Asynchronous Heterogeneous Linear Quadratic Regulator Design

### Leonardo F. Toso\*, Han Wang\*, James Anderson

### Department of Electrical Engineering, Columbia University

February 28, 2025



Conference on Decision and Control (CDC), 2024

\*equal contribution

1

### Single-agent model-free optimal control

$$x_{t+1} = A_{\star} x_t + B_{\star} u_t \quad t \ge 0$$

 $\begin{array}{l} \underline{\text{Objective:}} \\ \text{a cumulative cost } \mathcal{J}(K) = \mathbf{E}_{\mathbf{x}_0} \sum_{t=0}^{\infty} c_t(\mathbf{x}_t, u_t). \end{array}$ 

<u>Difficulty:</u> Do not have access to the ground-truth system dynamics, i.e.,  $(A_{\star}, B_{\star})$ .

### Single-agent model-free optimal control

$$x_{t+1} = A_{\star} x_t + B_{\star} u_t \quad t \ge 0$$

 $\begin{array}{l} \underline{\text{Objective:}} \text{ Design a control sequence } \{u_t = -Kx_t\}_t \text{ that minimizes} \\ \text{a cumulative cost } \mathcal{J}(K) = \mathbf{E}_{x_0} \sum_{t=0}^{\infty} c_t(x_t, u_t). \end{array}$ 

<u>Difficulty</u>: Do not have access to the ground-truth system dynamics, i.e.,  $(A_{\star}, B_{\star})$ .

Policy Gradient (PG)  $\downarrow$ <u>Collect data</u>:  $D_K := \{x_t, u_t\}_t \sim D_K$   $\downarrow$ <u>Estimation</u>:  $\hat{\nabla}_{D_K} = ZO(D_K)$   $\downarrow$ <u>Control</u>: PG updates  $K \leftarrow K - \eta \hat{\nabla}_{D_K}$   $\downarrow$ <u>Sample Complexity</u>:  $J(K) - J(K^*) \le e$ , # data points  $\approx O(1/e^2)$ 

### Single-agent model-free optimal control

$$x_{t+1} = A_{\star} x_t + B_{\star} u_t \quad t \ge 0$$

<u>Objective</u>: Design a control sequence  $\{u_t = -Kx_t\}_t$  that minimizes a cumulative cost  $\mathcal{J}(K) = \mathbf{E}_{x_0} \sum_{t=0}^{\infty} c_t(x_t, u_t)$ .

<u>Difficulty</u>: Do not have access to the ground-truth system dynamics, i.e.,  $(A_{\star}, B_{\star})$ .

Policy Gradient (PG)

Collect data: 
$$D_K := \{x_t, u_t\}_t \sim \mathcal{D}_K$$

Estimation:  $\hat{\nabla}_{D_K} = ZO(D_K)$ 



Necessity for a large amount of data samples



-

### Single-agent model-free optimal control

$$x_{t+1} = A_{\star} x_t + B_{\star} u_t \quad t \ge 0$$

<u>Objective</u>: Design a control sequence  $\{u_t = -Kx_t\}_t$  that minimizes a cumulative cost  $\mathcal{J}(K) = \mathbf{E}_{x_0} \sum_{t=0}^{\infty} c_t(x_t, u_t)$ .

<u>Difficulty</u>: Do not have access to the ground-truth system dynamics, i.e.,  $(A_{\star}, B_{\star})$ .

Policy Gradient (PG)

Collect data: 
$$D_K := \{x_t, u_t\}_t \sim \mathcal{D}_K$$

Estimation:  $\hat{\nabla}_{D_K} = ZO(D_K)$ 





Homogeneous  $x_{t+1} = A_{\star} x_t + B_{\star} u_t \quad t \ge 0$  $\hat{\nabla}_{D_{K}^{(M)}}$  $\hat{\nabla}_{D_{\nu}^{(1)}}$  $\frac{1}{M}\sum_{i=1}^{M}\hat{\nabla}_{D_{K}^{(i)}}$ Law of Large Numbers (LLN)  $\mathcal{J}(K) - \mathcal{J}(K^{\star}) \leq \epsilon$ # data points  $\approx \mathcal{O}(1/M\epsilon^2)$ 









### Challenges and Goals

Asynchronous LQR design



### Challenges and Goals

Asynchronous LQR design



**Question:** Can an asynchronous algorithm produce a controller that is near-optimal, even in the presence of staleness and heterogeneous system dynamics?

## Challenges and Goals

Asynchronous LQR design



**Question:** Can an asynchronous algorithm produce a controller that is near-optimal, even in the presence of staleness and heterogeneous system dynamics?

Linear Quadratic Regulator (LQR) Objective

Consider a discrete LTI dynamical system

$$x_{t+1}^{(i)} = A^{(i)}x_t^{(i)} + B^{(i)}u_t^{(i)}, \quad t = 0, 1, 2, \dots$$
 (sys-dyn)

where  $x_t^{(i)} \in \mathbb{R}^{n_x}$  and  $u_t^{(i)} \in \mathbb{R}^{n_u}$ .

Linear Quadratic Regulator (LQR) Objective

Consider a discrete LTI dynamical system

$$x_{t+1}^{(i)} = A^{(i)} x_t^{(i)} + B^{(i)} u_t^{(i)}, \quad t = 0, 1, 2, \dots$$
 (sys-dyn)

where  $x_t^{(i)} \in \mathbb{R}^{n_x}$  and  $u_t^{(i)} \in \mathbb{R}^{n_u}$ .

LQR Objective: Design a controller  $K_i^{\star}$   $(u_t^{(i)} = -K_i^{\star} x_t^{(i)})$  that solves

$$\mathcal{K}_{i}^{\star} = \operatorname{argmin}_{\mathcal{K} \in \mathcal{K}^{(i)}} \mathcal{J}^{(i)}(\mathcal{K}) = \mathbb{E}\left[\sum_{t=0}^{\infty} x_{t}^{(i)\top} \left(Q^{(i)} + \mathcal{K}^{\top} R^{(i)} \mathcal{K}\right) x_{t}^{(i)}\right],$$

subject to (sys-dyn).

Stabilizing set:  $\mathcal{K}^{(i)} = \{ \mathcal{K} \mid \rho(\mathcal{A}^{(i)} - \mathcal{B}^{(i)}\mathcal{K}) < 1 \}$ 

**Global convergence** despite of the non-convexity<sup>\*</sup> of  $\mathcal{J}^{(i)}(\mathcal{K})$  [FGKM,ICML 2018],

Initial Stabilizing Controller:  $K_0 \in \mathcal{K}^{(i)}$ 

**Controllability:**  $(A^{(i)}, B^{(i)})$  is controllable

Global convergence despite of the non-convexity<sup>\*</sup> of  $\mathcal{J}^{(i)}(\mathcal{K})$  [FGKM,ICML 2018], Initial Stabilizing Controller:  $\mathcal{K}_0 \in \mathcal{K}^{(i)}$ 

**Controllability:**  $(A^{(i)}, B^{(i)})$  is controllable

$$K_{n+1} = K_n - \eta \widehat{\nabla} \mathcal{J}^{(i)}(K_n), \text{ for } n \in \{0, 1, \dots, N-1\},$$

$$\mathcal{J}^{(i)}(K_N) - \mathcal{J}^{(i)}(K_i^{\star}) \leq \epsilon_i$$

after  $N = \mathcal{O}(\log(1/\epsilon))$  iterations

\* Non-convex for  $n_X \ge 3$ .

Zeroth-order Estimation:  $K_{n+1} = K_n - \eta \widehat{\nabla} \mathcal{J}^{(i)}(K_n)$ 

$$\operatorname{ZO}(m,r,K) \to \widehat{\nabla} \mathcal{J}^{(i)}(K) := \sum_{l=1}^{m} \frac{n_{X} n_{u} (\mathcal{J}^{(i)}(K+U_{l}) - \mathcal{J}^{(i)}(K-U_{l})) U_{l}}{2mr^{2}},$$

*m* (number of trajectories), *r* (smoothing radius) and  $||U_l||_F = r$ .



Zeroth-order Estimation:  $K_{n+1} = K_n - \eta \widehat{\nabla} \mathcal{J}^{(i)}(K_n)$ 

$$\operatorname{ZO}(m,r,K) \to \widehat{\nabla} \mathcal{J}^{(i)}(K) := \sum_{l=1}^{m} \frac{n_{X} n_{u} (\mathcal{J}^{(i)}(K+U_{l}) - \mathcal{J}^{(i)}(K-U_{l})) U_{l}}{2mr^{2}},$$

*m* (number of trajectories), *r* (smoothing radius) and  $||U_l||_F = r$ .

**Estimation Error**<sup>\*</sup>: Suppose  $m = O(1/\epsilon^2)$  and  $r = O(\epsilon)$ , it holds that

$$\|\nabla \mathcal{J}^{(i)}(K) - \widehat{\nabla} \mathcal{J}^{(i)}(K)\| \leq \epsilon,$$

with high probability.

<sup>\*</sup> Bernstein matrix inequality (Tropp (2012)).

## Multi-Agent Heterogeneous LQR design Objective

Consider *M* distinct systems (sys-dyn) with different LQR objectives, i.e.,

**Objective:** 
$$\bar{K}^{\star} := \operatorname{argmin}_{K \in \mathcal{S}} \left\{ \bar{\mathcal{J}}(K) := \frac{1}{M} \sum_{i=1}^{M} \mathcal{J}^{(i)}(K) \right\},$$

 $\bar{K}^{\star}$  should stabilize each (sys-dyn) and on *average* minimize their LQR objectives

### Multi-Agent Heterogeneous LQR design Objective

Consider *M* distinct systems (sys-dyn) with different LQR objectives, i.e.,

$$\textbf{Objective:} \ \bar{K}^{\star} := \operatorname{argmin}_{K \in \mathcal{S}} \left\{ \bar{\mathcal{J}}(K) := \tfrac{1}{M} \sum_{i=1}^{M} \mathcal{J}^{(i)}(K) \right\},$$

 $\bar{K}^{\star}$  should stabilize each (sys-dyn) and on *average* minimize their LQR objectives Heterogeneity:  $\max_{i \neq j} ||A^{(i)} - A^{(j)}|| \le \epsilon_A \rightarrow \text{same for } \epsilon_B, \epsilon_Q \text{ and } \epsilon_R$ 

### Multi-Agent Heterogeneous LQR design Objective

Consider *M* distinct systems (sys-dyn) with different LQR objectives, i.e.,

**Objective:** 
$$\bar{K}^{\star} := \operatorname{argmin}_{K \in \mathcal{S}} \left\{ \bar{\mathcal{J}}(K) := \frac{1}{M} \sum_{i=1}^{M} \mathcal{J}^{(i)}(K) \right\},$$

 $\bar{K}^{\star}$  should stabilize each (sys-dyn) and on *average* minimize their LQR objectives

**Heterogeneity:**  $\max_{i \neq j} ||A^{(i)} - A^{(j)}|| \le \epsilon_A \rightarrow \text{ same for } \epsilon_B, \epsilon_Q \text{ and } \epsilon_R$ 

Stabilizing sub-level set:  $S^{(i)} := \left\{ K \mid \mathcal{J}^{(i)}(K) - \mathcal{J}^{(i)}(K_i^{\star}) \leq \gamma(\mathcal{J}^{(i)}(K_0) - \mathcal{J}^{(i)}(K_i^{\star})) \right\},$ 



## Multi-Agent Heterogeneous LQR design

Gradient heterogeneity

**Gradient Heterogeneity:** [TZAW, L4DC 2024] For any two distinct systems with different LQR objectives, and given a stabilizing controller  $K \in S$ . It holds that,

 $\|\nabla \mathcal{J}^{(i)}(K) - \nabla \mathcal{J}^{(j)}(K)\|^2 \le f(\epsilon_A, \epsilon_B, \epsilon_Q, \epsilon_R) = \epsilon_{\mathsf{het}} \quad (\mathsf{grad-het})$ 

## Multi-Agent Heterogeneous LQR design

Gradient heterogeneity

**Gradient Heterogeneity:** [TZAW, L4DC 2024] For any two distinct systems with different LQR objectives, and given a stabilizing controller  $K \in S$ . It holds that,

 $\|\nabla \mathcal{J}^{(i)}(K) - \nabla \mathcal{J}^{(j)}(K)\|^2 \le f(\epsilon_A, \epsilon_B, \epsilon_Q, \epsilon_R) = \epsilon_{\mathsf{het}} \quad (\mathsf{grad-het})$ 



FedLQR [WTMA, 2023]

### Asynchronous updates



**Bounded staleness:**  $\tau_s(n) \leq \tau_{\max} \in \mathbb{N}$  denotes the staleness in the controller that system  $s \in [b_s]$  possesses when locally estimating its policy gradient at step n

# Asynchronous Policy Gradient LQR Algorithm

Input:  $\overline{K}_0 \in S$ , N,  $\eta$ , m, r,  $b_s$ Initialize:  $K_i = \overline{K}_0 \ \forall i \in [M]$  and  $\overline{\nabla} \leftarrow 0$  (1) Iteration and batch counters: s = n = 0

In parallel: compute and send  $\widehat{\nabla}_i = \text{ZO}(K_i, r, m)$  to the server  $\forall i \in [M]$  (2)

While n < NServer accumulates  $\overline{\nabla} = \overline{\nabla} + \nabla_i$ , s + = 1If  $s = b_s$  then  $\overline{K}_{n+1} = \overline{K}_n - \frac{\eta}{b_s}\overline{\nabla}$  (3) Send  $\overline{K}_{n+1}$  to the idle systems (4)

Output:  $\bar{K}_N$ 



Interplay between heterogeneity and staleness

Controlling the staleness effect:  $\mathbb{E} \| \bar{K}_n - \bar{K}_{n-\tau_i(n)} \|^2$ 



10

Interplay between heterogeneity and staleness

Controlling the staleness effect:  $\mathbb{E} \| \bar{K}_n - \bar{K}_{n-\tau_i(n)} \|^2$ 

$$\mathbb{E} \left\| \bar{K}_n - \bar{K}_{n-\tau_i(n)} \right\|^2 \le \tau_{\max} \sum_{l=n-\tau_i(n)}^{n-1} \mathbb{E} \left\| \bar{K}_{l+1} - \bar{K}_l \right\|^2$$



heterogeneity + staleness

Interplay between heterogeneity and staleness

Controlling the staleness effect:  $\mathbb{E} \| \bar{K}_n - \bar{K}_{n-\tau_i(n)} \|^2$ 

$$\mathbb{E} \left\| \bar{K}_{n} - \bar{K}_{n-\tau_{i}(n)} \right\|^{2} \leq \tau_{\max} \sum_{l=n-\tau_{i}(n)}^{n-1} \mathbb{E} \left\| \bar{K}_{l+1} - \bar{K}_{l} \right\|^{2}$$

$$\mathbb{E}\left\|\bar{K}_{l+1} - \bar{K}_{l}\right\|^{2} \leq \eta^{2}\tau_{\mathsf{max}}\mathcal{O}\left(\epsilon_{\mathsf{het}} + \mathbb{E}\|\nabla\mathcal{J}^{(i)}(\bar{K}_{n})\|^{2}\right)^{*}$$



heterogeneity + staleness

\* See proof of Lemma 4.

Interplay between heterogeneity and staleness

Controlling the staleness effect:  $\mathbb{E} \| \bar{K}_n - \bar{K}_{n-\tau_i(n)} \|^2$ 

$$\mathbb{E} \left\| \bar{K}_{n} - \bar{K}_{n-\tau_{i}(n)} \right\|^{2} \leq \tau_{\max} \sum_{l=n-\tau_{i}(n)}^{n-1} \mathbb{E} \left\| \bar{K}_{l+1} - \bar{K}_{l} \right\|^{2}$$

$$\mathbb{E} \left\| \bar{K}_{l+1} - \bar{K}_{l} \right\|^{2} \leq \eta^{2} \tau_{\max} \mathcal{O} \left( \epsilon_{\mathsf{het}} + \mathbb{E} \| \nabla \mathcal{J}^{(i)}(\bar{K}_{n}) \|^{2} \right)^{*}$$

$$\frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E} \left\| \bar{K}_{n} - \bar{K}_{n-\tau_{i}(n)} \right\|^{2} \leq \frac{\eta^{2} \tau_{\max}^{3} \epsilon_{\mathsf{het}}}{b_{s}} + \eta^{2} \tau_{\max}^{3} \mathbb{E} \| \nabla \bar{\mathcal{J}}(\bar{K}_{n}) \|^{2}$$

\* See proof of Lemma 4.

Convergence guarantees - ergodic convergence rate

Goal: Control  $\frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E} \|\nabla \bar{\mathcal{J}}(\bar{K}_n)\|^2 \to \text{find a stationary solution}$ 

• Initial stabilizing controller:  $\bar{K}_0 \in \mathcal{S}$ 

•  $\eta = \mathcal{O}\left(\sqrt{rac{b_s}{N}}
ight)$  and r sufficiently small

Convergence guarantees - ergodic convergence rate

Goal: Control  $\frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E} \|\nabla \bar{\mathcal{J}}(\bar{K}_n)\|^2 \to \text{find a stationary solution}$ 

• Initial stabilizing controller:  $\bar{K}_0 \in \mathcal{S}$ 

•  $\eta = \mathcal{O}\left(\sqrt{rac{b_s}{N}}
ight)$  and r sufficiently small

$$\frac{1}{N}\sum_{n=0}^{N-1} \mathbb{E} \|\nabla \bar{\mathcal{J}}(\bar{K}_n)\|_F^2 \leq \mathcal{O}\left(\frac{\bar{\Delta}_0}{\sqrt{Nb_s}} + \frac{\epsilon_{\mathsf{het}}}{\sqrt{Nb_s}} + \frac{\tau_{\mathsf{max}}^2 \epsilon_{\mathsf{het}}}{N}\right)$$

•  $\bar{\Delta}_0 = \mathbb{E}[\bar{\mathcal{J}}(\bar{\mathcal{K}}_0) - \bar{\mathcal{J}}(\bar{\mathcal{K}}^\star)] o$  initial optimality gap

 $\mathcal{O}\left(rac{ au_{\max} \epsilon_{\mathsf{het}}}{N}
ight) o$  staleness effect becomes negligible when  $N \gg b_s$ 

Convergence guarantees - optimality gap

 $\textbf{Goal: Control} ~ \mathbb{E} \left[ \mathcal{J}^{(i)}(\bar{K}_N) - \mathcal{J}^{(i)}(K_i^{\star}) \right] \rightarrow \textbf{system-specific optimality gap}$ 

- Initial stabilizing controller:  $\bar{K}_0 \in \mathcal{S}$
- ullet  $\eta = \mathcal{O}\left(1/ au_{\mathsf{max}}^{3/2}
  ight)$  and r sufficiently small

Convergence guarantees - optimality gap

 $\textbf{Goal: Control} ~ \mathbb{E} \left[ \mathcal{J}^{(i)}(\bar{K}_N) - \mathcal{J}^{(i)}(K_i^{\star}) \right] \rightarrow \textbf{system-specific optimality gap}$ 

• Initial stabilizing controller:  $\bar{K}_0 \in \mathcal{S}$ 

•  $\eta = \mathcal{O}\left(1/ au_{\mathsf{max}}^{\mathsf{3/2}}
ight)$  and r sufficiently small

$$\mathbb{E}\left[\mathcal{J}^{(i)}(\bar{K}_{\mathsf{N}}) - \mathcal{J}^{(i)}(K_{i}^{\star})\right] \leq \mathcal{O}\left(c^{\mathsf{N}}\Delta_{0}^{(i)} + \epsilon_{\mathsf{het}}\right)$$
•  $\Delta_{0}^{(i)} = \mathbb{E}[\mathcal{J}^{(i)}(\bar{K}_{0}) - \mathcal{J}^{(i)}(\bar{K}^{\star})], \ c = 1 - \frac{\eta\lambda}{4} \in (0, 1)$ 

Convergence guarantees - optimality gap

**Goal:** Control  $\mathbb{E}\left[\mathcal{J}^{(i)}(\bar{K}_N) - \mathcal{J}^{(i)}(K_i^{\star})\right] \rightarrow$  system-specific optimality gap

• Initial stabilizing controller:  $\bar{K}_0 \in \mathcal{S}$ 

ullet  $\eta = \mathcal{O}\left(1/ au_{\mathsf{max}}^{3/2}
ight)$  and r sufficiently small

$$\mathbb{E}\left[\mathcal{J}^{(i)}(ar{\kappa}_{N}) - \mathcal{J}^{(i)}(K_{i}^{\star})
ight] \leq \mathcal{O}\left(c^{N}\Delta_{0}^{(i)} + \epsilon_{\mathsf{het}}
ight)$$
  
•  $\Delta_{0}^{(i)} = \mathbb{E}[\mathcal{J}^{(i)}(ar{\kappa}_{0}) - \mathcal{J}^{(i)}(ar{\kappa}^{\star})], \ c = 1 - rac{\eta\lambda}{4} \in (0, 1)$ 

Within the number of iterations in the order of  $N = \mathcal{O}\left( au_{\max}^{3/2}\log(1/\epsilon)
ight)$  we have

$$\mathbb{E}\left[\mathcal{J}^{(i)}(\bar{K}_{N}) - \mathcal{J}^{(i)}(K_{i}^{\star})\right] \leq \mathcal{O}\left(\epsilon + \epsilon_{\mathsf{het}}\right)$$

Convergence guarantees - optimality gap

Within the number of iterations in the order of  $N = \mathcal{O}\left( au_{\max}^{3/2}\log(1/\epsilon)
ight)$  we have



Stability guarantees

**Goal:** Show that  $\bar{K}_n \in S$  for any iteration  $n \in \{0, 1, \dots, N-1\}$ 

- Initial stabilizing controller:  $\bar{K}_0 \in \mathcal{S}$
- $\epsilon_{\sf het}$  and r sufficiently small and  $\eta = \mathcal{O}(1/ au_{\sf max}^{3/2})$



## Numerical Validation

Asynchronous vs synchronous updates

M = 100 heterogeneous systems with  $n_x = 4$  states and  $n_u = 2$  inputs



### Numerical Validation

Convergence - batch size and heterogeneity



(left)  $\tau_{\max} = 1$ , (right)  $\tau_{\max} = 5$ ,  $b_s = 20$ 

## Conclusions

Æ

- We studied the problem of learning linear quadratic regulators from **heterogeneous** systems with **asynchronous** policy gradient updates
- The proposed asynchronous aggregation **fully exploits** the parallelism in the distributed computation while alleviating the impact of **straggler systems**
- We provided local and global convergence guarantees, i.e.,

$$\frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E} \|\nabla \bar{\mathcal{J}}(\bar{K}_n)\|_F^2 \leq \mathcal{O}\left(\frac{\bar{\Delta}_0}{\sqrt{Nb_s}} + \frac{\epsilon_{\mathsf{het}}}{\sqrt{Nb_s}} + \frac{\tau_{\mathsf{max}}^2 \epsilon_{\mathsf{het}}}{N}\right)$$
$$\left[\mathcal{J}^{(i)}(\bar{K}_N) - \mathcal{J}^{(i)}(K_i^{\star})\right] \leq \mathcal{O}\left(\epsilon + \epsilon_{\mathsf{het}}\right), \text{ with } N = \mathcal{O}\left(\frac{\tau_{\mathsf{max}}^{3/2} \log(1/\epsilon)}{2}\right)$$

• We also showed that  $ar{K}_n \in \mathcal{S}$  for any iteration  $n \in \{0, 1, \dots, N-1\}$ 

- **1** JA. Tropp. **"User-Friendly Tail Bounds for Sums of Random Matrices"**, FoCM 2012.
- M. Fazel, R. Ge, S. Kakade, M. Mesbahi. "Global Convergence of Policy Gradient Methods for the Linear Quadratic Regulator", ICML 2018.
- H. Wang, LF. Toso, A. Mitra, J. Anderson. "Model-free Learning with Heterogeneous Dynamical Systems: A Federated LQR Approach", 2023.
- LF. Toso, D. Zhan, J. Anderson, H. Wang. "Meta-Learning Linear Quadratic Regulators: A Policy Gradient MAML Approach for Model-free LQR", L4DC 2024.
- Sha, F. Zhao, K. You. "Asynchronous Parallel Policy Gradient Methods for the Linear Quadratic Regulator", 2024. → homogeneous setting

## Collaborators





Han Wang James Anderson

Funding **Capital**One Columbia Engineering The Fu Foundation School of Engineering and Applied Science

# To find out more:





My website

Full paper

Happy to take questions!

