}_t|\theta) \propto \exp(-\sigma^{-2} \sum_{i=1}^t (f(\theta; \mathbf{x}_i) - y_i)^2) = \exp(-L_t(\theta))$ , where  $L_t(\theta) = \sigma^{-2} \sum_{i=1}^t (f(\theta; \mathbf{x}_i) - y_i)^2$ . Therefore, we can approximate the posterior in each reverse step for conditional sampling as follows (refer to Appendix C.1 for details of derivation),

$$p(\theta_{\ell-1}|\theta_\ell, \mathcal{H}_t) \propto \exp(-L_t(\theta_{\ell-1}/\sqrt{\bar{\alpha}_{\ell-1}})) \cdot \mathcal{N}(\theta_{\ell-1}; \mu_\ell, \Sigma_\ell).$$

To sample from  $p(\theta_{\ell-1}|\theta_\ell, \mathcal{H}_t)$ , we apply Langevin Monte Carlo for approximate sampling. Specifically, at round  $t$  and diffusion reverse step  $\ell$ , we iteratively conduct the Langevin update,

$$\theta_{\ell, k+1} = \theta_{\ell, k} - \eta_\ell \left[ \nabla_{\theta} L_t(\theta_{\ell, k}/\sqrt{\bar{\alpha}_{\ell-1}}) + \Sigma_\ell^{-1}(\theta_{\ell, k} - \mu_\ell) \right] + \sqrt{2\eta_\ell \zeta_\ell^{-1}} \xi_k,$$

where  $\xi_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\eta_\ell$  is the step size, and  $\zeta_\ell$  is the temperature. We initialize  $\theta_{\ell, 0} = \mu_\ell$ , after  $K_\ell$  iterations, we set  $\theta_{\ell-1} = \theta_{\ell, K_\ell}$  as the sample from the reverse step.

Specially, when we only conduct one LMC update ( $K_\ell = 1$ ), we have the reverse update

$$\theta_{\ell-1} \leftarrow \left(1 - \frac{(1-\bar{\alpha}_\ell)\eta_\ell}{(1-\bar{\alpha}_{\ell-1})\beta_\ell}\right) \theta_\ell + \frac{(1-\bar{\alpha}_\ell)\eta_\ell}{(1-\bar{\alpha}_{\ell-1})\beta_\ell} \mu_\ell - \eta_\ell \nabla_{\theta} L_t(\theta_\ell/\sqrt{\bar{\alpha}_{\ell-1}}) + \sqrt{2\eta_\ell \zeta_\ell^{-1}} \xi. \quad (3.3)$$

Note that (3.3) aligns with our unified update (3.1). This confirms the rationality of our framework and the essence of diffusion-based posterior sampling. Intuitively, DLTS explicitly forms a conditional posterior by incorporating the history into the diffusion prior at each reverse step, then draws samples from it via approximate sampling.

### 3.3 DIFFUSION POSTERIOR SAMPLING WITH GUIDANCE

We now give a simpler and faster realization inspired by Diffusion Posterior Sampling (DPS) (Chung et al., 2023), a common approach for conditional sampling in inverse problems. DPS augments the unconditional reverse step with a likelihood score term and then runs the reverse chain, yielding a conditional sampler. With a pretrained unconditional process, drawing a parameter sample for arm selection is straightforward: run the unconditional reverse steps and add a likelihood drift for guidance. This requires only the unconditional score and a tractable proxy for the likelihood at the Tweedie estimate, so there is no need to construct a per-level conditional posterior and apply approximate sampling as in DLTS. Following this design, we propose **Diffusion Posterior Sampling with Guidance (DPSG)**, shown in Algorithm 2.

To derive diffusion reverse update (Line 8 in Algorithm 2), we follow Chung et al. (2023) and decompose the conditional score into two terms,

$$\nabla_{\theta} \log p_s(\theta_s|\mathcal{H}_t) = \underbrace{\nabla_{\theta} \log p_s(\theta_s)}_{\text{unconditional score}} + \underbrace{\nabla_{\theta} \log p_s(\mathcal{H}_t|\theta_s)}_{\text{likelihood score}}.$$

For the first term, we can use pretrained diffusion model score network  $s_{\psi^*}$  to approximate the unconditional score. For the second term, we approximate the likelihood score via Tweedie’s formula. Since  $p(\mathcal{H}_t|\theta) \propto \exp(-L_t(\theta))$ , based on Tweedie’s formula (Efron, 2011), we can obtain a tractable approximation for  $\nabla_{\theta} \log p_s(\mathcal{H}_t|\theta_s)$  (refer to Appendix C.2 for details of derivation):  $\nabla_{\theta} \log p_s(\mathcal{H}_t|\theta_s) \simeq -\nabla_{\theta} L_t(\hat{\theta}_0(\theta_s, s))$ , where  $\hat{\theta}_0(\theta_s, s) \simeq 1/\sqrt{\bar{\alpha}_s}(\theta_s + (1-\bar{\alpha}_s)s_{\psi^*}(\theta_s, s))$ . Therefore, we can further have the approximation of the conditional score of reverse dynamics,

$$\nabla_{\theta} \log p_s(\theta_s|\mathcal{H}_t) \simeq s_{\psi^*}(\theta_s, s) - \nabla_{\theta} L_t(\hat{\theta}_0(\theta_s, s)). \quad (3.4)$$

In practice, to obtain a discrete-time algorithm we use DDPM ancestral sampling (Ho et al., 2020) to implement conditional sampling via the conditional score (3.4). The reverse update is done by

<sup>1</sup>We define  $(\mu_{L+1}, \Sigma_{L+1}) = (\mathbf{0}, \mathbf{I})$  for notation simplicity.

**Algorithm 2** Diffusion Posterior Sampling with Guidance (DPSG)

**Input:** score network  $s_{\psi^*}(\cdot, \cdot)$ , noise schedule  $\{\beta_\ell\}_{\ell=1}^L$ , learning rate  $\{\eta_\ell\}_{\ell=1}^L$ .

```

1:  $\alpha_\ell \leftarrow 1 - \beta_\ell$  and  $\bar{\alpha}_\ell \leftarrow \prod_{j=1}^\ell \alpha_j$  for all  $\ell$ .
2: for  $t = 1, 2, \dots, T$  do
3:   Receive contextual vector  $\{\mathbf{x}_i(a)\}_{a \in \mathcal{A}}$ .
4:    $\boldsymbol{\theta}_L \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .
5:   for  $\ell = L, L-1, \dots, 1$  do
6:      $\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell) \leftarrow \frac{1}{\sqrt{\alpha_\ell}}(\boldsymbol{\theta}_\ell + (1 - \bar{\alpha}_\ell)s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell))$ .
7:      $\mathbf{z}_\ell \sim \mathcal{N}(\mathbf{0}, \beta_\ell \mathbf{I})$ .
8:      $\boldsymbol{\theta}_{\ell-1} \leftarrow \frac{1}{\sqrt{\alpha_\ell}}\boldsymbol{\theta}_\ell + \frac{\beta_\ell}{\sqrt{\alpha_\ell}}s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell) - \eta_\ell \nabla_{\boldsymbol{\theta}} L_t(\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell)) + \mathbf{z}_\ell$ .
9:   end for
10:  Let  $\tilde{\boldsymbol{\theta}}_t \leftarrow \boldsymbol{\theta}_0$ , take action  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x}_t(a)} f(\tilde{\boldsymbol{\theta}}_t; \mathbf{x}_t(a))$ .
11:  Receive reward  $y_t$  and update history  $\mathcal{H}_{t+1} = \{(\mathbf{x}_i, y_i)\}_{i=1}^t$ .
12: end for

```

two steps. First, we use the unconditional score  $s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell)$  through Tweedie’s estimate  $\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell)$  to form the Gaussian posterior mean, then we sample

$$\boldsymbol{\theta}'_{\ell-1} \leftarrow \frac{\sqrt{\alpha_\ell(1-\bar{\alpha}_{\ell-1})}}{1-\bar{\alpha}_\ell}\boldsymbol{\theta}_\ell + \frac{\sqrt{\bar{\alpha}_{\ell-1}\beta_\ell}}{1-\bar{\alpha}_\ell}\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell) + \mathbf{z}_\ell, \quad \mathbf{z}_\ell \sim \mathcal{N}(\mathbf{0}, \beta_\ell \mathbf{I}), \quad (3.5)$$

Based on (3.4), the second step is that we make this update conditional by subsequently adding one likelihood-based term. In practice, after sampling  $\boldsymbol{\theta}'_{\ell-1}$ , we take a gradient step using the likelihood  $L_t$  at the Tweedie estimate for guidance:  $\boldsymbol{\theta}_{\ell-1} \leftarrow \boldsymbol{\theta}'_{\ell-1} - \eta_\ell \nabla_{\boldsymbol{\theta}} L_t(\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell))$ . When we merge this with (3.5), we obtain line 8 in Algorithm 2,

$$\boldsymbol{\theta}_{\ell-1} \leftarrow \frac{1}{\sqrt{\alpha_\ell}}\boldsymbol{\theta}_\ell + \frac{\beta_\ell}{\sqrt{\alpha_\ell}}s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell) - \eta_\ell \nabla_{\boldsymbol{\theta}} L_t(\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell)) + \mathbf{z}_\ell. \quad (3.6)$$

Note that (3.6) aligns with our unified update (3.1). This confirms the rationality of our framework and the essence of diffusion-based posterior sampling. However, from the intuition perspective, DPSG is to use history information as guidance after unconditional update based on the diffusion prior. This is different from DLTS, which directly integrates history information into the unconditional update based on the diffusion prior to achieve conditional sampling.

**Empirical Variant of DPSG** DPSG is simple and fast, but it relies on the accuracy of likelihood-score at the Tweedie estimate. A recent work Xu et al. (2025) observes that DPS has the properties of high bias and low diversity, thus behaving like an implicit, but unstable, MAP estimator. Therefore, multi-step projection is used to solve this MAP optimization problem. Motivated by Xu et al. (2025), to stabilize learning, we introduce **DPSG with Multi-step Projection (DPSG-MP)** (Algorithm 3 in Appendix C.3), which is a practical DPSG variant that replaces the single correction with a short inner loop of gradient ascent on the likelihood score, and after each step projects the iterate onto the sphere where the reverse transition  $p_{\psi^*}(\boldsymbol{\theta}_{\ell-1}|\boldsymbol{\theta}_\ell)$  concentrates around. Specifically, we first initialize  $\boldsymbol{\theta}_{\ell-1,0} \leftarrow 1/\sqrt{\alpha_\ell}\boldsymbol{\theta}_\ell + \beta_\ell/\sqrt{\alpha_\ell}s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell) + \mathbf{z}_\ell$ , then we iteratively update

$$\boldsymbol{\theta}_{\ell-1,k+1} \leftarrow \operatorname{Proj}(\boldsymbol{\theta}_{\ell-1,k} - \eta_\ell \nabla_{\boldsymbol{\theta}} L_t(\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell))), \quad k = 0, 1, \dots, K_\ell,$$

where function  $\operatorname{Proj}(\cdot)$  projects the input onto sphere surface  $\mathcal{S}(p_{\psi^*}(\boldsymbol{\theta}_{\ell-1}|\boldsymbol{\theta}_\ell))$ , with radius  $r_\ell = \|\mathbf{z}_\ell\|_2$  and center at  $\mathbb{E}[\boldsymbol{\theta}_{\ell-1}|\boldsymbol{\theta}_\ell]$ . Finally, we assign  $\boldsymbol{\theta}_{\ell-1,K_\ell}$  to  $\boldsymbol{\theta}_{\ell-1}$  and finish the update. We provide the complete algorithm and further explanation in Appendix C.3.

## 4 THEORETICAL RESULTS

In this section, we analyze the connection between Thompson sampling and diffusion Thompson Sampling. For convenience of analysis, we first define two oracle algorithms.

**Definition 4.1.** We define the following oracle algorithm as **Oracle Thompson Sampling (OTS)**, which applies exact posterior sampling and greedy policy according to the true prior and likelihood:

- 1) Use true prior  $p_0$ .
- 2) Maintain exact posterior in each round:  $p_t(\boldsymbol{\theta}) \equiv p(\boldsymbol{\theta}|\mathcal{H}_t) \propto p_0(\boldsymbol{\theta}) \prod_{i=1}^{t-1} p(y_i|\boldsymbol{\theta}, \mathbf{x}_i)$ .

3) Sample  $\theta_t^{OTS} \sim p_t(\theta)$ , then select arm  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}_t} f(\theta_t^{OTS}; \mathbf{x})$ .

**Definition 4.2.** We define the following oracle algorithm as **Oracle Diffusion Thompson Sampling (ODTS)**, which performs posterior sampling via a reverse-diffusion chain, using the same pretrained diffusion model as DLTS (Algorithm 1).

- 1) Train a diffusion reverse process to learn prior  $p_0$  and obtain score/denoiser network  $s_{\psi^*}(\cdot, \cdot)$ .
- 2) The per-level conditional reverse process is implemented exactly:  $\tilde{p}_t(\theta_{0:L}) \equiv p(\theta_{0:L} | \mathcal{H}_t) = p(\theta_L | \mathcal{H}_t) \prod_{\ell=1}^L p(\theta_{\ell-1} | \theta_\ell, \mathcal{H}_t)$  and  $\tilde{p}_t(\theta) = \int \tilde{p}_t(\theta_{0:L}) d\theta_{1:L} = \int p(\theta_{0:L} | \mathcal{H}_t) d\theta_{1:L}$ .
- 3) Sample  $\theta_t^{ODTS} \sim \tilde{p}_t(\theta)$ , then selects arm  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}_t} f(\theta_t^{ODTS}; \mathbf{x})$ .

We make the following assumption to measure the accuracy of score estimation.

**Assumption 4.3.** For any diffusion time  $s \in [0, \mathcal{S}]$  and let  $p_s$  denote the  $s$ -marginal distribution of the forward diffusion over parameters  $\theta \in \Theta \subseteq \mathbb{R}^d$ . Let  $s_{\psi^*}(\theta, s)$  be the learned score estimator. Assume there exists  $\epsilon_{score}$  such that for all  $s$  and  $\theta$ ,  $\|s_{\psi^*}(\theta, s) - \nabla_{\theta} \log p_s(\theta)\|_2 \leq \epsilon_{score}$ .

Note that the key distinction between OTS and ODTS is whether the reverse-diffusion process is trained to recover the true prior  $p_0$ . The theorem below shows that if ODTS learns the prior exactly (i.e.  $\epsilon_{score} = 0$  in Assumption 4.3), then OTS and ODTS are equivalent.

**Theorem 4.4.** If ODTS perfectly train a diffusion reverse process to get true prior  $p_0$  and the score/denoiser network matches the true score at every diffusion level  $\ell$ :  $s_{\psi^*}(\cdot, \ell) \equiv \nabla_{\theta} \log p_{\ell}(\cdot)$ , then oracle algorithms OTS and ODTS are equivalent in the sense that they induce the same posterior distribution  $\tilde{p}_t(\theta) = p_t(\theta)$  and both algorithms induce the same arm selection distribution.

**Theorem 4.5.** For any round  $t \in [T]$ , denote  $p_t(\theta)$  as the exact posterior used by OTS, and  $\tilde{p}_t(\theta)$  as the marginal posterior used by ODTS. Assume the reward is bounded: for all  $(\theta, \mathbf{x})$ ,  $|f(\theta; \mathbf{x})| \leq f_{max}$ . Then the per-round expected regret gap between OTS and ODTS is bounded by  $\mathbb{E}[r_t^{ODTS} - r_t^{OTS}] \leq \Delta_t^{Score} = 2f_{max} \sum_{\ell=1}^L \kappa_{\ell} \epsilon_{score}$ , where coefficient  $\kappa_{\ell}$  is determined by noise schedule.

## 5 EXPERIMENTS

In this section, we conduct experiments to evaluate our proposed algorithms DLTS and a variant of DPSG named DPSG-MP (refer to Appendix C.3). We aim to demonstrate the effectiveness of our proposed algorithms in the different settings such as the linear and non-linear contextual bandit. We conduct these experiments on simulation experiments.

### 5.1 IMPLEMENTATIONS

Our diffusion prior is implemented within a Denoising Diffusion Probabilistic Model (DDPM) framework, using the standard  $\epsilon$ -prediction parameterization. The denoiser network is an MLP that incorporates time embeddings. For synthetic simulations, we use  $L = 100$  diffusion steps and a linear noise schedule with a constant  $\beta_{\ell}$ . A full description of the architectures and hyperparameters is available in Appendix H.

### 5.2 SIMULATION

**Prior Distribution** We follow Kveton et al. (2024) to design the prior over task parameters. We consider six prior distributions, including the ‘cross’, ‘rays’, ‘triangles’, ‘swirl’, ‘H’ and ‘corners’ in our simulation experiments. Each prior distribution has 10000 samples with  $d = 2$  parameters defined as  $\theta$ . We use 80% of the samples as the training samples, which refer to the previous task. For the remaining 20% samples, we use them as the test samples, which refer to the new tasks. We illustrate the prior distribution in Figure 6 in Appendix H.2. Based on the prior distribution, we design the linear and non-linear contextual bandit on different reward models.

**Linear Bandit** We first evaluate our methods in the linear contextual bandit setting. For each task, the ground-truth parameter  $\theta$  is drawn from the test set of a given prior distribution, and rewards are generated according to the model  $r = \mathbf{x}^{\top} \theta + \varepsilon$  where the  $\varepsilon \sim \mathcal{N}(0, 1)$ . We compare our proposed DLTS and DPSG-MP which is a more stable version of DPSG with several baseline methods, including the LinTS (Agrawal & Goyal, 2013b), LinUCB (Chu et al., 2011),  $\epsilon$ -greedy and DiffTS (Kveton et al., 2024). For LinTS, LinUCB, and  $\epsilon$ -greedy algorithms, we directly conduct them in the test bandit task. Our methods and DiffTS leverage a diffusion prior trained on the 80% training

partition of the parameter data. Performance is measured over 64 new tasks per trial, with each task running for a horizon of 200 steps. To ensure statistical robustness, all results are averaged over 8 independent trials. The quality of the learned diffusion prior is visualized in Figure 1, and the comparative regret performance is summarized in Figure 2.

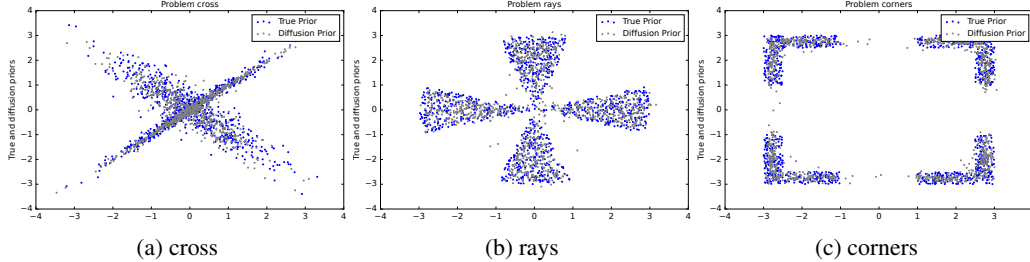


Figure 1: Visualization of the learned diffusion prior versus the true prior. We demonstrate three of them including cross, rays and corners. For each figure, samples from our trained model (in grey) are overlaid on the ground-truth samples (in blue). More results can be found in Appendix I.

First, Figure 1 validates the quality of our diffusion model, confirming that it accurately captures the complex structure of the true prior distributions. From the regret comparison in Figure 2, we observe that our methods DLTS and DPSG-MP achieve cumulative regret that is highly competitive with, and in some cases on par with, the specialized baseline algorithms. This is a noteworthy outcome, as methods like LinTS and DiffTS are designed specifically for the linear setting and leverage efficient closed-form posterior updates. In contrast, our framework relies on a more general, gradient-based sampling approach. The strong performance of our methods in a setting where they cannot exploit such closed-form solutions underscores the robustness of our framework and highlights its potential for nonlinear environments where those specialized updates are no longer applicable. The good performance of the algorithms also demonstrates the benefits of leveraging the diffusion prior learned on the previous task in solving the new tasks with the same prior distribution.

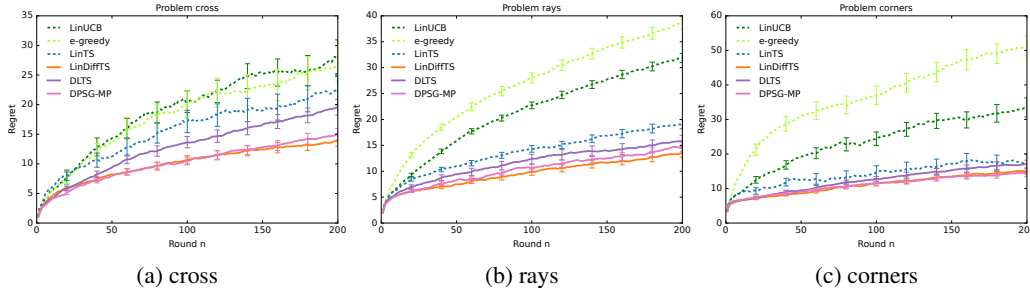


Figure 2: Performance of all algorithms on the linear contextual bandit tasks for the three prior distributions. Our proposed methods achieve cumulative regret comparable to LinTS and DiffTS. More results can be found in Appendix I.

**Nonlinear Bandit** We also evaluate our framework in scenarios where the reward function is inherently nonlinear. We construct three such environments where the reward is defined by  $r = f(\mathbf{x}^\top \boldsymbol{\theta}) + \varepsilon$ , with  $\varepsilon \sim \mathcal{N}(0, 1)$ . The nonlinear functions  $f(\cdot)$  we consider are a cosine model ( $f(z) = \cos(3z)$ ), a quadratic model ( $f(z) = z^2$ ), and a sigmoid-gated model ( $f(z) = 2z \cdot \text{sigmoid}(z) + 1$ ). Due to the space limit, the results of the cosine model and the quadratic model in Appendix I.2.

In this challenging setting, our proposed methods DLTS and DPSG-MP, use a neural network (MLP) to represent the reward function. A key distinction from prior sections is that the diffusion model is now trained to learn a prior directly over the parameters of this MLP, using network weights from past tasks as training data. This setup tests the ability of our framework to handle priors in the high-dimensional weight space of neural networks. We compare our proposed algorithms with several baseline methods, including the LinTS, DiffTS, NeuralTS (Zhang et al., 2021) and NeuralUCB (Zhou et al., 2020). Visualizations of the learned diffusion priors over the MLP parameters for each reward model are provided in Figure 3.

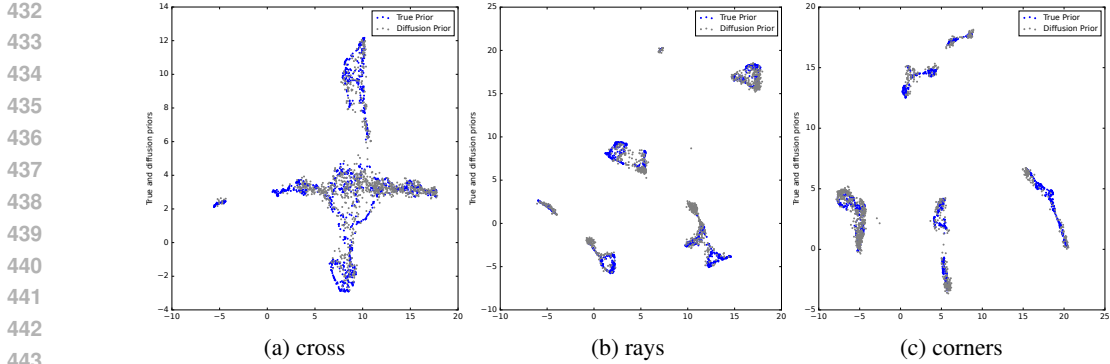


Figure 3: Visualization of the learned diffusion prior versus the true prior (under sigmoid-gated reward model). The true prior is the MLP parameters learned from the previous task. Samples from our trained model (grey) are overlaid on the ground-truth samples (blue). More results can be found in Appendix I.

When tested on the more complex sigmoid-gated reward model, our proposed methods demonstrate a clear performance advantage, as illustrated in Figure 4. Both DLTS and DPSG-MP consistently outperform the strong NeuralTS and NeuralUCB baselines, achieving the lowest cumulative regret. This superior performance underscores the primary benefit of our framework. By leveraging a powerful, pretrained diffusion prior, our algorithms can explore the parameter space more effectively and adapt more quickly in challenging nonlinear environments. This stands in contrast to the baseline methods, which must learn from scratch on each new task. The plotted results, which show the mean cumulative regret and standard error averaged over 64 tasks across 8 independent trials, confirm this significant performance gain.

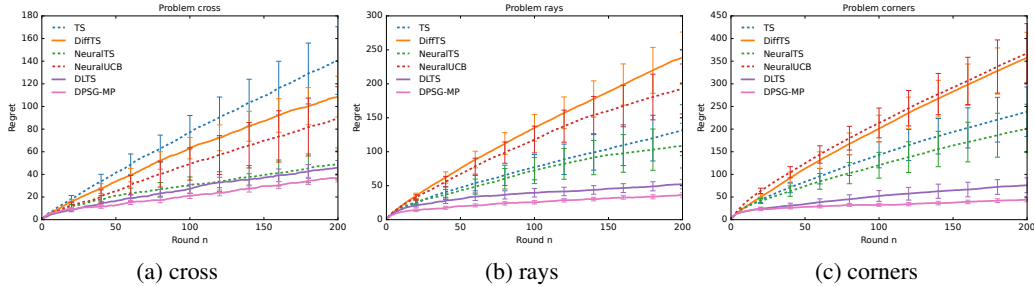


Figure 4: Performance of all algorithms on the sigmoid-gated nonlinear bandit tasks for the three prior distributions. Our proposed methods outperform strong baselines like NeuralTS and NeuralUCB, achieving the lowest cumulative regret. More results can be found in Appendix I.

Table 1: Results of ablation study on Diffusion Steps  $L$ . Reported values are Cumulative Regret at  $T = 200$ . The regrets are averaged over 64 tasks across 8 independent runs.

Environment	Algorithm	Diffusion Steps ( $L$ )		
		$L = 20$	$L = 50$	$L = 100$
Rays	DLTS	136.08 ± 3.18	53.51 ± 4.42	<b>51.20 ± 2.90</b>
	DPSG-MP	63.90 ± 0.72	50.28 ± 1.09	<b>36.31 ± 1.09</b>
Triangles	DLTS	433.16 ± 7.51	97.68 ± 5.95	<b>47.29 ± 1.79</b>
	DPSG-MP	86.57 ± 0.66	73.09 ± 1.19	<b>52.20 ± 3.54</b>
Swirl	DLTS	732.70 ± 7.29	125.55 ± 4.06	<b>53.95 ± 1.74</b>
	DPSG-MP	75.45 ± 0.59	52.47 ± 0.51	<b>38.34 ± 1.63</b>

### 5.3 ABLATION STUDY AND TIME ANALYSIS

**Sensitivity to Inner-Loop Updates  $K_\ell$ .** To analyze the sensitivity of our algorithms, we conduct ablation studies on the number of diffusion steps  $L$  and the number of inner-loop updates  $K_\ell$  in Appendix I.2. These experiments are performed using the sigmoid-gated nonlinear bandit benchmark, a challenging setting from Section 5.2. We report the cumulative regret at  $T = 200$ , averaged over 64 tasks across 8 independent runs.

We examined the impact of the diffusion depth  $L$  while keeping other parameters fixed. As shown in Table 1, increasing  $L$  consistently improves performance. A larger number of diffusion steps allows for a more refined reverse generation process, resulting in more accurate posterior sampling and significantly lower regret for both DLTS and DPSG-MP.

**Wall-Clock Time Analysis.** To rigorously assess the trade-off between computational cost and sample efficiency, we perform a time-to-accuracy analysis rather than relying solely on per-round inference latency. This metric effectively normalizes for sample efficiency, demonstrating that a computationally more intensive algorithm (such as ours) can achieve superior overall time efficiency if it requires significantly fewer interactions to reach a high-quality policy compared to faster but less data-efficient baselines.

We benchmark all algorithms in the nonlinear sigmoid-gated bandit setting using the ‘triangles’ prior distribution. We conduct our DLTS and DPSG-MP with inner-loop  $K_\ell$  to be 10 and 1 respectively. We define a target performance threshold of 0.5 average regret (calculated as cumulative regret divided by the number of rounds). For each algorithm, we measure the total wall-clock time required to achieve this threshold. A maximum time budget of 3000 seconds is enforced. For fair comparison, all evaluations are conducted on an Intel Xeon E5-2640 v4 CPU. The results, summarized in Table 2, show that our methods (DLTS and DPSG-MP) reach the performance threshold significantly faster than the baselines, despite having a higher per-step computational cost, due to their superior sample efficiency.

Table 2: Analysis of time to achieve the same average regret which is 0.5 in the sigmoid-gated nonlinear bandit with the Triangles problem in 64 task with 1 runs. The results show that our methods reach the performance threshold significantly faster than the baselines (NeuralTS and NeuralUCB).

Algorithm	Round	Total Time (s)
DiffTS	> 2000	> 3000
NeuralTS	1505	> 3000
NeuralUCB	799	~ 1967.71
<b>DLTS (Ours)</b>	63	~ <b>1823.81</b>
<b>DPSG-MP (Ours)</b>	54	~ <b>173.63</b>

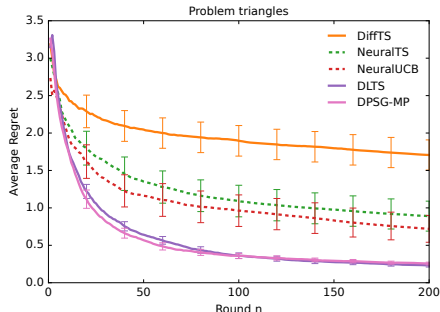


Figure 5: Results of average regret in Triangles with  $T = 200$ . It shows that our DLTS and DPSG-MP achieve the best performance.

## 6 CONCLUSION

We address multi-task nonlinear contextual bandits by learning a flexible diffusion prior from past tasks and performing posterior sampling via a conditional reverse-diffusion process. We provide a unified update framework that, at each step, combines (i) an unconditional drift from the diffusion prior, (ii) a likelihood-driven drift from the interaction history, and (iii) a noise term enabling randomized exploration. Built on this view, we instantiate two variants: **DLTS**, which integrates history into the diffusion prior at every reverse step to form a conditional posterior and draw approximate samples; and **DPSG**, which first performs unconditional reverse sampling from the pretrained diffusion prior and then applies a single history-guided gradient correction. In theory part, we analyzed the connection and expected regret gap between Thompson Sampling and diffusion Thompson Sampling by formalizing corresponding oracle algorithms. In simulations across linear and non-linear contextual bandits, our methods are competitive with closed-form baselines in the linear regime (e.g., LinTS and DiffTS) and achieve consistent gains in neural settings. These demonstrate the potential of our algorithms to solve more realistic problems.

540 REPRODUCIBILITY STATEMENT

541  
542 Our code will be made available in the supplementary material and will be publicly released upon  
543 acceptance. The implementation of all our experiments is based on the PyTorch framework. Our  
544 simulation environment, particularly the design of the prior distributions, follows the benchmark  
545 described in Kveton et al. (2024).

546 We provide detailed pseudocode for our proposed algorithms, DLTS, DPSG, and DPSG-MP, in Sec-  
547 tion 3 in Algorithm 1, Algorithm 2 and Algorithm 3. We provide implementation details, including  
548 the network architectures for the diffusion models in Table 3 in Appendix H.1. Furthermore, all  
549 hyperparameters used for the nonlinear bandit settings are detailed in Table 4 and Table 5 in Ap-  
550 pendix H.1.

552 REFERENCES

- 553  
554 Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs.  
555 In *International conference on machine learning*, pp. 127–135. PMLR, 2013a.
- 556 Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear payoffs.  
557 In Sanjoy Dasgupta and David McAllester (eds.), *Proceedings of the 30th International Con-*  
558 *ference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pp.  
559 127–135, Atlanta, Georgia, USA, 17–19 Jun 2013b. PMLR. URL [https://proceedings.](https://proceedings.mlr.press/v28/agrawal13.html)  
560 [mlr.press/v28/agrawal13.html](https://proceedings.mlr.press/v28/agrawal13.html).
- 561  
562 Shipra Agrawal and Navin Goyal. Near-optimal regret bounds for thompson sampling. *Journal of*  
563 *the ACM (JACM)*, 64(5):1–24, 2017.
- 564 Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B. Tenenbaum, Tommi S. Jaakkola, and Pulkit Agrawal.  
565 Is conditional generative modeling all you need for decision making? In *The Eleventh Interna-*  
566 *tional Conference on Learning Representations*, 2023. URL [https://openreview.net/](https://openreview.net/forum?id=sP1fo2K9DFG)  
567 [forum?id=sP1fo2K9DFG](https://openreview.net/forum?id=sP1fo2K9DFG).
- 568  
569 Imad Aouali. Diffusion models meet contextual bandits with large action spaces. *arXiv preprint*  
570 *arXiv:2402.10028*, 2024.
- 571 Sébastien Bubeck, Nicolo Cesa-Bianchi, and Gábor Lugosi. Bandits with heavy tail. *IEEE Trans-*  
572 *actions on Information Theory*, 59(11):7711–7717, 2013.
- 573  
574 Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. *Advances in neural*  
575 *information processing systems*, 24, 2011.
- 576 Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff func-  
577 tions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and*  
578 *Statistics*, pp. 208–214. JMLR Workshop and Conference Proceedings, 2011.
- 579  
580 Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul  
581 Ye. Diffusion posterior sampling for general noisy inverse problems. In *The Eleventh Interna-*  
582 *tional Conference on Learning Representations*, 2023. URL [https://openreview.net/](https://openreview.net/forum?id=OnD9zGAGT0k)  
583 [forum?id=OnD9zGAGT0k](https://openreview.net/forum?id=OnD9zGAGT0k).
- 584 Pierre Clavier, Tom Huix, and Alain Oliviero Durmus. VITS : Variational inference thompson  
585 sampling for contextual bandits. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian  
586 Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st*  
587 *International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning*  
588 *Research*, pp. 9033–9075. PMLR, 21–27 Jul 2024. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v235/clavier24a.html)  
589 [press/v235/clavier24a.html](https://proceedings.mlr.press/v235/clavier24a.html).
- 590 Zehao Dou and Yang Song. Diffusion posterior sampling for linear inverse problem solving: A  
591 filtering perspective. In *The Twelfth International Conference on Learning Representations*, 2024.
- 592  
593 Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Associa-*  
*tion*, 106(496):1602–1614, 2011.

- 594 Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. *Advances*  
595 *in neural information processing systems*, 31, 2018.  
596
- 597 Florence Forbes and Darren Wraith. A new family of multivariate heavy-tailed distributions with  
598 variable marginal amounts of tailweight: application to robust clustering. *Statistics and comput-*  
599 *ing*, 24(6):971–984, 2014.
- 600 Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
601 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information*  
602 *processing systems*, 27, 2014.  
603
- 604 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
605 *neural information processing systems*, 33:6840–6851, 2020.  
606
- 607 Joey Hong, Branislav Kveton, Manzil Zaheer, Mohammad Ghavamzadeh, and Craig Boutilier.  
608 Thompson sampling with a mixture prior. In *International Conference on Artificial Intelligence*  
609 *and Statistics*, pp. 7565–7586. PMLR, 2022.
- 610 Yu-Guan Hsieh, Shiva Kasiviswanathan, Branislav Kveton, and Patrick Blöbaum. Thompson sam-  
611 pling with diffusion generative prior. In Andreas Krause, Emma Brunskill, Kyunghyun Cho,  
612 Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th Inter-*  
613 *national Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning*  
614 *Research*, pp. 13434–13468. PMLR, 23–29 Jul 2023. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v202/hsieh23a.html)  
615 [press/v202/hsieh23a.html](https://proceedings.mlr.press/v202/hsieh23a.html).
- 616 Hao-Lun Hsu, Weixin Wang, Miroslav Pajic, and Pan Xu. Randomized exploration in cooperative  
617 multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 37:  
618 74617–74689, 2024.  
619
- 620 Haque Ishfaq, Qiwen Cui, Viet Nguyen, Alex Ayoub, Zhuoran Yang, Zhaoran Wang, Doina Precup,  
621 and Lin Yang. Randomized exploration in reinforcement learning with general value function  
622 approximation. In *International Conference on Machine Learning*, pp. 4607–4616. PMLR, 2021.  
623
- 624 Haque Ishfaq, Qingfeng Lan, Pan Xu, A. Rupam Mahmood, Doina Precup, Anima Anandkumar, and  
625 Kamyar Azizzadenesheli. Provable and practical: Efficient exploration in reinforcement learning  
626 via langevin monte carlo. In *The Twelfth International Conference on Learning Representations*,  
627 2024a. URL <https://openreview.net/forum?id=nfIAEJFiBZ>.
- 628 Haque Ishfaq, Yixin Tan, Yu Yang, Qingfeng Lan, Jianfeng Lu, A. Rupam Mahmood, Doina Precup,  
629 and Pan Xu. More efficient randomized exploration for reinforcement learning via approximate  
630 sampling. *Reinforcement Learning Journal*, 3:1211–1235, 2024b.  
631
- 632 Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for  
633 flexible behavior synthesis. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari,  
634 Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine*  
635 *Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 9902–9915. PMLR,  
636 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/janner22a.html>.
- 637 David Janz, Alexander Litvak, and Csaba Szepesvári. Ensemble sampling for linear bandits: small  
638 ensembles suffice. *Advances in Neural Information Processing Systems*, 37:23679–23704, 2024.  
639
- 640 Tianyuan Jin, Pan Xu, Jieming Shi, Xiaokui Xiao, and Quanquan Gu. Mots: Minimax optimal  
641 thompson sampling. In *International Conference on Machine Learning*, pp. 5074–5083. PMLR,  
642 2021.
- 643 Tianyuan Jin, Xianglin Yang, Xiaokui Xiao, and Pan Xu. Thompson sampling with less exploration  
644 is fast and optimal. In *International Conference on Machine Learning*, pp. 15239–15261. PMLR,  
645 2023.  
646
- 647 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
*arXiv:1312.6114*, 2013.

- 648 Branislav Kveton, Manzil Zaheer, Csaba Szepesvari, Lihong Li, Mohammad Ghavamzadeh, and  
649 Craig Boutilier. Randomized exploration in generalized linear bandits. In *International Confer-*  
650 *ence on Artificial Intelligence and Statistics*, pp. 2066–2076. PMLR, 2020.
- 651  
652 Branislav Kveton, Boris N. Oreshkin, Youngsuk Park, Aniket Anand Deshmukh, and Rui Song.  
653 Online posterior sampling with a diffusion prior. In *Advances in Neural Information Processing*  
654 *Systems*, 2024.
- 655 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 656  
657 Harin Lee and Min-hwan Oh. Improved regret of linear ensemble sampling. *Advances in Neural*  
658 *Information Processing Systems*, 37:92803–92831, 2024.
- 659 Ji Li and Chao Wang. Efficient diffusion posterior sampling for noisy inverse problems. *SIAM*  
660 *Journal on Imaging Sciences*, 18(2):1468–1492, 2025.
- 661  
662 Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and  
663 Mark D Plumbley. AudioLDM: Text-to-audio generation with latent diffusion models. In An-  
664 dreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan  
665 Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume  
666 202 of *Proceedings of Machine Learning Research*, pp. 21450–21474. PMLR, 23–29 Jul 2023.  
667 URL <https://proceedings.mlr.press/v202/liu23f.html>.
- 668 Xiuyuan Lu and Benjamin Van Roy. Ensemble sampling. *Advances in neural information processing*  
669 *systems*, 30, 2017.
- 670  
671 Eric Mazumdar, Aldo Pacchiano, Yian Ma, Michael Jordan, and Peter Bartlett. On approximate  
672 thompson sampling with langevin algorithms. In *international conference on machine learning*,  
673 pp. 6797–6807. PMLR, 2020.
- 674  
675 Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via  
676 posterior sampling. *Advances in Neural Information Processing Systems*, 26, 2013.
- 677  
678 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lak-  
679 shminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of Machine*  
680 *Learning Research*, 22(57):1–64, 2021.
- 681  
682 Chao Qin, Zheng Wen, Xiuyuan Lu, and Benjamin Van Roy. An analysis of ensemble sampling.  
683 *Advances in Neural Information Processing Systems*, 35:21602–21614, 2022.
- 684  
685 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
686 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*  
687 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 688  
689 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
690 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
691 text-to-image diffusion models with deep language understanding. *Advances in neural informa-*  
692 *tion processing systems*, 35:36479–36494, 2022.
- 693  
694 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
695 learning using nonequilibrium thermodynamics. In *International conference on machine learn-*  
696 *ing*, pp. 2256–2265. pmlr, 2015.
- 697  
698 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
699 Poole. Score-based generative modeling through stochastic differential equations. In *Internation-*  
700 *al Conference on Learning Representations*, 2021. URL [https://openreview.net/](https://openreview.net/forum?id=PxTIG12RRHS)  
701 [forum?id=PxTIG12RRHS](https://openreview.net/forum?id=PxTIG12RRHS).
- 702  
703 William R Thompson. On the likelihood that one unknown probability exceeds another in view of  
704 the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- 705  
706 Risto Vuorio, Shao-Hua Sun, Hexiang Hu, and Joseph J Lim. Multimodal model-agnostic meta-  
707 learning via task-aware modulation. *Advances in neural information processing systems*, 32,  
708 2019.

702 Zihui Wu, Yu Sun, Yifan Chen, Bingliang Zhang, Yisong Yue, and Katherine Bouman. Princi-  
703 pled probabilistic imaging using diffusion models as plug-and-play priors. *Advances in Neural*  
704 *Information Processing Systems*, 37:118389–118427, 2024.

705  
706 Pan Xu, Hongkai Zheng, Eric V Mazumdar, Kamyar Azizzadenesheli, and Animashree Anand-  
707 kumar. Langevin monte carlo for contextual bandits. In *International Conference on Machine*  
708 *Learning*, pp. 24830–24850. PMLR, 2022.

709  
710 Tongda Xu, Xiyan Cai, Xinjie Zhang, Xingtong Ge, Dailan He, Ming Sun, Jingjing Liu, Ya-Qin  
711 Zhang, Jian Li, and Yan Wang. Rethinking diffusion posterior sampling: From conditional score  
712 estimator to maximizing a posterior. In *The Thirteenth International Conference on Learning*  
713 *Representations*, 2025. URL <https://openreview.net/forum?id=GcvLoqOoXL>.

714  
715 Weitong Zhang, Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural thompson sampling. In  
716 *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=tkAtoZkcUnm>.

717  
718 Haoyang Zheng, Wei Deng, Christian Moya, and Guang Lin. Accelerating approximate thomp-  
719 son sampling with underdamped langevin monte carlo. In *International Conference on Artificial*  
720 *Intelligence and Statistics*, pp. 2611–2619. PMLR, 2024.

721  
722 Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with UCB-based ex-  
723 ploration. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International*  
724 *Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*,  
725 pp. 11492–11502. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/zhou20a.html>.

726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A THE USE OF LARGE LANGUAGE MODELS (LLMs)

We utilized large language models (LLMs) to assist with language editing and improve the clarity and readability of the manuscript. The authors reviewed and revised all LLM-generated suggestions and take full responsibility for the final content of this paper. The LLM’s role was strictly limited to editing and did not contribute to the core scientific ideas.

## B RELATED WORK

**Diffusion Models** Generative models have achieved strong results in modeling complex, multimodal distributions (Kingma & Welling, 2013; Goodfellow et al., 2014; Papamakarios et al., 2021; Ho et al., 2020). Among them, diffusion models are a good fit for inverse problems and decision-making because their reverse-time iterative sampling naturally supports conditioning on observed data (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021). An inverse problem is to seek unknown parameters or signals from indirect, noisy observations using a forward measurement model. In this area, Chung et al. (2023) is a cornerstone: it samples from a diffusion-model posterior by augmenting the score with the observation likelihood. Follow-up work improves stability (Xu et al., 2025), image quality (Dou & Song, 2024), learning efficiency (Li & Wang, 2025), and offers principled formulations with theoretical guarantees that avoid inaccurate approximations (Wu et al., 2024). While most of these focus on imaging, they provide useful insight for sequential decision making. In offline settings, conditional diffusion has been used to synthesize trajectories with strong results on standard benchmarks (Janner et al., 2022; Ajay et al., 2023). In online settings, Hsieh et al. (2023) studied multi-armed bandits, learn a flexible diffusion prior over task parameters, and couple it with Thompson Sampling, including variance calibration for reliable uncertainty. Kveton et al. (2024) adapted diffusion priors to linear and generalized linear contextual bandits via closed-form reverse updates. Aouali (2024) studied contextual bandits with a linear diffusion-model prior but the model reduces to a linear Gaussian case rather than a general diffusion model under this assumption. However, multi-task nonlinear contextual bandits remain underexplored. The core challenge is the absence of closed-form reverse updates for nonlinear reward models, which complicates algorithm design. Without closed forms, sampling error can be large; coupled with limited interactions per task and possible inaccuracies in the learned diffusion prior, posterior sampling can become unstable and exploration can be misled.

**Randomized Exploration** In sequential decision making problem such as bandits and reinforcement learning (RL), randomized exploration often outperforms deterministic strategies by preventing early convergence to suboptimal actions (Jin et al., 2021; 2023). Among such methods, Thompson Sampling (TS) (Thompson, 1933) is a key approach for multi-armed bandits (Agrawal & Goyal, 2017), contextual bandits (Agrawal & Goyal, 2013a), and RL (Osband et al., 2013). Unlike Upper Confidence Bound (UCB) algorithms, which rely on deterministic confidence intervals (Chu et al., 2011; Lattimore & Szepesvári, 2020), TS samples from posterior distributions, enabling robust exploration (Chapelle & Li, 2011). However, exact posterior sampling in TS is computationally intensive, particularly for non-conjugate priors. To address this, approximate sampling methods such as Langevin Monte Carlo (LMC) (Xu et al., 2022; Hsu et al., 2024), Stochastic gradient Langevin dynamics (SGLD) variants (Mazumdar et al., 2020; Zheng et al., 2024) and variational inference (Clavier et al., 2024) have been developed and applied to various problem settings, including multi-armed bandits with non-conjugate or highly nonlinear rewards, nonlinear contextual bandits and RL (Ishfaq et al., 2024a;b; Hsu et al., 2024). Perturb History Exploration (PHE) improves efficiency by perturbing historical data to approximate posterior sampling, making it applicable to complex reward distributions (Kveton et al., 2020; Ishfaq et al., 2021). Ensemble sampling keeps a small set of independently perturbed model replicas and selects actions using a randomly chosen replica (Lu & Van Roy, 2017; Qin et al., 2022; Janz et al., 2024; Lee & Oh, 2024).

## C MORE DETAILS ON DIFFUSION POSTERIOR SAMPLING ALGORITHMS

### C.1 DETAILED ALGORITHM INTERPRETATION FOR DLTS

We run DLTS for  $T$  rounds to solve the new bandit task ( $N + 1$ ). At each time  $t$ , with a pretrained diffusion model for unconditional sampling, we modify the reverse diffusion update to perform conditional sampling by injecting the history at every reverse step. We implement this conditional sampling by applying approximate sampling via Langevin Monte Carlo (LMC). After the reverse process ends, we take the final state  $\theta_0$  (i.e.,  $\tilde{\theta}_t$ ) as the parameter sample for arm selection.

Next, we provide a detailed derivation of the conditional sampling used in the reverse diffusion update in Algorithm 1. After training the diffusion denoiser  $\varepsilon_{\phi^*}(\cdot, \cdot)$  (see Appendix G for more details), we obtain the diffusion model parameter  $(\mu_\ell, \Sigma_\ell)_{\ell \in [L+1]}$ <sup>2</sup> for unconditional sampling in each reverse step, which is formulated as follow,

$$p(\theta_{\ell-1}|\theta_\ell) = \mathcal{N}(\theta_{\ell-1}; \mu_\ell, \Sigma_\ell), \quad (\text{C.1})$$

where  $\mu_\ell = \frac{1}{\sqrt{\alpha_\ell}}(\theta_\ell - \frac{1-\alpha_\ell}{\sqrt{1-\alpha_\ell}}\varepsilon_{\phi^*}(\theta_\ell, \ell))$ ,  $\Sigma_\ell = \frac{1-\alpha_{\ell-1}}{1-\alpha_\ell}\beta_\ell\mathbf{I}$ . Note that in (3.2), we can factorize  $p(\mathcal{H}_t|\theta_{\ell-1})$  as follows,

$$p(\mathcal{H}_t|\theta_{\ell-1}) = \int_{\theta_0} p(\mathcal{H}_t|\theta_0, \theta_{\ell-1})p(\theta_0|\theta_{\ell-1})d\theta_0 = \int_{\theta_0} p(\mathcal{H}_t|\theta_0)p(\theta_0|\theta_{\ell-1})d\theta_0,$$

where the second equality comes from that  $\mathcal{H}_t$  and  $\theta_{\ell-1}$  are conditionally independent on  $\theta_0$ . To derive the conditional diffusion reverse update, we need to approximate the following term  $\int_{\theta_0} p(\mathcal{H}_t|\theta_0)p(\theta_0|\theta_{\ell-1})d\theta_0$ . Based on Kveton et al. (2024), we approximate it by

$$\int_{\theta_0} p(\mathcal{H}_t|\theta_0)p(\theta_0|\theta_{\ell-1})d\theta_0 \approx p(\mathcal{H}_t|\theta_{\ell-1}/\sqrt{\alpha_{\ell-1}}). \quad (\text{C.2})$$

The motivation of approximation in (C.2) is as follows. Note that  $\theta_\ell = \sqrt{\alpha_\ell}\theta_0 + \sqrt{1-\alpha_\ell}\tilde{\xi}_\ell$ , where  $\tilde{\xi}_\ell \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is standard Gaussian noise. Rearranging gives  $\theta_0 = (\theta_\ell - \sqrt{1-\alpha_\ell}\tilde{\xi}_\ell)/\sqrt{\alpha_\ell}$ , so  $\theta_0$  can be viewed as a random variable with mean  $\theta_\ell/\sqrt{\alpha_\ell}$ . Replacing  $\theta_0$  in the LHS of (C.2) by this mean yields the approximation.

Note that at step  $t$ , given a pre-determined  $f$ , the nonlinear reward model is defined as

$$y_t = f(\theta^*; \mathbf{x}_t) + \varepsilon_t.$$

When we assume Gaussian noise  $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$ , then the likelihood becomes

$$p(\mathcal{H}_t|\theta) \propto \exp\left(-\sigma^{-2}\sum_{i=1}^t (f(\theta; \mathbf{x}_i) - y_i)^2\right) = \exp(-L_t(\theta)), \quad (\text{C.3})$$

where  $L_t(\theta) = \sigma^{-2}\sum_{i=1}^t (f(\theta; \mathbf{x}_i) - y_i)^2$ . Therefore, based on (3.2), (C.2) and (C.3), we approximate the posterior in each reverse step for conditional sampling as follows,

$$p(\theta_{\ell-1}|\theta_\ell, \mathcal{H}_t) \propto \exp(-L_t(\theta_{\ell-1}/\sqrt{\alpha_{\ell-1}})) \cdot \mathcal{N}(\theta_{\ell-1}; \mu_\ell, \Sigma_\ell). \quad (\text{C.4})$$

To sample from  $p(\theta_{\ell-1}|\theta_\ell, \mathcal{H}_t)$ , we apply Langevin Monte Carlo for approximate sampling. Specifically, at bandit round  $t$  and diffusion reverse step  $\ell$ , we iteratively conduct the Langevin update,

$$\theta_{\ell, k+1} = \theta_{\ell, k} - \eta_\ell \left[ \nabla_{\theta} L_t(\theta_{\ell, k}/\sqrt{\alpha_{\ell-1}}) + \Sigma_\ell^{-1}(\theta_{\ell, k} - \mu_\ell) \right] + \sqrt{2\eta_\ell \zeta_\ell^{-1}} \xi_k, \quad (\text{C.5})$$

where  $\xi_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\eta_\ell$  is step-size,  $\zeta_\ell$  is temperature and (C.5) is initialized at  $\theta_{\ell, 0} = \mu_\ell$  and after  $K_\ell$  iterations return  $\theta_{\ell, K_\ell}$  to be the diffusion reverse process sample  $\theta_{\ell-1}$ .

Specially, when we only conduct one LMC update, based on (C.5), we have the reverse update

$$\theta_{\ell-1} \leftarrow \left(1 - \frac{(1-\alpha_\ell)\eta_\ell}{(1-\alpha_{\ell-1})\beta_\ell}\right) \theta_\ell + \frac{(1-\alpha_\ell)\eta_\ell}{(1-\alpha_{\ell-1})\beta_\ell} \mu_\ell - \eta_\ell \nabla_{\theta} L_t(\theta_\ell/\sqrt{\alpha_{\ell-1}}) + \sqrt{2\eta_\ell \zeta_\ell^{-1}} \xi. \quad (\text{C.6})$$

Note that (C.6) aligns with our unified update (3.1). This confirms the rationality of our framework and the essence of diffusion-based posterior sampling. Intuitively, DLTS forms a conditional posterior by incorporating the history into the diffusion prior at each reverse step, then draws samples from it via approximate sampling.

<sup>2</sup>We define  $(\mu_{L+1}, \Sigma_{L+1}) = (\mathbf{0}, \mathbf{I})$  for notation simplicity.

## C.2 DETAILED ALGORITHM INTERPRETATION FOR DPSG

We run DPSG for  $T$  rounds to solve the new bandit task ( $N + 1$ ). At each time  $t$ , after using a pretrained diffusion model for unconditional sampling, we then turn the reverse diffusion into conditional sampling by adding a likelihood drift at every reverse step: at level  $\ell$  we form the Tweedie proxy  $\hat{\theta}_0$  from the current state and apply a single gradient update inside the reverse transition. The history interaction data enter as one likelihood-score step per level. After the reverse process ends, we take the final state  $\theta_0$  (i.e.,  $\hat{\theta}_t$ ) as the parameter sample for arm selection.

To derive diffusion reverse update (Line 8 in Algorithm 2), we follow Chung et al. (2023) and decompose the conditional score into two terms,

$$\nabla_{\theta} \log p_s(\theta_s | y) = \underbrace{\nabla_{\theta} \log p_s(\theta_s)}_{\text{unconditional score}} + \underbrace{\nabla_{\theta} \log p_s(y | \theta_s)}_{\text{likelihood score}}.$$

For the first term, we can use pretrained diffusion model score network  $s_{\psi^*}$  to approximate the unconditional score. For the second term, we can use Tweedie’s formula to provide a tractable approximation for  $p_s(y | \theta_s)$  such that one can use the surrogate function for approximate posterior sampling, which is formulated as,

$$\nabla_{\theta} \log p_s(y | \theta_s) \simeq \nabla_{\theta} \log p_s(y | \hat{\theta}_0), \quad (\text{C.7})$$

where  $\theta$  is the conditional posterior mean of  $\theta_0$ , which have the following close-form expression derived by Tweedie’s formula (Lemma F.2), that is,

$$\hat{\theta}_0 = \mathbb{E}_{\theta_0 \sim p(\theta_0 | \theta_s)}[\theta_0] = \frac{1}{\sqrt{\bar{\alpha}_s}} (\theta_s + (1 - \bar{\alpha}_s) \nabla_{\theta} \log p_s(\theta_s)).$$

In practice, given learned score network  $s_{\psi^*}$ , we can further approximate the posterior mean as

$$\hat{\theta}_0(\theta_s, s) \simeq \frac{1}{\sqrt{\bar{\alpha}_s}} (\theta_s + (1 - \bar{\alpha}_s) s_{\psi^*}(\theta_s, s)). \quad (\text{C.8})$$

At timestep  $t$ , the nonlinear reward model is defined as  $y_t = f(\theta^*; \mathbf{x}_t) + \varepsilon_t$ , where  $f$  is a pre-determined nonlinear function. Same as (C.3), when we assume  $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$ , then the likelihood becomes  $p(\mathcal{H}_t | \theta) \propto \exp(-\sigma^{-2} \sum_{i=1}^t (f(\theta; \mathbf{x}_i) - y_i)^2) = \exp(-L_t(\theta))$  where  $L_t(\theta) = \sigma^{-2} \sum_{i=1}^t (f(\theta; \mathbf{x}_i) - y_i)^2$ . Based on (C.7) and (C.8), then we can have the following tractable approximation of likelihood score,

$$\nabla_{\theta} \log p_s(\mathcal{H}_t | \theta_s) \simeq -\nabla_{\theta} L_t(\hat{\theta}_0(\theta_s, s)).$$

Therefore, we can further have the approximation of the conditional score of reverse dynamics,

$$\nabla_{\theta} \log p_s(\theta_s | \mathcal{H}_t) \simeq s_{\psi^*}(\theta_s, s) - \nabla_{\theta} L_t(\hat{\theta}_0(\theta_s, s)). \quad (\text{C.9})$$

In practice, to obtain a discrete-time algorithm we use DDPM ancestral sampling (Ho et al., 2020) to implement conditional sampling via the conditional score (C.9). The reverse update is done by two steps. First, we use the unconditional score  $s_{\psi^*}(\theta_{\ell}, \ell)$  through Tweedie’s estimate  $\hat{\theta}_0(\theta_{\ell}, \ell)$  to form the Gaussian posterior mean, we then sample

$$\theta'_{\ell-1} \leftarrow \frac{\sqrt{\alpha_{\ell}}(1 - \bar{\alpha}_{\ell-1})}{1 - \bar{\alpha}_{\ell}} \theta_{\ell} + \frac{\sqrt{\bar{\alpha}_{\ell-1}} \beta_{\ell}}{1 - \bar{\alpha}_{\ell}} \hat{\theta}_0(\theta_{\ell}, \ell) + \mathbf{z}_{\ell}, \quad \mathbf{z}_{\ell} \sim \mathcal{N}(\mathbf{0}, \beta_{\ell} \mathbf{I}), \quad (\text{C.10})$$

Based on (C.9), the second step is that we make this update conditional by subsequently adding one likelihood term. In practice, after sampling  $\theta'_{\ell-1}$ , we take a gradient step using the likelihood  $L_t$  at the Tweedie estimate for guidance, that is

$$\theta_{\ell-1} \leftarrow \theta'_{\ell-1} - \eta_{\ell} \nabla_{\theta} L_t(\hat{\theta}_0(\theta_{\ell}, \ell)). \quad (\text{C.11})$$

When we merge (C.10) and (C.11) together, we obtain line 8 in Algorithm 2,

$$\theta_{\ell-1} \leftarrow \frac{1}{\sqrt{\alpha_{\ell}}} \theta_{\ell} + \frac{\beta_{\ell}}{\sqrt{\alpha_{\ell}}} s_{\psi^*}(\theta_{\ell}, \ell) - \eta_{\ell} \nabla_{\theta} L_t(\hat{\theta}_0(\theta_{\ell}, \ell)) + \mathbf{z}_{\ell}. \quad (\text{C.12})$$

Note that (C.12) aligns with our unified update (3.1). This confirms the rationality of our framework and the essence of diffusion-based posterior sampling. However, from the intuition perspective, DPSG is to use history information as guidance after unconditional update based on the diffusion prior. This is different from DLTS, which directly integrates history information into the unconditional update based on the diffusion prior to achieve conditional sampling.

### 918 C.3 ALGORITHM INTERPRETATION FOR DPSG-MP

919  
920 DPSG is simple and fast, but it relies on the likelihood-score at the Tweedie estimate. A recent work  
921 Xu et al. (2025) observes that DPS has the properties of high bias and low diversity, thus behaving  
922 like an implicit, but unstable, MAP estimator. Therefore, multi-step projection is used to solve this  
923 MAP optimization problem. Xu et al. (2025) proposed Diffusion Maximize a Posterior (DMAP),  
924 which keeps the standard reverse diffusion prior step and replaces the single likelihood drift with  
925 several small gradient-ascent steps on the log posterior at each noise level, projecting after each step  
926 onto a sphere around the unconditional mean to stay consistent with the diffusion dynamics.

927 Therefore, we introduce DPSG with Multi-step Projection (DPSG-MP) shown in Algorithm 3. Com-  
928 pared to Algorithm 2, it replaces the single correction with a short inner loop of gradient ascent on  
929 the likelihood score, and after each step projects the iterate onto the sphere where the reverse tran-  
930 sition  $p_{\psi^*}(\theta_{\ell-1}|\theta_\ell)$  concentrates around. This avoids the need for an accurate conditional score,  
931 reduces drift bias, and yields stable, data-consistent parameter estimates for arm selection.

---

#### 932 **Algorithm 3** DPSG with Multi-step Projection (DPSG-MP)

---

933 **Input:** score network  $s_{\psi^*}(\cdot, \cdot)$ , noise schedule  $\{\beta_\ell\}$ , trajectory history  $\mathcal{H}_t$ , learning rate  $\{\eta_{\ell,k}\}$ .

934  
935 1:  $\alpha_\ell \leftarrow 1 - \beta_\ell$  and  $\bar{\alpha}_\ell \leftarrow \prod_{j=1}^\ell \alpha_j$  for all  $\ell$ .  
936 2: **for**  $t = 1, 2, \dots, T$  **do**  
937 3:   Receive contextual vector  $\{\mathbf{x}_t(a)\}_{a \in \mathcal{A}}$ .  
938 4:    $\boldsymbol{\theta}_L \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .  
939 5:   **for**  $\ell = L, L-1, \dots, 1$  **do**  
940 6:      $\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell) \leftarrow \frac{1}{\sqrt{\bar{\alpha}_\ell}}(\boldsymbol{\theta}_\ell + (1 - \bar{\alpha}_\ell)s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell))$ .  
941 7:      $\mathbf{z}_\ell \sim \mathcal{N}(\mathbf{0}, \beta_\ell \mathbf{I})$ ,  $r_\ell \leftarrow \|\mathbf{z}_\ell\|_2$ .  
942 8:      $\mathbf{m}_\ell \leftarrow \frac{\sqrt{\bar{\alpha}_\ell(1-\bar{\alpha}_\ell-1)}}{1-\bar{\alpha}_\ell}\boldsymbol{\theta}_\ell + \frac{\sqrt{\bar{\alpha}_\ell-1}\beta_\ell}{1-\bar{\alpha}_\ell}\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell)$ .  
943 9:      $\boldsymbol{\theta}_{\ell-1,0} \leftarrow \mathbf{m}_\ell + \mathbf{z}_\ell$ .  
944 10:     **for**  $k = 1, 2, \dots, K_\ell$  **do**  
945 11:        $\boldsymbol{\theta}'_{\ell-1,k} \leftarrow \boldsymbol{\theta}_{\ell-1,k-1} - \eta_{\ell,k} \nabla_{\boldsymbol{\theta}} L_t(\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell))$ .  
946 12:        $\boldsymbol{\theta}_{\ell-1,k} \leftarrow \mathbf{m}_\ell + r_\ell \frac{\boldsymbol{\theta}'_{\ell-1,k} - \mathbf{m}_\ell}{\|\boldsymbol{\theta}'_{\ell-1,k} - \mathbf{m}_\ell\|_2}$ .  
947 13:     **end for**  
948 14:      $\boldsymbol{\theta}_{\ell-1} \leftarrow \boldsymbol{\theta}_{\ell-1, K_\ell}$ .  
949 15:   **end for**  
950 16:   Let  $\tilde{\boldsymbol{\theta}}_t \leftarrow \boldsymbol{\theta}_0$ , take action  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x}_t(a)} f(\tilde{\boldsymbol{\theta}}_t; \mathbf{x}_t(a))$ .  
951 17:   Receive reward  $y_t$  and update history  $\mathcal{H}_{t+1} = \{(\mathbf{x}_i, y_i)\}_{i=1}^t$ .  
952 18: **end for**

---

## 953 D PROOF IN SECTION 4

### 954 D.1 PROOF OF THEOREM 4.4

955 *Proof of Theorem 4.4.* We show that, under the ideal ODTS assumptions, the ODTS marginal pos-  
956 terior over the clean parameter is same as the OTS posterior.

957 Let the diffusion-path variable be  $\boldsymbol{\theta}_{0:L} = (\boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_L)$ , where  $\boldsymbol{\theta}_0$  is the clean parameter used for  
958 decisions and  $\boldsymbol{\theta}_L$  is the noisiest level. The diffusion prior has the the Markov factorization

$$959 p(\boldsymbol{\theta}_{0:L}) = p(\boldsymbol{\theta}_L) \prod_{\ell=1}^L p(\boldsymbol{\theta}_{\ell-1}|\boldsymbol{\theta}_\ell), \quad (\text{D.1})$$

960 so that its marginal posterior is as follows,

$$961 p_0(\boldsymbol{\theta}_0) = \int p(\boldsymbol{\theta}_{0:L}) d\boldsymbol{\theta}_{1:L}, \quad (\text{D.2})$$

962 which is the true prior used by OTS according to the condition. Note that the history  $\mathcal{H}_t$  depends  
963 only on  $\boldsymbol{\theta}_0$ , then we have

$$964 p(\mathcal{H}_t|\boldsymbol{\theta}_{0:L}) = p(\mathcal{H}_t|\boldsymbol{\theta}_0). \quad (\text{D.3})$$

972 ODTs is assumed ideal: (i) the learned score/denoiser equals the true score so that (D.1) holds for  
 973 the implemented prior, and (ii) each per-level conditional reverse transition kernel is implemented  
 974 exactly, thus representing the true joint posterior over the path.

975 By Bayes' rule and (D.3), the joint posterior over the path satisfies

$$976 \tilde{p}_t(\boldsymbol{\theta}_{0:L}) \equiv p(\boldsymbol{\theta}_{0:L}|\mathcal{H}_t) \propto p(\mathcal{H}_t|\boldsymbol{\theta}_0)p(\boldsymbol{\theta}_{0:L}). \quad (\text{D.4})$$

977 Integrating (D.4) over the noisy layers  $\boldsymbol{\theta}_{1:L}$  and using (D.2) gives

$$\begin{aligned} 978 \int p(\boldsymbol{\theta}_{0:L}|\mathcal{H}_t)d\boldsymbol{\theta}_{1:L} &= \frac{1}{p(\mathcal{H}_t)} \int p(\mathcal{H}_t|\boldsymbol{\theta}_0)p(\boldsymbol{\theta}_{0:L})d\boldsymbol{\theta}_{1:L} \\ 979 &= \frac{1}{p(\mathcal{H}_t)} p(\mathcal{H}_t|\boldsymbol{\theta}_0) \underbrace{\int p(\boldsymbol{\theta}_{0:L})d\boldsymbol{\theta}_{1:L}}_{=p_0(\boldsymbol{\theta}_0)} \\ 980 &= \frac{p(\mathcal{H}_t|\boldsymbol{\theta}_0)p_0(\boldsymbol{\theta}_0)}{p(\mathcal{H}_t)} = p(\boldsymbol{\theta}_0|\mathcal{H}_t). \end{aligned} \quad (\text{D.5})$$

981 Note that  $\tilde{p}_t(\boldsymbol{\theta}) = \int \tilde{p}_t(\boldsymbol{\theta}_{0:L})d\boldsymbol{\theta}_{1:L} = \int p(\boldsymbol{\theta}_{0:L}|\mathcal{H}_t)d\boldsymbol{\theta}_{1:L}$ , then we have

$$982 \tilde{p}_t(\boldsymbol{\theta}) = p(\boldsymbol{\theta}|\mathcal{H}_t) \equiv p_t(\boldsymbol{\theta}). \quad (\text{D.6})$$

983 Therefore the ODTs posterior marginal  $\tilde{p}_t(\boldsymbol{\theta})$  and the OTS posterior  $p_t(\boldsymbol{\theta})$  are same. Note that both  
 984 oracle algorithms use greedy policy to select arm, thus share the same arm selection distribution.  $\square$

## 985 D.2 PROOF OF THEOREM 4.5

986 *Proof of Theorem 4.5.* Based on the definition of cumulative expected regret, we focus on the per-  
 987 round expected regret in round  $t$ ,

$$988 r_t = f(\boldsymbol{\theta}^*; \mathbf{x}^*) - f(\boldsymbol{\theta}^*; \mathbf{x}_t).$$

989 We then focus on the regret difference between ODTs and OTS,

$$\begin{aligned} 990 \mathbb{E}[r_t^{\text{ODTS}} - r_t^{\text{OTS}}] &= \mathbb{E}[f(\boldsymbol{\theta}^*; \mathbf{x}^*) - f(\boldsymbol{\theta}^*; \mathbf{x}_t^{\text{ODTS}})] - \mathbb{E}[f(\boldsymbol{\theta}^*; \mathbf{x}^*) - f(\boldsymbol{\theta}^*; \mathbf{x}_t^{\text{OTS}})] \\ 991 &= \mathbb{E}_{\boldsymbol{\theta} \sim p_t}[f(\boldsymbol{\theta}^*; \pi(\boldsymbol{\theta}))] - \mathbb{E}_{\boldsymbol{\theta} \sim \tilde{p}_t}[f(\boldsymbol{\theta}^*; \pi(\boldsymbol{\theta}))] \\ 992 &\leq 2f_{\max} \text{TV}(p_t, \tilde{p}_t), \end{aligned}$$

993 where the inequality comes from the dual definition of TV divergence.  $\text{TV}(p_t, \tilde{p}_t)$  contains the score  
 994 estimation error, we then further bound this term based on Assumption 4.3.

995 **Bounding  $\text{TV}(p_t, \tilde{p}_t)$ .** At reverse level  $\ell$ , the true and learned reverse Gaussian transition share  
 996 covariance  $\Sigma_\ell$  while their means differ because of the accuracy of denoiser. Denote that  $\mu_\ell^*$  and  $\varepsilon^*$   
 997 is the mean and noise of the true transition, then we have

$$998 \|\boldsymbol{\mu}_\ell - \boldsymbol{\mu}_\ell^*\|_2 \leq \frac{1 - \alpha_\ell}{\sqrt{\alpha_\ell \bar{\alpha}_\ell}} \|\varepsilon_{\phi^*} - \varepsilon^*\|_2 = c_\ell \|s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell) - \nabla_{\boldsymbol{\theta}} \log p_\ell(\boldsymbol{\theta}_\ell)\|_2 \leq c_\ell \epsilon_{\text{score}},$$

999 where coefficient  $c_\ell$  is determined by noise schedule, the second equality holds due to Tweedie's  
 1000 formula, the last inequality holds because of Assumption 4.3. For Gaussian transition with equal  
 1001 covariance, we have

$$1002 \text{TV}(\mathcal{N}(\boldsymbol{\mu}_\ell^*, \Sigma_\ell), \mathcal{N}(\boldsymbol{\mu}_\ell, \Sigma_\ell)) \leq \kappa_\ell \epsilon_{\text{score}},$$

1003 with  $\kappa_\ell = \frac{c_\ell}{2\sqrt{\lambda_{\min}(\Sigma_\ell)}} = \frac{1 - \alpha_\ell}{2\sqrt{\alpha_\ell} \sqrt{\lambda_{\min}(\Sigma_\ell)}}$ . Propagating this perturbation through the reverse Markov  
 1004 chain yields

$$1005 \text{TV}(p_t, \tilde{p}_t) \leq \sum_{\ell=1}^L \kappa_\ell \epsilon_{\text{score}}.$$

1006 Therefore, we have the score estimation error,

$$1007 \Delta_t^{\text{Score}} = 2f_{\max} \sum_{\ell=1}^L \kappa_\ell \epsilon_{\text{score}}.$$

1008 This completes the proof.  $\square$

## E INSTANTIATION UNDER SPECIFIC SETTINGS

In this section, we describe the bandits settings where we instantiate our proposed algorithms. The general frameworks for DLTS and DPSG can be readily adapted to various specific models by defining the reward function  $f(\boldsymbol{\theta}; \mathbf{x})$  and the corresponding loss function  $L_t(\boldsymbol{\theta})$ . The gradient of this loss,  $\nabla_{\boldsymbol{\theta}} L_t(\boldsymbol{\theta})$ , is the crucial model-dependent component that is plugged into the reverse diffusion update rules for both DLTS and DPSG. We now provide the concrete formulations for three widely-studied settings: linear, generalized linear, and neural contextual bandits.

**Linear Bandit** We first consider the linear contextual bandit settings. In linear contextual bandits, the reward is given by a linear function of the context and an unknown parameter  $\boldsymbol{\theta}^* \in \mathbb{R}^d$  plus some noise, that is,  $y_t = \boldsymbol{\theta}^{*\top} \mathbf{x}_t + \varepsilon_t$ , where  $\mathbf{x}_t \in \mathbb{R}^d$  is the context at time  $t$  and  $\varepsilon_t$  is the noise term. We assume that  $\varepsilon_t$  is i.i.d. Gaussian noise with mean zero and variance  $\sigma^2$ . The goal of the agent is to learn the unknown parameter  $\boldsymbol{\theta}^*$  and select actions that maximize the cumulative reward over time. The reward function is defined as follows

$$f(\boldsymbol{\theta}; \mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}.$$

In our proposed algorithms, we define the loss function over the observed history  $\mathcal{H}_t$  as follows

$$L_t(\boldsymbol{\theta}) = \sigma^{-2} \sum_{i=1}^t (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i)^2,$$

where  $\sigma^2$  is the variance,  $\mathbf{x}_i$  is the context at time  $i$  and  $y_i$  is the corresponding reward. For the linear model, the crucial gradient term  $\nabla_{\boldsymbol{\theta}} L_t(\boldsymbol{\theta})$  required for our diffusion updates has a convenient closed-form expression. It can be computed as:

$$\nabla_{\boldsymbol{\theta}} L_t(\boldsymbol{\theta}) = 2\sigma^{-2} \sum_{i=1}^t (\boldsymbol{\theta}^\top \mathbf{x}_i - y_i) \mathbf{x}_i.$$

For DLTS, the Langevin Monte Carlo update from (C.6), which is iterated within each reverse diffusion step  $\ell$ . We have

$$\boldsymbol{\theta}_{\ell-1} \leftarrow \left(1 - \frac{(1 - \bar{\alpha}_\ell)\eta_\ell}{(1 - \bar{\alpha}_{\ell-1})\beta_\ell}\right) \boldsymbol{\theta}_\ell + \frac{(1 - \bar{\alpha}_\ell)\eta_\ell}{(1 - \bar{\alpha}_{\ell-1})\beta_\ell} \boldsymbol{\mu}_\ell - 2\eta_\ell \sigma^{-2} \sum_{i=1}^t \left(\frac{\boldsymbol{\theta}_\ell^\top \mathbf{x}_i}{\sqrt{\bar{\alpha}_{\ell-1}}} - y_i\right) \frac{\mathbf{x}_i}{\sqrt{\bar{\alpha}_{\ell-1}}} + \sqrt{2\eta_\ell \zeta_\ell^{-1}} \boldsymbol{\xi}.$$

For DPSG, the reverse update from Line 8 in Algorithm 2 is instantiated as:

$$\boldsymbol{\theta}_{\ell-1} \leftarrow \frac{1}{\sqrt{\alpha_\ell}} \boldsymbol{\theta}_\ell + \frac{\beta_\ell}{\sqrt{\alpha_\ell}} s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell) - 2\eta_\ell \sigma^{-2} \sum_{i=1}^t ((\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell))^\top \mathbf{x}_i - y_i) \mathbf{x}_i + \mathbf{z}_\ell.$$

**Generalized Linear Bandit** Then we extend our proposed algorithms to generalized linear bandit settings. In generalized linear bandits, the reward is given by a generalized linear function of the context and an unknown parameter  $\boldsymbol{\theta}^* \in \mathbb{R}^d$  plus some noise, that is,  $y_t = g(\boldsymbol{\theta}^{*\top} \mathbf{x}_t) + \varepsilon_t$ , where  $\mathbf{x}_t \in \mathbb{R}^d$  is the context at time  $t$ ,  $g$  is a known link function and  $\varepsilon_t$  is the noise term. We assume that  $\varepsilon_t$  is i.i.d. Gaussian noise with mean zero and variance  $\sigma^2$ . The link function  $g$  is a known function that maps the linear combination of the context and parameter to the expected reward. We have the following formulations for the reward function

$$f(\boldsymbol{\theta}; \mathbf{x}) = g(\boldsymbol{\theta}^\top \mathbf{x}).$$

The goal of the agent is to learn the unknown parameter  $\boldsymbol{\theta}^*$  and select actions that maximize the cumulative reward over time. In our proposed algorithms, we define the loss function over the observed history  $\mathcal{H}_t$  as follows

$$L_t(\boldsymbol{\theta}) = \sigma^{-2} \sum_{i=1}^t (g(\boldsymbol{\theta}^\top \mathbf{x}_i) - y_i)^2,$$

where  $\sigma^2$  is the variance,  $\mathbf{x}_i$  is the context at time  $i$  and  $y_i$  is the corresponding reward. We use the logistic function as the example in our experiments. The logistic function is defined as  $g(z) = 1/(1 + e^{-z})$ . For the gradient, we have:

$$\nabla_{\boldsymbol{\theta}} L_t(\boldsymbol{\theta}) = 2\sigma^{-2} \sum_{i=1}^t (g(\boldsymbol{\theta}^\top \mathbf{x}_i) - y_i) g'(\boldsymbol{\theta}^\top \mathbf{x}_i) \mathbf{x}_i,$$

where  $g'(\cdot)$  is the derivative of the link function with respect to its input.

For DLTS, the single-step Langevin update (C.6) becomes:

$$\begin{aligned} \boldsymbol{\theta}_{\ell-1} \leftarrow & \left(1 - \frac{(1 - \bar{\alpha}_\ell)\eta_\ell}{(1 - \bar{\alpha}_{\ell-1})\beta_\ell}\right) \boldsymbol{\theta}_\ell + \frac{(1 - \bar{\alpha}_\ell)\eta_\ell}{(1 - \bar{\alpha}_{\ell-1})\beta_\ell} \boldsymbol{\mu}_\ell \\ & - 2\eta_\ell \sigma^{-2} \sum_{i=1}^t \left(g\left(\frac{\boldsymbol{\theta}_\ell^\top \mathbf{x}_i}{\sqrt{\bar{\alpha}_{\ell-1}}}\right) - y_i\right) g'\left(\frac{\boldsymbol{\theta}_\ell^\top \mathbf{x}_i}{\sqrt{\bar{\alpha}_{\ell-1}}}\right) \mathbf{x}_i + \sqrt{2\eta_\ell \zeta_\ell^{-1}} \boldsymbol{\xi}. \end{aligned}$$

For DPSG, the reverse update (Line 8 in Algorithm 2) is instantiated as:

$$\begin{aligned} \boldsymbol{\theta}_{\ell-1} \leftarrow & \frac{1}{\sqrt{\alpha_\ell}} \boldsymbol{\theta}_\ell + \frac{\beta_\ell}{\sqrt{\alpha_\ell}} s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell) \\ & - 2\eta_\ell \sigma^{-2} \sum_{i=1}^t (g((\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell))^\top \mathbf{x}_i) - y_i) g'((\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell))^\top \mathbf{x}_i) \mathbf{x}_i + \mathbf{z}_\ell. \end{aligned}$$

**Neural Bandit** Finally, we consider the neural contextual bandit settings. In neural contextual bandits, the reward is given by a neural network function of the context and an unknown parameter  $\boldsymbol{\theta}^* \in \mathbb{R}^p$  plus some noise, that is,  $y_t = f(\boldsymbol{\theta}^*; \mathbf{x}_t) + \varepsilon_t$ , where  $\mathbf{x}_t \in \mathbb{R}^d$  is the context at time  $t$  and  $\varepsilon_t$  is the noise term. We assume that  $\varepsilon_t$  is i.i.d. Gaussian noise with mean zero and variance  $\sigma^2$ . The goal of the agent is to learn the unknown parameter  $\boldsymbol{\theta}^*$  and select actions that maximize the cumulative reward over time. The reward function is defined as follows

$$f(\boldsymbol{\theta}; \mathbf{x}) = \sqrt{m} \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 \mathbf{x}))),$$

where  $\sigma(x) = \max\{x, 0\}$  is the rectified linear unit (ReLU) activation function,  $\mathbf{W}_1 \in \mathbb{R}^{m \times d}$ ,  $\mathbf{W}_i \in \mathbb{R}^{m \times m}$ ,  $2 \leq i \leq L-1$ ,  $\mathbf{W}_L \in \mathbb{R}^{m \times 1}$ , and  $\boldsymbol{\theta} = [\text{vec}(\mathbf{W}_1)^\top, \dots, \text{vec}(\mathbf{W}_L)^\top]^\top \in \mathbb{R}^p$  with  $p = m + md + m^2(L-1)$ . In our proposed algorithms, we define the loss function over the observed history  $\mathcal{H}_t$  as follows

$$L_t(\boldsymbol{\theta}) = \sigma^{-2} \sum_{i=1}^t (f(\boldsymbol{\theta}; \mathbf{x}_i) - y_i)^2,$$

where  $\sigma^2$  is the variance,  $\mathbf{x}_i$  is the context at time  $i$  and  $y_i$  is the corresponding reward. The gradient of this loss function is required for our diffusion updates. We obtain:

$$\nabla_{\boldsymbol{\theta}} L_t(\boldsymbol{\theta}) = 2\sigma^{-2} \sum_{i=1}^t (f(\boldsymbol{\theta}; \mathbf{x}_i) - y_i) \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}; \mathbf{x}_i).$$

For DLTS, the single-step Langevin update (C.6) becomes:

$$\begin{aligned} \boldsymbol{\theta}_{\ell-1} \leftarrow & \left(1 - \frac{(1 - \bar{\alpha}_\ell)\eta_\ell}{(1 - \bar{\alpha}_{\ell-1})\beta_\ell}\right) \boldsymbol{\theta}_\ell + \frac{(1 - \bar{\alpha}_\ell)\eta_\ell}{(1 - \bar{\alpha}_{\ell-1})\beta_\ell} \boldsymbol{\mu}_\ell \\ & - 2\eta_\ell \sigma^{-2} \sum_{i=1}^t \left(f\left(\frac{\boldsymbol{\theta}_\ell}{\sqrt{\bar{\alpha}_{\ell-1}}}; \mathbf{x}_i\right) - y_i\right) \nabla_{\boldsymbol{\theta}} f\left(\frac{\boldsymbol{\theta}_\ell}{\sqrt{\bar{\alpha}_{\ell-1}}}; \mathbf{x}_i\right) + \sqrt{2\eta_\ell \zeta_\ell^{-1}} \boldsymbol{\xi}. \end{aligned}$$

For DPSG, the reverse update (Line 8 in Algorithm 2) is instantiated as:

$$\boldsymbol{\theta}_{\ell-1} \leftarrow \frac{1}{\sqrt{\alpha_\ell}} \boldsymbol{\theta}_\ell + \frac{\beta_\ell}{\sqrt{\alpha_\ell}} s_{\psi^*}(\boldsymbol{\theta}_\ell, \ell) - 2\eta_\ell \sigma^{-2} \sum_{i=1}^t (f(\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell); \mathbf{x}_i) - y_i) \nabla_{\boldsymbol{\theta}} f(\hat{\boldsymbol{\theta}}_0(\boldsymbol{\theta}_\ell, \ell); \mathbf{x}_i) + \mathbf{z}_\ell.$$

## F AUXILIARY LEMMAS

The following lemma provides the conditional probability distributions for diffusion reverse process based on history  $\mathcal{H}_t = \{(\mathbf{x}_i, y_i)\}_{i=1}^{t-1}$  at timestep  $t$ .

**Lemma F.1** ((Kveton et al., 2024)). *For probability measure  $p$  over the reverse process, we have*

$$\begin{aligned} p(\boldsymbol{\theta}_L | \mathcal{H}_t) &\propto \left( \int_{\boldsymbol{\theta}_0} p(\mathcal{H}_t | \boldsymbol{\theta}_0) p(\boldsymbol{\theta}_0 | \boldsymbol{\theta}_L) d\boldsymbol{\theta}_0 \right) \cdot p(\boldsymbol{\theta}_L), \\ p(\boldsymbol{\theta}_{\ell-1} | \boldsymbol{\theta}_\ell, \mathcal{H}_t) &\propto \left( \int_{\boldsymbol{\theta}_0} p(\mathcal{H}_t | \boldsymbol{\theta}_0) p(\boldsymbol{\theta}_0 | \boldsymbol{\theta}_{\ell-1}) d\boldsymbol{\theta}_0 \right) \cdot p(\boldsymbol{\theta}_{\ell-1} | \boldsymbol{\theta}_\ell), \quad \text{for } \ell = 2, \dots, L, \\ p(\boldsymbol{\theta}_0 | \boldsymbol{\theta}_1, \mathcal{H}_t) &\propto p(\mathcal{H}_t | \boldsymbol{\theta}_0) \cdot p(\boldsymbol{\theta}_0 | \boldsymbol{\theta}_1). \end{aligned}$$

**Lemma F.2** (Tweedie’s formula (Efron, 2011)). *Let  $p(\mathbf{y} | \boldsymbol{\eta})$  belong to the exponential family distribution  $p(\mathbf{y} | \boldsymbol{\eta}) = p_0(\mathbf{y}) \exp(\boldsymbol{\eta}^\top T(\mathbf{y}) - \varphi(\boldsymbol{\eta}))$ , where  $\boldsymbol{\eta}$  is the canonical vector of the family,  $T(\mathbf{y})$  is some function of  $\mathbf{y}$ , and  $\varphi(\boldsymbol{\eta})$  is the cumulant generation function which normalizes the density, and  $p_0(\mathbf{y})$  is the density up to the scale factor when  $\boldsymbol{\eta} = \mathbf{0}$ . Then, the posterior mean  $\hat{\boldsymbol{\eta}} := \mathbb{E}[\boldsymbol{\eta} | \mathbf{y}]$  should satisfy*

$$(\nabla_{\mathbf{y}} T(\mathbf{y}))^\top \hat{\boldsymbol{\eta}} = \nabla_{\mathbf{y}} \log p(\mathbf{y}) - \nabla_{\mathbf{y}} \log p_0(\mathbf{y}).$$

## G DIFFUSION MODEL TRAINING

In this section, we introduce the training of diffusion model for bandit with nonlinear reward model, which is shown in Algorithm 4. The training process is a standard DDPM training process (Ho et al., 2020).

---

### Algorithm 4 Diffusion Model Training

---

**Input:** Training dataset  $\mathcal{D} = \{\boldsymbol{\theta}^{(i)}\}_{i=1}^N$ , total diffusion steps  $L$ , noise schedule  $\{\beta_\ell\}_{\ell=1}^L$ , learning rate  $\eta$ , mini-batch size  $B$ , epochs  $E$ .

- 1:  $\alpha_\ell \leftarrow 1 - \beta_\ell$  and  $\bar{\alpha}_\ell \leftarrow \prod_{j=1}^\ell \alpha_j$  for all  $\ell$ .
- 2: **Initialization:** denoiser parameters  $\phi \leftarrow \text{randn}()$ .
- 3: **for epoch**  $= 1, \dots, E$  **do**
- 4:   **for all** mini-batches  $\mathcal{B} \subset \mathcal{D}$  of size  $B$  **do**
- 5:     Sample time indices  $\ell_1, \dots, \ell_B \sim \text{Unif}\{1, \dots, L\}$ .
- 6:     Sample noises  $\varepsilon_1, \dots, \varepsilon_B \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .
- 7:     **for**  $b = 1, \dots, B$  **do**
- 8:        $X_0^{(b)} \leftarrow \boldsymbol{\theta}^{(b)}$ .
- 9:        $X_{\ell_b}^{(b)} \leftarrow \sqrt{\bar{\alpha}_{\ell_b}} X_0^{(b)} + \sqrt{1 - \bar{\alpha}_{\ell_b}} \varepsilon_b$ . {forward diffusion}
- 10:        $\hat{\varepsilon}_b \leftarrow \varepsilon_\psi(X_{\ell_b}^{(b)}, \ell_b)$ . {network prediction}
- 11:     **end for**
- 12:      $\mathcal{L} \leftarrow \frac{1}{B} \sum_{b=1}^B \|\hat{\varepsilon}_b - \varepsilon_b\|_2^2$ .
- 13:     **Gradient decent:**  $\phi \leftarrow \phi - \eta \nabla_\phi \mathcal{L}$ .
- 14:   **end for**
- 15: **end for**

**Output:** Trained diffusion model denoiser  $\varepsilon_{\phi^*}(\cdot, \cdot)$ .

---

## H EXPERIMENT DETAILS

In this section, we provide details about our experiment setups and implementations.

## H.1 IMPLEMENT DETAILS

**Diffusion Prior Implementation** Our diffusion prior is implemented within a Denoising Diffusion Probabilistic Model (DDPM) framework, using the standard  $\epsilon$ -prediction parameterization. The core of this model is a denoiser network responsible for predicting the noise at each diffusion step. We implement this denoiser using a Multi-Layer Perceptron (MLP) conditioned on a sinusoidal time embedding. The key architectural details of the MLP used at each diffusion step are summarized in Table 3.

Table 3: MLP Denoiser Configuration.

Parameter	Value
Hidden Layers	(32, 32)
Activation	ReLU
Optimizer	Adam

**Hyperparameters** We provide the hyperparameters for our proposed algorithm in Table 4 and Table 5.

Table 4: Hyperparameters of DLTS

Reward Model	Prior Distribution	Update Steps	Step size	Noise Scale
Cosine	cross	1	0.05	0.005
	rays	10	0.01	0.005
	triangles	1	0.005	0.05
	swirl	1	0.01	0.01
	H	1	0.05	0.01
	corners	1	0.05	0.05
Quadratic	cross	1	0.01	0.005
	rays	10	0.05	0.005
	triangles	10	0.05	0.01
	swirl	1	0.05	0.005
	H	1	0.005	0.1
	corners	10	0.01	0.05
Sigmoid-gated	cross	1	0.05	0.05
	rays	1	0.05	0.05
	triangles	10	0.05	0.005
	swirl	10	0.01	0.1
	H	1	0.1	0.01
	corners	1	0.05	0.1

## H.2 PRIOR DISTRIBUTION

The ‘cross’ problem is generated from a mixture of two highly correlated 2D Gaussian distributions, creating a distinct cross shape. The ‘rays’ distribution is formed using rejection sampling, where accepted points are constrained to lie close to the four cardinal axes, resembling rays emanating from the origin. The ‘triangles’ distribution, also generated via rejection sampling, consists of points within two triangular regions. The ‘swirl’ distribution is generated parametrically, creating a spiral pattern with added noise. The ‘H’ distribution uses rejection sampling to accept points within three rectangular regions that form the letter ‘H’. Finally, the ‘corners’ distribution consists of points sampled from four distinct rectangular areas located in the corners of the sampling domain.

## I ADDITIONAL EXPERIMENTAL RESULTS

In this section, we present additional experimental results to complement those in Section 5.

Table 5: Hyperparameters of DPSG-MP

Reward Model	Prior Distribution	Update Steps	Step Size
Cosine	cross	1	0.05
	rays	10	0.005
	triangles	10	0.01
	swirl	10	0.005
	H	1	0.05
	corners	1	0.05
Quadratic	cross	10	0.01
	rays	1	0.05
	triangles	1	0.1
	swirl	10	0.05
	H	10	0.005
	corners	10	0.001
Sigmoid-gated	cross	1	0.01
	rays	10	0.01
	triangles	10	0.1
	swirl	10	0.01
	H	10	0.05
	corners	10	0.05

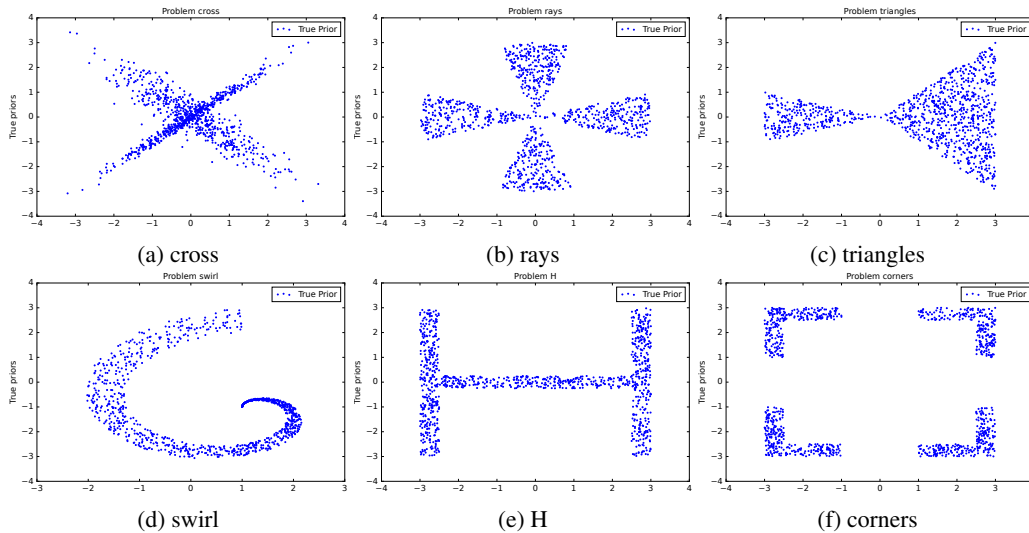


Figure 6: Prior Distribution. Each panel visualizes a prior over  $\theta \in \mathbb{R}^2$ . For each prior, we draw  $N = 10000$  parameter vectors; 80% are used as training tasks (previously seen tasks) and 20% as test tasks (newly encountered task). These parameters determine the true reward functions in our simulated bandit environments.

### I.1 RESULTS FOR THE TRIANGLES, SWIRL AND H PRIORS IN LINEAR BANDITS AND SIGMOID-GATED BANDITS

We present the triangles, swirl and H results in Linear bandits in Figure 7 and Figure 8. We also put the non-linear bandit results in Figure 9 and Figure 10.

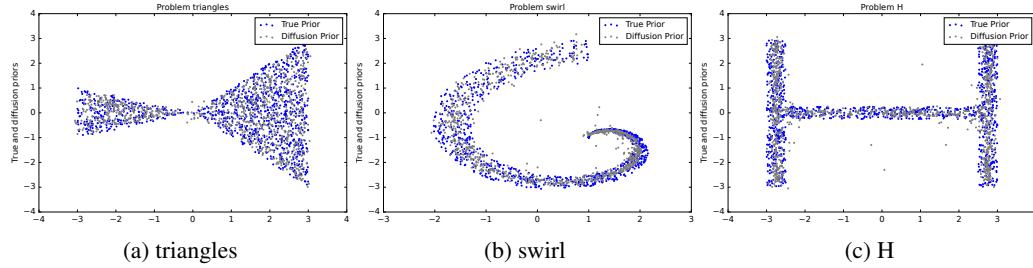


Figure 7: Visualization of the learned diffusion prior versus the true prior. We demonstrate three of them. For each figure, samples from our trained model (in grey) are overlaid on the ground-truth samples (in blue).

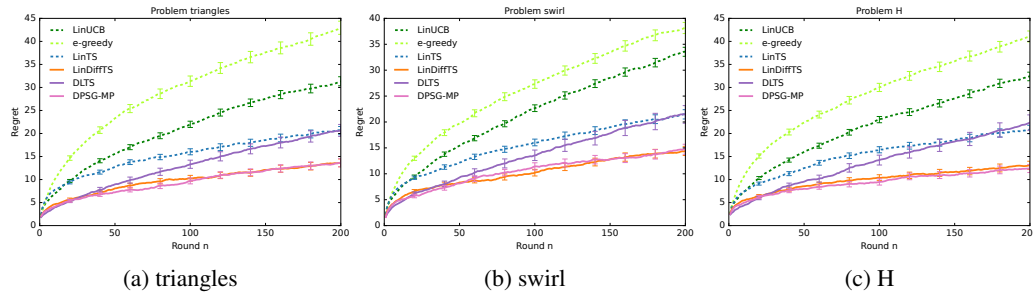


Figure 8: Performance of all algorithms on the linear contextual bandit tasks for the three prior distributions. Our proposed methods achieve cumulative regret comparable to LinTS and DiffTS.

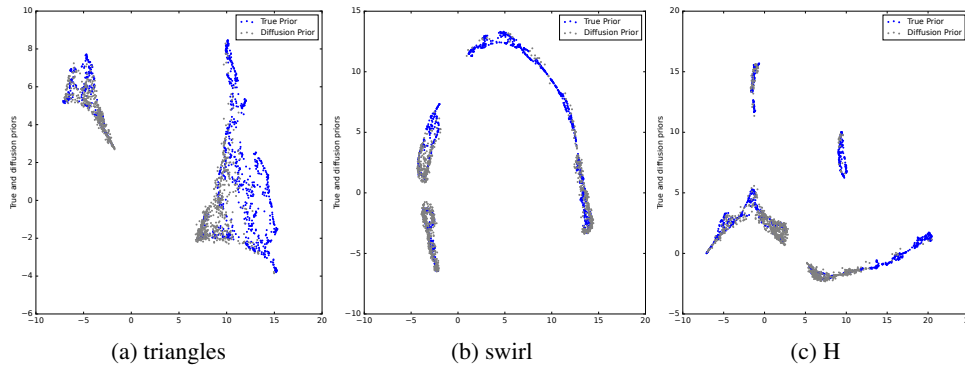


Figure 9: Visualization of the learned diffusion prior versus the true prior (under sigmoid-gated reward model). The true prior is the MLP parameters learned from the previous task. Samples from our trained model (grey) are overlaid on the ground-truth samples (blue).

## I.2 RESULTS FOR COSINE AND QUADRATIC NONLINEAR BANDITS

In the cosine reward setting, our proposed algorithms demonstrate strong performance, as shown in Figure 12. Both DLTS and the stabilized DPSG-MP achieve cumulative regret comparable to the state-of-the-art NeuralTS and NeuralUCB baselines. This result validates that a diffusion prior over neural network weights can effectively guide exploration in a complex, nonlinear environment. The reported results show the mean cumulative regret and standard error, averaged across 64 new tasks over 8 independent trials.

Similarly, for the quadratic reward model in Figure 14, our methods perform competitively against the strong neural baselines. Notably, for the challenging ‘corners’ prior, our algorithms outperform both NeuralTS and NeuralUCB. This highlights a key advantage of our approach: when the underlying parameter distribution has a complex structure, the guidance from the learned diffusion prior

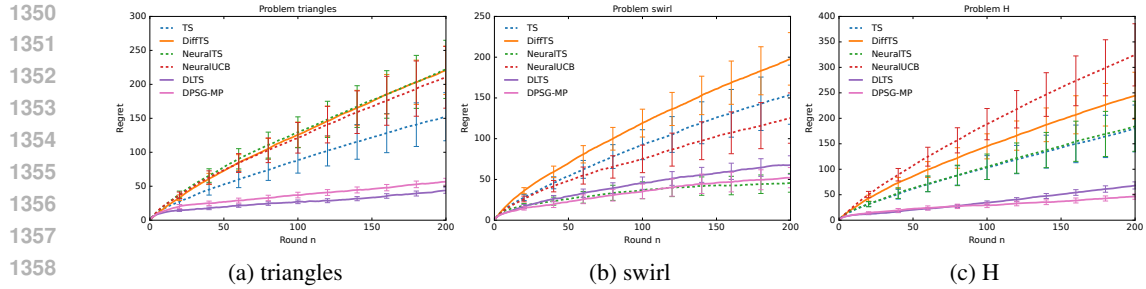


Figure 10: Performance of all algorithms on the sigmoid-gated nonlinear bandit tasks for the three prior distributions. Our proposed methods outperform strong baselines like NeuralTS and NeuralUCB, achieving the lowest cumulative regret.

becomes particularly beneficial, leading to more efficient exploration and lower regret. The evaluation protocol remains the same, with results averaged over 64 tasks across 8 independent trials.

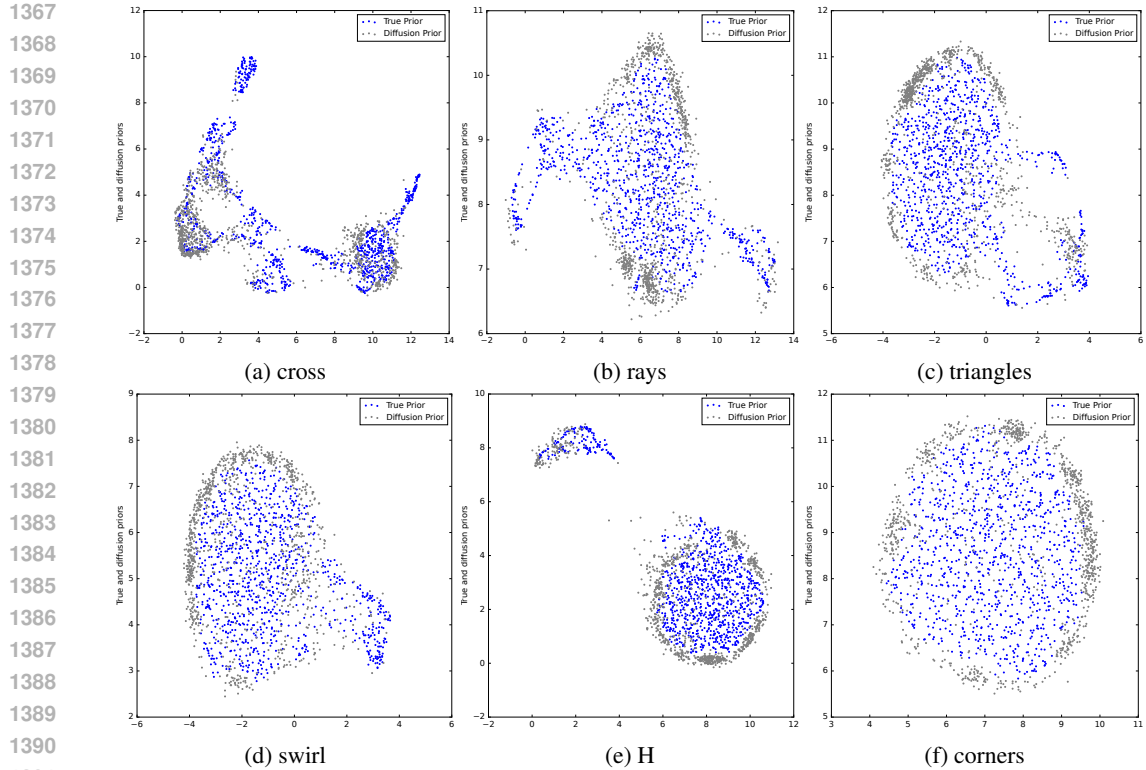


Figure 11: Visualization of the learned diffusion prior versus the true prior (under cosine reward model). The true prior is the MLP parameters learned from the previous task. Samples from our trained model (grey) are overlaid on the ground-truth samples (blue).

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

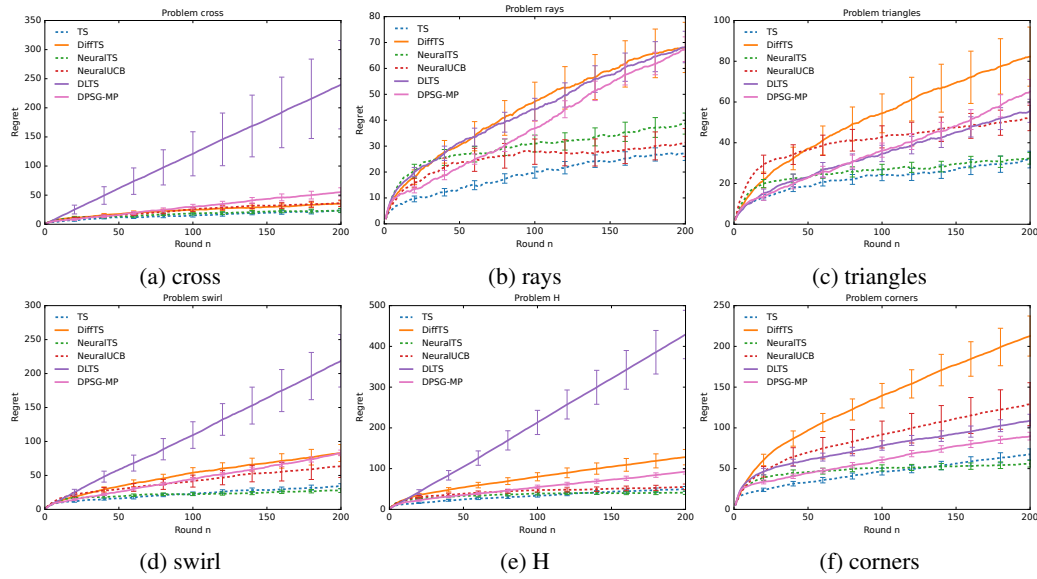


Figure 12: Performance of all algorithms on the cosine nonlinear bandit tasks for the six prior distributions.

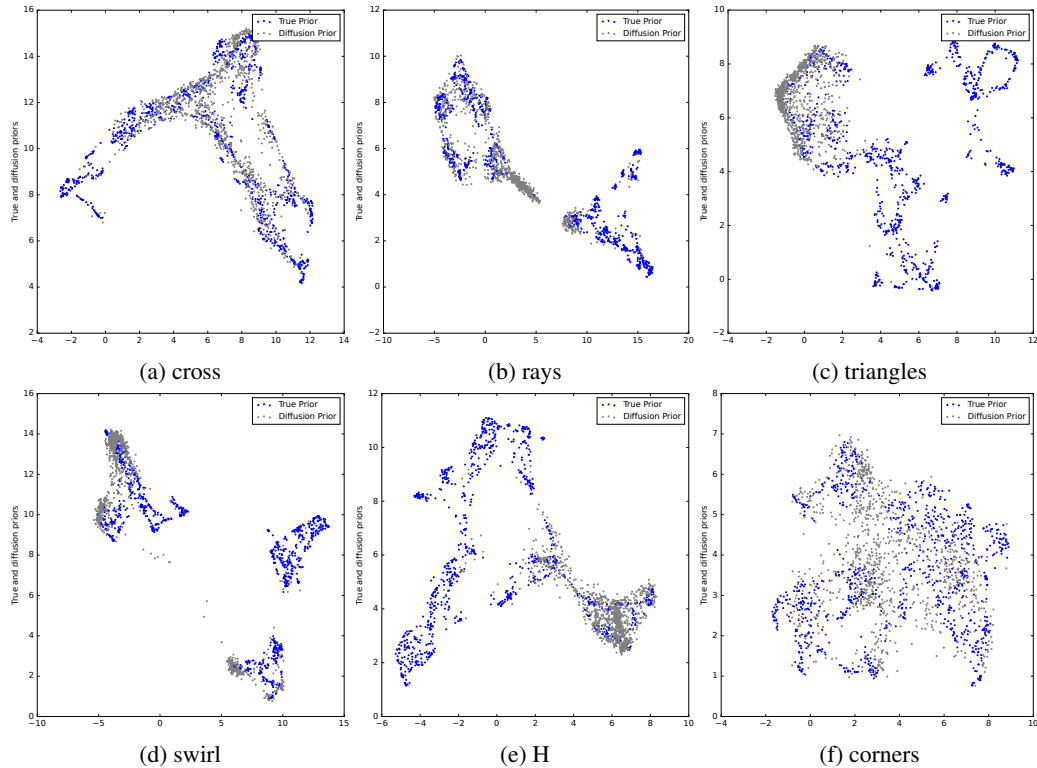


Figure 13: Visualization of the learned diffusion prior versus the true prior (under quadratic reward model). The true prior is the MLP parameters learned from the previous task. Samples from our trained model (grey) are overlaid on the ground-truth samples (blue).

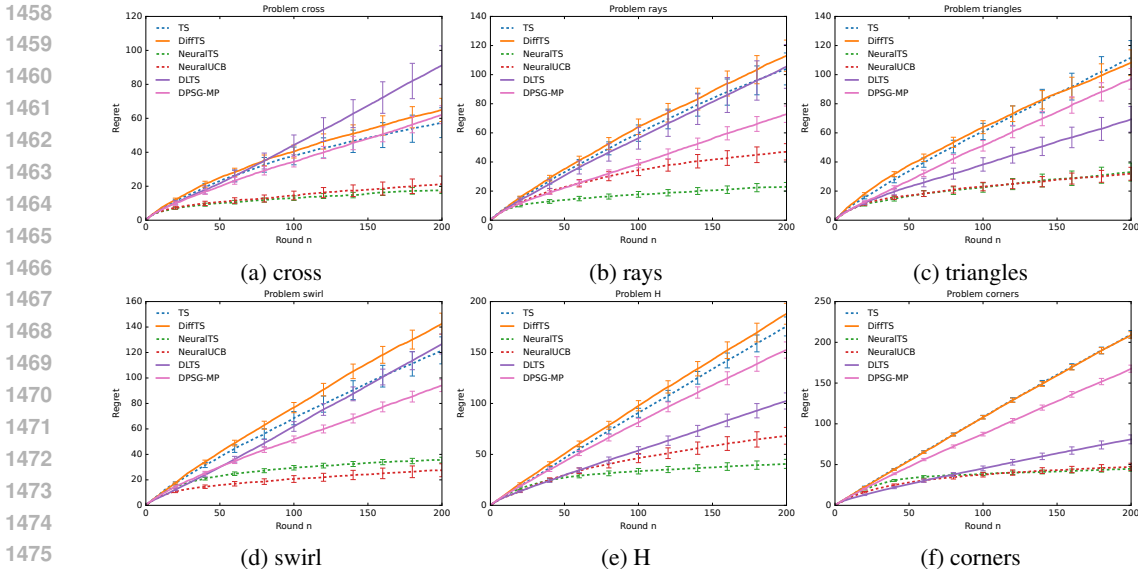


Figure 14: Performance of all algorithms on the quadratic nonlinear bandit tasks for the six prior distributions.

Table 6: Results of ablation study on Inner-Loop Updates  $K_\ell$  with fixed Diffusion Steps  $L = 100$ . Reported values are Cumulative Regret at  $T = 200$ . The regrets are averaged over 64 tasks across 8 independent runs.

Environment	Algorithm	Inner-Loop Steps ( $K_\ell$ )			
		$K_\ell = 1$	$K_\ell = 2$	$K_\ell = 5$	$K_\ell = 10$
Rays	DLTS	51.20 ± 2.90	<b>43.72 ± 1.58</b>	63.24 ± 3.48	55.87 ± 3.04
	DPSG-MP	39.59 ± 0.80	39.47 ± 0.85	37.36 ± 0.33	<b>36.31 ± 1.09</b>
Triangles	DLTS	79.60 ± 2.49	52.65 ± 5.32	160.99 ± 25.42	<b>47.29 ± 1.79</b>
	DPSG-MP	54.94 ± 6.35	51.28 ± 2.44	<b>49.51 ± 1.96</b>	52.20 ± 3.54
Swirl	DLTS	90.20 ± 3.63	78.81 ± 0.77	66.49 ± 2.34	<b>53.95 ± 1.74</b>
	DPSG-MP	<b>34.80 ± 1.43</b>	36.46 ± 0.87	34.81 ± 0.93	38.34 ± 1.63

### I.3 SENSITIVITY TO INNER-LOOP UPDATES $K_\ell$

We then investigate the effect of the number of inner-loop updates  $K_\ell$  which are LMC steps for DLTS and projection steps for DPSG-MP, fixing the Diffusion Steps  $L = 100$ . The results, summarized in Table 6, indicate that a small number of steps (e.g.,  $K_\ell \in \{1, 2\}$ ) is generally sufficient to achieve strong performance. Further increasing  $K_\ell$  yields diminishing returns and may even degrade performance in certain cases (e.g., DLTS on Triangles with  $K_\ell = 5$ ), likely due to overfitting the local likelihood approximation. Consequently, we recommend a small  $K_\ell$  as it offers the best trade-off between regret minimization and computational efficiency.

### I.4 LONG HORIZON EVALUATION

To validate the long-term robustness of our approach, we extend the evaluation of our most challenging nonlinear experiment, the sigmoid-gated reward model, to a horizon of  $T = 2000$  rounds

Table 7: Cumulative Regret at  $T = 2000$  on Sigmoid-Gated Nonlinear Bandits

Environment	LinTS	DiffTS	NeuralTS	NeuralUCB	DLTS (Ours)	DPSG-MP (Ours)
Cross	719.30 $\pm$ 329.99	899.33 $\pm$ 622.08	441.20 $\pm$ 264.00	413.51 $\pm$ 383.25	380.51 $\pm$ 343.59	<b>252.91 <math>\pm</math> 7.50</b>
Rays	1392.70 $\pm$ 545.47	1753.93 $\pm$ 883.84	501.88 $\pm$ 403.60	396.61 $\pm$ 304.10	<b>301.98 <math>\pm</math> 21.53</b>	355.43 $\pm$ 98.89
Triangles	942.77 $\pm$ 505.50	2873.07 $\pm$ 530.49	151.44 $\pm$ 110.73	1717.32 $\pm$ 532.37	513.40 $\pm$ 137.74	272.09 $\pm$ 82.04
Swirl	1388.62 $\pm$ 723.92	1012.04 $\pm$ 231.75	295.03 $\pm$ 304.34	997.91 $\pm$ 597.07	515.18 $\pm$ 72.08	<b>176.44 <math>\pm</math> 33.48</b>
H	1702.86 $\pm$ 1043.99	3120.56 $\pm$ 827.30	<b>114.73 <math>\pm</math> 25.09</b>	1303.26 $\pm$ 603.24	1293.67 $\pm$ 136.93	323.96 $\pm$ 36.16
Corners	1539.22 $\pm$ 635.22	2057.89 $\pm$ 1181.21	<b>244.09 <math>\pm</math> 105.71</b>	3515.52 $\pm$ 953.39	474.54 $\pm$ 56.57	258.20 $\pm$ 24.41

(increased from  $T = 200$ ). The cumulative regret results, summarized in the table below across various prior distributions.

The extended horizon experiments confirm the robustness of our framework. In the ‘Cross’ and ‘Rays’ environments, both our methods (DLTS and DPSG-MP) clearly outperform all baselines, achieving the lowest cumulative regret. In ‘Swirl’ environment, DPSG-MP clearly outperform all baselines. In the remaining three scenarios, our approach consistently achieves top-two performance. Specifically, in ‘Corners’, NeuralTS yields regret comparable to DPSG-MP (244.09 vs. 258.20) but suffers from significantly higher instability ( $\pm 105.71$  vs.  $\pm 24.41$ ). In general, the neural baselines exhibit much higher variance across tasks, whereas our diffusion-based methods provide a more stable and reliable exploration strategy over long horizons.