

# LQR Learning Pipelines

KIOS Graduate Training School 2024

Florian Dörfler

**ETH** zürich

# Context & acknowledgements

- collaboration with Claudio dePersis & Pietro Tesi to develop an explicit version of regularized DeePC  
→ **data-driven & regularized LQR**
- extension to adaptive LQR with Feiran Zhao, Keyou You, Linbin Huang, & Alessandro Chiuso  
→ **data-enabled policy optimization**
- revisit *old open problems with new perspectives*

Pietro Tesi (Florence)

Alessandro Chiuso (Padova)

Claudio de Persis (Groningen)

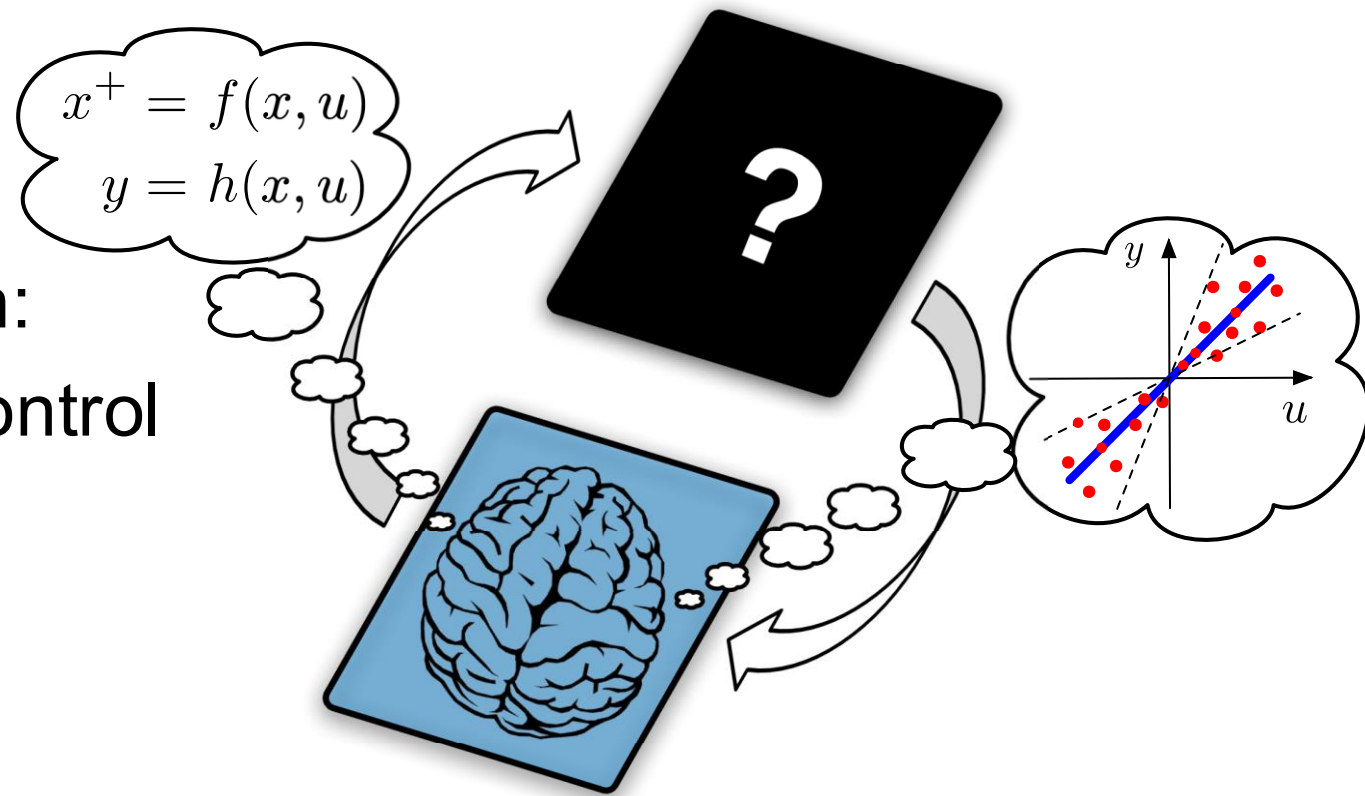
Feiran Zhao (Tsinghua)

Keyou You (Tsinghua)

Linbin Huang (Zhejiang)

# Data-driven pipelines

- **indirect** (model-based) approach:  
data  $\xrightarrow{ID}$  model + uncertainty  $\rightarrow$  control
- **direct** (model-free) approach:  
direct MRAC, RL, behavioral, ...
- **episodic & batch** algorithms:  
collect batch of data  $\rightarrow$  design policy  
 $\uparrow$  -----  
|
- **online & adaptive** algorithms:  
measure  $\rightarrow$  update policy  $\rightarrow$  actuate  
 $\uparrow$  -----  
|



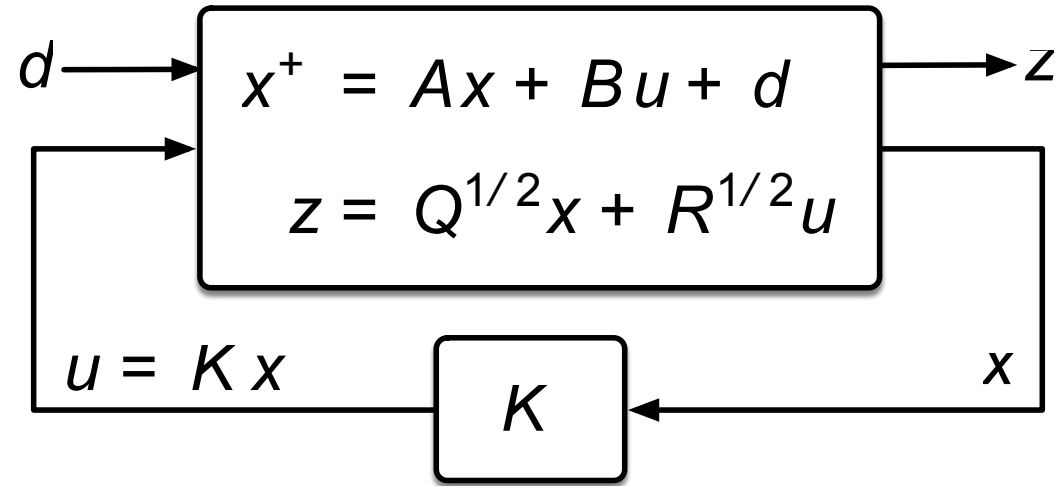
well-documented **trade-offs** concerning

- complexity: data, compute, & analysis
- goal: optimality vs (robust) stability
- practicality: modular vs end-to-end ...

$\rightarrow$  **gold(?) standard**: direct, adaptive, optimal yet robust, cheap, & tractable

# LQR

- **cornerstone** of automatic control



Equivalent LQR formulations:

- $J(K) = \sum_{t=0}^{\infty} x_t^T Q x_t + u_t^T R u_t = \sum_{t=0}^{\infty} x_t^T Q x_t + x_t^T K^T R K x_t$

- solution to  $x_{t+1} = (A + BK)x_t$  is  $x_t = (A + BK)^t x_0$

$$\leadsto J(K) = \sum_{t=0}^{\infty} x_0^T (A + BK)^{T^t} (Q + K^T R K) (A + BK)^t x_0$$

- Recall the closed-loop observability Gramian:  $W = \sum_{t=0}^{\infty} ((A + BK)^T)^t (Q + K^T R K) (A + BK)^t$



- $W$  can also be obtained as the unique positive definite solution to the Lyapunov equation:

$$(A+BK)^T W (A+BK) - W + Q + K^T R K = 0$$

⇒ equivalent reformulation

of  $J(u) = x_0^T W x_0 = \text{trace}(W x_0 x_0^T)$

- yet another reformulation

using  $\text{tr}(x_t^T Q x_t) = \text{tr}(Q x_t x_t^T)$

↳ covariance of (random)  $x_0$

$$J(K) = \text{tr}(Q \cdot P) + \text{tr}(K^T R K P)$$

where  $P = \sum_{t=0}^{\infty} x_t x_t^T$  (state covariance)

$$= \sum_{t=0}^{\infty} (A+BK)^t x_0 x_0^T ((A+BK)^T)^t$$

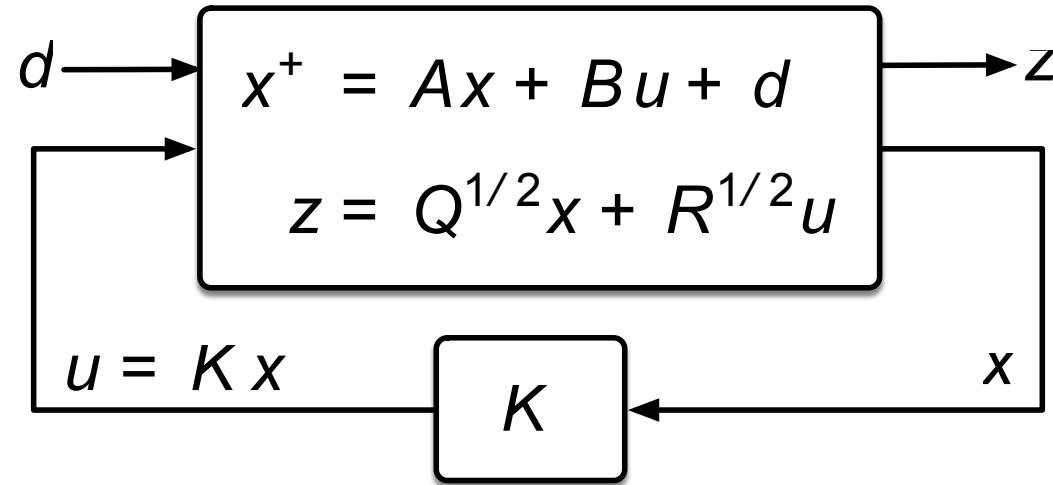
- recall that the above is the controllability Gramian which can be calculated uniquely as positive definite solution to

$$(A+BK) P (A+BK)^T - P + x_0 x_0^T = 0$$

side note: as it turns the actual value of  $x_0 x_0^T$  does not matter for the final result, and often one simply sets it to be identity

# LQR

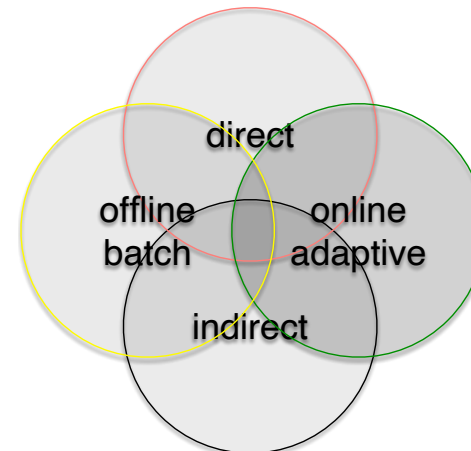
- **cornerstone** of automatic control



- $\mathcal{H}_2$  **parameterization**  
(can be posed as convex SDP,  
as differentiable program, as... )

$$\begin{aligned} & \text{minimize} && \text{trace}(QP) + \text{trace}(K^T R K P) \\ & P \succeq I, K \\ & \text{subject to} && (A + BK)P(A + BK)^T - P + I \preceq 0 \end{aligned}$$

- **the benchmark** for all data-driven control approaches in last decades but there is **no direct & adaptive LQR**



# Contents

- 1. model-based pipeline with model-free elements**  
→ data-driven parametrization & robustifying regularization
- 2. model-free pipeline with model-based elements**  
→ adaptive method: policy gradient & sample covariance
- 3. case studies: academic & power systems/electronics**  
→ LQR is academic example but can be made useful

# Contents

## 1. regularizations bridging direct & indirect data-driven LQR → story of a *model-based pipeline with model-free elements*

### On the Role of Regularization in Direct Data-Driven LQR Control




Florian Dörfler, Pietro Tesi, and Claudio De Persis

**Abstract**—The linear quadratic regulator (LQR) problem is a cornerstone of control theory and a widely studied benchmark problem. When a system model is not available, the conventional approach to LQR design is indirect, i.e., based on a model identified from data. Recently a suite of direct data-driven LQR design approaches has surfaced by-passing explicit system identification (SysID) and based on ideas from subspace methods and behavioral systems theory. In either approach, the data underlying the design can be taken at face value (certainty-equivalence) or the design is robustified to account for noise. An emerging topic in direct data-driven LQR design is to regularize the optimal control objective to account for implicit SysID (in a least-square or low-rank sense) or to promote robust stability. These regularized formulations are flexible, computationally attractive, and theoretically certifiable: they can interpolate

problems when identifying models from data. They facilitate finding solutions to optimization problems by rendering them unique or speeding up algorithms. Aside from such numerical advantages, a Bayesian interpretation of regularizations is that they condition models on prior knowledge [26], and they robustify problems to uncertainty [27], [28].

An emergent approach to data-driven control is borne out of the intersection of behavioral systems theory and subspace methods [29]. In particular, the so-called *Fundamental Lemma* characterizes the behavior of an LTI system by the range space of matrix time series data [30]. This perspective gave rise to direct data-driven predictive and

### On the Certainty-Equivalence Approach to Direct Data-Driven LQR Design

Florian Dörfler , Senior Member, IEEE, Pietro Tesi , Member, IEEE, and Claudio De Persis , Member, IEEE

**Abstract**—The linear quadratic regulator (LQR) problem is a cornerstone of automatic control, and it has been widely studied in the data-driven setting. The various data-driven approaches can be classified as indirect (i.e., based on an identified model) versus direct or as robust (i.e., taking uncertainty into account) versus certainty-equivalence. Here, we show how to bridge these different formulations and propose a novel, direct, and regularized formulation. We start from indirect certainty-equivalence LQR, i.e., least-square identification of state-space matrices followed by a nominal model-based design, formalized as a bilevel program. We show how to transform this problem into a single-level, regularized, and direct data-driven control formulation, where the regularizer accounts for the least-square data fitting criterion. For this novel formulation, we carry out a robustness and performance analysis in presence of noisy data. In a numerical case study, we compare regularizers promoting either robustness or certainty-equivalence, and we demonstrate the remarkable performance when blending both of them.

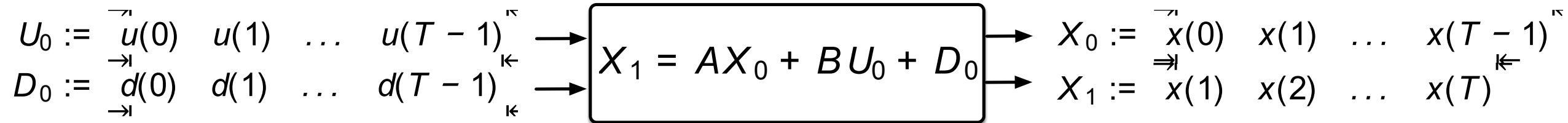
methods [10], [11], [12], reinforcement learning [13], behavioral methods [14], and Riccati-based methods [15] in the certainty-equivalence setting as well as [16], [17], [18] in the robust setting. We remark that the world is not black and white: a multitude of approaches have successfully bridged the direct and indirect paradigms, such as identification for control [19], [20], dual control [21], [22], control-oriented identification [23], and regularized data-enabled predictive control [24]. In essence, these approaches all advocate that the identification and control objectives should be blended to regularize each other.

An emergent approach to data-driven control is borne out of the intersection of behavioral systems theory and subspace methods; see the recent survey [25]. In particular, a result termed the *Fundamental Lemma* [26] implies that the behavior of an LTI system can be characterized by the range space of a matrix containing raw time series data. This perspective gave rise to implicit formulations (notably data-enabled predictive control [24], [27], [28]) as well as the design of explicit feedback policies [14], [15], [16], [17]. Both of these are direct

with Pietro Tesi (Florence) &  
Claudio de Persis (Groningen)

# Indirect & certainty-equivalence LQR

- collect I/O data  $(X_0, U_0, X_1)$  with  $D_0$  unknown & PE:  $\text{rank} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} = n + m$



- indirect & certainty-equivalence LQR**  
(optimal in MLE setting)

$$\text{minimize}_{P \succeq I, K} \text{trace}(QP) + \text{trace}(K^T R K P)$$

$$\text{subject to } (\hat{A} + \hat{B}K)P(\hat{A} + \hat{B}K)^T - P + I \preceq 0$$

$$[\hat{B} \quad \hat{A}] = \arg \min_{B, A} \left\| X_1 - \begin{bmatrix} B & A \end{bmatrix} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \right\|_F$$

certainty-equivalent LQR

least squares SysID

# Recall indirect approach on the board

- **I/O data**  $(X_0, U_0, X_1)$  with  $D_0$  unknown & PE:  $\text{rank} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} = n + m$

$$\begin{array}{l}
 U_0 := \begin{matrix} \overleftarrow{u(0)} & u(1) & \dots & u(T-1) \\ \overrightarrow{u(0)} & \overrightarrow{u(1)} & \dots & \overrightarrow{u(T-1)} \end{matrix} \\
 D_0 := \begin{matrix} \overleftarrow{d(0)} & d(1) & \dots & d(T-1) \\ \overrightarrow{d(0)} & \overrightarrow{d(1)} & \dots & \overrightarrow{d(T-1)} \end{matrix}
 \end{array}
 \rightarrow
 \boxed{X_1 = AX_0 + BU_0 + D_0}
 \rightarrow
 \begin{array}{l}
 X_0 := \begin{matrix} \overleftarrow{x(0)} & x(1) & \dots & x(T-1) \\ \overrightarrow{x(0)} & \overrightarrow{x(1)} & \dots & \overrightarrow{x(T-1)} \end{matrix} \\
 X_1 := \begin{matrix} x(1) & x(2) & \dots & x(T) \end{matrix}
 \end{array}$$

minimize  $\text{trace}(QP) + \text{trace}(K^T RKP)$   
 $P \succeq I, K$

subject to  $(\hat{A} + \hat{B}K)P(\hat{A} + \hat{B}K)^T - P + I \preceq 0$

$$[\hat{B} \quad \hat{A}] = \arg \min_{B, A} \left\| X_1 - [B \quad A] \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \right\|_F$$

$= X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^{\dagger}$  ← Moore-Penrose (right) inverse

$= X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^T \left( \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^T \right)^{-1}$  ← invertible due to PE

$\frac{\partial}{\partial C} = 0: 0 = 2(X_1 - [\hat{B} \quad \hat{A}]) \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}$

$\bullet \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^T$

$\rightsquigarrow [\hat{B} \quad \hat{A}] = X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^T \left( \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^T \right)^{-1}$

are uniquely determined

$\rightsquigarrow$  model-based design

# Derivation of a direct approach on the board

- **I/O data**  $(X_0, U_0, X_1)$  with  $D_0$  unknown & PE:  $\text{rank} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} = n + m$

$$\begin{array}{l}
 U_0 := \begin{matrix} \overrightarrow{\leftarrow} \\ \leftarrow \end{matrix} \begin{matrix} u(0) & u(1) & \dots & u(T-1) \\ \leftarrow & & & \leftarrow \end{matrix} \\
 D_0 := \begin{matrix} \overrightarrow{\leftarrow} \\ \leftarrow \end{matrix} \begin{matrix} d(0) & d(1) & \dots & d(T-1) \\ \leftarrow & & & \leftarrow \end{matrix}
 \end{array}
 \rightarrow \boxed{X_1 = AX_0 + BU_0 + D_0}
 \begin{array}{l}
 \rightarrow X_0 := \begin{matrix} \overrightarrow{\leftarrow} \\ \leftarrow \end{matrix} \begin{matrix} x(0) & x(1) & \dots & x(T-1) \\ \leftarrow & & & \leftarrow \end{matrix} \\
 \rightarrow X_1 := \begin{matrix} \overrightarrow{\leftarrow} \\ \leftarrow \end{matrix} \begin{matrix} x(1) & x(2) & \dots & x(T) \\ \leftarrow & & & \leftarrow \end{matrix}
 \end{array}$$

• PE implies that  $\forall K \exists G$  so that  $\begin{bmatrix} K \\ I \end{bmatrix} = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G$

• subspace relations for closed-loop matrix

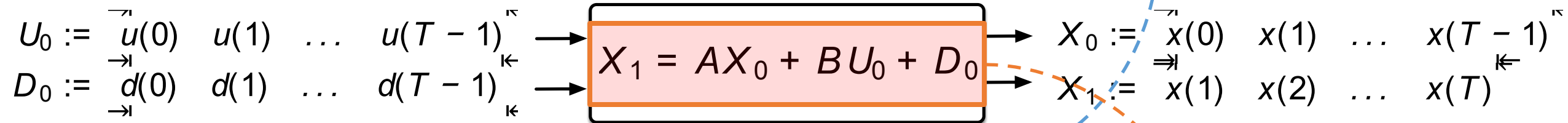
$$A + BK = [B \ A] \begin{bmatrix} K \\ I \end{bmatrix} = [B \ A] \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G = [AX_0 + BU_0] G = (X_1 - D_0) G$$

$\leadsto$  can replace  $A + BK$  in any LMI by  $(X_1 - D_0)G$

$\leadsto$  data driven parameterization of linear control design

# Direct approach from subspace relations in data

- **PE data:**  $\text{rank} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} = n + m \Rightarrow \forall K \exists G \text{ s.t. } \begin{bmatrix} K \\ I \end{bmatrix} = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G$



- **subspace relations**

$$A + BK = [B \quad A] \begin{bmatrix} K \\ I \end{bmatrix} \stackrel{\text{blue}}{=} [B \quad A] \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G \stackrel{\text{orange}}{=} (X_1 - D_0)G$$

- **data-driven LQR** LMIs by substituting  $A + BK = (X_1 - D_0)G$

→ certainty equivalence by neglecting noise  $D_0$ :  $A + BK = X_1G$



# Indirect

vs

# direct

minimize  $\text{trace}(QP) + \text{trace}(K^T R K P)$   
 $P \succeq I, K$

subject to  $(\hat{A} + \hat{B}K)P(\hat{A} + \hat{B}K)^T - P + I \preceq 0$

$$[\hat{B} \quad \hat{A}] = \arg \min_{B, A} \left\| X_1 - [B \quad A] \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \right\|_F$$

minimize  $\text{trace}(QP) + \text{trace}(K^T R K P)$   
 $P \succeq I, K, G$

subject to  $X_1 G P G^T X_1^T - P + I \preceq 0$

$$\begin{bmatrix} K \\ I \end{bmatrix} = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G$$

$$= X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^+$$

$\leadsto$  minimize  $\text{trace}(QP) + \text{trace}(K^T R K P)$

$P, K$

s.t.

$$X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^+ \begin{bmatrix} K \\ I \end{bmatrix} - P - (\dots)^T - P + I \preceq 0$$

$$[\hat{B} \quad \hat{A}]$$

$$\hat{A} + \hat{B}K$$

• issue  $\begin{bmatrix} U_0 \\ X_0 \end{bmatrix}$  has a big null space  
 $\leadsto$  solution  $G$  is not unique

• pick least norm solution satisfying

$$(I - \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^+ \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}) G = 0$$

$\leadsto$  projects orthogonal to null space

# Equivalence: direct + xxx $\Leftrightarrow$ indirect

- **direct** approach

→ optimizer has

nullspace  $\ker \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}$

→ orthogonality constraint

$$\text{minimize}_{P \succeq I, K, G} \quad \text{trace}(QP) + \text{trace}(K^\top RKP)$$

$$\text{subject to} \quad X_1 G P G^\top X_1^\top - P + I \preceq 0$$

$$\begin{bmatrix} K \\ I \end{bmatrix} = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G$$

$$\left( I - \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\dagger \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \right) G = 0$$

$$G = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\dagger \begin{bmatrix} K \\ I \end{bmatrix}$$

**equivalent constraints:**

$$\left( X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\dagger \begin{bmatrix} K \\ I \end{bmatrix} \right) P \left( \dots \right)^\top - P + I \preceq 0$$

- **indirect** approach

$$\text{minimize}_{P \succeq I, K} \quad \text{trace}(QP) + \text{trace}(K^\top RKP)$$

$$\text{subject to} \quad (\hat{A} + \hat{B}K)P(\hat{A} + \hat{B}K)^\top - P + I \preceq 0$$

$$\begin{bmatrix} \hat{B} & \hat{A} \end{bmatrix} = \arg \min_{B, A} \left\| X_1 - \begin{bmatrix} B & A \end{bmatrix} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \right\|_F$$

$$\begin{bmatrix} \hat{B} & \hat{A} \end{bmatrix} = X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\dagger$$

# Convex reformulation of the control design problem

minimize  
 $P \succeq I, K, G$

subject to

$$\text{trace}(QP) + \text{trace}(K^T R K P)$$

$$X_1 G P G^T X_1^T - P + I \leq 0$$

$$\begin{bmatrix} K \\ I \end{bmatrix} = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G$$

$$\Pi G = 0$$

$\hookrightarrow \text{trace}(R^{1/2} K P K^T R^{1/2})$

② can be pushed  
 to constraint  
 via epigraph  
 formulation

①  $K = U_0 G$   
 can be  
 eliminated

④ remove  $P = X_0 Y$

⑤ interpret  $\geq \mathcal{Q} \leq$   
 as Schur complement

③ substitute  $Y = G P$  or  $G = Y \cdot P^{-1} \Rightarrow K = U_0 G = U_0 Y P^{-1}$

minimize  $Y, X, P$

$$\text{trace}(QP) + \text{trace}(X)$$

$$X_1 Y P^{-1} Y^T X_1 - P + I \leq 0$$

$$X - R^{1/2} U_0 Y P^{-1} P P^{-1} U_0 Y^T U_0^T R^{1/2} \geq 0$$

$$I = X_0 G = X_0 Y \overset{P^{-1}}{\underbrace{P^{-1}}} \Leftrightarrow P = X_0 Y$$

$$\Pi G = 0$$

e.g.:

$$\begin{bmatrix} a & b \\ b^T & d \end{bmatrix} \geq 0$$

$$\Leftrightarrow d - b^T a^{-1} b \geq 0$$

if  $a > 0$

# Regularized, direct, & certainty-equivalent LQR

- orthogonality constraint

$$\Pi = I - \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\dagger \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}$$

**lifted** to regularizer

$$\begin{array}{ll} \text{minimize} & \text{trace}(QP) + \text{trace}(K^\top RKP) + \lambda \cdot \|\Pi G\| \\ P \succeq I, K, G & \\ \text{subject to} & X_1 G P G^\top X_1^\top - P + I \preceq 0 \\ & \begin{bmatrix} K \\ I \end{bmatrix} = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G \end{array}$$

- **equivalent** to indirect certainty-equivalent LQR design for  $\lambda$  suff. large
- $\lambda$  **interpolates** between direct & indirect approaches
- **multi-criteria interpretation**:  $\lambda$  interpolates control & SysID objectives
- however, certainty-equivalence formulation may not be **robust (?)**

# Robustness-promoting regularization

- effect of noise** entering data:  $A + BK = (X_1 - D_0)G$
  - Lyapunov constraint  $X_1 G P G^\top X_1^\top - P + I \preceq 0$
  - becomes  $(X_1 - D_0) G P G^\top (X_1 - D_0)^\top - P + I \preceq 0$
- } for robustness  $G P G^\top$  should be small
- previous **certainty-equivalence regularizer**  $\|\Pi G\|$  achieves small  $\|G\|$

- robustness-promoting regularizer** [de Persis & Tesi, '21]

$$\begin{aligned}
 & \underset{P \succeq I, K, G}{\text{minimize}} \quad \text{trace}(QP) + \text{trace}(K^\top R K P) \\
 & \quad + \rho \cdot \text{trace}(G P G^\top) \\
 & \text{subject to} \quad X_1 G P G^\top X_1^\top - P + I \preceq 0 \\
 & \quad \begin{bmatrix} K \\ I \end{bmatrix} = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G
 \end{aligned}$$

# Performance & robustness analysis

- **SNR** (signal-to-noise-ratio)  $\frac{\sigma_{min}([X_0 \ U_0])}{\sigma_{max}(D_0)}$

- **relative performance** metric

*realized cost from regularized design with  $\lambda$  &  $\rho$*

*if exact system matrices  $A$  and  $B$  were known*

$$\frac{\{\text{regularized data-driven LQR performance}\} - \{\text{ground-truth performance}\}}{\{\text{ground-truth performance}\}}$$

**certificate:** optimal control problem is **always feasible & stabilizing** for suff. large SNR & **relative performance**  $\sim \mathcal{O}(\text{SNR}^{-1}) + \text{const.}$   $\rho$  robust reg.

*proof* bounds Lyapunov constraint  $(X_1 - D_0)GPG^\top (X_1 - D_0)^\top - P + I \preceq 0$  <sub>19</sub>

# FYI: another regularization promoting low-rank

- de-noising of data-matrices via **low-rank approximation**

$$\begin{aligned} & \underset{\hat{U}_0, \hat{X}_0, \hat{X}_1}{\text{minimize}} \left\| \begin{bmatrix} \hat{U}_0 \\ \hat{X}_0 \\ \hat{X}_1 \end{bmatrix} - \begin{bmatrix} U_0 \\ X_0 \\ X_1 \end{bmatrix} \right\| \\ & \text{subject to } \text{rank} \begin{bmatrix} \hat{U}_0 \\ \hat{X}_0 \\ \hat{X}_1 \end{bmatrix} = \text{rank} \begin{bmatrix} \hat{U}_0 \\ \hat{X}_0 \end{bmatrix} = n + m \end{aligned}$$

Let PE hold:  $\text{rank} \begin{bmatrix} u_0 \\ x_0 \end{bmatrix} = n + m$

The following are equivalent:

$$(i) \text{rank} \begin{bmatrix} u_0 \\ x_0 \\ x_1 \end{bmatrix} = \text{rank} \begin{bmatrix} u_0 \\ x_0 \end{bmatrix} = n + m$$

$$(ii) \exists \text{ unique } B \& A \text{ so that } x_1 = AX_0 + BU_0$$

Proof: (ii)  $\Rightarrow$  (i) follows since  $x_1 = AX_0 + BU_0$  implies that  $x_1$  is dependent

(i)  $\Rightarrow$  (ii):  $n$  rows of  $\begin{bmatrix} u_0 \\ x_0 \\ x_1 \end{bmatrix}$  are dependent

due to PE, the rows of  $\begin{bmatrix} u_0 \\ x_0 \end{bmatrix}$  are independent

$\Rightarrow \exists [B \ A]$  so that  $x_1 = [B \ A] \begin{bmatrix} u_0 \\ x_0 \end{bmatrix}$

$\Rightarrow$  uniqueness due to PE

# Surrogate for low-rank pre-processing

minimize  $\text{trace}(QP) + \text{trace}(K^T R K P)$   
 $P, K, G$

$$A_{ce} P A_{ce}^T - P + I \leq 0$$

$$\begin{bmatrix} K \\ I \\ 0 \end{bmatrix} = \begin{bmatrix} x_0 \\ x_0 \\ x_1 \end{bmatrix} G \quad \text{or} \quad x_1 G = A + B K = A_{ce}$$

①

new constraint

without loss of  
 generality since  
 rank of  $\begin{bmatrix} x_0 \\ x_0 \\ x_1 \end{bmatrix} = n+m$

number of non-zero  
 entries of every column  
 $G_i$  of  $G$  is less than  $n+m$

① relax new constraint as  $\|G_i\|_1 \leq \alpha_i$  for suitable  $\alpha_i$

② relax as  $\|G\|_1 \leq \max \alpha_i$

③ lift to cost function as a penalty  $\|G\|_1$



# $l_1$ regularization as low-rank surrogate

- de-noising of data-matrices via **low-rank approximation** (low rank is equivalent to uniqueness of  $(A, B)$  matrices)
- $l_1$  **regularizer** as surrogate of pre-processing by low-rank approximation: bias solution  $G$  towards sparsity  $\sim$  low-rank

$$\begin{aligned} & \underset{\hat{U}_0, \hat{X}_0, \hat{X}_1}{\text{minimize}} \quad \left\| \begin{bmatrix} \hat{U}_0 \\ \hat{X}_0 \\ \hat{X}_1 \end{bmatrix} - \begin{bmatrix} U_0 \\ X_0 \\ X_1 \end{bmatrix} \right\| \\ & \text{subject to} \quad \text{rank} \begin{bmatrix} \hat{U}_0 \\ \hat{X}_0 \\ \hat{X}_1 \end{bmatrix} = \text{rank} \begin{bmatrix} \hat{U}_0 \\ \hat{X}_0 \end{bmatrix} = n + m \end{aligned}$$

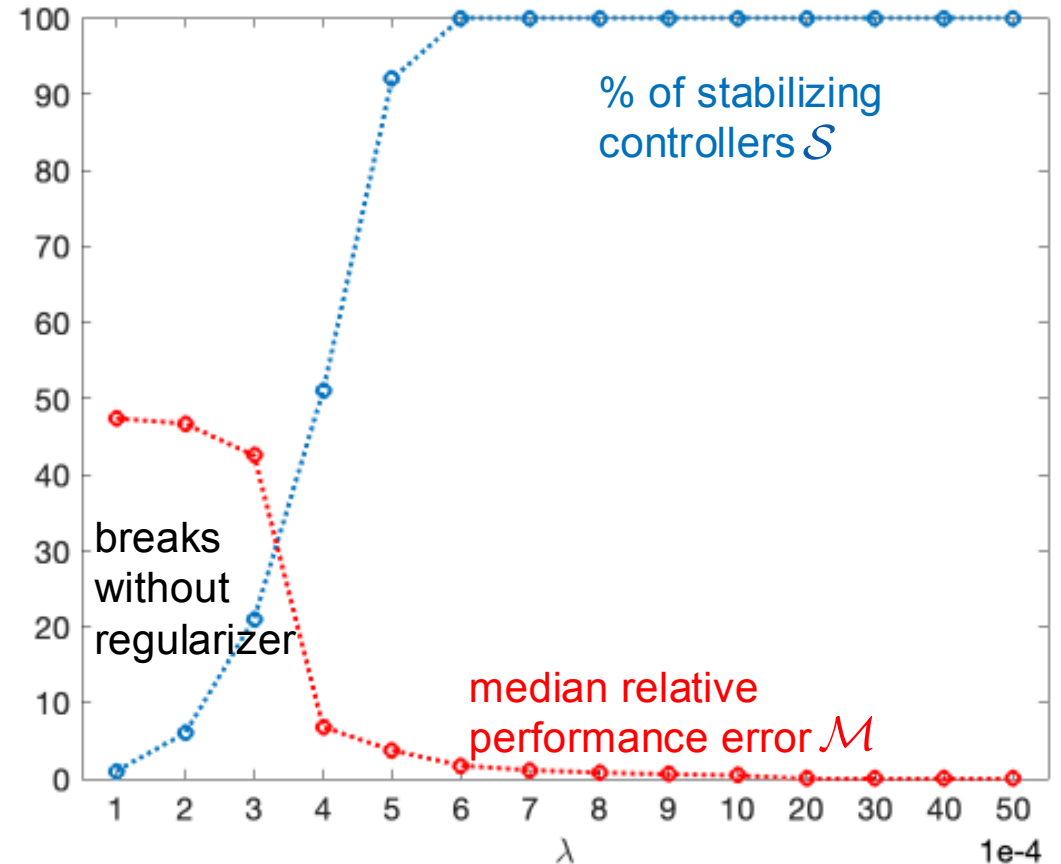
$$\begin{aligned} & \underset{K, P \succeq I, G}{\text{minimize}} \quad \text{trace}(QP) + \text{trace}(K^\top RKP) + \lambda \|G\|_1 \\ & \text{subject to} \quad X_1 G P G^\top X_1^\top - P + I \preceq 0 \\ & \quad \quad \quad \begin{bmatrix} K \\ I \end{bmatrix} = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G \end{aligned}$$

# Numerical case study

- **case study** [Dean et al. '19]: discrete-time marginally unstable Laplacian system subject to noise of variance  $\sigma^2 = 0.01$

- **take-home message 1:**  
*regularization is needed!*  
prior work without regularizer  
has no robustness margin

$$A = \begin{bmatrix} 1.01 & 0.01 & 0 \\ 0.01 & 1.01 & 0.01 \\ 0 & 0.01 & 1.01 \end{bmatrix}, \quad B = I$$



# Numerical case study cont'd

- **take-home message 2:** different regularizers promote different features: robustness vs. certainty-equivalence (performance)

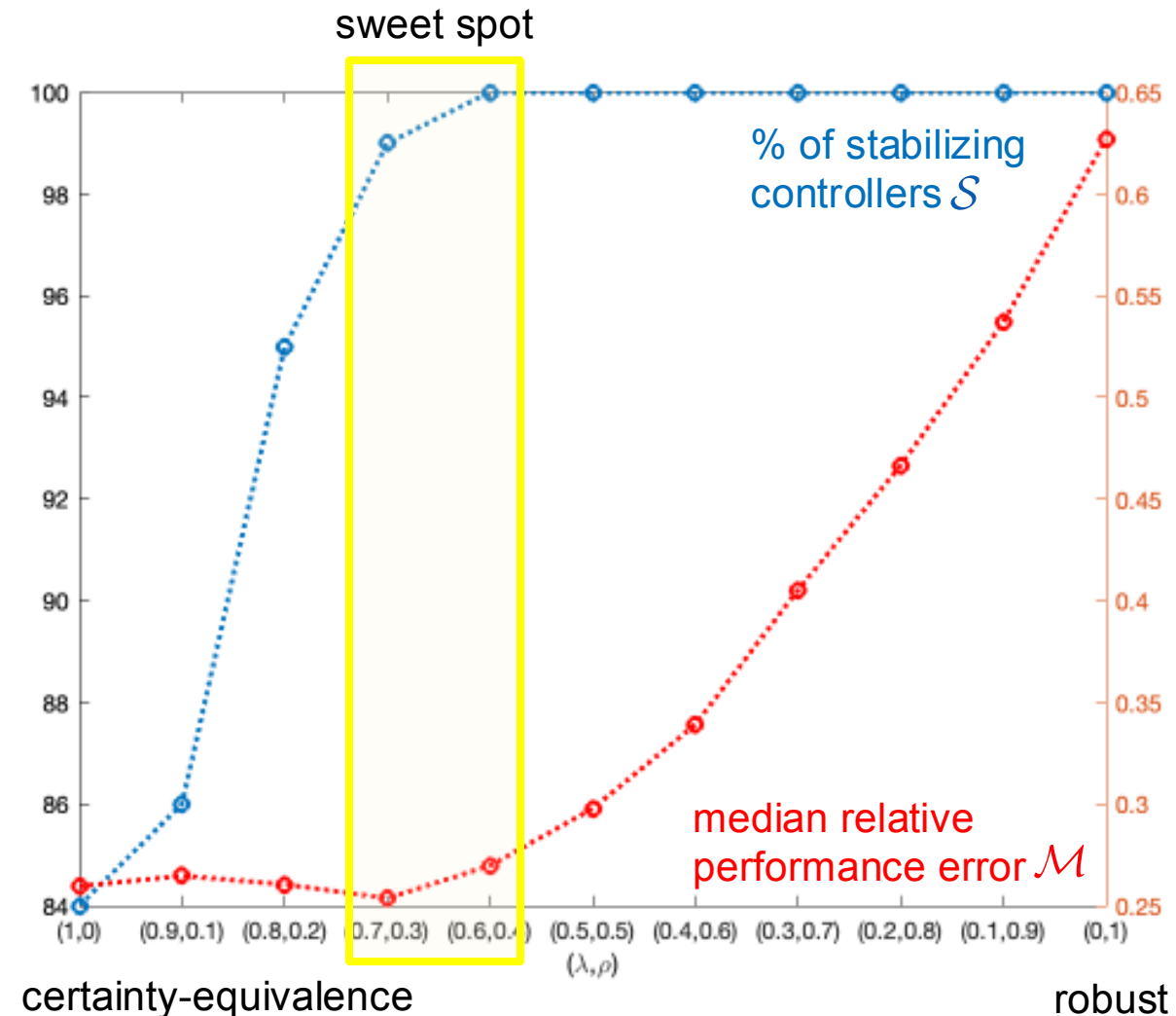
	$\sigma = 0.01$ (SNR > 15dB)	$\sigma = 0.1$ (SNR $\in [5, 10]$ dB)	$\sigma = 0.3$ (SNR $\in [0, 5]$ dB)	$\sigma = 0.7$ (SNR $\approx 0$ dB)	$\sigma = 1$ (SNR < -5dB)
Certainty-equivalence ( $\lambda = 1, \rho = 0$ )	$\mathcal{S} = 100\%$ $\mathcal{M} = 2.5599e-05$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.0026$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.0237$	$\mathcal{S} = 97\%$ $\mathcal{M} = 0.1366$	$\mathcal{S} = 84\%$ $\mathcal{M} = 0.2596$
Robust approach ( $\lambda = 0, \rho = 1$ )	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.0035$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.0074$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.0369$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.2350$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.6270$

- **take-home message 3:** mixed regularization achieves best of both

Mixed regularization ( $\lambda = \rho = 0.5$ )	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.0010$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.0035$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.0235$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.1262$	$\mathcal{S} = 100\%$ $\mathcal{M} = 0.2978$
--	---	---	---	---	---

# Intermediate conclusions... so far

- **interpolation** of different regularizers with high noise:  $\sigma^2 = 1$  (SNR < -5db)
  - **flexible multi-criteria formulation** trading off different objectives by regularizers (best of all is attainable)
  - **classification direct vs. indirect** is less relevant:  $\lambda$  interpolates
- works... but lame: **learning is offline**



# Contents

## 2. data-enabled policy optimization for online adaptation → story of a *model-free pipeline with model-based elements*

### Data-enabled Policy Optimization for the Linear Quadratic Regulator

Feiran Zhao, Florian Dörfler, Keyou You

**Abstract**—Policy optimization (PO), an essential approach of reinforcement learning for a broad range of system classes, requires significantly more system data than indirect (identification-followed-by-control) methods or behavioral-based direct methods even in the simplest linear quadratic regulator (LQR) problem. In this paper, we take an initial step towards bridging this gap by proposing the data-enabled policy optimization (DeePO) method, which requires only a finite number of sufficiently exciting data to iteratively solve the LQR problem via PO. Based on a data-driven closed-loop parameterization, we are able to directly compute the

a considerable gap in the sample complexity between PO and indirect methods, which have proved themselves to be more sample-efficient [9], [10] for solving the LQR problem. This gap is due to the exploration or trial-and-error nature of RL, or more specifically, that the cost used for gradient estimate can only be evaluated *after* a whole trajectory is observed. Thus, the existing PO methods require numerous system trajectories to find an optimal policy, even in the simplest LQR setting.

### Data-Enabled Policy Optimization for Direct Adaptive Learning of the LQR

Feiran Zhao, Florian Dörfler, Alessandro Chiuso, Keyou You

**Abstract**—Direct data-driven design methods for the linear quadratic regulator (LQR) mainly use offline or episodic data batches, and their online adaptation has been acknowledged as an open problem. In this paper, we propose a direct adaptive method to learn the LQR from online closed-loop data. First, we propose a new policy parameterization based on the sample covariance to formulate a direct data-driven LQR problem, which is shown to be equivalent to the certainty-equivalence LQR with optimal non-asymptotic guarantees. Second, we design a novel data-enabled policy optimization (DeePO) method to directly update the policy, where the gradient is explicitly computed using only a batch of persistently exciting (PE) data. Third, we establish its global convergence via a projected gradient dominance property. Importantly, we efficiently use DeePO to adaptively learn the LQR by performing only one-step projected gradient descent per sample of the closed-loop system, which also leads to an explicit recursive update of the policy. Under PE inputs and for bounded noise, we show that the average regret of the LQR cost is upper-bounded by two terms signifying a sublinear decrease in time  $\mathcal{O}(1/\sqrt{T})$  plus a bias scaling inversely with signal-to-

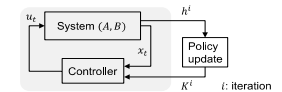


Fig. 1. An illustration of episodic approaches, where  $h^i = (x_0, u_0, \dots, x_{T-1})$  denotes the trajectory of the  $i$ -th episode.

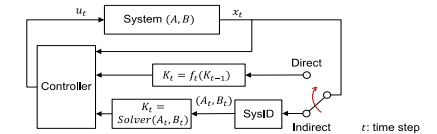


Fig. 2. An illustration of indirect and direct adaptive approaches in closed-loop, where  $f_t$  is some explicit function.

with Alessandro Chiuso (Padova),  
Feiran Zhao, Keyou You (Tsinghua),  
& Linbin Huang (Zhezjiang)

# Online & adaptive solutions

- **shortcoming** of separating offline learning & online control  
→ cannot improve policy **online** & cheaply / rapidly **adapt** to changes

Adaptive Control:  
Towards a Complexity-Based General Theory\*  
G. ZAMES-

*“adaptive = improve over best control with a priori info”*

- (elitist) **desired adaptive** solution: direct, online (non-episodic/non-batch) algorithms, with closed-loop data, & recursive algorithmic implementation
- “best” way to improve policy with new data → **go down the gradient !**

---

\* disclaimer: a large part of the adaptive control community focuses on stability & not optimality

# Ingredient 1: policy gradient methods

- LQR viewed as smooth program (many formulations)

$$\begin{aligned} & \text{minimize}_{P \succeq I, K} \quad \text{trace}(QP) + \text{trace}(K^\top RKP) \\ & \text{subject to} \quad (A + BK)P(A + BK)^\top - P + I \preceq 0 \end{aligned}$$

after eliminating  
(unique)  $P$ ,  
denote this  
as  $J(K)$

- $J(K)$  is not convex ...

but on the set of stabilizing gains  $K$ , it's

- coercive with compact sublevel sets,
- smooth with bounded Hessian, &
- degree-2 gradient dominated

$$J(K) - J^* \leq \text{const.} \cdot \|\nabla J(K)\|^2$$

**Fact:** policy gradient descent

$$K^+ = K - \eta \nabla J(K)$$

initialized from a stabilizing policy converges linearly to  $K^*$ .

# Insights into the proof

•  $J(k)$  is smooth with  $\|\nabla^2 J(k)\| \leq L$  : By Taylor or mean-value theorem:

$$J(k') \leq J(k) + \nabla J(k)^T (k' - k) + \frac{L}{2} \|k' - k\|_F^2 \quad (1)$$

• gradient dominance :  $J(k) - J(k^*) \leq \frac{1}{2\mu} \|\nabla J(k)\|_F^2 \quad (2)$

• gradient descent :  $k^T = k - \eta \nabla J(k)$

$$\leadsto J(k^T) = J(k - \eta \nabla J(k)) \stackrel{(1)}{\leq} J(k) + \nabla J(k)^T (k - \eta \nabla J(k) - k) + \frac{L}{2} \eta^2 \|\nabla J(k)\|^2$$

$$= J(k) - \left(\eta - \frac{L\eta^2}{2}\right) \|\nabla J(k)\|^2 - \eta \nabla J(k)^T$$

$$\stackrel{(2)}{=} J(k) - \left(\eta - \frac{L\eta^2}{2}\right) 2\mu (J(k) - J(k^*)) \Rightarrow J(k^T) - J(k^*) \leq \left(1 - \eta - \frac{L\eta^2}{2}\right) 2\mu (J(k) - J(k^*))$$



# Explicit formulae for model-based gradient

• For these results we need the equivalent LQR formulations (see beginning)

$$J(k) = \text{tr}(PQ) + \text{tr}(K^T R K P) \quad \text{where } P > 0 \text{ solves } (A+Bk)P(A+Bk)^T - P + Q = 0$$
$$= \text{tr}(W X) \quad \text{where } W > 0 \text{ solves } (A+Bk)^T W (A+Bk) - W + Q + k^T R k = 0$$

where  $X = x_0 x_0^T$  is the initial state covariance, though its particular value is irrelevant

• To calculate the gradient, we recognize  $\nabla J(k) = \frac{\partial}{\partial k} \text{tr}(W(k) \cdot X)$

... as you can see, the math for such derivatives can get cumbersome. For these reasons, we will work with differentials which will simplify the derivations.

$$= \begin{bmatrix} \vdots & \text{tr}\left(\frac{\partial W}{\partial k_{ij}} \cdot X\right) & \vdots \end{bmatrix}$$

The differential  $dx$  is the linear part (Jacobian) of the function  $f(x+dx) - f(x)$

• mathematical preliminaries on differentials

•  $d \text{Tr}(A) = \text{Tr}(dA)$

•  $d(A \cdot B) = dA \cdot B + dB \cdot A$

•  $d(A^T) = dA^T$

• Let  $J$  be a function of  $x$ . If  $dJ = \text{Tr}(C \cdot dx)$ , then  $\nabla_x J = C^T$

this one is constant  
↓  
with zero differential

• Derivation of  $\nabla_k J^{(k)}$ : (1) since  $dJ = d \text{tr}(W \cdot X) = \text{tr}(dW \cdot X)$

(2) to obtain  $dW$ , we evaluate ●:

$$\leadsto (A+Bk)^T \delta W (A+Bk) - \delta W + \overbrace{\delta k^T (A+Bk)^T W B + k^T R}^{\Pi} + (\dots \text{same term} \dots)^T$$

$\Pi^T$

$\leadsto$  this is a Lyapunov equation and thus

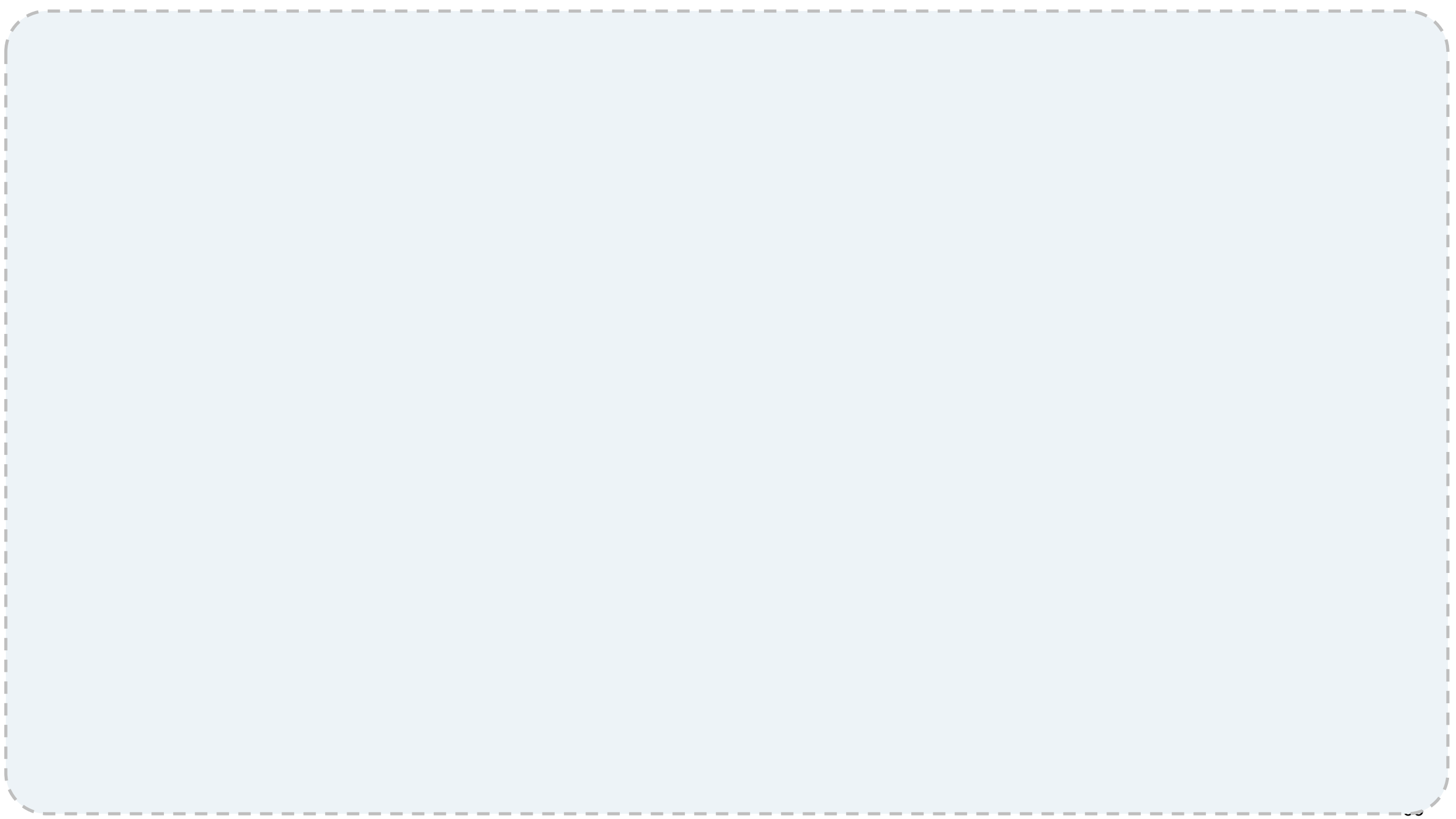
$$\frac{\partial W}{\partial k} = \sum_{t=0}^{\infty} ((A+Bk)^T)^t (\Pi + \Pi^T) (A+Bk)^t$$

$$\rightarrow \text{Hence, } \text{tr}(SW \cdot X) = \text{tr} \left( 2 \pi^T \underbrace{\sum_{t=0}^{\infty} (A+Bk)^t X ((A+Bk)^t)^T}_{= P \text{ (controllability Gramian)}} \right)$$

$$= \text{tr} \left( \delta k^T 2 (B^T W (A+Bk) + Rk) \cdot P \right)$$

$\rightarrow$  Last, using that  $dJ = \text{Tr}(C \cdot dk) \Rightarrow \nabla_k J = C^T$ , we obtain

$$\nabla_k J(k) = 2 (B^T W (A+Bk) + Rk) \cdot P$$



# Model-free policy gradient methods

- **model-based setting**: explicit formulae for  $\nabla J(K)$  based on closed-loop controllability + observability Gramians [Levine & Athans, '70]
- **model-free 0<sup>th</sup> order methods** constructing two-point gradient estimate

conceptual for a scalar function:  $\nabla f(x) = \lim_{\epsilon \rightarrow 0} \frac{1}{2\epsilon} (f(x+\epsilon) - f(x-\epsilon)) = \lim_{\eta \rightarrow 0} E_{\eta \sim \text{uniform in } [-1,1]} \frac{\eta}{\epsilon} f(x+\eta\epsilon)$   
 $\rightarrow$  can be approximated sampling function, but scales very poorly for high dimension

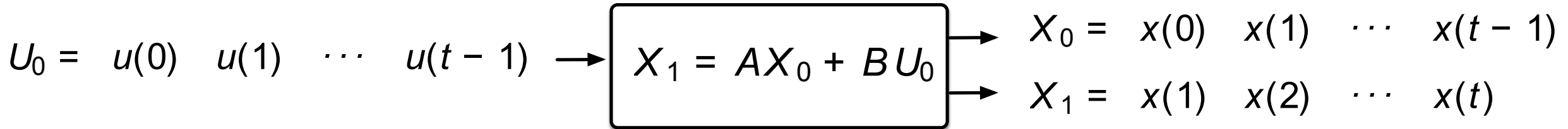
from numerous & very long trajectories  $\rightarrow$  extremely sample inefficient

relative performance gap	$\epsilon = 1$	$\epsilon = 0.1$	$\epsilon = 0.01$
# trajectories (100 samples)	1414	43850	142865

$\sim 10^7$  samples

- IMO: policy gradient is a **potentially great** candidate for direct adaptive control **but sadly useless in practice**: sample-inefficient, episodic, ...

# Ingredient 2: sample covariance parameterization



## prior parameterization

- PE condition: full row rank  $\begin{bmatrix} U_0 \\ X_0 \end{bmatrix}$
- $A + BK = [B \ A] \begin{bmatrix} K \\ I \end{bmatrix} = [B \ A] \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} G = X_1 G$
- robustness:  $G = \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\top (\cdot) \leftrightarrow$  regularization
- dimension of all matrices grows with  $t$

## covariance parameterization

- sample covariance  $\Lambda = \frac{1}{t} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\top > 0$
- $A + BK = [B \ A] \begin{bmatrix} K \\ I \end{bmatrix} = [B \ A] \Lambda V = \frac{1}{t} X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\top V$
- robustness for free without regularization
- dimension of all matrices is constant  
+ cheap rank-1 updates for online data

# Covariance parameterization of the LQR

- state / input sample covariance  $\Lambda = \frac{1}{t} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix} \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\top$  &  $\bar{X}_1 = \frac{1}{t} X_1 \begin{bmatrix} U_0 \\ X_0 \end{bmatrix}^\top$

- closed-loop matrix  $A + BK = \bar{X}_1 V$  with  $\begin{bmatrix} K \\ I \end{bmatrix} = \Lambda V = \begin{bmatrix} \bar{U}_0 \\ \bar{X}_0 \end{bmatrix} V$

- LQR covariance parameterization after eliminating  $K$  with variable  $V$ , Lyapunov eqn (explicitly solvable), smooth cost  $J(V)$  (after removing  $P$ ), & linear parameterization constraint

$$\begin{aligned} \min_{V, P > 0} & \text{trace}(QP) + \text{trace}(V^T \bar{U}_0^T R \bar{U}_0 V P) \\ \text{s.t. } & P = I + \bar{X}_1 V P V^T \bar{X}_1^T, V = \bar{X}_0 V \end{aligned}$$

details are not important

# Projected policy gradient with sample covariances

- **data-enabled policy optimization (DeePO)**

$$V^+ = V - \eta \Pi_{\bar{X}_0}(\nabla J(V))$$

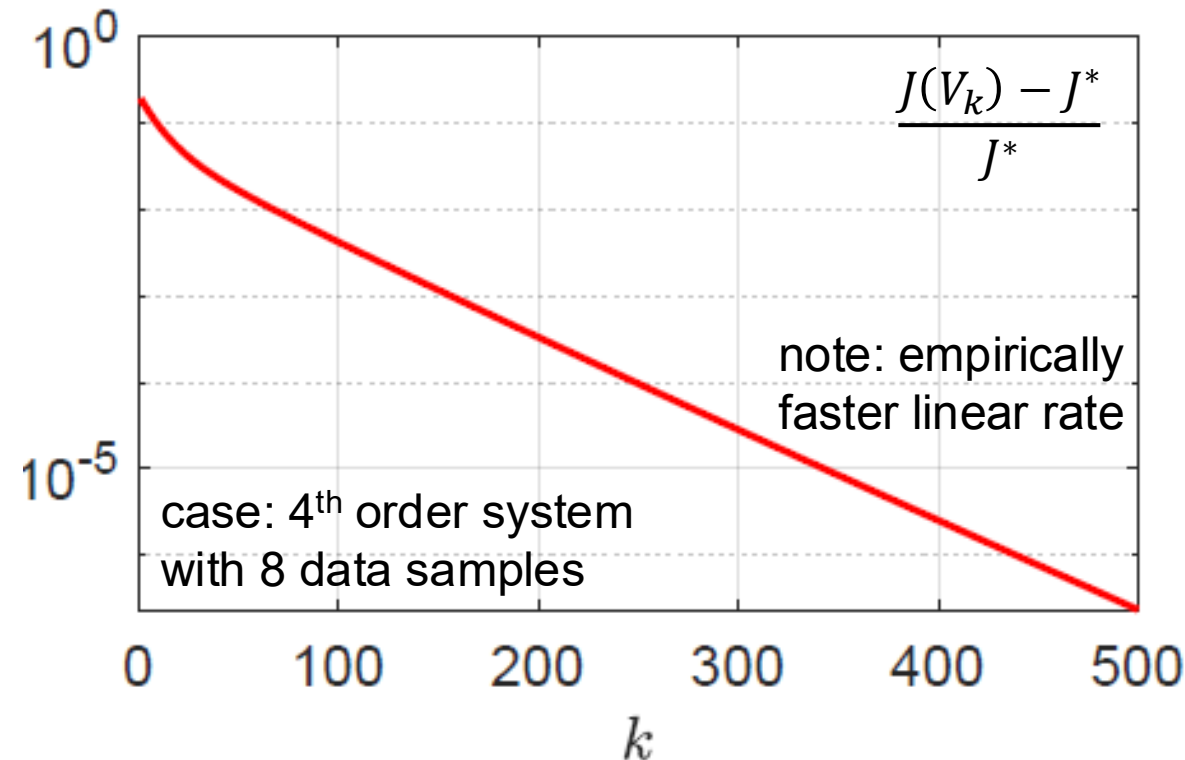
$\Pi_{\bar{X}_0}$  projects on parameterization constraint  $I = \bar{X}_0 V$  & gradient  $\nabla J(V)$  is computed from two Lyapunov equations with sample covariances

- **optimization landscape:** smooth, degree-1 proj. grad dominance

$$J(V) - J^* \leq \text{const.} \cdot \left\| \Pi_{\bar{X}_0}(\nabla J(V)) \right\|$$

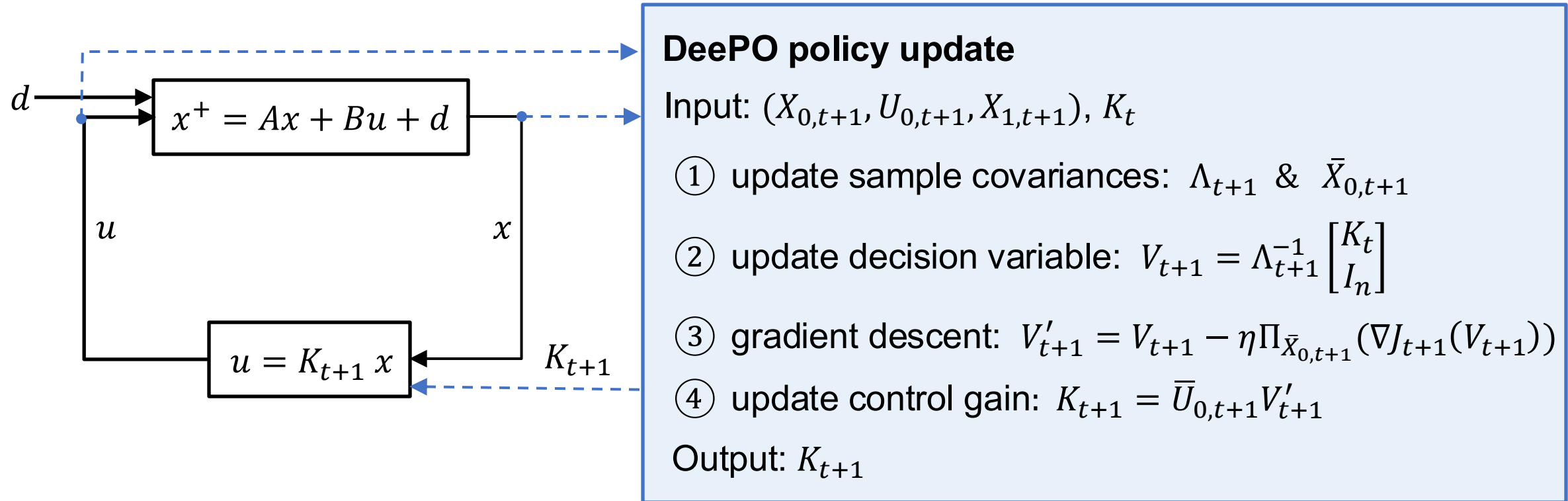
- warm-up: offline data & no disturbance

**Sublinear convergence** for feasible initialization  $J(V^k) - J^* \leq \mathcal{O}(1/k)$ .





# Online, adaptive, & closed-loop DeePO



where  $X_{0,t+1} = [x(0), x(1), \dots, x(t), x(t+1)]$  & similar for other matrices

- **cheap & recursive implementation:** rank-1 update of (inverse) sample covariances, cheap computation, & no memory needed to store old data

# Underlying assumptions for theoretic certificates

- **initially stabilizing controller:** the LQR problem parameterized by offline data  $(X_{0,t_0}, U_{0,t_0}, X_{1,t_0})$  is feasible with stabilizing gain  $K_{t_0}$ .
- **persistence of excitation** due to process noise or probing:  
$$\underline{\sigma}(\mathcal{H}_{n+1}(U_{0,t})) \geq \gamma \cdot \sqrt{t} \quad \text{with Hankel matrix } \mathcal{H}_{n+1}(U_{0,t})$$
- **bounded noise:**  $\|d(t)\| \leq \delta \quad \forall t \rightarrow$  **signal-to-noise** ratio  $SNR := \gamma/\delta$
- **BIBO:** there are  $\bar{u}, \bar{x}$  such that  $\|u(t)\| \leq \bar{u} \quad \& \quad \|x(t)\| \leq \bar{x}$   
( $\exists$  common Lyapunov function ?)

# Bounded regret of DeePO in adaptive setting

- **average regret** performance metric  $\text{Regret}_T := \frac{1}{T} \sum_{t=t_0}^{t_0+T-1} (J(K_t) - J^*)$

**Sublinear regret:** Under the assumptions, there are  $\nu_1, \nu_2, \nu_3, \nu_4 > 0$  such that for  $\eta \in (0, \nu_1]$  &  $SNR \geq \nu_2$ , it holds that  $\{K_t\}$  is stabilizing &

$$\text{Regret}_T \leq \frac{\nu_3}{\sqrt{T}} + \frac{\nu_4}{\sqrt{SNR}} .$$

- **comments** on the qualitatively expected result:
  - analysis is independent of the noise statistics & **consistent**  $\text{Regret}_{T \rightarrow \infty} \rightarrow 0$
  - **favorable sample complexity:** sublinear decrease term matches best rate  $\mathcal{O}(1/\sqrt{T})$  of first-order methods in online convex optimization
  - empirically observe smaller **bias term:**  $\mathcal{O}(1/SNR^2)$  & not  $\mathcal{O}(1/\sqrt{SNR})$

# Comparison case studies

- **same case study** [Dean et al. '19]

- **case 1: offline LQR**

vs direct adaptive DeePO

vs indirect adaptive: rls + dlqr

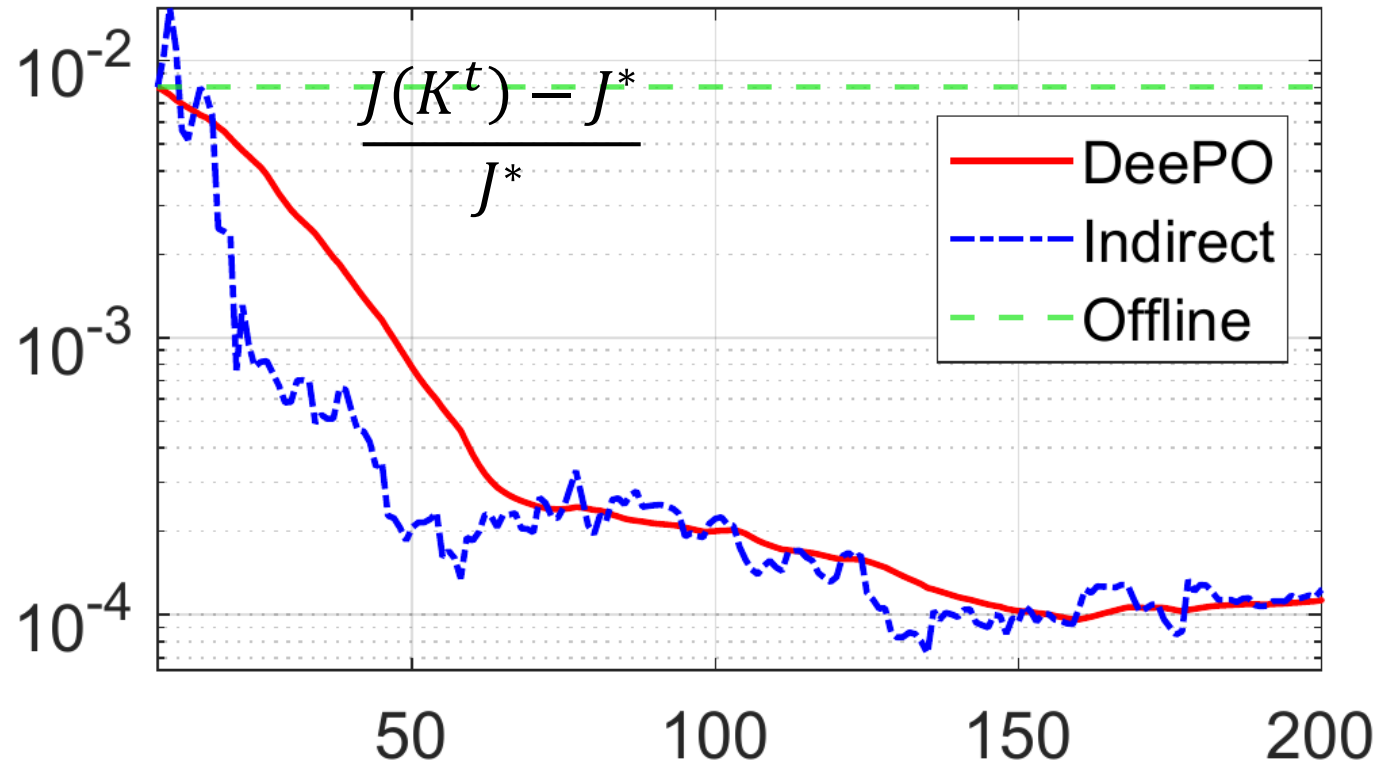
→ **adaptive outperforms offline**

→ direct/indirect **rates matching**  
but **direct is much(!) cheaper**

- **case 2: adaptive DeePO**

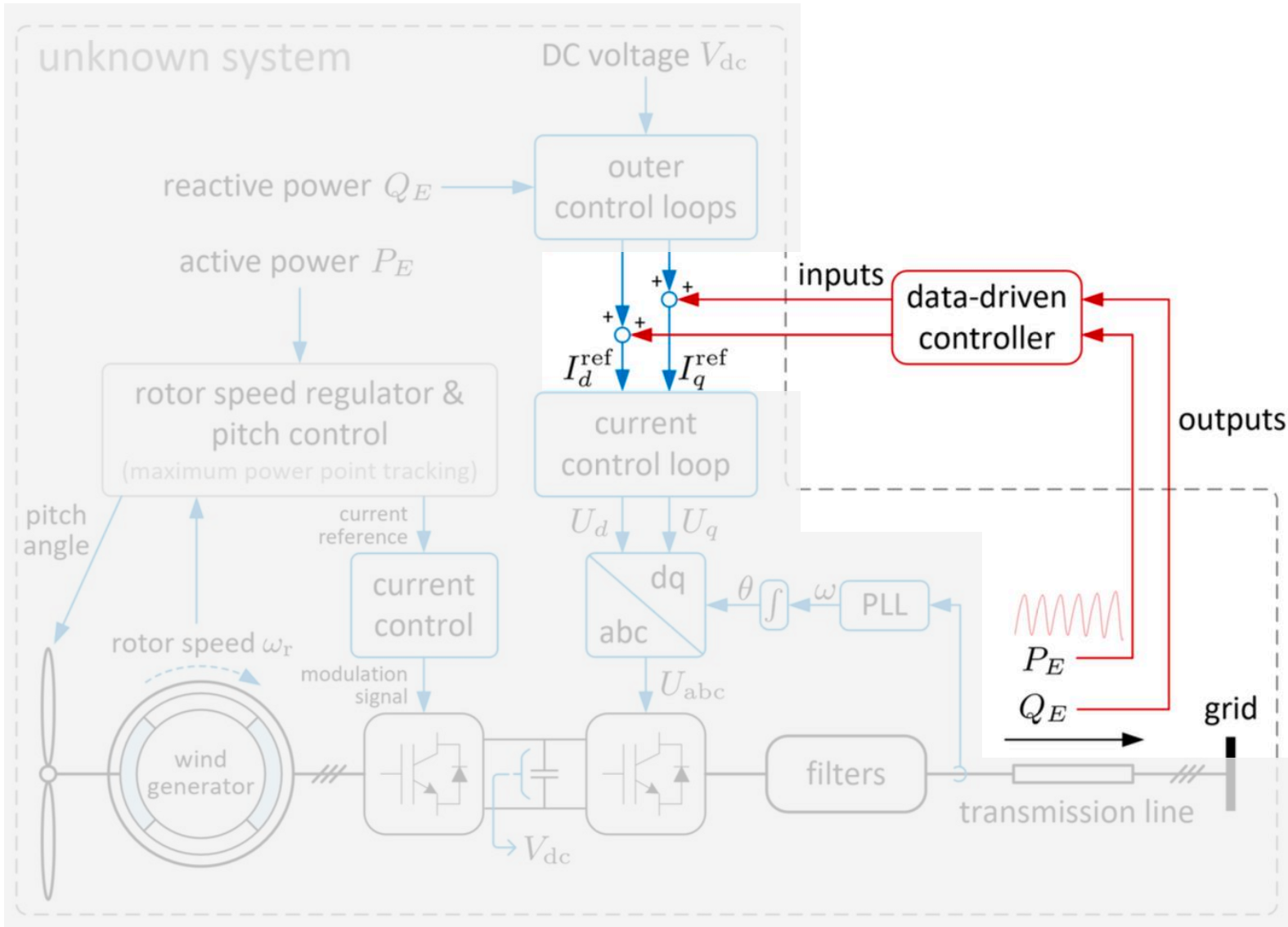
vs  $0^{th}$  order methods

→ **significantly less data**



relative performance gap	$\epsilon = 1$	$\epsilon = 0.1$	$\epsilon = 0.01$
# long trajectories ( <b>100</b> samples) for $0^{th}$ order LQR	1414	43850	<b>142865</b>
DeePO (# I/O samples)	10	24	<b>48</b>

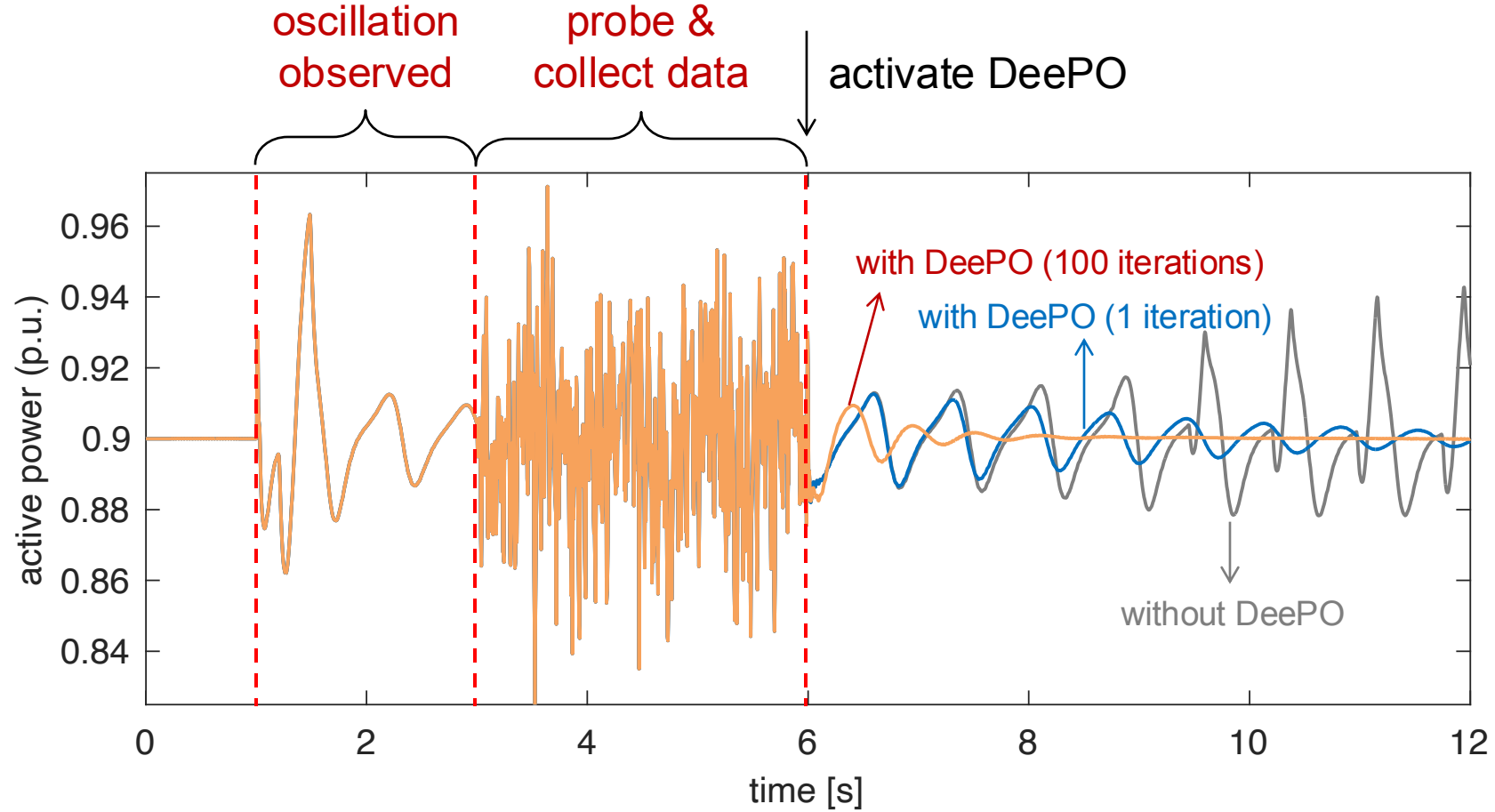
# Power systems / electronics case study



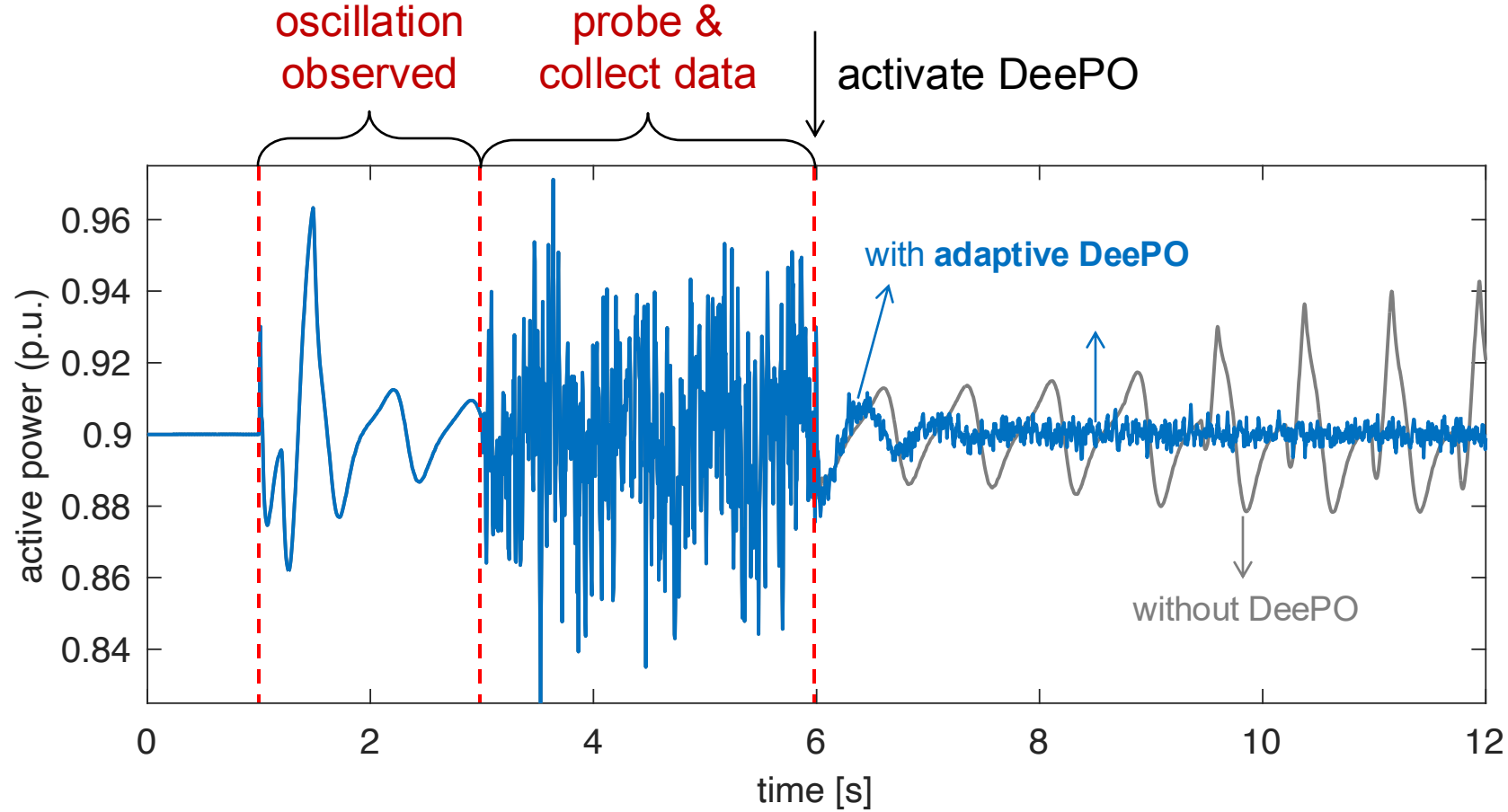
synchronous generator & full-scale converter

- wind turbine becomes **unstable** in weak grids with nonlinear oscillations
- converter, turbine, & grid are a **black box** for the commissioning engineer
- construct state from time shifts (5ms sampling) of  $(y(t), u(t))$  & use **DeePO**

# Power systems / electronics case study



# ... same in the adaptive setting with excitation



# Conclusions

- **Summary**

- model-based pipeline with model-free block: data-driven LQR parametrization  
→ works well when regularized (note: further flexible regularizations available)
- model-free pipeline with model-based block: policy gradient & sample covariance  
→ DeePO is adaptive, online, with closed-loop data, & recursive implementation
- academic case studies & can be made useful in power systems/electronics

- **Future work**

- technicalities: weaken assumptions & improve rates
- control: based on output feedback & for other objectives
- further system classes: stochastic, time-varying, & nonlinear
- open questions: online vs episodic? “best” batch size? triggered?



# Papers

## 1. model-based pipeline with model-free elements

### On the Role of Regularization in Direct Data-Driven LQR Control

Florian Dörfler, Pietro Tesi, and Claudio De Persis

**Abstract**—The linear quadratic regulator (LQR) problem is a cornerstone of control theory and a widely studied benchmark problem. When a system model is not available, the conventional approach to LQR design is indirect, i.e., based on a model identified from data. Recently a suite of direct data-driven LQR design approaches has surfaced by-passing explicit system identification (SysID) and based on ideas from subspace methods and behavioral systems theory. In either approach, the data underlying the design can be taken at face value (certainty-

problems when identifying models from data. They facilitate finding solutions to optimization problems by rendering them unique or speeding up algorithms. Aside from such numerical advantages, a Bayesian interpretation of regularizations is that they condition models on prior knowledge [26], and they robustify problems to uncertainty [27], [28].

An emergent approach to data-driven control is borne out of the intersection of behavioral systems theory and

## 2. model-free pipeline with model-based elements




### Data-enabled Policy Optimization for the Linear Quadratic Regulator

Feiran Zhao, Florian Dörfler, Keyou You

**Abstract**—Policy optimization (PO), an essential approach of reinforcement learning for a broad range of system classes, requires significantly more system data than indirect (identification-followed-by-control) methods or behavioral-based direct methods even in the simplest linear quadratic regulator (LQR) problem. In this paper, we take an initial step towards bridging this gap by proposing the data-enabled policy optimization (DeePO) method, which requires only a finite number of sufficiently exciting data to iteratively solve the LQR problem via PO. Based on a data-driven closed-loop parameterization, we are able to directly compute the

a considerable gap in the sample complexity between PO and indirect methods, which have proved themselves to be more sample-efficient [9], [10] for solving the LQR problem. This gap is due to the exploration or trial-and-error nature of RL, or more specifically, that the cost used for gradient estimate can only be evaluated *after* a whole trajectory is observed. Thus, the existing PO methods require numerous system trajectories to find an optimal policy, even in the simplest LQR setting.

### On the Certainty-Equivalence Approach to Direct Data-Driven LQR Design

Florian Dörfler , Senior Member, IEEE, Pietro Tesi , Member, IEEE, and Claudio De Persis , Member, IEEE

**Abstract**—The linear quadratic regulator (LQR) problem is a cornerstone of automatic control, and it has been widely studied in the data-driven setting. The various data-driven approaches can be classified as indirect (i.e., based on an identified model) versus direct or as robust (i.e., taking uncertainty into account) versus certainty-equivalence. Here, we show how to bridge these different formulations and propose a novel, direct, and regularized formulation. We start from indirect certainty-equivalence LQR, i.e., least-square identification of state-space matrices followed by a nominal model-based design, formalized as a bilevel program. We show how to transform this problem into a single-level, regularized, and direct data-driven control formulation, where the regularizer accounts for the least-square data fitting criterion. For this novel formulation, we carry out a robustness and performance analysis in presence of noisy data. In a numerical case study, we compare regularizers promoting either robustness or certainty-equivalence, and we demonstrate the remarkable performance when blending both of them.

methods [10], [11], [12], reinforcement learning [13], behavioral methods [14], and Riccati-based methods [15] in the certainty-equivalence setting as well as [16], [17], [18] in the robust setting. We remark that the world is not black and white: a multitude of approaches have successfully bridged the direct and indirect paradigms, such as identification for control [19], [20], dual control [21], [22], control-oriented identification [23], and regularized data-enabled predictive control [24]. In essence, these approaches all advocate that the identification and control objectives should be blended to regularize each other.

An emergent approach to data-driven control is borne out of the intersection of behavioral systems theory and subspace methods; see the recent survey [25]. In particular, a result termed the *Fundamental Lemma* [26] implies that the behavior of an LTI system can be characterized by the range space of a matrix containing raw time series data. This perspective gave rise to implicit formulations (notably data-enabled predictive control [24], [27], [28]) as well as the design of explicit feedback policies [14], [15], [16], [17]. Both of these are direct

### Data-Enabled Policy Optimization for Direct Adaptive Learning of the LQR

Feiran Zhao, Florian Dörfler, Alessandro Chiuso, Keyou You

**Abstract**—Direct data-driven design methods for the linear quadratic regulator (LQR) mainly use offline or episodic data batches, and their online adaptation has been acknowledged as an open problem. In this paper, we propose a direct adaptive method to learn the LQR from online closed-loop data. First, we propose a new policy parameterization based on the sample covariance to formulate a direct data-driven LQR problem, which is shown to be equivalent to the certainty-equivalence LQR with optimal non-asymptotic guarantees. Second, we design a novel data-enabled policy optimization (DeePO) method to directly update the policy, where the gradient is explicitly computed using only a batch of persistently exciting (PE) data. Third, we establish its global convergence via a projected gradient dominance property. Importantly, we efficiently use DeePO to adaptively learn the LQR by performing only one-step projected gradient descent per sample of the closed-loop system, which also leads to an explicit recursive update of the policy. Under PE inputs and for bounded noise, we show that the average regret of the LQR cost is upper-bounded by two terms signifying a sublinear decrease in time  $\mathcal{O}(1/\sqrt{T})$ , plus a bias scaling inversely with signal-to-

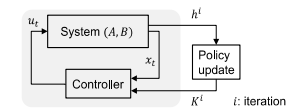


Fig. 1. An illustration of episodic approaches, where  $h^i = (x_0, u_0, \dots, x_{T_i})$  denotes the trajectory of the  $i$ -th episode.

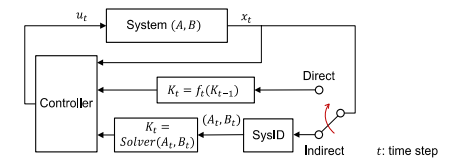


Fig. 2. An illustration of indirect and direct adaptive approaches in closed-loop, where  $f_t$  is some explicit function.

**thanks**