

Robust application-oriented exploration in LQ control

Jack Umenberger
Thomas B. Schön

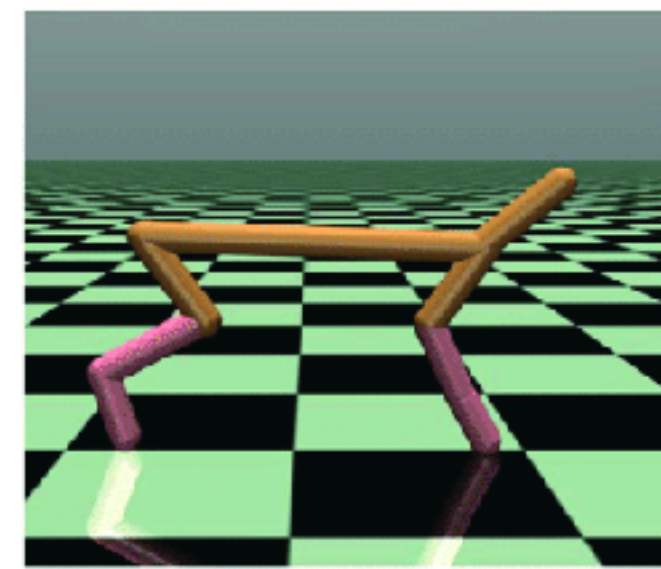
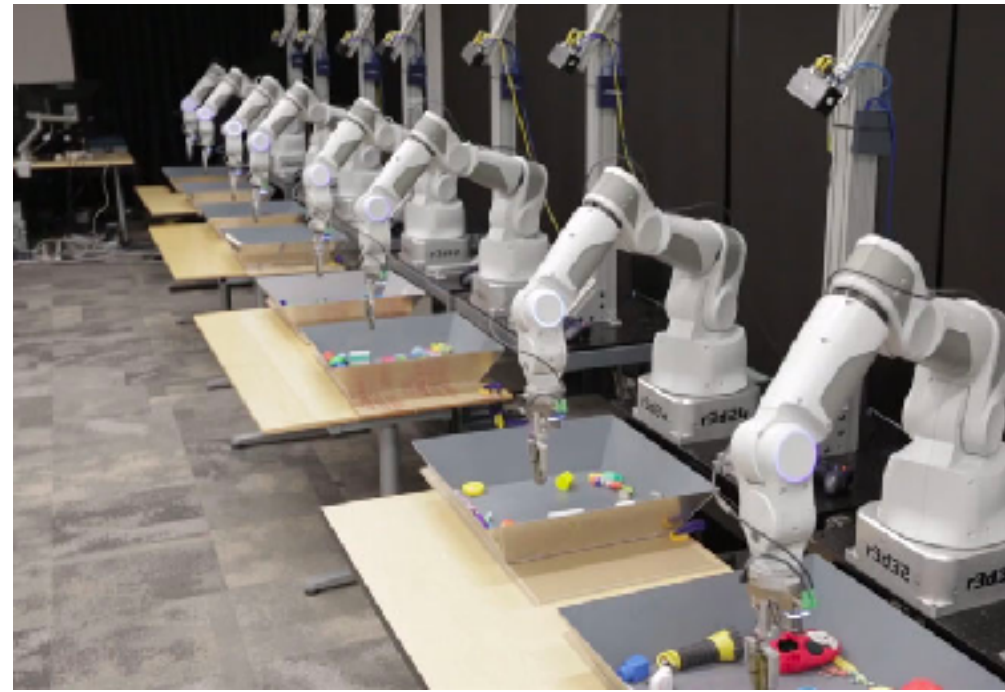


Mina Ferizbegovic
Håkan Hjalmarsson

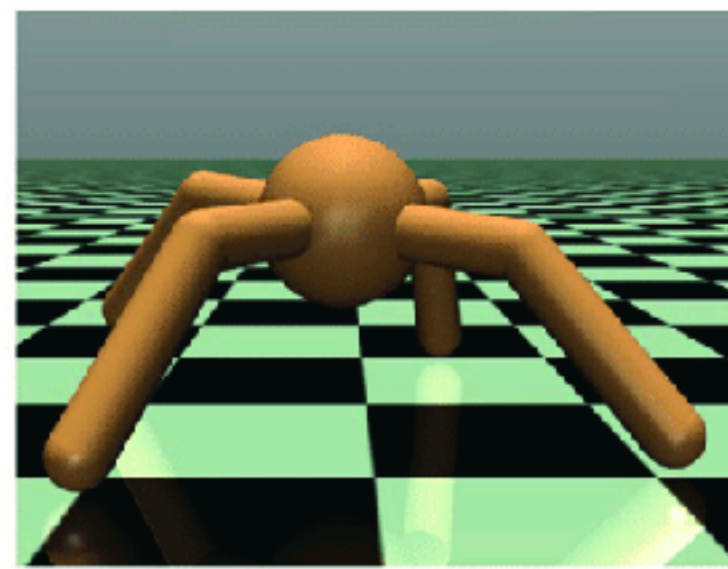




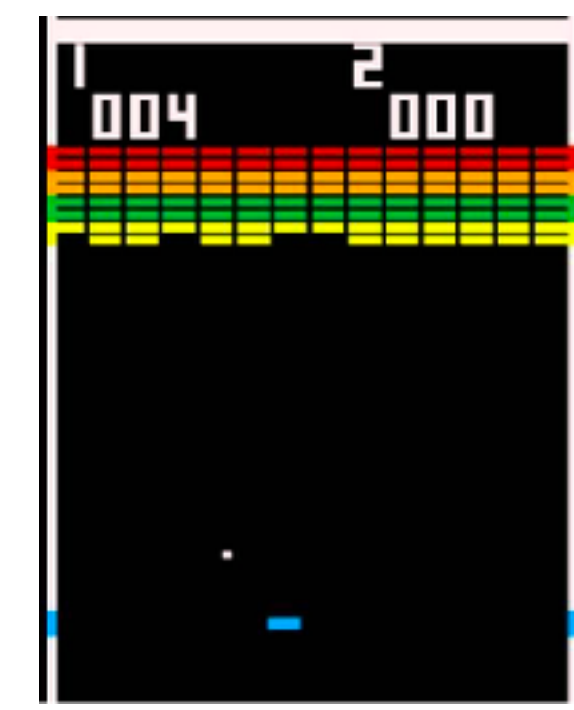
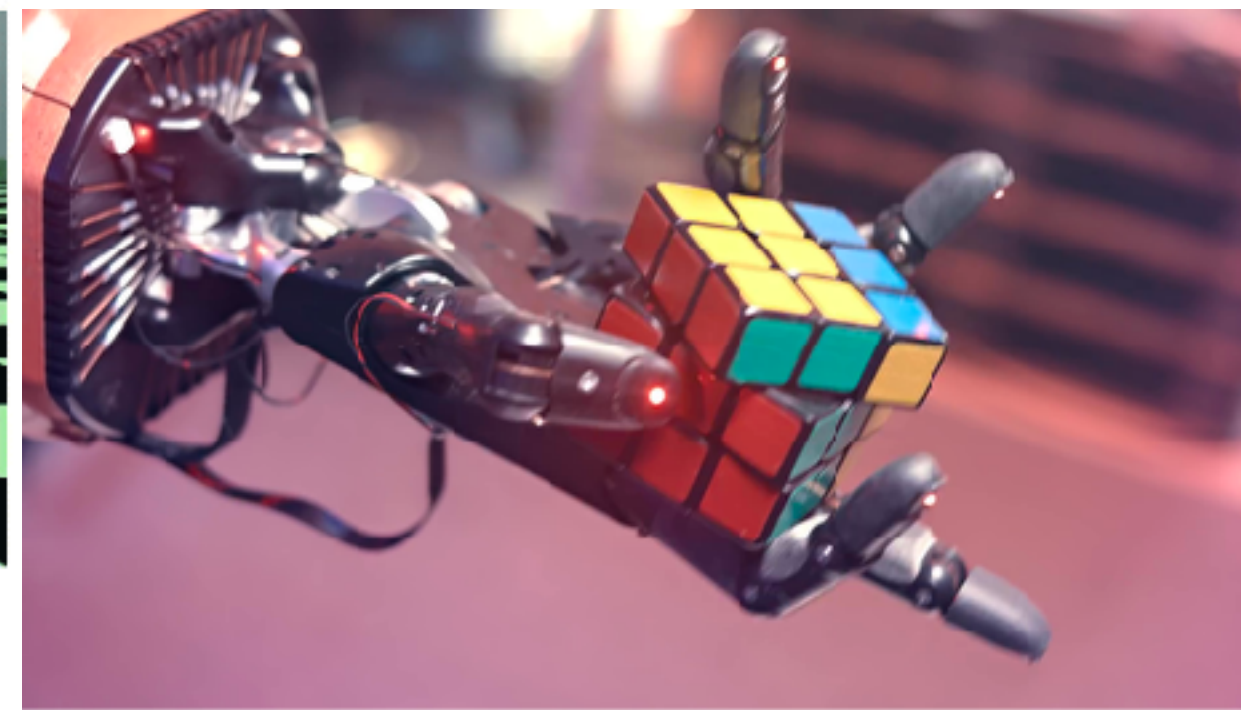
Goals



Half-Cheetah



Ant



- In reinforcement learning, agents need to explore their environment to learn which actions maximize rewards
- exploration is *often* random trial and error
- we want to do exploration that is **robust** and **targeted**
- **robust**: does not destabilize the system or cause failure
- **targeted**: provides knowledge that helps complete the task



Balancing the tradeoff...



- There typically exists a **tradeoff** between exploration and exploitation
- actions that provide the most information about the environment may incur high short term costs



Balancing the tradeoff...



- There typically exists a **tradeoff** between exploration and exploitation
 - actions that provide the most information about the environment may incur high short term costs
- We achieve the **optimal tradeoff** between exploration and exploitation by formulating the search for a policy as a **convex optimization** problem, that can be solved to **global optimality**



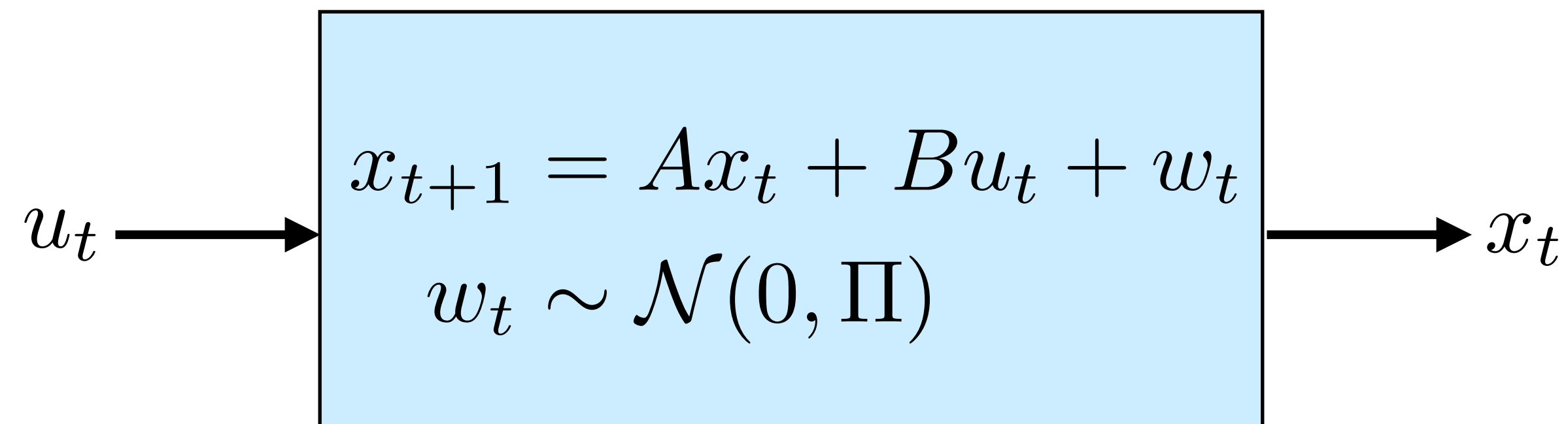
Balancing the tradeoff...



- There typically exists a **tradeoff** between exploration and exploitation
 - actions that provide the most information about the environment may incur high short term costs
- We achieve the **optimal tradeoff** between exploration and exploitation by formulating the search for a policy as a **convex optimization** problem, that can be solved to **global optimality**
- Relies only strong assumptions about the environment...



Linear quadratic control

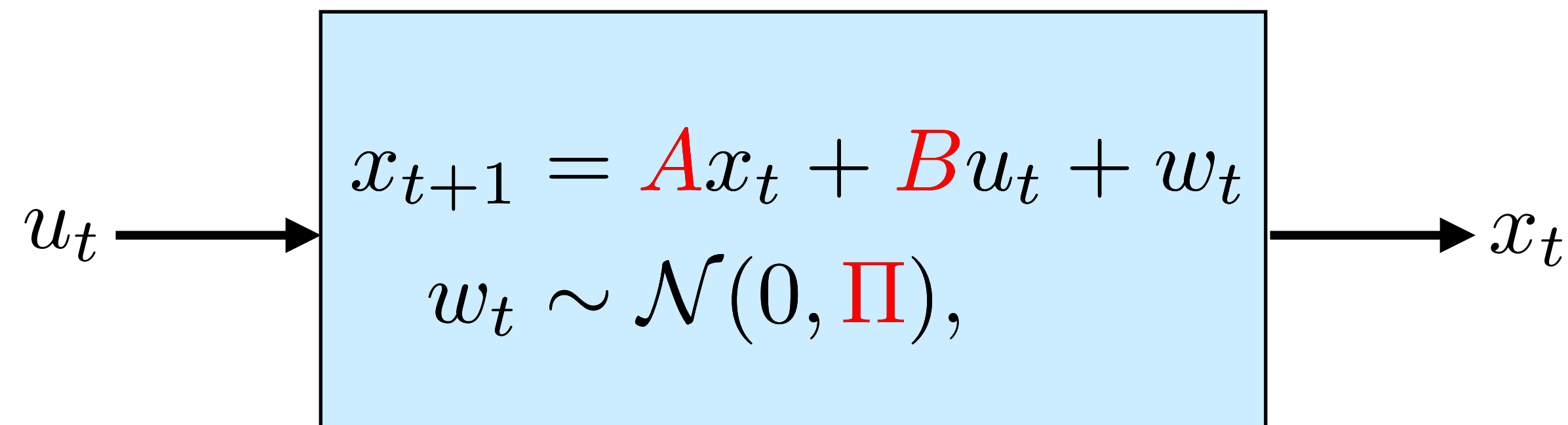


Task: find a state-feedback controller to minimize the quadratic cost

$$\sum_{t=0}^T x_t^\top Q x_t + u_t^\top R u_t$$



Linear quadratic control



Task: find a state-feedback controller to minimize the quadratic cost

$$\sum_{t=0}^T x_t^\top Q x_t + u_t^\top R u_t$$

Challenge: we don't know the system parameters A, B, Π



Key quantity: empirical covariance



- The key quantity in our formulation is the empirical covariance $D = \sum_{t=0}^T \begin{bmatrix} x_t \\ u_t \end{bmatrix} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^\top$
- Shows up in both



Key quantity: empirical covariance



- The key quantity in our formulation is the empirical covariance $D = \sum_{t=0}^T \begin{bmatrix} x_t \\ u_t \end{bmatrix} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^\top$
- Shows up in both
 - the cost, trace $\begin{bmatrix} Q & 0 \\ 0 & R \end{bmatrix} D$, where “**smaller**” is **better** for lower cost



Key quantity: empirical covariance



- The key quantity in our formulation is the empirical covariance $D = \sum_{t=0}^T \begin{bmatrix} x_t \\ u_t \end{bmatrix} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^\top$
- Shows up in both
 - the cost, trace $\begin{bmatrix} Q & 0 \\ 0 & R \end{bmatrix} D$, where “**smaller**” is **better** for lower cost
 - the system uncertainty (i.e. inverse variance of the posterior), $D \otimes I$, where “**bigger**” is **better** for reduced uncertainty
- This clearly illustrates the exploration (“big”) and exploitation (“small”) tradeoff



Preserving structure in uncertainty



- recent works consider **absolute** measures of uncertainty

$$\|A_{\text{est}} - A_{\text{true}}\| \leq \epsilon_A, \quad \|B_{\text{est}} - B_{\text{true}}\| \leq \epsilon_B$$

- all uncertainty information is collapsed into a single **scalar**; structure is lost



Preserving structure in uncertainty



- recent works consider **absolute** measures of uncertainty

$$\|A_{\text{est}} - A_{\text{true}}\| \leq \epsilon_A, \quad \|B_{\text{est}} - B_{\text{true}}\| \leq \epsilon_B$$

- all uncertainty information is collapsed into a single **scalar**; structure is lost

- we preserve structure by working with the **matrix** $D = \sum_{t=0}^T \begin{bmatrix} x_t \\ u_t \end{bmatrix} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^\top$

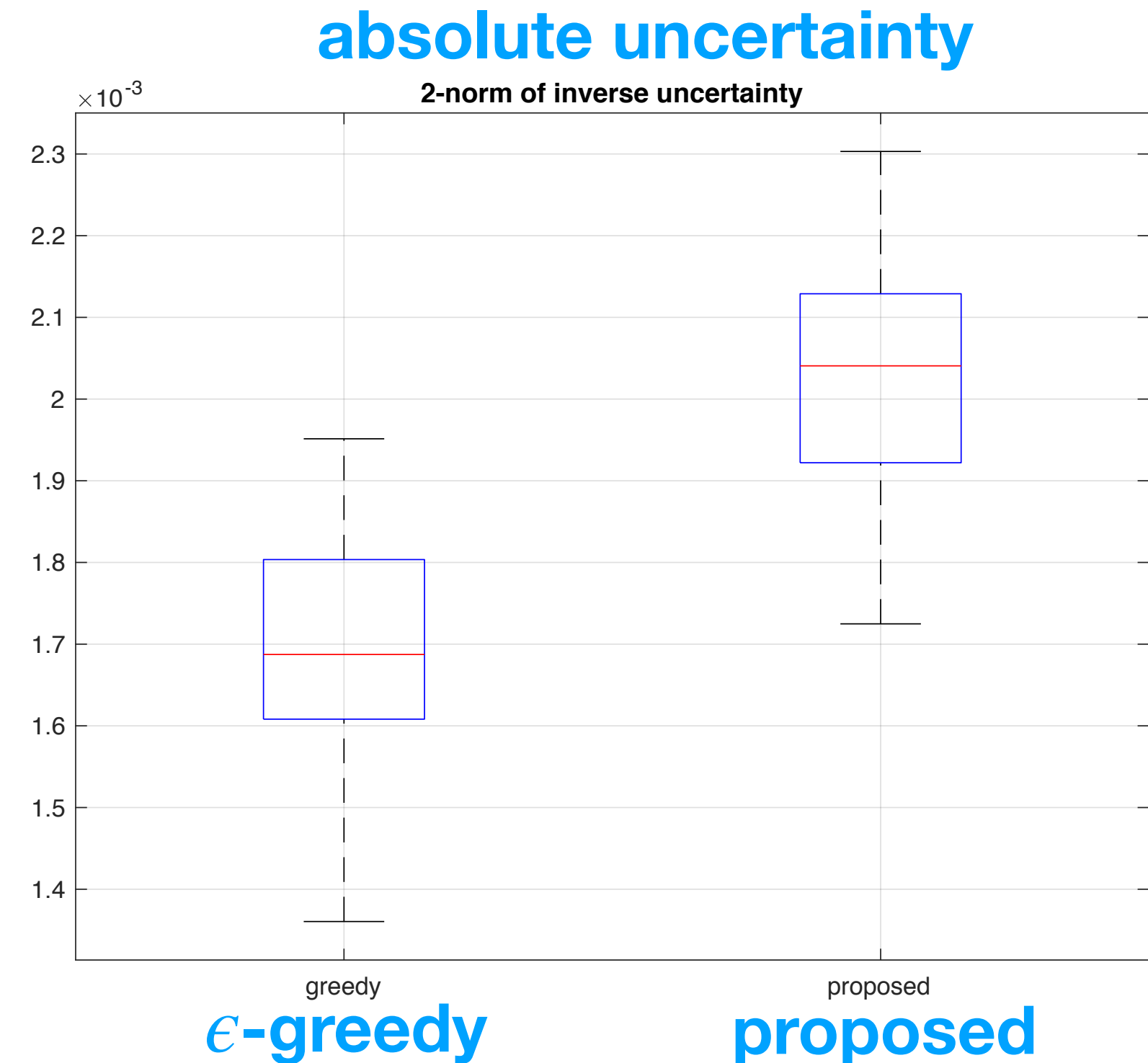
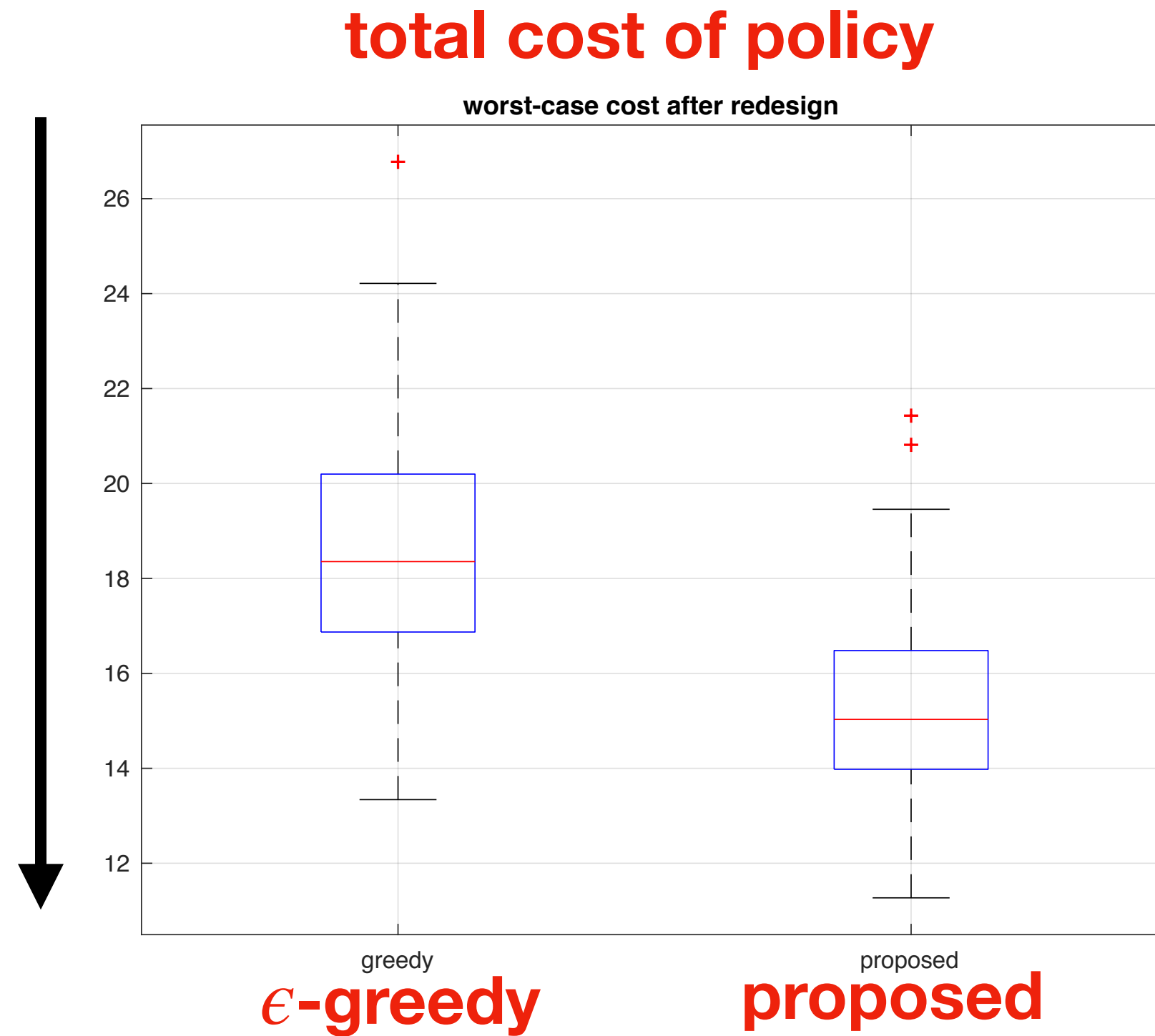
- allows to optimize for reduction of uncertainty in the parameters that **matter for the task**



less uncertainty \neq lower cost



lower cost
(better
performance)



lower absolute
uncertainty

- **performance** of proposed method is better, even though absolute **uncertainty** is larger
- the **structure** of the uncertainty matters
- uncertainty has been reduced in the parameters that matter for the control task



Poster presentation



Poster #177

Today 05:00 -- 07:30 PM @ East Hall B+C



SWEDISH FOUNDATION *for*
STRATEGIC RESEARCH