

**Department of Advanced Control**

**COMBINING MODEL-BASED  
AND MODEL-FREE OPTIMAL  
CONTROL FOR DYNAMIC  
SYSTEMS**

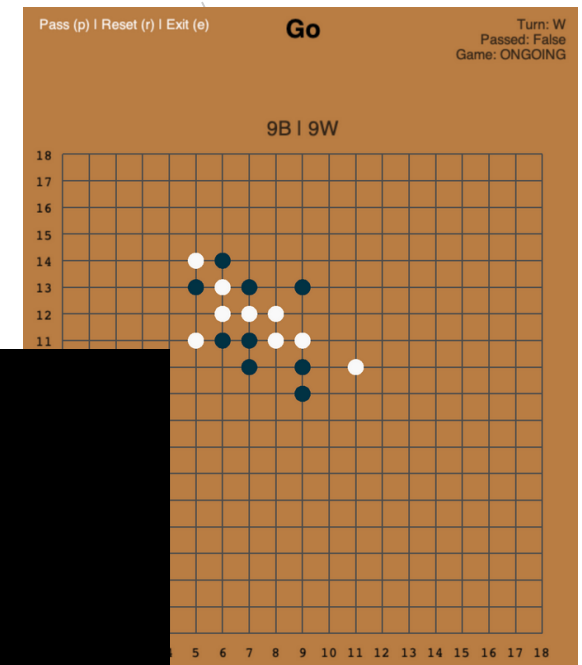
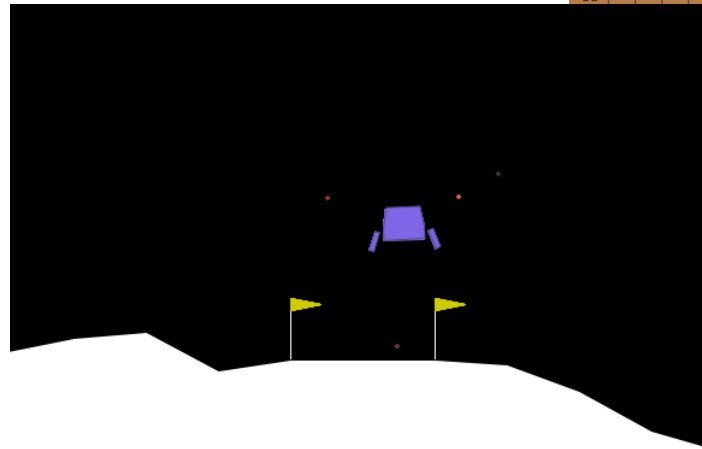
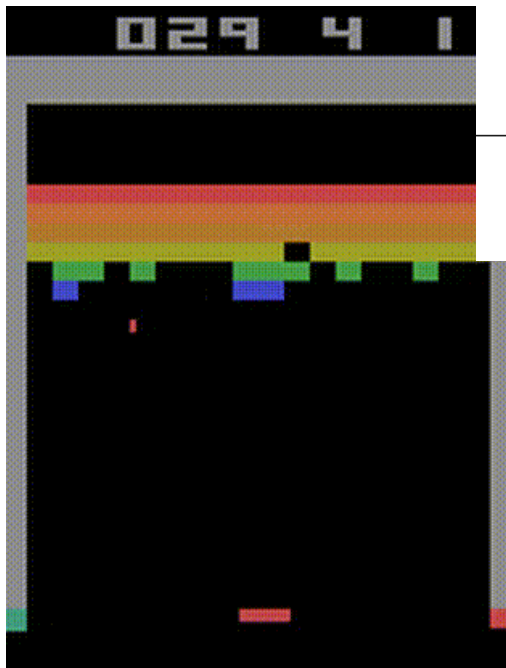
**MASTER'S THESIS PRESENTATION - HANDOUT**

Pascal Peters

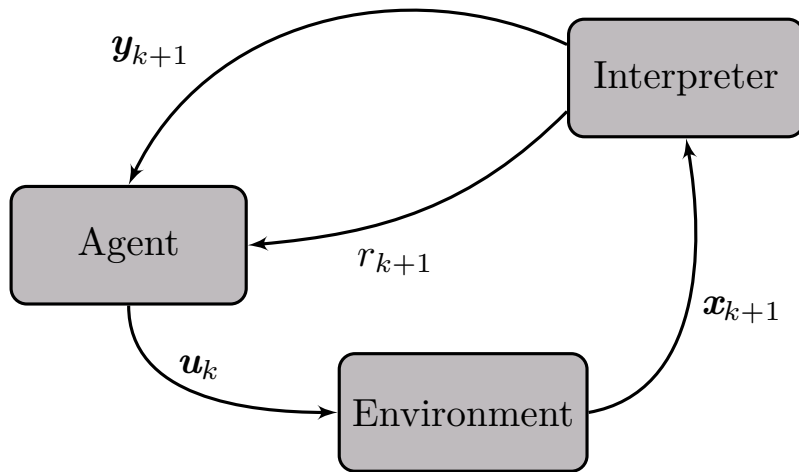
Paderborn, 04.05.2022



## Motivation



# Reinforcement Learning



Combining model-based and model-free control

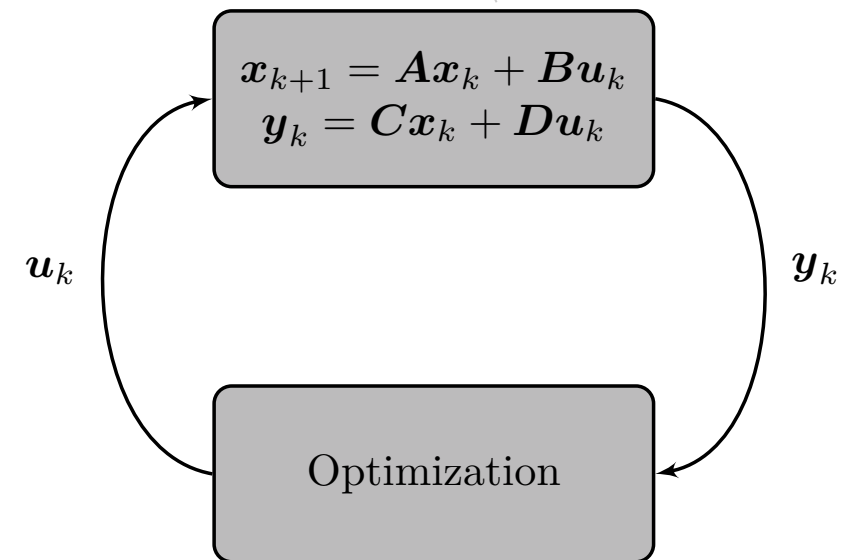
- Model-free
- Maximize the return:  $g_k = \sum_{i=0}^N \gamma^i r_{k+i+1}$
- Solves Bellman equation for a horizon  $N \rightarrow \infty$ :

$$q_{\pi}(\mathbf{x}_k, \mathbf{u}_k) = \mathbb{E}_{\pi} \left[ R_{k+1} + \gamma q_{\pi}(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) \middle| \mathbf{x}_k, \mathbf{u}_k \right]$$

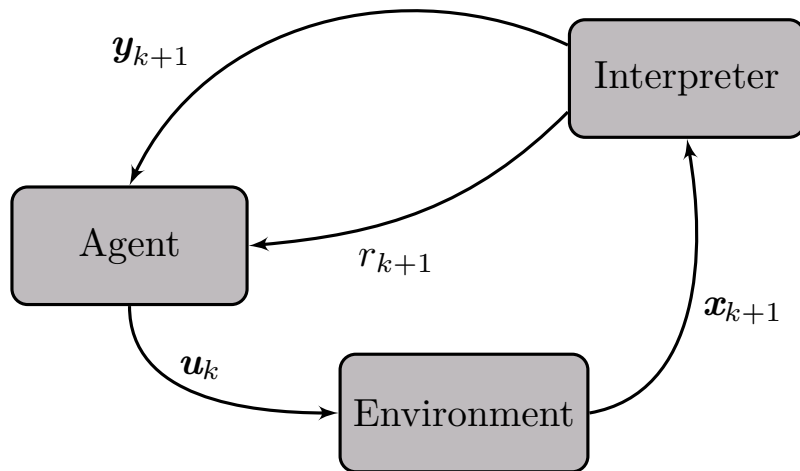
- Explicit control law to determine action  $\mathbf{u}_k$
- Learn through experience and feedback

## Model Predictive Control

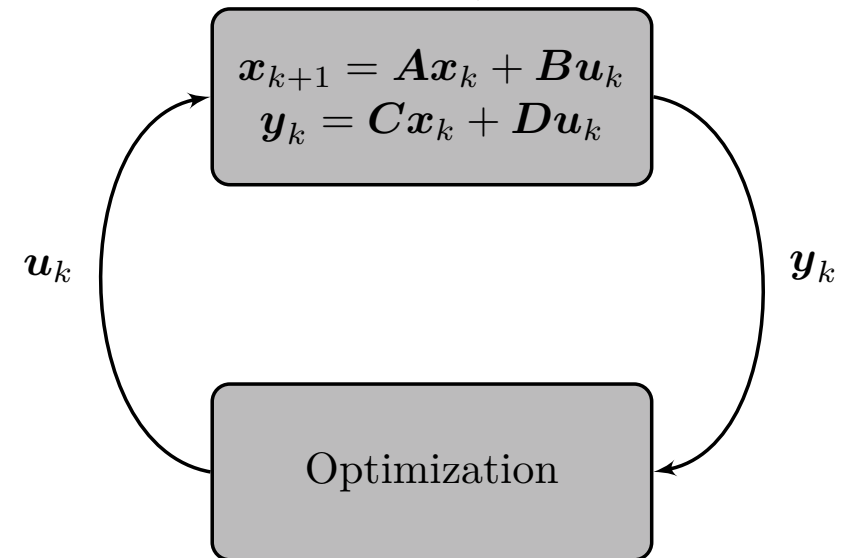
- Model-based
- Minimize costs
- Solve optimization problem for finite  $N$ -step prediction horizon
- Implicit control law with complexity growing with horizon length  $N$
- Constrained optimization for safe control



## Motivation



+

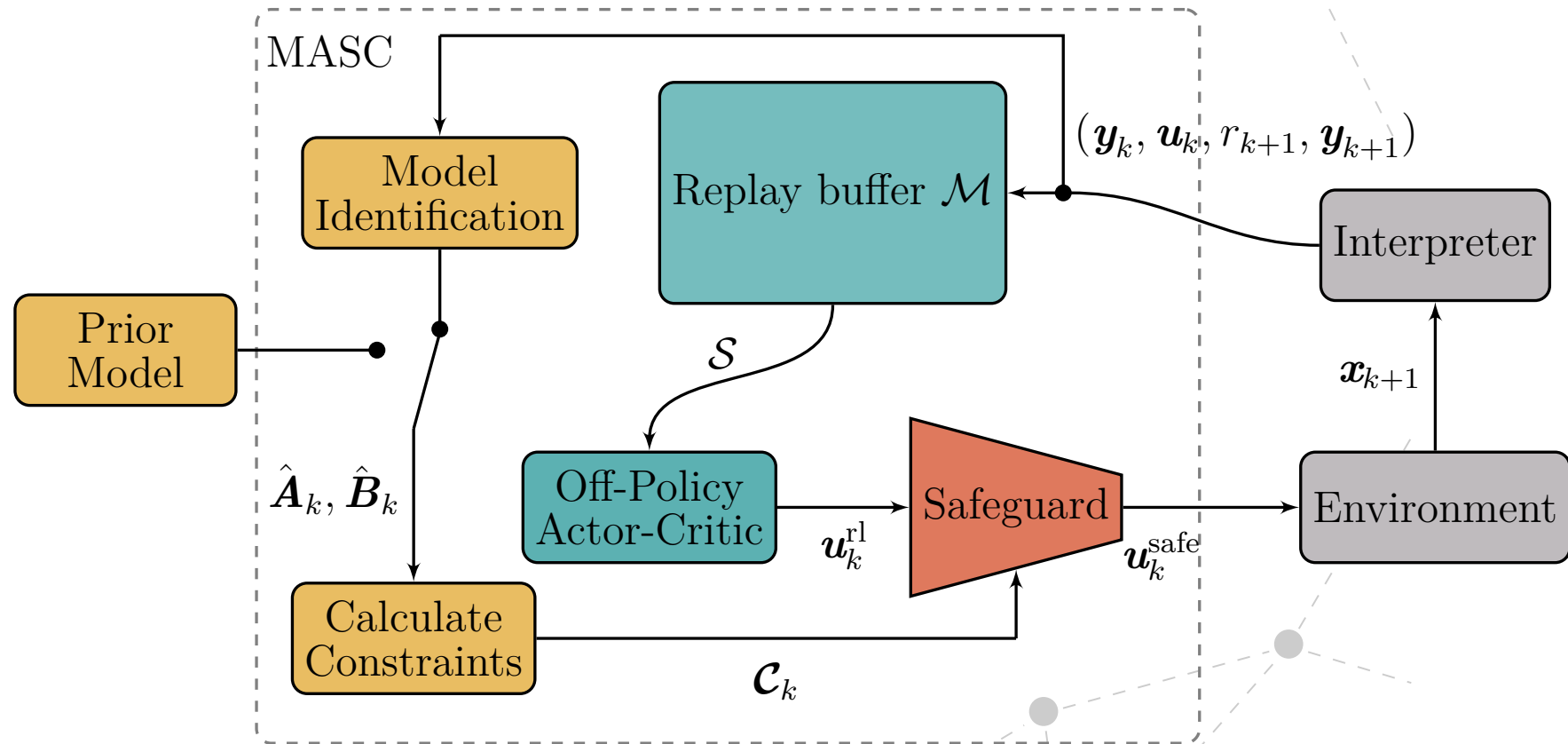


Combining model-based and model-free control

## Developed System

- Using a model to extend an off-policy actor-critic controller to achieve safe behavior
  - By monitoring actions for constraint violation
  - Modifying unsafe actions such that the system stays within the constraints
- Eliminate the dependency on prior available model knowledge by identifying model parameters online

## Model Adaptive Safeguard Controller



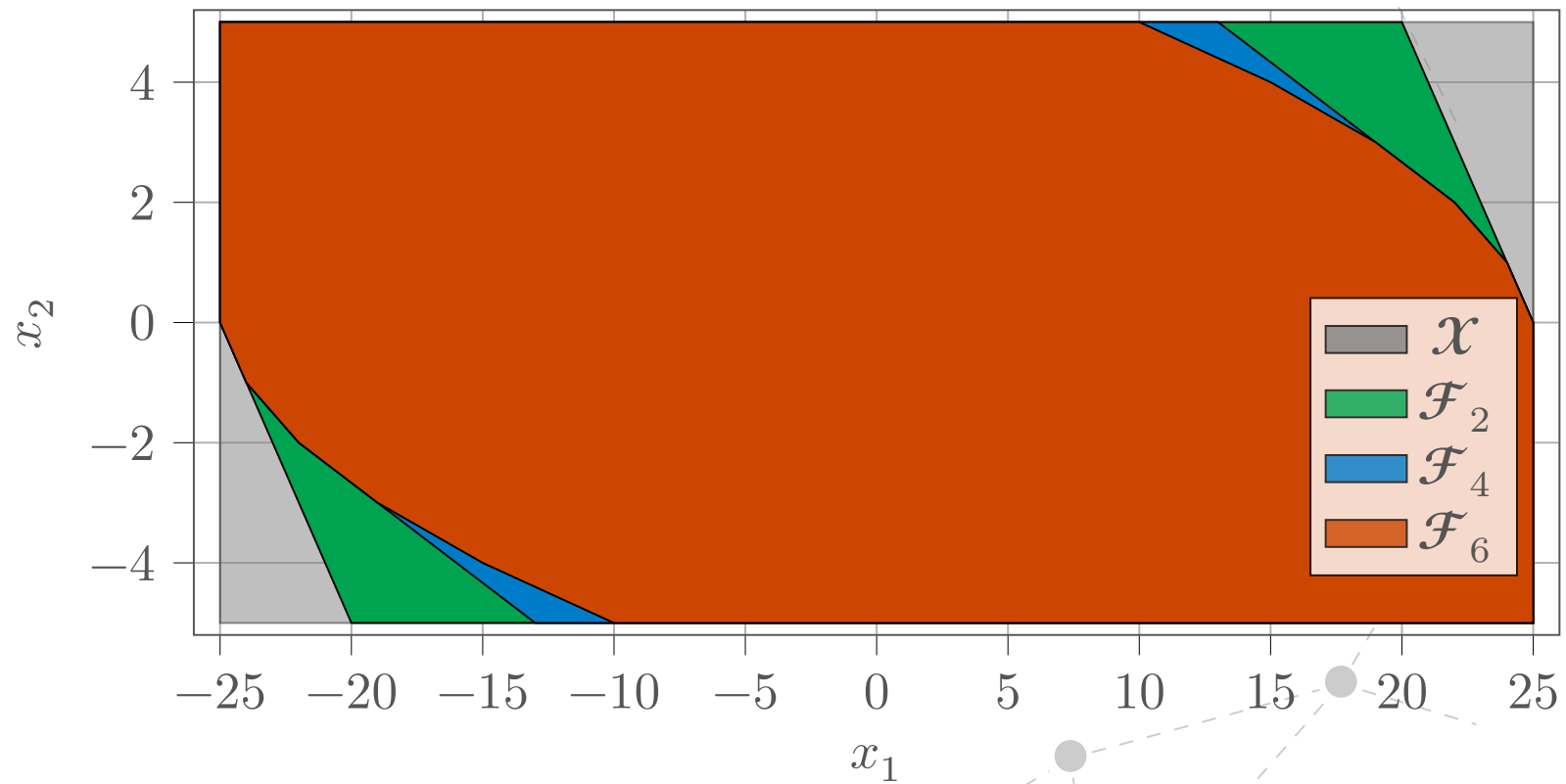
## Safeguard

- Check for possible constraint violation
  - $\mathbf{u}_k^{\text{RL}}$  safe if:  $\begin{bmatrix} \mathbf{x}_k^{\text{T}} & \mathbf{u}_k^{\text{RL T}} \end{bmatrix}^{\text{T}} \in \mathcal{C}$
  - $\mathbf{u}_k^{\text{RL}}$  unsafe if:  $\begin{bmatrix} \mathbf{x}_k^{\text{T}} & \mathbf{u}_k^{\text{RL T}} \end{bmatrix}^{\text{T}} \notin \mathcal{C}$
- Solving optimization problem for unsafe actions:

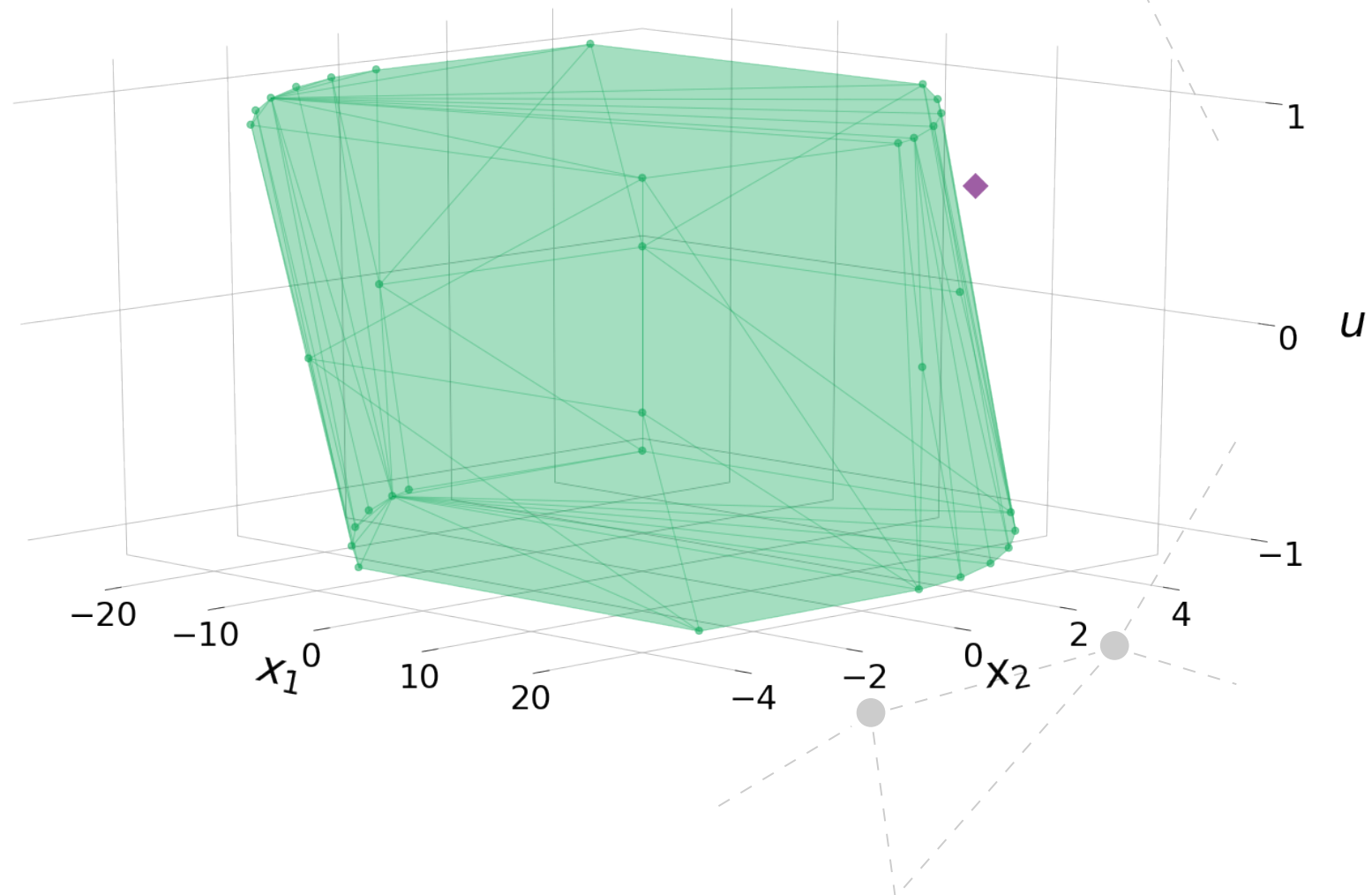
$$\mathbf{u}_k^{\text{safe}} = \mathbf{u}_k^* = \underset{\mathbf{u}_k}{\operatorname{argmin}} \|\mathbf{u}_k - \mathbf{u}_k^{\text{RL}}\|^2,$$
$$\text{s.t. } \mathbf{u}_k \in \mathcal{C}_u$$



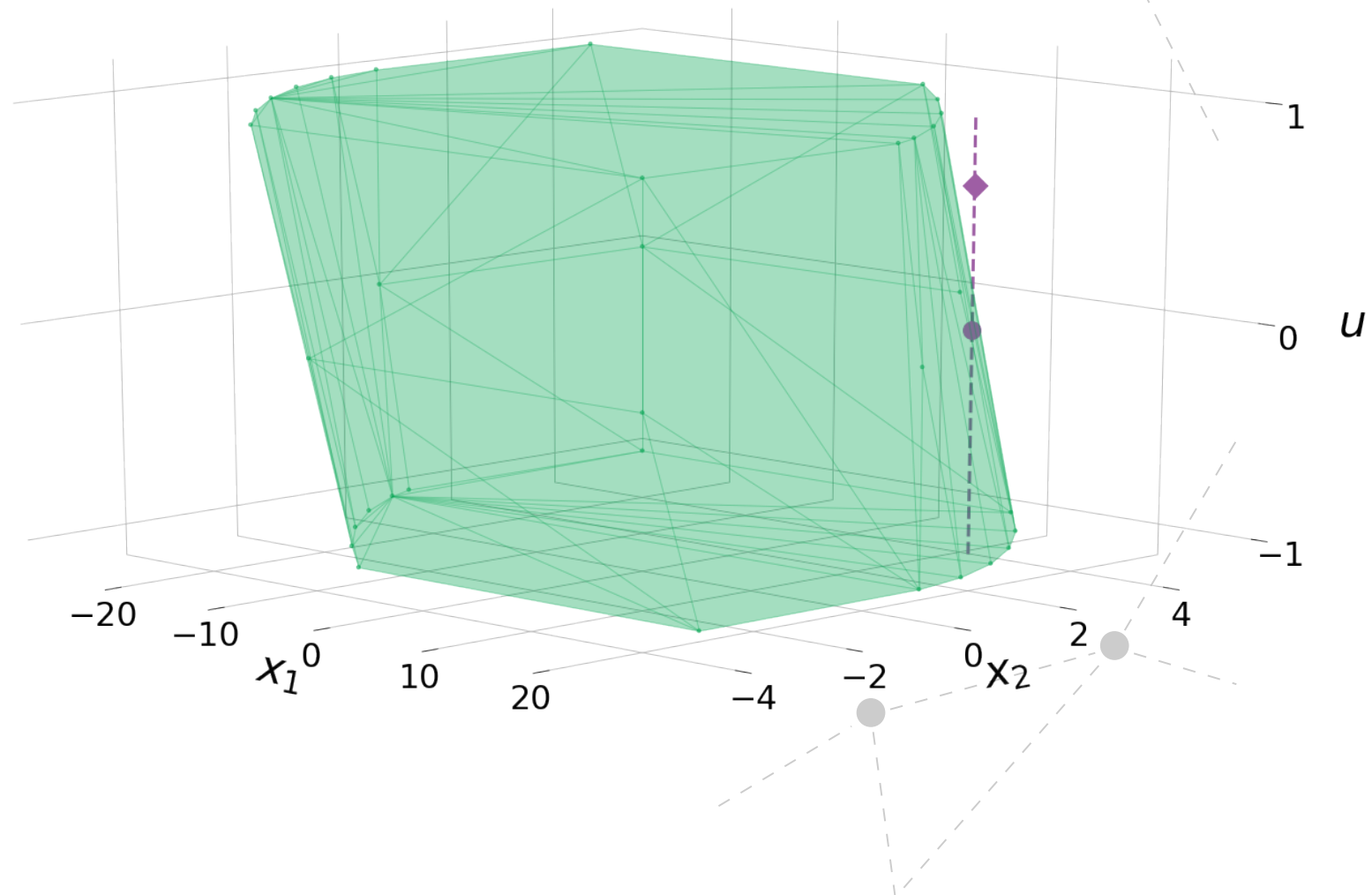
## Safeguard Constraints



## Safeguard Constraints



## Safeguard Constraints



## Model Identification

- Online model identification with the recursive least squares algorithm

$$\begin{aligned}\kappa_k &= \frac{\mathbf{P}_k \xi_{k+1}}{\lambda_{k+1} + \xi_{k+1}^\top \mathbf{P}_k \xi_{k+1}} \\ \hat{\theta}_{k+1} &= \hat{\theta}_k + \kappa_k \left( \psi_{k+1} - \xi_{k+1}^\top \hat{\theta}_k \right) \\ \mathbf{P}_{k+1} &= \left( \mathbf{I} - \kappa_k \xi_{k+1}^\top \right) \mathbf{P}_k \frac{1}{\lambda_{k+1}}\end{aligned}$$

- Update feasible set only if change of new estimates is significant:

$$\Delta \mathbf{A}_k = \frac{\|\hat{\mathbf{A}}_k - \hat{\mathbf{A}}_{k-1}\|_F}{\|\hat{\mathbf{A}}_k\|_F}, \Delta \mathbf{B}_k = \frac{\|\hat{\mathbf{B}}_k - \hat{\mathbf{B}}_{k-1}\|_F}{\|\hat{\mathbf{B}}_k\|_F}$$

## Evaluation

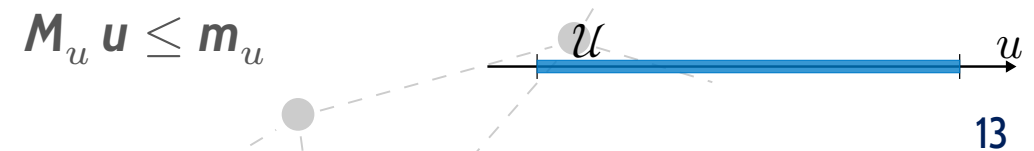
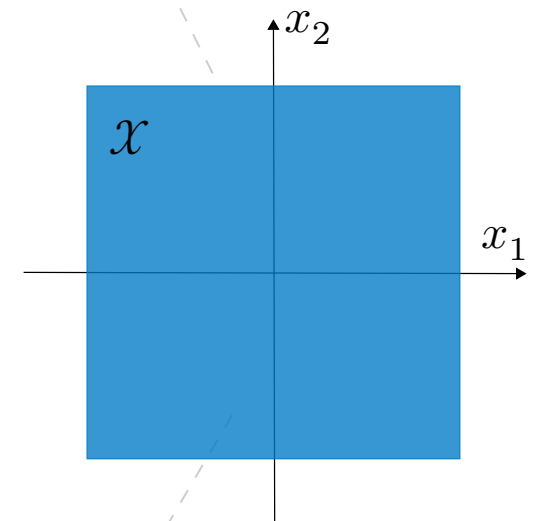
- Double Integrator as test environment

$$\dot{\mathbf{x}}(t) = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{x}(t) + \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mathbf{u}(t),$$

$$\mathbf{y}(t) = \begin{bmatrix} 1 & 0 \end{bmatrix} \mathbf{x}(t), \quad \mathbf{x}(0) = \mathbf{x}_0$$

$$\text{s.t. } \mathbf{M}_x \mathbf{x} \leq \mathbf{m}_x,$$

$$\mathbf{M}_u \mathbf{u} \leq \mathbf{m}_u$$



## Evaluation

- Scaled weighted sum of errors (SWSE):

$$r_{k+1}^{\text{WSE}} = \begin{cases} -\frac{1}{f} \left( \sum_n w_n \left| \frac{x_{k,n} - x_{k,n}^{\text{ref}}}{2x_n^{\text{lim}}} \right| \right) & , \mathbf{x}_{k+1} \in \mathcal{X} \\ r^{\text{violation}} & , \mathbf{x}_{k+1} \notin \mathcal{X} \end{cases}$$

- Incorporate safeguard-penalty  $r^{\text{SG}}$  for actuated safeguard

## Metrics

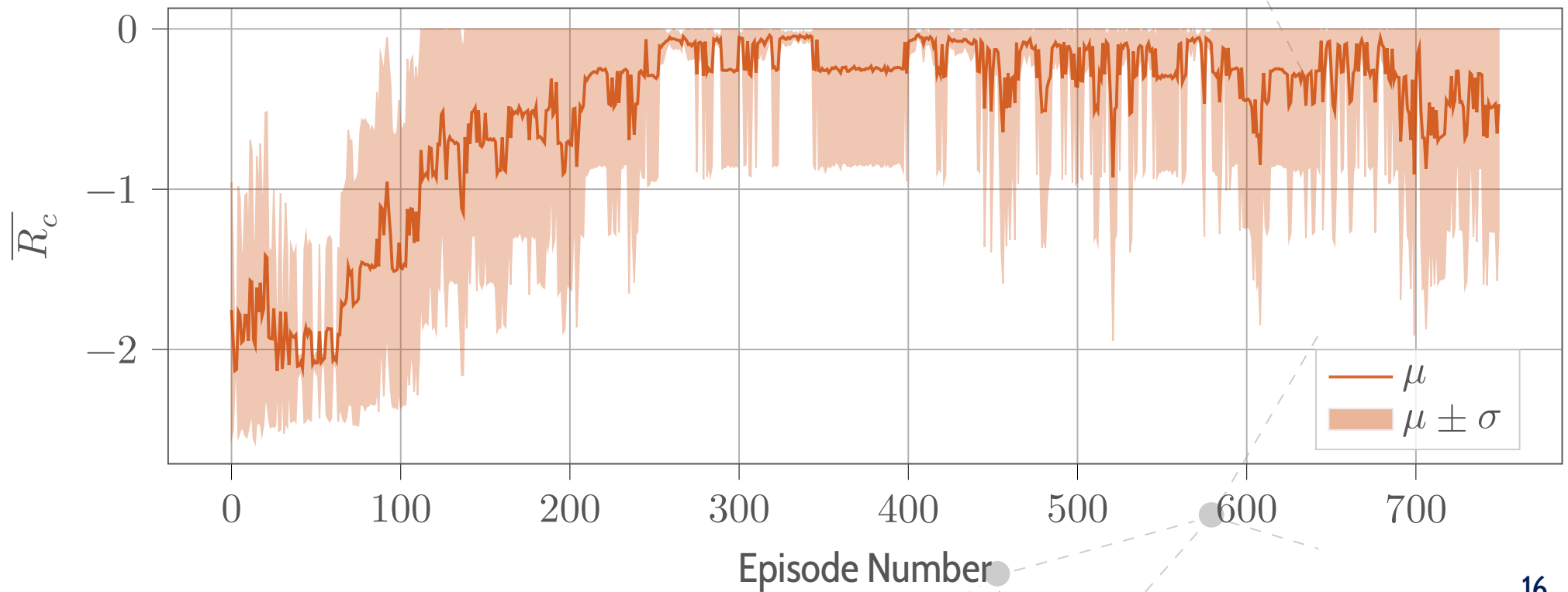
- Averaged over multiple seeds
- Cumulated reward per episode:

$$R_c = \sum_{k=1}^N r_k$$

- Cumulated constraint violations (CCV):

$$\text{CCV} = \sum_{e=1}^E c_e(\mathbf{x}), \quad \text{with } c_e(\mathbf{x}) = \begin{cases} 0 & , \text{if } \mathbf{x} \in \mathcal{X} \\ 1 & , \text{otherwise} \end{cases}$$

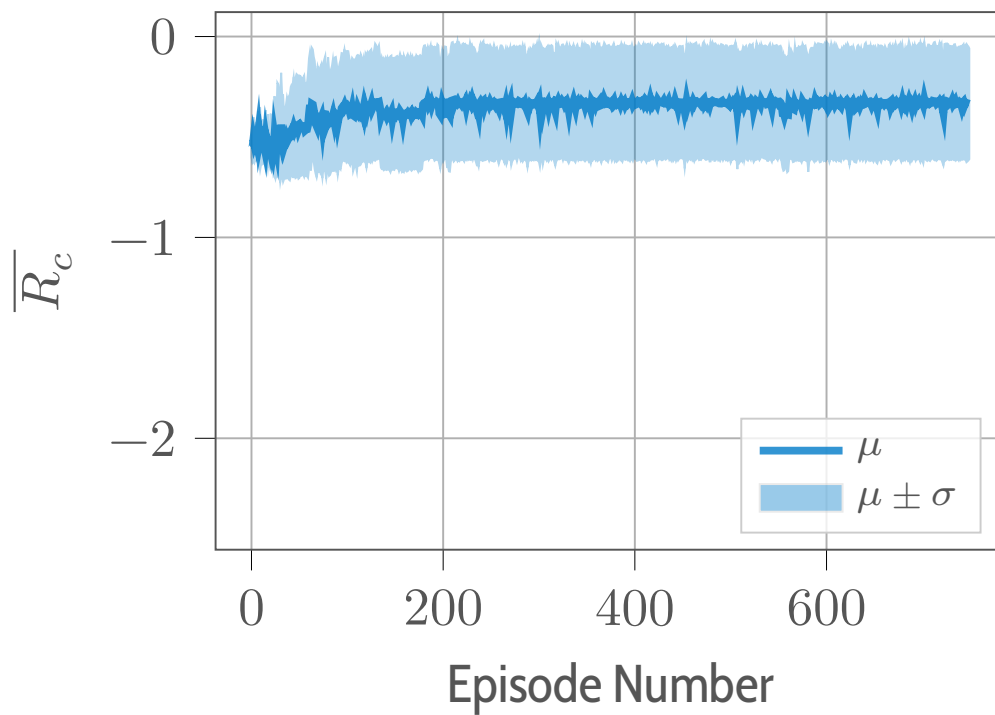
## Reinforcement Learning Controller



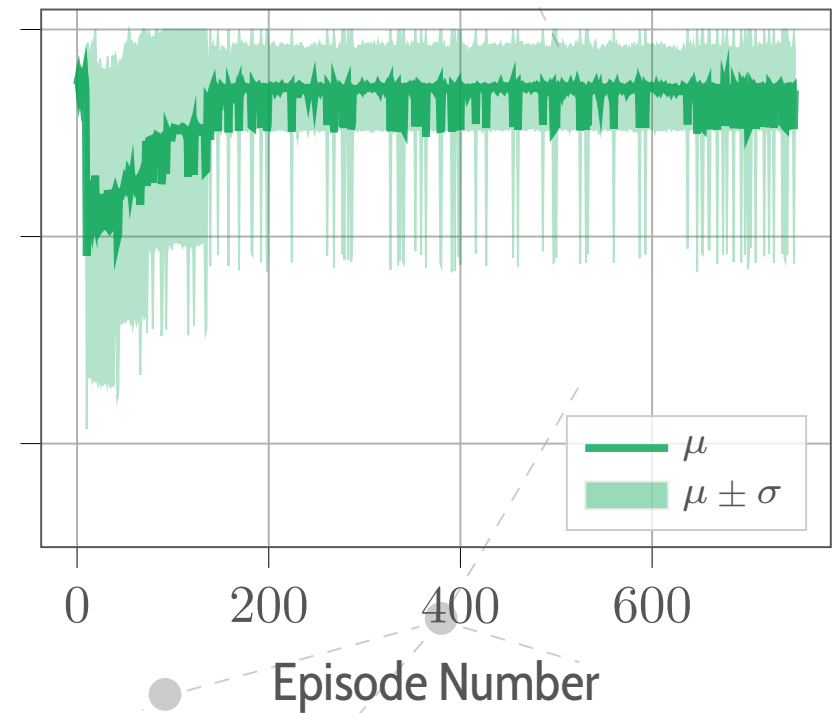


## Safeguard Controller

○ Prior available model:

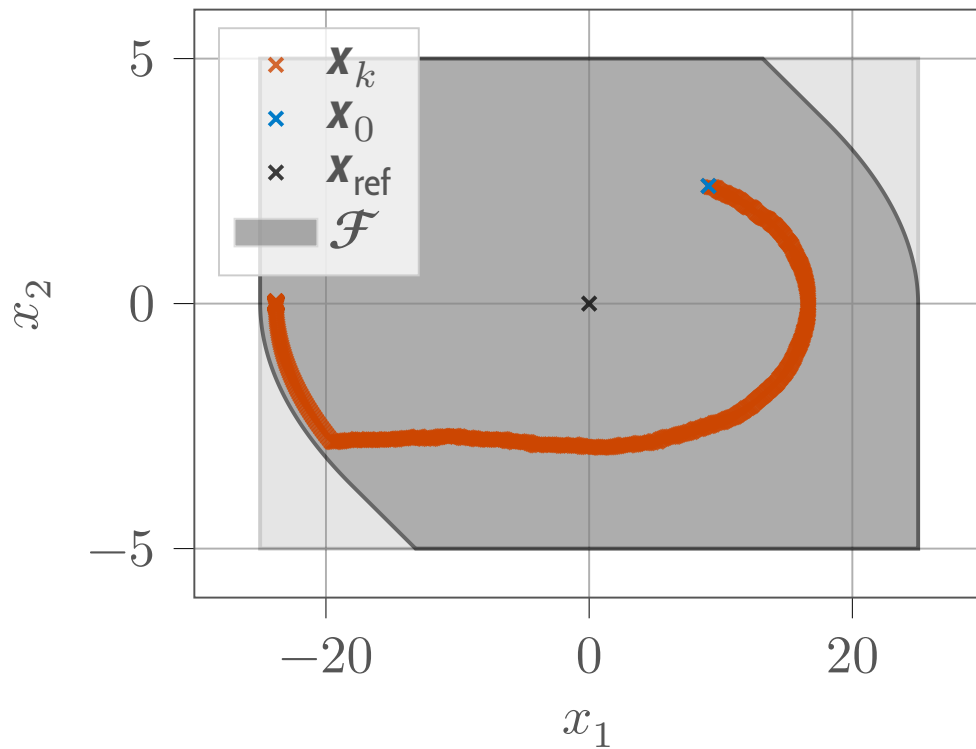


○ Estimated model:

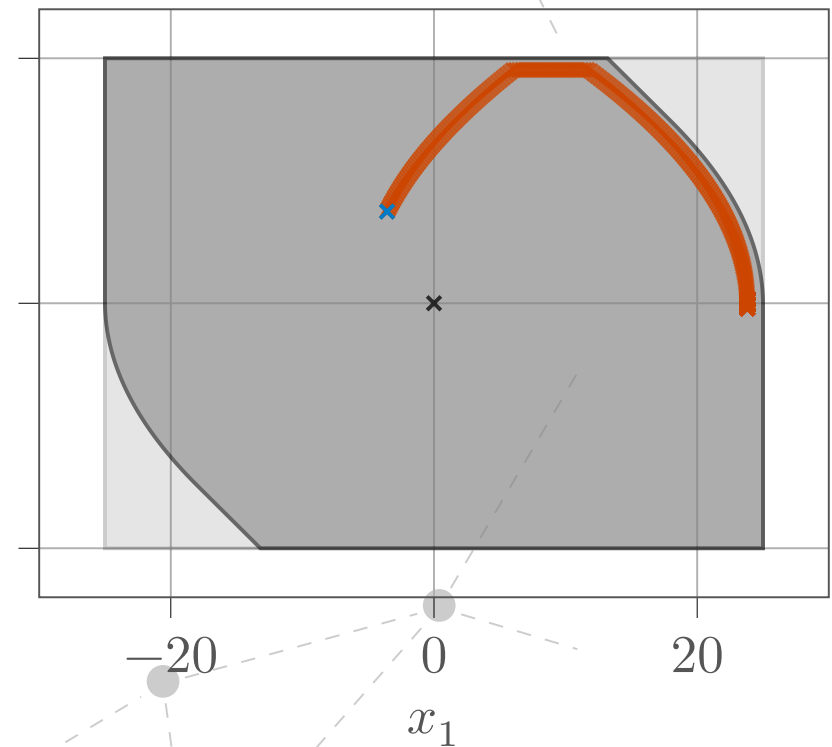


## Exemplary Safeguard Controller Trajectories

○ Training start:

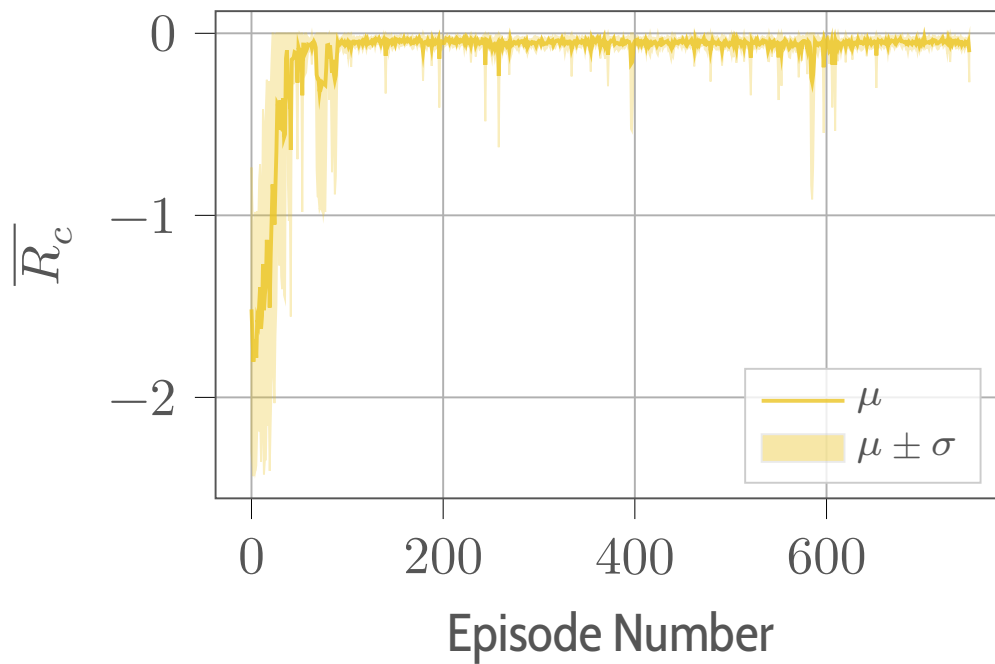


○ Training end:

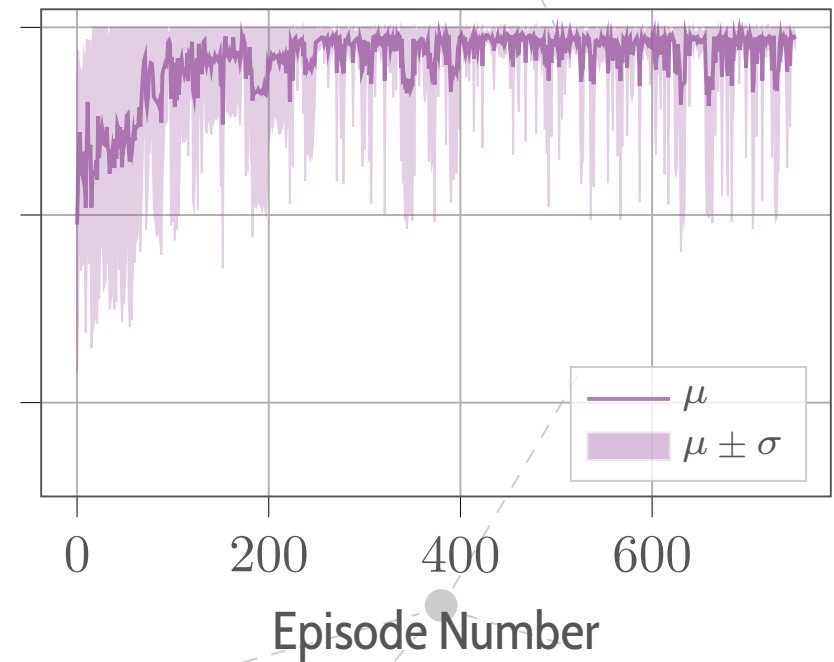


## Safeguard Penalty

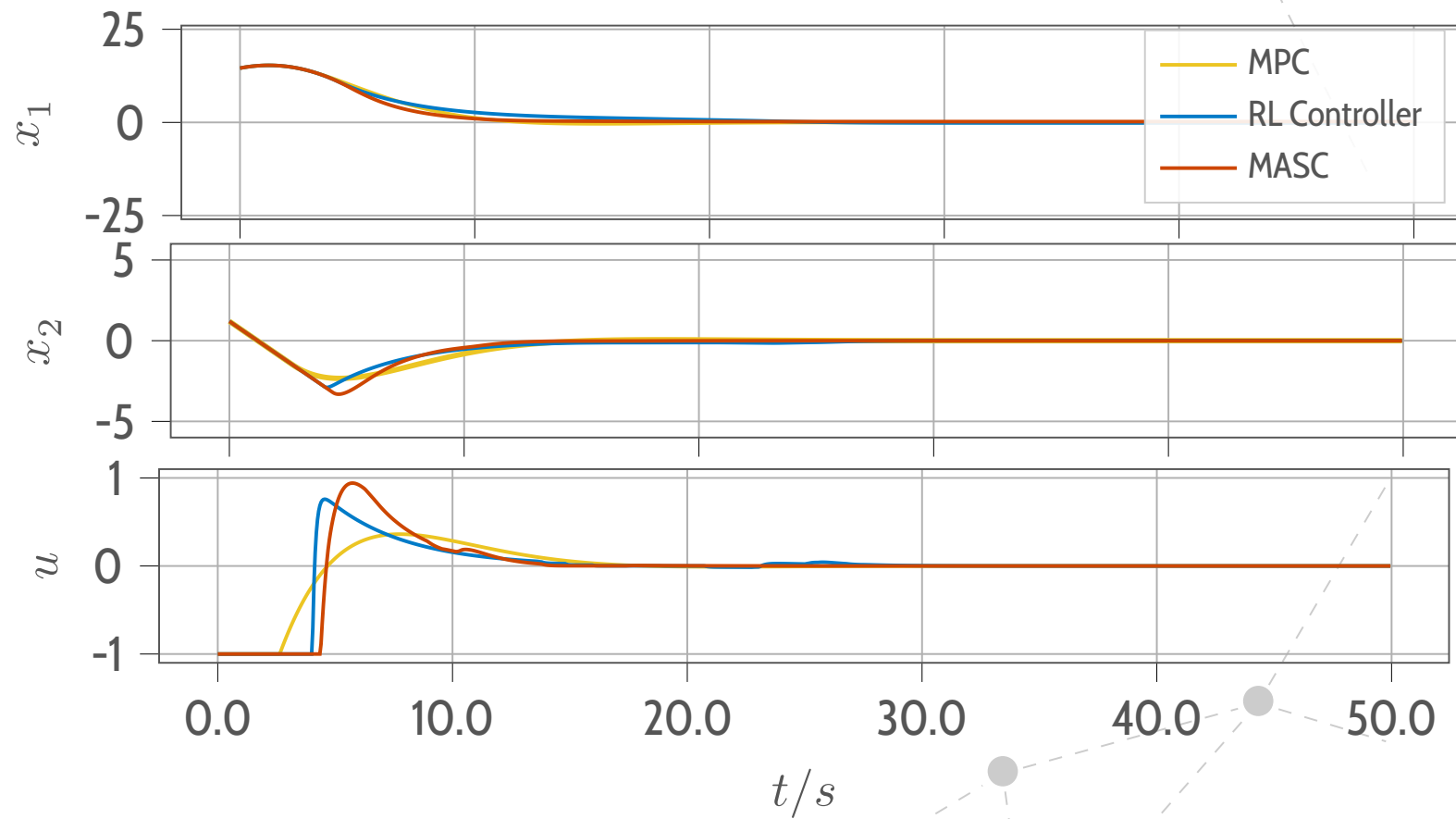
○ Prior available model with safeguard penalty:



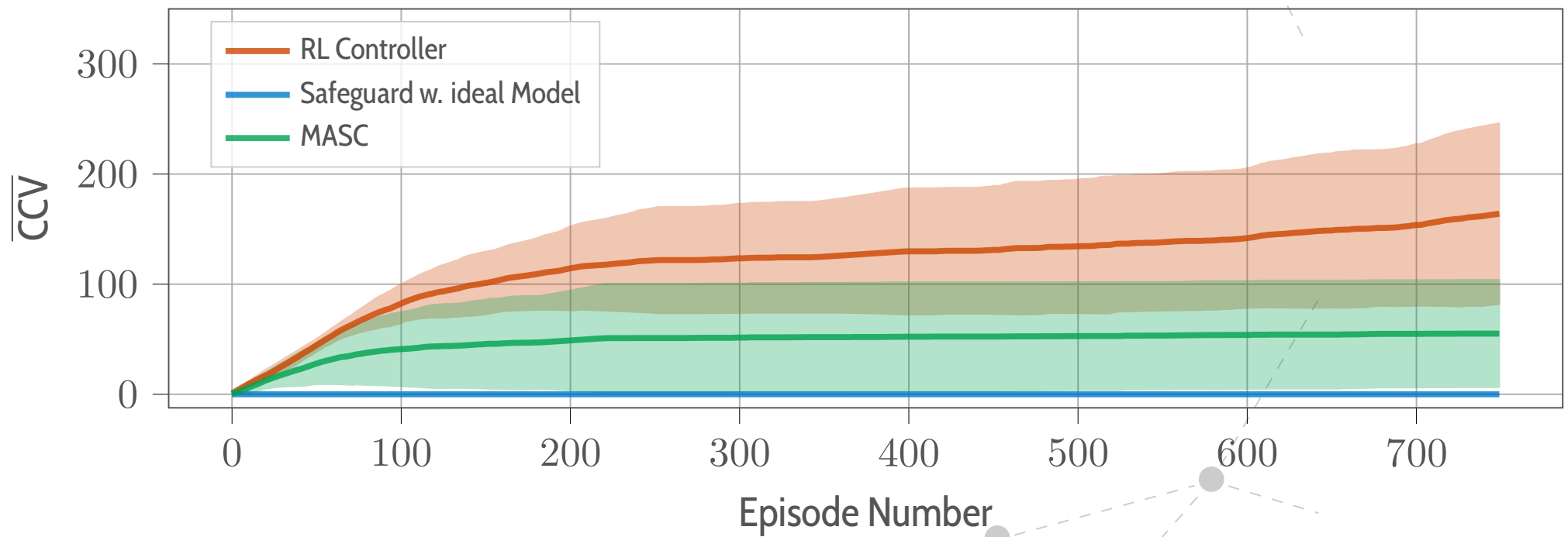
○ Estimated model with safeguard penalty:



## Exemplary Control Trajectories



## Cumulated Constraint Violations



## Summary

Combining model-based and model-free control

- Extension of an off-policy actor-critic controller with the so called safeguard to enforce constraint satisfaction
- Utilize the feasible-set as safeguard constraints to restrain the RL controller to move in the feasible state-action space
- Use the RLS algorithm for online model estimation

## Future Work

- Use computational less expensive methods to determine the feasible set
- Extend exploration mechanism of the RL controller:
  - Utilize estimated model to employ MPC for directional, controlled exploration on basis of risk and curiosity criteria
  - providing safer exploration during the starting model identification phase
- Extensive hyperparameter optimization of full setup MASC

Thank you for your attention.