



# Non-Asymptotic System Identification

CDC'23 and <https://arxiv.org/abs/2309.03873>

Anastasios Tsiamis (ETH)

# System Identification

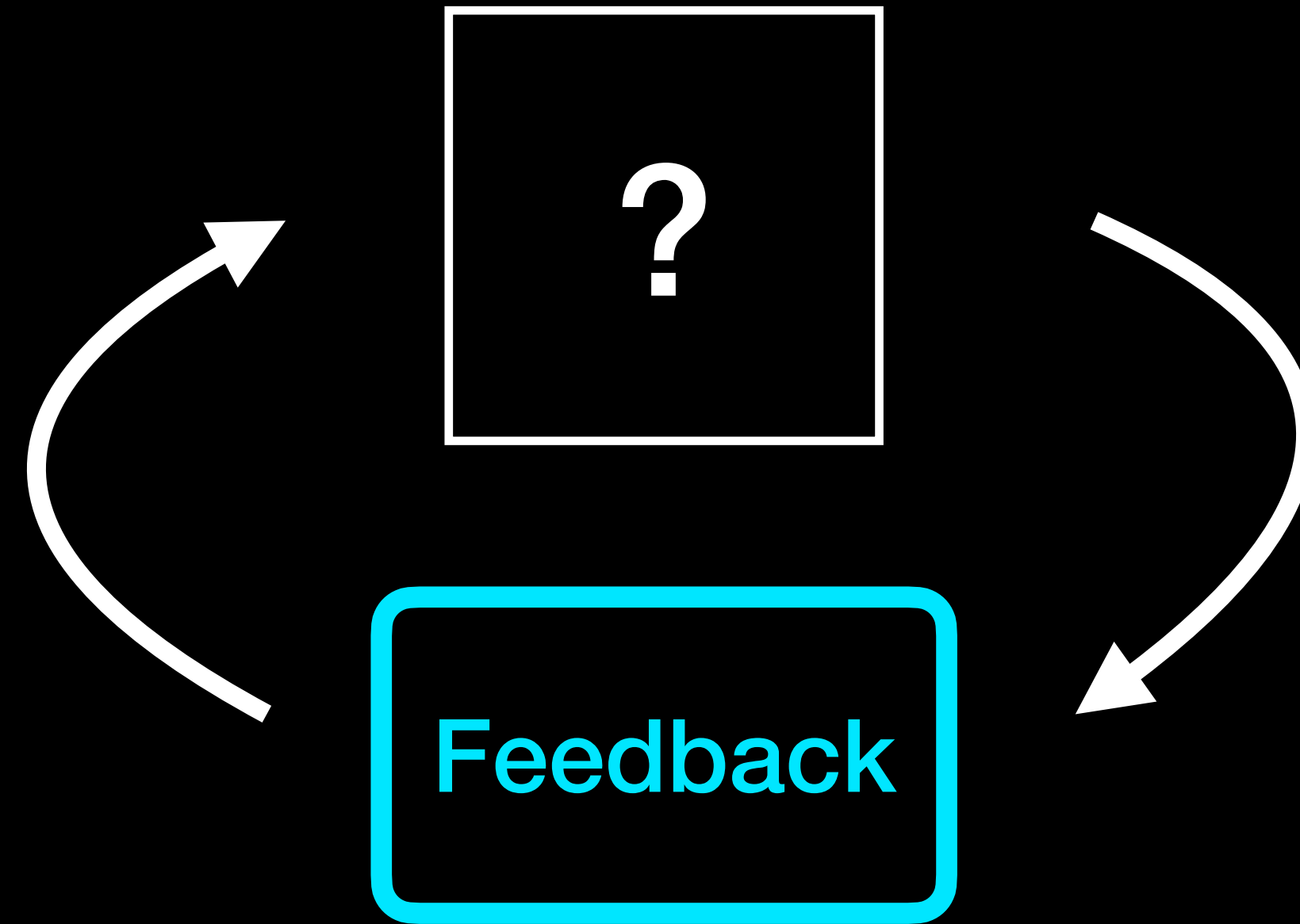


**Experiments:** excite system

**SysID:** learn model

Today's focus: **SysID**

# Why SysID?

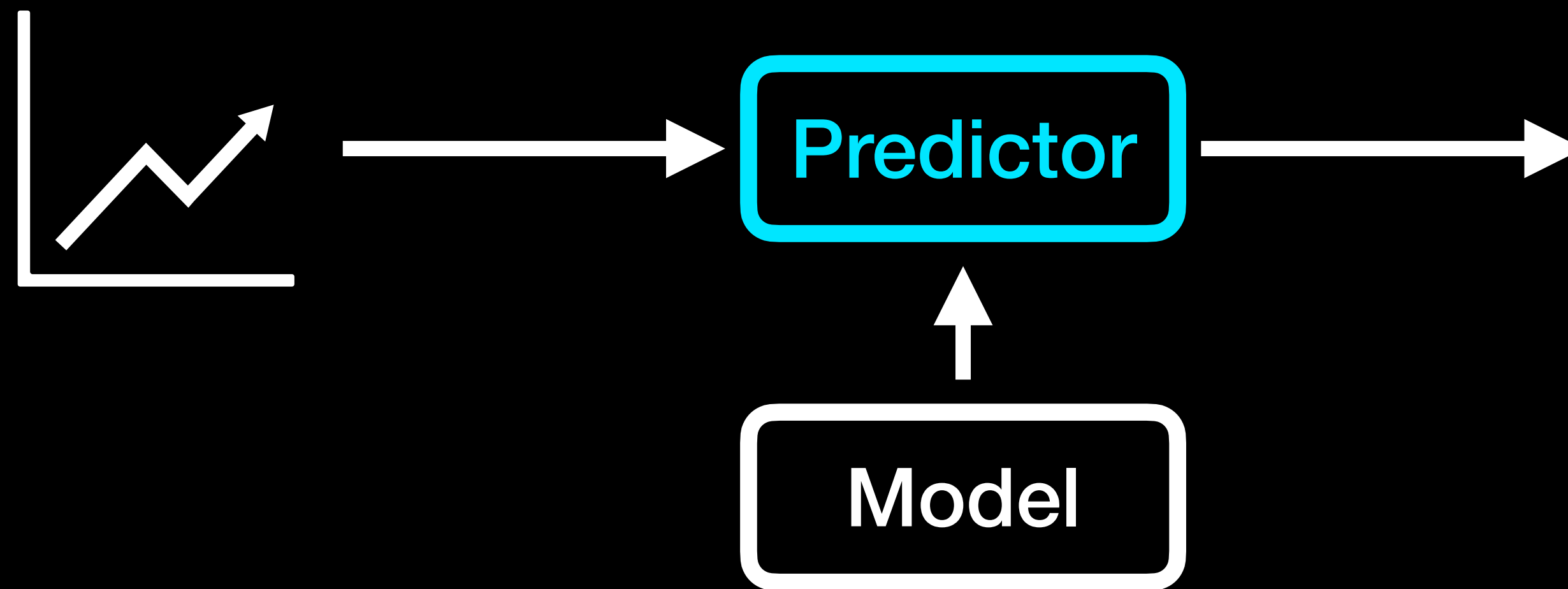


Classical **control** design: requires model

Unknown or complex system?

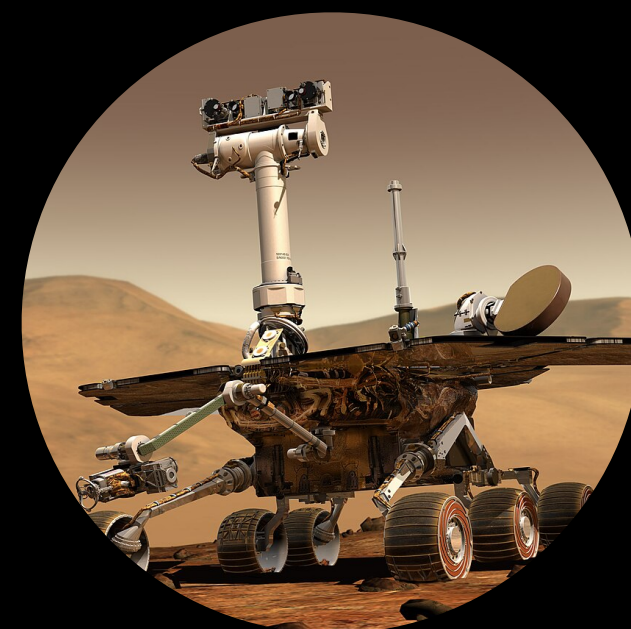
Learn from data

# Why SysID?

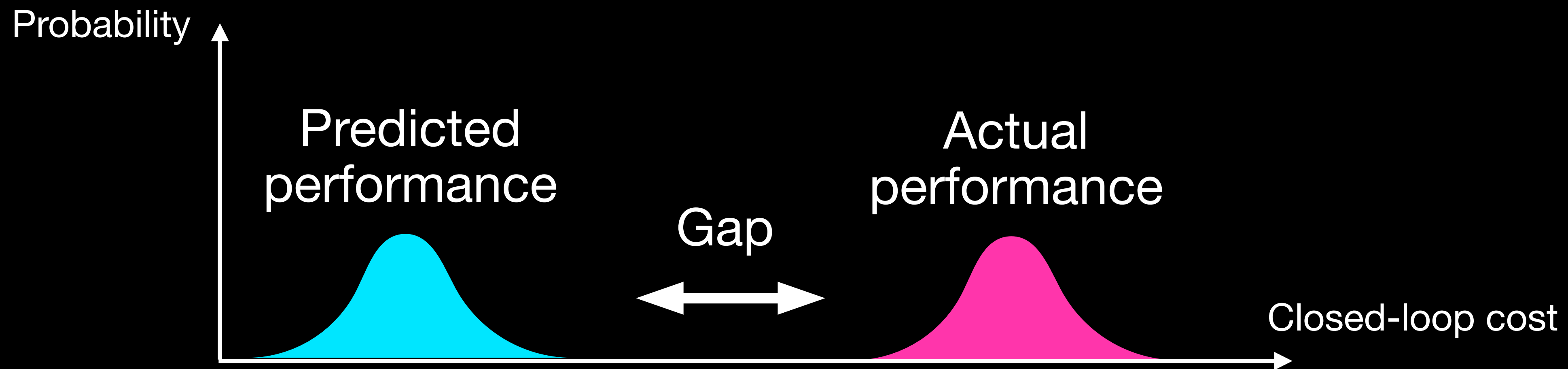
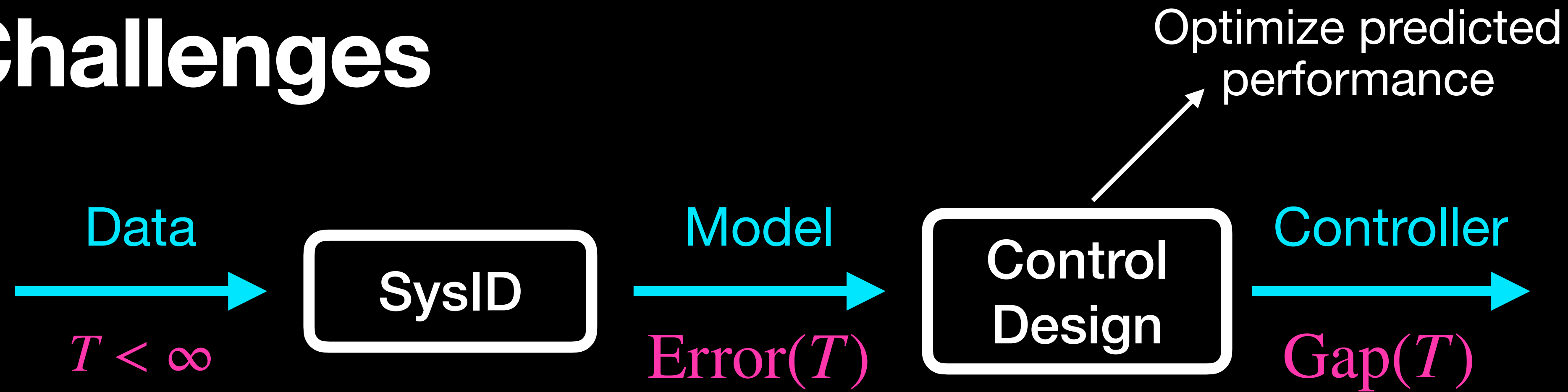


Also: forecasting, fault detection

# Ubiquitous Applications



# Challenges



Many factors (**statistical**, modeling bias, experiment design, sim2real, ...)

Today: **Finite** Data  $T < \infty$



# Finite Sample Error of SysID



How fast does the **error decay** with number  $T$  of data?

How does it depend on system properties (size, structure, noise)?

# Linear System Identification

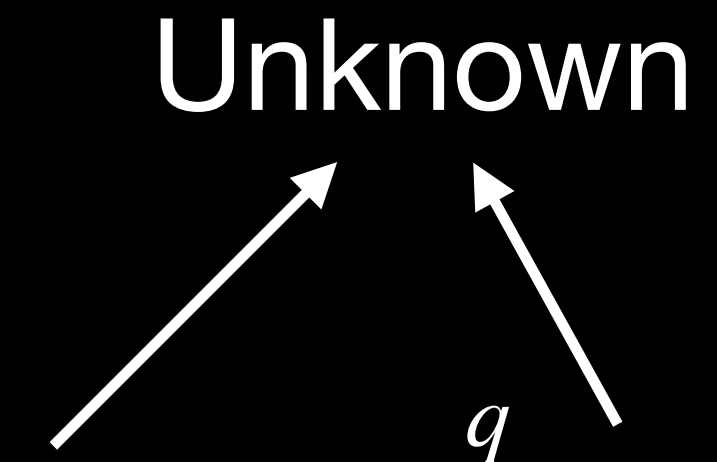


# Linear Systems

- ARX:

$$Y_t = \sum_{i=1}^p A_i^* Y_{t-i} + \sum_{j=1}^q B_j^* U_{t-j} + W_t$$

Unknown



- Building block for learning state-space systems

$$Z_{t+1} = A^* Z_t + B^* U_t + K^* W_t$$
$$Y_t = C^* Z_t + W_t$$

- Linear: **simple** but **non-trivial** class

Stay tuned for non-linear!

# Assumptions

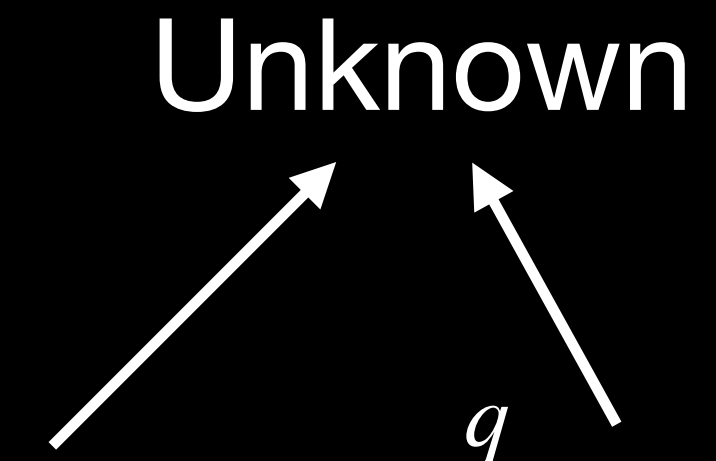
$$Y_t = \sum_{i=1}^p A_i^* Y_{t-i} + \sum_{j=1}^q B_j^* U_{t-j} + W_t$$

1. Marginally **Stable** system (poles on or inside unit circle)
2. **Known**  $p, q$  (upper bounds)
3.  $\sigma^2$ -**sub-Gaussian** noise, **zero-mean**, **independent**
4. White-noise **inputs**
5. Single trajectory **data**:  $Y_0, U_0, \dots, Y_T, U_T$

# Problem

$$Y_t = \sum_{i=1}^p A_i^* Y_{t-i} + \sum_{j=1}^q B_j^* U_{t-j} + W_t$$

Unknown



Given data:  $Y_0, U_0, \dots, Y_T, U_T$

Return  $\hat{A}_1^*, \dots, \hat{A}_p^*, \hat{B}_1^*, \dots, \hat{B}_q^*$  and finite-sample bounds for the error

# ARX in the Least Squares Framework

$$Y_t = \sum_{i=1}^p \theta_i^* X_{t-i} + \sum_{j=1}^q B_j^* U_{t-j} + W_t$$

Instance of least squares setup

$$\begin{bmatrix} Y_{t-1} \\ \vdots \\ Y_{t-p} \end{bmatrix} \quad \begin{bmatrix} U_{t-1} \\ \vdots \\ U_{t-q} \end{bmatrix}$$



$$X_t = \begin{bmatrix} Y_{t-1:t-p}^\top & U_{t-1:t-q}^\top \end{bmatrix}^\top$$

$$[A_{1:p}^* \quad B_{1:q}^*]$$



$$\theta^*$$

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^{d_Y \times d_X}} \left\{ \frac{1}{T} \sum_{t=1}^T \|Y_t - \theta X_t\|_2^2 \right\}$$

# ARX Least Squares error

$$Y_t = \theta^* X_t + W_t$$

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathbb{R}^{d_Y \times d_X}} \left\{ \frac{1}{T} \sum_{t=1}^T \|Y_t - \theta X_t\|_2^2 \right\}$$

$$\hat{\theta} - \theta^* = \left( \sum_{t=1}^T W_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2}$$

Sample Covariance  
Lower Tail Methods

Self-normalized  
martingale methods

# Bounding the terms

$$\Sigma_t \triangleq \mathbf{E}X_t X_t^\top$$

$$\frac{1}{T} \sum_{t=1}^T X_t X_t^\top \geq c \Sigma_\tau, \quad \text{for } T \geq T_{\text{burn}}(\tau, \delta)$$

Empirical cov.    Covariance

} With prob.  
1 -  $\delta$

$$\left( \sum_{t=1}^T W_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \leq c' \sigma^2 (d_X + \log \det \Sigma_T \Sigma_\tau^{-1} + \log \frac{1}{\delta})$$

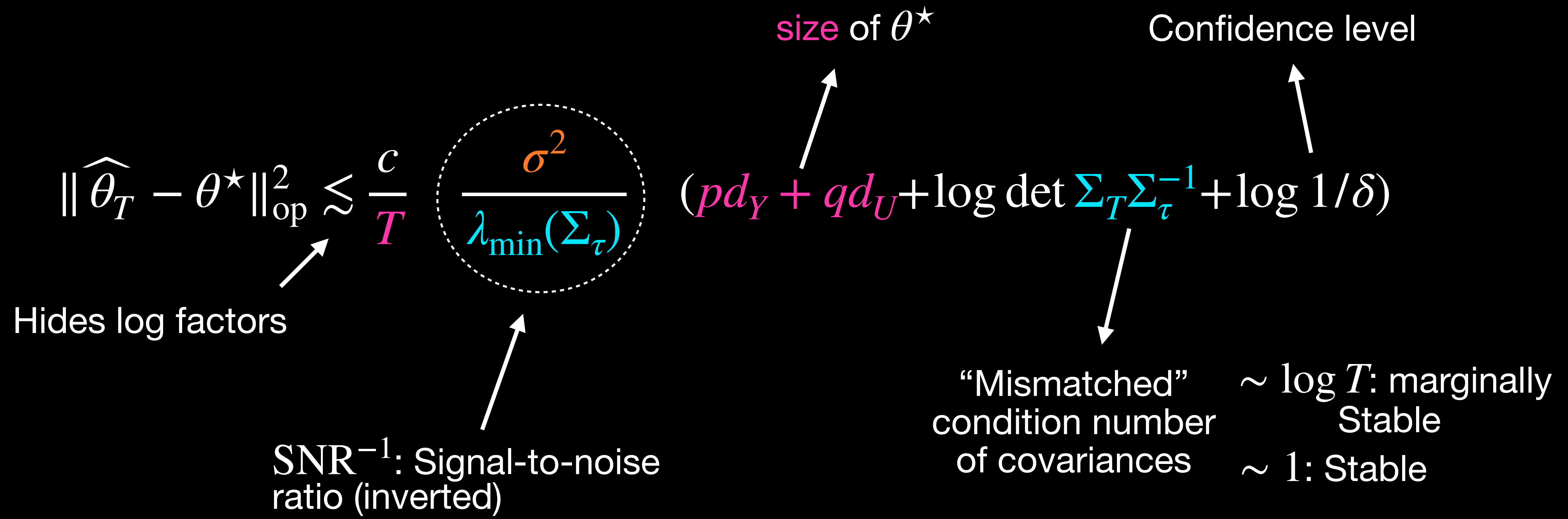
Mismatch?    Empirical Cov. at  $T$   
Covariance at some  $\tau < T$

$$\hat{\theta} - \theta^* = \left( \sum_{t=1}^T W_t X_t^\top \right) \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2} \left( \sum_{t=1}^T X_t X_t^\top \right)^{-1/2}$$

$\tau$ : tunable to reduce mismatch  
**But:** It increases burn-in time

# Finite-sample bound

$$\Sigma_t \triangleq \mathbf{E}X_t X_t^\top$$



Non-asymptotic rate of  $1/T$

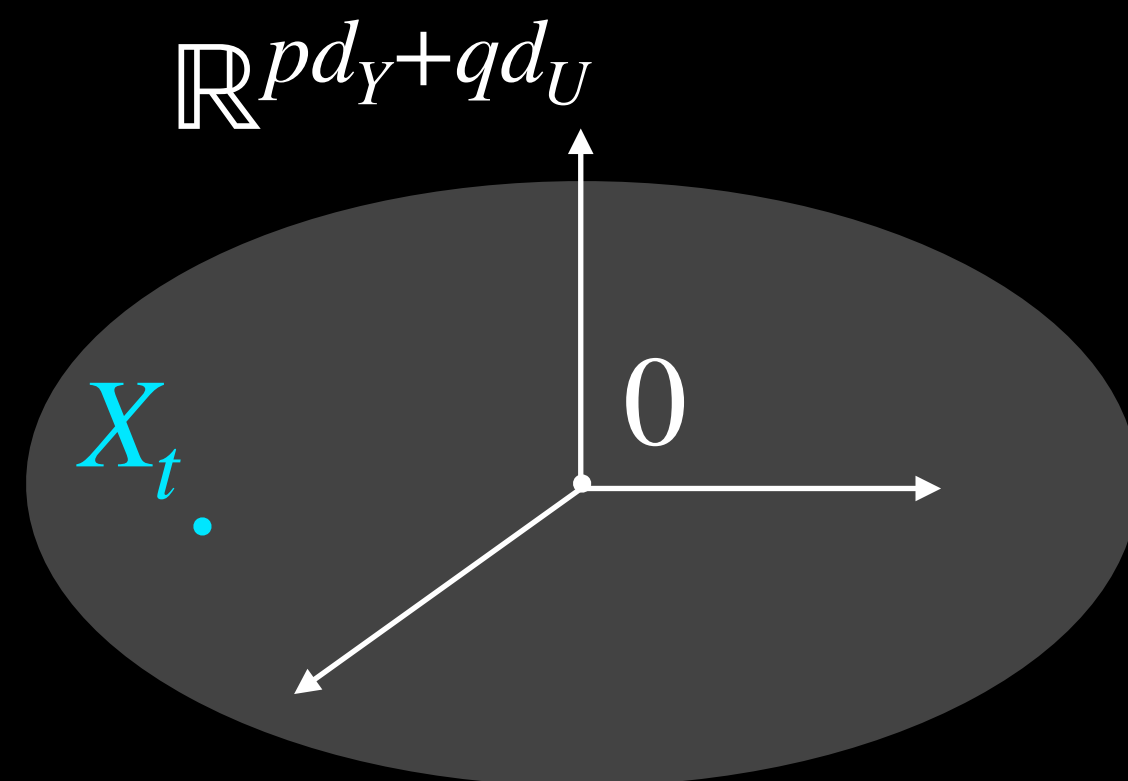
Agrees with asymptotics (Central Limit Theorem)

# Persistency of excitation

$$\Sigma_t \triangleq \mathbf{E}X_t X_t^\top$$

$$\frac{1}{T} \sum_{t=1}^T X_t X_t^\top \geq c \Sigma_\tau, \quad \text{for } T \geq T_{\text{burn}}(\tau, \delta)$$

Before  $T_{\text{burn}}$  matrix  
can be singular



After  $T_{\text{burn}}$  matrix invertible &  
persistently away from zero

$$T_{\text{burn}} \simeq c' \tau (\log 1/\delta + (pd_Y + qd_U) \dots)$$

Trades excitation with  
larger Burn-in

Dimension of  $X_t$



# Tuning the bounds

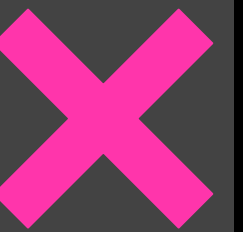
$$\Sigma_t \triangleq \mathbf{E}X_t X_t^\top$$

$$\frac{1}{T} \sum_{t=1}^T X_t X_t^\top \geq c \Sigma_\tau, \quad \text{for } T \geq T_{\text{burn}}(\tau, \delta)$$

$$T_{\text{burn}} \simeq c\tau (\log 1/\delta + (pd_Y + qd_U)\dots)$$

Tune  $\tau$ ?

$$\tau = 1 < p, q$$



$$X_1 = \begin{bmatrix} Y_0^\top & 0 & \dots & U_0^\top & 0 & \dots & 0 \end{bmatrix}^\top$$

Singular  $\Sigma_\tau$

$$\tau = \max\{p, q\}$$



$$X_\tau = \begin{bmatrix} Y_{\tau-1:\tau-p}^\top & U_{\tau-1:\tau-q}^\top \end{bmatrix}^\top$$

Invertible  $\Sigma_\tau$  (details in paper)

$\tau$  arbitrary large



only if  $T \geq T_{\text{burn}}(\tau, \delta) \simeq c\tau(\dots)$

Feasible if  $\tau = o(T)$ , e.g.  $\tau = \sqrt{T}$

But large burn-in!

# Tuning for Stable Systems

$$\Sigma_t \triangleq \mathbf{E}X_t X_t^\top$$

$$\|\widehat{\theta}_T - \theta^\star\|_{\text{op}}^2 \lesssim \frac{c}{T} \frac{\sigma^2}{\lambda_{\min}(\Sigma_\tau)} (pd_Y + qd_U + \log \det \Sigma_T \Sigma_\tau^{-1} + \log 1/\delta)$$

For **stable plants** we have  $\Sigma_t \rightarrow \Sigma$

$$\|\widehat{\theta}_T - \theta^\star\|_{\text{op}}^2 \lesssim \frac{c'}{T} \frac{\sigma^2}{\lambda_{\min}(\Sigma)} (pd_Y + qd_U + \log 1/\delta)$$

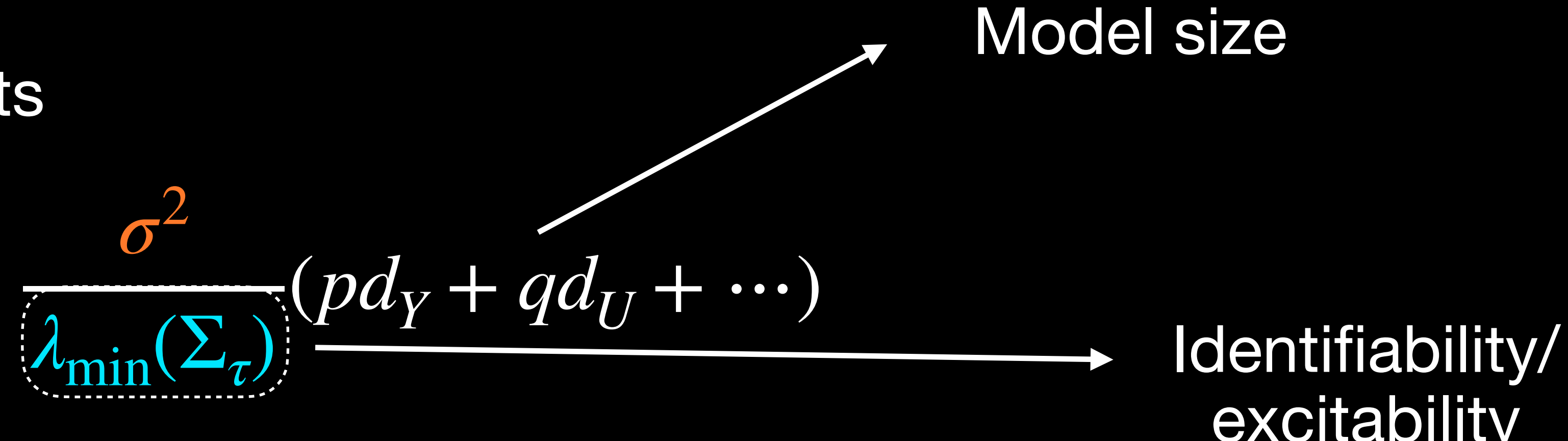
where  $\tau = \tau_{\text{mix}}$ : **mixing** time (**settling** time of dominant pole)

# Remarks

Guarantees valid even if  $\tau = \max\{p, q\} < \tau_{\text{mix}}$

- Stability or mixing is not required
- Such tradeoffs not captured by asymptotics

System theoretic constants

$$\bullet \quad \|\widehat{\theta}_T - \theta^*\|_{\text{op}}^2 \lesssim \frac{c}{T} \frac{\sigma^2}{\lambda_{\min}(\Sigma_\tau)} (pd_Y + qd_U + \dots)$$


- SNR depends on system's **controllability structure**
- Both **size** and **structure** affect difficulty

# Remarks

$$\|\widehat{\theta}_T - \theta^*\|_{\text{op}}^2 \lesssim \frac{c}{T} \frac{\sigma^2}{\lambda_{\min}(\Sigma_\tau)} (\dots)$$

operator norm: **worst-case** error vs **worst-case** excitation

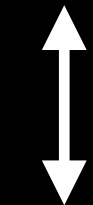
Alternative statements are possible, e.g.

$$(\widehat{\theta}_T - \theta^*) \Sigma_T (\widehat{\theta}_T - \theta^*)^\top \preceq \frac{c}{T} \sigma^2 (\dots)$$

# Generalization to State Space

$$\begin{aligned} Z_t &= A^* Z_{t-1} + B^* U_{t-1} + K^* W_{t-1} \\ Y_{t-1} &= C^* Z_{t-1} + W_{t-1} \end{aligned}$$

Innovation Form



Under conditions, e.g. Gaussianity

$$\begin{aligned} Z_t &= A^* Z_{t-1} + B^* U_{t-1} + W_{t-1} \\ Y_{t-1} &= C^* Z_{t-1} + V_{t-1} \end{aligned}$$

Standard form

1. (Marginally) **Stable** system
2.  $\sigma^2$ -**sub-Gaussian** noise, **zero-mean**, **independent**
3. White-noise **inputs**
4. **Minimum-phase**:  $A^* - K^* C^*$  is stable (linked to observability)

# Reformulation

Can we approximate it as ARX?

$$\begin{aligned} \boxed{Z_t} &= A^* Z_{t-1} + B^* U_{t-1} + K^* W_{t-1} \\ Y_{t-1} &= C^* Z_{t-1} + \boxed{W_{t-1}} \end{aligned}$$

$$\begin{aligned} A_{cl}^* &\triangleq A^* - K^* C^* \\ A_i^* &\triangleq C^* (A_{cl}^*)^{i-1} K^* \\ B_i^* &\triangleq C^* (A_{cl}^*)^{i-1} B^* \end{aligned}$$

$$\begin{aligned} Y_t &= C^* \boxed{Z_t} + W_t \\ &= C^* A^* Z_{t-1} + C^* B^* U_{t-1} + C^* K^* \boxed{W_{t-1}} + W_t \\ &= C^* A_{cl}^* Z_{t-1} + C^* B^* U_{t-1} + C^* K^* Y_{t-1} + W_t \\ &= C^* (A_{cl}^*)^2 Z_{t-2} + C^* A_{cl}^* B^* U_{t-2} + C^* A_{cl}^* K^* Y_{t-2} \\ &\quad + C^* B^* U_{t-1} + C^* K^* Y_{t-1} + W_t \\ &= \sum_{i=1}^p A_i^* Y_{t-i} + \sum_{j=1}^p B_j^* U_{t-j} + W_t + C^* (A_{cl}^*)^p Z_{t-p} \end{aligned}$$

# ARX Approximation

$$Y_t = \sum_{i=1}^p A_i^* Y_{t-i} + \sum_{j=1}^p B_j^* U_{t-j} + W_t + C^* (A_{cl}^*)^p Z_{t-p} \xrightarrow{\text{Bias term}}$$

$$A_{cl}^* \triangleq A^* - K^* C^*$$

$$A_i^* \triangleq C^* (A_{cl}^*)^{i-1} K^*$$

$$B_i^* \triangleq C^* (A_{cl}^*)^{i-1} B^*$$

Minimum-phase

$$\|(A_{cl}^*)^p\| \leq c\rho^p$$

For  $p = c' \log T$

$$Y_t \approx \theta_p^* X_t + W_t$$

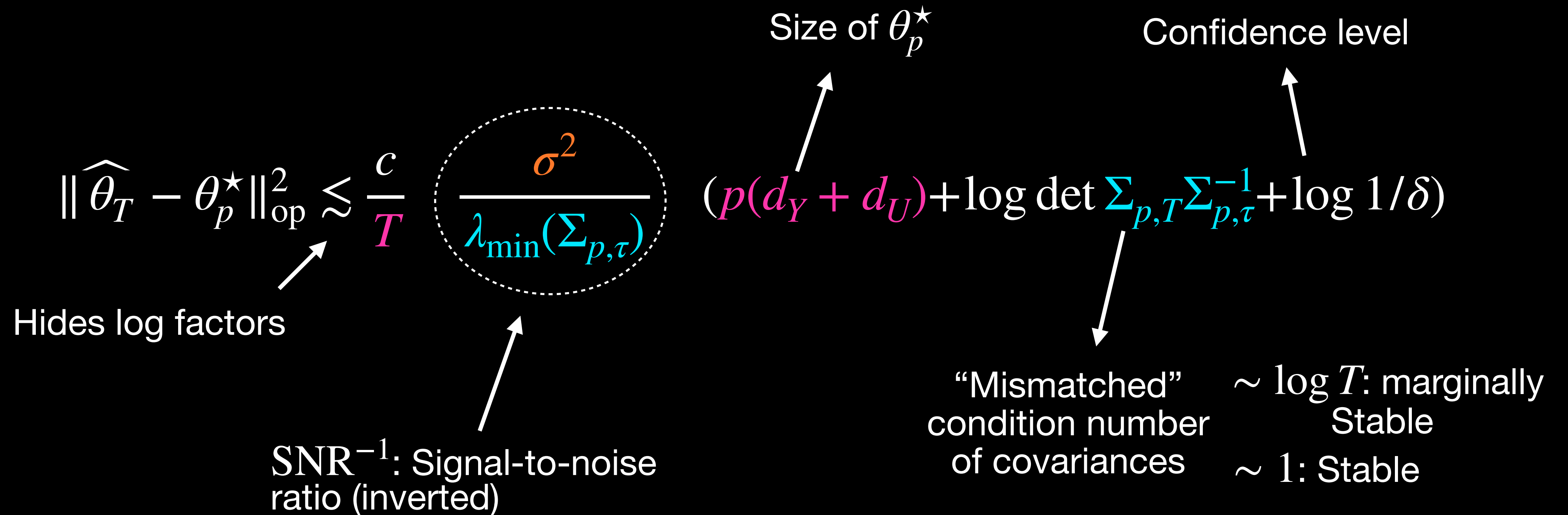
$$X_t = \begin{bmatrix} Y_{t-1:t-p}^\top & U_{t-1:t-p}^\top \end{bmatrix}^\top$$

$$\theta_p^* = \begin{bmatrix} A_{1:p}^* & B_{1:p}^* \end{bmatrix}$$

Learn Markov Parameters  $C^* A_{cl}^{i} B^*$   
instead of original  $A^*, B^*, C^*, K^*$

# State-Space bound

$$\Sigma_{p,t} \triangleq \mathbf{E} X_t X_t^\top$$



Non-asymptotic rate of  $1/T$

But  $p = c \log T$



# Persistency of excitation

$$\Sigma_{p,t} \triangleq \mathbf{E}X_t X_t^\top$$

$$\frac{1}{T} \sum_{t=1}^T X_t X_t^\top \geq c \Sigma_{p,\tau}, \quad \text{for } T \geq T_{\text{burn}}(\tau, \delta) \quad p = c' \log T$$

$$T_{\text{burn}} \simeq c\tau \left( \log 1/\delta + (p(d_Y + d_U)) \dots \right)$$

$$= c''\tau \left( \log 1/\delta + (\log T(d_Y + d_U)) \dots \right)$$

$$= c'''\log T \left( \log 1/\delta + (\log T(d_Y + d_U)) \dots \right)$$

$$\tau = p$$



$$X_\tau = \begin{bmatrix} Y_{\tau-1:\tau-p}^\top & U_{\tau-1:\tau-p}^\top \end{bmatrix}^\top$$

Invertible  $\Sigma_{p,\tau}$  (details in paper)

# Comments

Guarantees valid even if  $\tau = \max\{p, q\} < \tau_{\text{mix}}$

- (Strict) Stability or mixing (of  $A^*$ ) is not required

However we need minimum phase

- Stability of  $A^* - K^*C^*$  (linked to observability, not mixing)

Improper learning

- Learn Markov Parameters  $C^*A_{c_1}^iB^*$  instead of original  $A^*, B^*, C^*, K^*$
- Recover original using realization theory

# Moving forward

Improve universal constants  $c$

Rates are upper bounds. Are they optimal? What about lower bounds?

*Statistical Learning for Control  
A finite sample perspective*



Nonlinear systems

# Thank you!

*Statistical Learning for Control  
A finite sample perspective*



*A Tutorial on the Non-  
Asymptotic Theory of System  
Identification*

