

TL;DR

Meta-learning for lifelong kernelized bandit optimization with oracle optimal performance guarantees

Overview

- ▷ We solve a sequence of kernelized bandit optimization tasks, where we assume the kernel to be **unknown**, but **shared** across all tasks.
- ▷ We develop LiBO, an algorithm that sequentially meta-learns an approximate kernel and solves the incoming tasks with the latest kernel estimate.
- ▷ Our method pairs with any kernelized bandit algorithm, ensuring oracle optimal performance, meaning that the LiBO's task-specific regret approaches the regret of an algorithm with oracle knowledge of the true kernel over time.
- ▷ We also propose F-LiBO, which solves the lifelong problem in a federated manner.

Problem Setting

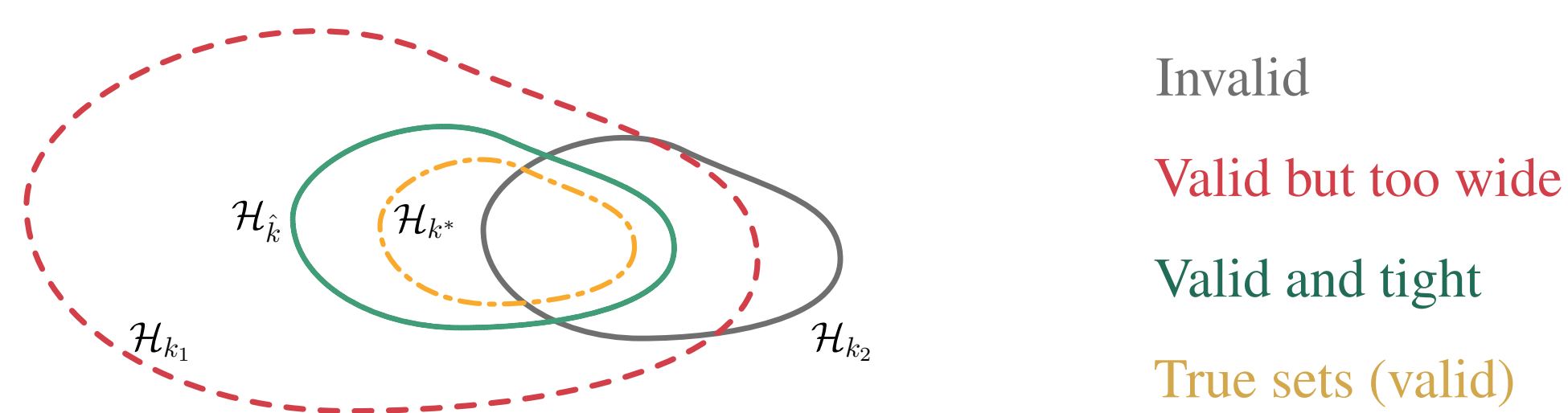
- ▷ Sequentially interacting with the environments

$$y_{s,i} = f_s(\mathbf{x}_{s,i}) + \epsilon_{s,i} \quad 1 \leq i \leq n \text{ and } 1 \leq s \leq m$$

$\epsilon_{s,i}$: also σ^2 sub-Gaussian, i.i.d.

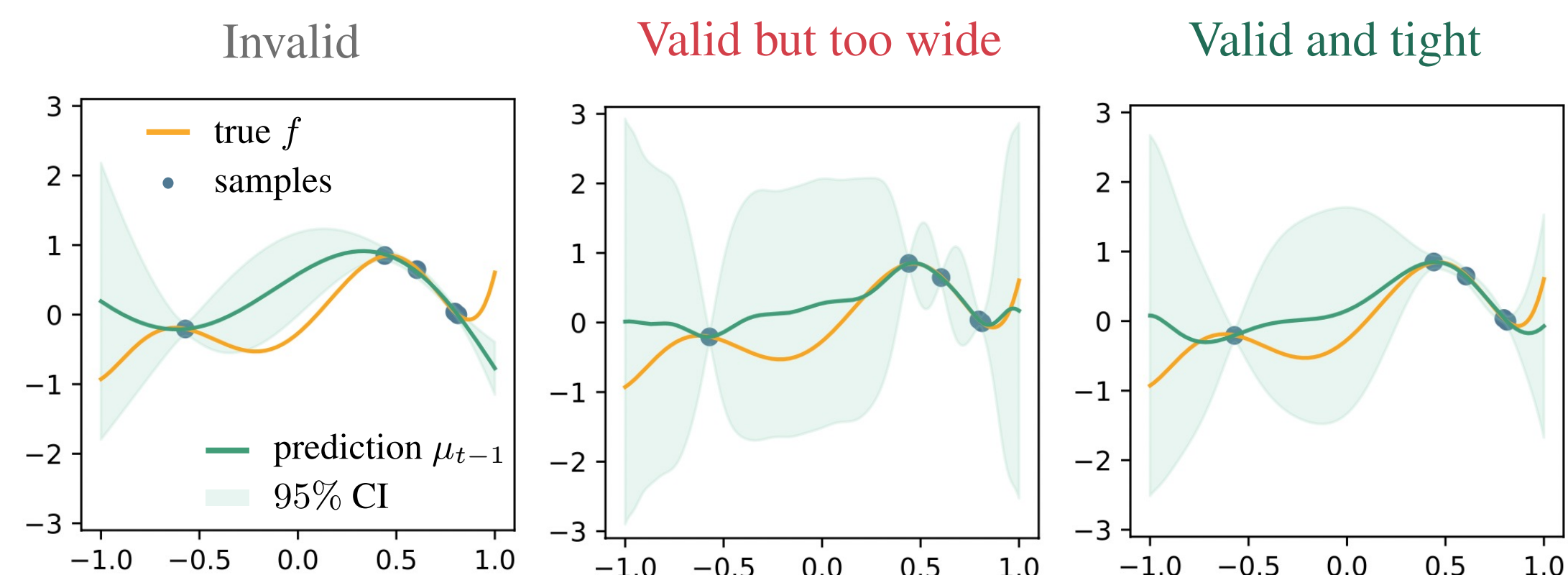
$$f_s : \mathcal{X} \rightarrow \mathbb{R}, f_s \in \mathcal{H}_{k^*}, \|f_s\|_{k^*} \leq B$$

- ▷ Approximating the kernel



Invalid
Valid but too wide
Valid and tight
True sets (valid)

- ▷ We want valid and tight confidence bands for the reward functions such that the bandit solvers converge properly
- ▷ We achieve this with a good kernel approximation



- ▷ Applications: Bandits, Safe BO, Bayesian Optimization, Model-Based RL

Meta-Learning the Kernel

- ▷ Assume the kernel is a linear combination of known base kernels

$$k^*(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^p \eta_j^* k_j(\mathbf{x}, \mathbf{x}') \quad p < \infty$$

η_j^* : unknown, non-negative k_j : known, finite-dimensional

- ▷ We use group Lasso to find the sparsity pattern with offline data

$$\hat{\beta} \implies \hat{J} \implies \hat{k} \implies \mathcal{H}_{\hat{k}}$$

$$\hat{\beta} := \arg \min_{\beta \in \mathbb{R}^{dm}} \left(\frac{1}{mn} \|y - \Phi\beta\|_2^2 + \lambda \sum_{j=1}^p \|\beta^{(j)}\|_2 \right)$$

$$\hat{k} = \frac{1}{|\hat{J}|} \sum_{j \in \hat{J}} k_j \quad \hat{J} := \left\{ j \in \{1, \dots, p\} \mid \|\hat{\beta}^{(j)}\|_2 / \sqrt{m} > \omega \right\}$$

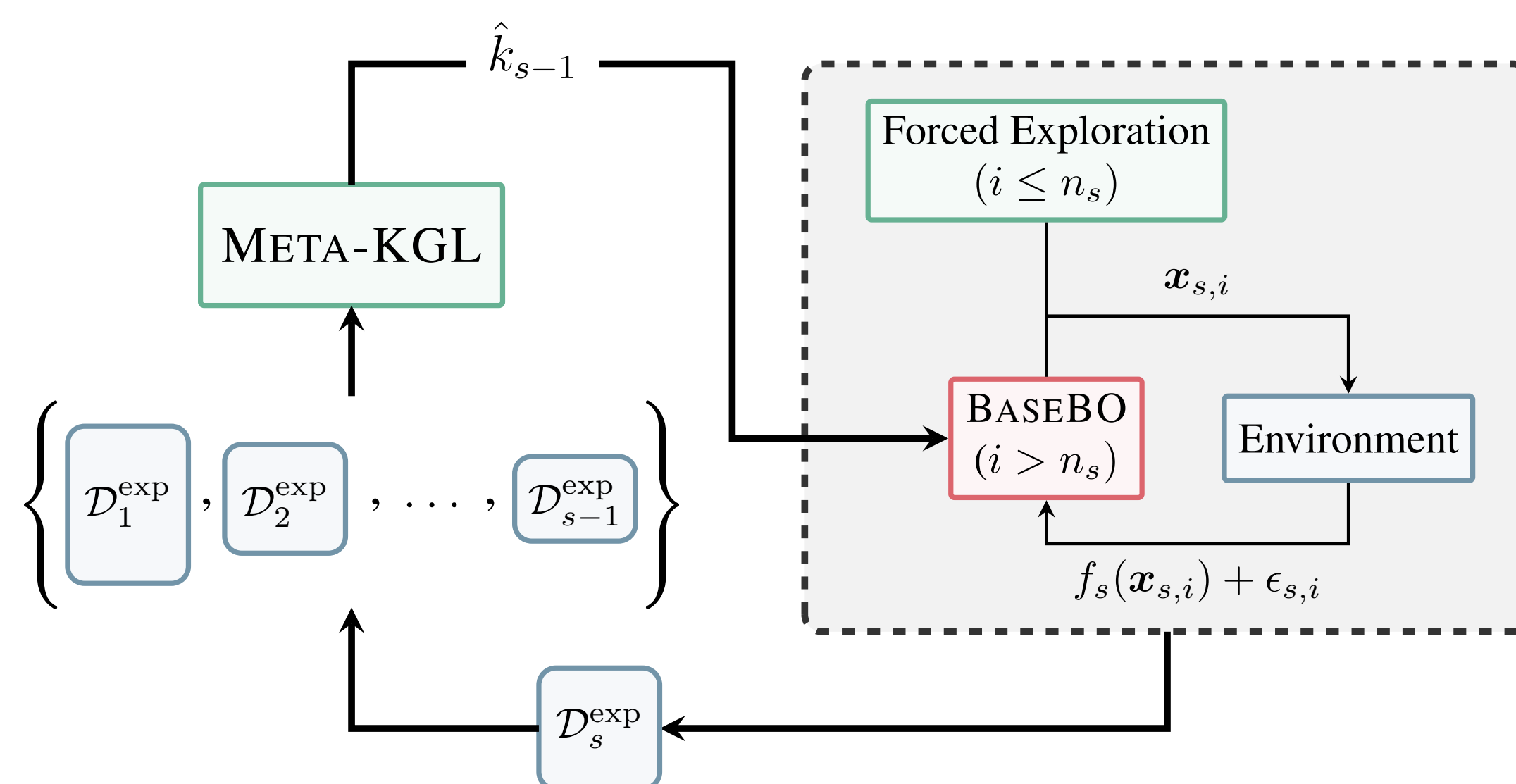
- ▷ Estimator is consistent

Theorem (Informal)

Under mild assumptions on the data model, $\mathcal{H}_{\hat{k}}$ is with high probability a **consistent estimator** in the number of samples n and the number of tasks m . That is,

$$\lim_{n \rightarrow \infty} \mathbb{P}[\mathcal{H}_{\hat{k}} = \mathcal{H}_{k^*}] = 1, \text{ and } \lim_{m \rightarrow \infty} \mathbb{P}[\mathcal{H}_{\hat{k}} = \mathcal{H}_{k^*}] = 1.$$

Lifelong Bandit Optimization



- ▷ Key features
 - Solves a sequence of kernelized bandit tasks in lifelong setting
 - Pairs with any kernelized bandit algorithm
 - Uses forced exploration to improve convergence
- ▷ Guarantee

Theorem (Informal)

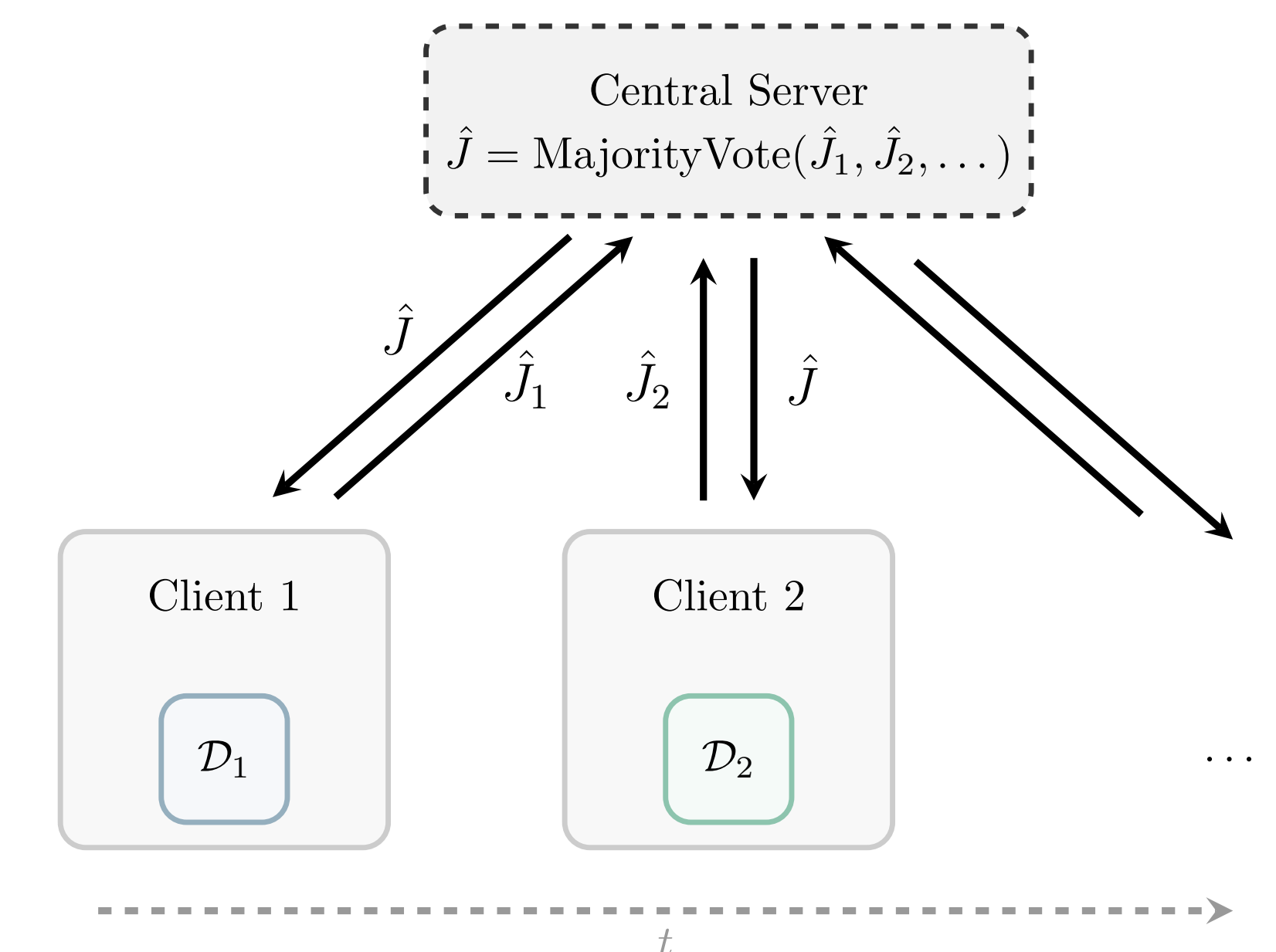
Under mild assumptions, LiBO paired with BASEBO achieves oracle optimal performance with high probability, i.e.

$$R(n, m) = \mathcal{O}(B\sqrt{nm} + R^*(n, m)) = \mathcal{O}(R^*(n, m))$$

where $R^*(n, m)$ is the lifelong regret of BASEBO, if given knowledge of the kernel.

Federated Setting

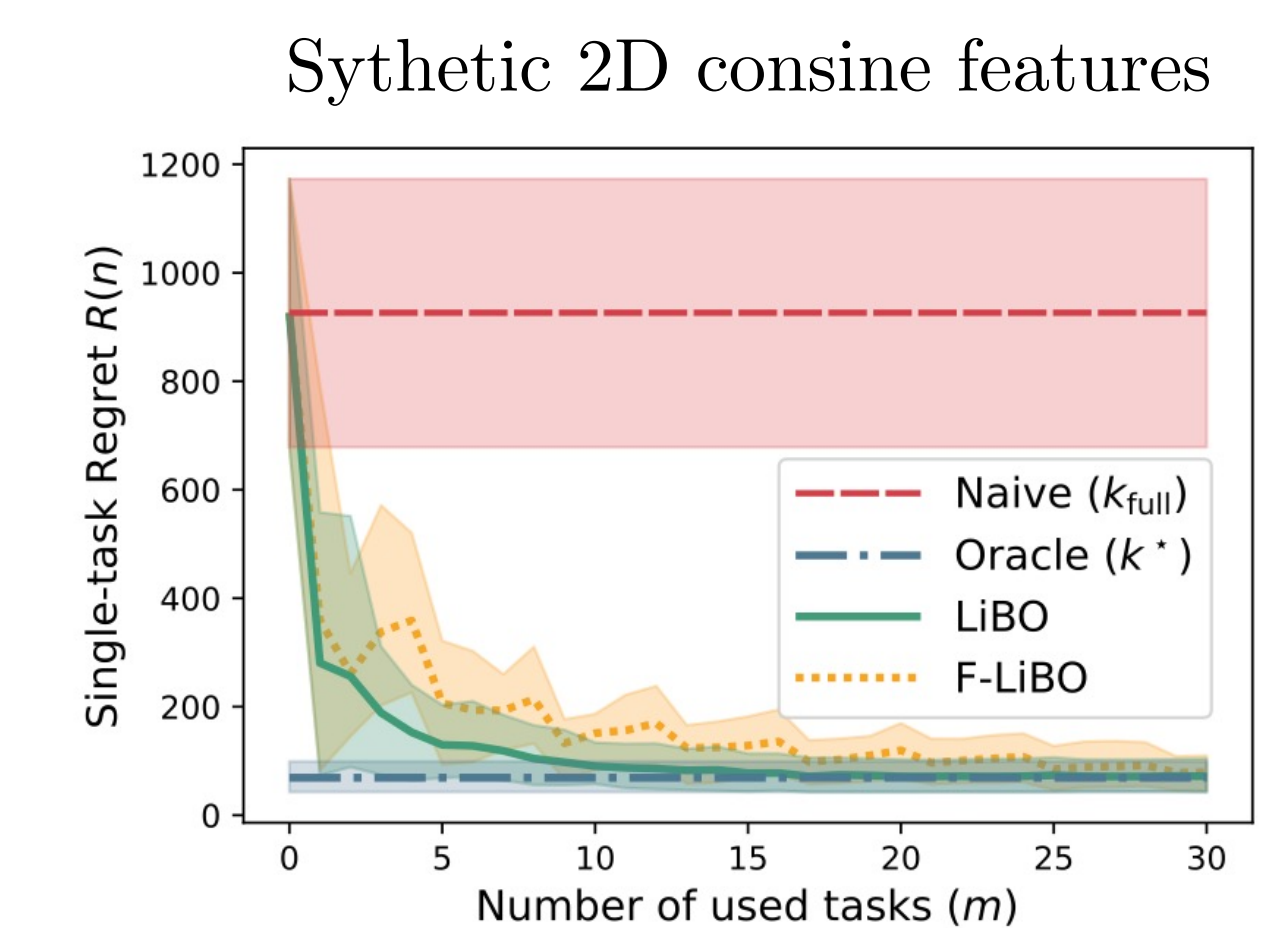
- ▷ We develop F-LiBO for the federated setting:



- ▷ Under similar assumptions we can prove the same convergence guarantees as for the non-federated setting

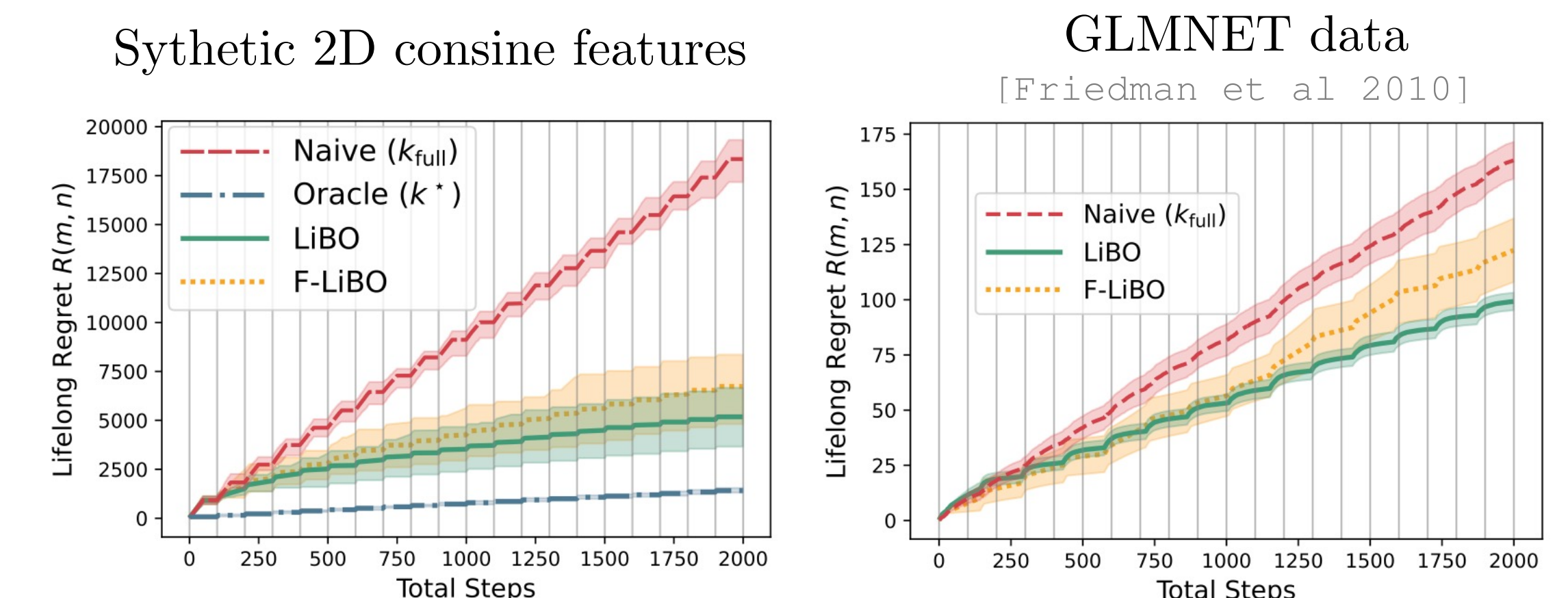
Experiments

- ▷ Offline kernel convergence experiment



Regret of UCB using the estimated kernel quickly converges to the regret achieved with oracle knowledge of the kernel

- ▷ Lifelong bandit optimization



- ▷ Meta-learning with LiBO improves the performance on both synthetic and real world data
- ▷ In the beginning, when kernel estimate is invalid, the regret behaves the same as the naive bandit.
- ▷ Eventually, the regret becomes oracle optimal.