

Gaussian Process Planning with Lipschitz Continuous Reward Functions: Towards Unifying Bayesian Optimization, Active Learning, and Beyond

Chun Kai Ling
Bryan Kian Hsiang Low
Patrick Jaillet

chunkai@comp.nus.edu.sg
lowkh@comp.nus.edu.sg
jaillet@mit.edu

1 Motivation



- Bayesian Optimization (BO) and some Active Learning (AL) problems often assume Gaussian Process (GP) priors.
- Identifying a common *planning* framework allows us to tackle both problems more effectively by utilizing planning techniques
- Provide basis for **theoretical guarantees** and **non-myopic** decision making
- Identify and solve **novel problems** with similar rewards and priors
- Potential applications: robot exploration in spatially correlated fields, robotic energy harvesting, hyperparameter tuning

2 Gaussian Process Planning (GPP)

Offers Flexibility in Reward functions:

- Weak restriction of R to be Lipschitz in Z
- Encompasses several existing formulations
- Generalizes to new interesting tasks

Examples:

- Maximum Entropy Sampling (Shewry, 1987)
- UCB selection criterion (Srinivas et al, 2010)
- Diminishing Rewards: $\log(Z)$ for $Z > 1$, 0 otherwise

$$V_t^\pi(d_t) \triangleq Q_t^\pi(\pi(d_t), \boxed{d_t}) \text{ Information state}$$

$$Q_t^\pi(s_{t+1}, d_t) \triangleq \mathbb{E}[R(Z_{t+1}, s_{t+1}) + V_{t+1}^\pi(\langle s_{t+1}, z_t \oplus Z_{t+1} \rangle) | s_{t+1}, d_t]$$

Immediate reward
Future reward
Dependency on past observations

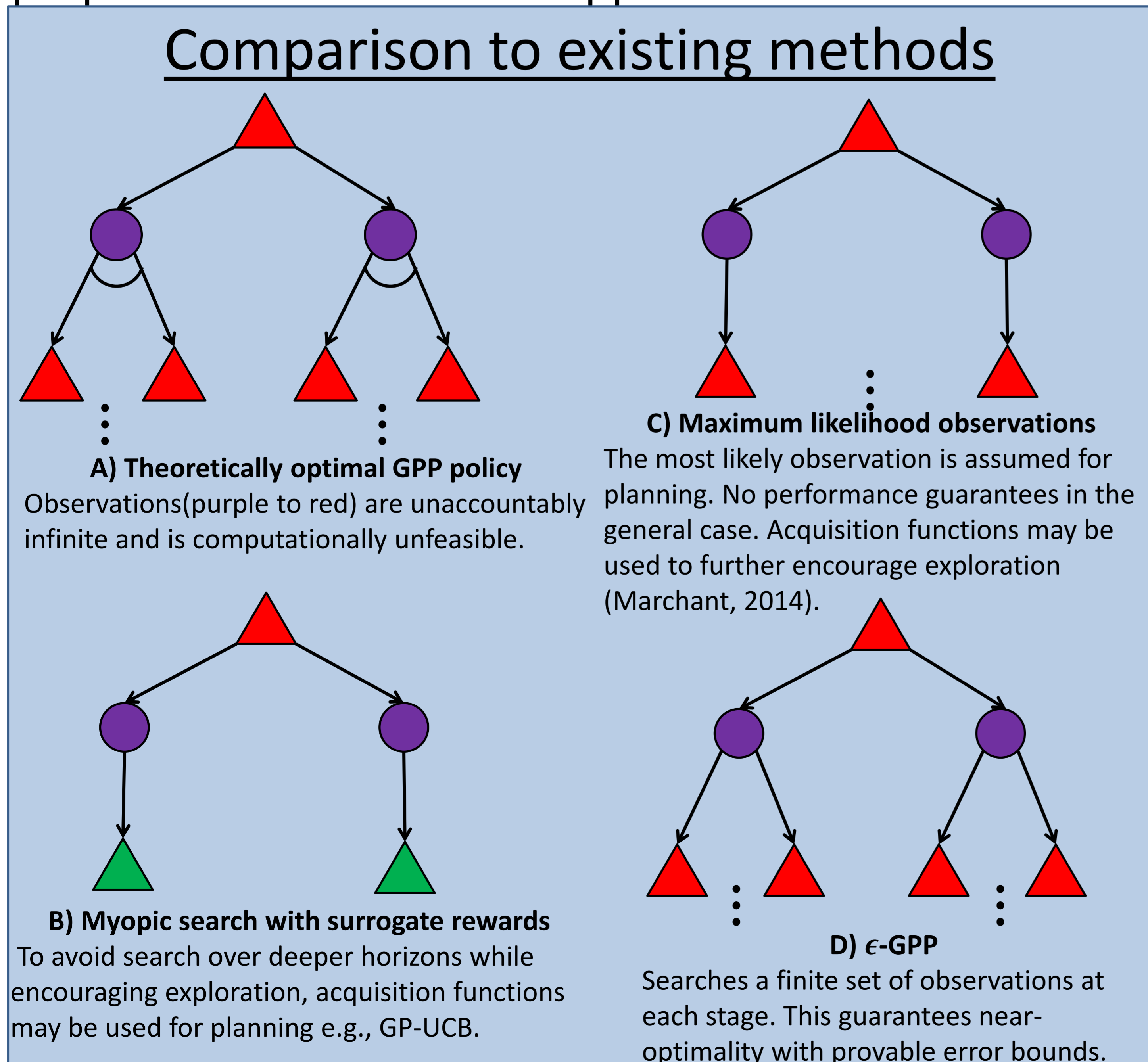
H-stage Bellman Equations

- Resolves exploration-exploitation tradeoff
- Removes the need to explicitly encourage exploration e.g., acquisition functions

3 ϵ -GPP

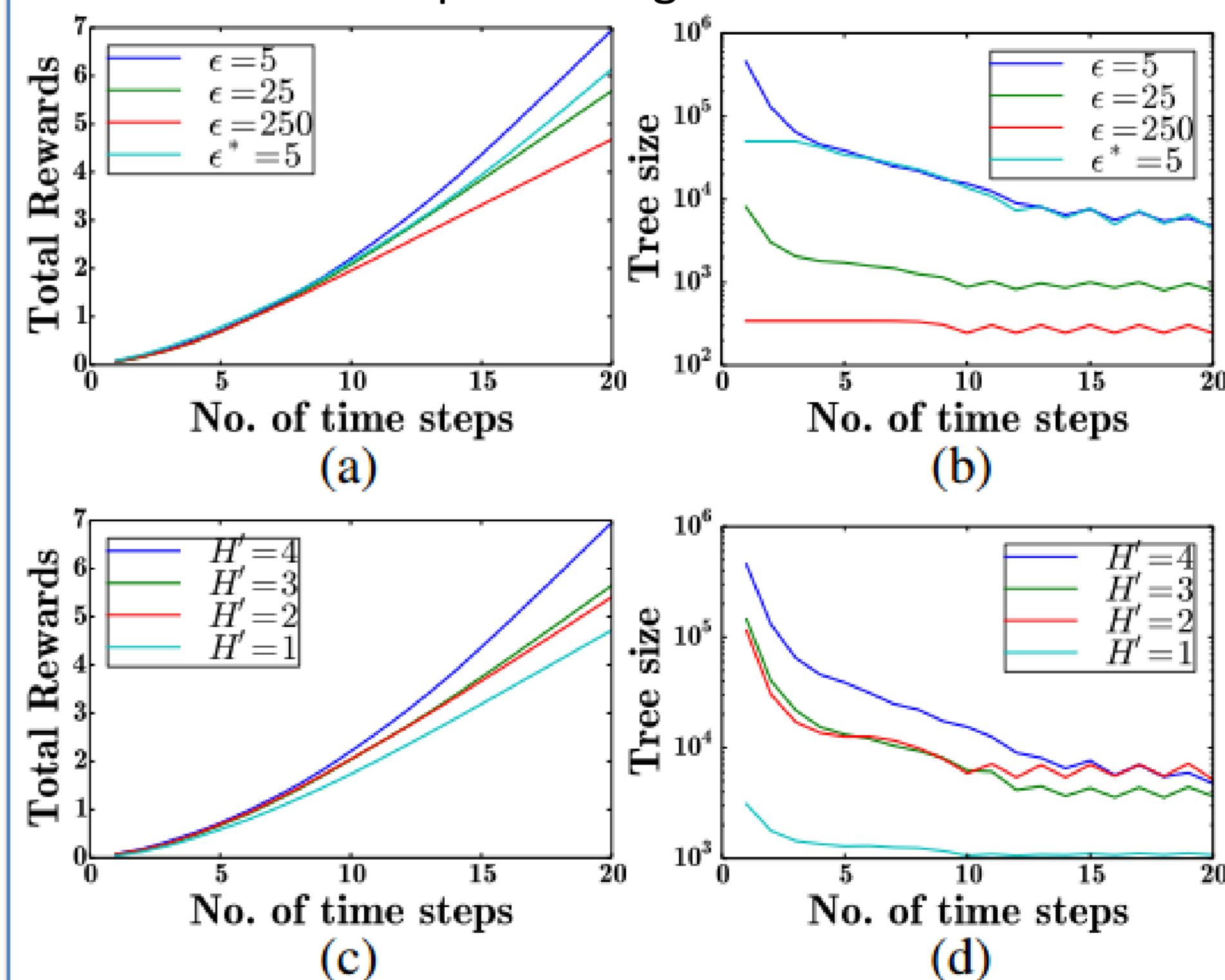
Key idea: By exploiting L1-continuous reward functions, a *finite* tree search may be used to approximate the search over all possible sample observations/actions

Guarantees policy is ϵ -optimal in specified horizon H . An anytime branch and bound extension of ϵ -GPP is proposed to suit real-time applications.



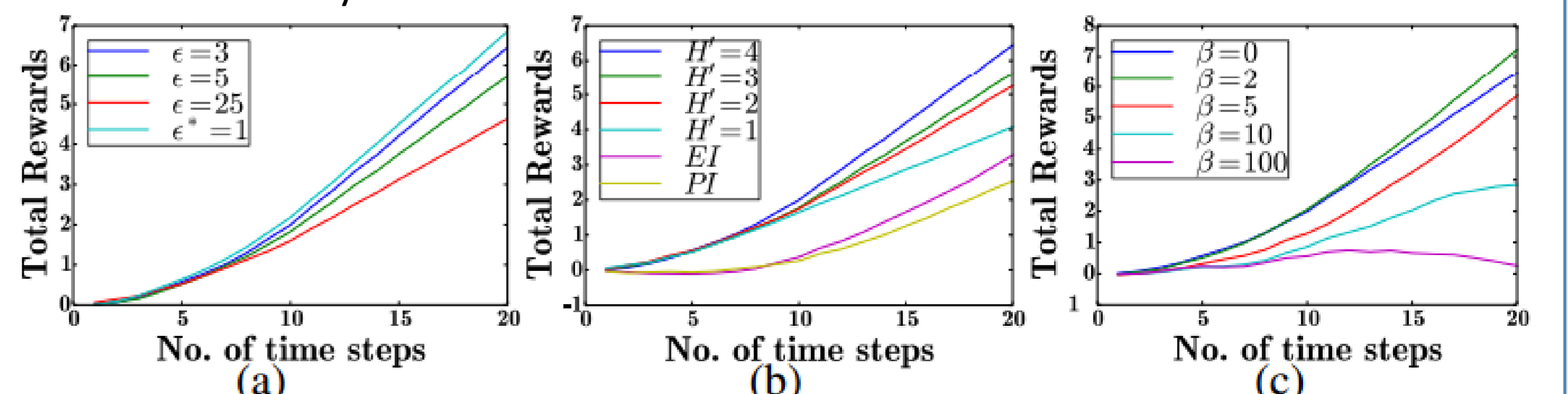
4 Empirical Results

Figure 1: Optimal planning on synthetically generated environments. Graphs of total rewards and tree size of ϵ -GPP policies with (a-b) online planning horizon $H' = 4$ and varying ϵ and (c-d) varying $H' = 1, 2, 3, 4$ (respectively, $\epsilon = 0.002, 0.06, 0.8, 5$) vs. no. of time steps with logarithmic rewards.



The plot of $\epsilon^* = 5$ uses our anytime variant with a maximum tree size of 50000 nodes while the plot of $\epsilon = 250$ effectively assumes maximum likelihood observations during planning (Marchant, 2014).

Figure 2: BO on real world log potassium concentration field. Graphs of total normalized rewards of ϵ -GPP policies using UCB-based rewards with (a) $H' = 4$, $\beta = 0$, and varying ϵ , (b) varying $H' = 1, 2, 3, 4$ (respectively, $\epsilon = 0.002, 0.003, 0.4, 2$) and $\beta = 0$, and (c) $H' = 4$, $\epsilon = 1$, and varying β vs. no. of time steps. The plot of $\epsilon^* = 1$ uses our anytime variant with a maximum tree size of 30000 nodes while the plot of $\epsilon = 25$ effectively assumes maximum likelihood observations.



Observations: Nonmyopic, ϵ -optimal planning improves performance significantly over myopic rewards employing EI/PI. Setting a small value of beta improves performance slightly as exploration is encouraged. However, ϵ -GPP is competitive and does not require tuning of the parameter β .