

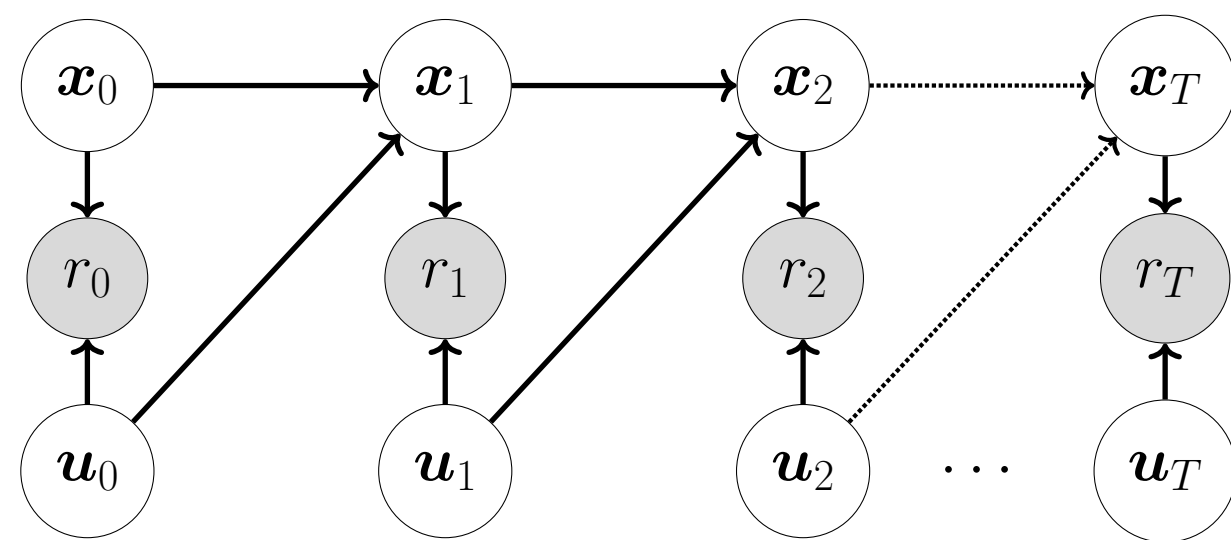
CLASSICAL BANDITS vs. OUR BANDITS



Setting: Bandits with Latent Dynamics

$$r_t = \mathbf{u}_t^\top \mathbf{C} \mathbf{x}_t + z_t, \quad \mathbf{x}_{t+1} = \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{u}_t + \mathbf{w}_t.$$

- Only rewards are observed; the state \mathbf{x}_t is latent.
- $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are unknown and $\rho(\mathbf{A}) < 1$.
- Action \mathbf{u}_t changes both current reward and future rewards.



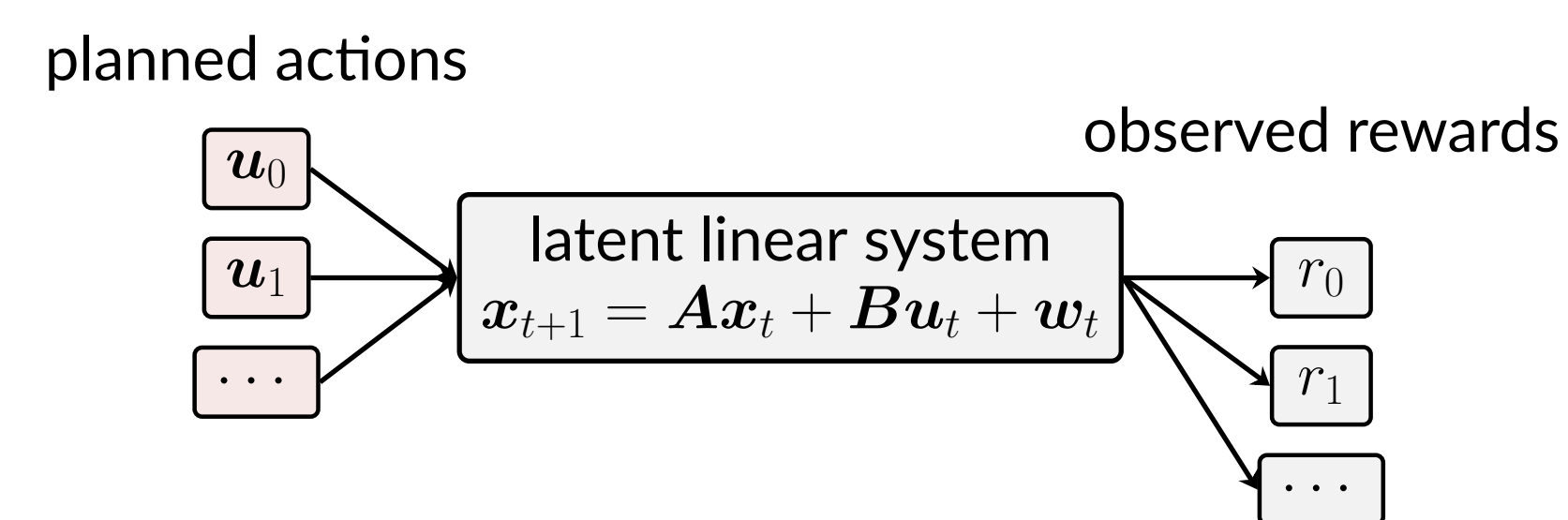
Objective & Benchmark

- Goal:** choose $\mathbf{u}_t \in [-1, +1]^p$ to maximize $\mathbb{E} \left[\sum_{t=0}^T r_t \right]$
- Regret** compares against the best open-loop sequence:

$$R_T(\pi) = \mathbb{E} \left[\sum_{t=0}^T r_t^* - \sum_{t=0}^T r_t^\pi \right].$$

Open-loop Policy

- Open-loop:** Choose a full action sequence $\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_T$ without reacting to realized r_t .



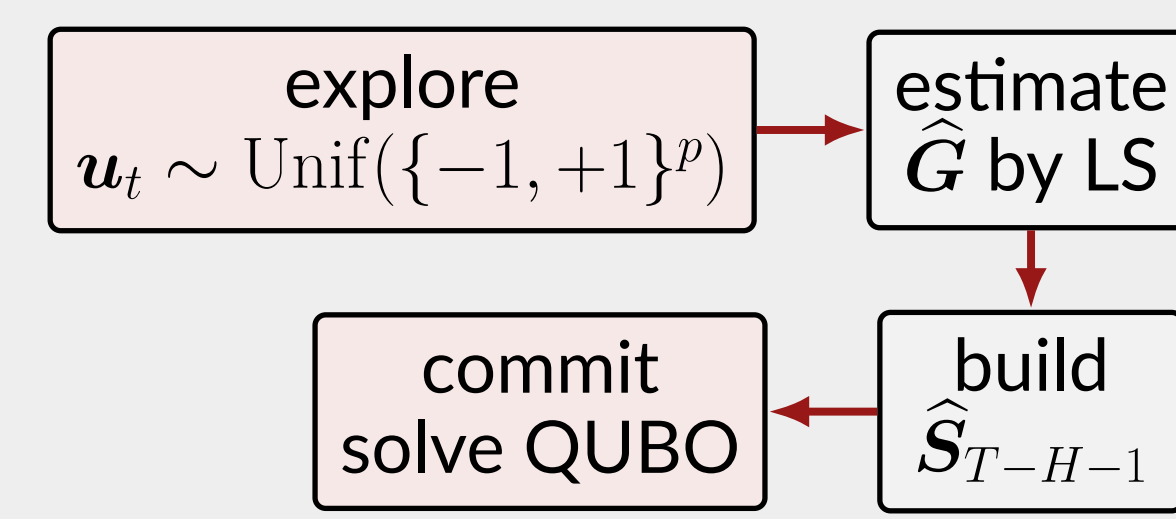
Optimal Open-loop Action Sequence from QUBO

Unrolling the dynamics gives a quadratic value:

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right] = \frac{1}{2} \mathbf{u}_{0:T}^\top \mathbf{S}_T \mathbf{u}_{0:T}, \quad \mathbf{u}_{0:T}^* \in \arg \max_{\mathbf{u}_t \in [-1, +1]^p} \frac{1}{2} \mathbf{u}_{0:T}^\top \mathbf{S}_T \mathbf{u}_{0:T}.$$

- An optimal sequence exists at a vertex: $\mathbf{u}_t^* \in \{-1, +1\}^p, \forall t$.
- The resulting QUBO is **NP-hard**.

Algorithm: Explore-then-Commit



- Require:** Horizon T , exploration length H , truncation length L
- Play random actions for $t = 0, \dots, H$ and observe rewards.
 - Estimate the first L Markov parameters $\hat{\mathbf{G}}$.
 - Optimize the estimated commit objective

$$\max_{\mathbf{u}_{H+1:T} \in \{-1, +1\}^{p(T-H)}} \frac{1}{2} \mathbf{u}_{H+1:T}^\top \hat{\mathbf{S}}_{T-H-1} \mathbf{u}_{H+1:T}.$$

- Play the chosen sequence for $t = H + 1, \dots, T$.

Learning How Actions Affect Rewards

Markov Parameters characterize past action-reward relationship.

$$\mathbf{G} = [\mathbf{C}\mathbf{B} \quad \mathbf{C}\mathbf{A}\mathbf{B} \quad \dots \quad \mathbf{C}\mathbf{A}^{L-1}\mathbf{B}]$$

- We only consider first L of them due to stability assumption.

For $t \geq L$, the reward r_t depends on the past L actions [3]:

$$r_t = \text{vec}(\mathbf{G})^\top (\bar{\mathbf{u}}_{t-1} \otimes \mathbf{u}_t) + \text{noise} + \mathcal{O}(\rho^L).$$

where $\bar{\mathbf{u}}_{t-1} = [\mathbf{u}_{t-1}^\top \mathbf{u}_{t-2}^\top \dots \mathbf{u}_{t-L}^\top]^\top$. From $\{(\mathbf{u}_t, r_t)\}_{t=0}^H$,

$$\hat{\mathbf{G}} = \arg \min_{\mathbf{G} \in \mathbb{R}^{p \times pL}} \sum_{t=L+1}^H (r_t - \text{vec}(\mathbf{G})^\top (\bar{\mathbf{u}}_{t-1} \otimes \mathbf{u}_t))^2.$$

Estimation guarantee With high probability,

$$H - L \gtrsim \tilde{\mathcal{O}}(p^2 L) \implies \|\hat{\mathbf{G}} - \mathbf{G}\|_F \lesssim \tilde{\mathcal{O}} \left(\sqrt{\frac{p^2 L}{H - L}} \right)$$

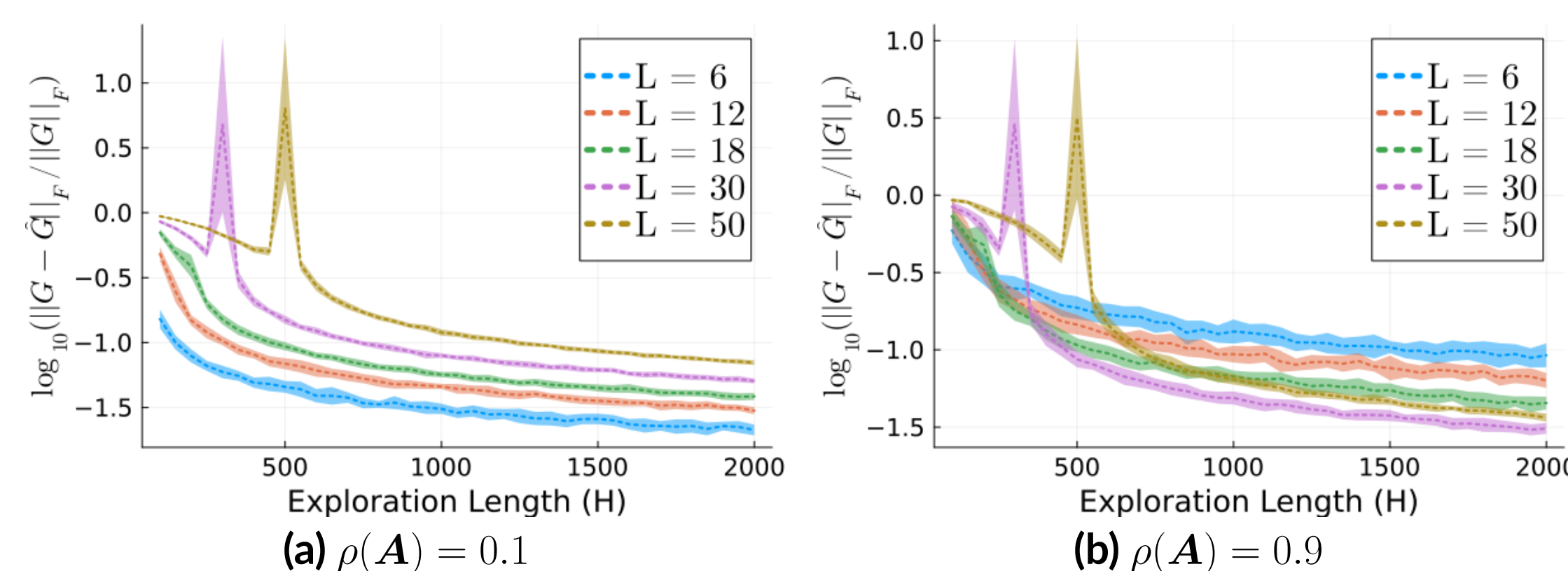


Figure 1. Markov Parameter Estimation Error

Regret Analysis Outline

Suboptimality Lemma. If $\|\hat{\mathbf{G}} - \mathbf{G}\|_F \leq \epsilon$, the estimated QUBO chooses a near-optimal open-loop sequence:

$$\text{commit suboptimality} \lesssim p(T - H)(\epsilon + \rho^L).$$

Regret decomposition:

$$R_T(\pi) \lesssim \underbrace{pH}_{\text{exploration}} + \underbrace{pT\epsilon}_{\text{estimation}} + \underbrace{pT\rho^L}_{\text{truncation}} \quad \text{w.h.p.}$$

Taking $\epsilon \simeq \tilde{\mathcal{O}}(H^{-1/2})$ and $L = \Theta(\log T)$:

$$R_T(\pi) \lesssim p \left(H + \frac{T}{\sqrt{H}} \right).$$

Main Regret Guarantee

Theorem. With $H = \tilde{\mathcal{O}}(T^{2/3})$ and $L = \Theta(\log T)$, the explore-then-commit policy satisfies

$$R_T(\pi) = \tilde{\mathcal{O}}(pT^{2/3})$$

with high probability.

- H controls the explore/exploit tradeoff.
- L makes the truncation negligible.

Practical QUBO Approximation

The commit problem is an NP-hard QUBO

$$\max_{\mathbf{x} \in \{\pm 1\}^d} \mathbf{x}^\top \mathbf{W} \mathbf{x} \quad \text{with } d = p(T - H).$$

Semidefinite Relaxation with Goemans-Williamson Rounding:

$$\max_{\mathbf{X} \succeq 0} \text{tr}(\mathbf{W} \mathbf{X}) \quad \text{s.t. } \text{rank}(\mathbf{X}) = 1, \mathbf{X}_{ii} = 1.$$

Solve SDP with Mosek [1]; then random-hyperplane round [2].

SignIter: Fast heuristic method; useful when the SDP is too large.

$$\mathbf{x}_i^{k+1} = \text{sign}((\mathbf{W} \mathbf{x}^k)_i).$$

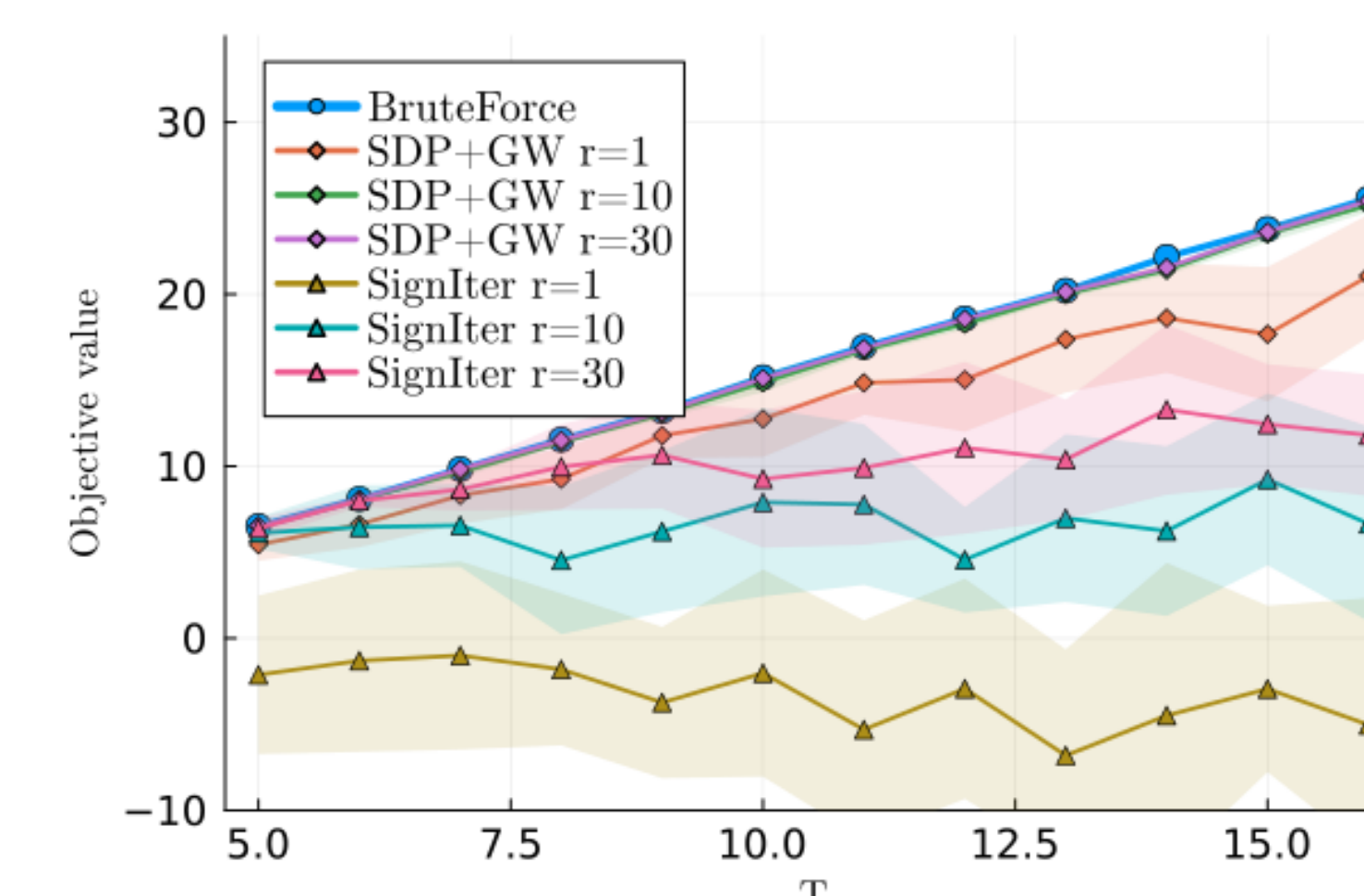


Figure 2. Small-dimensional comparison against brute force.

Numerical Results

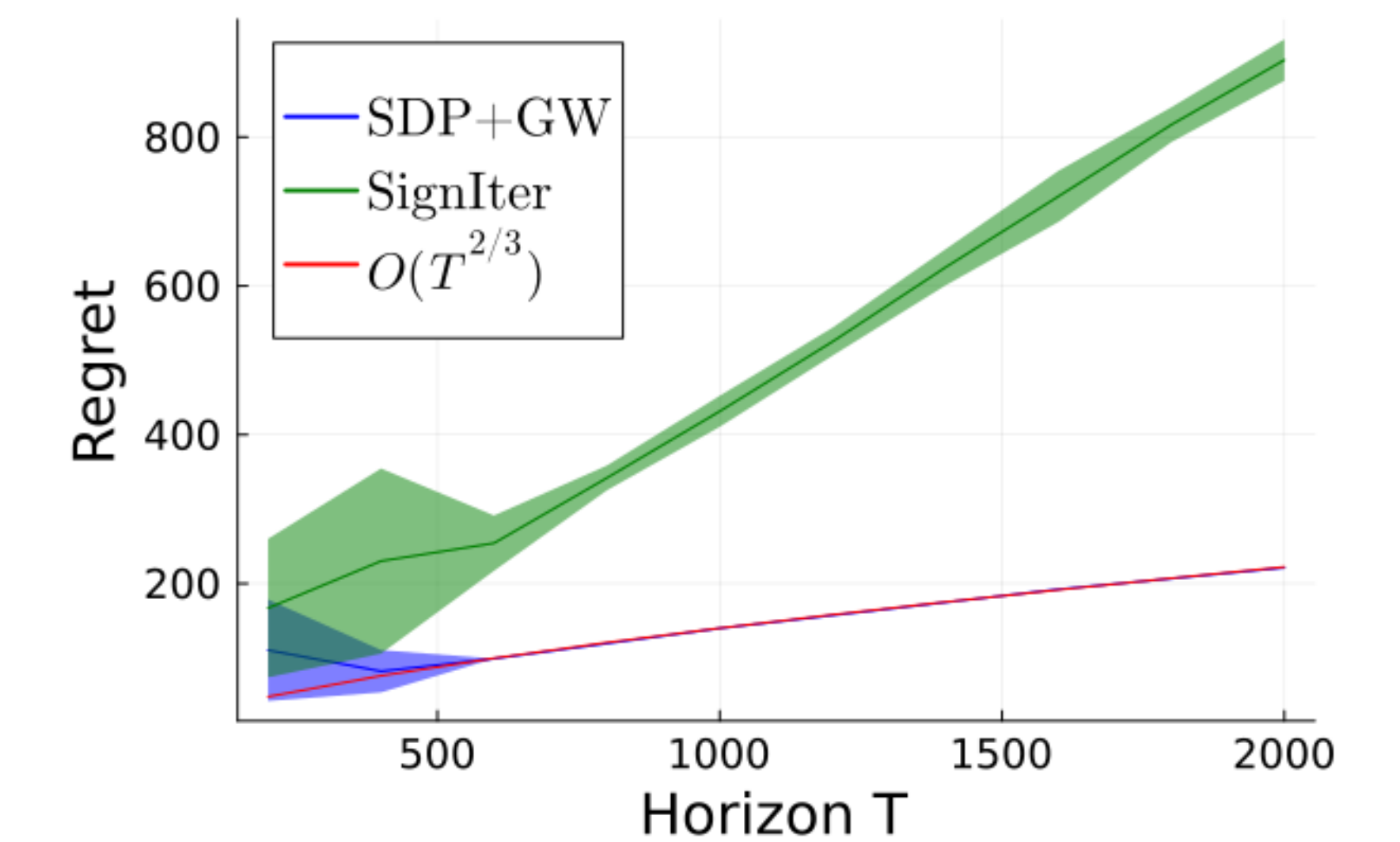


Figure 3. Regret of ETC using SDP+GW and SignIter.

- SDP+GW gives the stronger commit-phase benchmark.
- SignIter is scalable but can be more suboptimal.
- Empirical growth is consistent with theoretical regret.

Takeaways

- Regret is measured against the best open-loop action sequence.
- Random exploration learns the action-to-reward history.
- Approximate planning is enough for $\tilde{\mathcal{O}}(pT^{2/3})$ regret.

Future Work

- Adaptive closed-loop policies that use reward feedback.
- Faster QUBO solvers for larger horizons and actions.

References

- MOSEK ApS. MOSEK Optimizer API for Julia 11.0.29, 2025.
- Michel X Goemans and David P Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. Journal of the ACM (JACM), 42(6):1115–1145, 1995.
- Yahya Sattar, Yassir Jedra, and Sarah Dean. Learning linear dynamics from bilinear observations. In 2025 American Control Conference (ACC), pages 3109–3115. IEEE, 2025.

