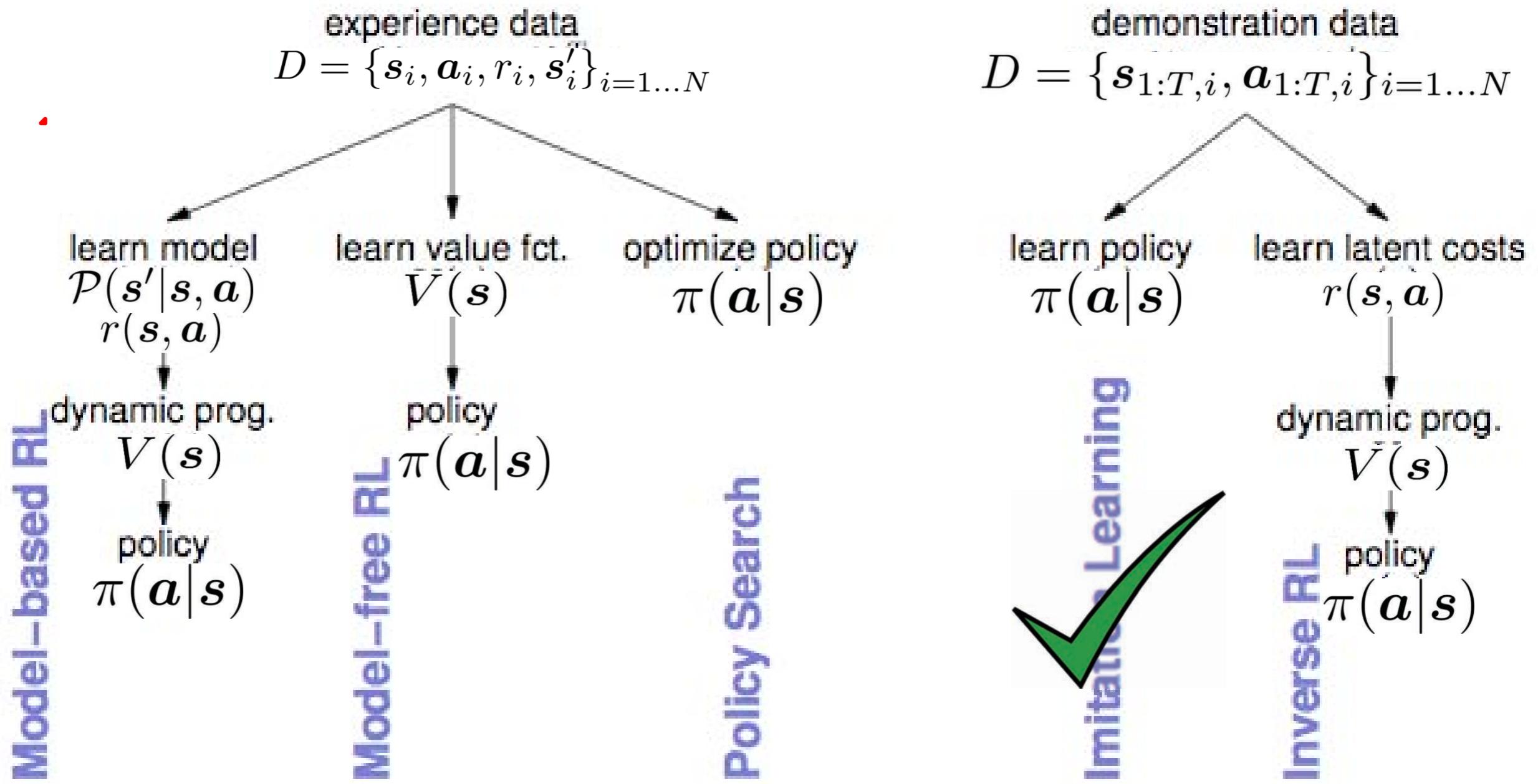


# Reinforcement Learning Part I: Continuous State-Action Optimal Control ...with Learned Models

**Jan Peters**  
**Gerhard Neumann**

# The Bigger Picture: How to learn policies



1. This Lecture

2.

3.

4.



# Motivation

Today we want to use **optimal decision making** for a specific system

- Linear system, Quadratic Reward, Gaussian Noise
  - The optimal policy is called **Linear Quadratic Regulator (LQR)**
- Optimal Decision making for continuous dynamical systems is also called **Optimal Control**

Why? Its the only **continuous case** where we can do it analytically...

If we do not know the model, we **can learn it!**

If the (learned) model is not linear, **we can linearize it!**

# Outline of the Lecture



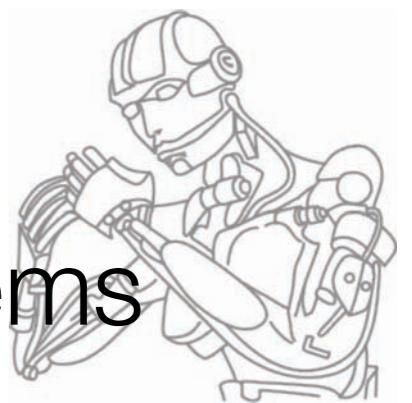
## **1. Optimal Control**

## **2. Solving the Optimal Control for LQR systems**

## **3. Approximating Non-Linear Systems**

## **4. Optimal Control with Learned Models**

## **5. Final Remarks**



# Optimal Decision Making in Continuous Systems

**In continuous systems we call it Optimal Control**

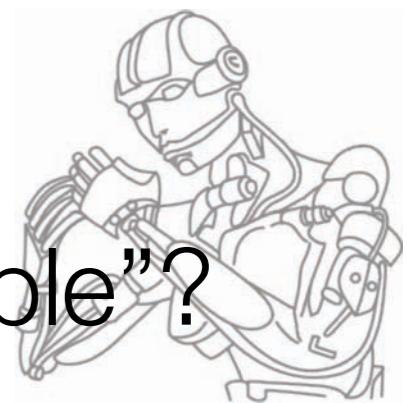
- continuous state space  $s \in \mathbb{R}^n$  (note: will be called  $\boldsymbol{x}$  )
- continuous action space  $a \in \mathbb{R}^m$  (note: same as  $\boldsymbol{u}$  )
- its transition dynamics as density

$$\mathcal{P}_t(s_{t+1}|s_t, a_t) = p_t(x_{t+1}|x_t, u_t)$$

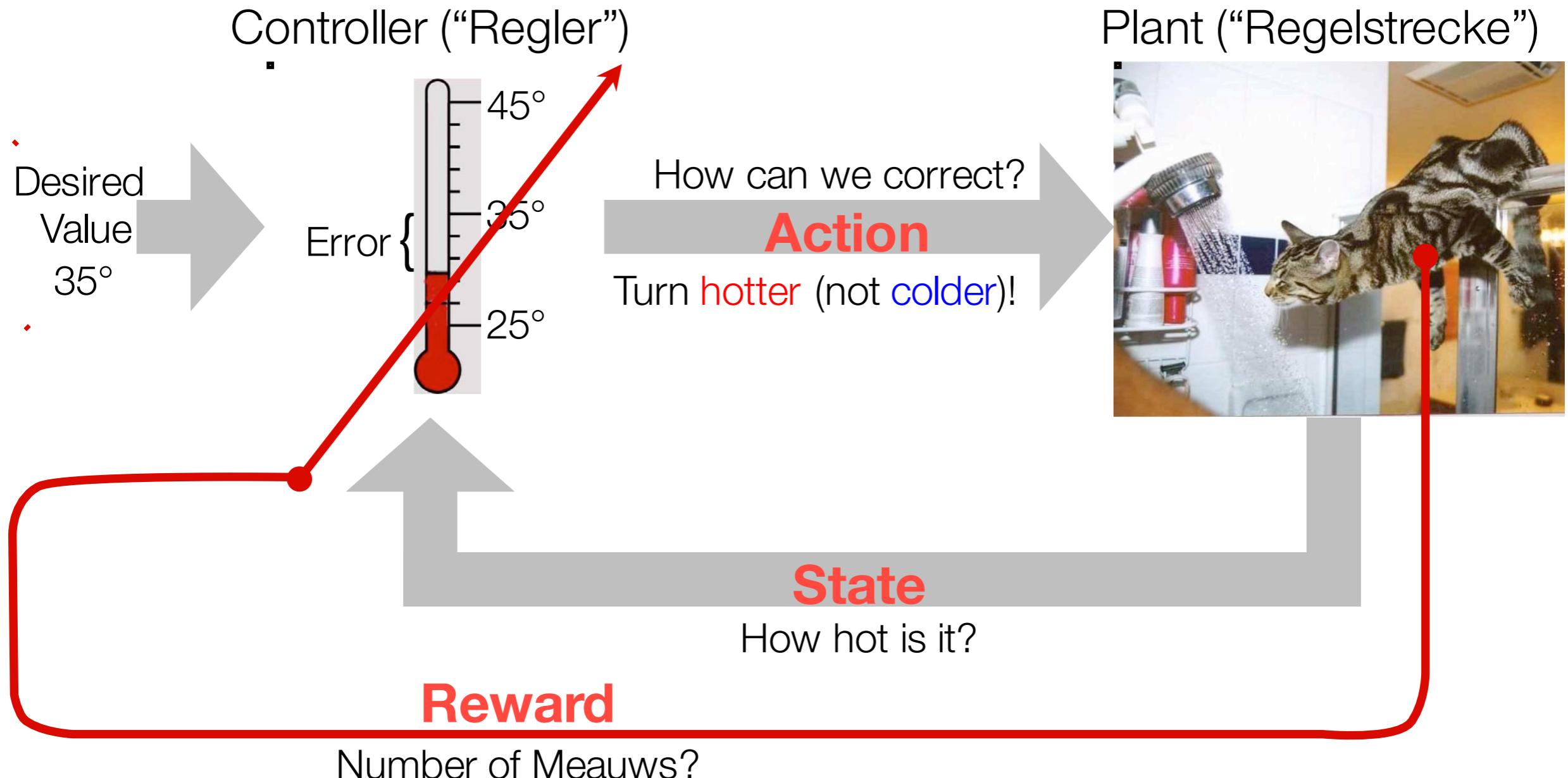
→ We use the more common **optimal control notation**

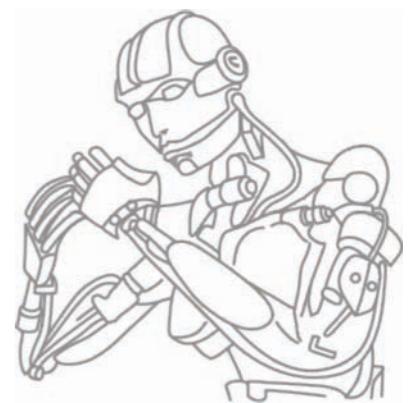
**First question:**

→ **How to define a reward for continuous systems?**



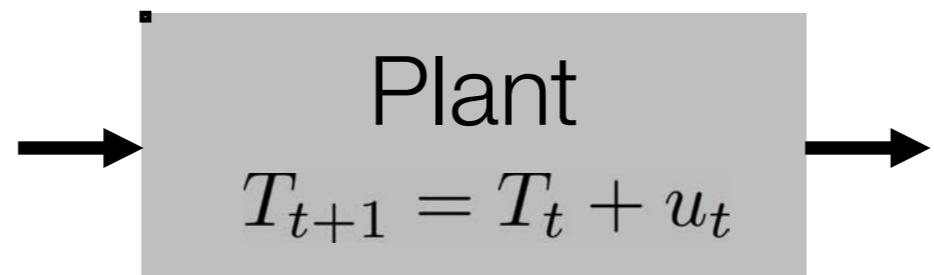
# Illustration: Remember our “Showering Example”?





# Let's Model the System

The system



can be modeled as with

$$p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t) = p(T_{t+1}|T_t, u_t) = \mathcal{N}(T_{t+1}|T_t + u_t, \sigma^2)$$

with  $\mathbf{x} = T, \dots, \mathbf{u} = u$

What kind of rewards induce which behavior?

# Let us find a reward!

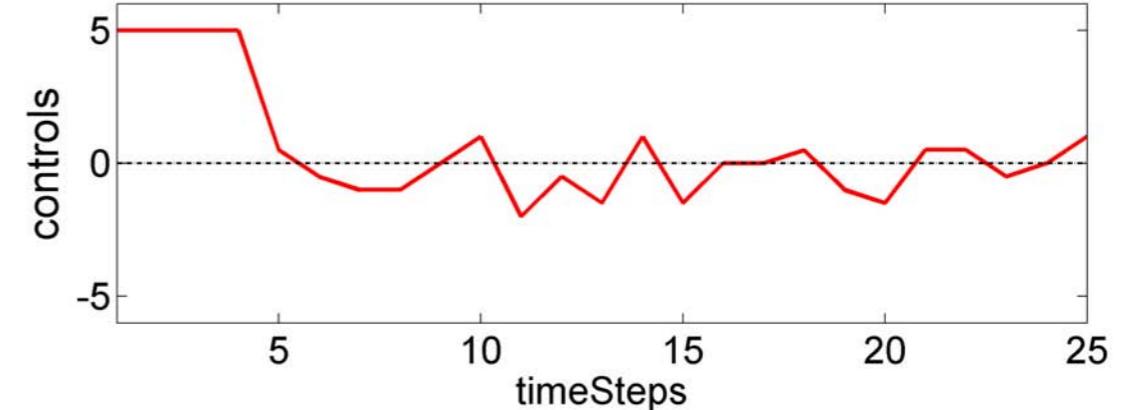
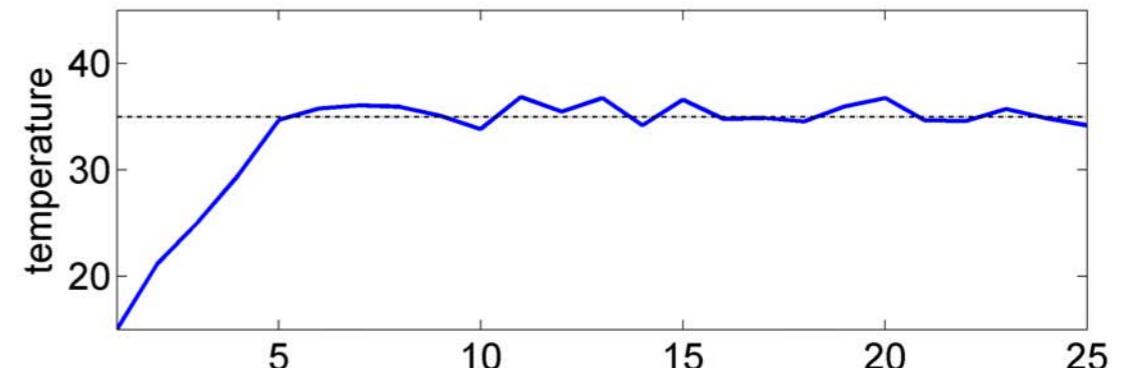
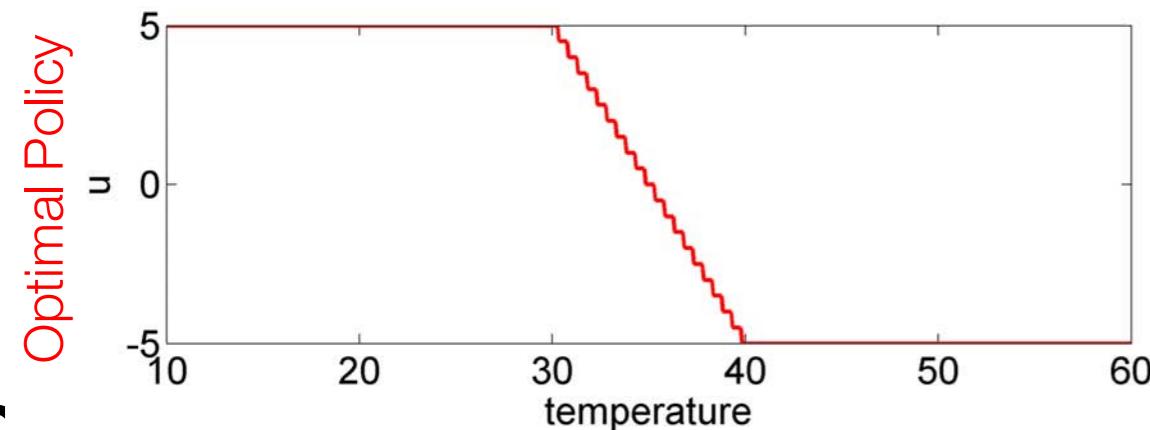
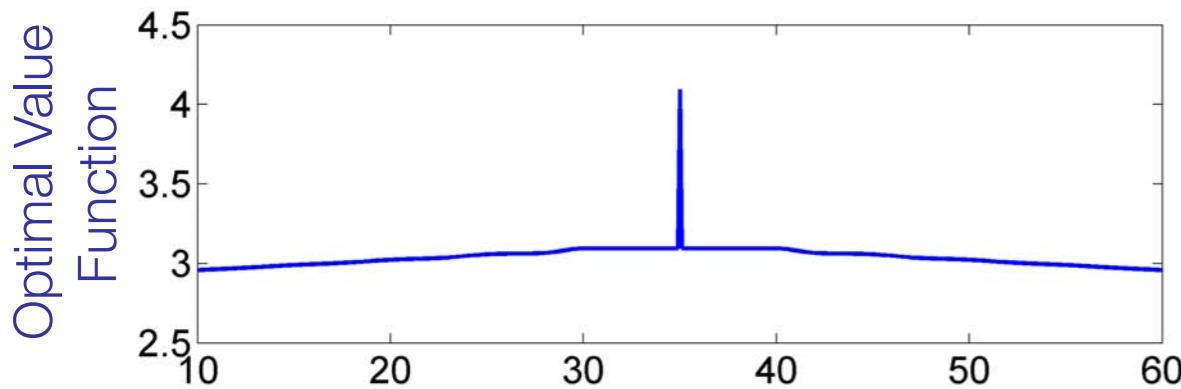


**Only give rewards to good temperatures:**

$$r(T, u) = \begin{cases} 1, & \text{if } T = T_{\text{des}} \\ 0, & \text{otherwise} \end{cases}$$

**How does the controller look?**

→ Rather jerky controls



# Let us find a reward!

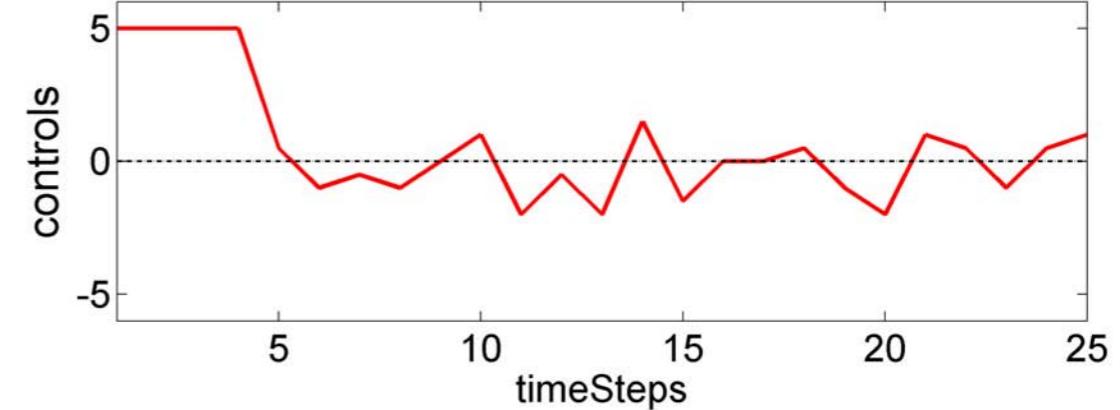
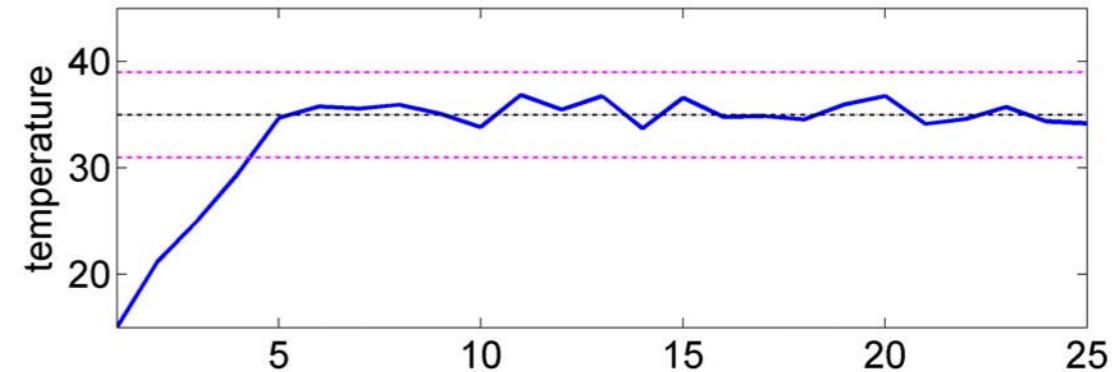
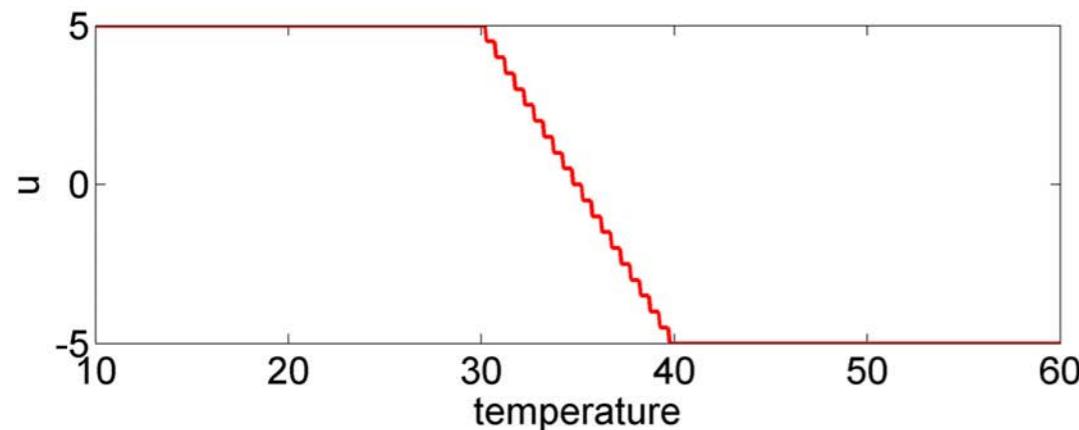
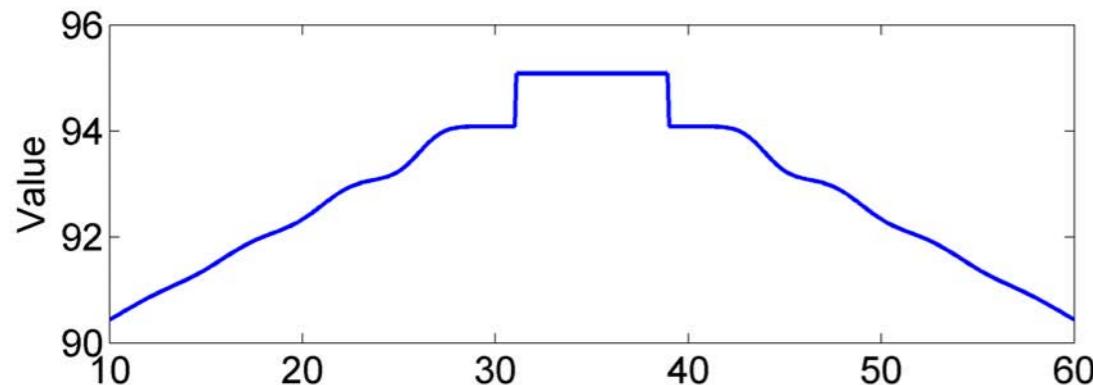


**We could enlarge the region:**

$$r(T, u) = \begin{cases} 1, & \text{if } |T - T_{\text{des}}| < 4 \\ 0, & \text{otherwise} \end{cases}$$

**How does the controller look?**

→ Rather jerky controls



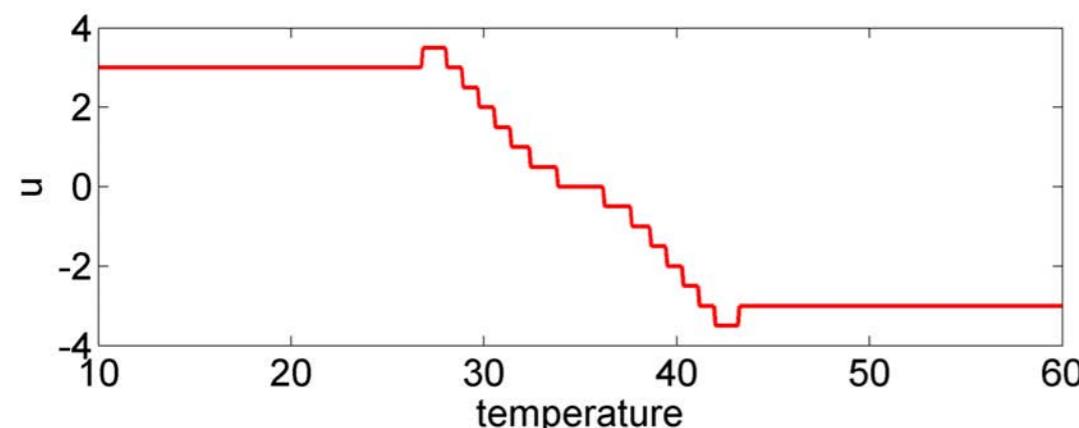
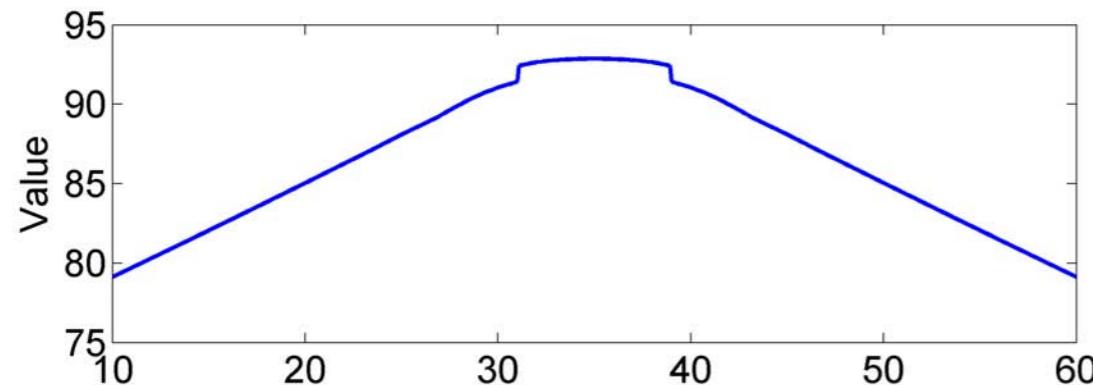


Let us find a reward!

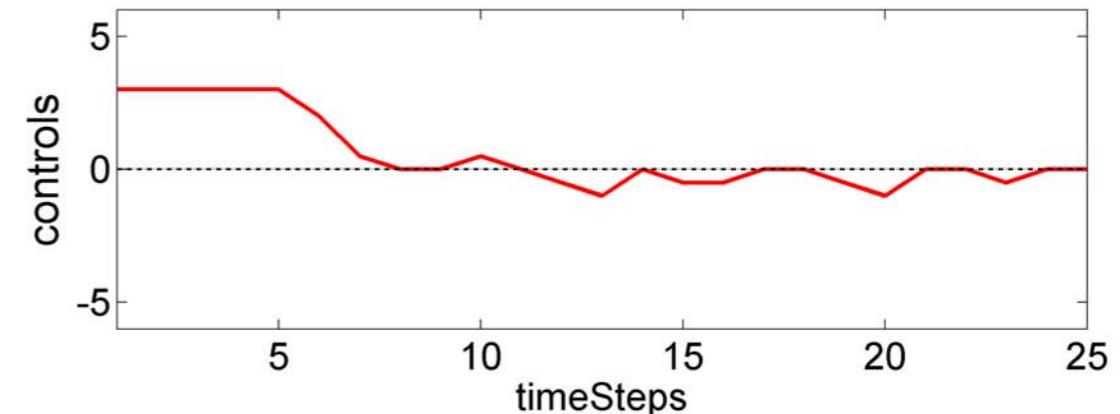
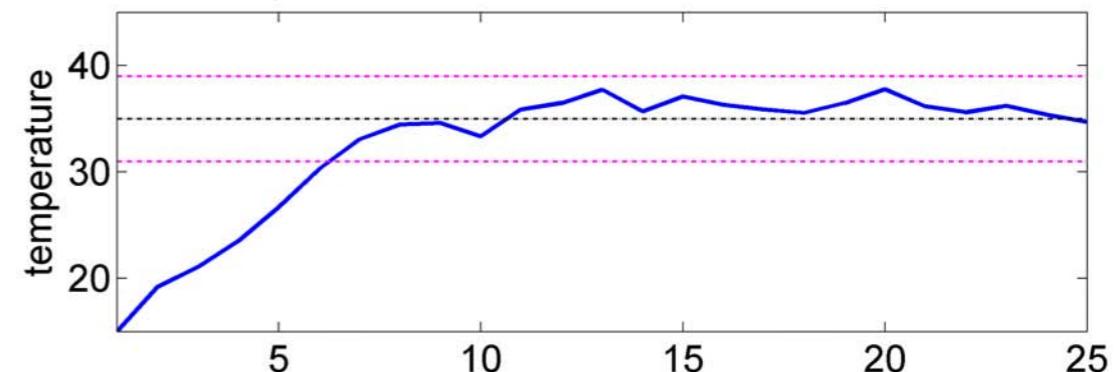
**Punish if we need to turn the nob too much:**

$$r(T, u) = \begin{cases} 1, & \text{if } |T - T_{\text{des}}| < 4 \\ -0.1u^2, & \text{otherwise} \end{cases}$$

**How does the controller look?**



→ Still complex value function and policy



# Let us find a reward!

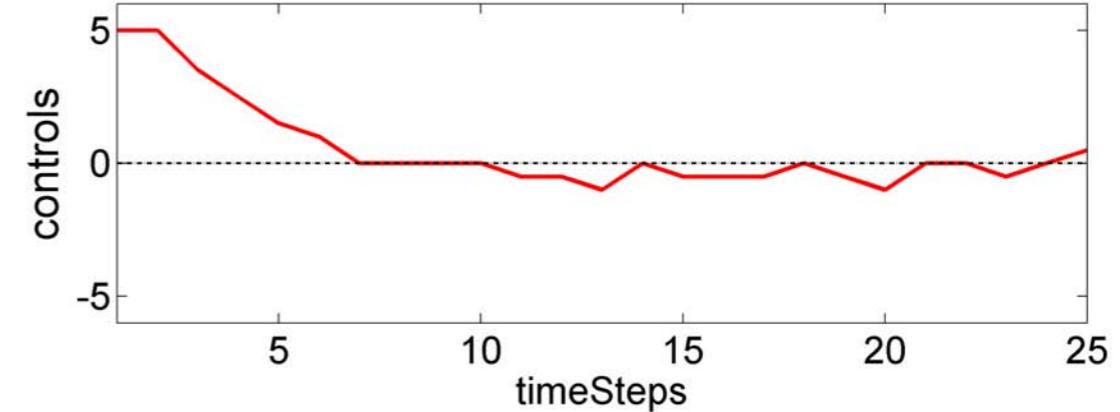
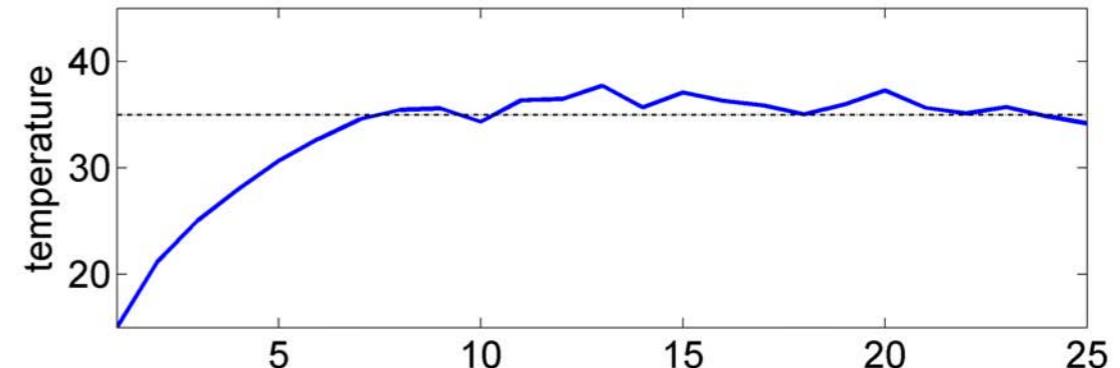
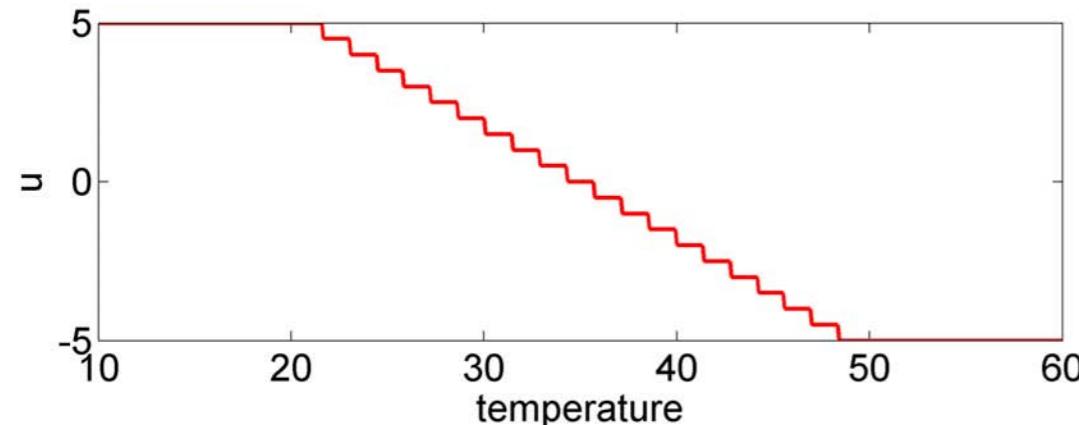
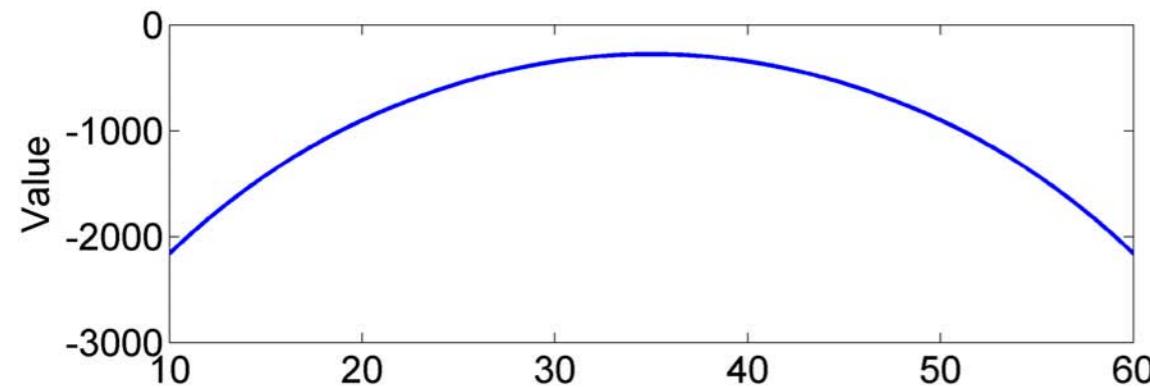


## Punish turning the nob and high deviations:

$$r(T, u) = -(T - T_{\text{des}})^2 - 5u^2$$

### How does the controller look?

- Linear optimal controller
- Quadratic Value Function



# Outline of the Lecture



- 1. Optimal Control**
- 2. Solving the Optimal Control for LQR systems**
- 3. Approximating Non-Linear Systems**
- 4. Optimal Control with Learned Models**
- 5. Final Remarks**



# Finite Horizon Objectives

The goal of the agent is to find a policy  $\pi(a|s)$  that maximizes its expected return  $J_\pi$  **for a finite time horizon**

**Finite Horizon T:** Accumulated expected reward for T steps

$$J_\pi = \mathbb{E}_{\mu_0, \mathcal{P}, \pi} \left[ \sum_{t=1}^{T-1} r_t(s_t, a_t) + r_T(s_T) \right]$$

$r_T(s_T)$  ... final reward



# Linear Quadratic Gaussian systems

An **LQR** system is defined as

- its state space  $\mathbf{x} \in \mathbb{R}^n$  (note: same as  $s$ )
- its action space  $\mathbf{u} \in \mathbb{R}^m$  (note: same as  $a$ )
- its (possibly time-dependent) **linear transition dynamics with Gaussian noise**

$$p_t(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \Sigma_t)$$

- its **quadratic** reward function

$$r_t(\mathbf{x}, \mathbf{u}) = (\mathbf{x} - \mathbf{r}_t)^T \mathbf{R}_t (\mathbf{x} - \mathbf{r}_t) + \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t$$

$$r_T(\mathbf{x}) = (\mathbf{x} - \mathbf{r}_T)^T \mathbf{R}_T (\mathbf{x} - \mathbf{r}_T)$$

- and its initial state density

$$\mu_0(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_0, \Sigma_0)$$



# Optimal Control for LQR systems

## Linear systems with Gaussian Noise

$$p_t(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \Sigma_t)$$

$\mathbf{A}_t$  ... system matrix,  $\mathbf{B}_t$  ... control matrix,  $\mathbf{b}_t$  ... drift term

$\Sigma_t$  ... system noise

## Quadratic reward functions

$$r_t(\mathbf{x}, \mathbf{u}) = -(\mathbf{x} - \mathbf{r}_t)^T \mathbf{R}_t (\mathbf{x} - \mathbf{r}_t) - \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t$$

$$r_T(\mathbf{x}) = -(\mathbf{x} - \mathbf{r}_T)^T \mathbf{R}_T (\mathbf{x} - \mathbf{r}_T)$$

$\mathbf{r}_t$  ... desired state,  $\mathbf{R}_t$  ... state metric for reward

$\mathbf{H}_t$  ... control metric for reward

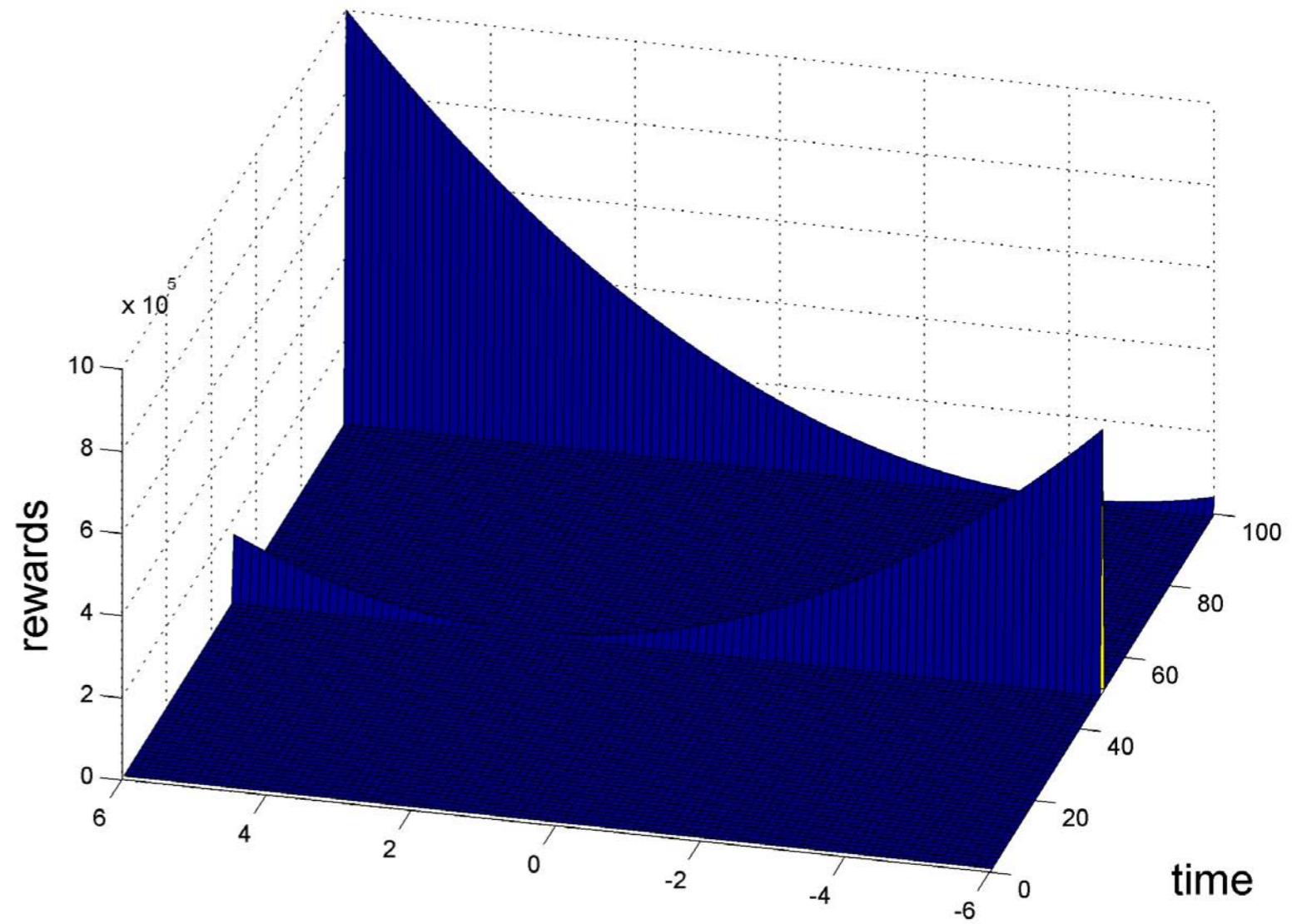
# Example



**RewardFunction:** Reach 2 Via-Points at  $t_1 = 50$  and  $t_2 = 100$

$$\mathbf{R}_{t_1, t_2} = \begin{bmatrix} 10^4 & 0 \\ 0 & 10^{-6} \end{bmatrix}, \text{ for all other } \tilde{t}, \mathbf{R}_{\tilde{t}} = \begin{bmatrix} 10^{-6} & 0 \\ 0 & 10^{-6} \end{bmatrix}$$

$$\mathbf{r}_{t_1} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}, \mathbf{r}_{t_2} = \begin{bmatrix} -4 \\ 0 \end{bmatrix}$$





# Optimal Control for Finite Horizon Objectives

We will look at the **simpler finite horizon case**

**Short refresher** from last lecture:

**Start with last layer...**

$$V_T^*(\mathbf{x}) = r_T(\mathbf{x})$$

Iterate **backwards in time**

$$V_t^*(\mathbf{x}) = \max_{\mathbf{u}} (r_t(\mathbf{x}_t, \mathbf{u}_t) + \mathbb{E}_p [V_{t+1}^*(\mathbf{x}_{t+1}) | \mathbf{x}_t, \mathbf{u}_t])$$

The optimal value function/policy for time step  $t$  is obtained after  $T - t + 1$  iterations

$$V_T^*(\mathbf{x}) \rightarrow V_{T-1}^*(\mathbf{x}) \rightarrow \dots \rightarrow V_1^*(\mathbf{x})$$



# Optimal Control for LQR systems

**We have to solve...**

- **Expectation** over the next value:

$$\mathbb{E}_{p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)} [V_{t+1}^*(\mathbf{x}_{t+1})|\mathbf{x}_t, \mathbf{u}_t]$$

- **Maximum** operator in continuous action spaces:

$$\max_{\mathbf{u}} (r_t(\mathbf{x}_t, \mathbf{u}_t) + \mathbb{E}_p [V_{t+1}^*(\mathbf{x}_{t+1})|\mathbf{x}_t, \mathbf{u}_t])$$

**When can we do that?**

In continuous systems: **only for LQR systems!**



Ok, lets solve the optimal control problem

For illustration, lets make it simpler (without any linear terms)...

$$p_t(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \Sigma_t)$$
$$r_t(\mathbf{x}, \mathbf{u}) = -(\mathbf{x} - \mathbf{r}_t)^T \mathbf{R}_t (\mathbf{x} - \mathbf{r}_t) - \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t$$
$$r_T(\mathbf{x}) = -(\mathbf{x} - \mathbf{r}_T)^T \mathbf{R}_T (\mathbf{x} - \mathbf{r}_T)$$

$$p_t(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t, \Sigma_t)$$



$$r_t(\mathbf{x}, \mathbf{u}) = -\mathbf{x}^T \mathbf{R}_t \mathbf{x} - \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t$$
$$r_T(\mathbf{x}) = -\mathbf{x}^T \mathbf{R}_T \mathbf{x}$$

For the derivation of the full problem including the drift and linear terms in the reward, see  
<http://ipvs.informatik.uni-stuttgart.de/mlr/marc/notes/soc.pdf>



# Bellman's Recipe

**1. At the last step, the value function is given as**

$$V_T^*(\mathbf{x}) = r_T(\mathbf{x}) = -\mathbf{x}^T \mathbf{R}_T \mathbf{x} = -\mathbf{x}^T \mathbf{V}_T \mathbf{x}$$

→  $\mathbf{V}_T = \mathbf{R}_T$

**2. To get from  $t+1$  to  $t$ , first compute the Q-Function**

$$Q_t^*(\mathbf{x}_t, \mathbf{u}_t) = r_t(\mathbf{x}_t, \mathbf{u}_t) + \mathbb{E}_p [V_{t+1}^*(\mathbf{x}_{t+1}) | \mathbf{x}_t, \mathbf{u}_t]$$

**3. then compute optimal policy  $\pi_t^*$**

$$\pi_t^*(\mathbf{x}) = \operatorname{argmax}_{\mathbf{u}} Q_t^*(\mathbf{x}, \mathbf{u})$$

**4. compute optimal value function for time step t**

$$V_t^*(\mathbf{x}) = Q_t^*(\mathbf{x}_t, \pi^*(\mathbf{x}))$$



# Bellman's Recipe

## Step 2a: **compute expectation for value of the next state**

→ Lets assume for now a quadratic structure for V-function of next time step

$$V_{t+1}^*(\mathbf{x}) = -\mathbf{x}^T \mathbf{V}_{t+1} \mathbf{x}$$

→ We need to compute:

$$\mathbb{E}_p [V_{t+1}^*(\mathbf{x}_{t+1}) | \mathbf{x}_t, \mathbf{u}_t] = - \int \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t, \Sigma_t) \mathbf{x}_{t+1}^T \mathbf{V}_{t+1} \mathbf{x}_{t+1} d\mathbf{x}_{t+1}$$

**Yet another useful Gaussian identity:** 2<sup>nd</sup> order expectation

if  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  than  $\mathbb{E}_p[\mathbf{x}^T \mathbf{M} \mathbf{x}] = \boldsymbol{\mu}^T \mathbf{M} \boldsymbol{\mu} + \text{Tr}(\mathbf{M} \boldsymbol{\Sigma})$

→ This identity yields

$$\mathbb{E}_p [V_{t+1}^*(\mathbf{x}_{t+1}) | \mathbf{x}_t, \mathbf{u}_t] = -(\mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t)^T \mathbf{V}_{t+1} (\mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t) + \text{Tr}(\mathbf{V}_{t+1} \boldsymbol{\Sigma}_t)$$



# Bellman's Recipe

## Step 2b: Compute **Q-function**

$$\begin{aligned} Q_t^*(\mathbf{x}_t, \mathbf{u}_t) &= r_t(\mathbf{x}_t, \mathbf{u}_t) + \mathbb{E}_p [V_{t+1}^*(\mathbf{x}_{t+1}) | \mathbf{x}_t, \mathbf{u}_t] \\ &= -\mathbf{x}^T \mathbf{R}_t \mathbf{x} - \mathbf{u}^T \mathbf{H}_t \mathbf{u} \\ &\quad - (\mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t)^T \mathbf{V}_{t+1} (\mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t) + \text{Tr}[\Sigma_{t+1} \Sigma_t] \end{aligned}$$

Not state or action dependent

→ Also the Q-function is quadratic in state and action!



# Bellman's Recipe

**Step 3: compute optimal policy**

$$\pi_t^*(x) = \operatorname{argmax}_u Q_t^*(x, u)$$

**Set derivation to zero...**

$$0^T = \frac{d}{du} (-x_t^T R_t x_t - u^T H_t u - (A_t x_t + B_t u)^T V_{t+1} (A_t x_t + B_t u))$$

**Remember matrix calculus...?**

$$0^T = -2u^T H_t - 2(A_t x_t + B_t u)^T V_{t+1} B_t$$

$$0^T = -u^T (H_t + B_t^T V_{t+1} B_t) - x_t^T A^T V_{t+1} B_t$$

**And solve for u**

$$\pi^*(x_t) = u^* = -(H_t + B_t^T V_{t+1} B_t)^{-1} B_t^T V_{t+1} A_t x_t = K_t x_t$$

→ The optimal policy a time-varying linear (PD) controller!



# Bellman's Recipe

## Step 4: compute value function

$$\begin{aligned}
 V_t^*(\mathbf{x}) &= -\mathbf{x}_t^T \mathbf{R}_t \mathbf{x}_t - \mathbf{u}_t^{*T} \mathbf{H}_t \mathbf{u}_t^* - (\mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t^*)^T \mathbf{V}_{t+1} (\mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t^*) \\
 &= -\mathbf{x}_t^T (\mathbf{R}_t + \mathbf{A}_t^T \mathbf{V}_{t+1} \mathbf{A}_t) \mathbf{x}_t - \mathbf{u}_t^{*T} (\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t) \mathbf{u}_t^* \\
 &\quad - 2\mathbf{u}_t^{*T} \mathbf{B}_t \mathbf{V}_{t+1} \mathbf{A}_t \mathbf{x}_t
 \end{aligned}$$

We first set in the optimal action for  $\mathbf{u}_t^* = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t \mathbf{x}_t$

$$V_t^*(\mathbf{x}) = -\mathbf{x}_t^T (\mathbf{R}_t + \mathbf{A}_t^T \mathbf{V}_{t+1} \mathbf{A}_t) \mathbf{x}_t - \mathbf{u}_t^{*T} \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t \mathbf{x}_t$$

Now we can substitute  $\mathbf{u}_t^* = \mathbf{K}_t \mathbf{x}_t$

$$V_t^*(\mathbf{x}) = -\mathbf{x}_t^T (\mathbf{R}_t + \mathbf{A}_t^T \mathbf{V}_{t+1} \mathbf{A}_t + \mathbf{K}_t^T \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t) \mathbf{x}_t$$

**Note:** this derivation works only if the matrices  $\mathbf{V}_{t+1}, \mathbf{H}_t$  and  $\mathbf{R}_t$  are positive definite (and hence symmetric), can always be guaranteed



# Bellman's Recipe

## Step 4: compute value function

→ The optimal value function  $V_t$  for time step  $t+1$  is also quadratic

$$V_t^*(x) = -x_t^T V_t x_t$$

→ we ended up in a recursive update equation for

$$\begin{aligned} V_t &= R_t + A_t^T V_{t+1} A_t + K_t^T B_t^T V_{t+1} A_t \\ &= R_t + (A_t + B_t K_t)^T V_{t+1} A_t \end{aligned}$$

with  $K_t = -(H_t + B_t^T V_{t+1} B_t)^{-1} B_t^T V_{t+1} A_t$

→ if  $V_{t+1}(s)$  is in quadratic form,  $V_t(s)$  also is

→ since  $V_T(s)$  is quadratic, all  $V_t(s)$  are quadratic



# Solving optimal control

**So how does the **full case** look like?**

$$p_t(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \Sigma_t)$$

$$r_t(\mathbf{x}, \mathbf{u}) = -(\mathbf{x} - \mathbf{r}_t)^T \mathbf{R}_t (\mathbf{x} - \mathbf{r}_t) - \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t$$

**The optimal value function has a **quadratic and linear** form**

$$V_t(\mathbf{x}_t) = -\mathbf{x}_t^T \mathbf{V}_t \mathbf{x}_t + 2\mathbf{v}_t^T \mathbf{x}_t + \text{const}$$

**With the update rules:**

$$\mathbf{V}_t = \mathbf{R}_t + (\mathbf{A}_t + \mathbf{B}_t \mathbf{K}_t)^T \mathbf{V}_{t+1} \mathbf{A}_t \quad (\text{same as before})$$

$$\mathbf{v}_t = \tilde{\mathbf{r}}_t + (\mathbf{A}_t + \mathbf{B}_t \mathbf{K}_t)^T (\mathbf{v}_{t+1} - \mathbf{V}_{t+1} \mathbf{b}_t)$$

$$\text{with } \tilde{\mathbf{r}}_t = \mathbf{r}_t^T \mathbf{R}_t \quad \text{and} \quad \mathbf{K}_t = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t$$

For the derivation of the full problem including the drift and linear terms in the reward, see  
<http://ipvs.informatik.uni-stuttgart.de/mlr/marc/notes/soc.pdf>



# Solving optimal control

**So how does the **full case** look like?**

$$p_t(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \Sigma_t)$$

$$r_t(\mathbf{x}, \mathbf{u}) = -(\mathbf{x} - \mathbf{r}_t)^T \mathbf{R}_t (\mathbf{x} - \mathbf{r}_t) - \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t$$

**The optimal policy is given by**

$$\begin{aligned} \mathbf{u}^* &= -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T (\mathbf{V}_{t+1} (\mathbf{A}_t \mathbf{x}_t + \mathbf{b}_t) - \mathbf{v}_{t+1}) \\ &= \mathbf{K}_t \mathbf{x}_t + \mathbf{k}_t \end{aligned}$$

$$\text{with } \mathbf{K}_t = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t$$

$$\text{and } \mathbf{k}_t = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T (\mathbf{V}_{t+1} \mathbf{b}_t - \mathbf{v}_{t+1})$$

**I.e. the optimal policy is a **time-dependent linear feedback controller with time dependent offset****

For the derivation of the full problem including the drift and linear terms in the reward, see  
<http://ipvs.informatik.uni-stuttgart.de/mlr/marc/notes/soc.pdf>

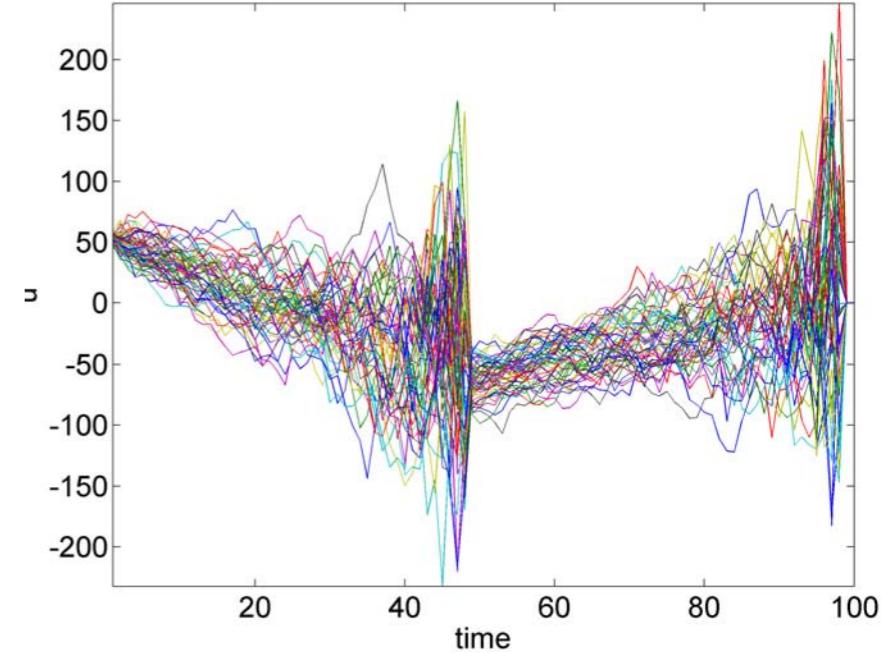
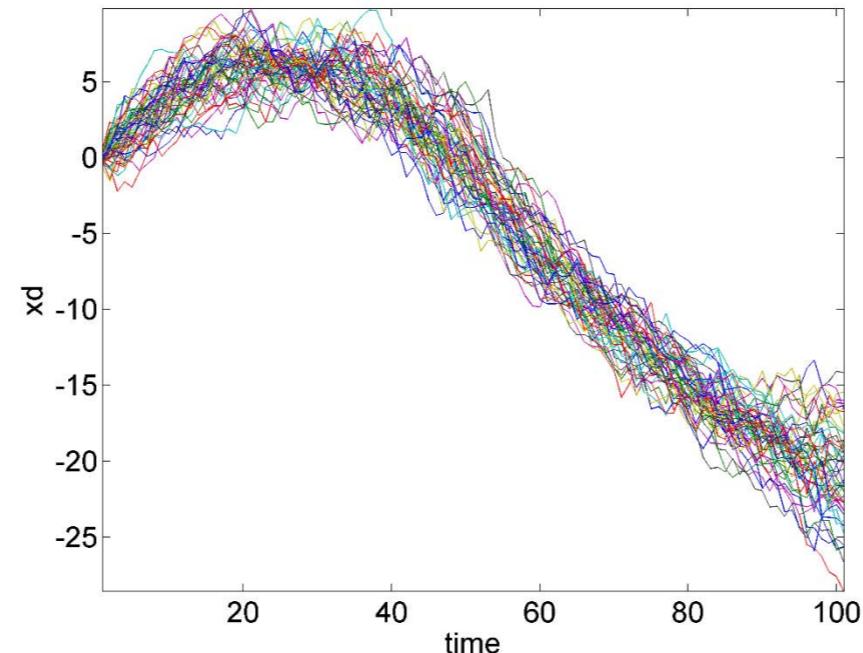
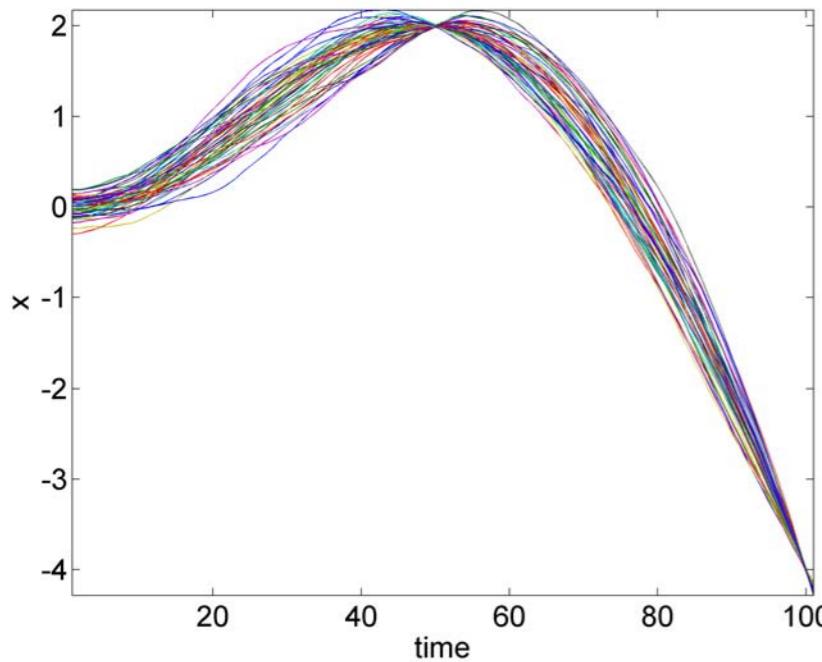


# Back to the Example:

## System:

Second order integrator (we directly set accelerations)

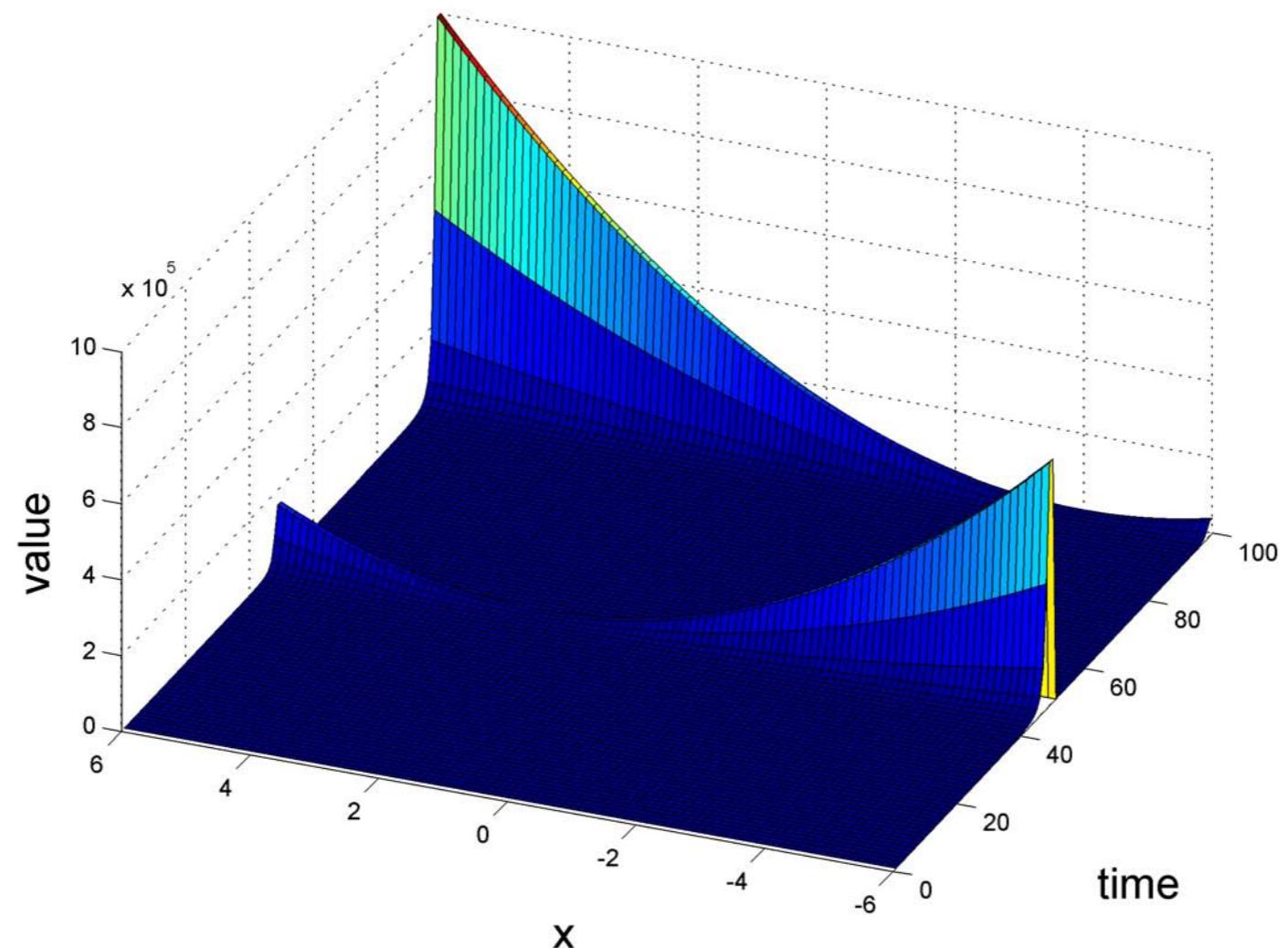
$$\boldsymbol{x}_{t+1} = \underbrace{\begin{bmatrix} 1 & dt \\ 0 & 1 \end{bmatrix}}_A \boldsymbol{x}_t + \underbrace{\begin{bmatrix} 0 \\ dt \end{bmatrix}}_B + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} 0 & 0 \\ 0 & 0.5dt^2 \end{bmatrix}\right)$$





# Optimal Control for LQR systems

## Illustration of the Value Function

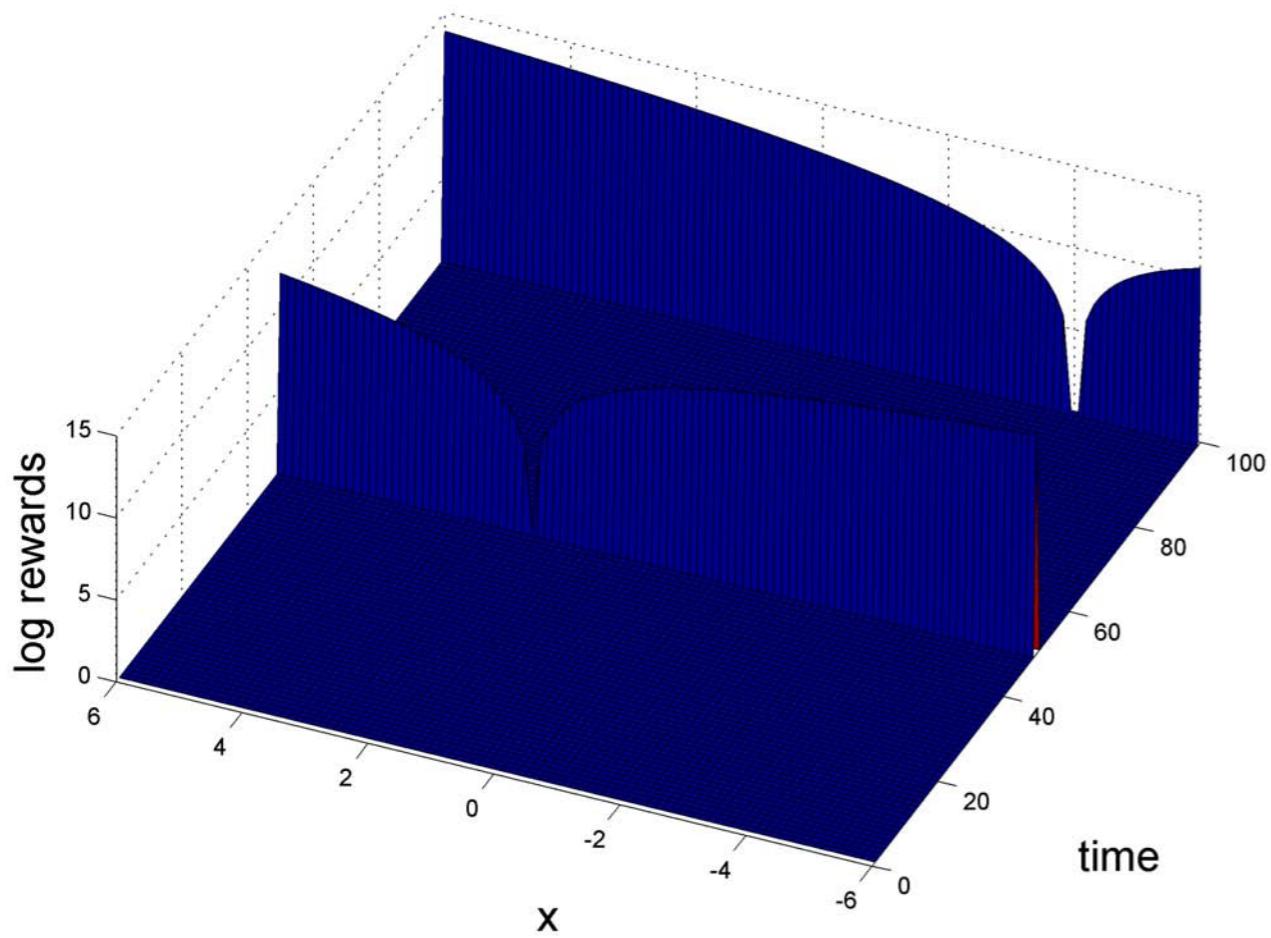


# Optimal Control for LQR systems

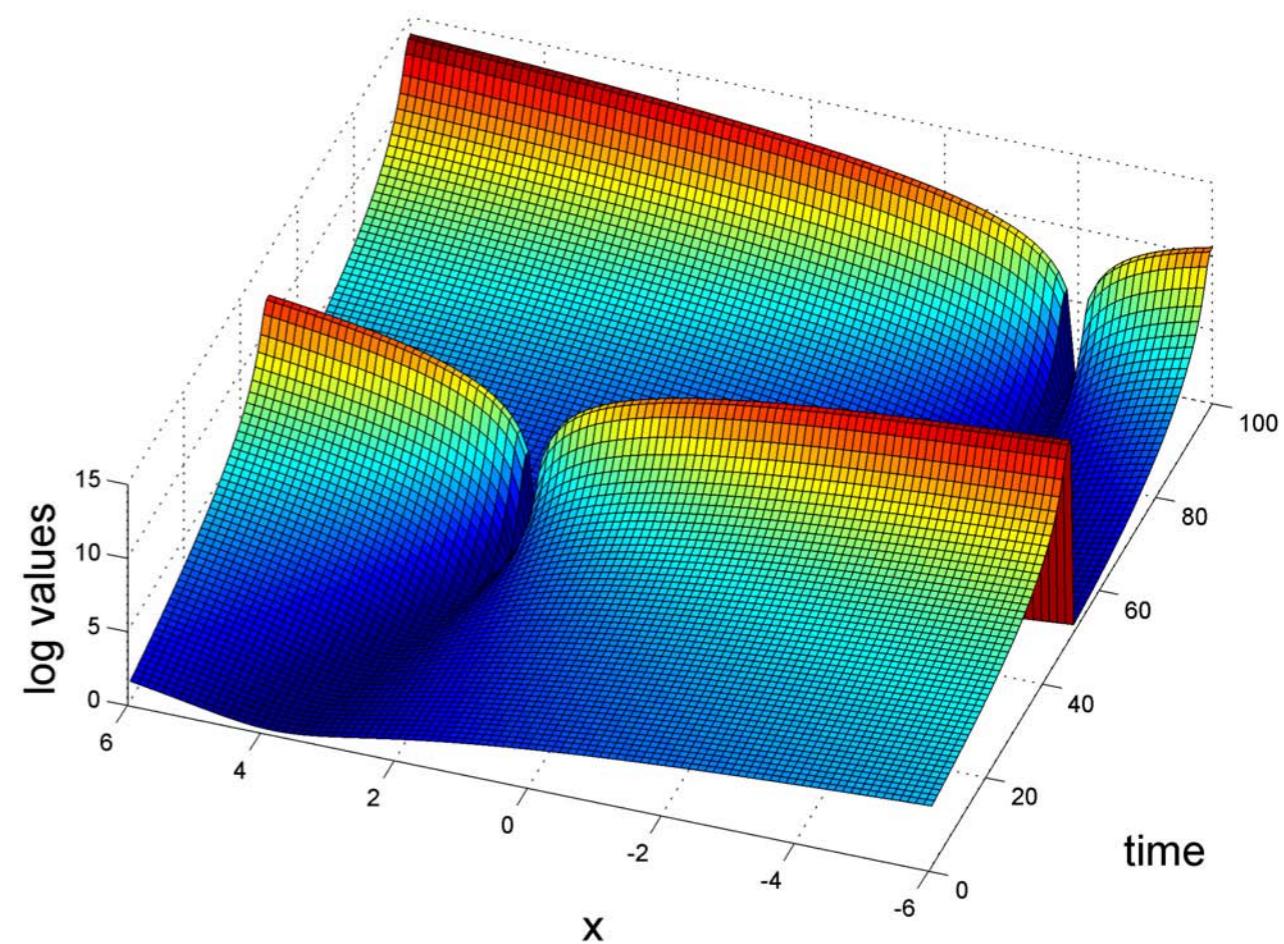


Comparison of Value and Reward Function (log domain)

Rewardfunction



Valuefunction

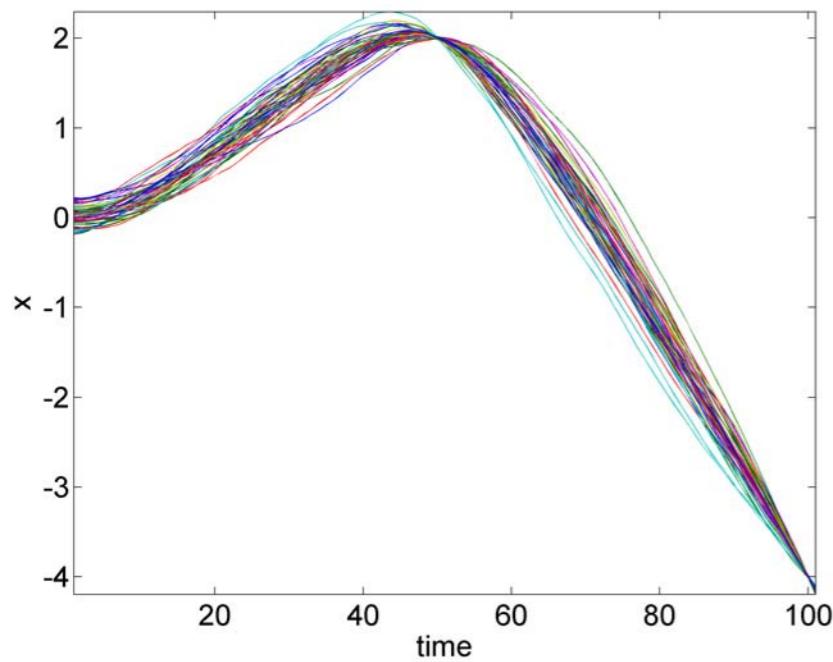


# Optimal Control for LQR systems

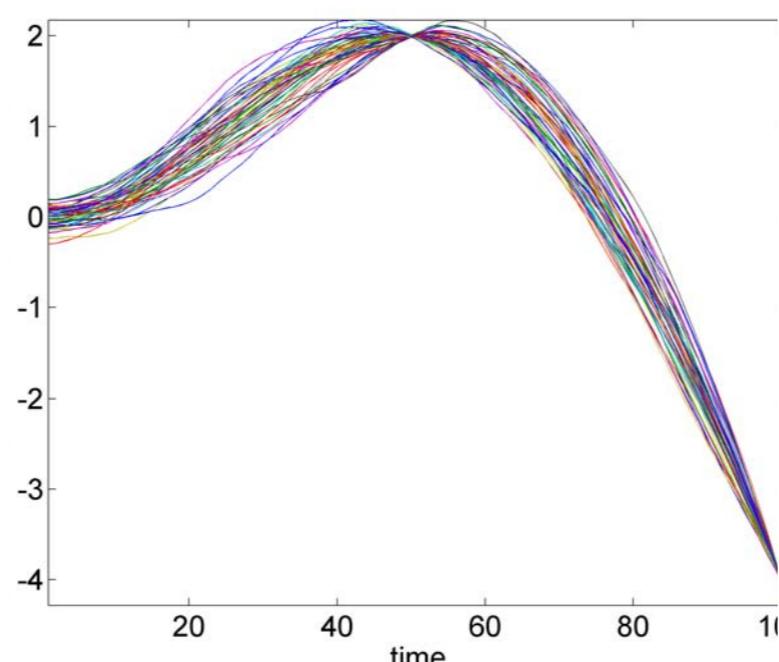


## Different Control Costs

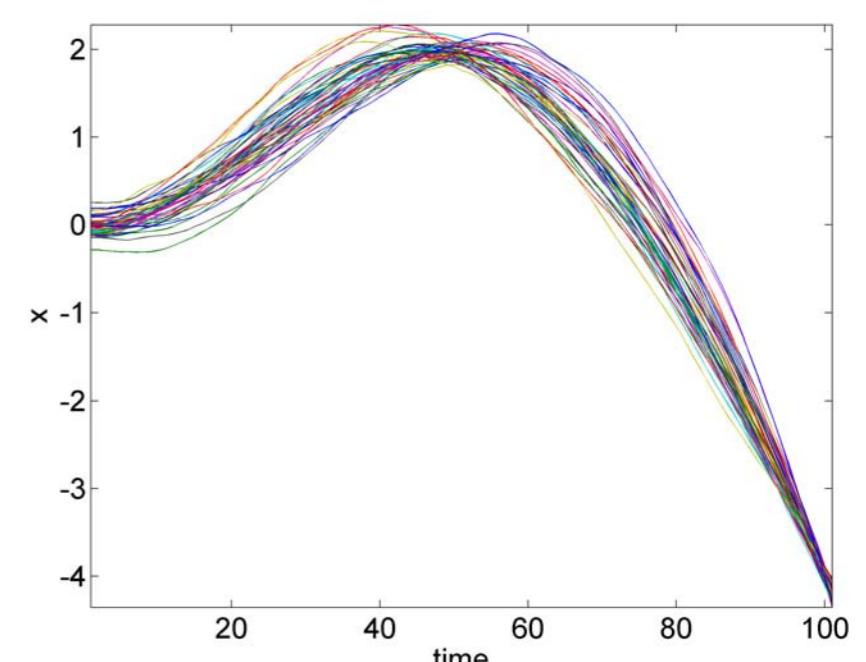
$$\mathbf{H} = 10^{-6}$$



$$\mathbf{H} = 10^{-4}$$



$$\mathbf{H} = 10^{-1}$$



# Outline of the Lecture



**1. Optimal Control**

**2. Solving the Optimal Control for LQR systems**

**3. Approximating Non-Linear Systems**

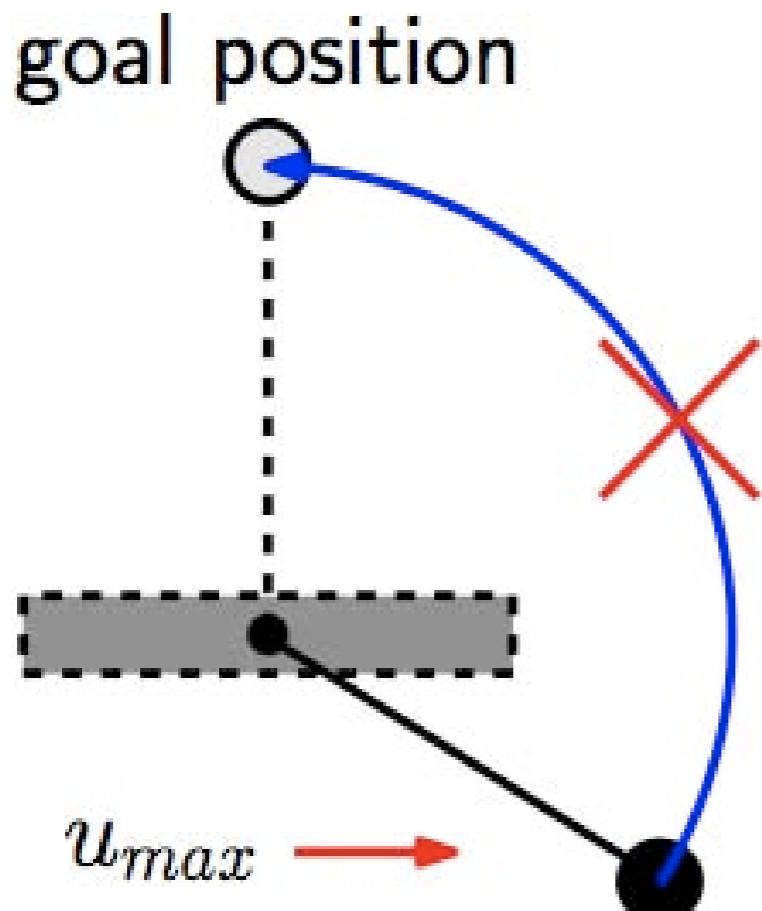
**4. Optimal Control with Learned Models**

**5. Final Remarks**



Example for non-linear dynamics: Swing-Up

## System

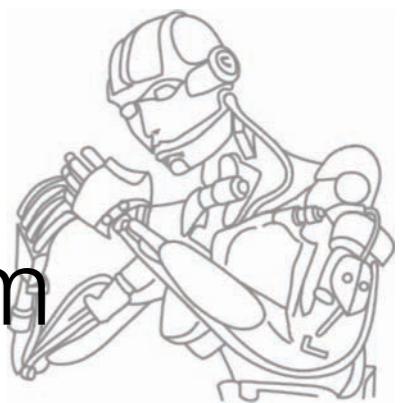


$$\ddot{\varphi}(t) = \frac{-\mu\dot{\varphi}(t) + mgl \sin(\varphi(t)) + u(t)}{ml^2}$$

$$\mathbf{x}_{k+1} := \begin{bmatrix} \varphi_{k+1} \\ \dot{\varphi}_{k+1} \end{bmatrix} = \begin{bmatrix} \varphi_k + \Delta_t \dot{\varphi}_k + \frac{\Delta_t^2}{2} \ddot{\varphi}_k \\ \dot{\varphi}_k + \Delta_t \ddot{\varphi}_k \end{bmatrix}$$

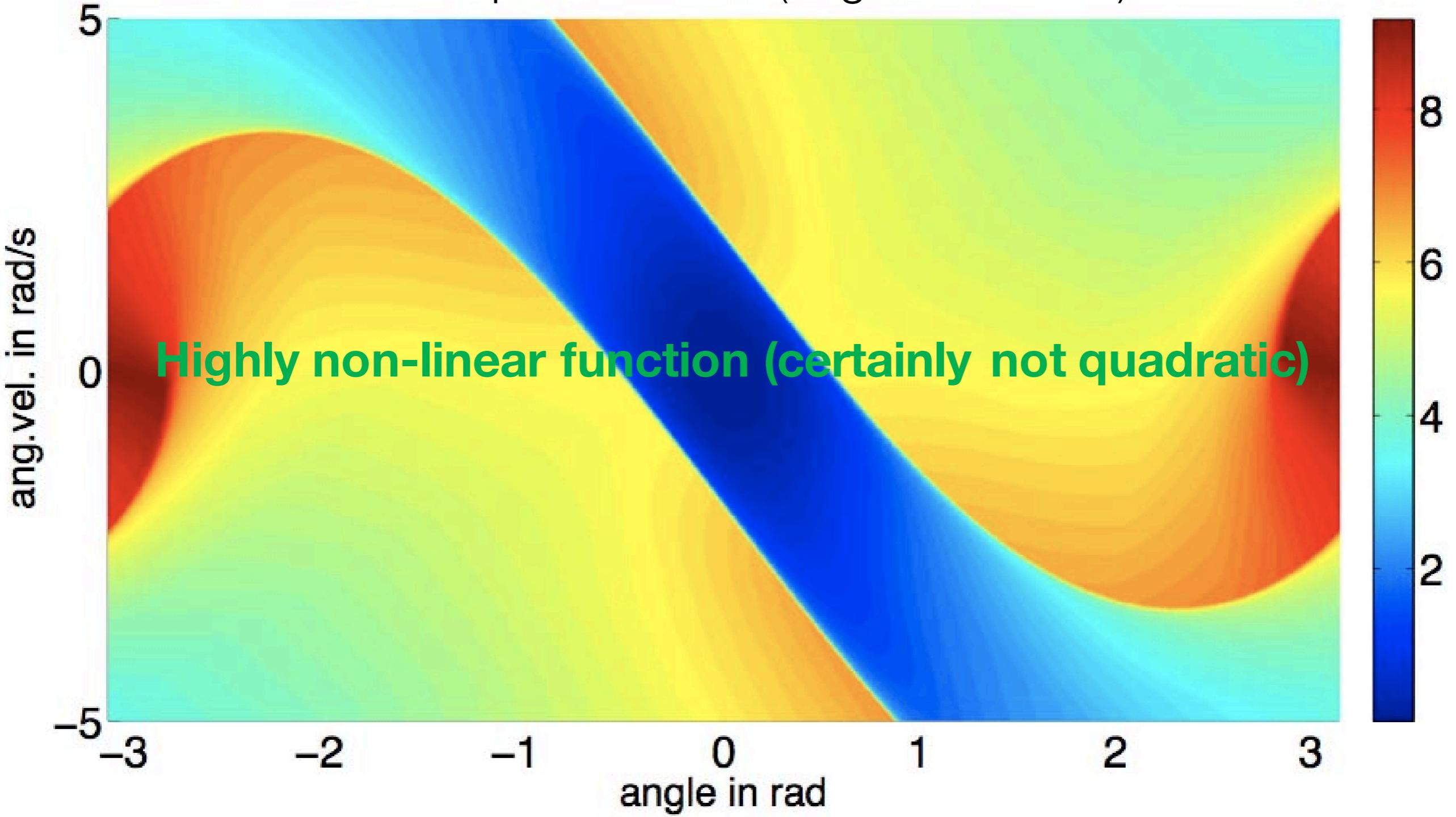
## Reward

$$r(\mathbf{s}, a) = -\mathbf{s}^T \text{diag}(1, 0.1)\mathbf{s} - 0.2a^2$$

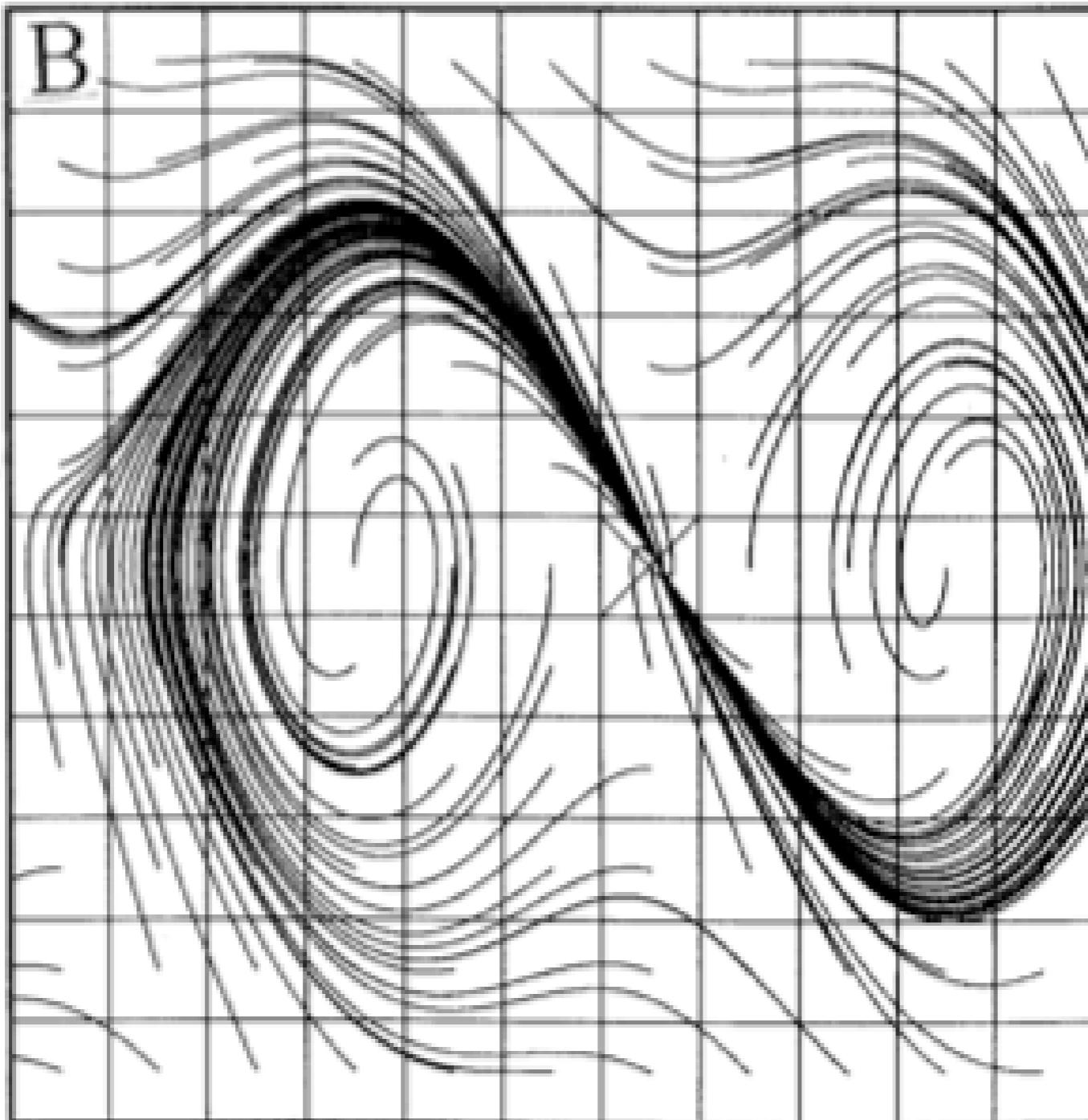


# Example: Value Function of Inverted Pendulum

Value function for the expected costs (negative reward)



# Possible: Learn Solutions only where needed!



**If you know places  
where we start...**

**... we can just look  
ahead and  
approximate the  
solution locally  
around an initial  
trajectory**



# Local Solutions by Linearizations

Every smooth function can be modeled with a Taylor expansion

$$f(\mathbf{x}) = f(\mathbf{a}) + \frac{df}{d\mathbf{x}} \Big|_{\mathbf{x}=\mathbf{a}} (\mathbf{x} - \mathbf{a}) + \frac{1}{2} (\mathbf{x} - \mathbf{a})^T \frac{d^2 f}{d\mathbf{x}^2} \Big|_{\mathbf{x}=\mathbf{a}} (\mathbf{x} - \mathbf{a}) + \dots$$

Hence, we can also **approximate the (learned) forward dynamics by linearizing** at the point  $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{u}}_t)$

$$\begin{aligned} \mathbf{x}_{t+1} &= f_t(\mathbf{x}_t, \mathbf{u}_t) \approx f(\tilde{\mathbf{x}}_t, \tilde{\mathbf{u}}_t) + \frac{df}{ds}(\mathbf{x}_t - \tilde{\mathbf{x}}_t) + \frac{df}{du}(\mathbf{u}_t - \tilde{\mathbf{u}}_t) \\ &= \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t \end{aligned}$$

with  $\mathbf{A}_t = \frac{df}{dx} \Big|_{x=\tilde{x}_t, u=\tilde{u}_t}$  and  $\mathbf{B}_t = \frac{df}{du} \Big|_{x=\tilde{x}_t, u=\tilde{u}_t}$

and  $\mathbf{b}_t = f(\tilde{\mathbf{x}}_t, \tilde{\mathbf{u}}_t) - \mathbf{A}_t \tilde{\mathbf{x}}_t - \mathbf{B}_t \tilde{\mathbf{u}}_t$



# Local Solutions by Linearizations

Similarly, we can **approximate the (learned) reward function by a second order approximation** at the point  $(\tilde{s}_t, \tilde{a}_t)$  (only shown for states )

$$\begin{aligned} r_t(s_t, a_t) &\approx r(\tilde{x}_t, \tilde{u}_t) + \frac{dr}{d\mathbf{x}}(\mathbf{x}_t - \tilde{\mathbf{x}}_t) + (\mathbf{x}_t - \tilde{\mathbf{x}}_t)^T \frac{dr}{d\mathbf{x}d\mathbf{x}}(\mathbf{x}_t - \tilde{\mathbf{x}}_t) - \mathbf{u}_t^T \mathbf{H}_t \mathbf{u}_t \\ &= -\mathbf{x}^T \mathbf{R}_t \mathbf{x} + 2\mathbf{r}_t^T \mathbf{x} - \mathbf{u}^T \mathbf{H}_t \mathbf{u} + \text{const} \end{aligned}$$

with  $\mathbf{R}_t = -\frac{dr}{d\mathbf{x}d\mathbf{x}}$  and  $\mathbf{r}_t = 0.5 \frac{dr}{d\mathbf{x}} - \frac{dr}{d\mathbf{x}d\mathbf{x}} \tilde{\mathbf{x}}_t$



# Local Solutions by Linearizations

So we are back to the **full linear optimal control case with...**

$$p_t(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = \mathcal{N}(\mathbf{x}_{t+1} | \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{b}_t, \Sigma_t)$$

$$r(\mathbf{x}, \mathbf{u}) = -\mathbf{x}^T \mathbf{R}_t \mathbf{x} + 2\mathbf{r}_t^T \mathbf{x} - \mathbf{u}^T \mathbf{H}_t \mathbf{u} + \text{const}$$

that we know how to solve...

Hence our algorithm for **solving non-linear optimal control** is...

**1. Backward Solution:** Compute optimal control law (i.e. Gains  $\mathbf{K}_t$  and offsets  $\mathbf{k}_t$ )

**2. Forward Propagation:** Run simulator with optimal control law to obtain linearization points  $(\tilde{\mathbf{x}}_{1:T}, \tilde{\mathbf{u}}_{1:T})$

1. If not converged, go to 1.

# Some interesting results (only in simulation)



Work by Emo Todorov  
and Yuval Tassa  
(They call basically the  
same algorithm  
incremental LQG, iLQG)

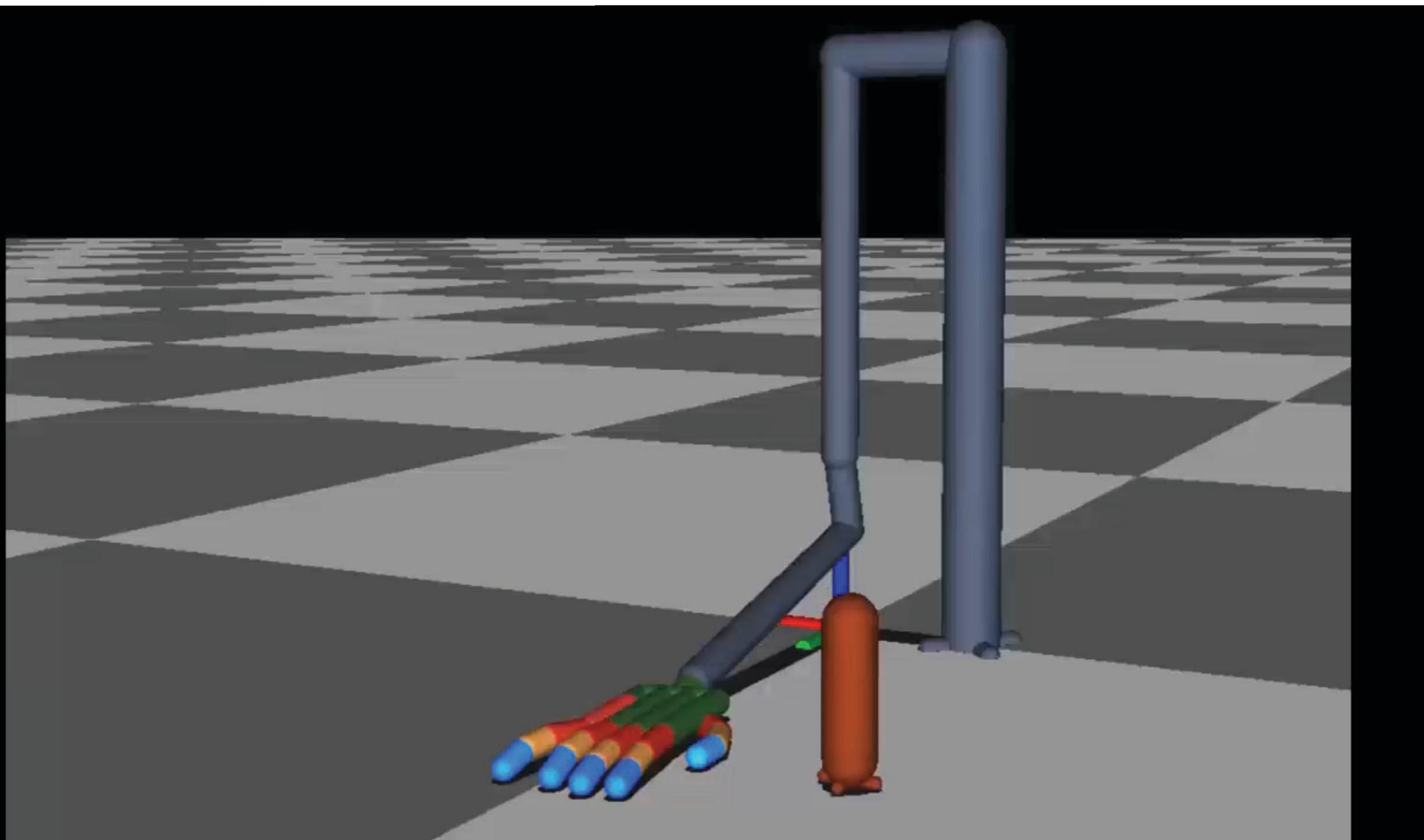
Synthesis of Complex Behaviors  
with  
Online Trajectory Optimization

(under review)

# Application to the Swing-Up



Some interesting results (only in simulation)



# Outline of the Lecture



- 1. Optimal Control**
- 2. Solving the Optimal Control for LQR systems**
- 3. Approximating Non-Linear Systems**
- 4. Optimal Control with Learned Models**
- 5. Final Remarks**



# Model Learning...

**Why does this work only in simulation?**

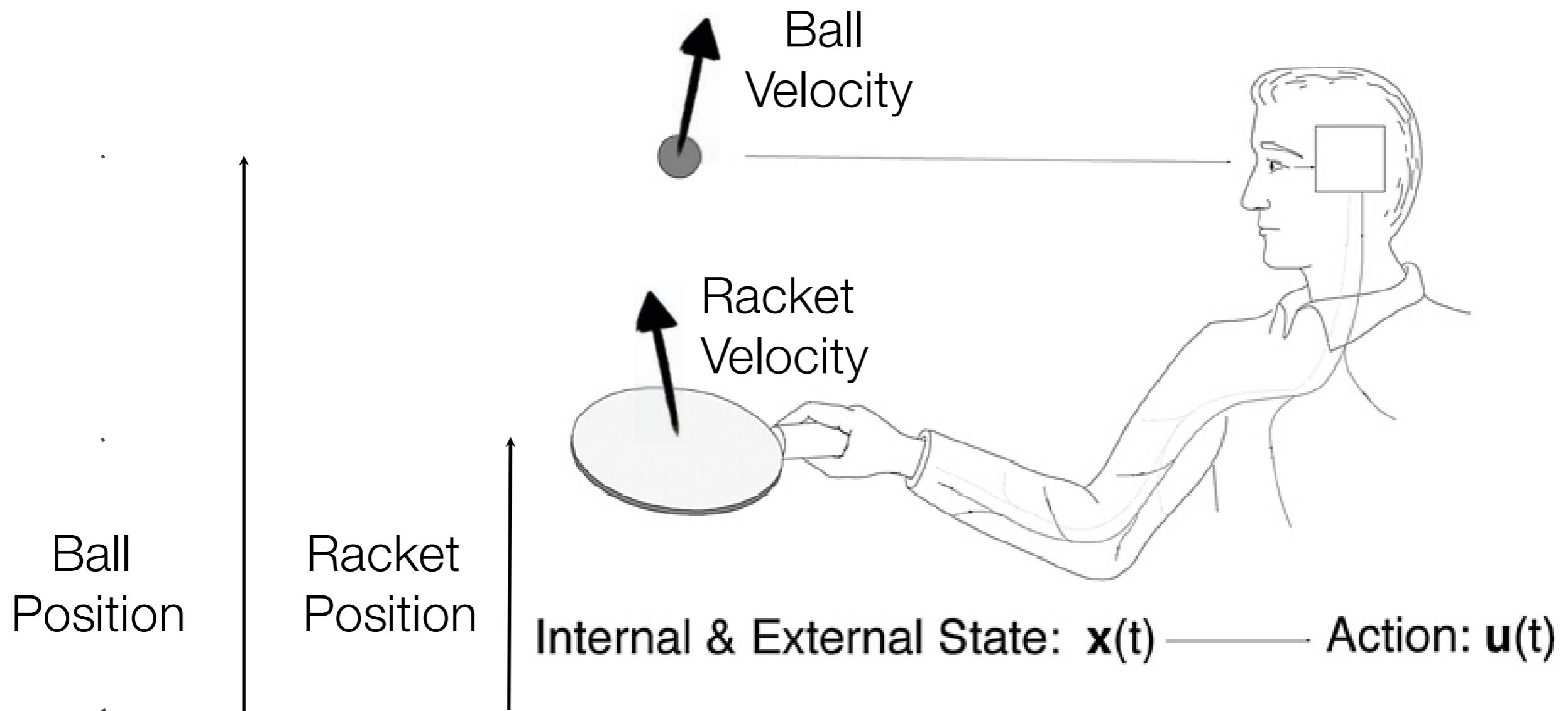
The models we have for such complex robots are... crap

We need **to learn the models !**

# Example: Ball Paddling



What are the states  $\mathbf{x}$ ?



# Example: Ball Paddling



**What are the actions  $u$ ?**

**All motor torques?**

If you do not have an inverse model ...

**Joint Accelerations?**

Perfect, if you have a good inverse model ...

Maybe identify the proper degrees of freedom?

**Accelerations in Task Space?**

**Ideally!**

... but only if you have a good operational space control law!



# Example: Ball Paddling



What are good rewards  $r$ ?

## Task knowledge or success/failure?

- For some algorithms rewards in  $\{1,0\}$  are perfect ...
- Real problems often require *reward shaping*...

## What's a good reward for our problem?

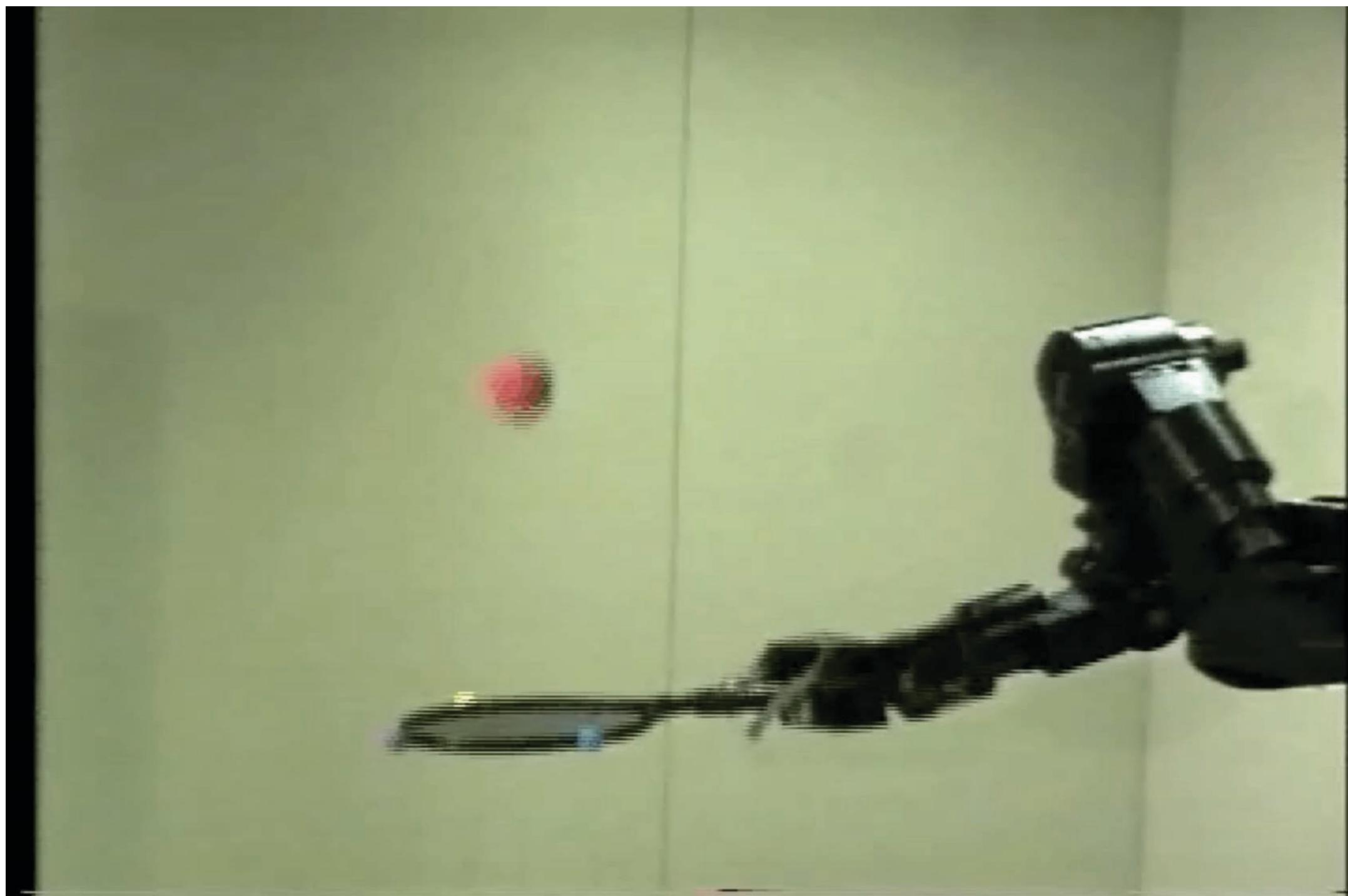
- Height of the ball?
- Distance between ball and the paddle?
- Ball needs to move in a certain region?
- All of the above?
- Additional punishments?



→ All of these together do the job!



Example: Real world application...

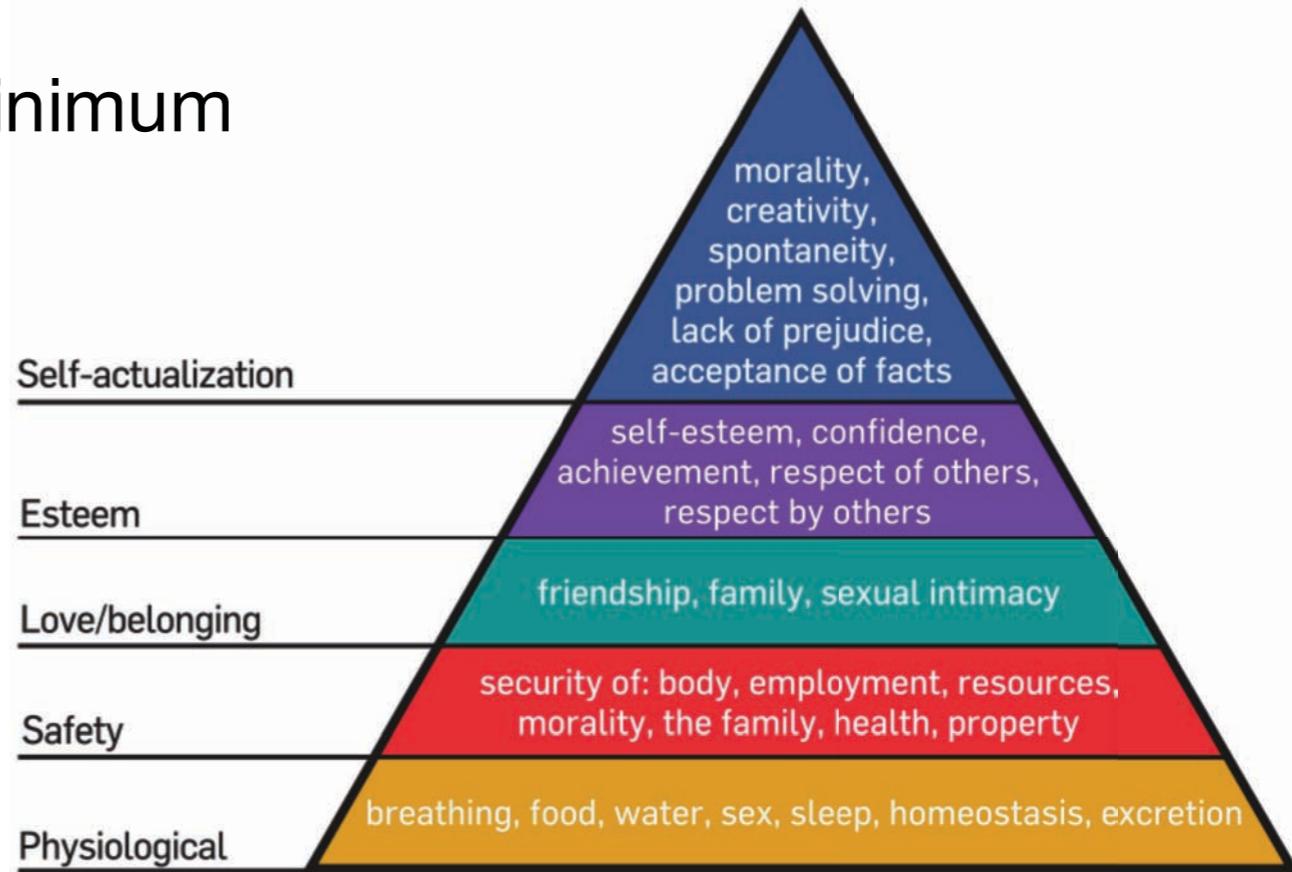




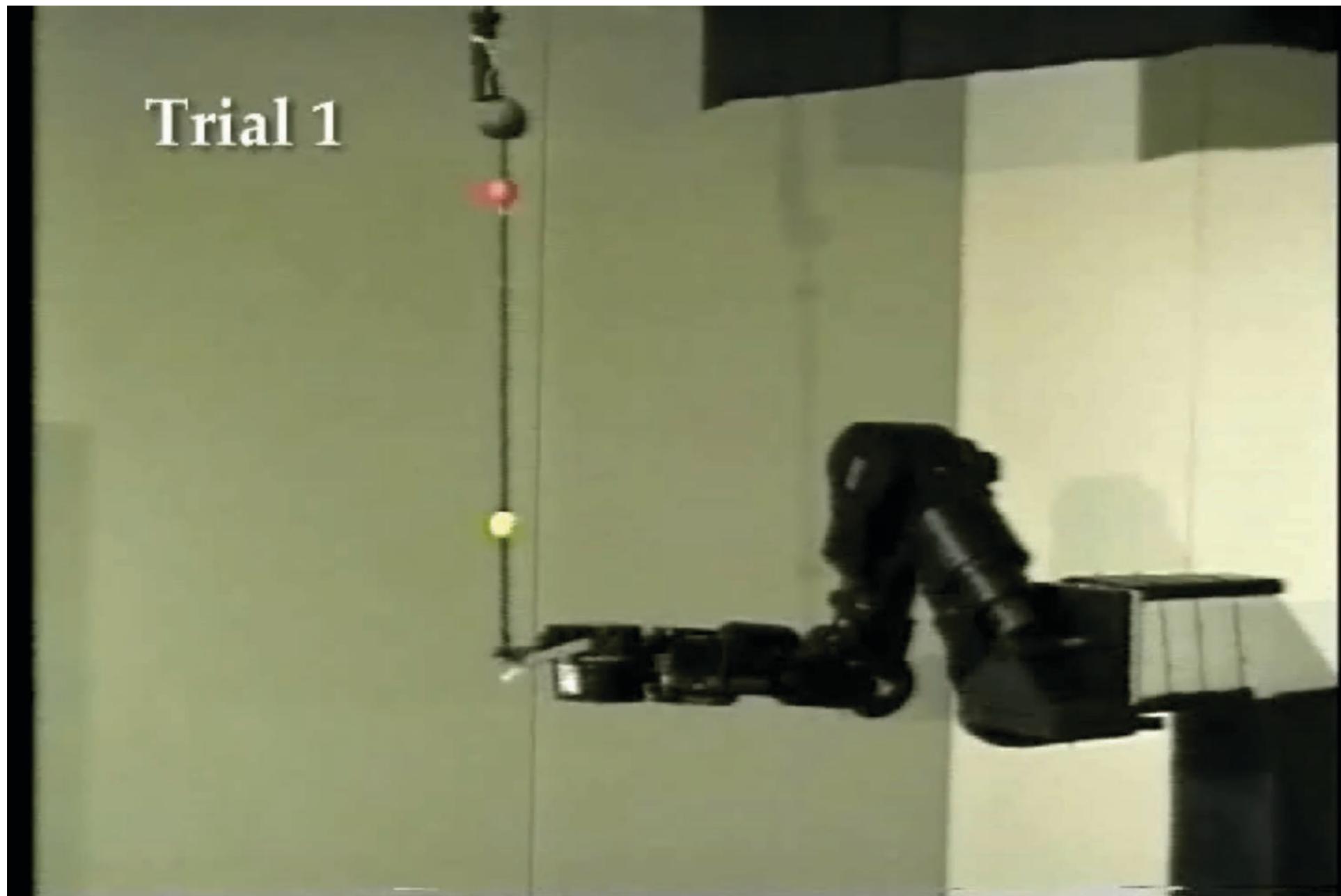
# Human Motor Cost Functions?

## Does this relate to Human Learning?

- **Maybe!** Many models from cognitive science are cost function based...
- *Reaching movements* can be explained by
  - Minimum jerk, Minimum torque change, Minimum end-point variance
- *Locomotion* can be explained by minimum metabolic energy consumption.
- Maslow's *Hierarchy of rewards* psychology ...

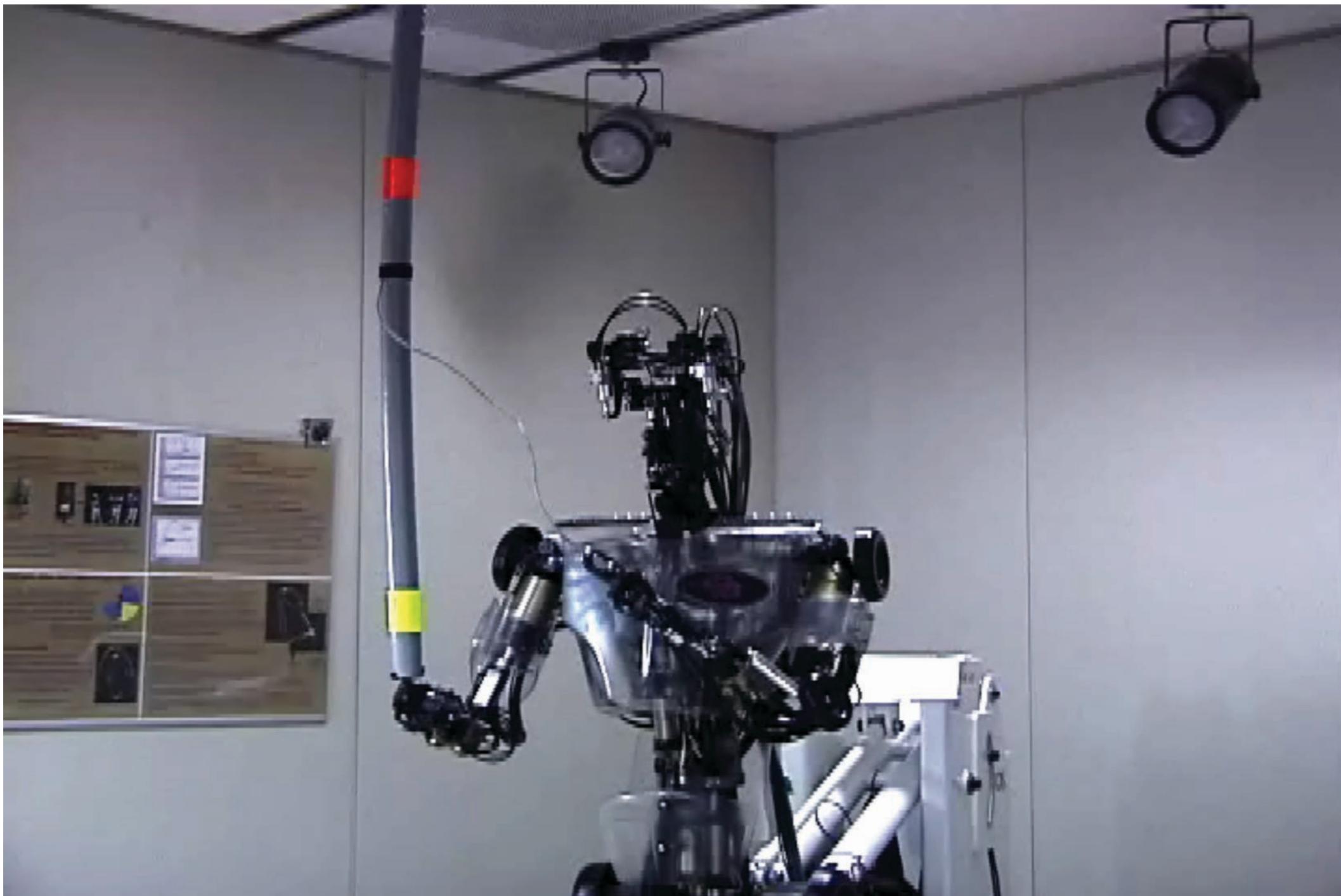


# Model Learning with subsequent Policy Optimization





# Model Learning with subsequent Policy Optimization



# Challenges in model-based policy learning



**This mainly works for balancing tasks where we live in a restricted state space**

**For more complex problems, using learned models becomes really hard**

- The model **is likely to be inaccurate**
- Inaccuracies could be exploited by optimizer such that the policy on the real system performs bad
- If we fully exploit an inaccurate model, we might jump into an area of the state space that we have not seen before
- **inherently unstable**



2 recent approaches...

## 1. PILCO (probabilistic inference for learning control)

Learn GP forward models

Use uncertainty of the GP-model for the long-term reward prediction

Policy Optimization with analytic gradient of expected reward

# Policy Optimization with PILCO



**Marc Peter Deisenroth, Carl Edward Rasmussen, Dieter Fox**  
**Learning to Control a Low-Cost Manipulator**  
**using Data-efficient Reinforcement Learning**



2 state of the art approaches...

## 2. Policy Search guided by trajectory optimization

Learn time dependent linear forward models

**Trajectory Optimization:** LQR like algorithm, additional constraint  
that new trajectory should stay close to the data → increase stability

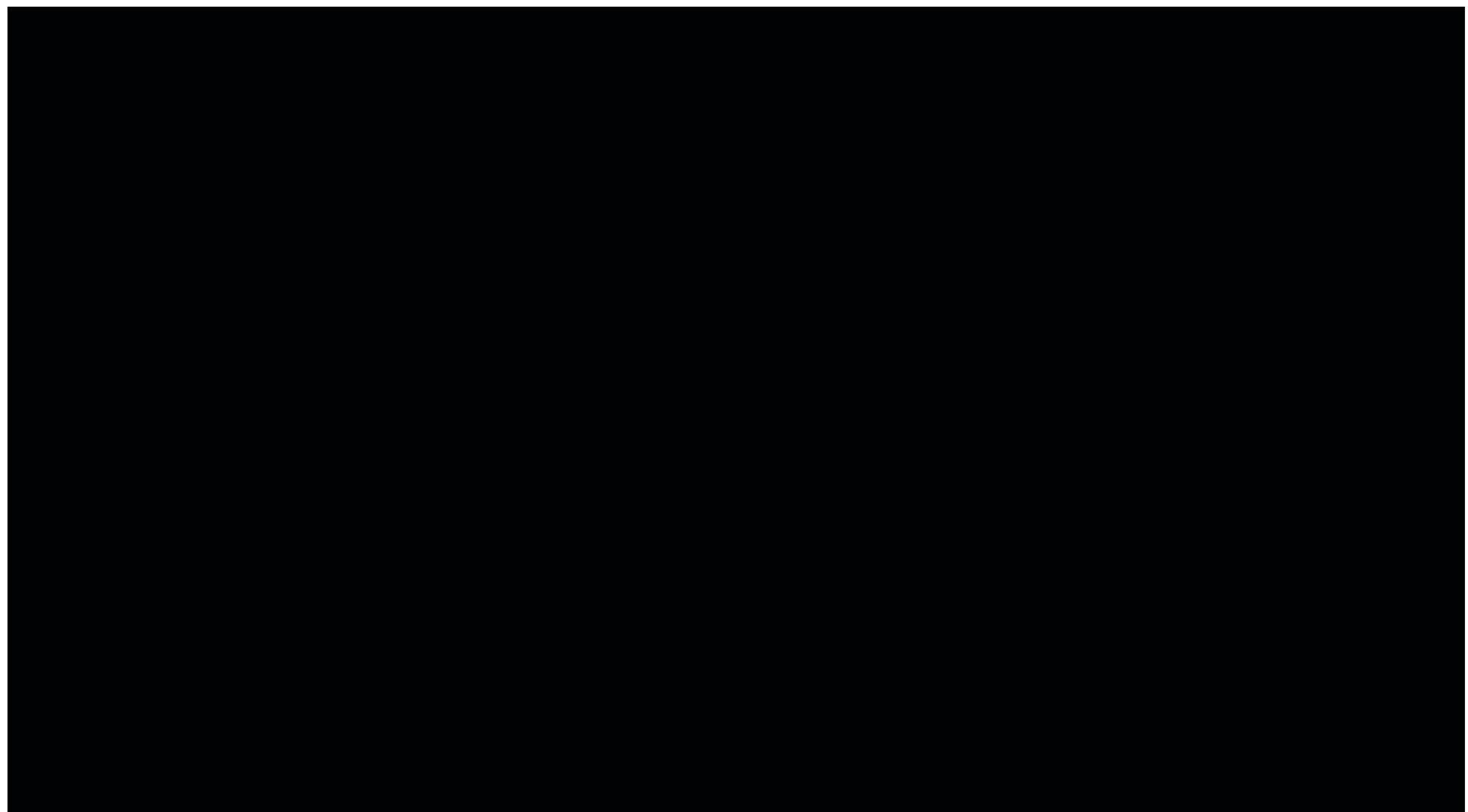
Use optimized trajectories to learn a generalizing neural network  
policy

# Guided Policy Search



Learning Neural Network  
Policies with Guided Policy Search  
under Unknown Dynamics

# Guided Policy Search





# Conclusions

- You have solved an **(stochastic) optimal control** problem today!
- Only two cases are solvable: **linear & discrete!**
- The **optimal policy** for a **LQG system** is a **time-varying linear feedback** controller
- Linearizations can be problematic → lead to **oscillations** (but can be made more stable)
- Works well if the system is **not too non-linear and model can be learned accurately!**
- We will continue with **Value Function** and **Policy Search Methods.**