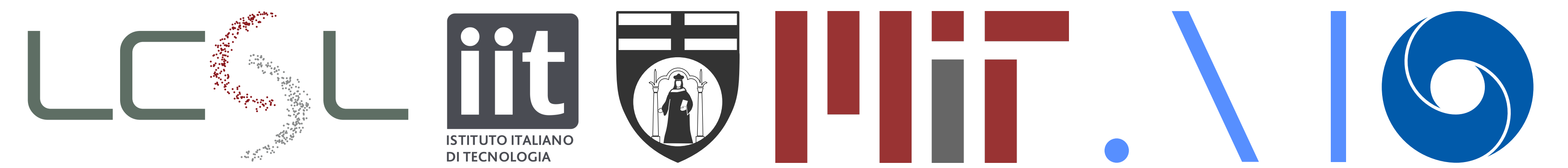


Near-linear Time Gaussian Process Optimization with Adaptive Batching and Resparsification

Daniele Calandriello^{*,†,¶}, Luigi Carratino^{*,‡}, Alessandro Lazaric[§], Michal Valko[¶], Lorenzo Rosasco^{†,‡,¶}

^{*}Equal contribution, [†]LCSL - Istituto Italiano di Tecnologia, [‡]Università degli Studi di Genova,

[§]Facebook AI Research Paris, [¶]DeepMind Paris, ^{||}Massachusetts Institute of Technology

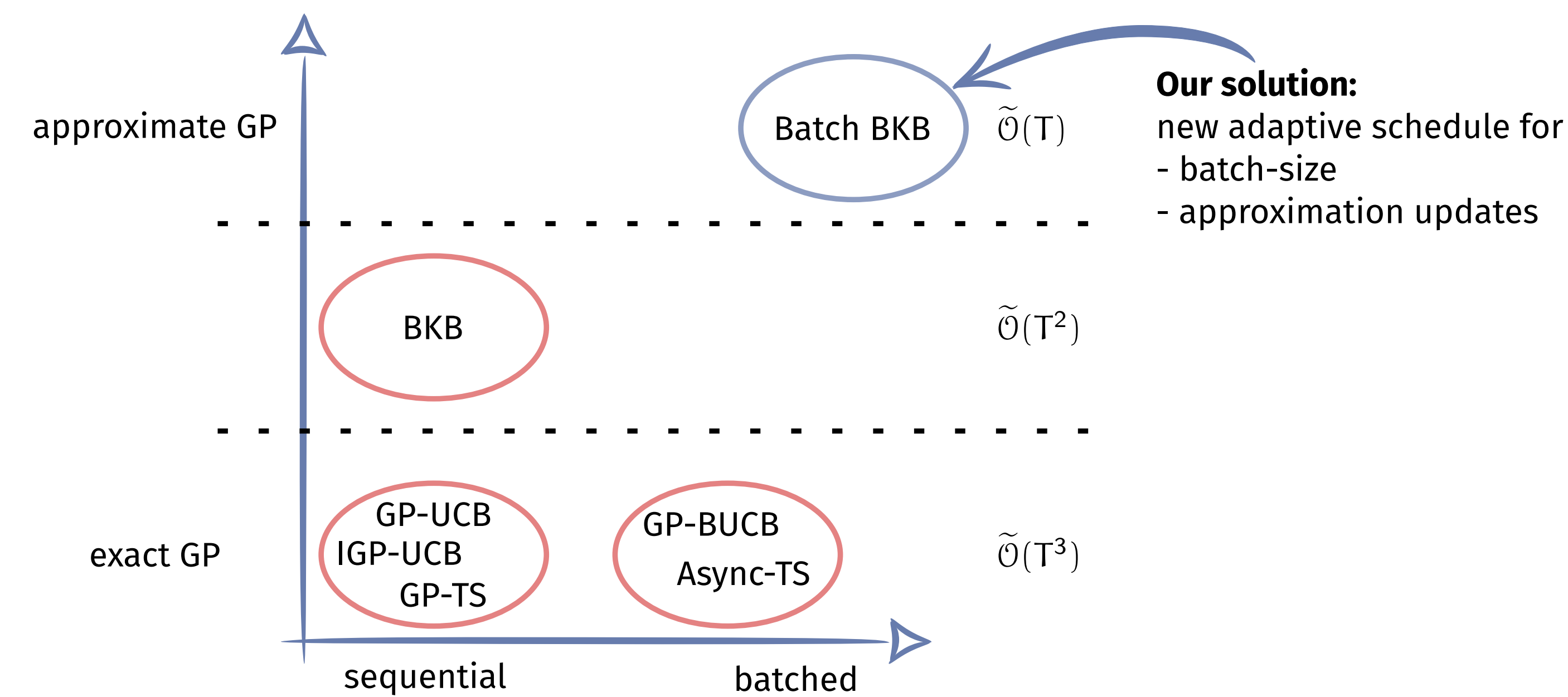


In a nutshell

Gaussian process optimization (GP-Opt) is emerging as a valuable tool for stochastic black-box optimization (e.g. hyperparameter tuning), slowly replacing grad student descent. However provably convergent GP-Opt algorithms (e.g. GP-UCB) suffer from:

- **computational bottleneck**: $\mathcal{O}(T^3)$ time and $\mathcal{O}(T^2)$ space
- **experimental bottleneck**: need feedback at each iteration

We introduce the first **general** GP optimization algorithm (BBKB) that is provably scalable, with **near-linear runtime** $\mathcal{O}(Td_{\text{eff}}^2)$, **no regret**, and also maintains **valid posterior variance estimates** at all steps.



Batch Budgeted Kernelized Bandits (BBKB)

BKB **for** $t = \{1, \dots, T-1\}$ **do**
for $i = \{1, \dots, A\}$ **do**
 $\tilde{u}_t(\mathbf{x}_i) = \tilde{\mu}_t(\mathbf{x}_i, \mathcal{D}_t) + \tilde{\beta}_t \tilde{\sigma}_t^2(\mathbf{x}_i, \mathcal{D}_t)$;
end
 Select $\mathbf{x}_{t+1} \leftarrow \arg \max_{\mathbf{x}_i \in A} \tilde{u}_t(\mathbf{x}_i)$;
 Set $\tilde{\mathbf{p}}_{t+1} \propto [\tilde{\sigma}_t^2(\mathbf{x}_1, \mathcal{D}_t), \dots, \tilde{\sigma}_t^2(\mathbf{x}_{t+1}, \mathcal{D}_t)]$;
 Sample $\mathcal{D}_{t+1} \sim \tilde{\mathbf{p}}_{t+1}$;

end

BBKB At step t start batch and for $t' > t$ continue selecting

$$\tilde{u}_{t'}(\mathbf{x}_i) = \tilde{\mu}_{t'}(\mathbf{x}_i, \mathcal{D}_t) + \tilde{\alpha}_t \tilde{\sigma}_{t'}^2(\mathbf{x}_i, \mathcal{D}_t)$$

while $\sum_{s=t}^{t'} \tilde{\sigma}_s^2(\mathbf{x}_s, \mathcal{D}_t) \leq 1$.

Main result: BBKB is scalable and no regret

Theorem. Let $\tilde{\alpha}_t \approx \sqrt{\sum_{s=1}^t \log(1 + 2\tilde{\sigma}_{s-1}^2(\tilde{\mathbf{x}}_s, \mathcal{D}_{s-1}))} + \sqrt{\lambda}F$. Assume $f \in \mathcal{H}$ arbitrary and $\|f\| \leq F$. Then w.h.p., $\forall t \in [T]$ and all $\mathbf{x} \in \mathcal{A}$,
 $\sigma_t^2(\mathbf{x})/2 \leq \tilde{\sigma}_t^2(\mathbf{x}) \leq 2\sigma_t^2(\mathbf{x})$ and $|\mathcal{D}_t| \leq \mathcal{O}(d_{\text{eff}} \log(t/\delta))$,
 and BBKB incurs at most regret $R_T^{\text{BBKB}} \leq 4R_T^{\text{GP-UCB}}$
 and runs in $\mathcal{O}(TAd_{\text{eff}}^2)$ time with average batch size $P \geq \Omega(T/d_{\text{eff}})$.

- 😊 adapts to GP's effective dimension/rank: $d_{\text{eff}} \triangleq \sum_{i=1}^T \sigma_T^2(\tilde{\mathbf{x}}_i)$
- 😊 $\tilde{\beta}_t$ **computable** in $\tilde{\mathcal{O}}(Ad_{\text{eff}}^2)$ time replacing worst-case bounds
- 😊 **No assumptions** on GP structure (e.g., not only stationary GPs)
- 😊 **No free lunch**: worst-case falls back to GP-UCB

- 😊 DTC is **not a GP** (not consistent), but now a **justified** heuristic
- 😊 Easy **extension to infinite \mathcal{A}** , but **how to optimize posterior?**

🤖 Big batches and rare resparsification with no loss of accuracy?

"Not too big" Lemma:

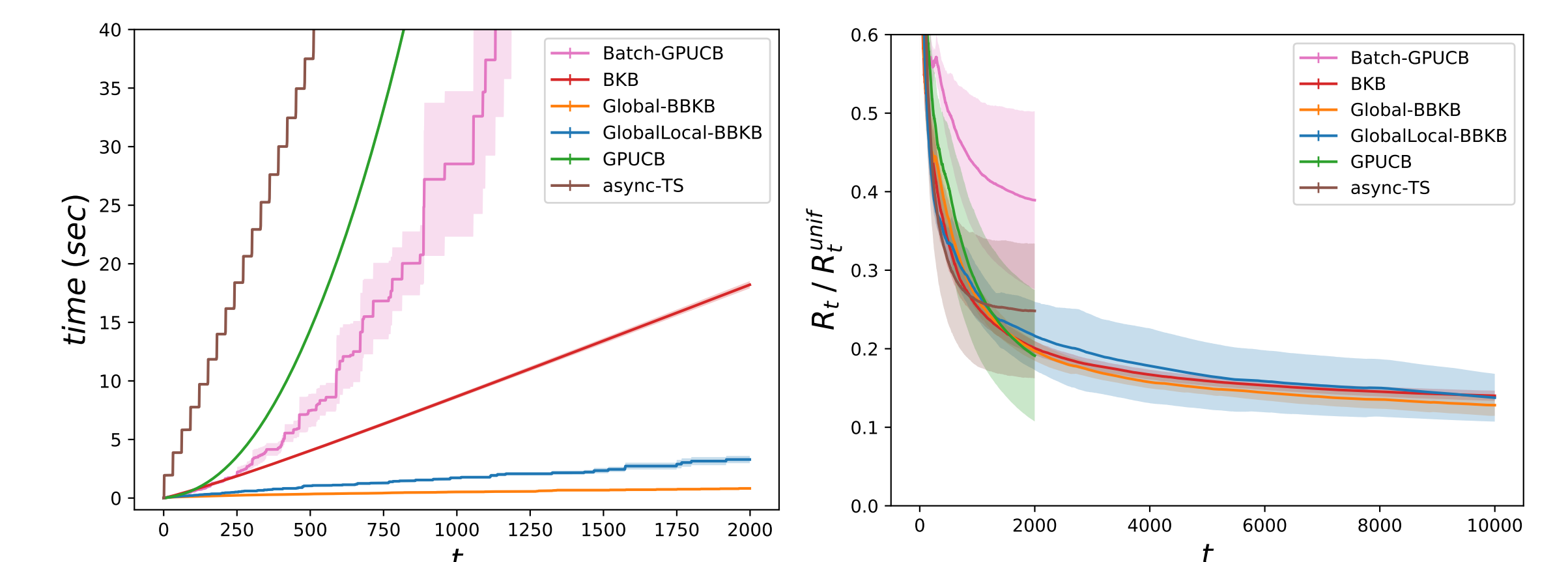
↳ valid UCBs

vs "Not too small" Lemma:

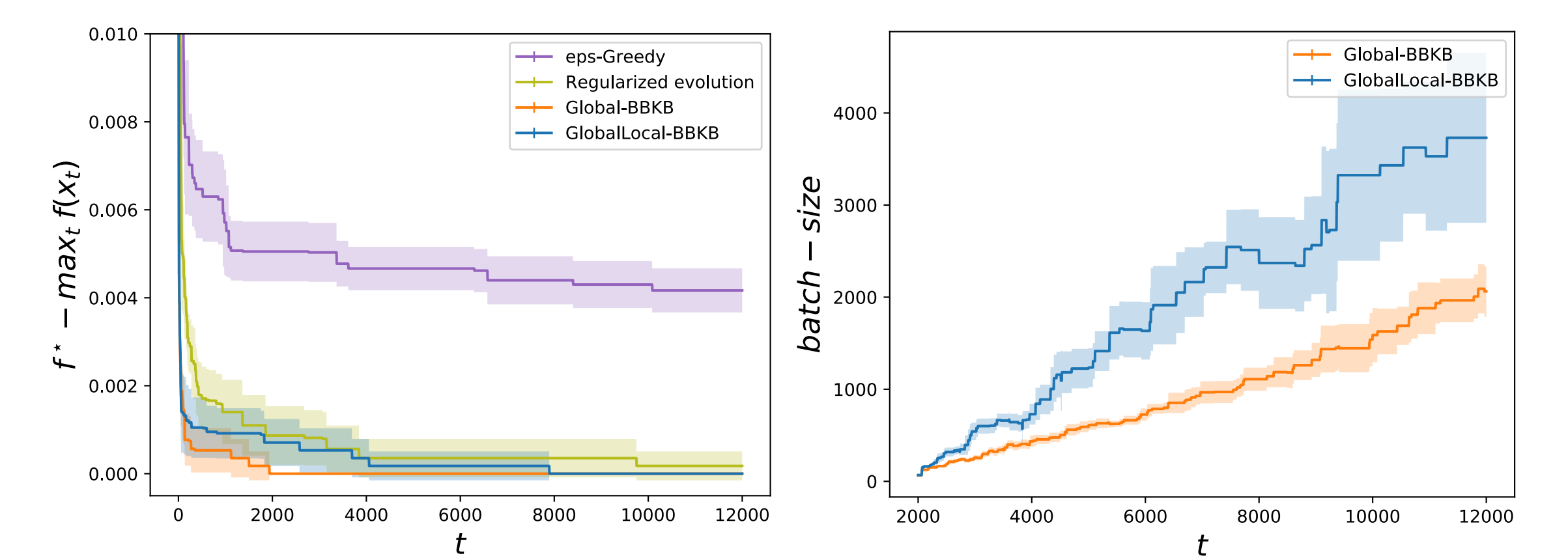
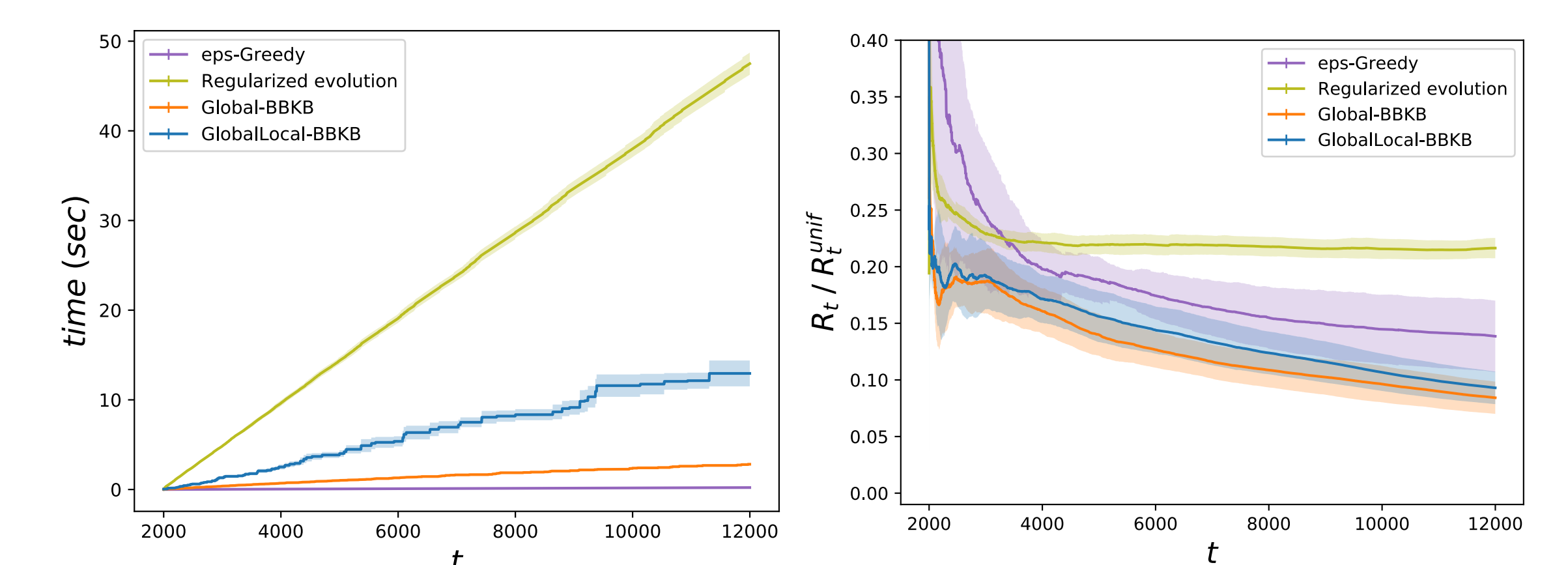
↳ batch-size = $\Omega(t)$

Experiments

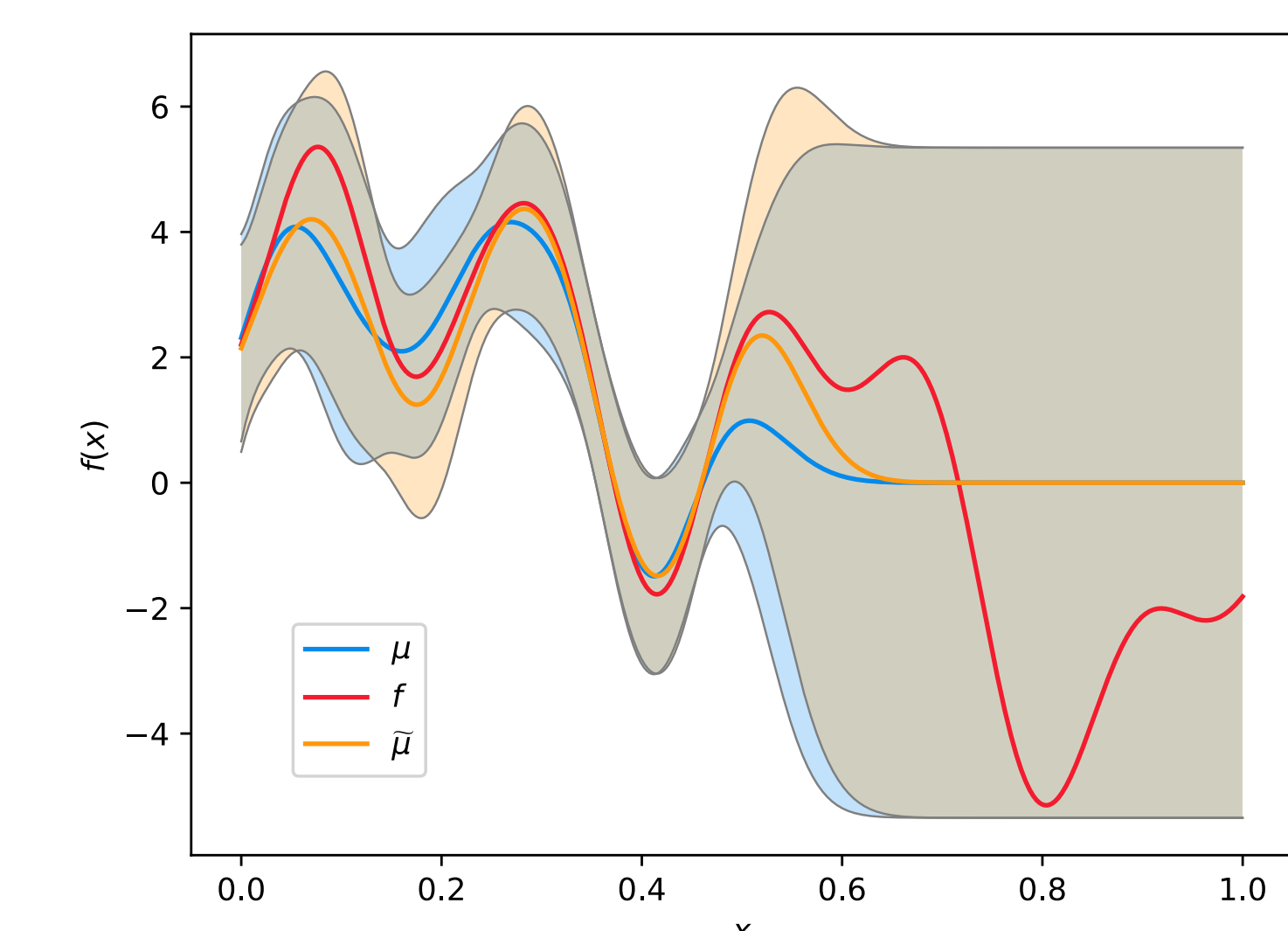
Cadata dataset ($A = 20640, d = 8, T = 10000$)



NAS-bench-101 dataset ($A = 12416, d = 19, T = 12000$)



Not simply an approximate GP-UCB



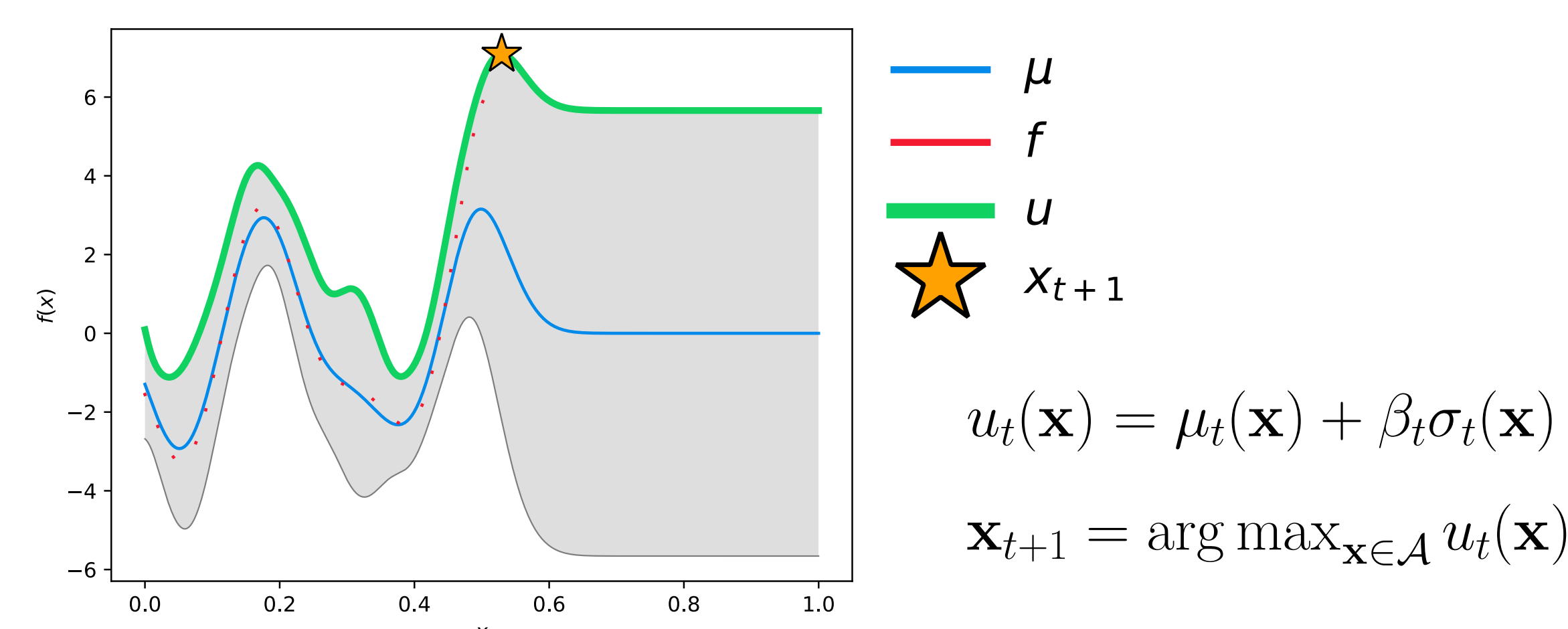
Sparse approximation **regularizes!** New bias-variance trade-off!

Gaussian process optimization and GP-UCB

Candidates $\mathcal{A} = \{\mathbf{x}_i\}_{i=1}^A$ with $\mathbf{x}_i \in \mathbb{R}^D$ (or RKHS \mathcal{H})

- For $t \in [1, \dots, T]$:
- (1) select $\mathbf{x}_{t+1} = \arg \max_{\mathbf{x}_i} u_t(\mathbf{x}_i)$
 - (2) Receive noisy feedback $y_{t+1} = f(\mathbf{x}_{t+1}) + \eta_{t+1}$
 - (3) Improve u_{t+1} for next time

Goal: minimize **regret** $R_T = \sum_{t=1}^T f(\mathbf{x}_*) - f(\mathbf{x}_t)$ vs. $\mathbf{x}_* = \arg \max_{\mathbf{x}_i} f(\mathbf{x}_i)$



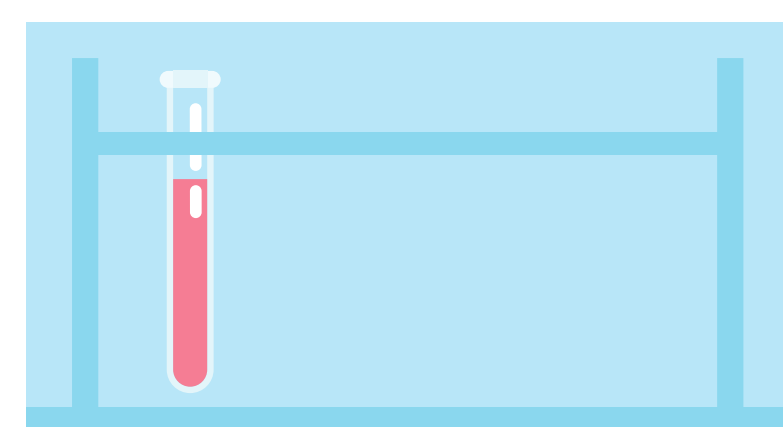
— μ
 — f
 — u
 ★ \mathbf{x}_{t+1}

$$u_t(\mathbf{x}) = \mu_t(\mathbf{x}) + \beta_t \sigma_t(\mathbf{x})$$

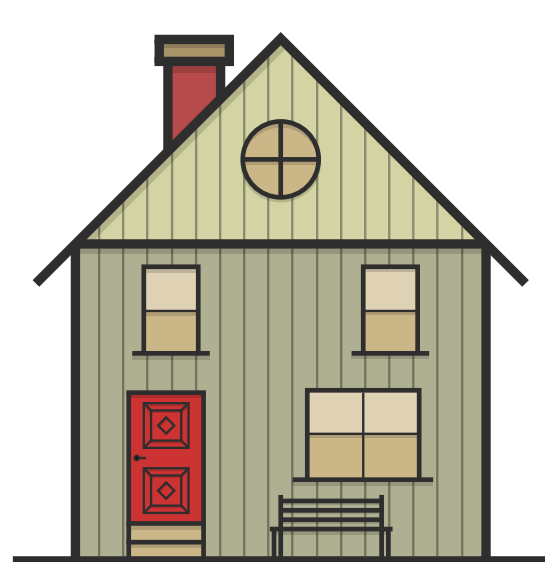
$$\mathbf{x}_{t+1} = \arg \max_{\mathbf{x} \in \mathcal{A}} u_t(\mathbf{x})$$

$$\mu_t(\mathbf{x}) = \mathbf{x}^T (\mathbf{X}_t^T \mathbf{X}_t + \lambda \mathbf{I})^{-1} \mathbf{X}_t^T \mathbf{y}_t, \quad \sigma_t^2(\mathbf{x}) = \mathbf{x}^T (\mathbf{X}_t^T \mathbf{X}_t + \lambda \mathbf{I})^{-1} \mathbf{x}$$

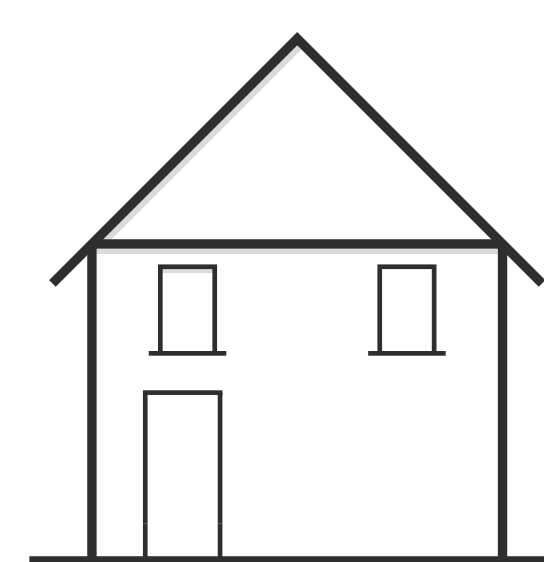
🤖 How to improve bottlenecks and still achieve low regret? 🤖



sequential vs batch



exact GP vs approximate GP



💡 Sparse GP with dictionary $\mathcal{D} = \{\mathbf{x}_j\}_{j=1}^m$ of m inducing points