



Off-Policy Reinforcement Learning for Adaptive Optimal Output Tracking of Unknown Linear Discrete-Time Systems

Ci Chen

2020. 09. 25



LUND
UNIVERSITY

Content

- ▲ **Background & Motivation:** Optimal output tracking
- ▲ **Problem Formulation:** Output regulation theory
- ▲ **Method:** Off-policy reinforcement learning
for output-feedback-based optimal control
- ▲ **Simulation**
- ▲ **Conclusion**

1. Background & Motivation

Challenge

How to design tracking control systems with **satisfactory performance without exact model knowledge?**

- Rapid response
- Stability/robustness/safety guarantee
- Optimality for reduced fuel consumption

Aeronautics



Power systems

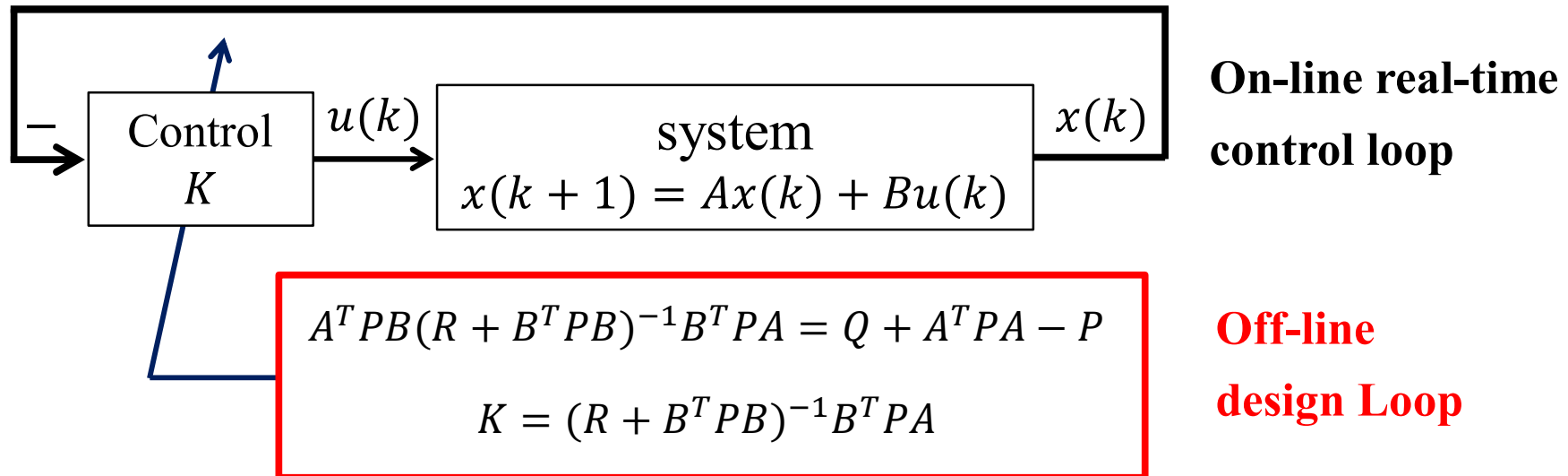


Transportation



1. Background & Motivation

Optimal Control--The Linear Quadratic Regulator (LQR)
if full system dynamics are available.



General goal

We want to find optimal control solutions

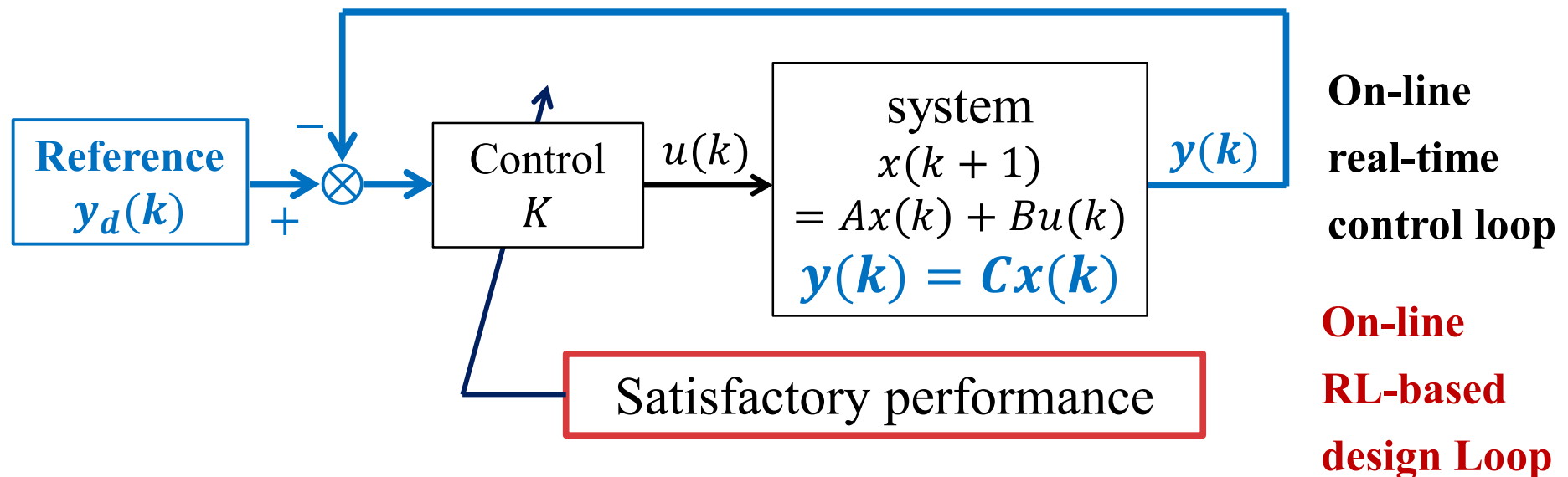
- Online in real-time
- Using adaptive control techniques
- Without knowing the full dynamics

Reinforcement learning (RL)
turns out to be
the key to **this goal!**

1. Background & Motivation

Problem to be studied

How to achieve **Optimal Output Tracking for DT systems** via **Output-Feedback-based Reinforcement learning?**



Content

 Background & Motivation

 **Problem Formulation**

 Method

 Simulation

 Conclusion

2. Problem Formulation (From Tracking to Regulation)

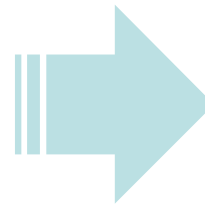
System to be controlled	Reference	Tracking error
$x(k+1) = Ax(k) + Bu(k)$	$x_d(k+1) = Sx_d(k)$	$y_e(k) = y(k) - y_d(k)$
$y(k) = Cx(k)$ $r_n \times r_m$ $r_p \times r_n$	$y_d(k) = Rx_d(k)$	

Control design by the standard output regulation (**full system dynamics**)

Standard controller

$$z(k+1) = Fz(k) - Gy_e(k)$$

$$u(k) = -Kx(k) - Hz(k)$$



Augmented system

$$\begin{cases} e(k+1) = \underline{A}e(k) + \bar{B}u_e(k) \\ y_{e(k)} = \bar{C}e(k) \end{cases}$$

$$u_e(k) = -\bar{K}e(k) \text{ with } \bar{K} = [K, H]$$



$$\underline{A} = \begin{bmatrix} A & O \\ -GC & F \end{bmatrix}$$

$$\bar{C} = [C, O] \quad \bar{B} = \begin{bmatrix} B \\ O \end{bmatrix}$$

Problem:

Feedback gains to be designed

(\underline{A}, \bar{C}) may not be detectable (observable).

A fundamental problem of detectability (observability) is resulted for output-feedback-based design.



2. Problem Formulation (From Tracking to Regulation)

System to be controlled

$$\begin{aligned} x(k+1) &= Ax(k) + Bu(k) \\ y(k) &= Cx(k) \end{aligned} \quad \begin{matrix} r_n \times r_m \\ r_p \times r_n \end{matrix}$$

Reference

$$\begin{aligned} x_d(k+1) &= Sx_d(k) \\ y_d(k) &= Rx_d(k) \end{aligned}$$

Tracking error

$$y_e(k) = y(k) - y_d(k)$$

To ensure the **detectability (observability)**, we design

Our controller

$$\begin{aligned} z(k+1) &= Fz(k) - Gy_e(k) \\ u(k) &= -Kx(k) - Hz(k) - Tz(k) \end{aligned}$$

Augmented system

$$\begin{cases} e(k+1) = \underline{A}e(k) + \bar{B}u_e(k) \\ y_e(k) = \bar{C}e(k) \end{cases}$$

$$u_e(k) = -\bar{K}e(k) \text{ with } \bar{K} = [K, H]$$

Our Solution:

- Given (F, T) detectable (observable) and $r_p \geq r_m$, then (\underline{A}, \bar{C}) is detectable (observable).
- Given any matrix T , then (\underline{A}, \bar{B}) is stabilizable (controllable).

New feedforward gain Tz

$$\underline{A} = \begin{bmatrix} A & -BT \\ -GC & F \end{bmatrix}$$

$$\bar{C} = [C, O] \quad \bar{B} = \begin{bmatrix} B \\ O \end{bmatrix}$$

2. Problem Formulation

A tracking control design problem is now transformed into a regulation-based optimization problem.

Augmented system

$$(*) \begin{cases} e(k+1) = \underline{A}e(k) + \bar{B}u_e(k) \\ y_e(k) = \bar{C}e(k) \end{cases}$$

Problem

$$\begin{cases} \min \sum_{i=k}^{\infty} (y_e^T(i) Q y_e(i) + u_e^T(i) \bar{R} u_e(i)) \\ \text{subject to } (*) \end{cases}$$

Question: How to use **i/o data to **learn the optimal controller** that solves the optimization problem **without exact model knowledge?****

Content

 Background & Motivation

 Problem Formulation

 **Method**

 Simulation

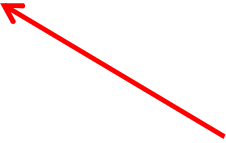
 Conclusion

3. Method (State-Feedback Case)

Design the behavior policy $\left\{ \begin{array}{l} \bar{u}(k) = -\bar{K}^0 r(k) + \xi(k) - T\bar{z}(k) \\ \bar{z}(k+1) = F\bar{z}(k) - G y(k) + G\vartheta(k) \end{array} \right.$

Policy-iteration-based Bellman equation solver in state-feedback form

$$\begin{aligned} & r^T(k+1)P^{j+1}r(k+1) - r^T(k)P^{j+1}r(k) \\ &= -r^T(k)(\bar{Q} + (\bar{K}^j)^T \bar{R} \bar{K}^j)r(k) + \vartheta^T(k)\bar{G}^T P^{j+1}\bar{G}\vartheta(k) \\ & \quad + (-\bar{K}^j r(k) + \bar{u}(k))^T \bar{B}^T P^{j+1} \bar{B} (\bar{K}^j r(k) + \bar{u}(k)) \\ & \quad + 2\vartheta^T(k)\bar{G}^T P^{j+1} \bar{B} u(k) + 2r^T(k)\underline{A}^T P^{j+1} \bar{G}\vartheta(k) \\ & \quad + 2r^T(k)\underline{A}^T P^{j+1} \bar{B} (\bar{K}^j r(k) + \bar{u}(k)). \end{aligned}$$



$$r = [x^T, \bar{z}^T]^T$$

unknown

We seek to reconstruct the state using input and output data.

3. Method (System State Reconstruction)

Reconstruct the state using **input and output data**.

$$\begin{aligned}\zeta_{\bar{u}}(k+1) &= (I_m \otimes A_\zeta)\zeta_{\bar{u}}(k) + \bar{u}(k) \otimes b, \\ \zeta_y(k+1) &= (I_p \otimes A_\zeta)\zeta_y(k) + y(k) \otimes b, \\ \zeta_{\vartheta}(k+1) &= (I_p \otimes A_\zeta)\zeta_{\vartheta}(k) + \vartheta(k) \otimes b,\end{aligned}\quad \bar{\zeta}^T = [\zeta_u^T, \zeta_y^T, \zeta_{\vartheta}^T]^T$$

where $b = [0, 0, \dots, 0, 1]^T$ and A_ζ is a companion matrix¹

1. G. Tao, Adaptive control design and analysis. John Wiley & Sons, 2003

Theorem 1: If the matrix pair (\underline{A}, \bar{B}) is controllable and (\underline{A}, \bar{C}) is observable, then the system state satisfies

$$r = \bar{M}\bar{\zeta} + (\underline{A} - \bar{L}\bar{C})^k r(0)$$

where \bar{M} is a full row rank matrix and $\bar{\zeta}$ is a known vector. The re-expression error $r - \bar{M}\bar{\zeta}$ converges to zero asymptotically.

3. Method (System State Reconstruction)

Theorem 1: If the matrix pair (\underline{A}, \bar{B}) is controllable and (\underline{A}, \bar{C}) is observable, then the system state satisfies

$$r = \bar{M}\bar{\zeta} + (\underline{A} - \bar{L}\bar{C})^k r(0)$$

where \bar{M} is a full row rank matrix and $\bar{\zeta}$ is a known vector. The re-expression error $r - \bar{M}\bar{\zeta}$ converges to zero asymptotically.

Policy iteration-based Bellman equation solver in state-feedback form

$$\begin{aligned} & r^T(k+1)P^{j+1}r(k+1) - r^T(k)P^{j+1}r(k) \\ &= -r^T(k)(\bar{Q} + (\bar{K}^j)^T \bar{R} \bar{K}^j)r(k) + \vartheta^T(k)\bar{G}^T P^{j+1}\bar{G}\vartheta(k) \\ & \quad + (-\bar{K}^j r(k) + \bar{u}(k))^T \bar{B}^T P^{j+1} \bar{B} (\bar{K}^j r(k) + \bar{u}(k)) \\ & \quad + 2\vartheta^T(k)\bar{G}^T P^{j+1} \bar{B} u(k) + 2r^T(k)\underline{A}^T P^{j+1} \bar{G}\vartheta(k) \\ & \quad + 2r^T(k)\underline{A}^T P^{j+1} \bar{B} (\bar{K}^j r(k) + \bar{u}(k)). \end{aligned}$$

$r = [x^T, \bar{z}^T]^T$
unknown

3. Method (Output-Feedback Case)

Solve the optimal control gain **through output feedback**

$$\text{Solve } \mathcal{Q}_o^j \begin{bmatrix} \text{vecs}(\bar{M}^T P^{j+1} \bar{M}) \\ \text{vec}(\bar{M}^T \underline{A}^T P^{j+1} \bar{B}) \\ \text{vecs}(\bar{B}^T P^{j+1} \bar{B}) \\ \text{vec}(\bar{M}^T \underline{A}^T P^{j+1} \bar{G}) \\ \text{vec}(\bar{G}^T P^{j+1} \bar{B}) \\ \text{vecs}(\bar{G}^T P^{j+1} \bar{G}) \end{bmatrix} = v_o^j + \mathcal{D}_{\bar{x}^{j+1}}$$

Collect the input-output data over $[k_i, k_{i+1}]$, $i=0, \dots, f$

$$\mathcal{C}_r = [\text{vecv}(r(k_1)) - \text{vecv}(r(k_0)), \dots, \text{vecv}(r(k_f)) - \text{vecv}(r(k_{f-1}))]^T$$

$$\mathcal{D}_{\bar{K}^j r} = [\text{vecv}(\bar{K}^j r(k_0)), \text{vecv}(\bar{K}^j r(k_1)), \dots, \text{vecv}(\bar{K}^j r(k_{f-1}))]^T$$

$$\mathcal{D}_{\bar{u}} = [\text{vecv}(\bar{u}(k_0)), \text{vecv}(\bar{u}(k_1)), \dots, \text{vecv}(\bar{u}(k_{f-1}))]^T$$

$$\mathcal{D}_r = [\text{vecv}(r(k_0)), \text{vecv}(r(k_1)), \dots, \text{vecv}(r(k_{f-1}))]^T$$

$$\mathcal{D}_g = [\text{vecv}(g(k_0)), \text{vecv}(g(k_1)), \dots, \text{vecv}(g(k_{f-1}))]^T$$

$$\mathcal{D}_{g_r} = [g(k_0) \otimes r(k_0), g(k_1) \otimes r(k_1), \dots, g(k_{f-1}) \otimes r(k_{f-1})]^T$$

$$\mathcal{D}_{rr} = [r(k_0) \otimes r(k_0), r(k_1) \otimes r(k_1), \dots, r(k_{f-1}) \otimes r(k_{f-1})]^T$$

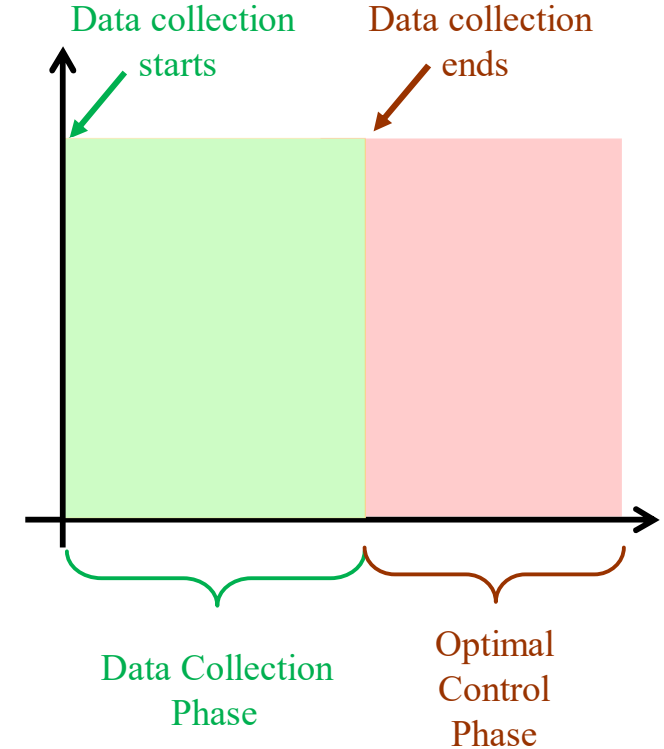


Fig. 1 data processing for off-policy learning

3. Method (Learning From Input-Output Data)

Solve the optimal control gain **through output feedback**

$$\text{Solve } \begin{bmatrix} \text{vecs}(\bar{M}^T P^{j+1} \bar{M}) \\ \text{vec}(\bar{M}^T \underline{A}^T P^{j+1} \bar{B}) \\ \text{vecs}(\bar{B}^T P^{j+1} \bar{B}) \\ \text{vec}(\bar{M}^T \underline{A}^T P^{j+1} \bar{G}) \\ \text{vec}(\bar{G}^T P^{j+1} \bar{B}) \\ \text{vecs}(\bar{G}^T P^{j+1} \bar{G}) \end{bmatrix} = v_o^j + \mathcal{D}_{\bar{x}^{j+1}}$$

Verifiable criterion for checking how much data.

$$\begin{aligned} & \text{rank}([\mathcal{D}_{\bar{\zeta}\bar{\zeta}}, \mathcal{D}_{\bar{u}\bar{\zeta}}, \mathcal{D}_{\bar{u}}, \mathcal{D}_{g\bar{\zeta}}, \mathcal{D}_{\bar{u}g}, \mathcal{D}_g]) \\ &= \frac{1}{2}(n_{\bar{\zeta}}(n_{\bar{\zeta}} + 1) + r_m(r_m + 1) + r_p(r_p + 1)) \\ & \quad + n_{\bar{\zeta}}r_m + r_p n_{\bar{\zeta}} + r_p r_m \end{aligned}$$

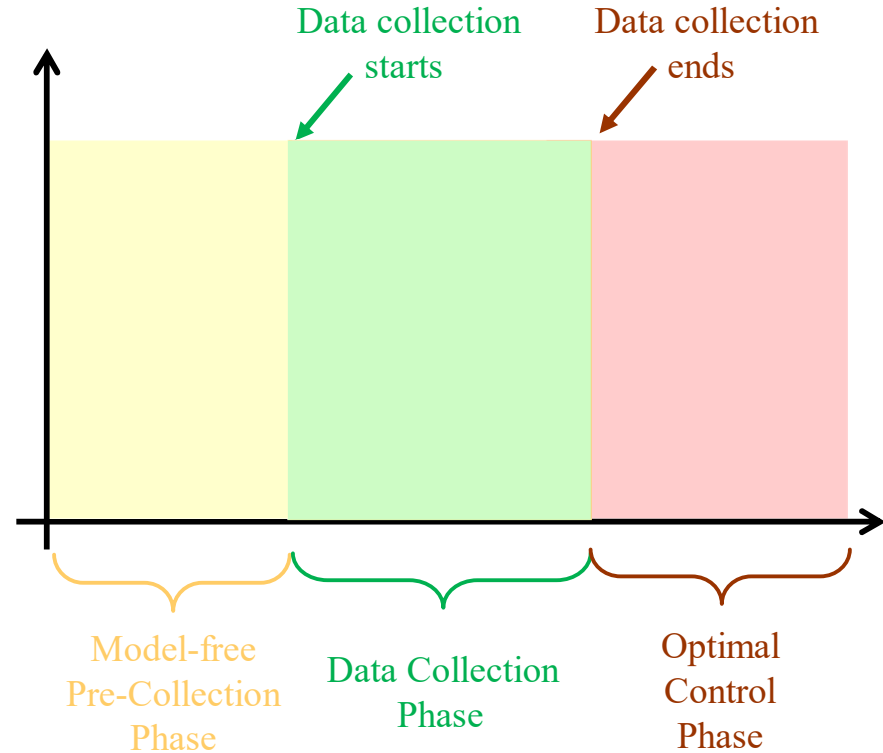


Fig. 2 The **proposed** data processing for off-policy learning

The optimal control gain $(\bar{R} + \bar{B}^T P^{j+1} \bar{B})^{-1} (\bar{M}^T \underline{A}^T P^{j+1} \bar{B})^T$ is uniquely learned from input-output data.

Content

 **Background & Motivation**

 **Problem Formulation**

 **Method**

 **Simulation**

 **Conclusion**

4. Simulation

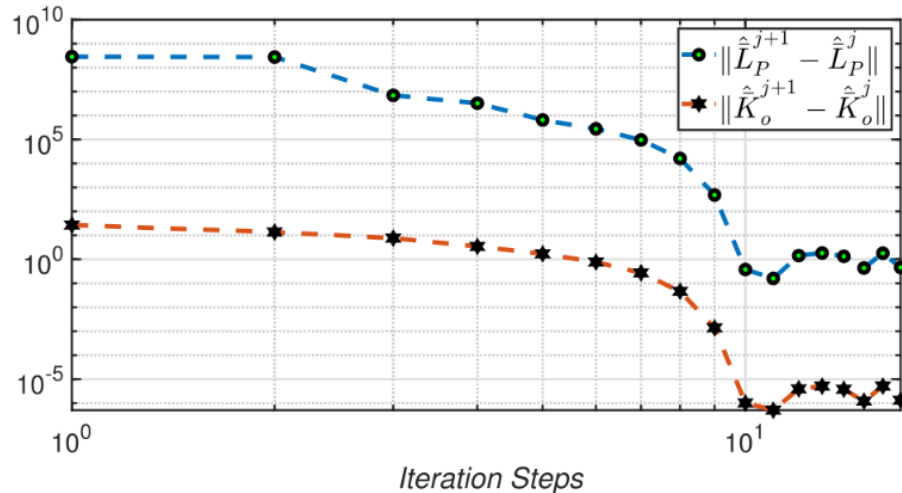
F-16 aircraft dynamics after the discretization

$$A = \begin{bmatrix} 0.887086 & -0.00423047 & -0.00281863 \\ 0.282727 & 0.999573 & -0.00043168 \\ 0.0284386 & 2.93148 & 0.994194 \end{bmatrix},$$

$$B = \begin{bmatrix} 1.77695 \\ 0.272081 \\ 2.82003 \end{bmatrix}, \quad C = [0 \ 57.2958 \ 0].$$

Reference dynamics

$$S = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad R = [1 \ 0].$$



(b) $\|\hat{L}_P^{j+1} - \bar{L}_P^j\|$ and $\|\hat{K}_o^{j+1} - \bar{K}_o^j\|$

Fig. 3 Convergence of the learned control gain

$x(1)$ angle attack; $x(3)$ elevator actuator;
 $x(2)$ pitch rate; $y = x(2)$ pitch rate.

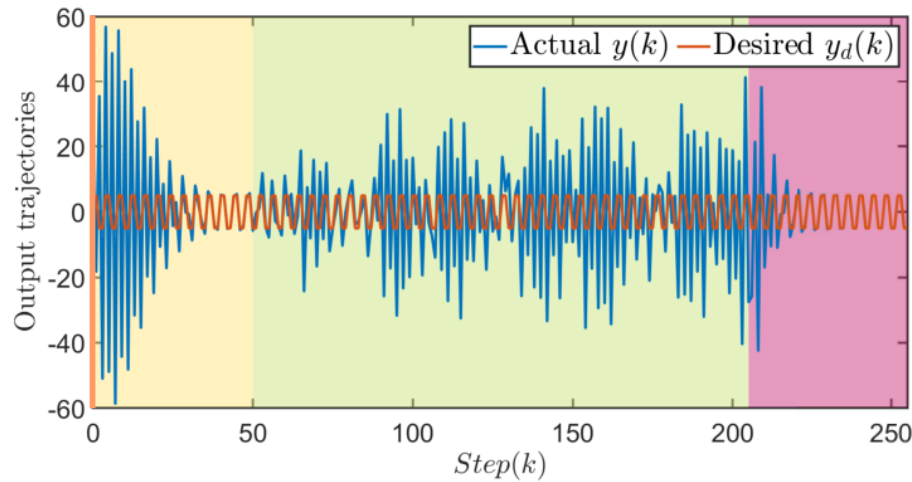
The initial stabilizing control gain as

$$\bar{K}_o^0 = \begin{bmatrix} 16.0669 & 20.11 & 20.3783 \\ 20.8936 & 3.93085 & 7.08518 \\ 19.9149 & 10.5774 & 4.29378 \\ -0.31044 & -8.03354 & -18.0887 \\ -12.2209 & -2.57785 & -0.414603 \end{bmatrix}$$

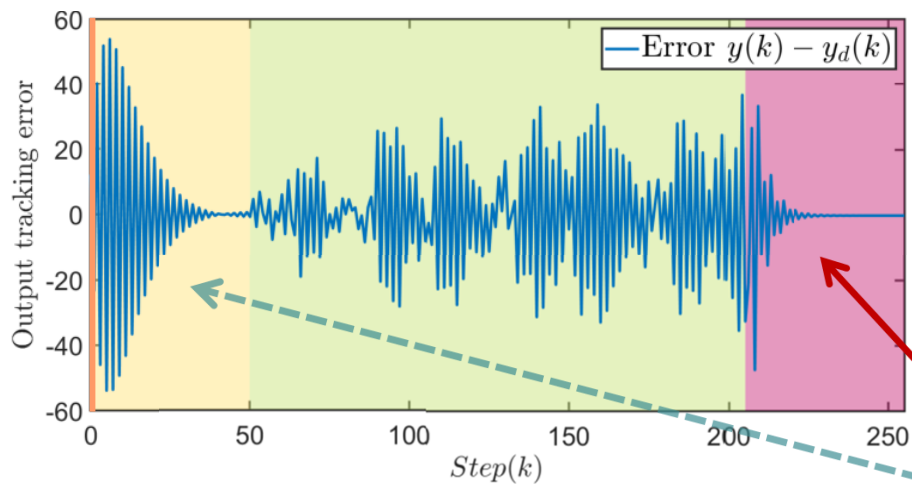
The learned optimal control gain as

$$\bar{K}_o^{10} = \begin{bmatrix} -3.14906 & -1.06988 & 0.068666 \\ -0.161422 & 2.73117 & -1.38856 \\ -2.63657 & 1.70597 & 1.54444 \\ 0.842932 & 1.57443 & 2.10917 \\ -1.08647 & -2.09949 & -0.47113 \end{bmatrix}$$

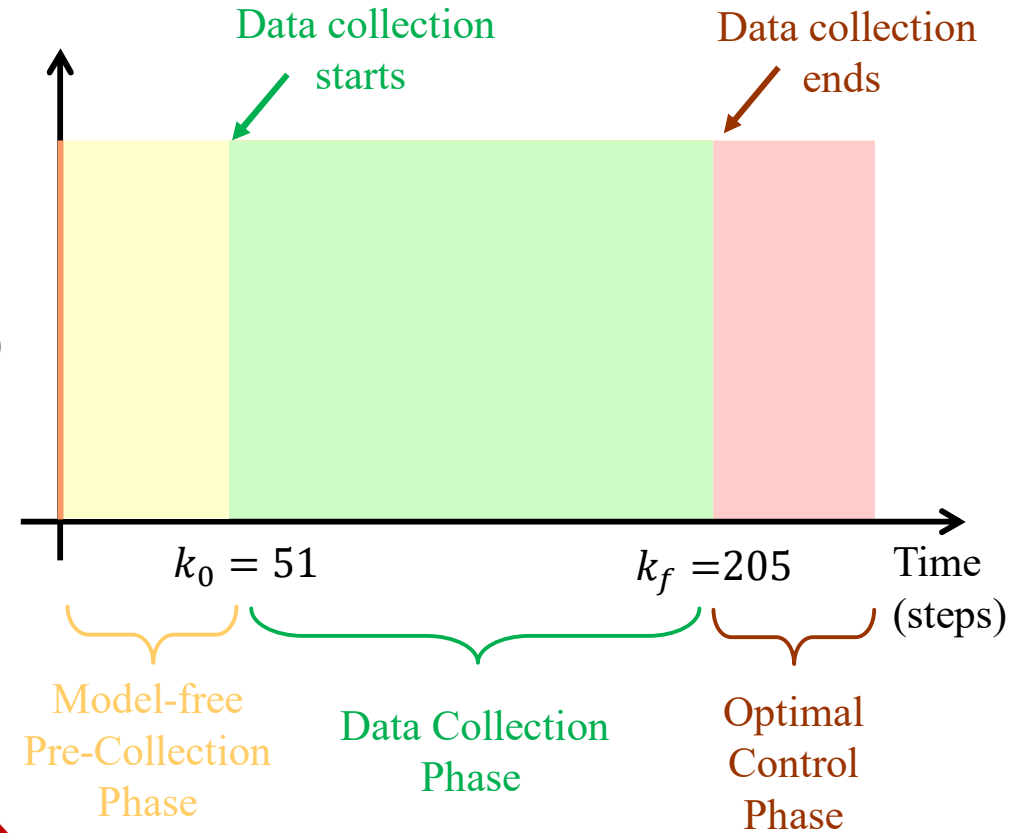
4. Simulation



(a) Output performance $y(k), y_d(k)$



(b) Tracking error $y(k) - y_d(k)$



Performance Improved

Fig. 4 Output trajectories over the time

Content

 **Background & Motivation**

 **Problem Formulation**

 **Method**

 **Simulation**

 **Conclusion**

5. Conclusion

Adaptive Optimal Output Tracking for DT systems via Output-Feedback-Based Reinforcement Learning

1. We proposed **an output regulation and off-policy RL-based controller to formulate adaptive optimal output tracking problem for DT systems.**
2. We derived a **verifiable rank condition** to ensure the uniqueness of the optimal control gain learned from **input-output data.**
3. We proposed **a model-free pre-collection phase** to supplement the off-policy learning for DT systems.



Thank you



LUND
20
UNIVERSITY