

Diffusion Prior for Online Decision Making: A Case Study of Thompson Sampling

Yu-Guan Hsieh¹, Patrick Bloebaum², Shiva Kasiviswanathan², Branislav Kveton² ¹Université Grenoble Alpes ²AWS AI Labs

Introduction

TL;DR: We show that diffusion model is able to learn a prior that can help reduce the regret in multi-armed bandit problems

Multi-Armed Bandits

A model for online decision making:

- Learner pulls arm a_t at round t
- Learner receives rewards r_t drawn from the arm's distribution
- The goal is to maximize the cumulative rewards, i.e., $\sum_t r_t$

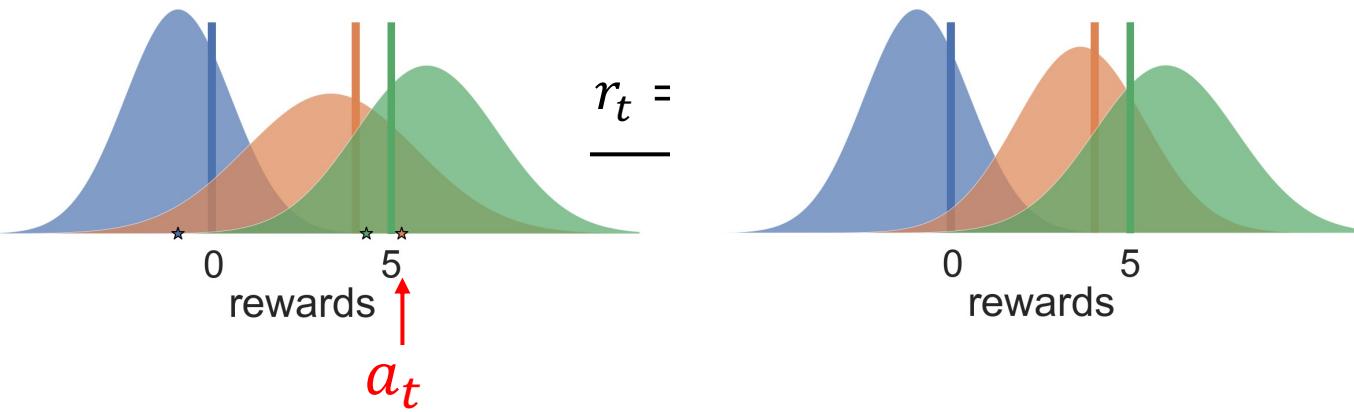
Thompson Sampling (TS)

A Bayesian approach to tackle exploration-exploitation tradeoff

- Given:
 - A prior $p(\mu)$ over mean reward vector μ
 - Interaction history $\mathcal{H}_t = (a_s, r_s)_{s \in \{1, \dots, t\}}$
- Maintain posterior distribution

$$p(w | \mathcal{H}_t) \propto p(\mathcal{H}_t | w)p(w)$$

- Sample $\tilde{\mu}_t$ from the posterior distribution
- Pull $a_t \in \operatorname{argmax}_{a \in \mathcal{A}} \tilde{\mu}_t^a$ and update posterior

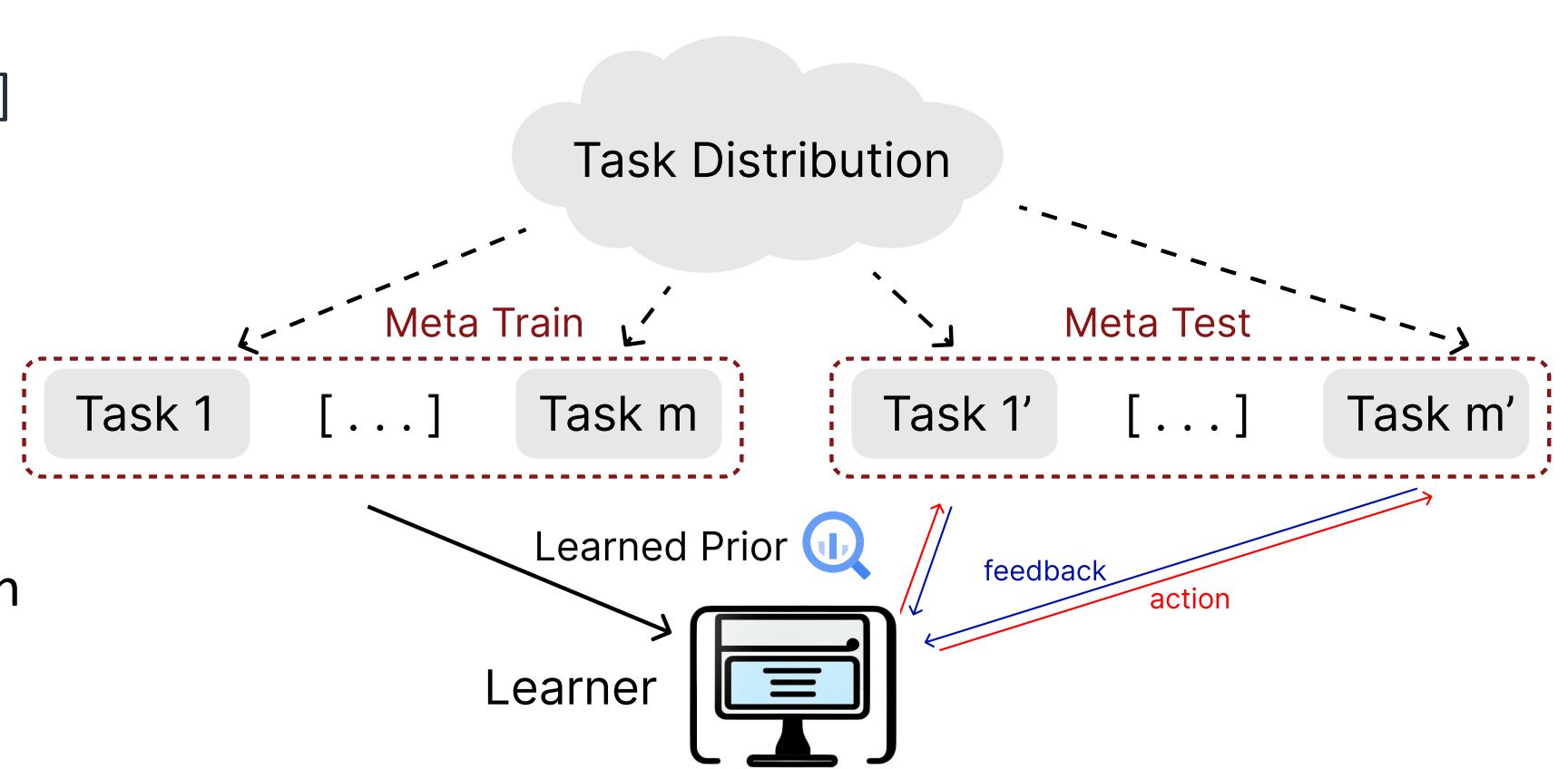


TS is sensitive to the choice of prior:
Can we learn a good prior to minimize regret?

Meta-Learning for Bandits^[1]

Different bandit instances can have similar patterns

- Recommend items to different customers
- Assign price to different items using an online pricing algorithm



Goal: Meta-learn a prior from similar tasks using diffusion models

Diffusion Models^[2]

- Noise is gradually added in the forward diffusion process

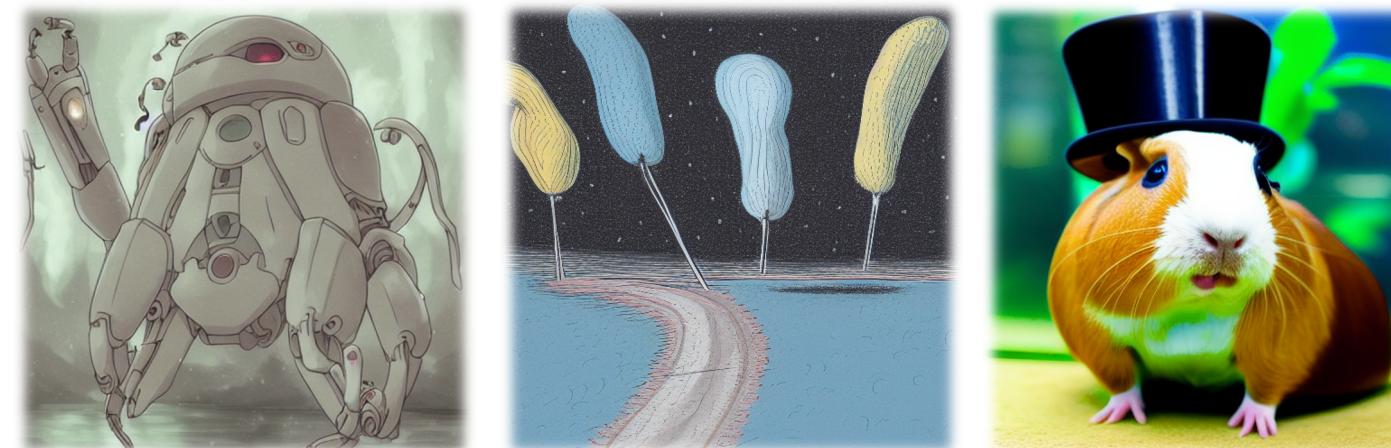
$$q(X_{\ell+1} | x_\ell) = \mathcal{N}(X_{\ell+1}; \sqrt{\alpha_{\ell+1}} x_\ell, (1 - \alpha_{\ell+1})I)$$
- The model is trained to reverse the process

$$p_\theta(X_\ell | x_{\ell+1}) = q(X_\ell | x_{\ell+1}, X_0 = h_\theta(x_{\ell+1}, \ell + 1))$$
where h_θ is a denoiser and predicts x_0
- Variance calibration

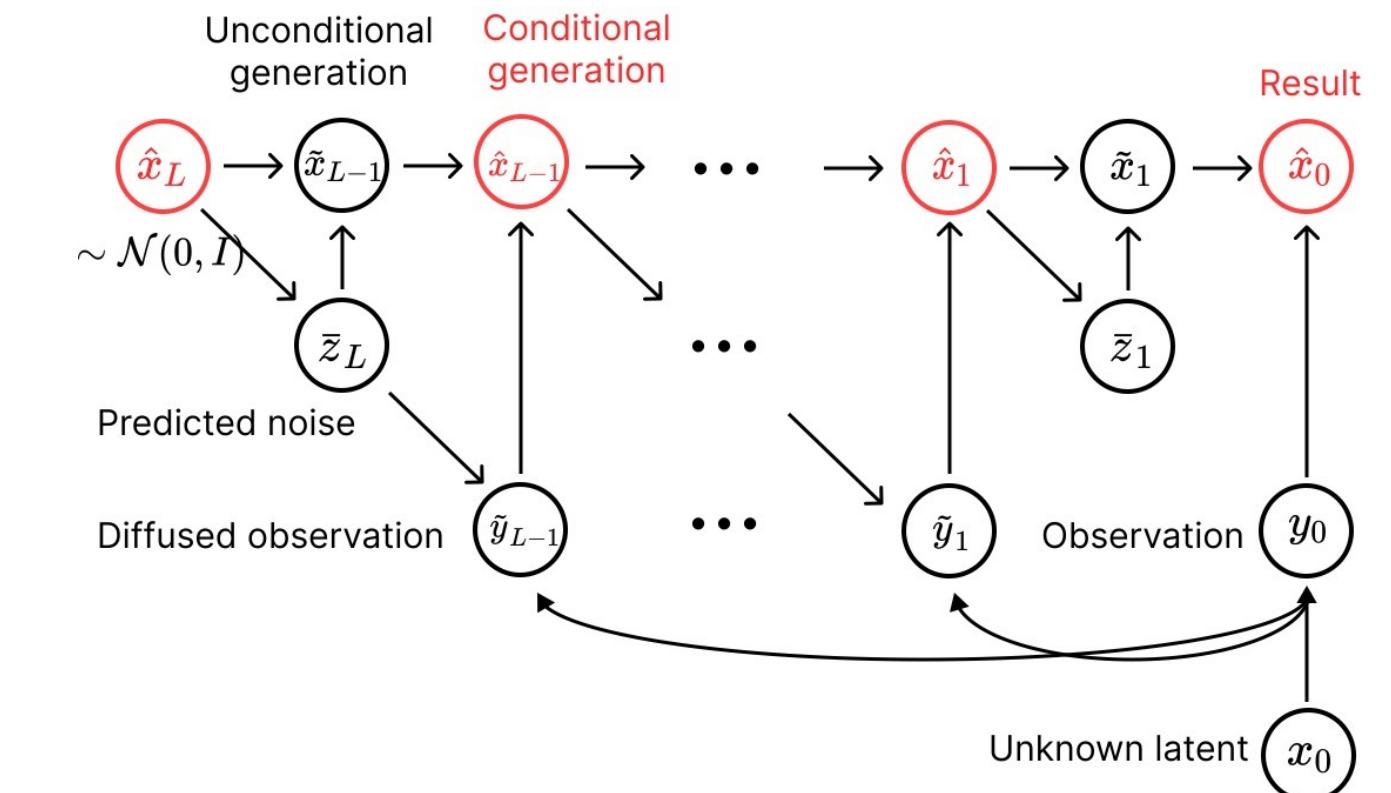
$$p_{\theta, \tau}(X_\ell | x_{\ell+1}) = \int q(X_\ell | x_{\ell+1}, x_0) p_{\theta, \tau}(x_0 | x_{\ell+1}) dx_0$$

with $p_{\theta, \tau}(x_0 | x_{\ell+1}) = \mathcal{N}(h_\theta(x_{\ell+1}, \ell + 1), \tau_{\ell+1}^2)$
and mean squared reconstruction error

$$\tau_\ell^2 = \mathbb{E}_{x_0, x_\ell} [\|x_0 - h_\theta(x_\ell, \ell)\|^2]$$



TS with Diffusion Prior



Let X_0 denote the vector of expected rewards

The goal is to sample $\tilde{\mu}_t$ from $X_0 | \mathcal{H}_t$

1. Sample \hat{x}_L from $\mathcal{N}(0, I)$
2. Sample \hat{x}_ℓ from $X_\ell | \hat{x}_{\ell+1}, y_0$
 - Arm never pulled: $\hat{x}_\ell^a \sim p_{\theta, \tau}(X_\ell^a | \hat{x}_{\ell+1})$
 - Arm pulled: compute empirical mean $\hat{\mu}_t^a$ and $\tilde{y}_\ell^a = \sqrt{\bar{\alpha}_\ell} \hat{\mu}_t^a + \sqrt{1 - \bar{\alpha}_\ell} \bar{z}_\ell^a$

$$\tilde{q}(X_\ell^a | \hat{x}_{\ell+1}, \mathcal{H}_t) \propto p_{\theta, \tau}(X_\ell^a | \hat{x}_{\ell+1}) \mathcal{N}(\tilde{y}_\ell^a, \sigma_{t, \ell}^a)$$

Sample : $\hat{x}_\ell^a \sim \tilde{q}(X_\ell^a | \hat{x}_{\ell+1}, \mathcal{H}_t)$

Numerical Experiments

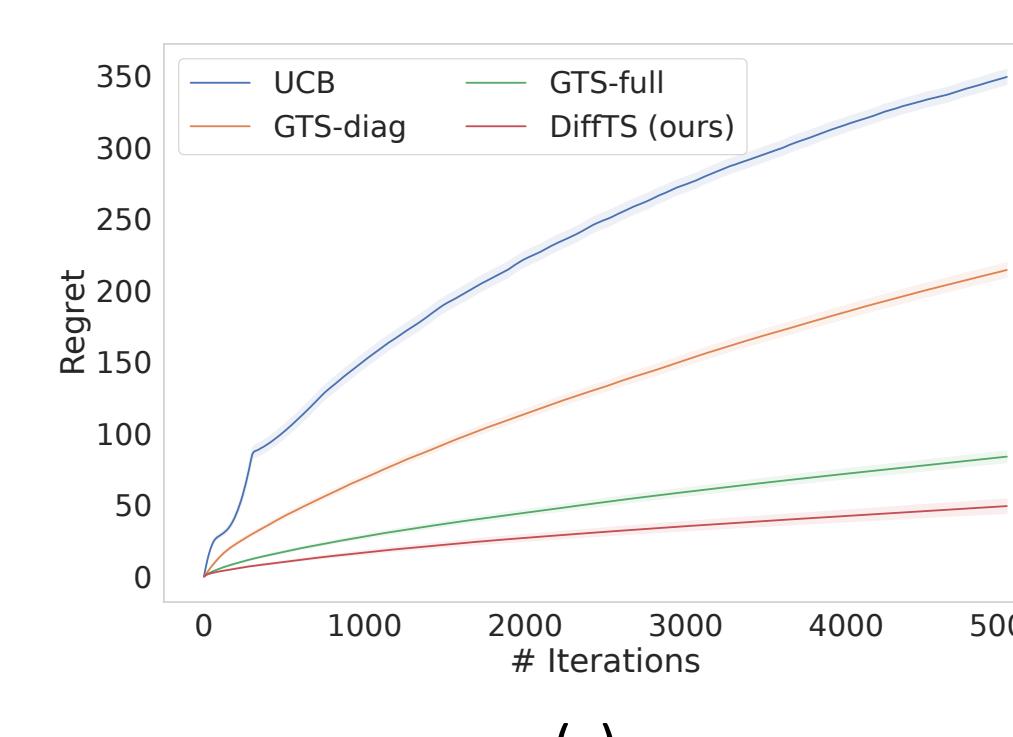
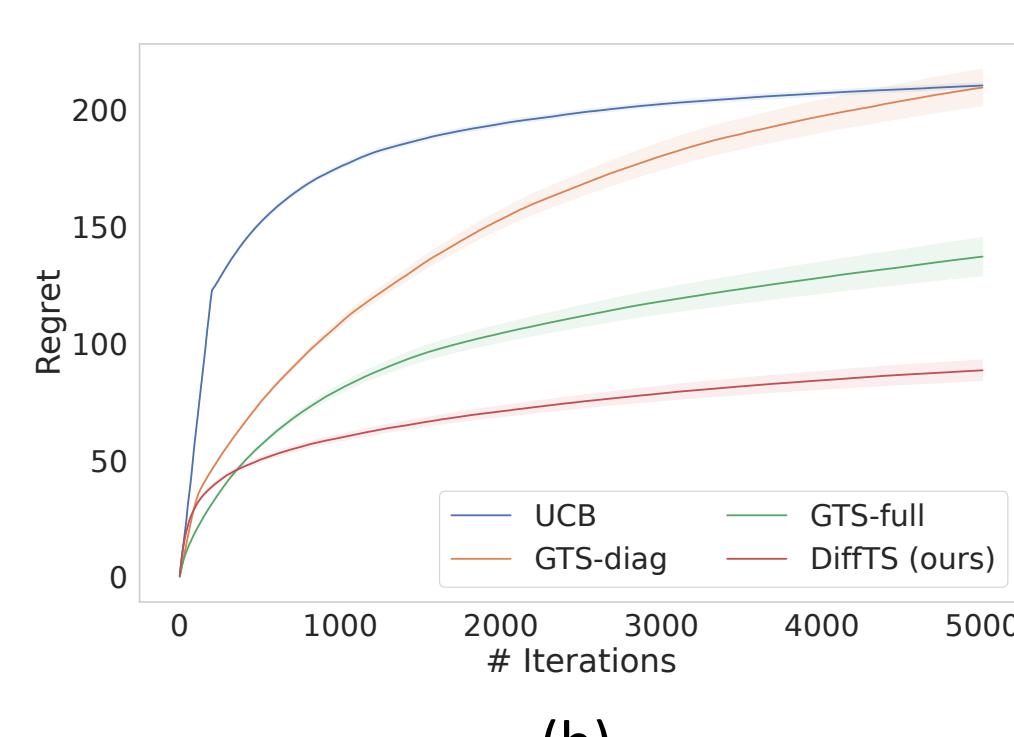
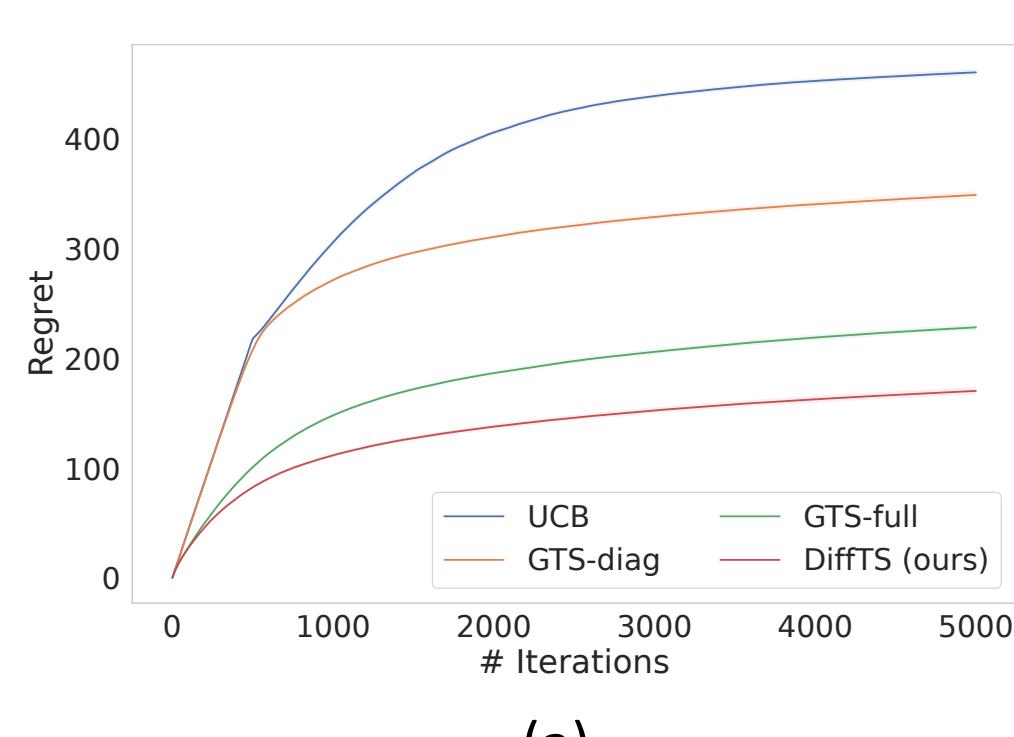
Handcrafted tasks with underlying patterns. We consider three settings:

- Arms are labeled and the rewards are related to these labels
- Arms are separated into different groups, some groups have higher means but are never optimal (hard for Gaussian)
- Modeling the reward of a simple bidding problem with a real-world auction dataset

Algorithm Meta Learning for Bandits with Diffusion Models

- 1: **Meta Training**
- 2: **Input:** A set of expected means $(\mu_B)_B$ from different tasks $B \sim \mathcal{T}$
- 3: Train a diffusion model (a denoiser) h_θ to model the distribution of the mean rewards
- 4: **Calibration**
- 5: **Input:** A set of expected means $(\mu_B)_B$ from different tasks $B \sim \mathcal{T}$
- 6: Compute $(\tau_\ell^2)_{\ell \in \{1, \dots, L\}}$ for the denoiser h_θ at different noise levels
- 7: **Meta Test / Deployment**
- 8: For any new task B , run Thompson sampling with diffusion prior using the trained model

Regret = difference of cumulative rewards between an algorithm and the one that consistently chooses the best action
The algorithms are provided presumed noise variance σ'
The means and covariance of gaussian prior are estimated using training + calibration data



Real data

