



Data-Driven Control and Data-Poisoning attacks in Buildings: the KTH Live-In Lab case study

Alessio Russo*, Marco Molinari and Alexandre Proutiere

Mediterranean Conference on Control and Automation (MED), 2021

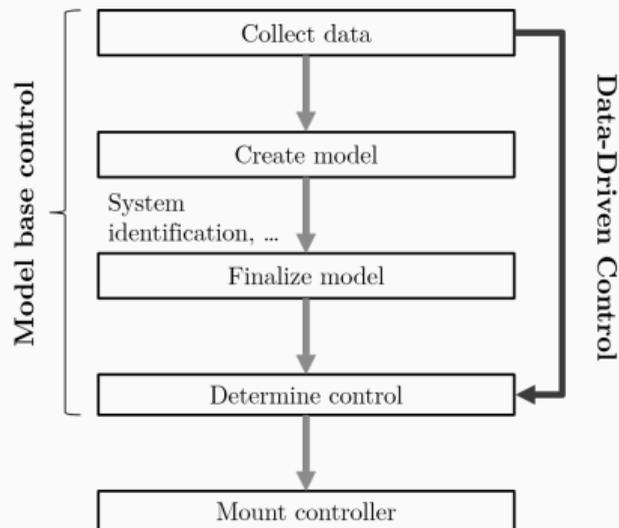
KTH, Royal Institute of Technology, Stockholm

Problem motivation and background



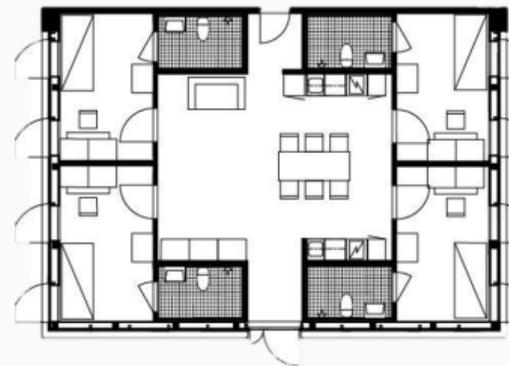
- **Temperature control in buildings may be complicated.**
- Data-driven control approaches: use data to directly compute a control law.
 1. Model-reference based methods: Virtual Reference Feedback Tuning (VRFT) [1], Iterative Feedback Tuning [2], correlation-based [3]...
 2. Methods based on Willems' et al. lemma [4,5].
- **The data can be poisoned.**
- We focus on VRFT, a popular model-reference based method.

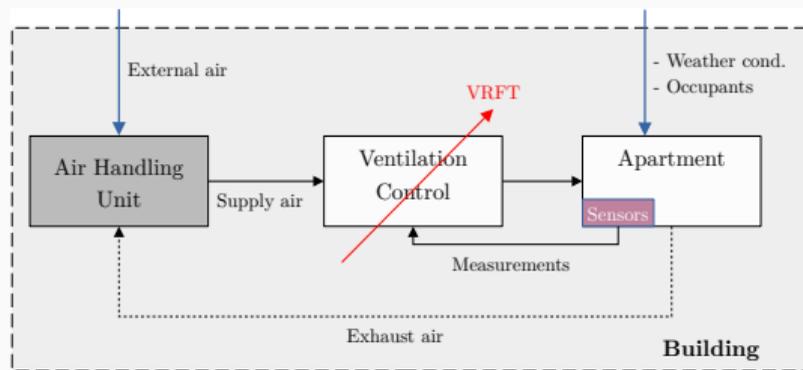
Problem motivation



- **Temperature control in buildings may be complicated.**
- Data-driven control approaches: use data to directly compute a control law.
 1. Model-reference based methods: Virtual Reference Feedback Tuning (VRFT) [1], Iterative Feedback Tuning [2], correlation-based [3]...
 2. Methods based on Willems' et al. lemma [4,5].
- **The data can be poisoned.**
- We focus on VRFT, a popular model-reference based method.

KTH Live-in Lab Testbed

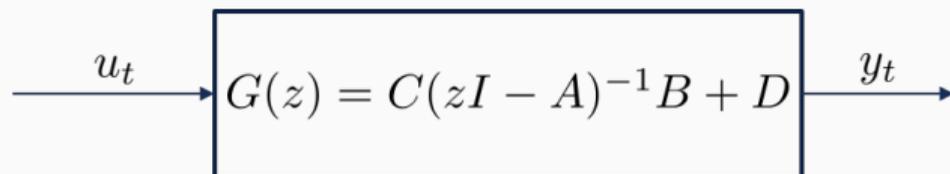




1. We modeled the building using **IDA-ICE**, a building performance simulation software [6].
2. We focused on the problem of **ventilation control of a single apartment**.
3. We applied **VRFT** to derive a control law, directly from the data of an (empty) apartment.
4. Finally, we tested whether VRFT is susceptible to **data poisoning attacks**.

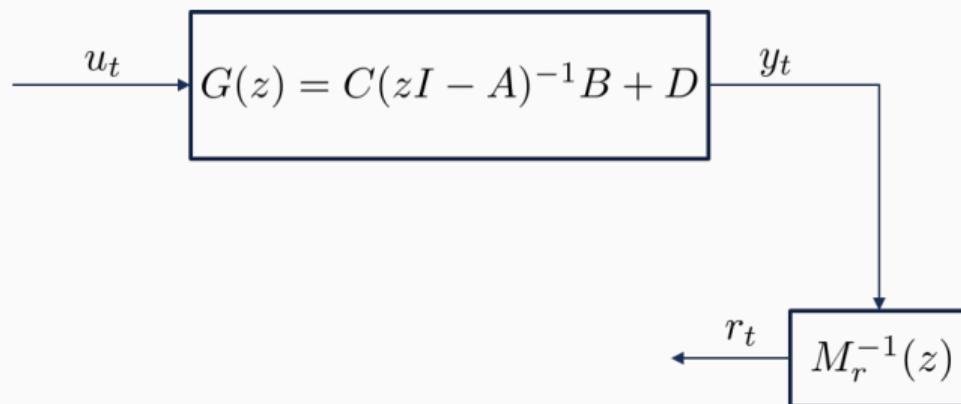
Temperature control

Virtual Reference Feedback Tuning [1]



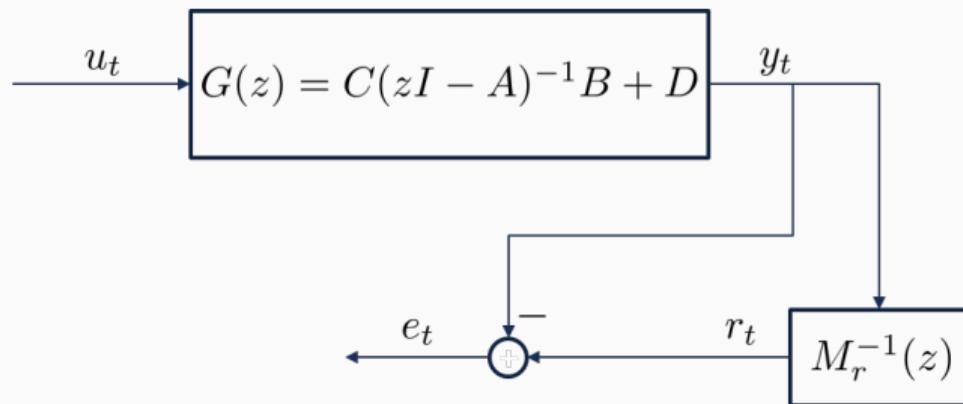
1. Feed a pre-designed signal u_t and measure y_t .

Virtual Reference Feedback Tuning [1]



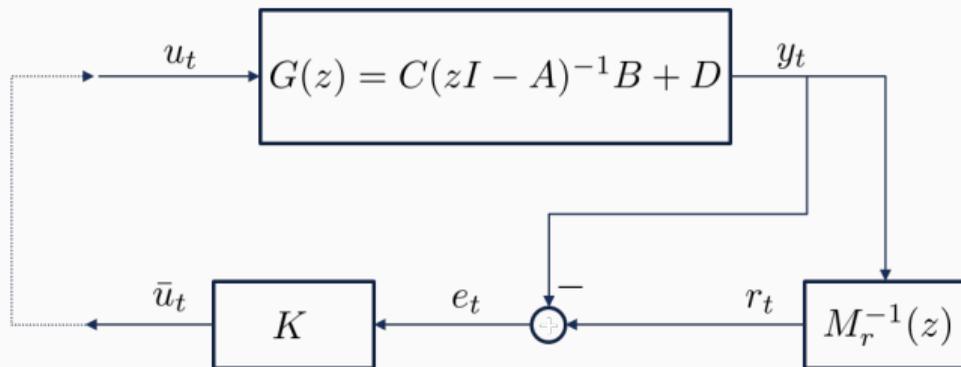
1. Feed a pre-designed signal u_t and measure y_t .
2. Given a reference model $M_r(z)$, compute the reference signal r_t .

Virtual Reference Feedback Tuning [1]



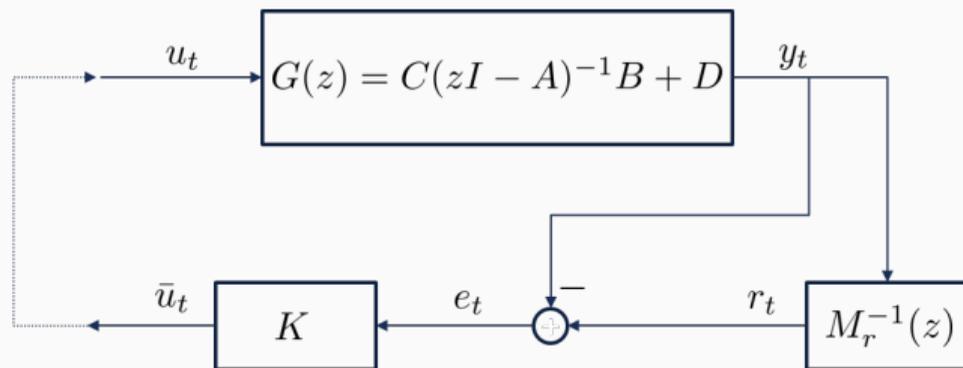
1. Feed a pre-designed signal u_t and measure y_t .
2. Given a reference model $M_r(z)$, compute the reference signal r_t .
3. Compute the *virtual error* $e_t = r_t - y_t$.

Virtual Reference Feedback Tuning [1]



1. Feed a pre-designed signal u_t and measure y_t .
2. Given a reference model $M_r(z)$, compute the reference signal r_t .
3. Compute the *virtual error* $e_t = r_t - y_t$.
4. Design a control law K that outputs a signal \bar{u}_t that is *close* to u_t .

Virtual Reference Feedback Tuning [1]

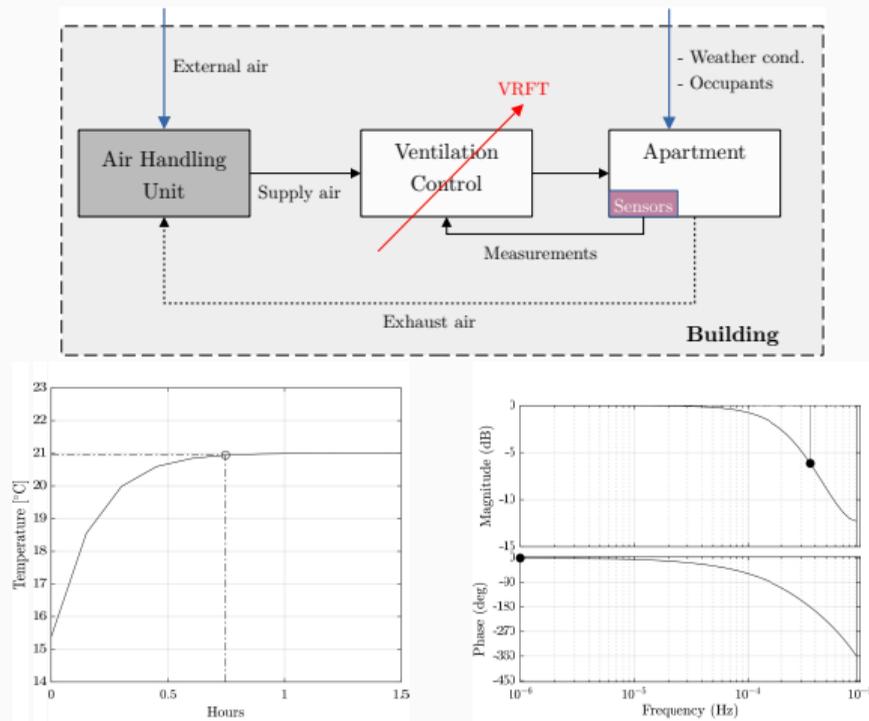


1. Feed a pre-designed signal u_t and measure y_t .
2. Given a reference model $M_r(z)$, compute the reference signal r_t .
3. Compute the *virtual error* $e_t = r_t - y_t$.
4. Design a control law K that outputs a signal \bar{u}_t that is *close* to u_t .

Under some assumptions, it is possible to show that minimizing $\frac{1}{N} \sum_{t=1}^N (\bar{u}_t - u_t)^2$, for $N \rightarrow \infty$, yields a law K that converges to the minimum of

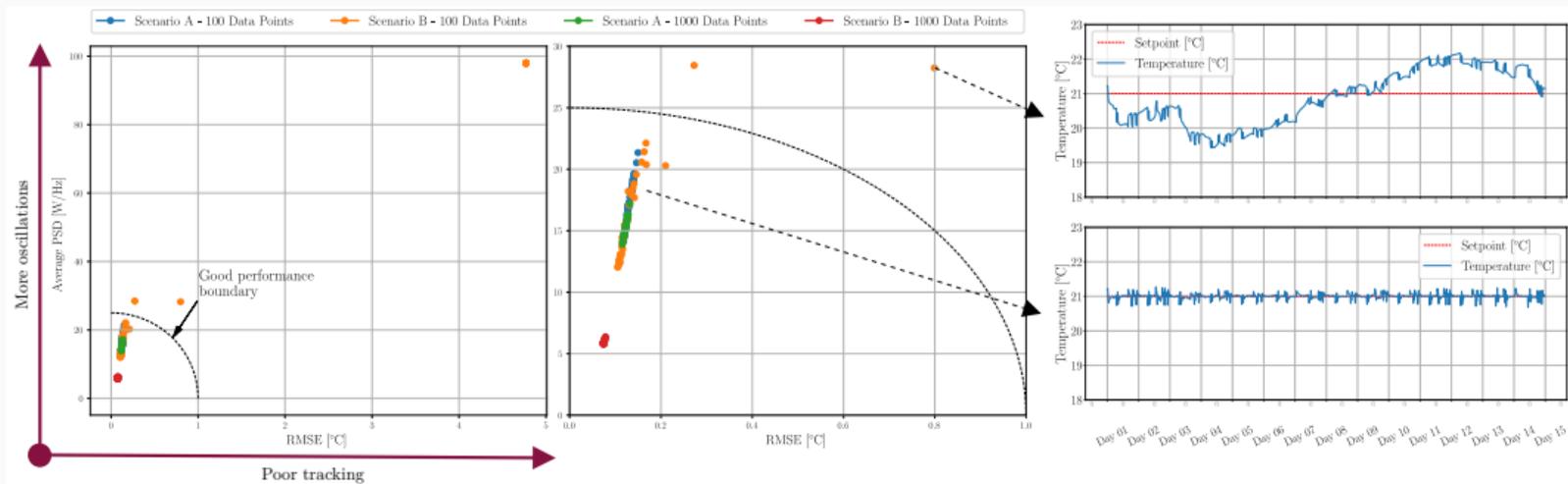
$$\min_K \|M_r(z) - (1 - M_r)KG(z)\|_2.$$

Temperature control: method



1. Data was sampled every 540 [s].
2. The control signal is a real number in $[0, 1]$. We designed 2 experiments for VRFT.
 - **Scenario A:** $u_t \sim \mathcal{N}(\frac{1}{2}, \frac{1}{6})$.
 - **Scenario B:** $u_t \sim \mathcal{N}(\frac{1}{2}, 1)$.
3. Goal of VRFT: compute $K_\theta(z)$, where $K_\theta(z) = \sum_{k=1}^3 \theta_k \frac{z^{-k+2}}{z-1}$ (PID-like controller).
4. We used a 2nd order reference model (see plot on the left).

Temperature control: results



1. **Scenario A:** $u_t \sim \mathcal{N}(\frac{1}{2}, \frac{1}{6})$; **Scenario B:** $u_t \sim \mathcal{N}(\frac{1}{2}, 1)$.
2. January was used for training of VRFT (empty apartment); February for evaluation of the control law (1 person). For each case we run 50 simulations.

Data poisoning

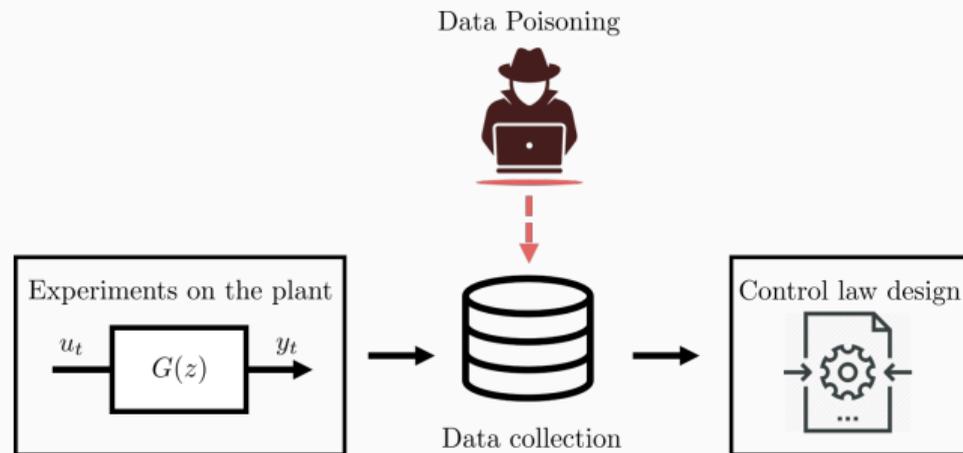


Figure 1: Data poisoning scheme.

Attack Formulation

We can cast the attacker's problem as a bi-level optimization problem.

$$\begin{aligned} \max_{\mathbf{u}', \mathbf{y}'} \quad & \mathcal{A}(\mathbf{u}, \mathbf{y}, K(\mathbf{u}', \mathbf{y}')) \\ \text{s.t.} \quad & K(\mathbf{u}', \mathbf{y}') \in \arg \min_K \mathcal{L}(\mathbf{u}', \mathbf{y}', K) \\ & \|\mathbf{u}' - \mathbf{u}\|_2 \leq \varepsilon_u \|\mathbf{u}\|_2, \quad \|\mathbf{y}' - \mathbf{y}\|_2 \leq \varepsilon_y \|\mathbf{y}\|_2, \end{aligned}$$

- We denote by $u'_t = u_t + a_{u,t}$ the poisoned input, where \mathbf{a}_u is the poisoning signal (similarly for y'_t).
- We denote by \mathcal{L} the learner's criterion (e.g., the MSE loss of VRFT).
- Similarly, \mathcal{A} is the attacker's criterion.

Attack based on Russo, A., Proutiere, A.. *Poisoning attacks against data-driven control methods*. *American Control Conference, 2021*.

VRFT: Attack Formulation

1. **Remember the VRFT criterion** $\frac{1}{N} \sum_{t=1}^N (u_t - \bar{u}_t)^2$, where $\bar{u}_t = K_\theta(z)(M_r^{-1}(z) - 1)y_t$.
2. The learner's criterion under attack can be rewritten in matrix form as

$$\mathcal{L}(\mathbf{u}', \mathbf{y}', \theta) = \frac{1}{N} \|\mathbf{u}' - \Phi(\mathbf{y}')\theta\|_2^2$$

where Φ is a matrix that includes the effect of $M_r(z)$ (ref. model) and $K_\theta(z)$.

3. **How do we choose the attacker's criterion? Simplest choice is to just maximize the original VRFT criterion!**

$$\begin{aligned} \max_{\mathbf{u}', \mathbf{y}'} \quad & \mathcal{A}(\mathbf{u}, \mathbf{y}, \hat{\theta}(\mathbf{u}', \mathbf{y}')) = \frac{1}{N} \left\| \mathbf{u} - \Phi(\mathbf{y})\hat{\theta}(\mathbf{u}', \mathbf{y}') \right\|_2^2 \\ \text{s.t.} \quad & \hat{\theta}(\mathbf{u}', \mathbf{y}') = (\Phi^\top(\mathbf{y}')\Phi(\mathbf{y}'))^{-1} \Phi^\top(\mathbf{y}')\mathbf{u}' \\ & \|\mathbf{u}' - \mathbf{u}\|_2 \leq \varepsilon_u \|\mathbf{u}\|_2, \quad \|\mathbf{y}' - \mathbf{y}\|_2 \leq \varepsilon_y \|\mathbf{y}\|_2. \end{aligned}$$

The problem is concave in the input signal \mathbf{u}' , and non-convex in the output signal \mathbf{y}' .

VRFT: Attack Formulation

1. **Remember the VRFT criterion** $\frac{1}{N} \sum_{t=1}^N (u_t - \bar{u}_t)^2$, where $\bar{u}_t = K_\theta(z)(M_r^{-1}(z) - 1)y_t$.
2. The learner's criterion under attack can be rewritten in matrix form as

$$\mathcal{L}(\mathbf{u}', \mathbf{y}', \theta) = \frac{1}{N} \|\mathbf{u}' - \Phi(\mathbf{y}')\theta\|_2^2$$

where Φ is a matrix that includes the effect of $M_r(z)$ (ref. model) and $K_\theta(z)$.

3. **How do we choose the attacker's criterion? Simplest choice is to just maximize the original VRFT criterion!**

$$\begin{aligned} \max_{\mathbf{u}', \mathbf{y}'} \quad & \mathcal{A}(\mathbf{u}, \mathbf{y}, \hat{\theta}(\mathbf{u}', \mathbf{y}')) = \frac{1}{N} \left\| \mathbf{u} - \Phi(\mathbf{y})\hat{\theta}(\mathbf{u}', \mathbf{y}') \right\|_2^2 \\ \text{s.t.} \quad & \hat{\theta}(\mathbf{u}', \mathbf{y}') = (\Phi^\top(\mathbf{y}')\Phi(\mathbf{y}'))^{-1} \Phi^\top(\mathbf{y}')\mathbf{u}' \\ & \|\mathbf{u}' - \mathbf{u}\|_2 \leq \varepsilon_u \|\mathbf{u}\|_2, \quad \|\mathbf{y}' - \mathbf{y}\|_2 \leq \varepsilon_y \|\mathbf{y}\|_2. \end{aligned}$$

The problem is concave in the input signal \mathbf{u}' , and non-convex in the output signal \mathbf{y}' .

VRFT: Attack Formulation

Input: Data-set (\mathbf{u}, \mathbf{y}) ; objective function \mathcal{A} ;

parameters $\varepsilon_u, \varepsilon_y, \eta$

Output: Attack vectors $\mathbf{a}_u, \mathbf{a}_y$

$i \leftarrow 0, (\mathbf{a}_u^{(i)}, \mathbf{a}_y^{(i)}) \leftarrow (\mathbf{0}, \mathbf{0})$

$\hat{\theta}^{(i)} \leftarrow \hat{\theta}(\mathbf{u} + \mathbf{a}_u^{(i)}, \mathbf{y} + \mathbf{a}_y^{(i)})$

$J^{(i)} \leftarrow \mathcal{A}(\mathbf{u}, \mathbf{y}, \hat{\theta}^{(i)})$

do

$\mathbf{a}_u^{(i+1)} \leftarrow$ solve attacker's problem in \mathbf{a}_u

using CCP [9]

$\mathbf{a}_y^{(i+1)} \leftarrow \text{PGA}(\varepsilon_y, \hat{\theta}(\mathbf{u} + \mathbf{a}_u^{(i+1)}, \mathbf{y} + \mathbf{a}_y^{(i)}))$

$\hat{\theta}^{(i+1)} \leftarrow \hat{\theta}(\mathbf{u} + \mathbf{a}_u^{(i+1)}, \mathbf{y} + \mathbf{a}_y^{(i+1)})$

$J^{(i+1)} \leftarrow \mathcal{A}(\mathbf{u}, \mathbf{y}, \hat{\theta}^{(i+1)})$

$i \leftarrow i + 1$

while $|J^{(i+1)} - J^{(i)}| > \eta$

-Remember that $\mathbf{u}' = \mathbf{u} + \mathbf{a}_y$ (resp. \mathbf{y}').

-The attacker wants to solve

$$\max_{\mathbf{u}', \mathbf{y}'} \frac{1}{N} \left\| \mathbf{u} - \Phi(\mathbf{y}) \hat{\theta}(\mathbf{u}', \mathbf{y}') \right\|_2^2$$

$$\text{s.t. } \hat{\theta}(\mathbf{u}', \mathbf{y}') = (\Phi^\top(\mathbf{y}') \Phi(\mathbf{y}'))^{-1} \Phi^\top(\mathbf{y}') \mathbf{u}'$$

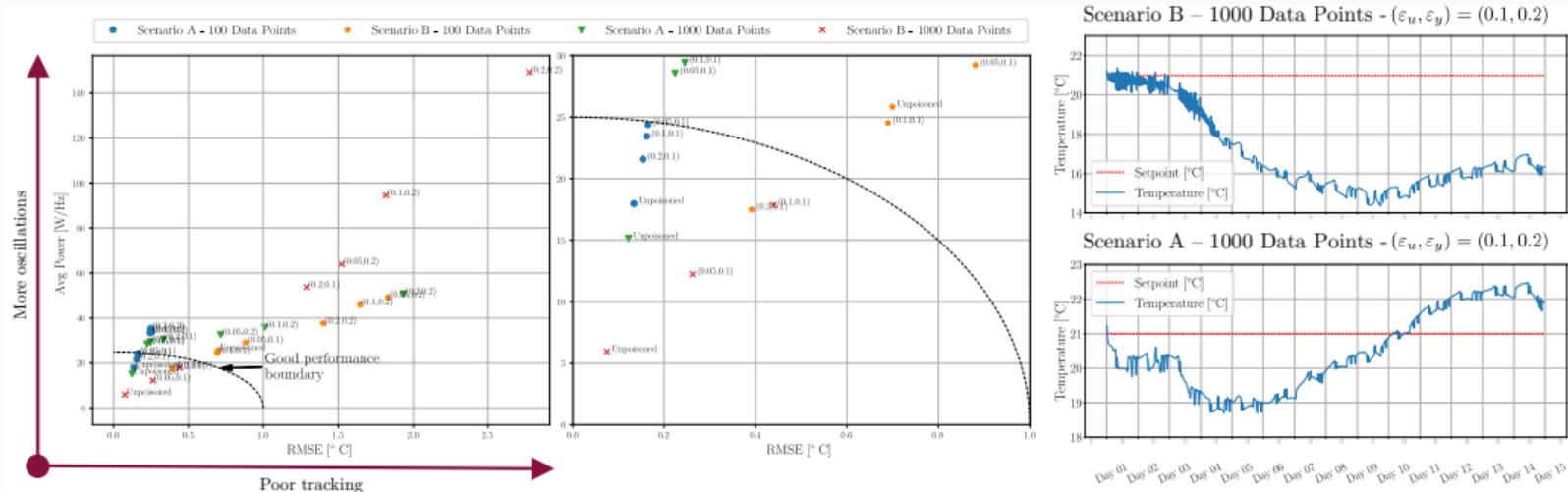
$$\|\mathbf{u}' - \mathbf{u}\|_2 \leq \varepsilon_u \|\mathbf{u}\|_2, \quad \|\mathbf{y}' - \mathbf{y}\|_2 \leq \varepsilon_y \|\mathbf{y}\|_2$$

-The problem is concave in the input signal \mathbf{u}' :

we use convex-concave programming techniques.

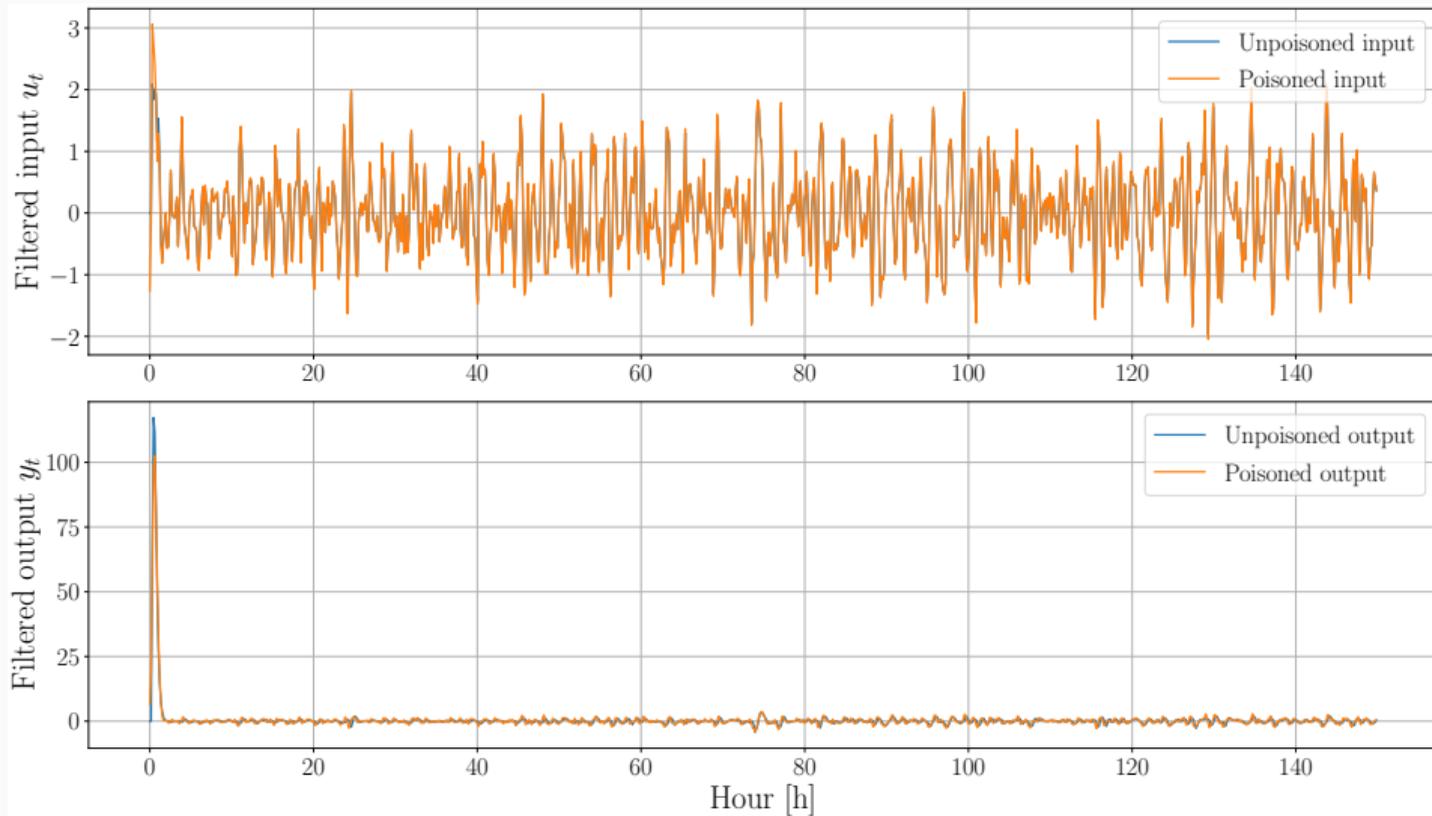
-The problem is non-convex in the output signal \mathbf{y}' : **we use projected gradient ascent.**

Data poisoning: results



1. **Scenario A:** $u_t \sim \mathcal{N}(\frac{1}{2}, \frac{1}{6})$; **Scenario B:** $u_t \sim \mathcal{N}(\frac{1}{2}, 1)$.
2. Each point on the left plots represents the average across 50 simulations for a specific set of values $(\epsilon_u; \epsilon_y)$, displayed on the top of each point (also the unpoisoned cases are depicted in the plots).

Data poisoning: original vs poisoned data



Conclusions

- Data-driven control methods can be used to derive control laws directly from data.
- Data Poisoning is not a new concept in Machine Learning (see Biggio et al. [10]).
- We must pay attention to the security aspects of data-driven methods!

Thank you for listening!

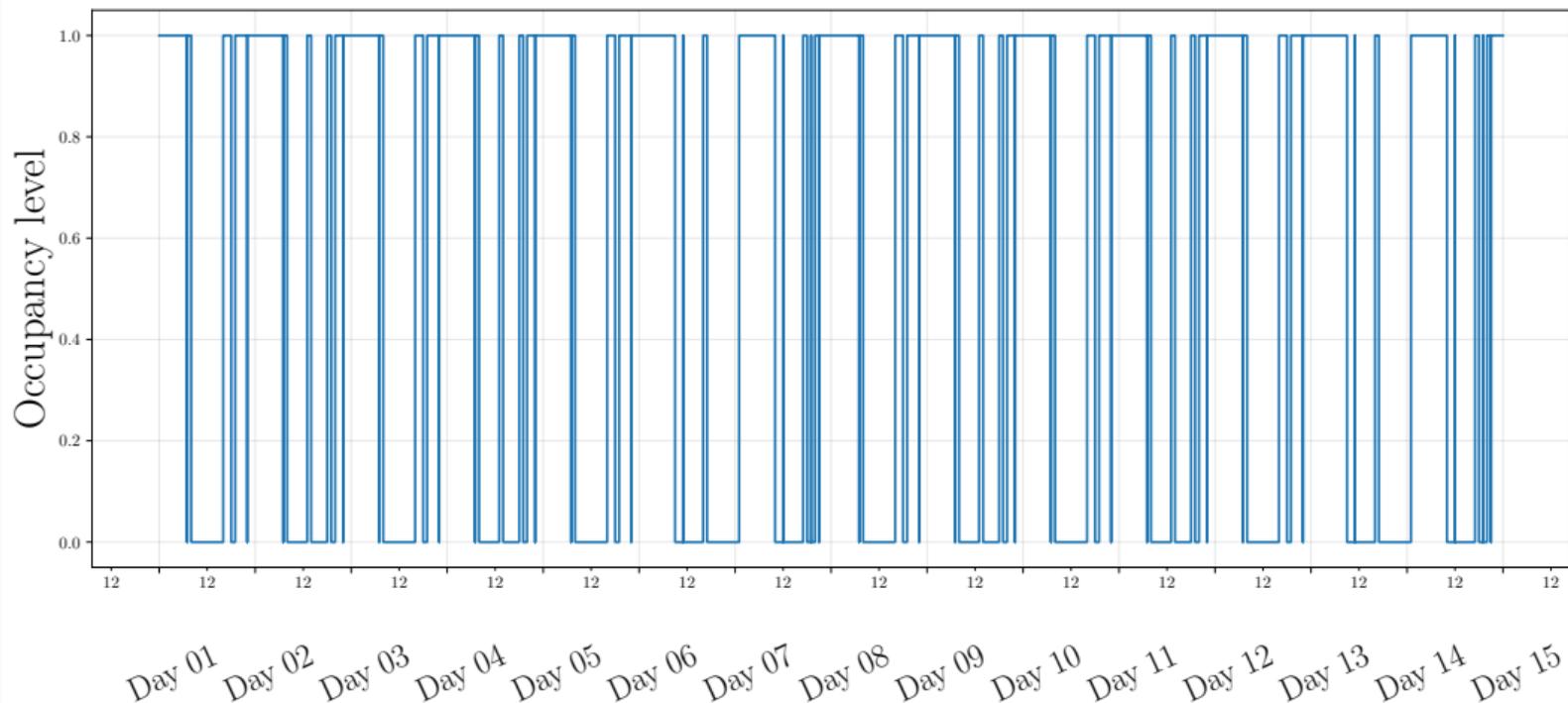
References

1. Campi, Marco C., Andrea Lecchini, and Sergio M. Savaresi. "Virtual reference feedback tuning: a direct method for the design of feedback controllers." *Automatica* 38.8 (2002): 1337-1346.
2. Hjalmarsson, Hakan, et al. "Iterative feedback tuning: theory and applications." *IEEE control systems magazine* 18.4 (1998): 26-41.
3. Karimi, A., L. Mikovi, and D. Bonvin. "Iterative correlationbased controller tuning." *International journal of adaptive control and signal processing* 18.8 (2004): 645-664.
4. Willems, Jan C., et al. "A note on persistency of excitation." *Systems & Control Letters* 54.4 (2005): 325-329.
5. De Persis, Claudio, and Pietro Tesi. "Formulas for data-driven control: Stabilization, optimality, and robustness." *IEEE Transactions on Automatic Control* 65.3 (2019): 909-924.
6. EQUA Simulation AB, IDA Indoor Climate and Energy (IDA-ICE), 2020, version 5.0. [Online]. Available: <https://www.equa.se/en/ida-ice>.
7. Russo, A., Proutiere, A. (2021). Poisoning attacks against data-driven control methods. American Control Conference, 2021.
8. Sinha, Ankur, Pekka Malo, and Kalyanmoy Deb. "A review on bilevel optimization: from classical to evolutionary approaches and applications." *IEEE Transactions on Evolutionary Computation* 22.2 (2017): 276-295.

References

9. Shen, Xinyue, et al. "Disciplined convex-concave programming." 2016 IEEE 55th Conference on Decision and Control (CDC). IEEE, 2016.
10. Biggio, Battista, Blaine Nelson, and Pavel Laskov. "Poisoning attacks against support vector machines." arXiv preprint arXiv:1206.6389 (2012).

Backup



Attack Formulation

$$\begin{aligned} \max_{\mathbf{u}', \mathbf{y}'} \quad & \mathcal{A}(\mathbf{u}, \mathbf{y}, K(\mathbf{u}', \mathbf{y}')) \\ \text{s.t.} \quad & K(\mathbf{u}', \mathbf{y}') \in \arg \min_K \mathcal{L}(\mathbf{u}', \mathbf{y}', K) \\ & \|\mathbf{u}' - \mathbf{u}\|_{q_u} \leq \delta_u, \quad \|\mathbf{y}' - \mathbf{y}\|_{q_y} \leq \delta_y, \end{aligned}$$

1. Assume the inner problem $K(\mathbf{u}', \mathbf{y}') \in \arg \min_K \mathcal{L}(\mathbf{u}', \mathbf{y}', K)$ is convex and sufficiently regular.
 - We can perform single-level reduction [6] and replace the inner problem with its KKT conditions.

2. Then, assume K is parameterized by θ (we will write K_θ). We can conclude that

$$\nabla_{\theta} \mathcal{L}(\mathbf{u}', \mathbf{a}', K_\theta) = 0 \Rightarrow \nabla_{\mathbf{a}_u} \theta = -(\nabla_{\mathbf{a}_u} \nabla_{\theta} \mathcal{L})(\nabla_{\theta}^2 \mathcal{L})^{-1}$$

(similarly also for \mathbf{a}_y).

3. **This allows us to find approximate attacks** by using gradient ascent methods.

Attack Formulation

$$\begin{aligned} \max_{\mathbf{u}', \mathbf{y}'} \quad & \mathcal{A}(\mathbf{u}, \mathbf{y}, K(\mathbf{u}', \mathbf{y}')) \\ \text{s.t.} \quad & K(\mathbf{u}', \mathbf{y}') \in \arg \min_K \mathcal{L}(\mathbf{u}', \mathbf{y}', K) \\ & \|\mathbf{u}' - \mathbf{u}\|_{q_u} \leq \delta_u, \quad \|\mathbf{y}' - \mathbf{y}\|_{q_y} \leq \delta_y, \end{aligned}$$

1. Assume the inner problem $K(\mathbf{u}', \mathbf{y}') \in \arg \min_K \mathcal{L}(\mathbf{u}', \mathbf{y}', K)$ is convex and sufficiently regular.
 - We can perform single-level reduction [6] and replace the inner problem with its KKT conditions.

2. Then, assume K is parameterized by θ (we will write K_θ). We can conclude that

$$\nabla_{\theta} \mathcal{L}(\mathbf{u}', \mathbf{a}', K_\theta) = 0 \Rightarrow \nabla_{\mathbf{a}_u} \theta = -(\nabla_{\mathbf{a}_u} \nabla_{\theta} \mathcal{L})(\nabla_{\theta}^2 \mathcal{L})^{-1}$$

(similarly also for \mathbf{a}_y).

3. **This allows us to find approximate attacks** by using gradient ascent methods.