

# Optimal and Learning Control

for

# Autonomous Robots

## Lecture 12



**A D R L**

Jonas Buchli  
Agile & Dexterous Robotics Lab



# Class logistics

Exercise 3 hand out today!

Due: Tue. May 26 - 18.00

Interviews: Thu/Fri. May 28/29

<https://ethz.doodle.com/bsi7gvkycvrmht6t>



# L11 Recap



# Forward diffusion / sampling

## Importance sampling

## Function approximation



# Linear Markov Decision Process

Three conditions on the optimal control problem:

1) Quadratic control cost

$$J = E \left\{ \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u} dt \right\}$$

2) Control affine system

$$d\mathbf{x} = \mathbf{f}(t, \mathbf{x})dt + \mathbf{g}(t, \mathbf{x}) (\mathbf{u}dt + d\mathbf{w}), \quad d\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma dt)$$

$$\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x}) (\mathbf{u} + \varepsilon), \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$$

3)  $\mathbf{R}\Sigma = \lambda \mathbf{I}$



# Linear Markov Decision Process (cnt)

nonlinear PDE

$$-\partial_t V^* = q - \frac{1}{2} \nabla_x^T V^* \Xi \nabla_x V^* + \nabla_x^T V^* \mathbf{f} + \frac{\lambda}{2} \text{Tr}[\nabla_{xx} V^* \Xi]$$



$$V^*(t, \mathbf{x}) = -\lambda \log \Psi(t, \mathbf{x})$$

$$-\partial_t \Psi = -\frac{1}{\lambda} q \Psi + \mathbf{f}^T \nabla_x \Psi + \frac{\lambda}{2} \text{Tr}[\Xi \nabla_{xx} \Psi]$$

$$-\partial_t \Psi = \mathbf{H}[\Psi] \quad \mathbf{H} = -\frac{1}{\lambda} q + \mathbf{f}^T \nabla_x + \frac{\lambda}{2} \text{Tr}[\Xi \nabla_{xx}]$$

$$\Psi(t_f, \mathbf{x}) = \exp\left(-\frac{1}{\lambda} \Phi(\mathbf{x})\right)$$

Final Value problem

$$\mathbf{g} \Sigma \mathbf{g}^T = \lambda \Xi$$

The effective Covariance



# Path Integral

$$\Psi(s, \mathbf{y}) = \int \mathbb{P}_{uc}(\tau \mid s, \mathbf{y}) e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}_N) + \sum_{i=0}^{N-1} q(t_i, \mathbf{x}(t_i)) dt \right)} d\mathbf{x}(t_1) \dots d\mathbf{x}(t_{N-1}) d\mathbf{x}_N$$

Equivalently

$$\begin{aligned} \Psi(s, \mathbf{y}) &= \mathbb{E}_{\tau_{uc}} \left\{ e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_N)) + \sum_{i=0}^{N-1} q(t_i, \mathbf{x}(t_i)) dt \right)} \right\} \\ &= \mathbb{E}_{\tau_{uc}} \left\{ e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)} \right\} \end{aligned}$$

Samples can be generated by

$$d\mathbf{x} = \mathbf{f}(t, \mathbf{x})dt + \mathbf{g}(t, \mathbf{x})d\mathbf{w}, \quad d\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma dt), \quad \mathbf{x}(t = s) = \mathbf{y}$$



# Closer look at Path Integral formula

- For calculating the Desirability function at each point

$$\Psi(s, \mathbf{y}) = E_{\tau_{uc}} \left\{ e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)} \right\}$$

$$d\mathbf{x} = \mathbf{f}(t, \mathbf{x})dt + \mathbf{g}(t, \mathbf{x})d\mathbf{w}, \quad d\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma dt), \quad \mathbf{x}(t = s) = \mathbf{y}$$

- 1) Forward simulate the uncontrolled system from  $(s, \mathbf{y})$  up to  $t_f$
- 2) Integrate the cost over the generated path



# Path Integral: Optimal Control

- Using the white noise formulation  $\varepsilon = \frac{d\mathbf{w}}{dt}$

$$\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x})\varepsilon, \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \Sigma), \quad \mathbf{x}(t = s) = \mathbf{y}$$

$$\mathbf{u}^*(s, \mathbf{y}) = \frac{\mathbb{E}_{\tau_{uc}} \left\{ \varepsilon e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)} \right\}}{\mathbb{E}_{\tau_{uc}} \left\{ e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)} \right\}}$$

# Example: Naive sampling

$$M(\theta) \cdot \ddot{\theta} + C(\theta, \dot{\theta}) = \tau$$

$$\ddot{\theta} = M(\theta)^{-1} \cdot (-C(\theta, \dot{\theta}) + \tau)$$

$$M(\theta) = \begin{pmatrix} d_1 + 2d_2 \cos(\theta_2) & d_3 + d_2 \cos(\theta_2) \\ d_3 + d_2 \cos(\theta_2) & d_3 \end{pmatrix} \quad (37)$$

$$C(\dot{\theta}, \theta) = \begin{pmatrix} -\dot{\theta}_2 (2\dot{\theta}_1 + \dot{\theta}_2) \\ \dot{\theta}_1^2 \end{pmatrix} d_2 \sin(\theta_2) \quad (38)$$

$$d_1 = I_1 + I_2 + m_2 l_1^2, \quad d_2 = m_2 l_1 s_2, \quad d_3 = I_2 \quad (39)$$

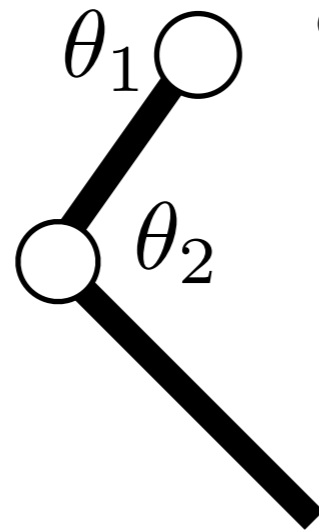
Compare to affine control assumption

$$\dot{x} = \Phi(x) + G(x) \cdot \tau$$

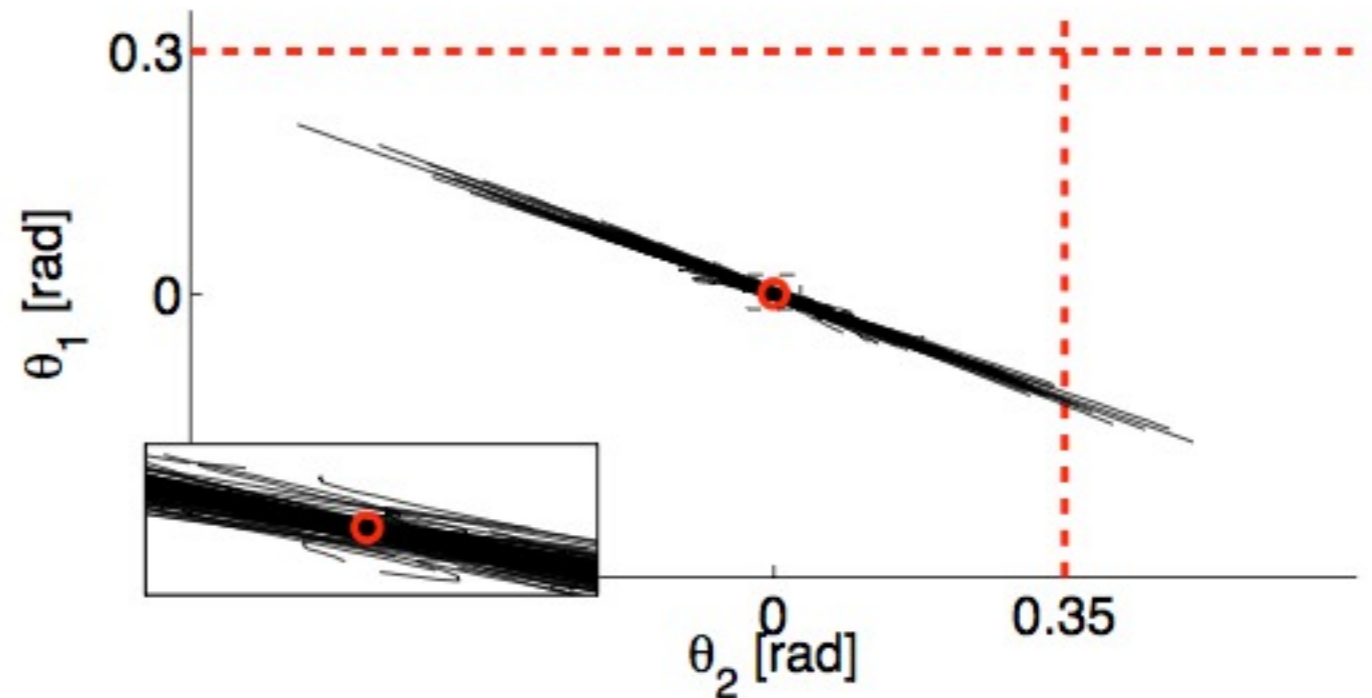
$$x = \begin{pmatrix} \theta_1 & \theta_2 & \dot{\theta}_1 & \dot{\theta}_2 \end{pmatrix}^T \quad \tau = \begin{pmatrix} \tau_1 & \tau_2 \end{pmatrix}$$

$$\Phi(x) = \begin{pmatrix} \dot{\theta}_1 \\ \dot{\theta}_2 \\ -M(\theta)^{-1} \cdot C(\theta, \dot{\theta}) \end{pmatrix} \quad G(x) = \begin{pmatrix} O_{2 \times 2} \\ M(\theta)^{-1} \end{pmatrix}$$

Symbol	Value	Unit
$m_1$	1.4	Kg
$m_2$	1	Kg
$s_1$	0.11	m
$s_2$	0.16	m
$I_1$	0.3	Kg m <sup>2</sup>
$I_2$	0.33	Kg m <sup>2</sup>
$l_1$	0.025	m
$l_2$	0.045	m



Path integral SOC requires sampling of passive dynamics with gaussian mean-free noise



# Improved sampling...?

i.d. - controller

$$\tau_u^i = M(\theta) \cdot (\alpha_i + \epsilon_i) + C(\theta, \dot{\theta})$$

$$M(\theta) \cdot \ddot{\theta} + C(\theta, \dot{\theta}) = \tau_u$$

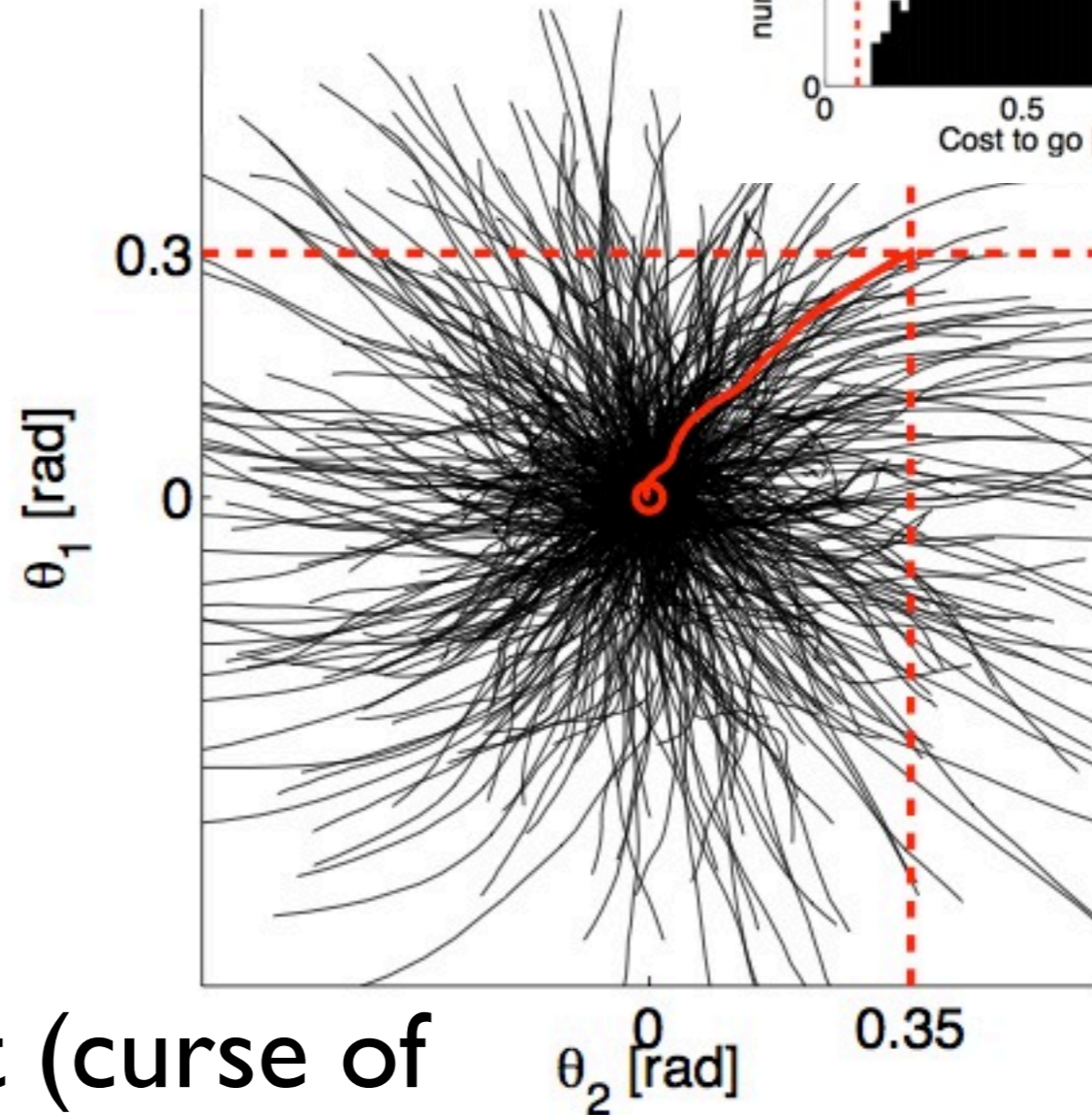
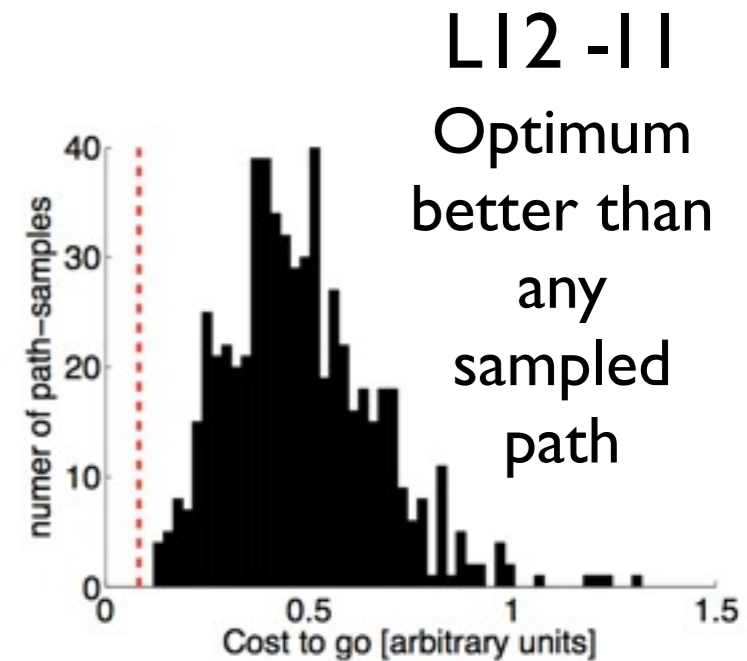
plant dynamics

Sample in acceleration space, use inverse dynamics controllers to find torques:

... still not very efficient (curse of dimensionality still strikes, needle in a haystack!)



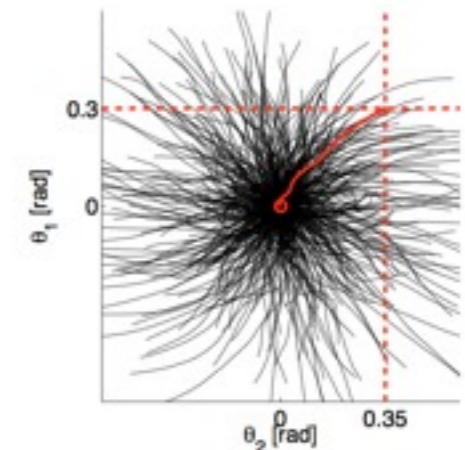
Buchli - OLCAR - 2013



# Path Integral: issues (1)

- Inefficient sampling

$$\mathbf{u}^*(s, \mathbf{y}) = \mathbb{E}_{\tau_{uc}} \left\{ \varepsilon \frac{e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)}}{\mathbb{E}_{\tau_{uc}} \left\{ e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)} \right\}} \right\}$$



Soft Max

It just has significant value for near optimal solution

What are the chances to hit the optimal solution by a random walk?

## Importance Sampling

# Path Integral: issues (2)

- Point-wise estimation of the optimal controls

$$\mathbf{u}^*(s, \mathbf{y}) = \mathbb{E}_{\tau_{uc}} \left\{ \frac{e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)}}{\mathbb{E}_{\tau_{uc}} \left\{ e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)} \right\}} \right\}$$

The optimal control is estimated independently for each point

Does the optimal control change drastically from one point to the other?

**Function Approximation**



# Reward weighting

$$\mathbf{u}^*(s, \mathbf{y}) = \mathbb{E}_{\tau_{uc}} \left\{ \varepsilon \frac{e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)}}{\mathbb{E}_{\tau_{uc}} \left\{ e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)} \right\}} \right\}$$

$$\alpha(\tau_{uc}; s, \mathbf{y}) = \frac{e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)}}{\mathbb{E}_{\tau_{uc}} \left\{ e^{-\frac{1}{\lambda} \left( \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) dt \right)} \right\}}$$

$$\boldsymbol{\theta}_i^* = \boldsymbol{\theta}_{i,c} + \operatorname{argmin}_{\Delta \boldsymbol{\theta}_i} \int \frac{e^{-\frac{1}{\lambda} R(\tau; s, \mathbf{y})}}{\mathbb{E}_{\tau_c} \left\{ e^{-\frac{1}{\lambda} R(\tau; s, \mathbf{y})} \right\}} \|\boldsymbol{\Upsilon}_i^T(s, \mathbf{y}) \Delta \boldsymbol{\theta}_i - \varepsilon\|_2^2 \mathbb{P}_{\tau_c}(\tau | s, \mathbf{y}) p(s, \mathbf{y}) d\tau d\mathbf{y} ds$$

# EOF Recap



# L12





# Lecture 12 Goals

- ★ Policy Improvement with Path Integrals - PI2
- ★ Combination of optimal and learning control



# General Path Integral Algorithm

## Importance Sampling

Goal: Sample more efficiently

Update sampling distribution to account what is already known about good solutions

Use current estimate of optimal controls for sampling

## Function approximation

Goal: Reduce complexity by generalizing

Function approximation reduces open parameters

Each parameter covers a neighborhood



---

**Algorithm 8** General Path Integral Algorithm
 

---

**given**

The cost function:

$$J = \Phi(\mathbf{x}(t_f)) + \int_s^{t_f} (q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u}) dt$$

A PDF defining the quality of approximation of optimal control at each time-state pair:  $p(t, \mathbf{x})$

An initial policy and a Linear Model:  $\mathbf{u}(t, \mathbf{x}) = [u_i(t, \mathbf{x})] = [\Upsilon_i^T(t, \mathbf{x}) \boldsymbol{\theta}_i]$

**repeat**

(a) Randomly choose a time-state pair from  $p(t, \mathbf{x})$ :  $(s, \mathbf{y})$

(b) Forward simulate the controlled system for  $K$  different rollouts:  $\{\tau^k\}_{k=1}^K$  sampling repeatedly

$$\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x}) (\mathbf{u} + \boldsymbol{\varepsilon})$$

$\mathbf{u}(t, \mathbf{x}) = [\Upsilon_i^T(t, \mathbf{x}) \boldsymbol{\theta}_i]$ ,  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ ,  $\mathbf{x}(t=s) = \mathbf{y}$  using current estimate of opt. control to sample

(c) Calculate the return for each rollout:  $\{R^k\}_{k=1}^K$

$$R(\tau; s, \mathbf{y}) = \Phi(\mathbf{x}(t_f)) + \int_s^{t_f} (q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u}) dt + \int_s^{t_f} \mathbf{u}^T \mathbf{R} d\mathbf{w}$$

(d) Calculate  $\{\alpha^k\}_{k=1}^K$

$$\alpha^k(s, \mathbf{y}) = \exp(-\frac{1}{\lambda} R^k) / \frac{1}{K} \sum_{j=1}^K \exp(-\frac{1}{\lambda} R^j)$$

weighting by reward  
reward weighted regression

(e) Solve the following linear regression problem for each control input  $i$ :

$$\Delta \boldsymbol{\theta}_i = \operatorname{argmin} \sum_{k=1}^K \alpha^k \|\Upsilon_i^T(s, \mathbf{y}) \Delta \boldsymbol{\theta}_i - \varepsilon_i^k(s)\|_2^2$$

(f) Update the parameter vector for each control input  $i$ :

$$\boldsymbol{\theta}_i \leftarrow \boldsymbol{\theta}_i + \omega \Delta \boldsymbol{\theta}_i$$

**until** convergence

---

# Policy improvement with Path Integrals - PI<sup>2</sup>

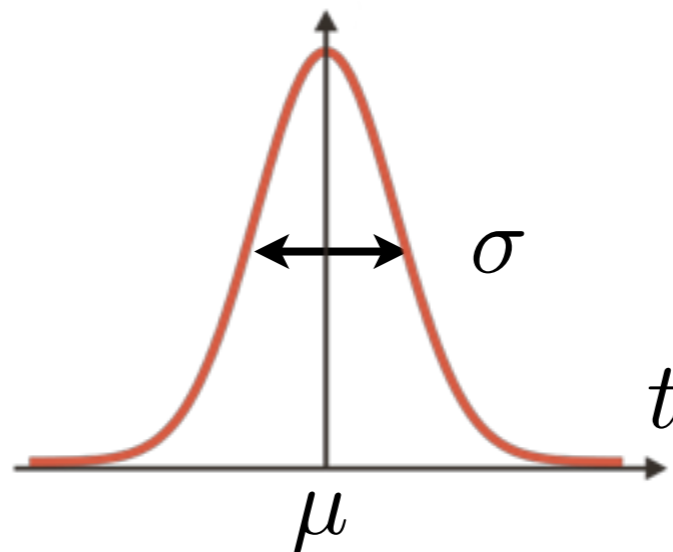
In contrast to the General Path Integral Algorithm, PI<sup>2</sup> assumes that samples are extracted over the entire time horizon, not only for a single time step. Therefore the regression problem should be over different time steps as well as the rollout batch.



# PI<sup>2</sup> w. time dependent policy

$$u_i(t) = \mathbf{r}^T(t) \boldsymbol{\theta}_i$$

$$\mathbf{r}(t) = [\mathbf{r}_n(t)]_{N \times 1} = \begin{bmatrix} e^{-\frac{1}{2} \frac{(t-\mu_n)^2}{\sigma_n^2}} \end{bmatrix}_{N \times 1}$$



# PI<sup>2</sup> Assumptions

1. The weighting function of MSE,  $p(t, \mathbf{x})$ , is assumed to be the probability distribution of the state under the latest estimation of the optimal control. Therefore the rollout trajectory's states are considered to be extracted from  $p(t, \mathbf{x})$ .
2. The return function,  $R(\tau; s, \mathbf{y})$ , won't be a function of state, if the system has been initialized from a similar initial condition. An immediate result of this assumption is that  $\alpha^k(s, \mathbf{y})$  is only a function of time i.e.  $\alpha^k(s)$ . Therefore for a batch of trajectories extracted from a similar initial condition,  $\alpha^k(s)$  can be estimated for each time step by the following

$$\alpha^k(s) = \frac{\exp(-\frac{1}{\lambda}R^k(s))}{\frac{1}{K} \sum_{j=1}^K \exp(-\frac{1}{\lambda}R^j(s))} \quad (3.60)$$

# PI<sup>2</sup> Assumptions (cont'd)

3. The basis function vector for all the control inputs is the same. Hence we will drop the  $i$  subscription of the basis function as  $[u_i(t)] = [\mathbf{r}^T(t)\boldsymbol{\theta}_i]$ .
4. Instead of adding noise to the control input, the noise is added directly to the parameter vector. Therefore the input noise will be as  $\varepsilon_i = \mathbf{r}(t)\boldsymbol{\epsilon}_i$ , where  $\boldsymbol{\epsilon}_i$  is the noise that is added to the parameter vector of the  $i$ th control input.
5. The PI2 regression problem should be modified as follows

$$\Delta\boldsymbol{\theta}_i = \operatorname{argmin} \sum_{s=t_0}^{t_f} \sum_{k=1}^K \alpha^k \|\mathbf{r}_i^T(s)\Delta\boldsymbol{\theta}_i - \varepsilon_i^k(s)\|_2^2$$

In contrast to the General Path Integral Algorithm, PI2 assumes that samples are extracted over the entire time horizon, not only for a single time step. Therefore the regression problem should be over different time steps as well as the rollout batch.



# PI<sup>2</sup> - Regression step

In order to solve this regression problem, PI2 breaks it into two separate optimizations. This method finds the optimal solution as long as we can assume the regression error has a zero mean over the samples. In this method first the optimization is solved for each time step separately. Therefore the first optimization will find a time-dependent parameter vector increment that has the minimum error over the rollouts at each time step. Finally in the second optimization, we will find a parameter vector increment that approximates the time-dependent one.

## I) min. error at each time step

The first optimization is defined as follows for each time step  $s$

$$\Delta\theta_i^*(s) = \operatorname{argmin} \sum_{k=1}^K \alpha^k(s) \|\mathbf{r}_i^T(s) \Delta\theta_i - \varepsilon_i^k(s)\|_2^2$$





# PI<sup>2</sup> - Regression step

(cont'd)

The first optimization is defined as follows for each time step  $s$

$$\Delta\boldsymbol{\theta}_i^*(s) = \operatorname{argmin} \sum_{k=1}^K \alpha^k(s) \|\boldsymbol{\Upsilon}_i^T(s) \Delta\boldsymbol{\theta}_i - \boldsymbol{\epsilon}_i^k(s)\|_2^2$$

$$\Delta\boldsymbol{\theta}_i^*(s) = \sum_{k=1}^K \alpha^k(s) \frac{\boldsymbol{\Upsilon}_i(s)}{\boldsymbol{\Upsilon}_i^T(s) \boldsymbol{\Upsilon}_i(s)} \boldsymbol{\epsilon}_i^k(s)$$

$$\Delta\boldsymbol{\theta}_i^*(s) = \sum_{k=1}^K \alpha^k(s) \frac{\boldsymbol{\Upsilon}_i(s) \boldsymbol{\Upsilon}_i^T(s)}{\boldsymbol{\Upsilon}_i^T(s) \boldsymbol{\Upsilon}_i(s)} \boldsymbol{\epsilon}_i^k(s)$$

# PI<sup>2</sup> - Regression step

(cont'd)

## 2) parameter vect. increment approximating time dependent update

The second optimization for finding the optimal  $\Delta\theta_i^*$  is defined in equation (3.65). The index  $n$  refers to the  $n$ th element of the vector  $\Delta\theta_i^* = [\Delta\theta_{i,n}^*]$

$$\Delta\theta_{i,n}^* = \operatorname{argmin}_{\Delta\theta_{i,n}} \sum_{s=t_0}^{t_f} (\Delta\theta_{i,n} - \Delta\theta_{i,n}^*(s))^2 \Upsilon_n(s) \quad (3.65)$$

$\Upsilon_n(t)$  is the  $n$ th element of the basis function vector  $\Upsilon(t)$

$$\Delta\theta_{in}^* = \frac{\sum_{t=t_0}^{t_f} \Delta\theta_{in}^*(s) \Upsilon_n(s)}{\sum_{t=t_0}^{t_f} \Upsilon_n(s)}$$



$$\Delta\theta_{in}^* = \frac{\sum_{t=t_0}^{t_f} \Delta\theta_{in}^*(s) \Upsilon_n(s)}{\sum_{t=t_0}^{t_f} \Upsilon_n(s)}$$

Compact notation: Use element-wise multiplication  
replace sum w. integral

$$\Delta\theta_i^* = \left( \int_{t_0}^{t_f} \Delta\theta_i^*(s) \circ \Upsilon(s) ds \right) \cdot / \int_{t_0}^{t_f} \Upsilon(s) ds$$

where  $\circ$  and  $\cdot /$  are element-wise multiplication and division.

**Algorithm 9** PI2 Algorithm for time-dependent policy**given**

The cost function:

$$J = \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} (q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u}) dt$$

A Linear Model for function approximation:  $\mathbf{u}(t) = [u_i(t)] = [\mathbf{\Upsilon}^T(t) \boldsymbol{\theta}_i]$ Initialize  $[\boldsymbol{\theta}_i]$  with a sophisticated guessInitialize exploration noise standard deviation:  $c$ **repeat**Create  $K$  rollouts of the system with the perturbed parameter  $[\boldsymbol{\theta}_i] + [\boldsymbol{\epsilon}_i]$ ,  $\boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, c^2 \mathbf{I})$ **for** the  $i$ th control input **do****for** each time,  $s$  **do**Calculate the Return from starting time  $s$  for the  $k$ th rollout:

$$R(\tau^k(s)) = \Phi(\mathbf{x}(t_f)) + \int_s^{t_f} (q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u}) dt$$

Calculate  $\alpha$  from starting time  $s$  for the  $k$ th rollout:

$$\alpha^k(s) = \frac{\exp(-\frac{1}{\lambda} R(\tau^k(s)))}{\sum_{k=1}^K \exp(-\frac{1}{\lambda} R(\tau^k(s)))}$$

time dependent parameter increment

Calculate the time varying parameter increment  $\Delta \boldsymbol{\theta}_i(s)$ :

$$\Delta \boldsymbol{\theta}_i(s) = \sum_{k=1}^K \alpha^k(s) \frac{\boldsymbol{\Upsilon}(s) \boldsymbol{\Upsilon}^T(s)}{\boldsymbol{\Upsilon}^T(s) \boldsymbol{\Upsilon}(s)} \boldsymbol{\epsilon}_i^k(s)$$

**end for**

Time-averaging the parameter vector

$$\Delta \boldsymbol{\theta}_i = \left( \int_{t_0}^{t_f} \Delta \boldsymbol{\theta}_i(s) \circ \boldsymbol{\Upsilon}(s) ds \right) \cdot \left/ \int_{t_0}^{t_f} \boldsymbol{\Upsilon}(s) ds \right.$$

time averaging of increment

Update parameter vector for control input  $i$ ,  $\boldsymbol{\theta}_i$ :

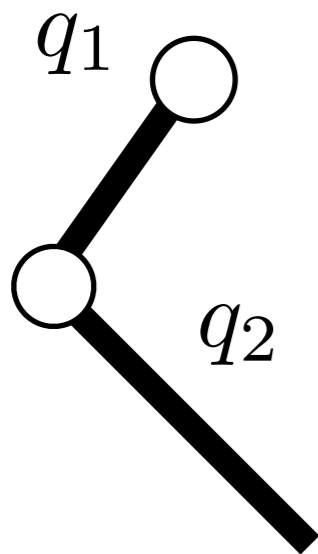
$$\boldsymbol{\theta}_i \leftarrow \boldsymbol{\theta}_i + \omega \Delta \boldsymbol{\theta}_i$$

**end for**- Decrease  $c$  for noise annealing**until** maximum number of iterations

# Simple example

Assume position PID  
controller in each joint

$$\tau_i = -K_{p,i}(q_i - q_{i,des}) - K_{d,i}(\dot{q}_i - \dot{q}_{i,des})$$

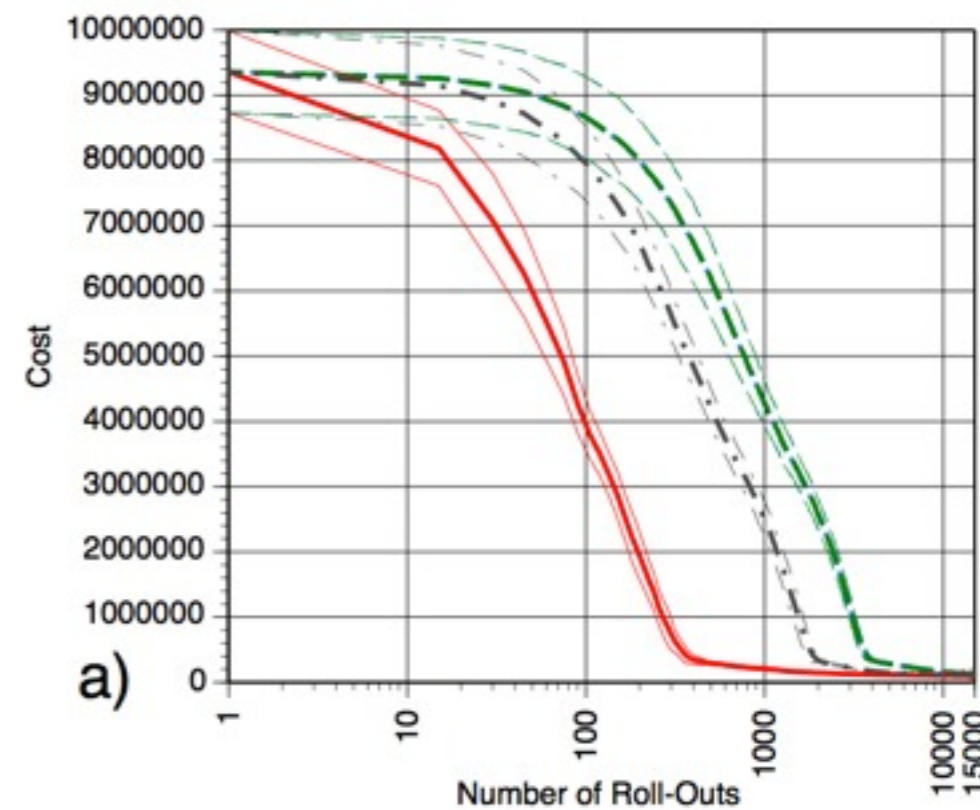
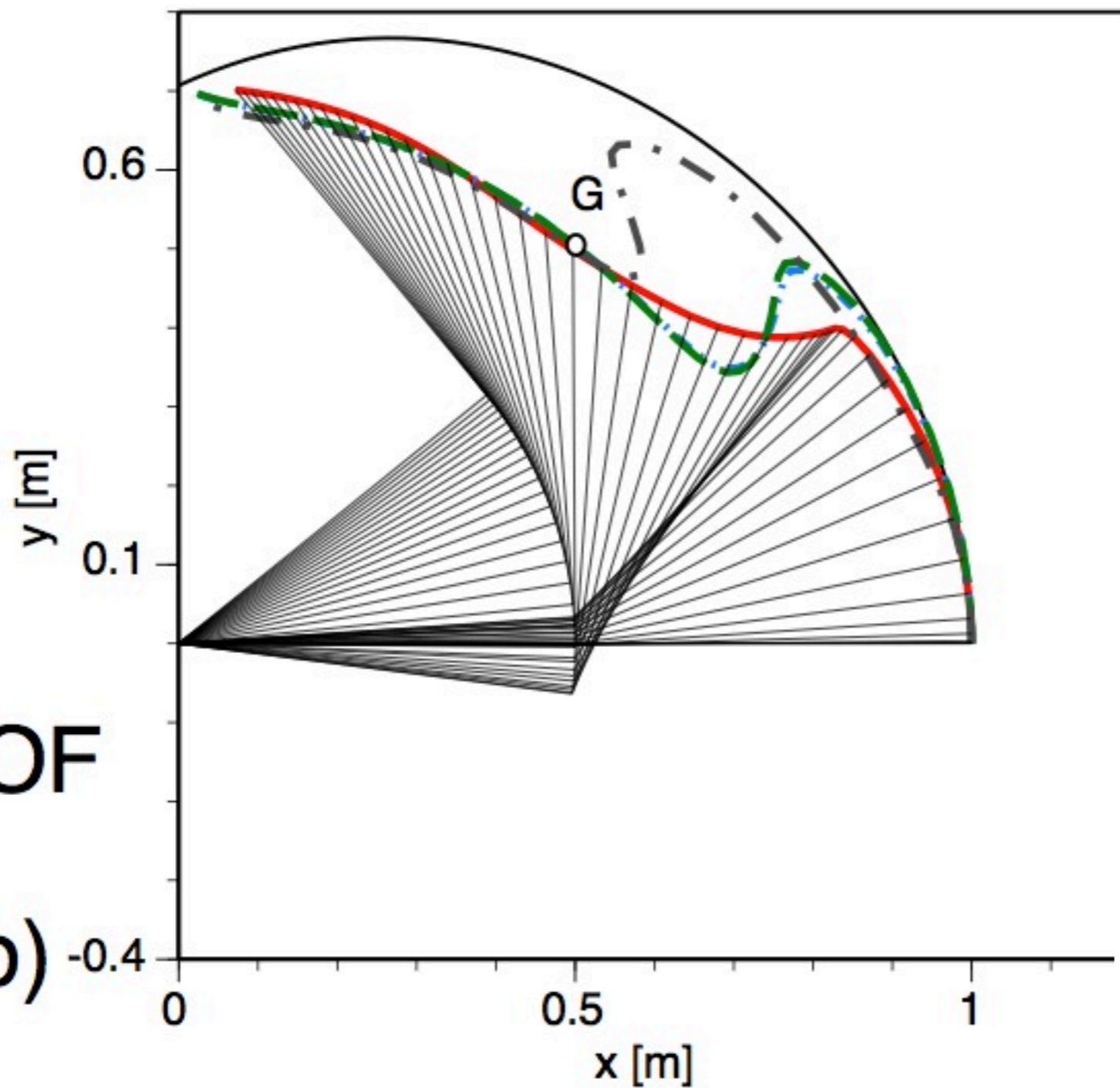


Parametrize learning problem via  
reference trajectories

$$q_{des,i}(t) = \Upsilon(t)^T \theta_{q,i}$$

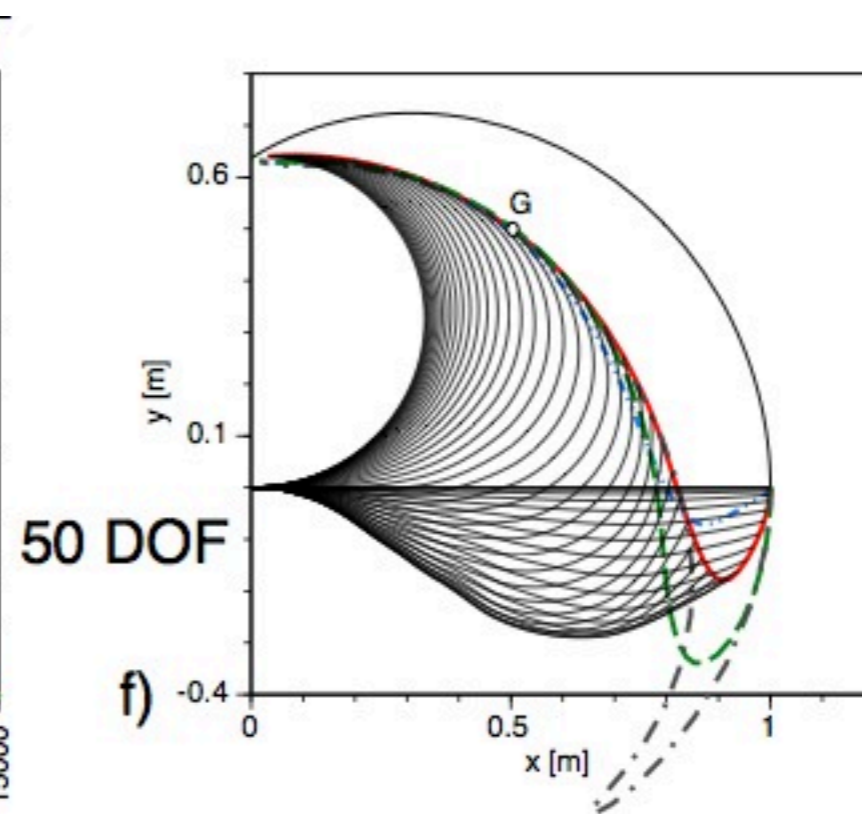
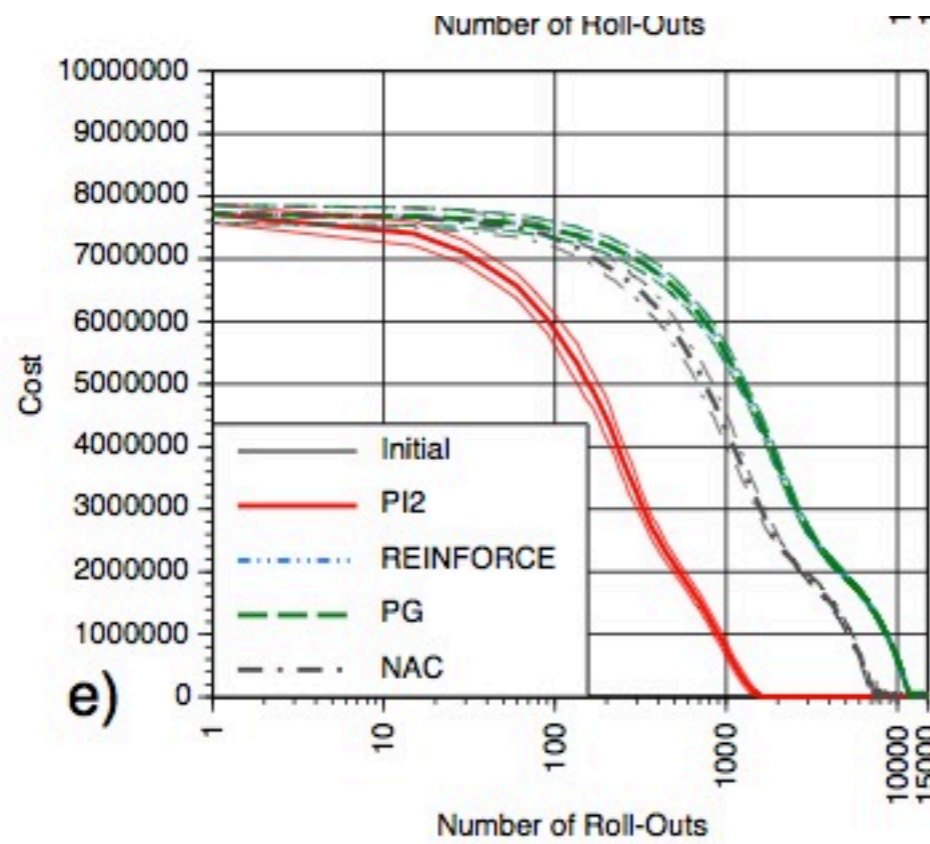
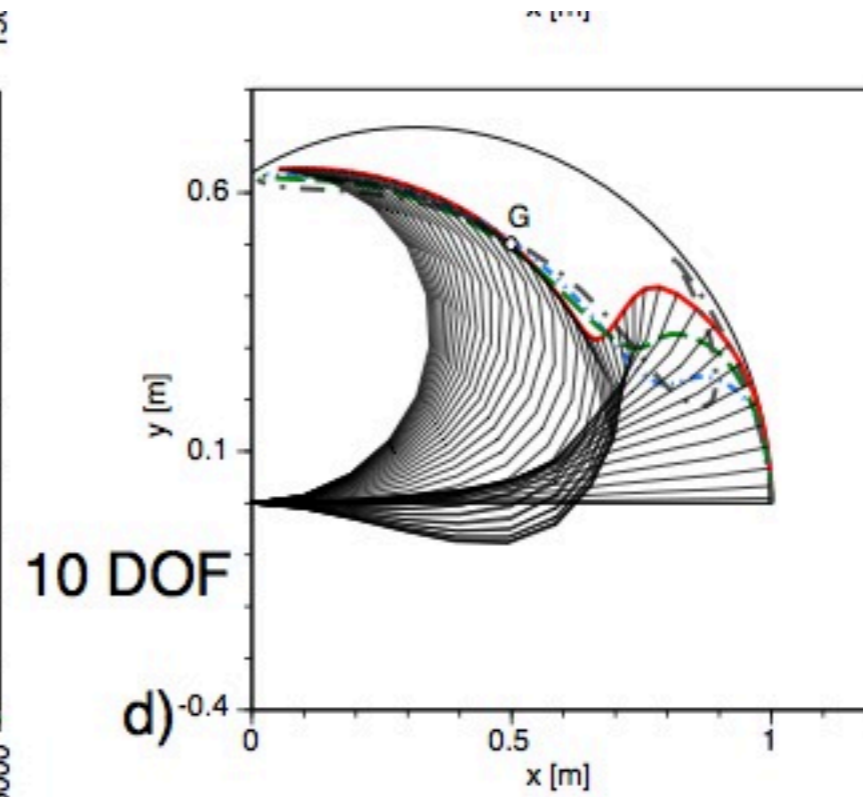
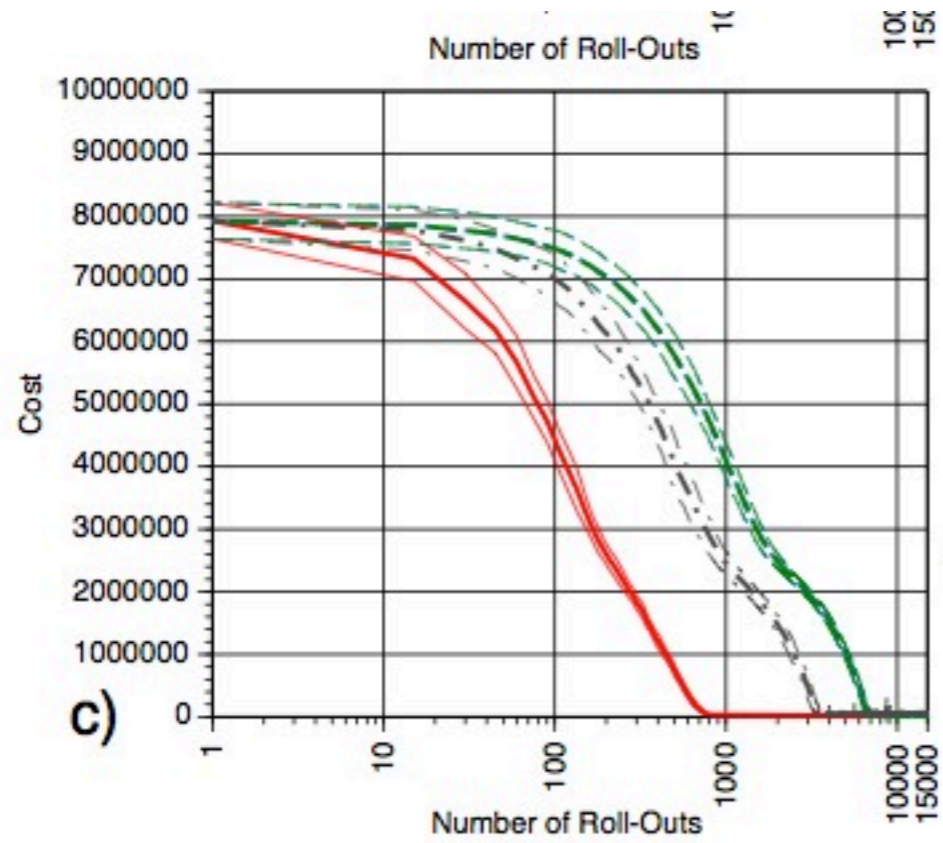
2 DOF

b)



a)





# Learning Feedback Gains with PI2

Consider system with tracking controller

$$\begin{cases} \dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x})\mathbf{u} \\ \mathbf{u} = \mathbf{K}^T(t)(\mathbf{x} - \mathbf{x}_{ref}) \end{cases}$$

$$\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x})(\mathbf{x} - \mathbf{x}_{ref})^T \mathbf{K}(t)$$

Treat gains  $\mathbf{K}$  as controls - use PI2 for time indexed policies





# Learning variable impedance control

Application of gain learning to rigid body  
dynamics / Torque controlled robots

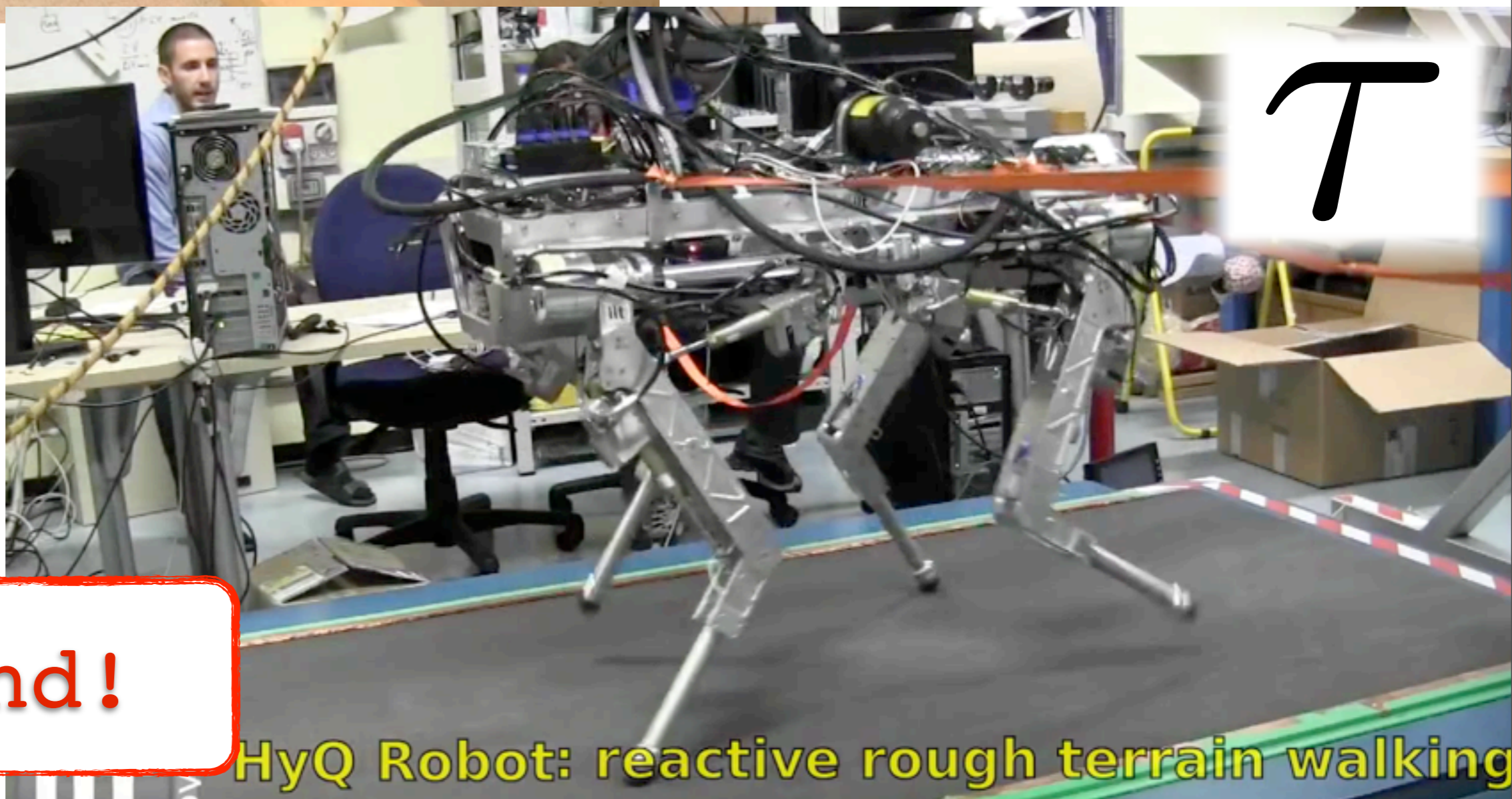


$\mathcal{X}$

High Gain PD control  
performs poorly on  
stochastic terrain

Position control  
only kinematic model  
high gains / 'stiff'

low gains / 'soft'  
Torque control  
+ dynamics model



$\mathcal{T}$

Blind!

HyQ Robot: reactive rough terrain walking

# Why?

Consider a robot with arms and legs:

- Low gains help to make the robot more robust in unknown surroundings
- Reduce energy consumption
- But position tracking is reduced

Thus, try to chose appropriate gain for the given situation. The requirements changes over time. Need a gain schedule.



# Model of a simple robot

Measurements: Joint angles:  $q = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix}$

Actuation: joint torques

$$\tau = \begin{bmatrix} \tau_1 \\ \tau_2 \\ \tau_3 \end{bmatrix}$$

Governing physics:  
Rigid body dynamics

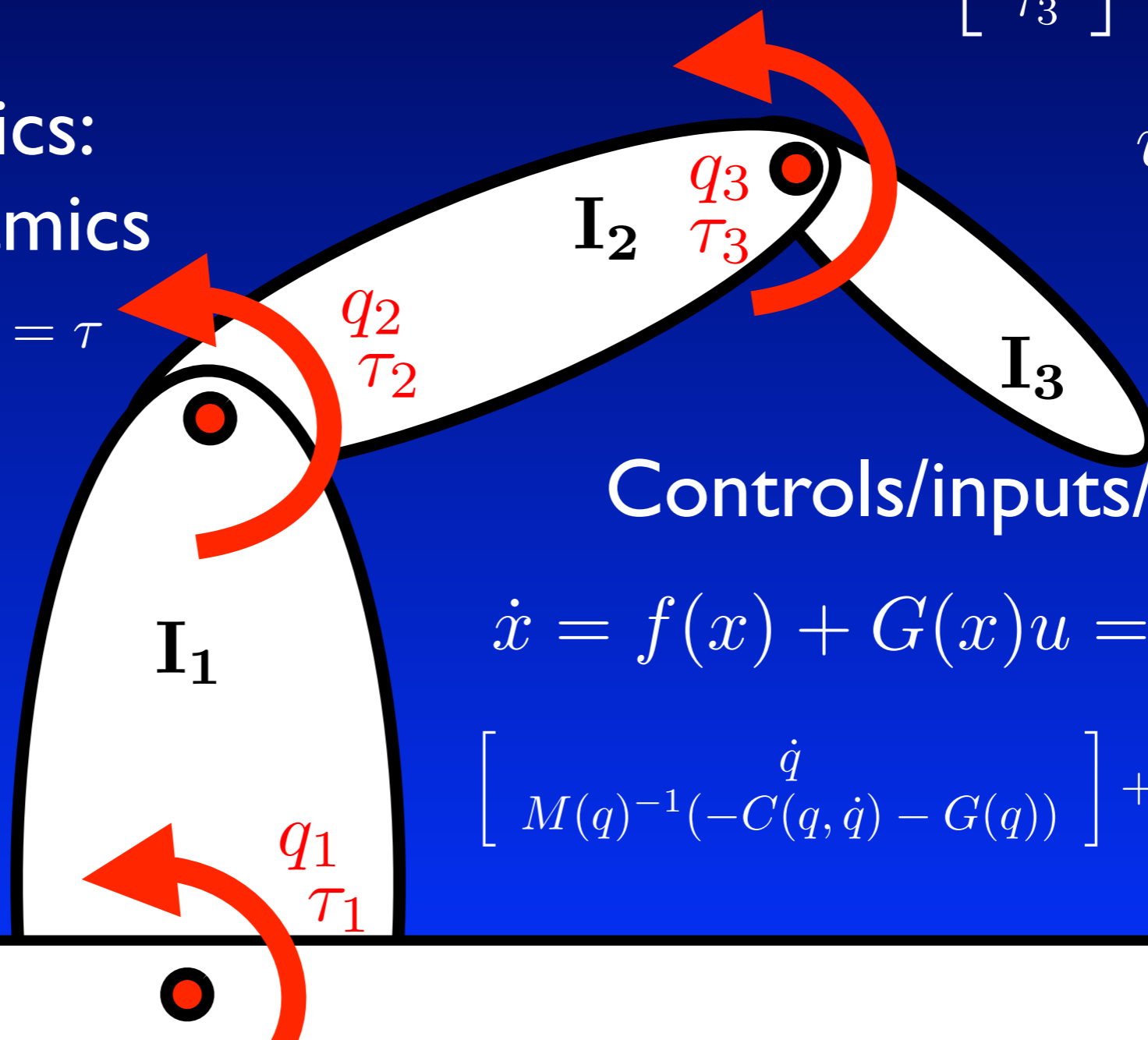
$$M\ddot{q} + C(q, \dot{q}) + G(q) = \tau$$

System states:

$$x = \begin{bmatrix} q \\ \dot{q} \end{bmatrix}$$

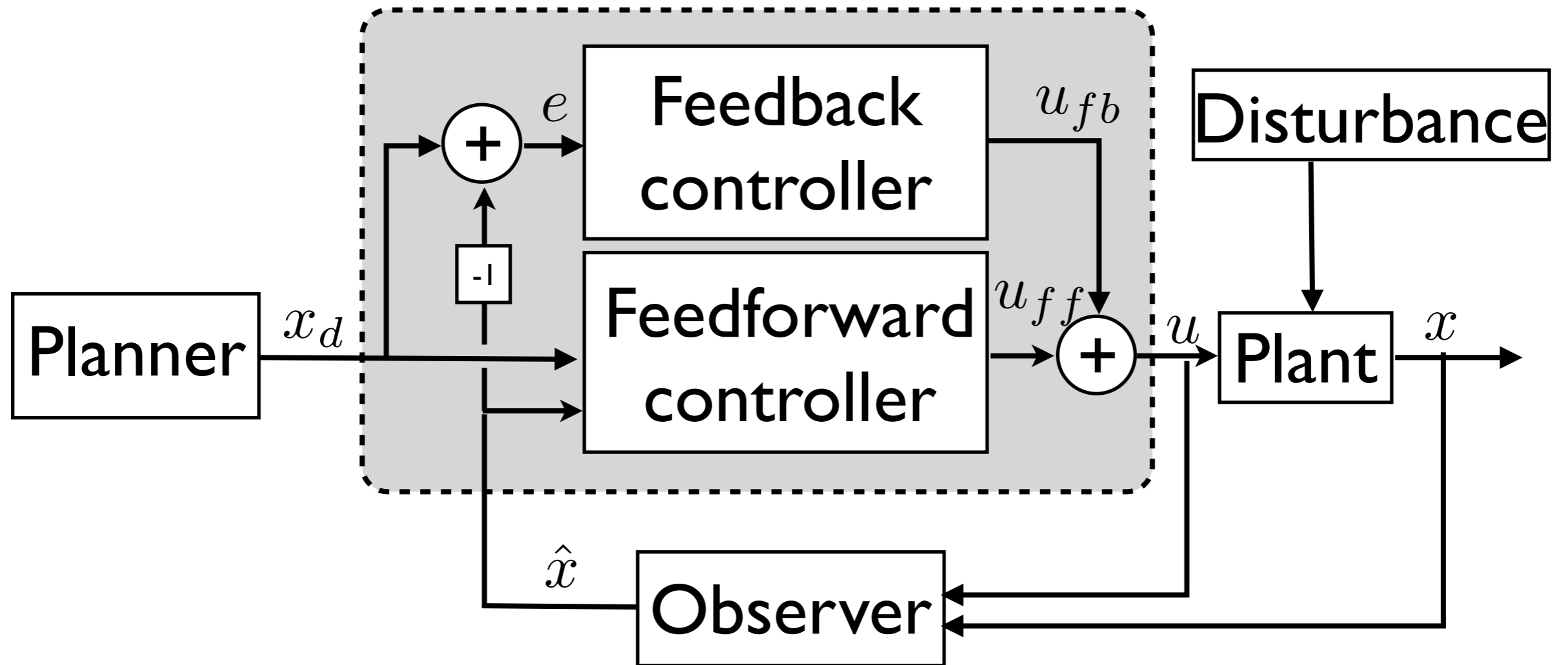
Controls/inputs/actions:

$$\dot{x} = f(x) + G(x)u = \begin{bmatrix} \dot{q} \\ M(q)^{-1}(-C(q, \dot{q}) - G(q)) \end{bmatrix} + \begin{bmatrix} 0 \\ M(q)^{-1} \end{bmatrix}^T u$$

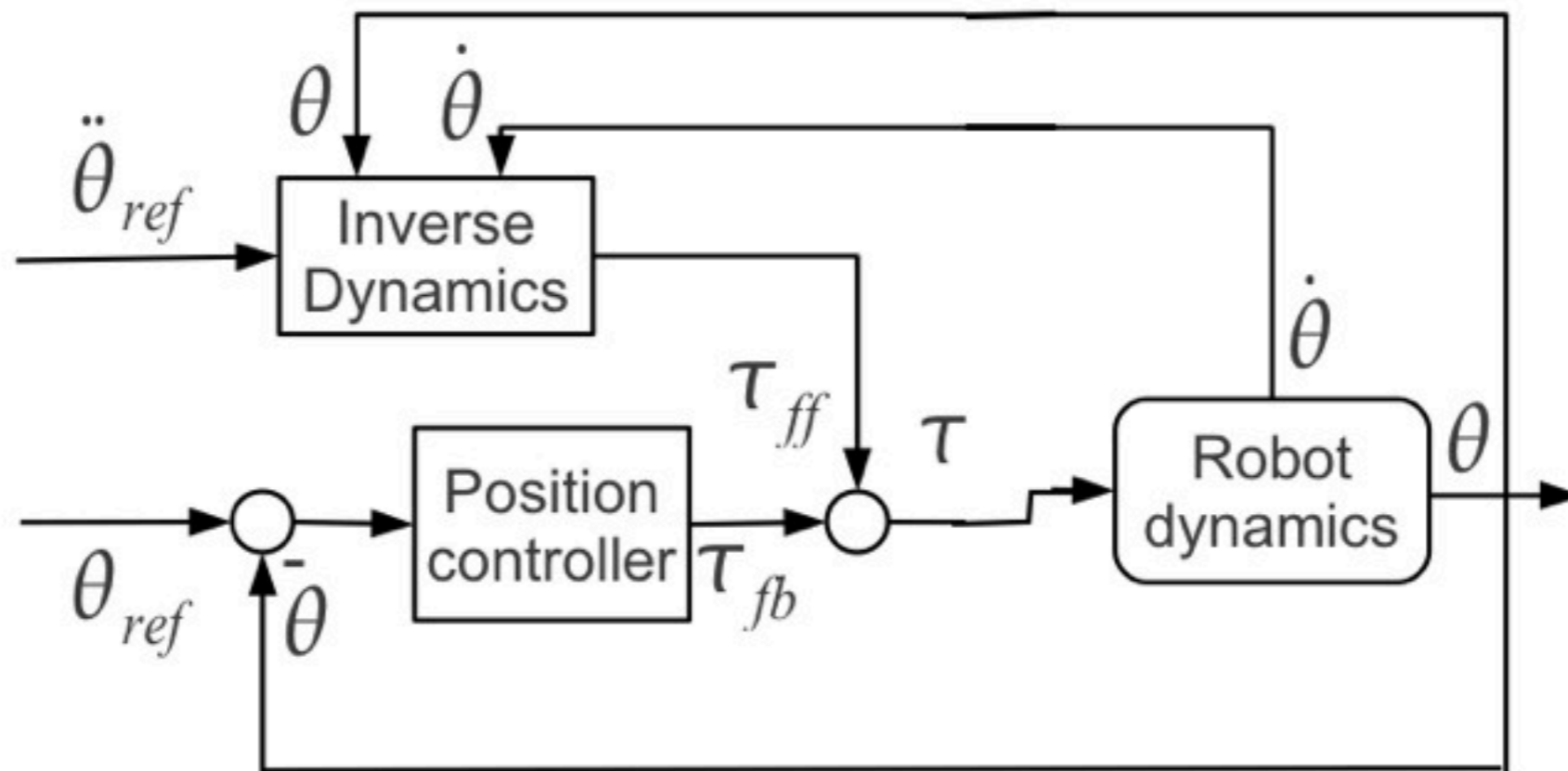


$$u = \tau$$

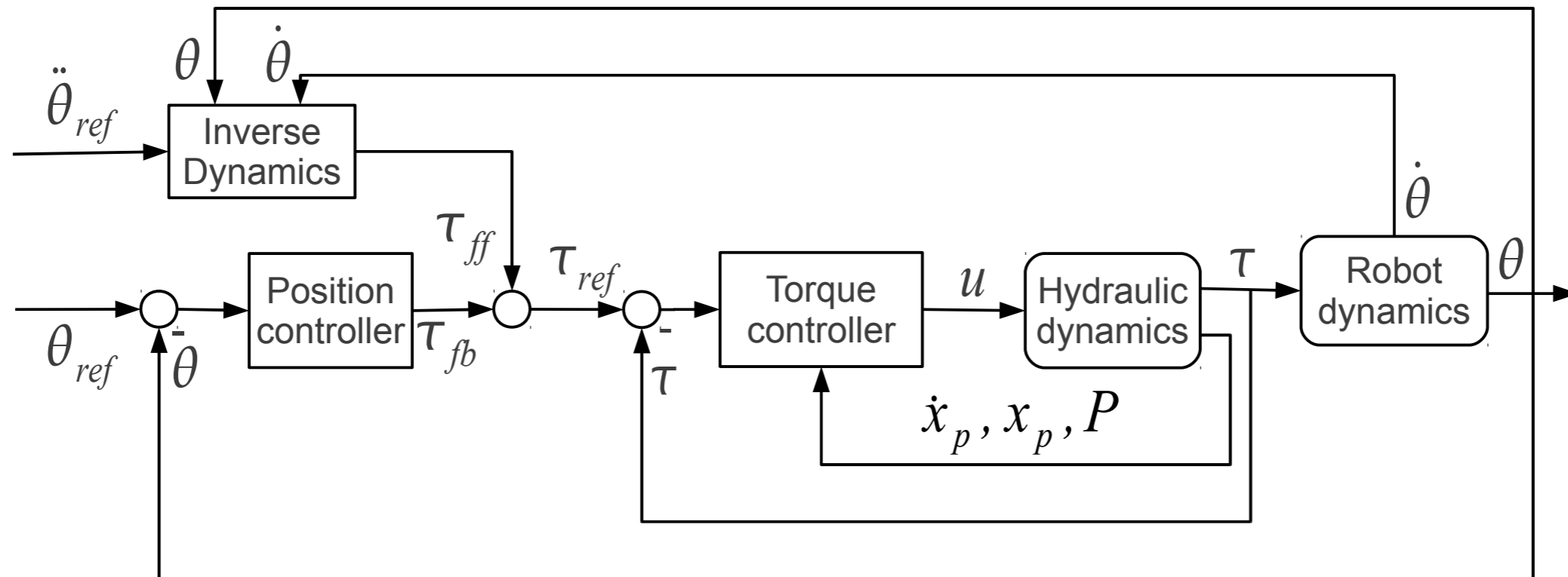
# General control structure



# Torque controlled robot - Control diagram



# Impedance control

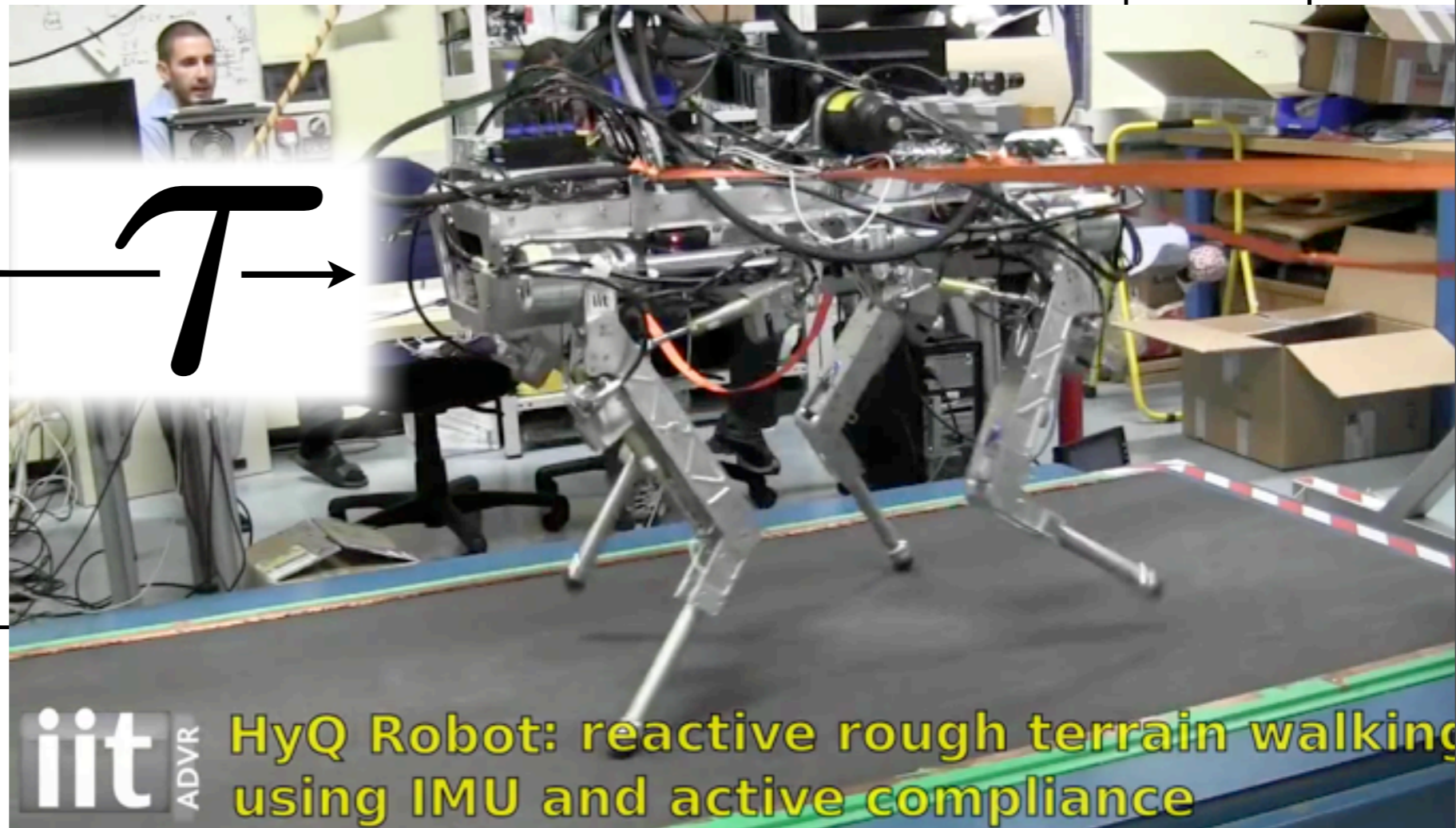
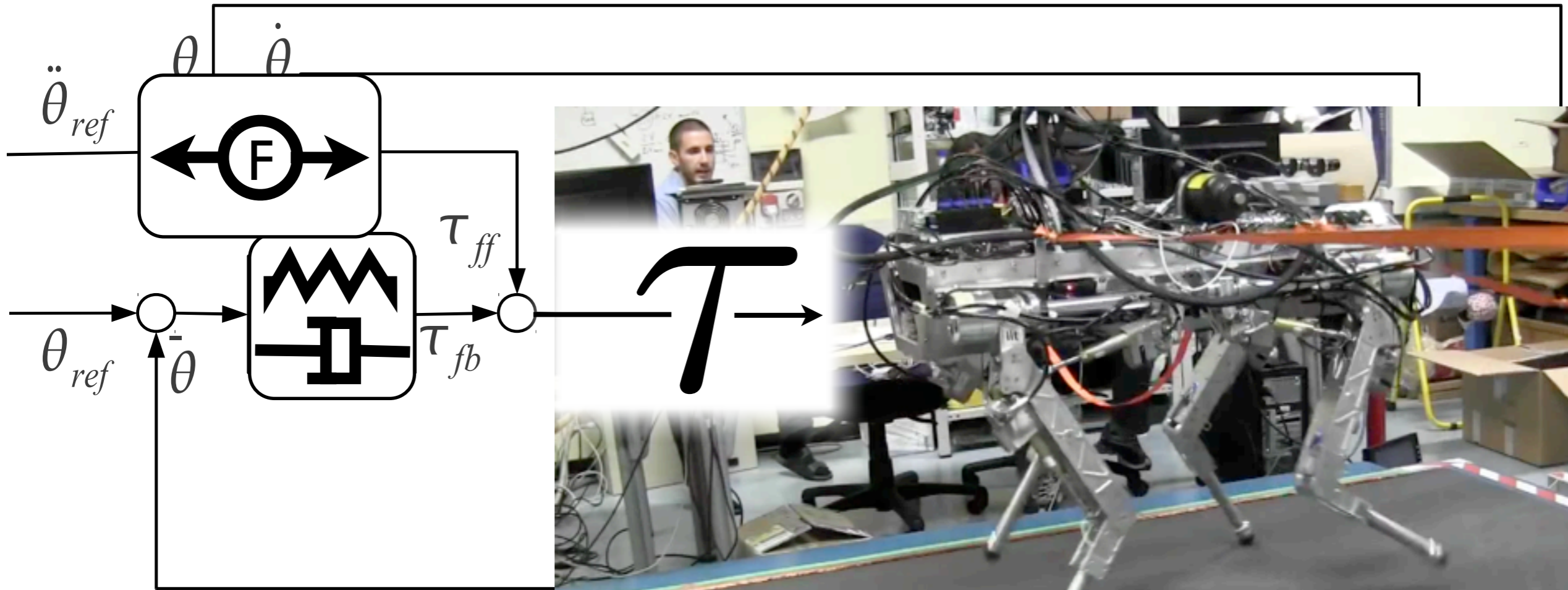


Compare to LQR control structure:

$$\mathbf{T} = -\mathbf{K}_P(\mathbf{q} - \mathbf{q}_d) - \mathbf{K}_D(\dot{\mathbf{q}} - \dot{\mathbf{q}}_d) + \mathbf{T}_{ff}$$

units?

# Torque controlled legged robots

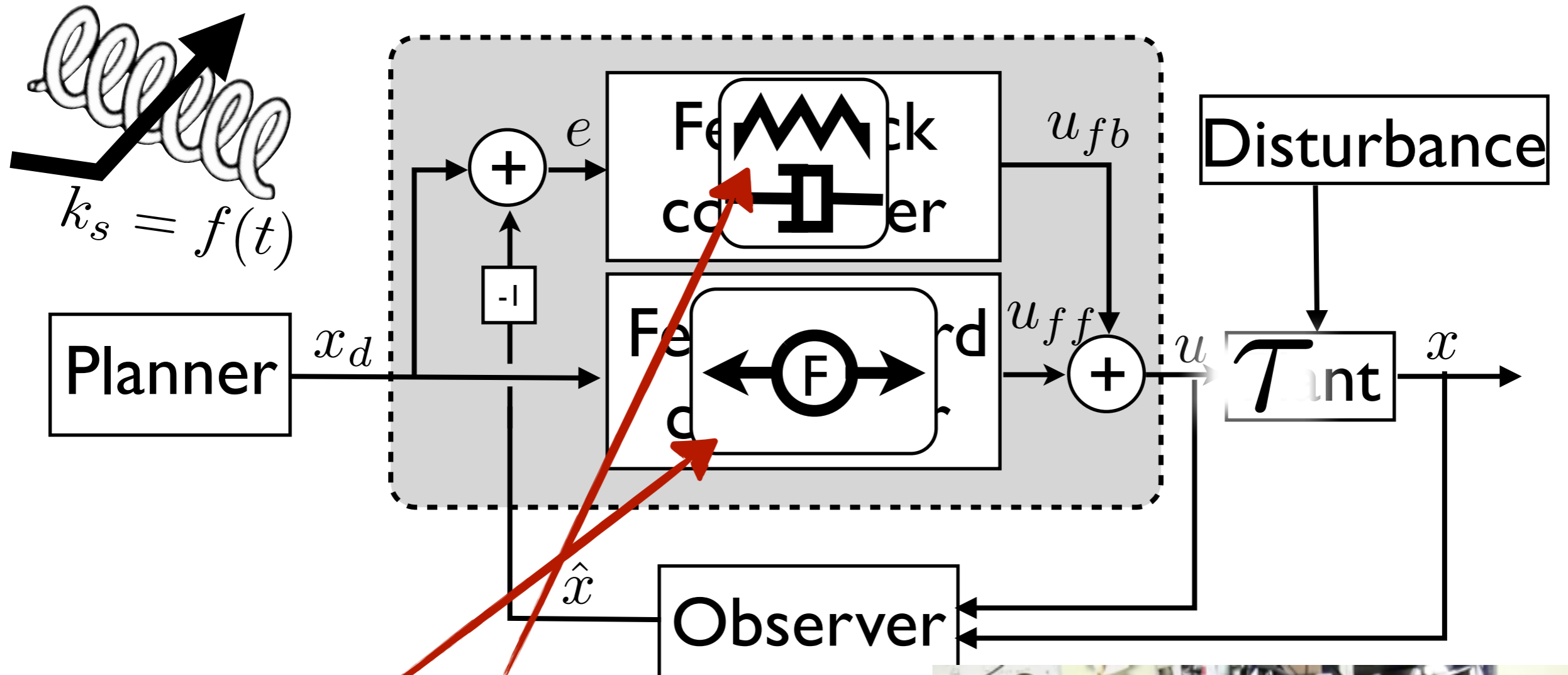


Torque control  
+ dynamics model

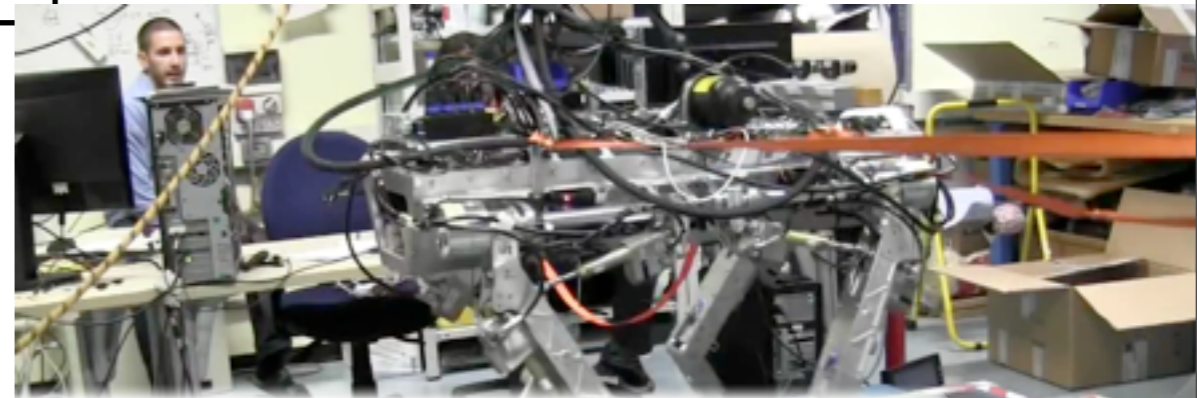




# Optimal Impedance Control



$$\begin{aligned}
 \mathbf{u}(t) &= \mathbf{u}^*(t) + \Delta \mathbf{u}^*(t) \\
 &= \mathbf{u}^*(t) - \mathbf{C}(t)[\mathbf{x}(t) - \mathbf{x}^*(t)]
 \end{aligned}$$



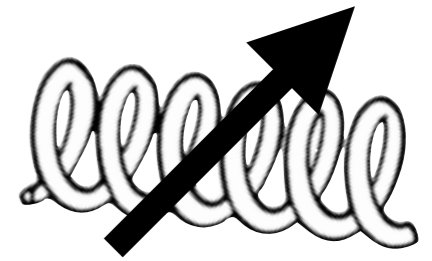
Principle of optimality: Time varying stiffness -  $\mathbf{K}(t)$

ain walkin



# Virtual compliance



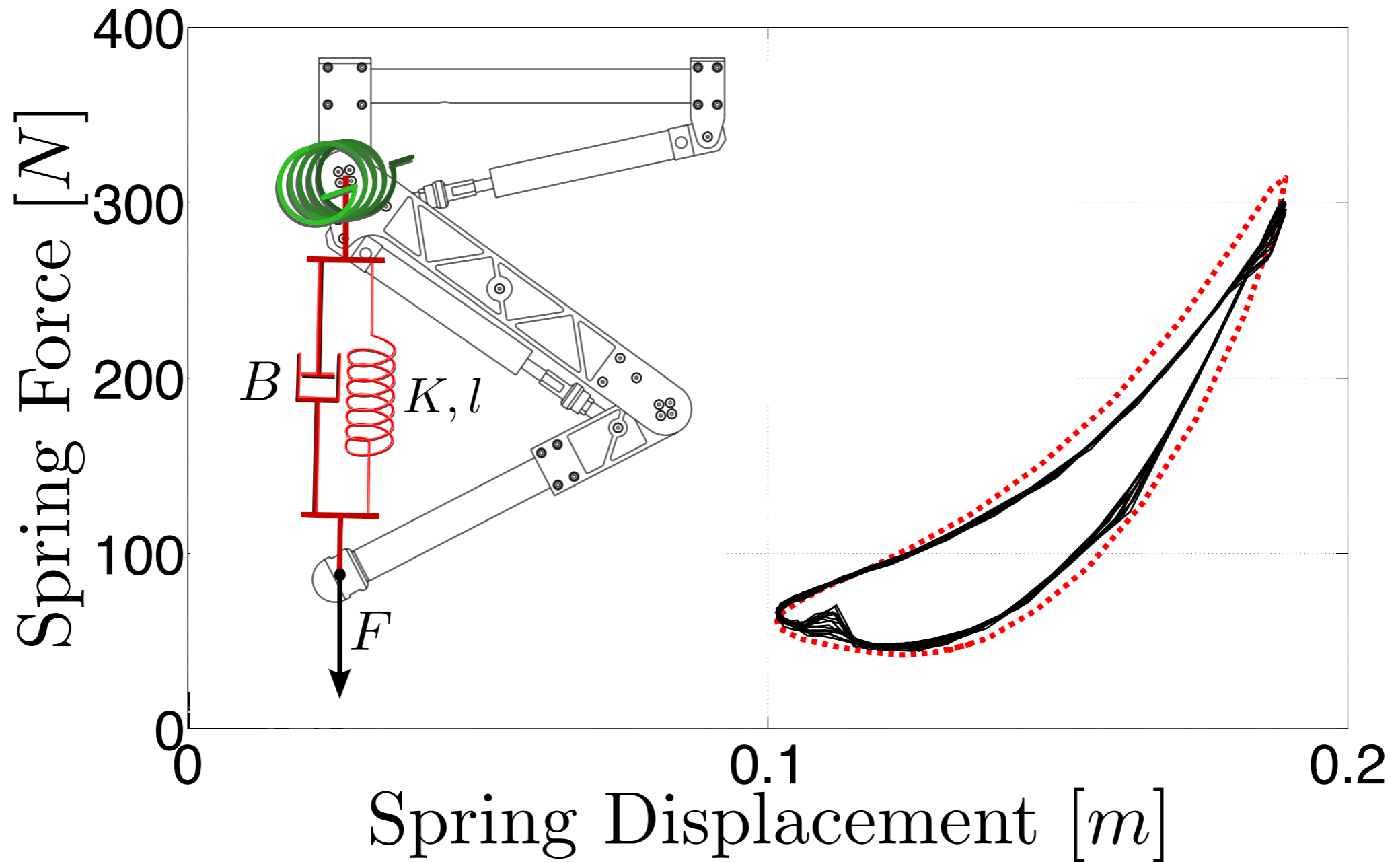


$\dot{k}_s = f(t)$   
Active springs!

# Virtual nonlinear springs



$$\tau = J^T F$$



# Active vs passive compliance

Active Compliance for highly-dynamic legged locomotion

Thiago Boaventura, Gustavo A. Medrano-Cerda, Ioannis Havoutis,  
Claudio Semini, Jonas Buchli, Darwin G. Caldwell



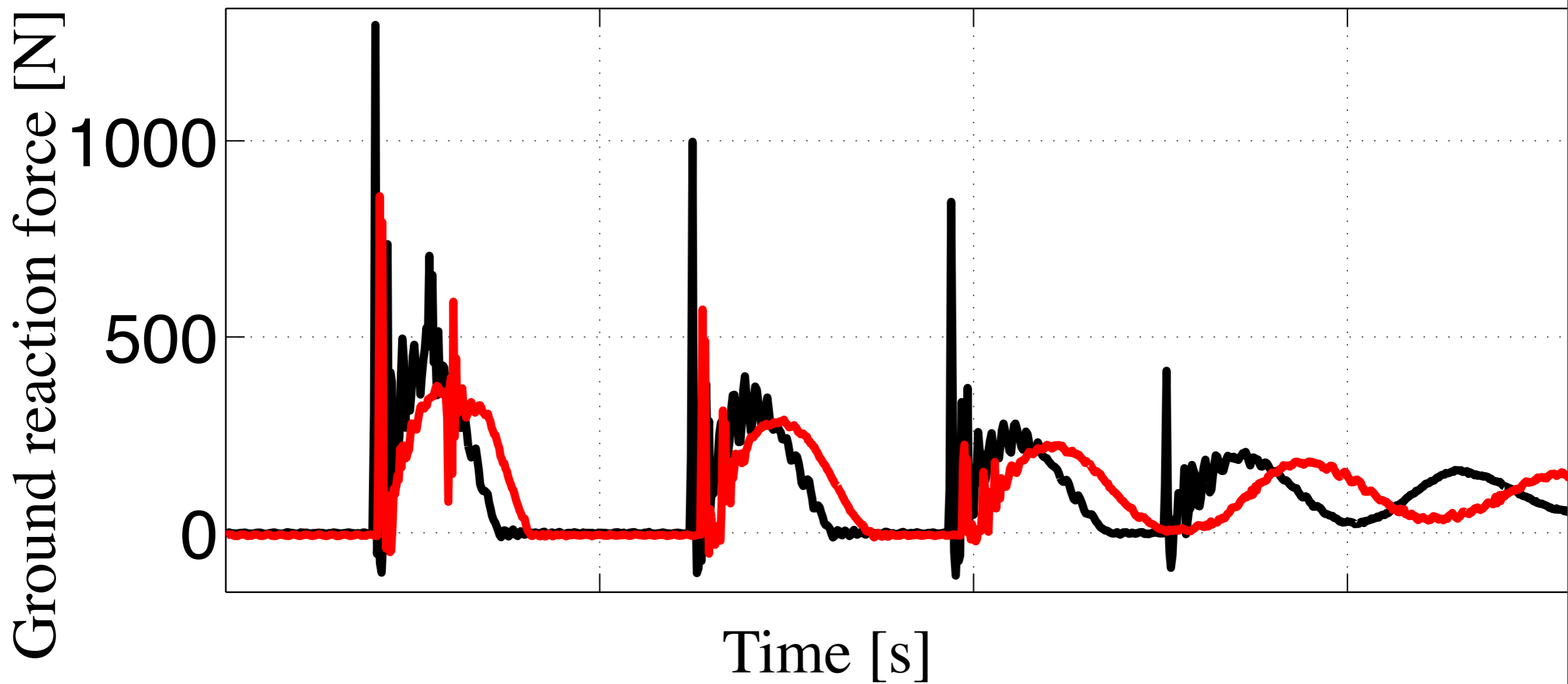
ICRA 2013

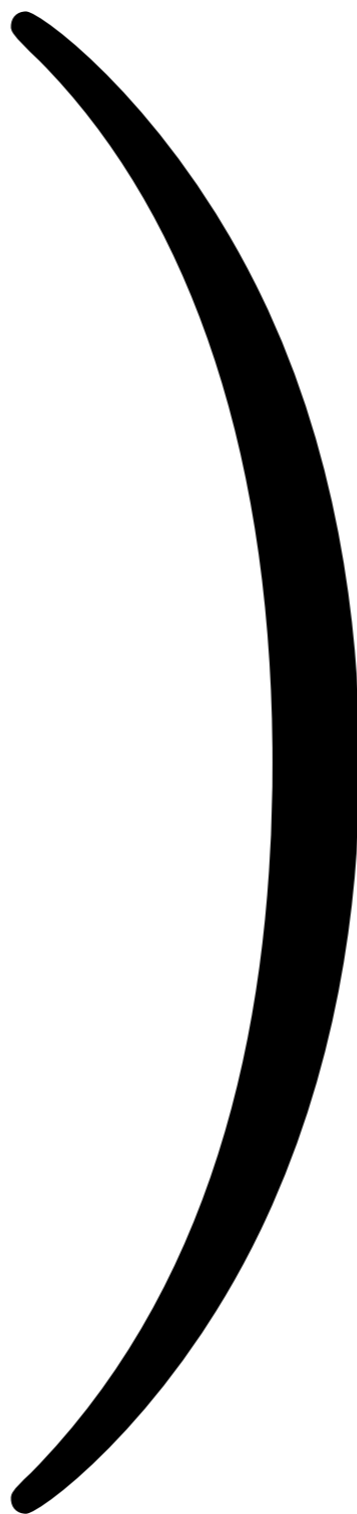


Controller: [Boaventura, ICRA 2012]



# Active vs passive compliance





Buchli - OLCAR - 2015

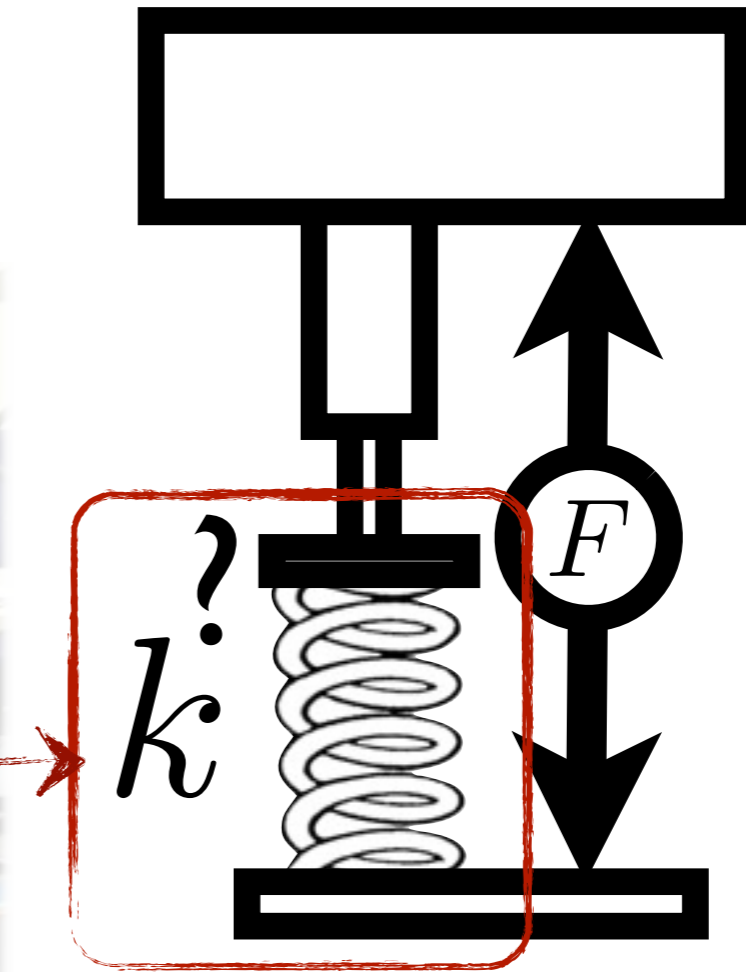
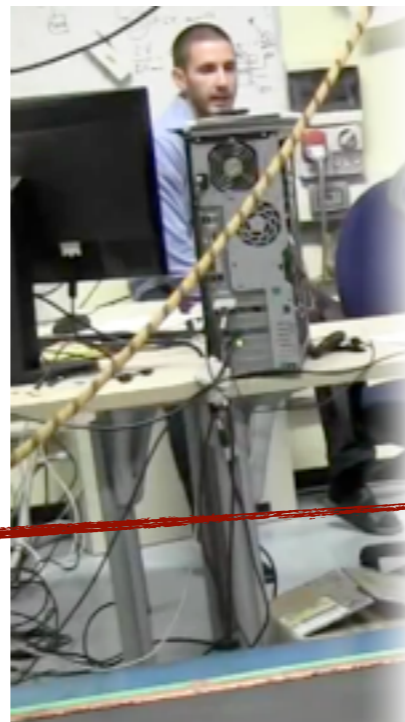
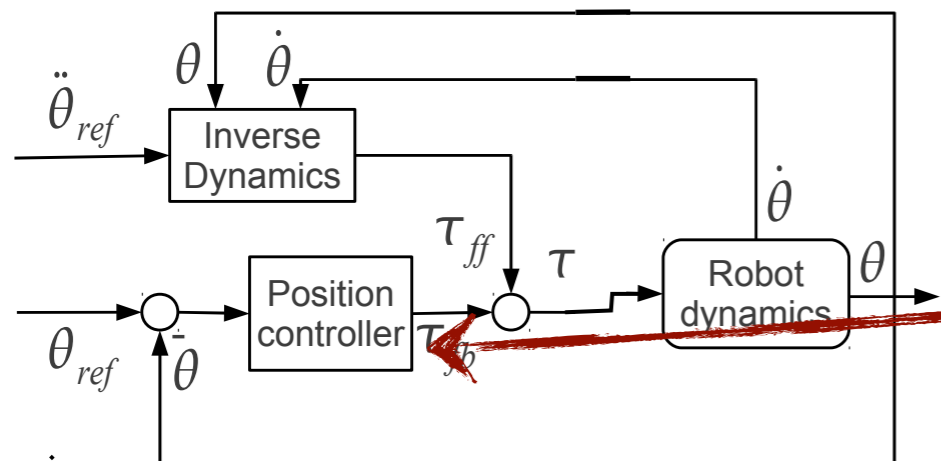


# Inverse dynamics

Position control  
only kinematic model



$$M(q)\ddot{q} + h(q, \dot{q}) = S^T \tau + J_C^T(q) \lambda$$



iit ADVR HyQ Robot: reactive rough terrain walking using IMU and active compliance



Torque control  
+ dynamics model



# Parametrization

## Policies to be learned

Position reference:  $\mathbf{X}_D(\mathbf{t}) = \phi(\mathbf{t})^T \boldsymbol{\theta}_x$

Spring resting length

Gain schedules:  $\mathbf{K}_P(\mathbf{t}) = \phi(\mathbf{t})^T \boldsymbol{\theta}_p$

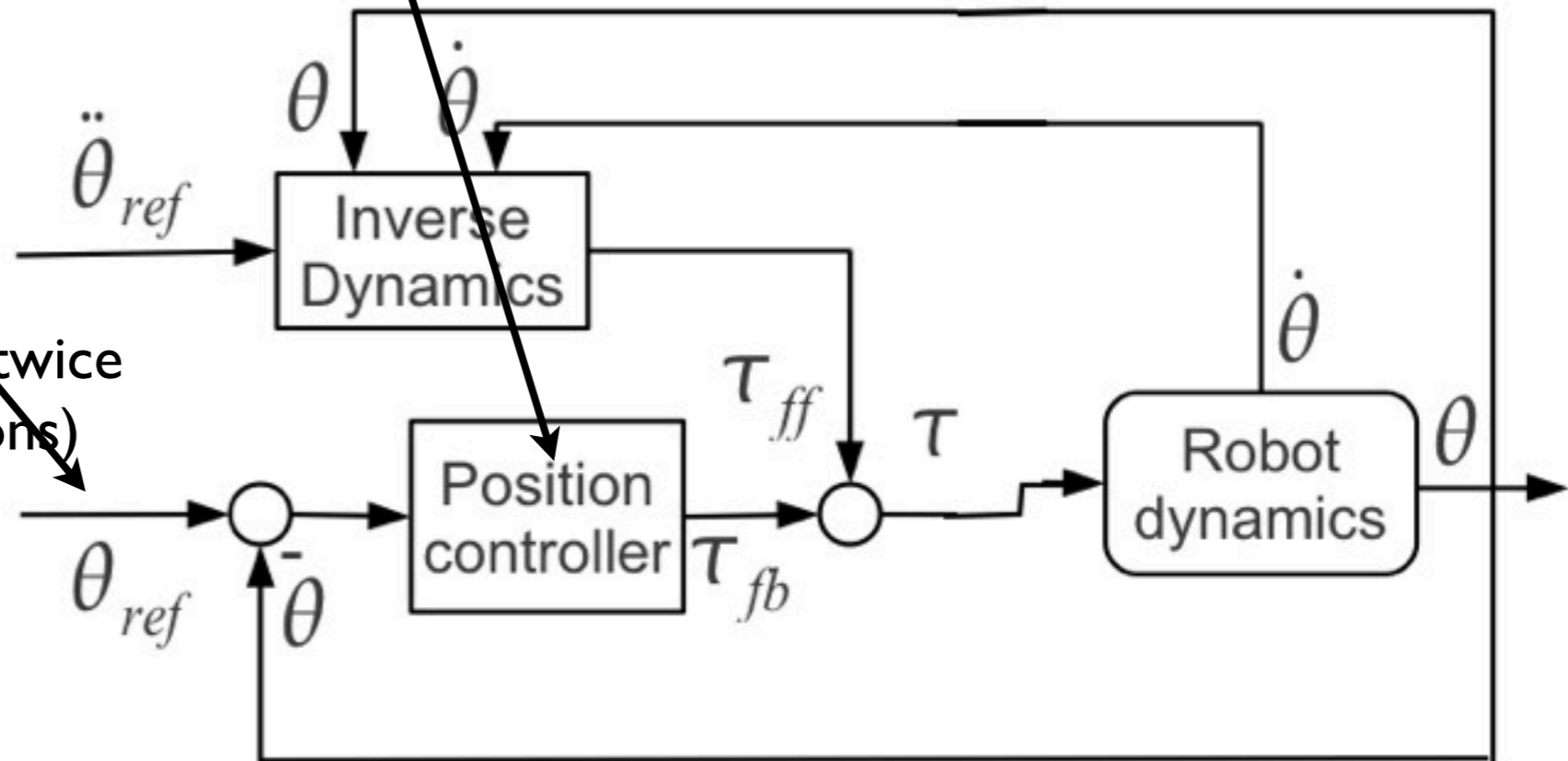
Spring Damper

$\mathbf{K}_D(\mathbf{t}) = \phi(\mathbf{t})^T \boldsymbol{\theta}_d$

Gain schedules:

$\mathbf{r}(t)^T$

make sure position reference is twice differentiable (finite accelerations)



# Via-point task

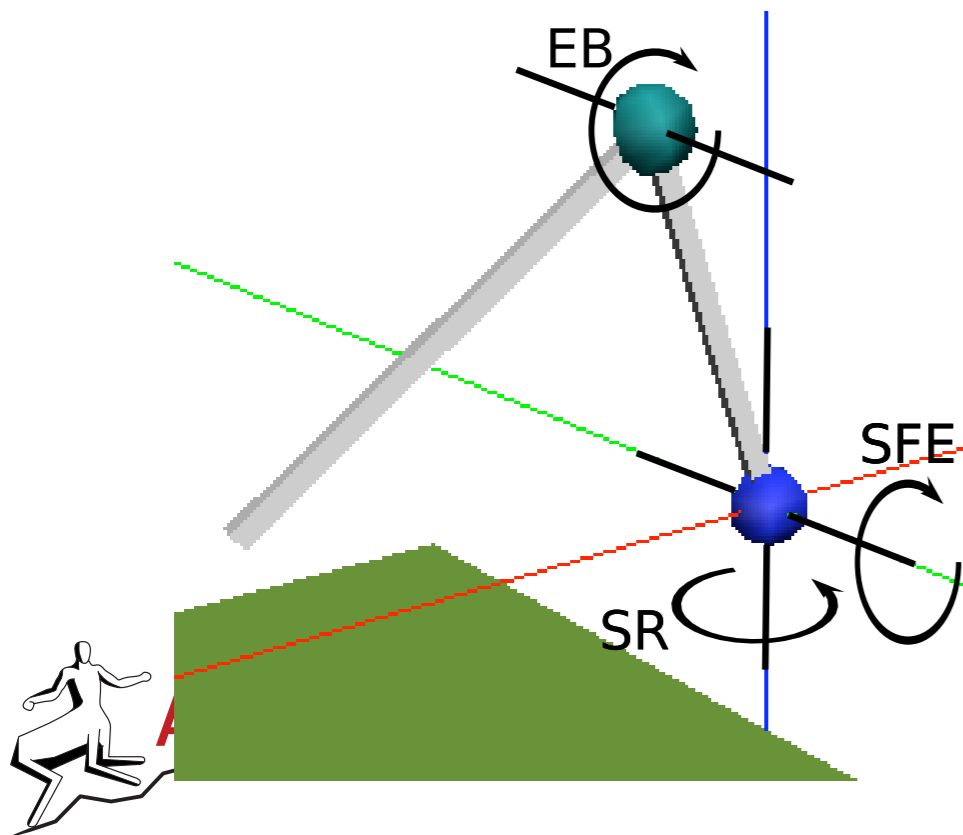
$$r_t = W_{gain} \sum_i K_{P,t}^i + W_{acc} \|\ddot{\mathbf{x}}\| + W_{via-point} C(t)$$

‘keep gains low / stay soft!’

‘don’t wiggle too much’

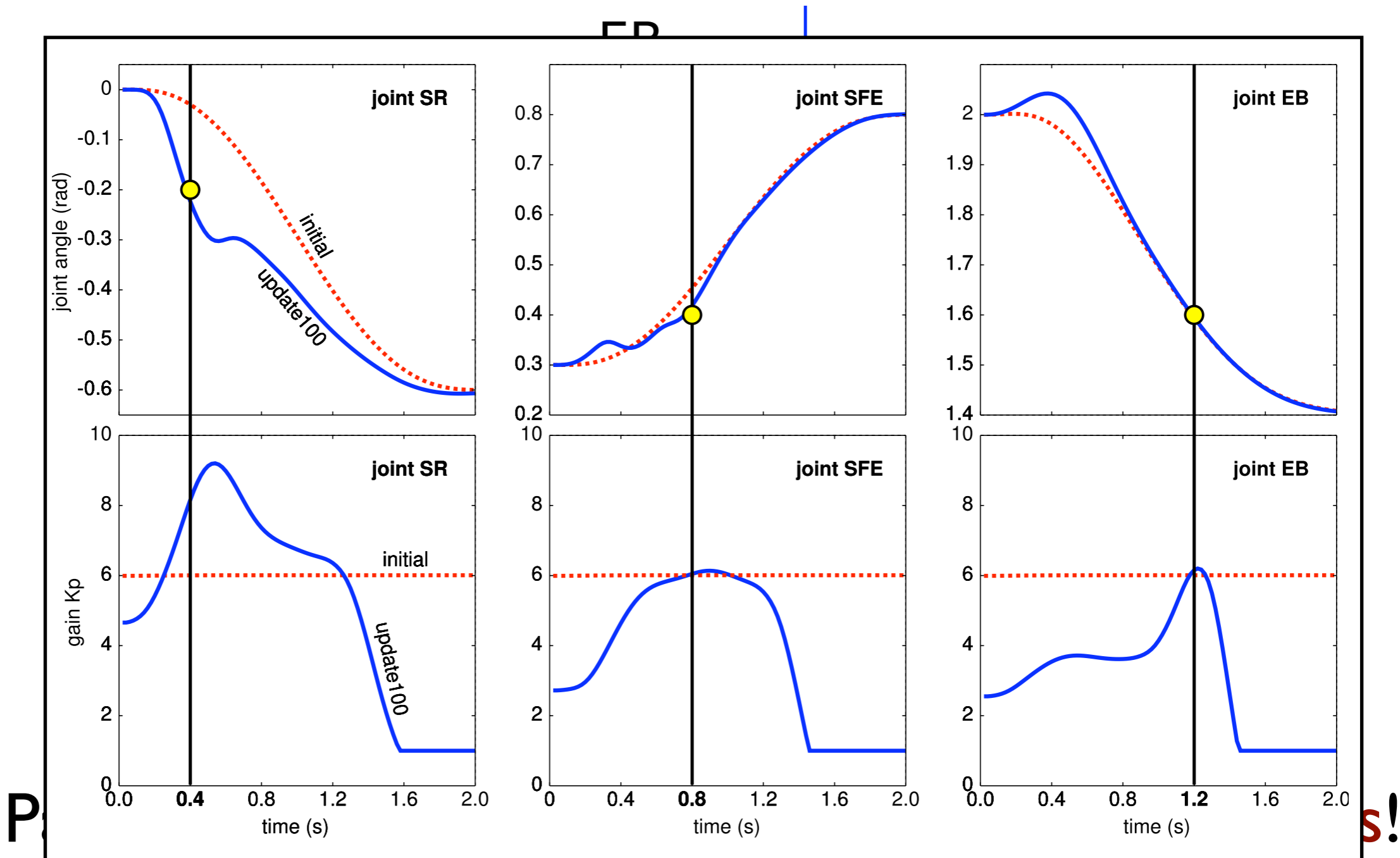
go through given joint angles at given times:

$$C(t) = \delta(t - 0.4) \cdot |q_{SR}(t) + 0.2| \\ + \delta(t - 0.8) \cdot |q_{SFE}(t) - 0.4| \\ + \delta(t - 1.2) \cdot |q_{EB}(t) - 1.5|$$



# Gain scheduling

## 3 DOF - Phantom robot



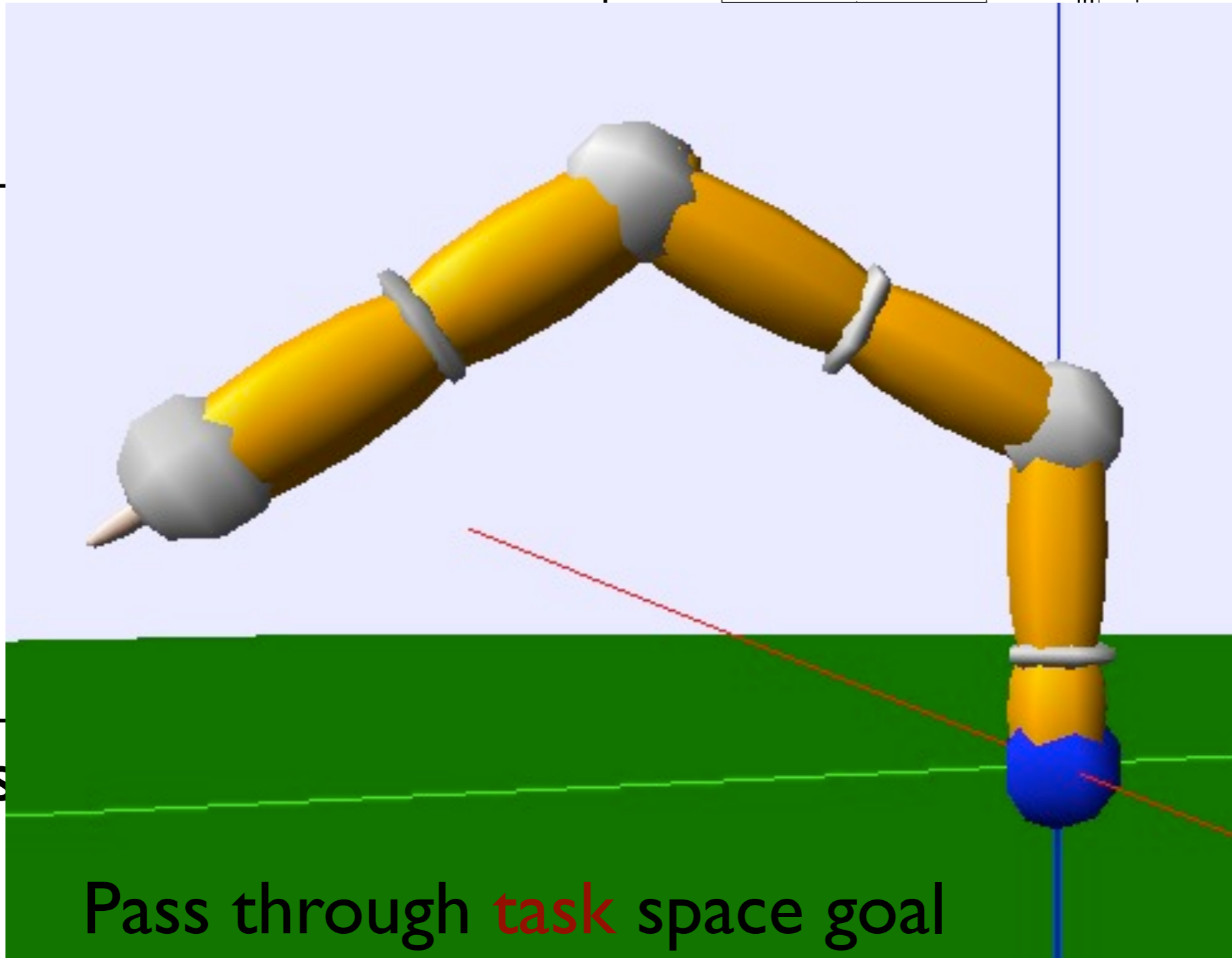
[Buchli et al, RSS 2010]

# Gain scheduling

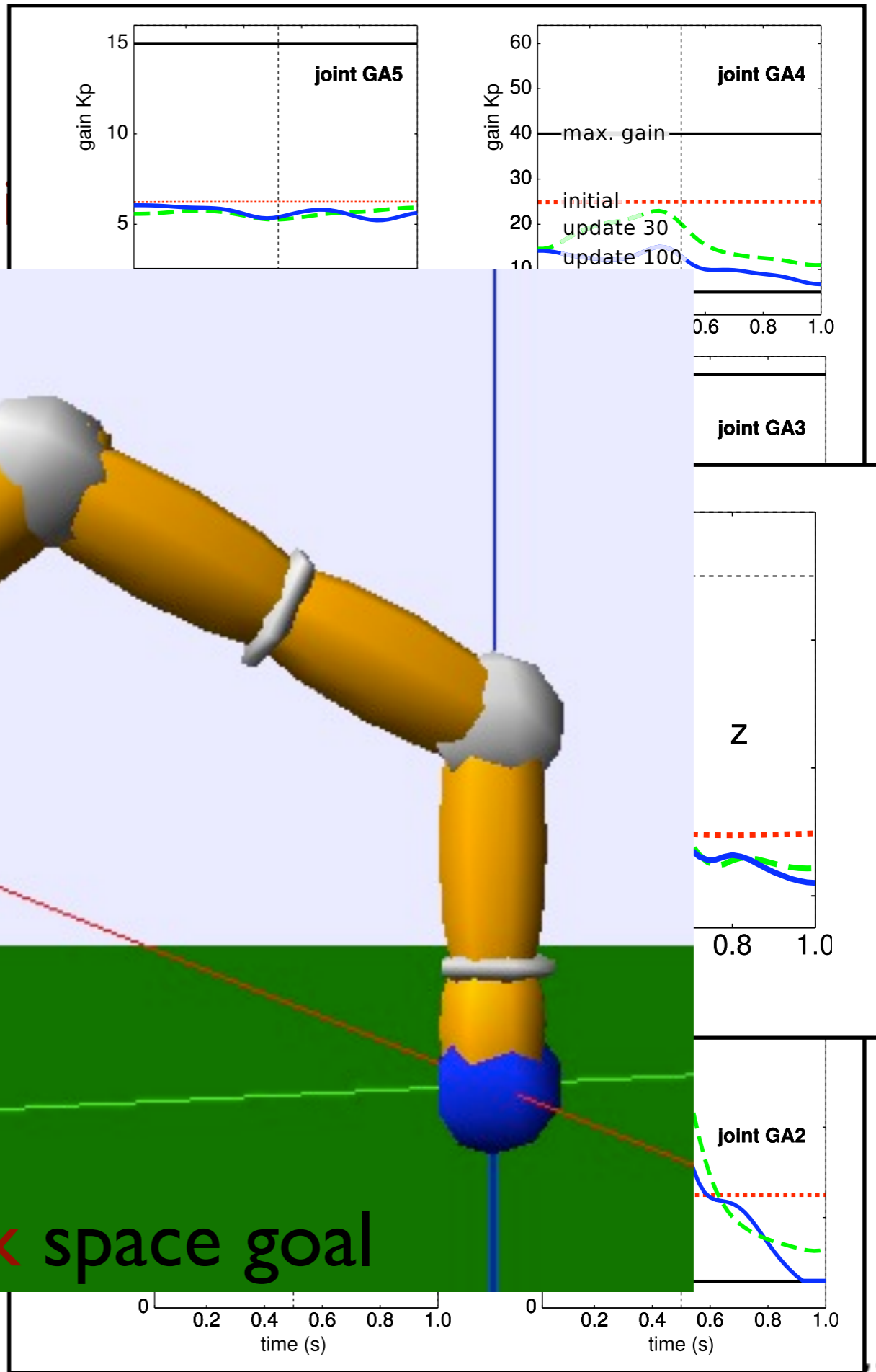
## 6DOF - Kuka I

end-effector position (m)

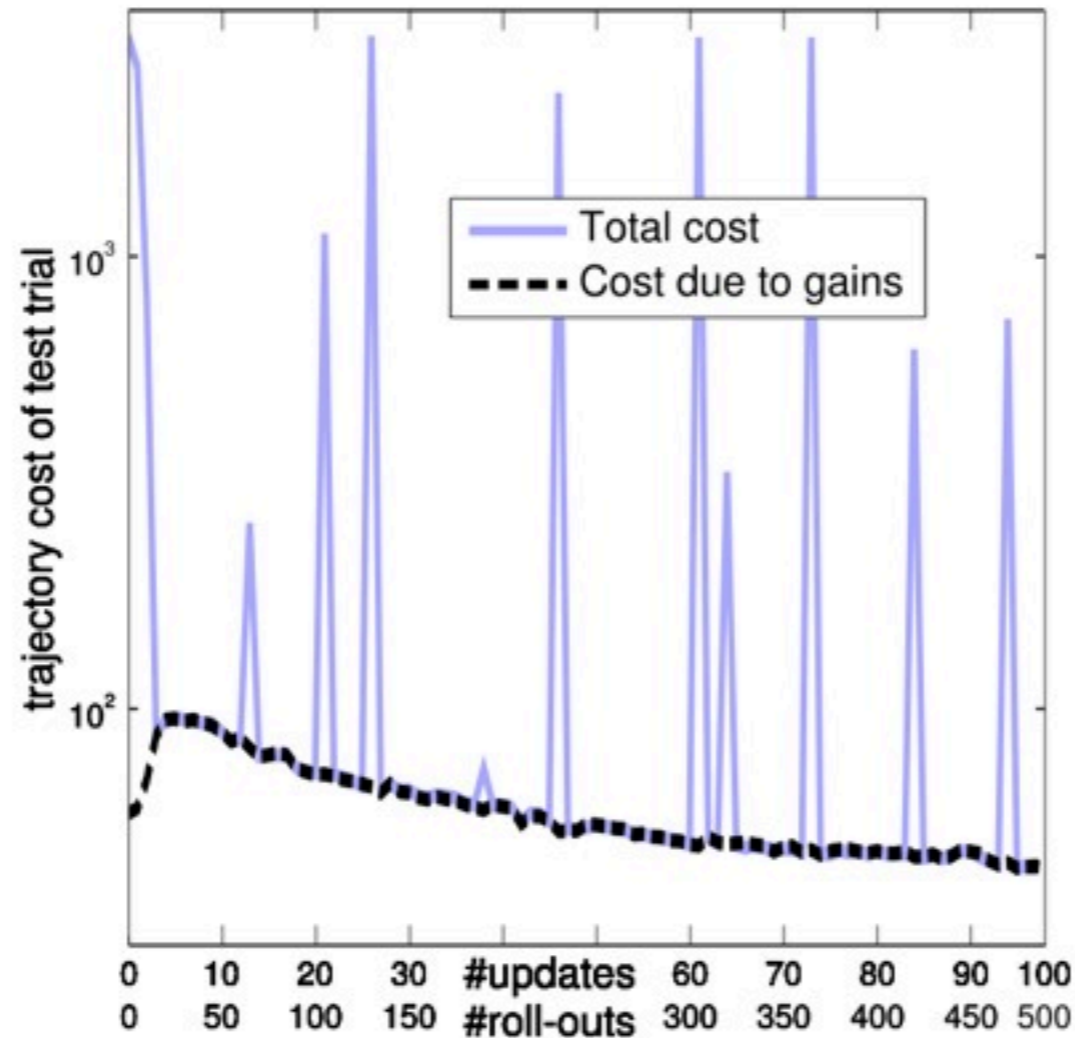
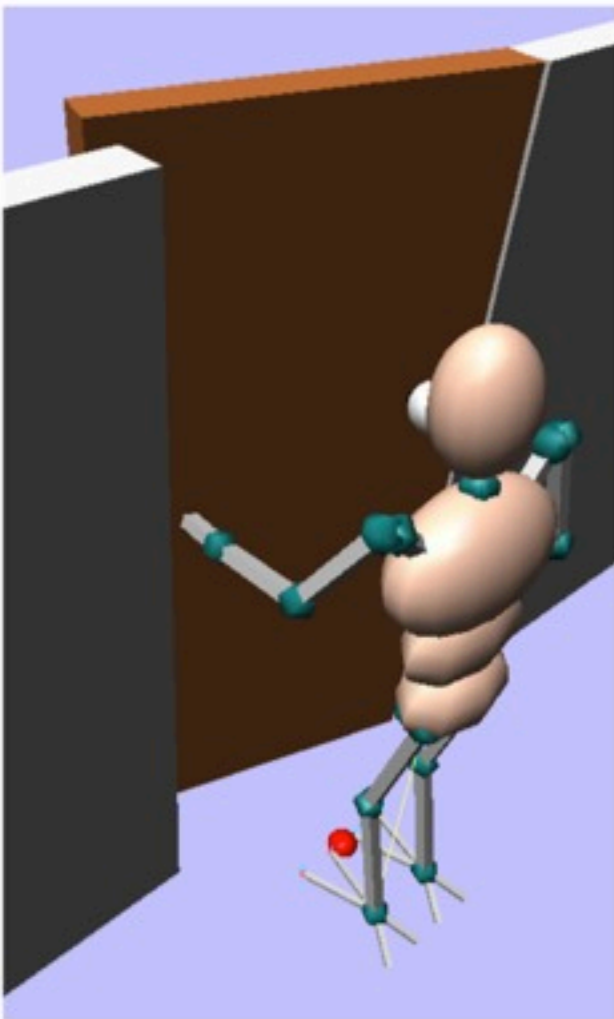
Pass



Pass through **task** space goal



Goal: Open door as far as possible and use lowest amount of effort



$$r_t = \frac{1}{N} \sum_{i=1}^3 K_{P,t}^i$$

$$\phi_{t_N} = 10^4 \cdot (\psi_{max} - \psi_N)$$



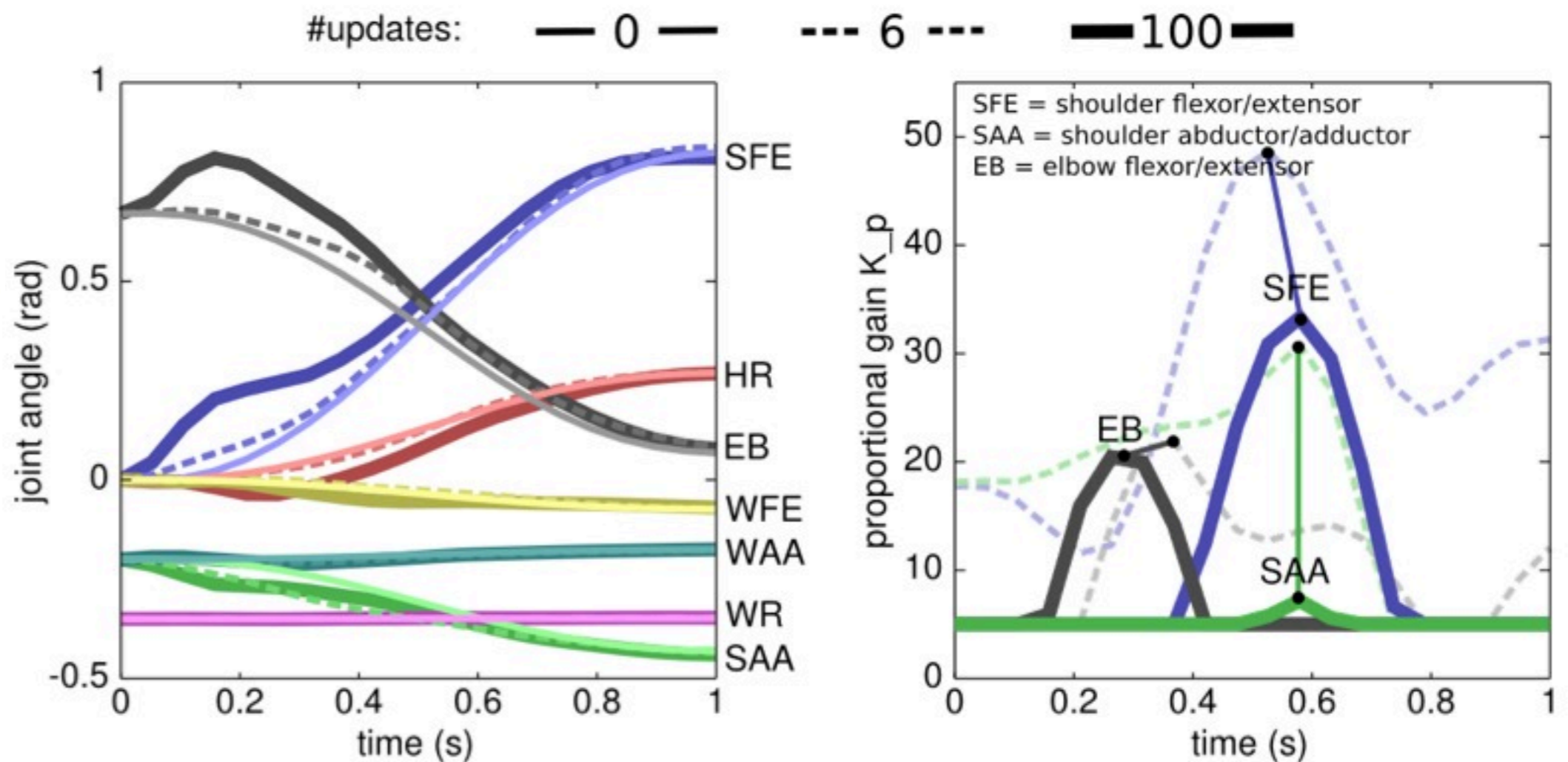
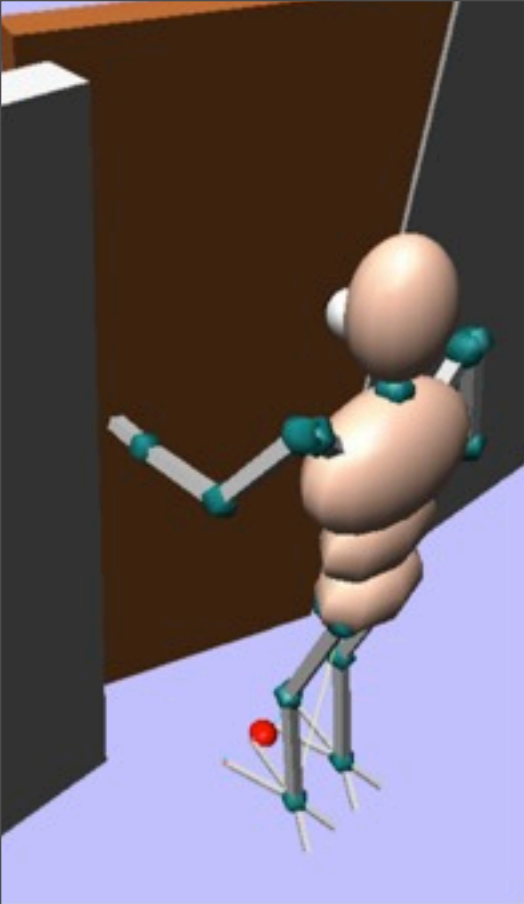
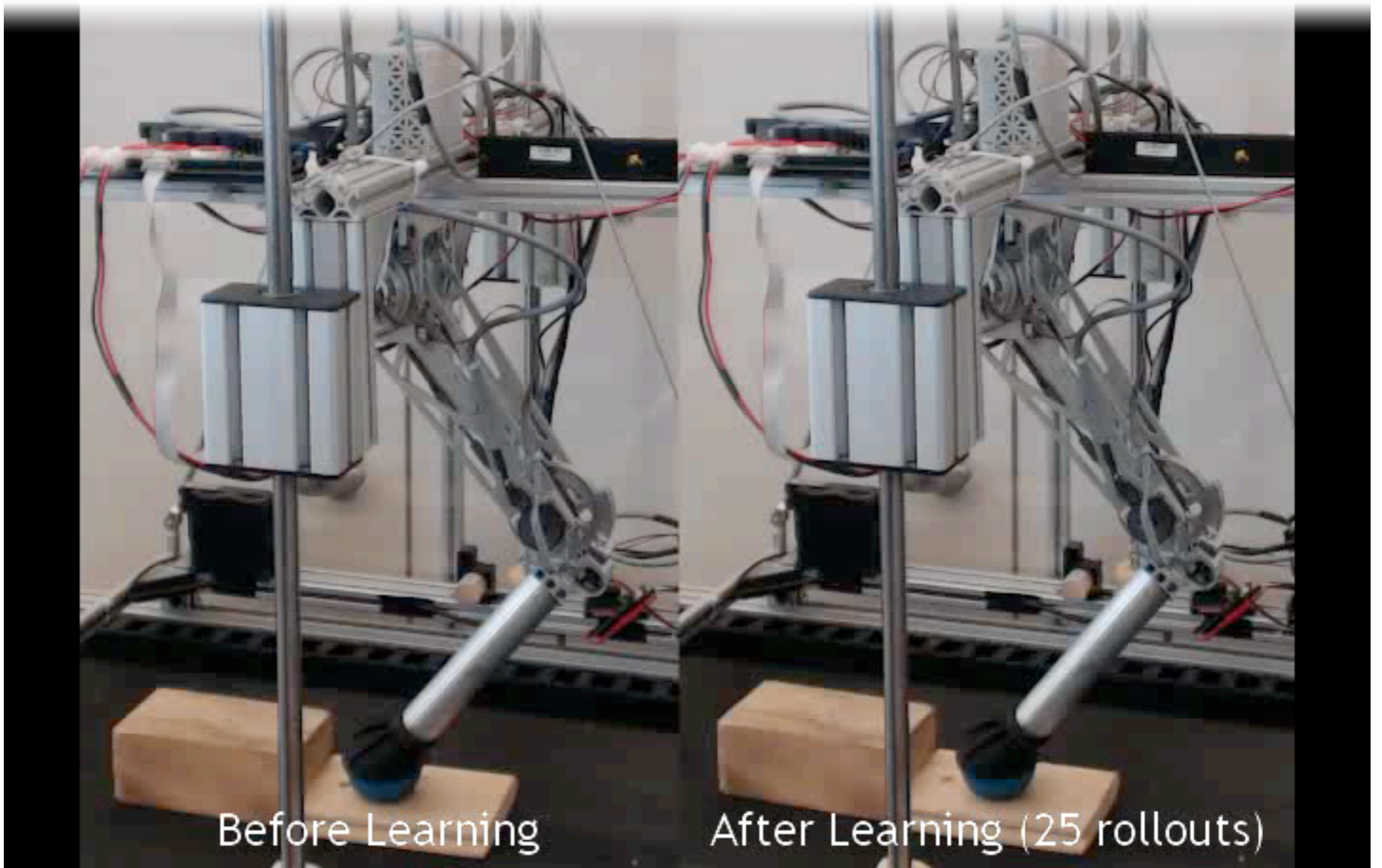
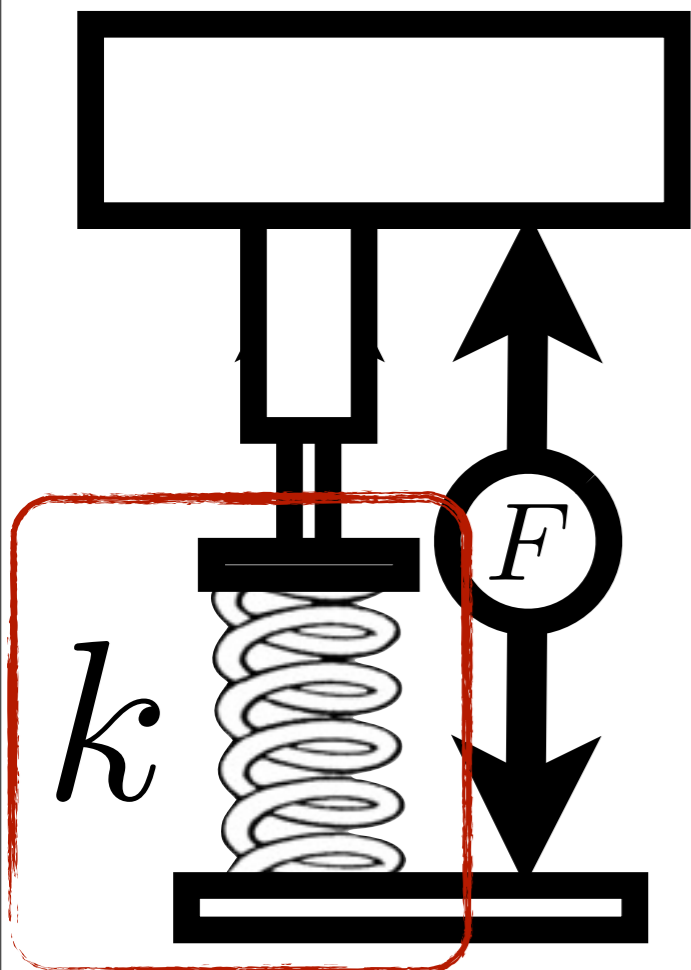
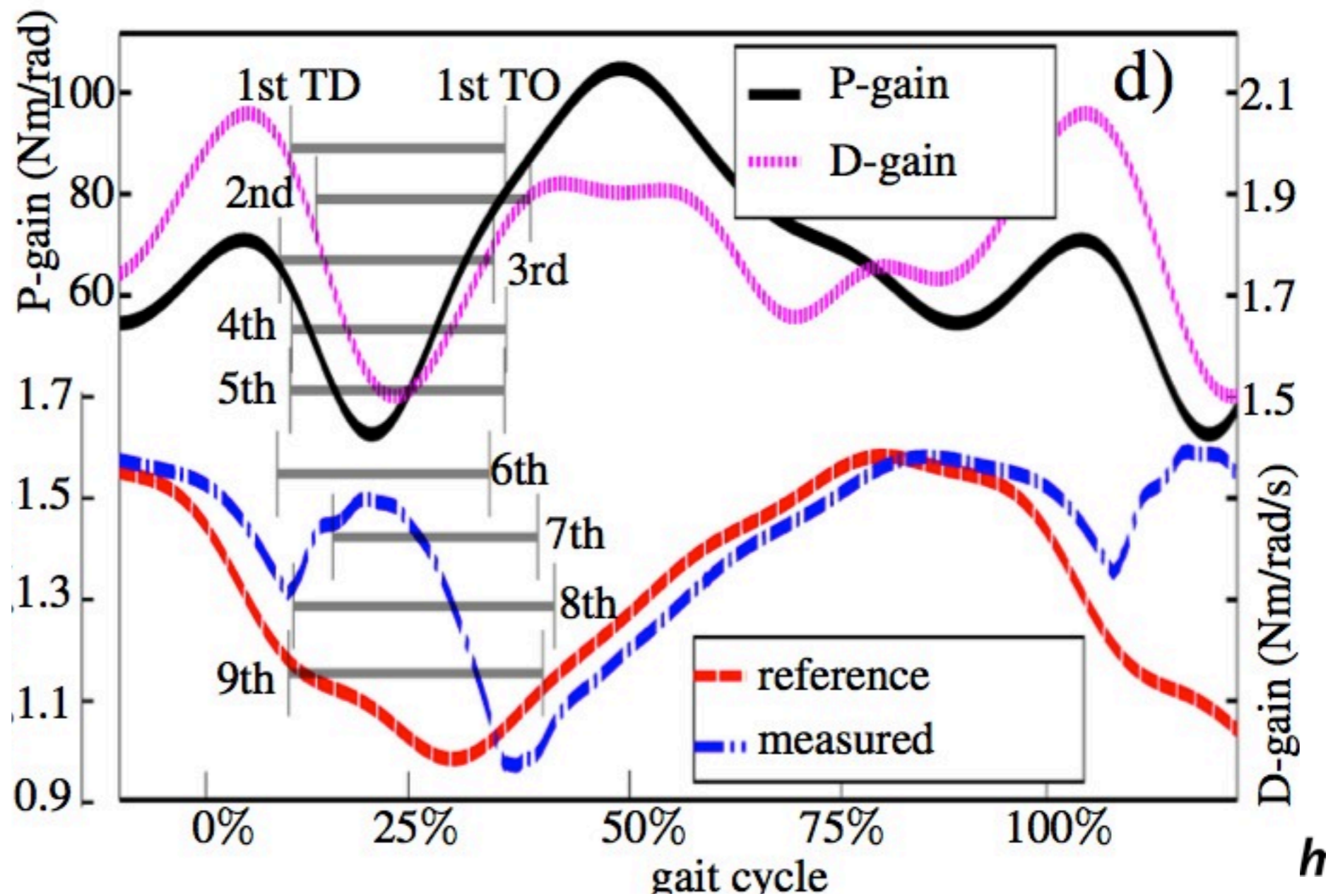
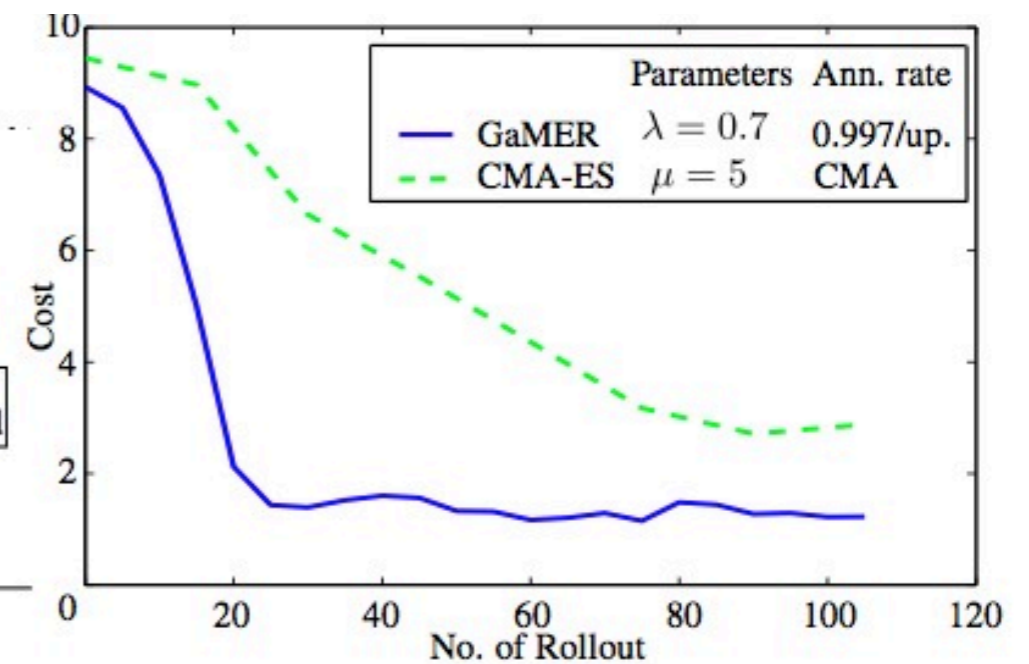
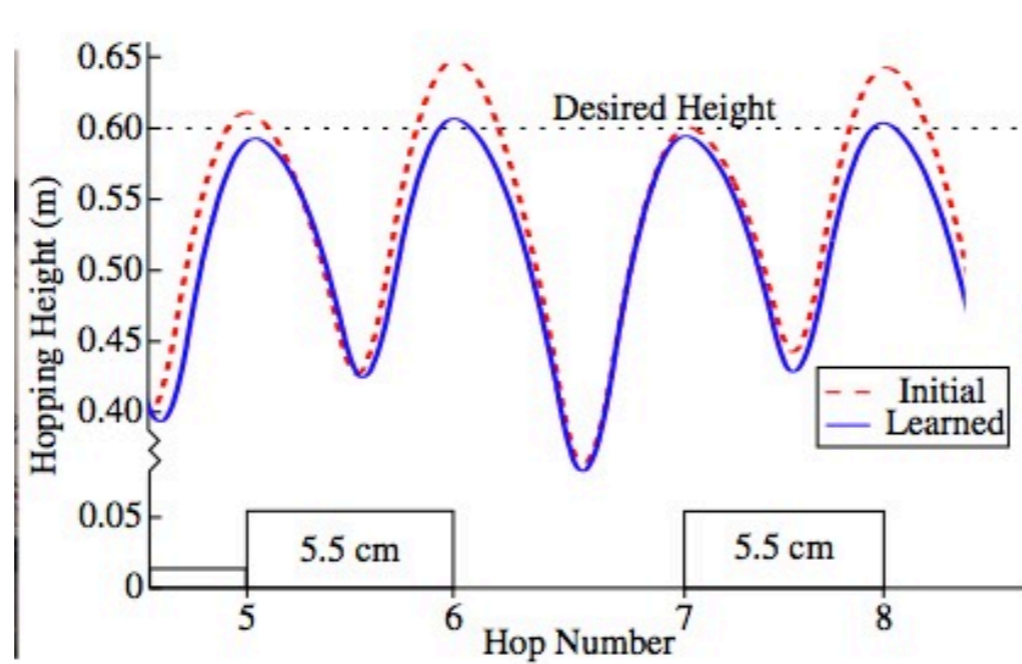
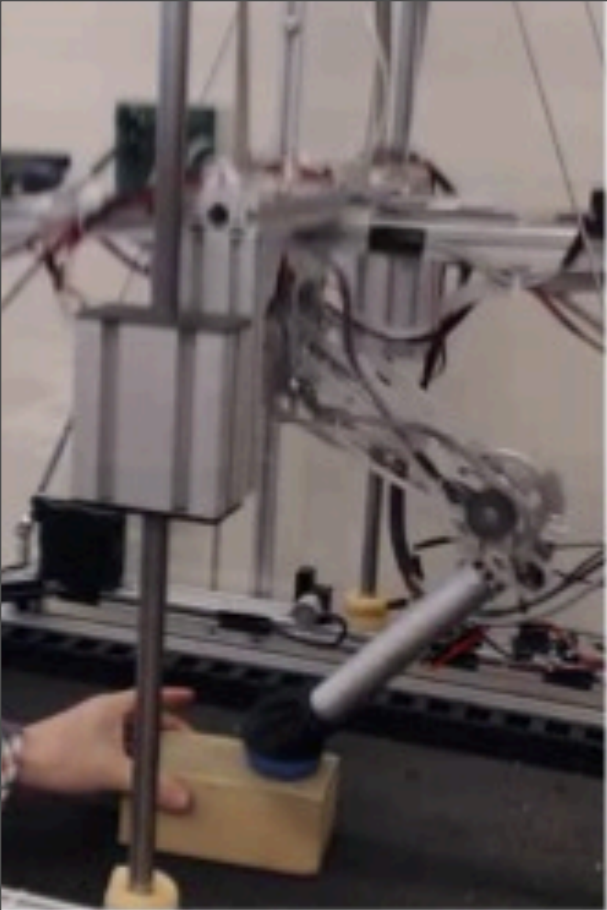


Fig. 4. Learned joint angle trajectories (center) and gain schedules (right) of the CBi arm after 0/6/100 updates. The gain schedules of only three joints have been depicted for sake of clarity.



Before Learning

After Learning (25 rollouts)





# PI<sup>2</sup> - General Algorithm

## PI2 for state feedback control

$$u_i(t, \mathbf{x}) = \text{grand sum } [\bar{\Upsilon}(t, \mathbf{x}) \circ \theta_i]$$

$$\bar{\Upsilon}(t, \mathbf{x}) = \Upsilon(t) \begin{bmatrix} 1 & \mathbf{x}^T \end{bmatrix} = \begin{bmatrix} e^{-\frac{1}{2} \frac{(t-\mu_n)^2}{\sigma_n^2}} & \begin{bmatrix} 1 & \mathbf{x}^T \end{bmatrix} \end{bmatrix}_{N \times (1+\dim[\mathbf{x}])}$$

‘feedforward’
‘state-feedback’

$\theta_i$  is a  $N \times (1+\dim[\mathbf{x}])$  parameter matrix for  $i$ th control input approximation

$\bar{\Upsilon}(t, \mathbf{x})$  basis function matrix  $N \times (1+\dim[\mathbf{x}])$

**Policy nonlinear in time, linear in states**



**Algorithm 10** General PI2 Algorithm**given**

The cost function:

$$J = \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} (q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u}) dt$$

general policy: parameter matrix

A Linear Model for function approximation:  $\mathbf{u}(t, \mathbf{x}) = [u_i(t, \mathbf{x})] = [\text{grand sum } [\tilde{\Upsilon}(t, \mathbf{x}) \circ \theta_i]]$ Initialize  $\{\theta_i\}$  with a sophisticated guessInitialize exploration noise standard deviation:  $c$ **repeat**Create  $K$  rollouts of the system with the perturbed parameter  $\{\theta_i\} + \{\epsilon_i\}$ ,  $\{\epsilon_{i,j}\} \sim \mathcal{N}(\mathbf{0}, c^2 \mathbf{I})$ **for** the  $i$ th control input **do****for** each time,  $s$  **do**Calculate the Return from starting time  $s$  for the  $k$ th rollout:

$$R(\tau^k(s)) = \Phi(\mathbf{x}(t_f)) + \int_s^{t_f} (q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u}) dt$$

Calculate  $\alpha$  from starting time  $s$  for the  $k$ th rollout:

$$\alpha^k(s) = \frac{\exp(-\frac{1}{\lambda} R(\tau^k(s)))}{\sum_{k=1}^K \exp(-\frac{1}{\lambda} R(\tau^k(s)))}$$

Calculate the time varying parameter increment  $\Delta\theta_i(s)$ :

$$\Delta\theta_i(s) = \sum_{k=1}^K \alpha^k(s) \frac{\Upsilon(s) \Upsilon^T(s)}{\Upsilon^T(s) \Upsilon(s)} \epsilon_i^k(s)$$

**end for****for** the  $j$ th column of  $\Delta\theta_i$  matrix,  $\Delta\theta_{i,j}$  **do**

Time-averaging the parameter vector

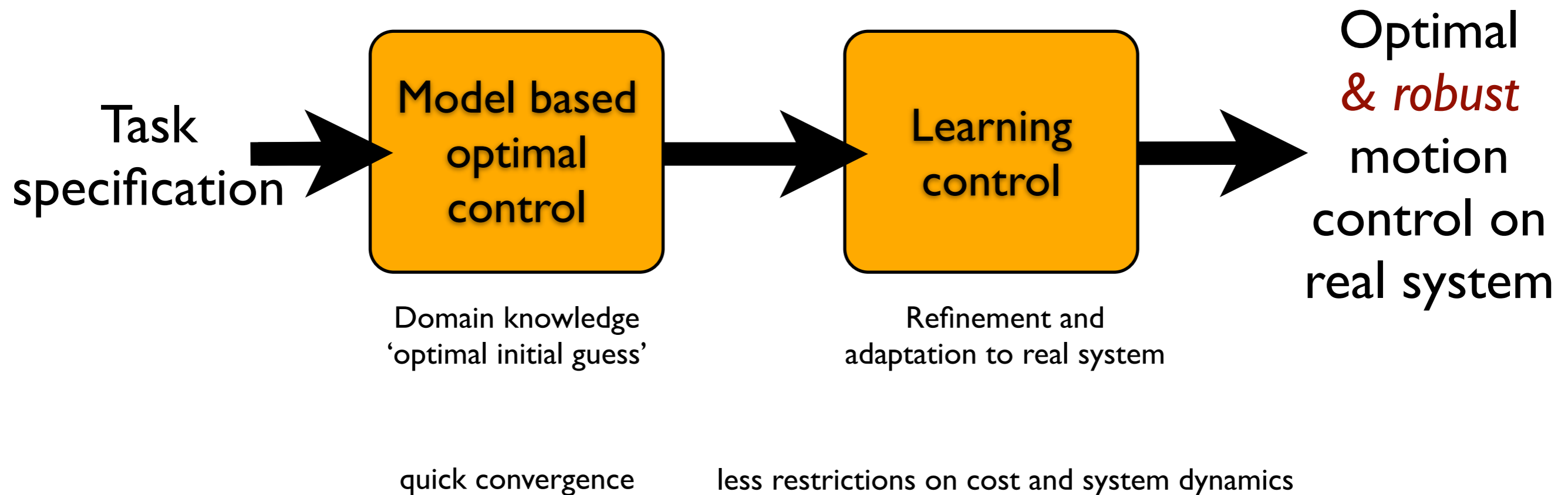
$$\Delta\theta_{i,j} = \left( \int_{t_0}^{t_f} \Delta\theta_{i,j}(s) \circ \Upsilon(s) ds \right) \cdot \left/ \int_{t_0}^{t_f} \Upsilon(s) ds \right.$$

**end for**Update parameter vector for control input  $i$ ,  $\theta_i$ :

$$\theta_i \leftarrow \theta_i + \omega \Delta\theta_i$$

**end for**- Decrease  $c$  for noise annealing**until** maximum number of iterations

# Combining optimal and learning control



# Example



# Rezero

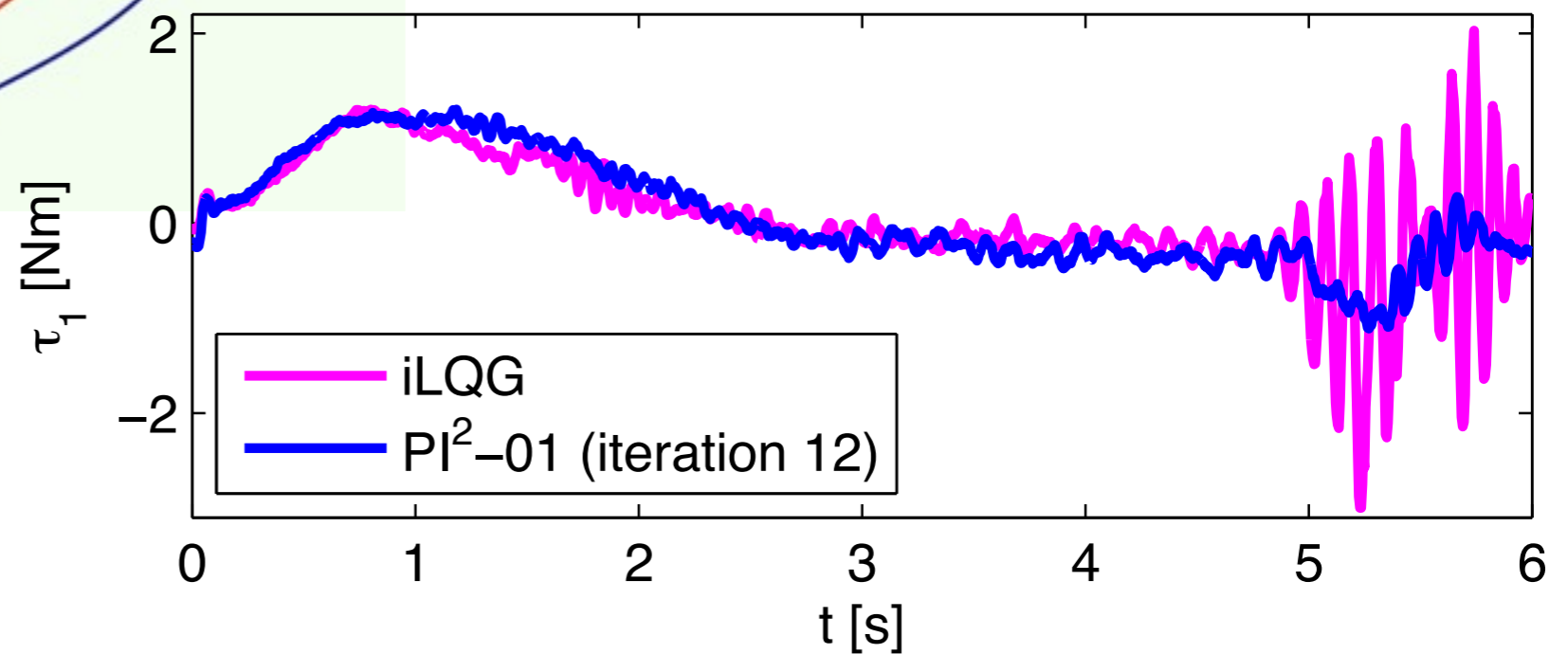
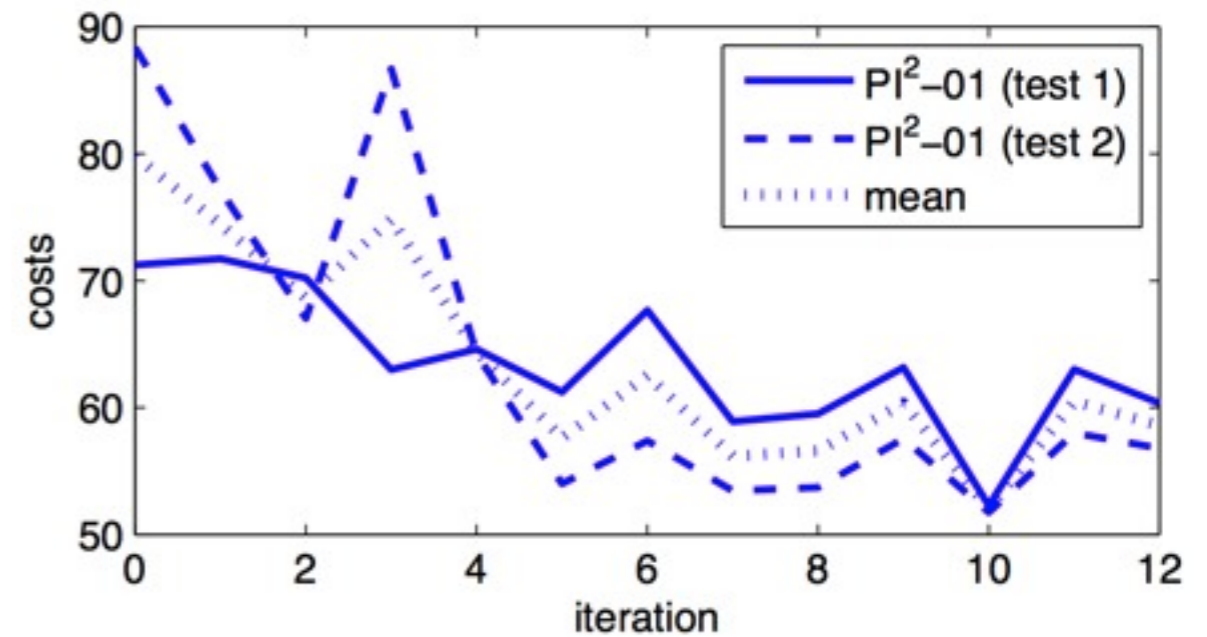
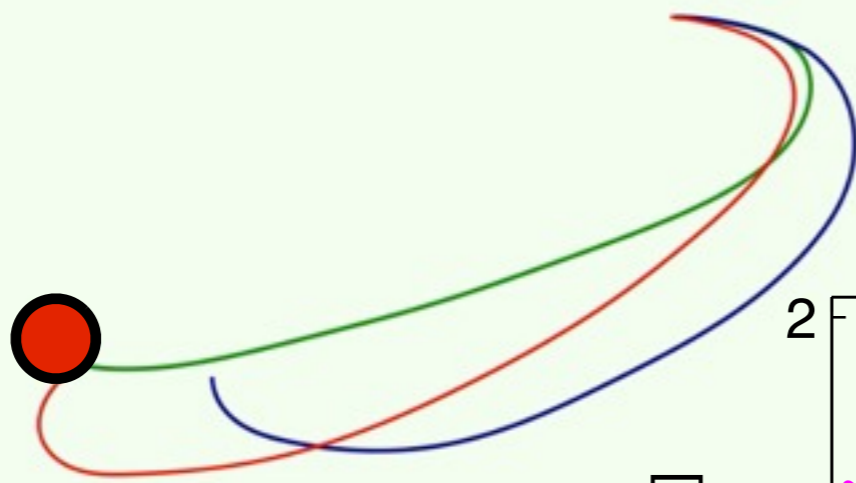


**Task: 'reach goal, using minimum torque'**

**Init: iLQG      Learning: PI<sup>2</sup>-01**

[Li & Todorov, 2005]

— Initialization - iLQG  
—  $PI^2-01$  Learning - Trial 5  
—  $PI^2-01$  Learning - Trial 12



Cost: 'reach goal, w. min. torque'

Init: iLQG

Learning:  $PI^2-01$

