# Optimal Control
# with Learned Forward Models

Pieter Abbeel
UC Berkeley

Jan Peters
TU Darmstadt

1

# Where we are?

**Reinforcement Learning Data**

$$\mathscr{D} = \{(\mathbf{x}_i, \mathbf{u}_i, \mathbf{x}_{i+1}, r_i)\}\}$$

$$\mathbf{x}' \sim \mathscr{P}^{\mathbf{u}}_{\mathbf{xx}'}$$
$$r \approx \mathscr{R}^{\mathbf{u}}_{\mathbf{xx}'}$$

$V^*(\mathbf{x})$

$\pi^*(\mathbf{u}|\mathbf{x})$

Now

$V^*(\mathbf{x})$

Policy Search

$\pi^*(\mathbf{u}|\mathbf{x})$

Value Function Methods

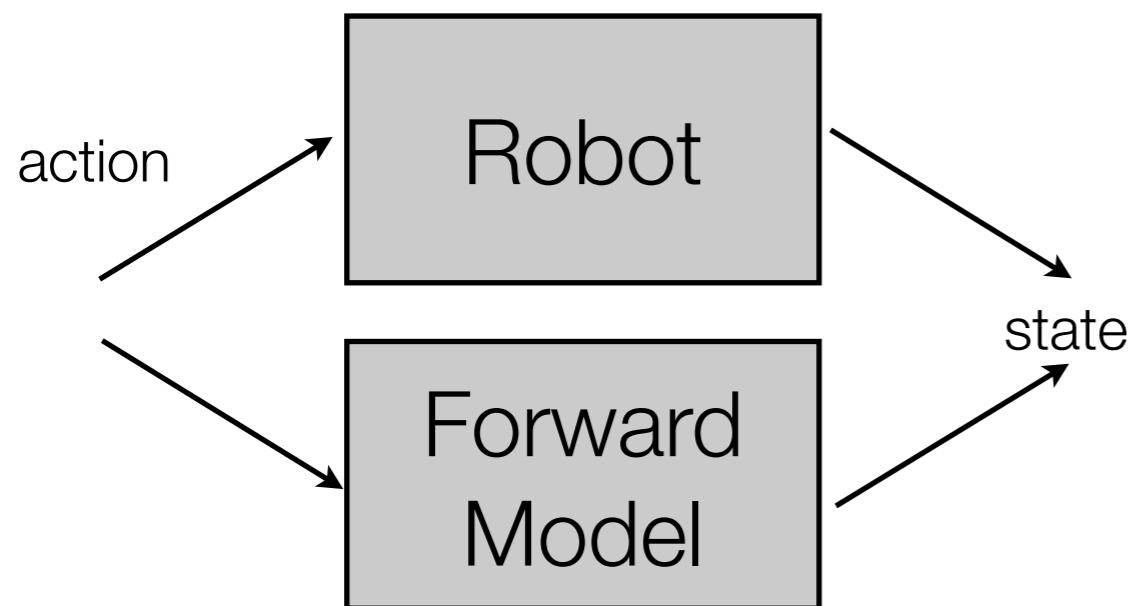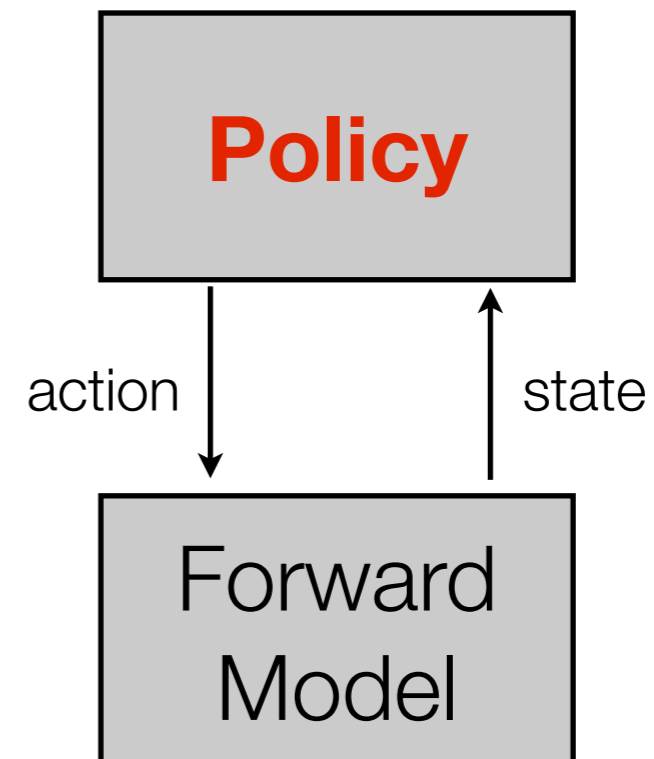$\pi^*(\mathbf{u}|\mathbf{x})$

Optimal Control with
Model Learning

2

# Goal of this Lecture

**1. Step:** Learn an Forward Model

**2. Step:** Use your favorite *Optimal Control Method* to get an optimal policy

action → | Robot |

| Forward Model | → state

| **Policy** |

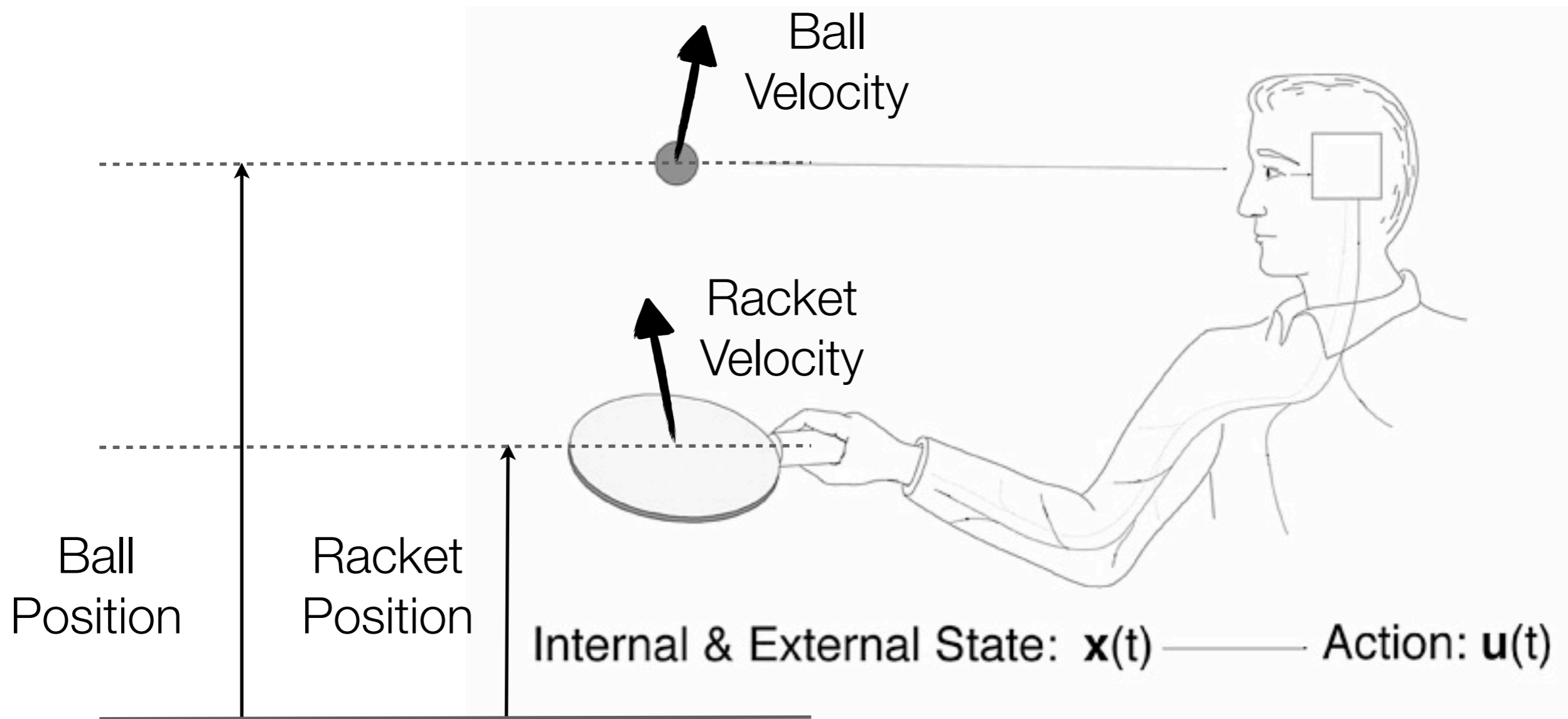action ↓   ↑ state

| Forward Model |

3

# Outline of the Lecture

▶ 1. Introduction to Optimal Control

2. Solving Linear-Quadratic Optimal Control Problems

3. Optimal Control with Learned Models

4. Hot story: Marc Deisenroth's PILCO Approach

5. Final Remarks

4

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



Ball Velocity

Racket Velocity

Ball Position

Racket Position

Internal & External State: $\mathbf{x}(t)$ ———— Action: $\mathbf{u}(t)$
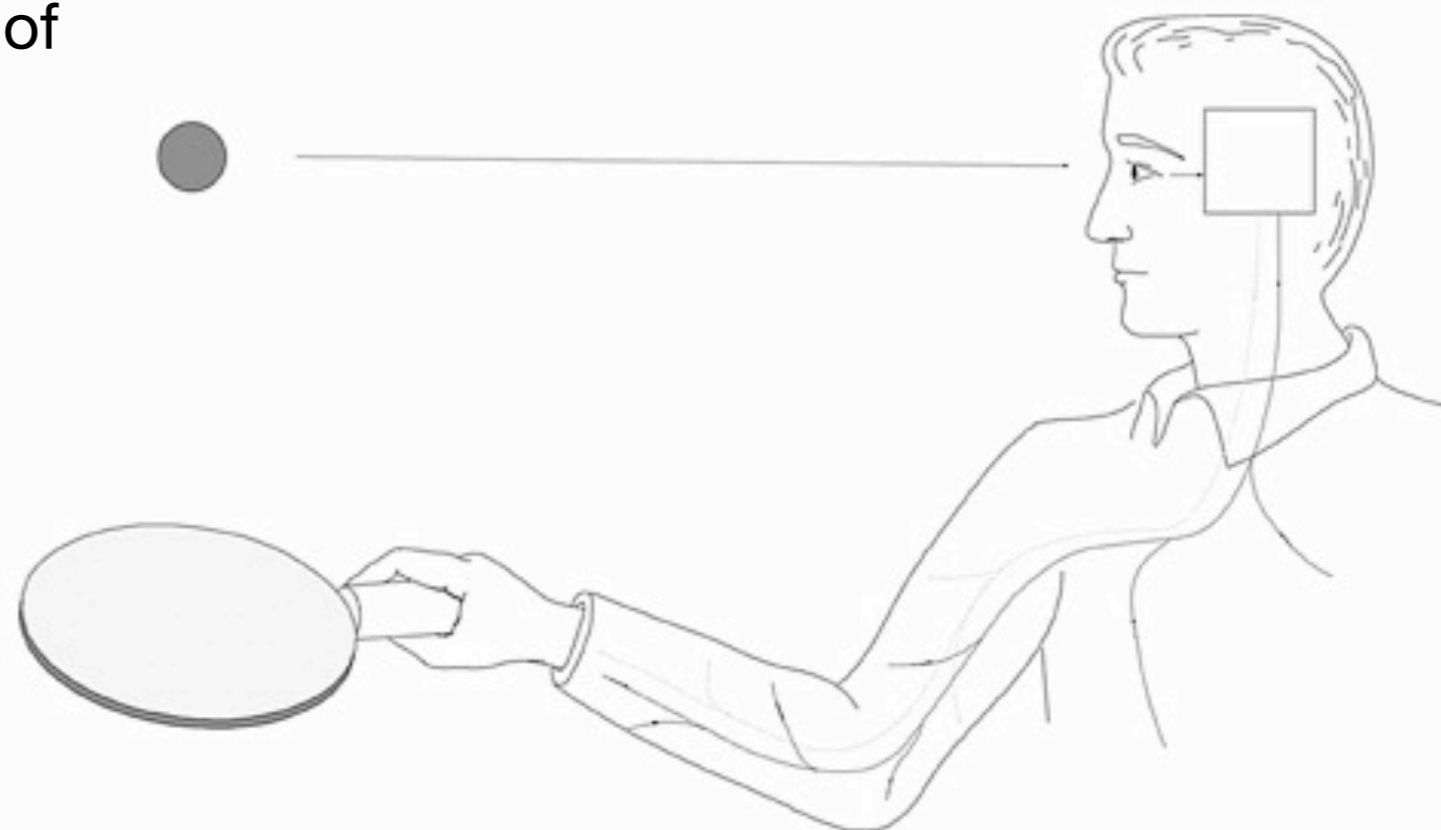
5

# Example: Ball Paddling

## What are the actions **u**?

- **All motor torques?**
  - If you do not have a model ...

- **Joint Accelerations?**
  - Perfect, if you have a good model ...
  - Maybe identify the proper degrees of freedom?

- **Accelerations in Operational Space?**
  - **Ideally!**
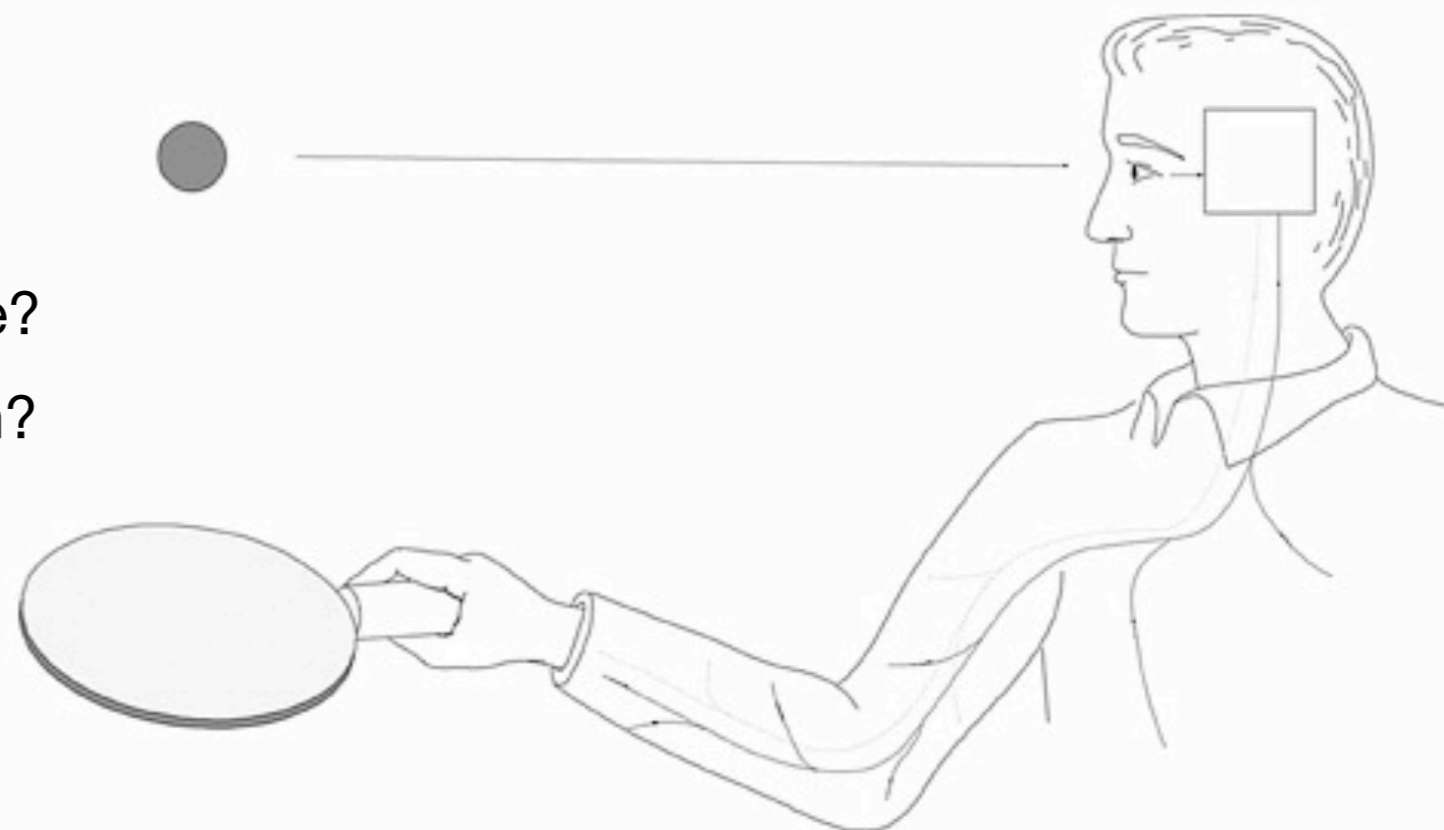  - ... but only if you have a good operational space control law!

# 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?
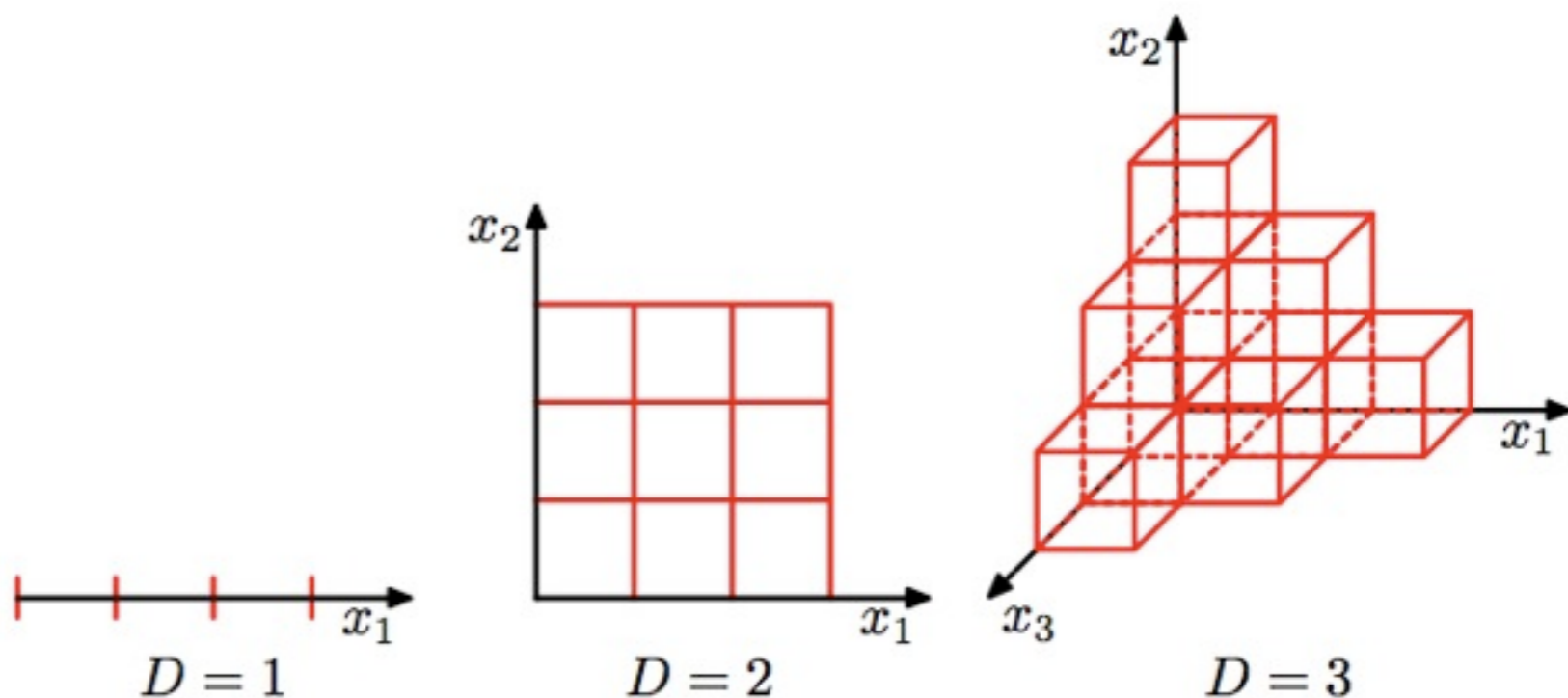
  - Additional punishments?

➡**All of these together do the job!**

7

# So can we get this to work?

- The state space has at least 12 dimensions.

- The action space has at least 3 dimensions.

- Can discretizations deal with such spaces?

- No!  Finding an Optimal Value Function is limited by the curse of dimensionality.



$D = 1$

$D = 2$

$D = 3$

**So how can you get this by RL?**

9

# Outline of the Lecture

1. Introduction to Optimal Control

2. Solving Linear-Quadratic Optimal Control Problems

3. Optimal Control with Learned Models

4. Hot story: Marc Deisenroth's PILCO Approach

5. Final Remarks

# Optimal Control Goal

- The goal of optimal control is find a policy $\boldsymbol{u} \sim \pi(\boldsymbol{x})$ such that

$$J(\pi) = \lim_{T \to \infty} \frac{1}{T} E \left\{ \sum_{k=0}^{T} r(\mathbf{x}_t, \mathbf{u}_t) \right\}$$

is maximal for a given

reward function such as $r(\mathbf{x}, \mathbf{u}) = -\mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{u}^T \mathbf{R} \mathbf{u}$

system: $\mathbf{x}' = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} + \mathbf{c}$ For Simplicity of Derivation! Nothing Changes!
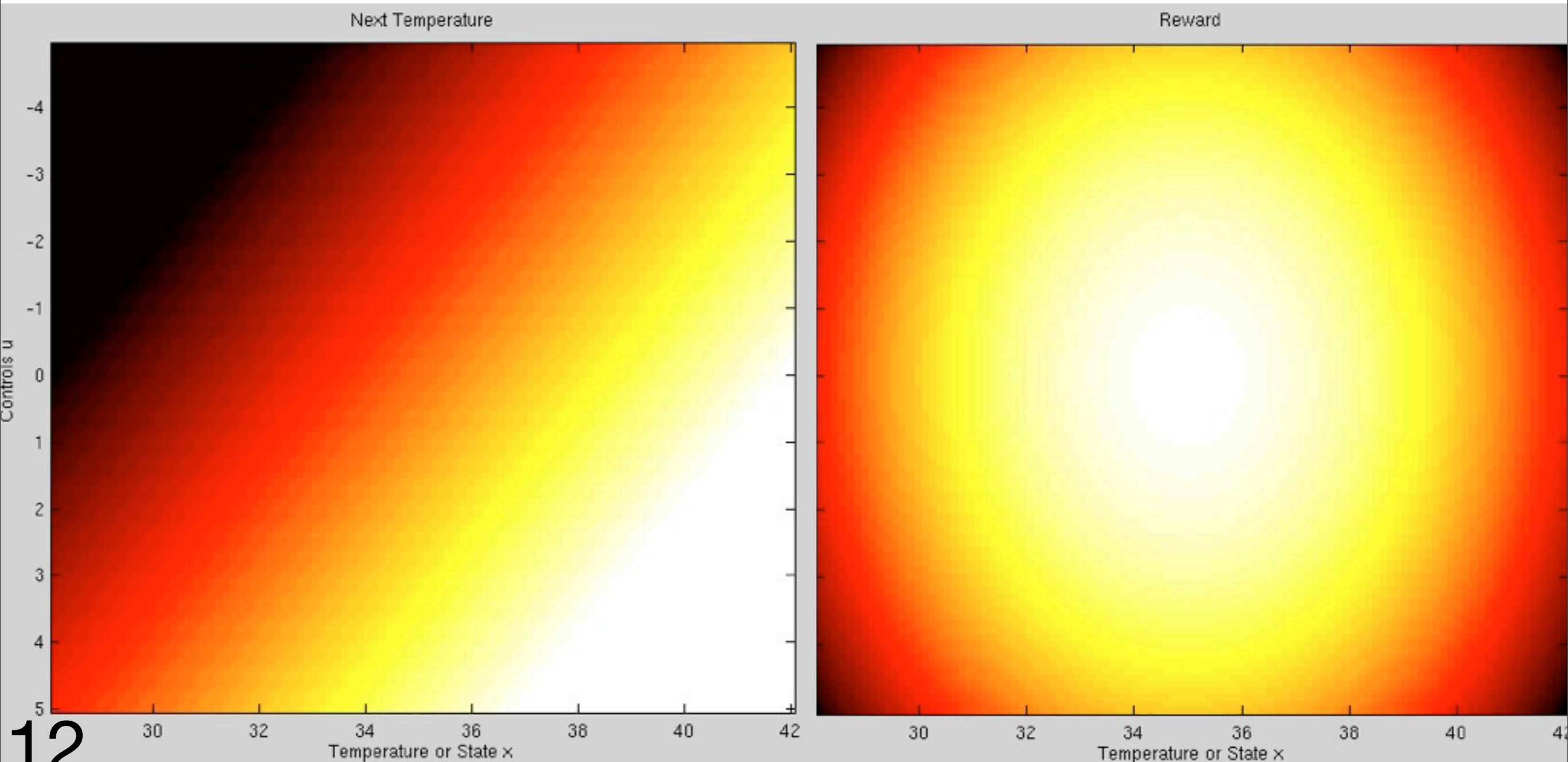
## ...so how do we solve this?

11

$$\mathbf{x}' = \mathbf{Ax} + \mathbf{Bu} + \cancel{\mathbf{c}} \qquad r(\mathbf{x}, \mathbf{u}) = -\mathbf{x}^T \mathbf{Qx} - \mathbf{u}^T \mathbf{Ru}$$



**12**

# Bellman' Recipe: Steps 1+2

1. At the last step, we have the value function

$$V_T^*(\mathbf{x}) = 0$$

2. For *t=T-1*, compute optimal policy such that

$$\pi_t^*(u|x) = \operatorname{argmax}_\pi \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\}$$

determined by

$$\frac{d}{d\mathbf{u}} \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\} = 0$$

$$\frac{d}{d\mathbf{u}} \left\{ -\mathbf{x}^T\mathbf{Q}\mathbf{x} - \mathbf{u}^T\mathbf{R}\mathbf{u}) \right\} = 0$$

$$\mathbf{u}^* = 0$$



Value Function

13

3.Obtain next value function $\quad V_t^*(x) = \max_\pi \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\}$

$$
\begin{aligned}
V_{t+1}^*(\mathbf{x}) &= r(\mathbf{x}, \mathbf{u}^*) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u}^*)) \\
&= -\mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{u}^{*T} \mathbf{R} \mathbf{u}^* \\
&= -\mathbf{x}^T \mathbf{Q} \mathbf{x}
\end{aligned}
$$

4.**As not converged, go back to Step 2.**

14

# Bellman' Recipe: Step 2 *again*!

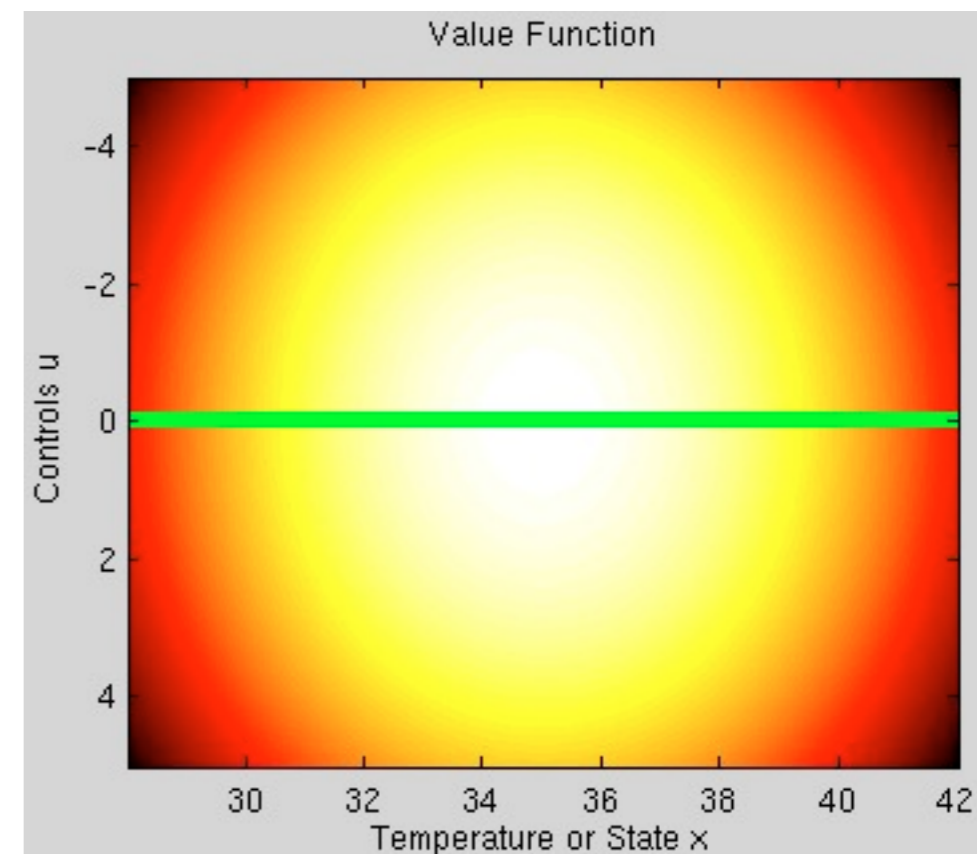2.For *t<T-1*, compute optimal policy such that

$$\pi_t^*(u|x) = \mathrm{argmax}_\pi \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\}$$

determined by

$$\frac{d}{d\mathbf{u}} \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\} = 0$$

$$\frac{d}{d\mathbf{u}} \left\{ -\mathbf{x}^T\mathbf{Q}\mathbf{x} - \mathbf{u}^T\mathbf{R}\mathbf{u} - f(\mathbf{x}, \mathbf{u})^T\mathbf{P}_{t+1}f(\mathbf{x}, \mathbf{u}) \right\} = 0$$

$$\frac{d}{d\mathbf{u}} \left\{ -\mathbf{x}^T\mathbf{Q}\mathbf{x} - \mathbf{u}^T\mathbf{R}\mathbf{u} - (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u})^T\mathbf{P}_{t+1}(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}) \right\} = 0$$

$$\frac{d}{d\mathbf{u}} \left\{ -\mathbf{R}\mathbf{u} - \mathbf{B}^T\mathbf{P}_{t+1}(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}) \right\} = 0$$

which implies

Policy Parameters!

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

15

3.Obtain next value function $\quad V_t^*(x) = \max_\pi \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\}$

$$
\begin{aligned}
V_t^*(\mathbf{x}) \quad &= \quad r(\mathbf{x}, \mathbf{u}^*) + V_{t+1}^*(\mathbf{Ax} + \mathbf{Bu}^*) \\
&= \quad -\mathbf{x}^T \mathbf{Qx} - (\boldsymbol{\theta}_t \mathbf{x})^T \mathbf{R} (\boldsymbol{\theta}_t \mathbf{x}) \\
&\quad -(\mathbf{Ax} + \mathbf{B}\boldsymbol{\theta}_{t+1} \mathbf{x})^T \mathbf{P}_{t+1} + (\mathbf{Ax} + \mathbf{B}\boldsymbol{\theta}_t \mathbf{x}) \\
&= \quad -\mathbf{x}^T [\mathbf{Q} - \boldsymbol{\theta}_t^T \mathbf{R} \boldsymbol{\theta}_t \\
&\quad -(\mathbf{A} + \mathbf{B}\boldsymbol{\theta}_t)^T \mathbf{P}_{t+1} + (\mathbf{A} + \mathbf{B}\boldsymbol{\theta}_t)] \mathbf{x} \\
&= \quad -\mathbf{x}^T \mathbf{P}_t \mathbf{x}
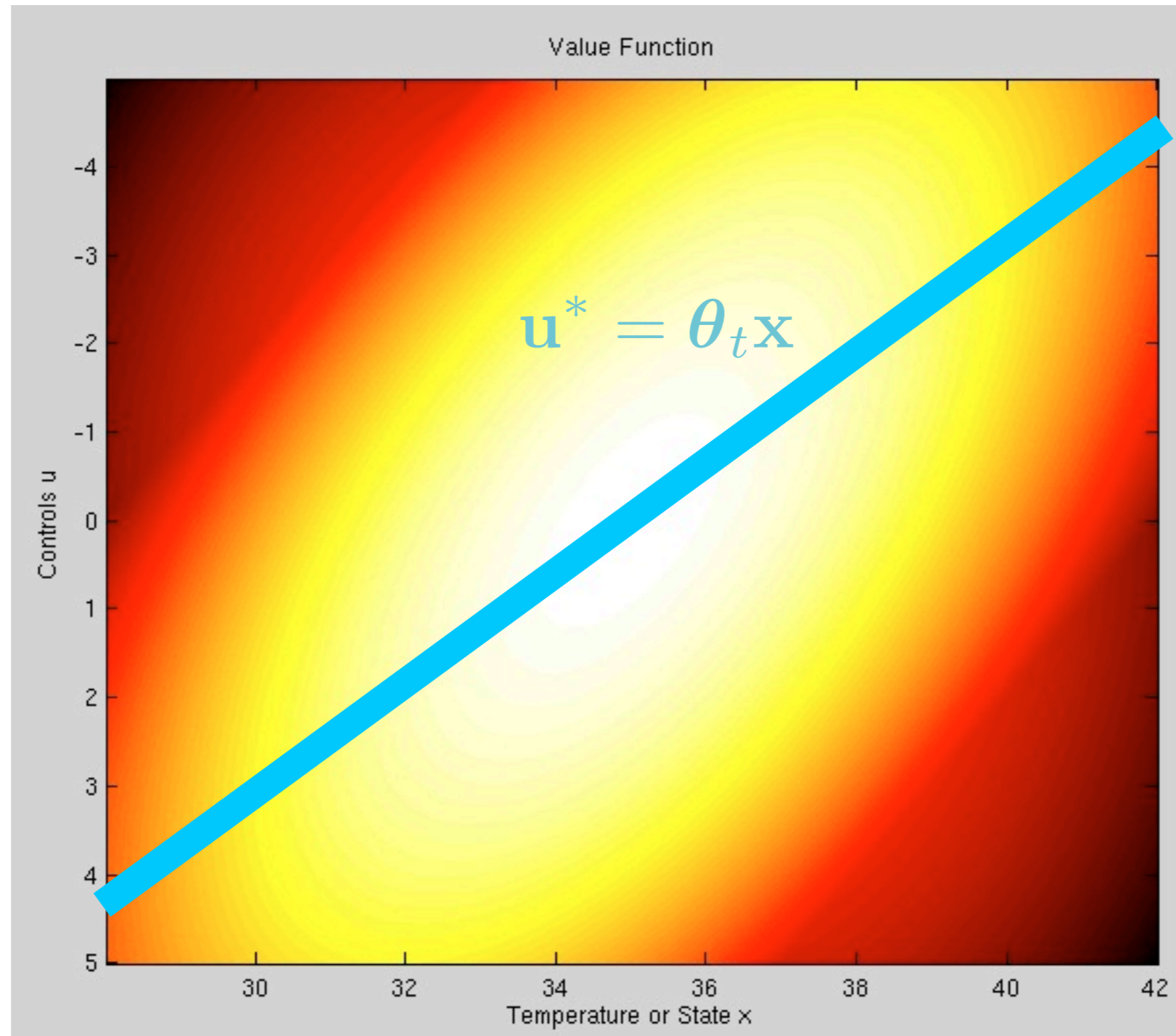\end{aligned}
$$

4. **We have a converged to _Recursion_:**

$$
\mathbf{P}_t = -\mathbf{Q} - \boldsymbol{\theta}_{t+1}^T \mathbf{R} \boldsymbol{\theta}_{t+1} - (\mathbf{A} + \mathbf{B}\boldsymbol{\theta}_{t+1})^T \mathbf{P}_{t+1} + (\mathbf{A} + \mathbf{B}\boldsymbol{\theta}_{t+1})
$$

16 $\quad \boldsymbol{\theta}_t = -(\mathbf{R} + \mathbf{B}^T \mathbf{P}_{t+1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{P}_{t+1} \mathbf{A}$

Value Function

$$\mathbf{u}^* = \boldsymbol{\theta}_t \mathbf{x}$$

# Outline of the Lecture

1. Introduction to Optimal Control

2. Solving Linear-Quadratic Optimal Control Problems
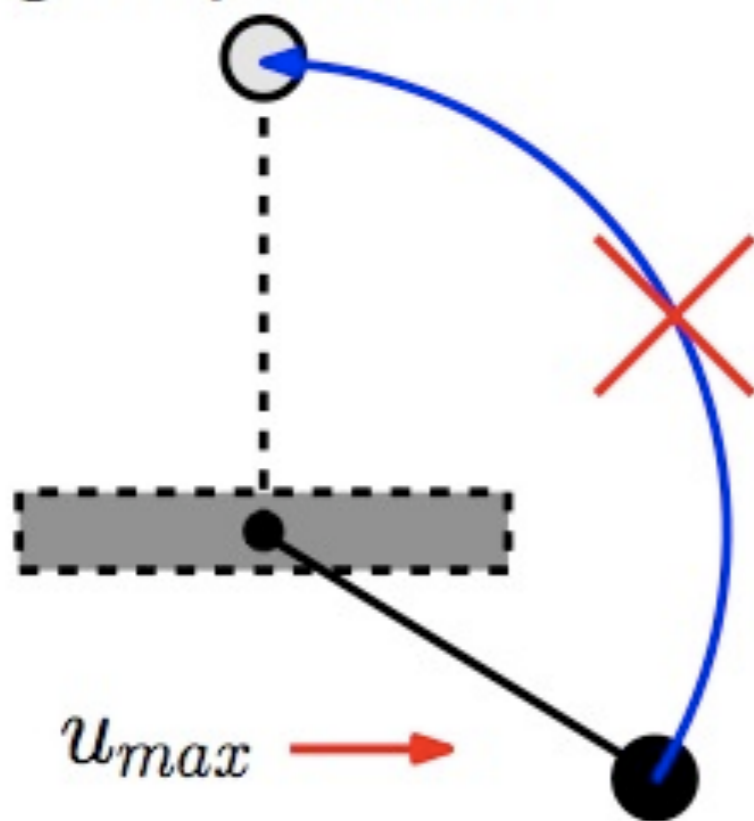
3. Optimal Control with Learned Models

4. Hot story: Marc Deisenroth's PILCO Approach

5. Final Remarks

Thursday, May 17, 2012

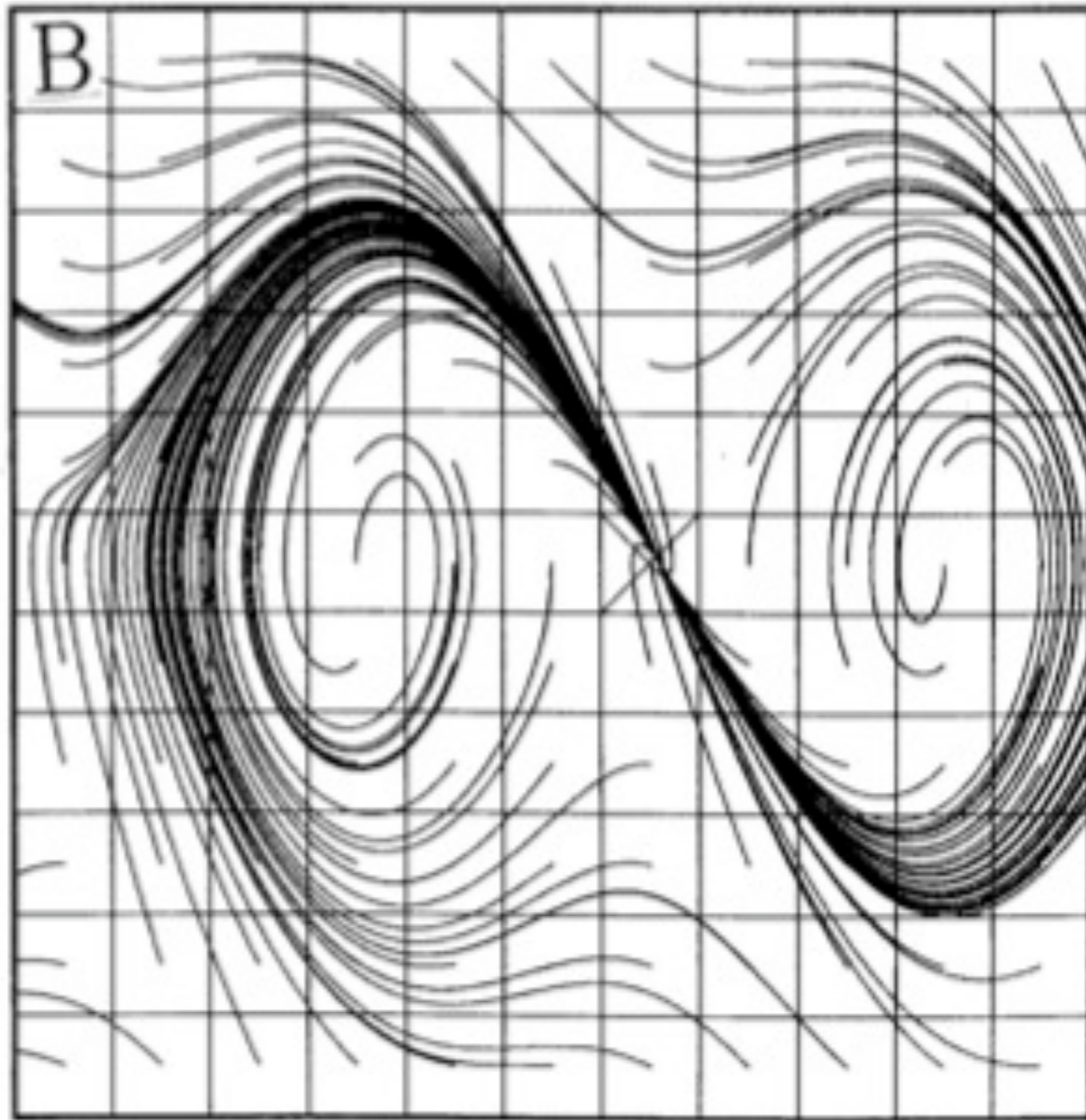# Example: Swing-Up

goal position



$u_{max}$ →

## 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

$$r(\mathbf{x}, \mathbf{u}) = -\mathbf{x}_k^T \operatorname{diag}(1, 0.1)\mathbf{x}_k - 0.2\, u_k^2$$

19

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

**... we can just look ahead!**

Atkeson & Schaal, 1995

# Local Solutions

- Every smooth function can be modeled with a Taylor expansion

$$f(\mathbf{x}) = f(\mathbf{a}) + \left.\frac{df}{d\mathbf{x}}\right|_{\mathbf{x}=\mathbf{a}} (\mathbf{x} - \mathbf{a}) + \frac{1}{2}(\mathbf{x} - \mathbf{a})^T \left.\frac{d^2 f}{d\mathbf{x}^2}\right|_{\mathbf{x}=\mathbf{a}} (\mathbf{x} - \mathbf{a}) + \ldots$$

- Hence, we can also approximate:

$$\mathbf{x}' \approx f(\hat{\mathbf{x}}, \hat{\mathbf{u}}) + \left.\frac{df}{d\mathbf{x}}\right|_{\hat{\mathbf{x}},\hat{\mathbf{u}}} (\mathbf{x} - \hat{\mathbf{x}}) + \left.\frac{df}{d\mathbf{u}}\right|_{\hat{\mathbf{x}},\hat{\mathbf{u}}} (\mathbf{u} - \hat{\mathbf{u}})$$

$$= \mathbf{a}_t^0 + \mathbf{A}_t(\mathbf{x} - \hat{\mathbf{x}}) + \mathbf{B}_t(\mathbf{u} - \hat{\mathbf{u}})$$

$$r \approx r(\hat{\mathbf{x}}, \hat{\mathbf{u}}) + \left[\begin{array}{c} \frac{dr}{d\mathbf{x}} \\ \frac{dr}{d\mathbf{u}} \end{array}\right]^T \left[\begin{array}{c} \mathbf{x} - \hat{\mathbf{x}} \\ \mathbf{u} - \hat{\mathbf{u}} \end{array}\right] + \frac{1}{2} \left[\begin{array}{c} \mathbf{x} - \hat{\mathbf{x}} \\ \mathbf{u} - \hat{\mathbf{u}} \end{array}\right]^T \left[\begin{array}{cc} \frac{d^2 r}{d\mathbf{x}^2} & \frac{d^2 r}{d\mathbf{x}d\mathbf{u}} \\ \frac{d^2 r}{d\mathbf{x}d\mathbf{u}} & \frac{d^2 r}{d\mathbf{u}^2} \end{array}\right] \left[\begin{array}{c} \mathbf{x} - \hat{\mathbf{x}} \\ \mathbf{u} - \hat{\mathbf{u}} \end{array}\right]$$

Atkeson & Schaal, 1995

21

# Bellman' Recipe: Steps 1-4

1. At the last step, we have the value function

$$V_T^*(\mathbf{x}) = 0$$

2. For *t=T-1*, compute optimal policy such that

$$\pi_t^*(u|x) = \mathrm{argmax}_\pi \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\}$$

gives $\mathbf{u} = \hat{\mathbf{u}}$.

3. Obtain next value function $V_t^*(x) = \max_\pi \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\}$

$$V_t^*(\mathbf{x}) = -r_{t+1} - (\mathbf{x} - \hat{\mathbf{x}}_t)^T \mathbf{Q}_{t+1}(\mathbf{x} - \hat{\mathbf{x}}_t)$$

4. **As not converged, go back to Step 2.**

22

# Bellman' Recipe: Step 2-4 *again*!

2. For *t<T-1*, compute optimal policy such that

$$\pi_t^*(u|x) = \mathrm{argmax}_\pi \left\{ r(\mathbf{x}, \mathbf{u}) + V_{t+1}^*(f(\mathbf{x}, \mathbf{u})) \right\}$$
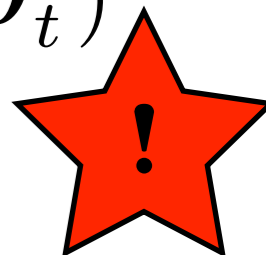
determined by

$$\mathbf{u}^* = \hat{\mathbf{u}}_t - (\mathbf{R}_t + \mathbf{B}_t^T \mathbf{P}_t \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{P}_{t+1}(\mathbf{a}_t^0 + \mathbf{A}_t(\mathbf{x} - \hat{\mathbf{x}}_t)) = \boldsymbol{\theta}_t^1(\mathbf{x} - \hat{\mathbf{x}}_t) + \boldsymbol{\theta}_t^0$$

3. Obtain the recursions

$$\boldsymbol{\theta}_t^1 = \hat{\mathbf{u}}_t - (\mathbf{R}_t + \mathbf{B}_t^T \mathbf{P}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{P}_{t+1} \mathbf{A}_t(\mathbf{x} - \hat{\mathbf{x}}_t)$$

$$\boldsymbol{\theta}_t^0 = \hat{\mathbf{u}}_t - (\mathbf{R}_t + \mathbf{B}_t^T \mathbf{P}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{P}_{t+1} \mathbf{a}_t^0$$

$$\mathbf{P}_t = -\mathbf{Q}_t - \boldsymbol{\theta}_t^T \mathbf{R}_t \boldsymbol{\theta}_t^1 + (\mathbf{A} + \mathbf{B}_t \boldsymbol{\theta}_t^1) \mathbf{P}_{t+1}(\mathbf{A} + \mathbf{B}_t \boldsymbol{\theta}_t^1)$$
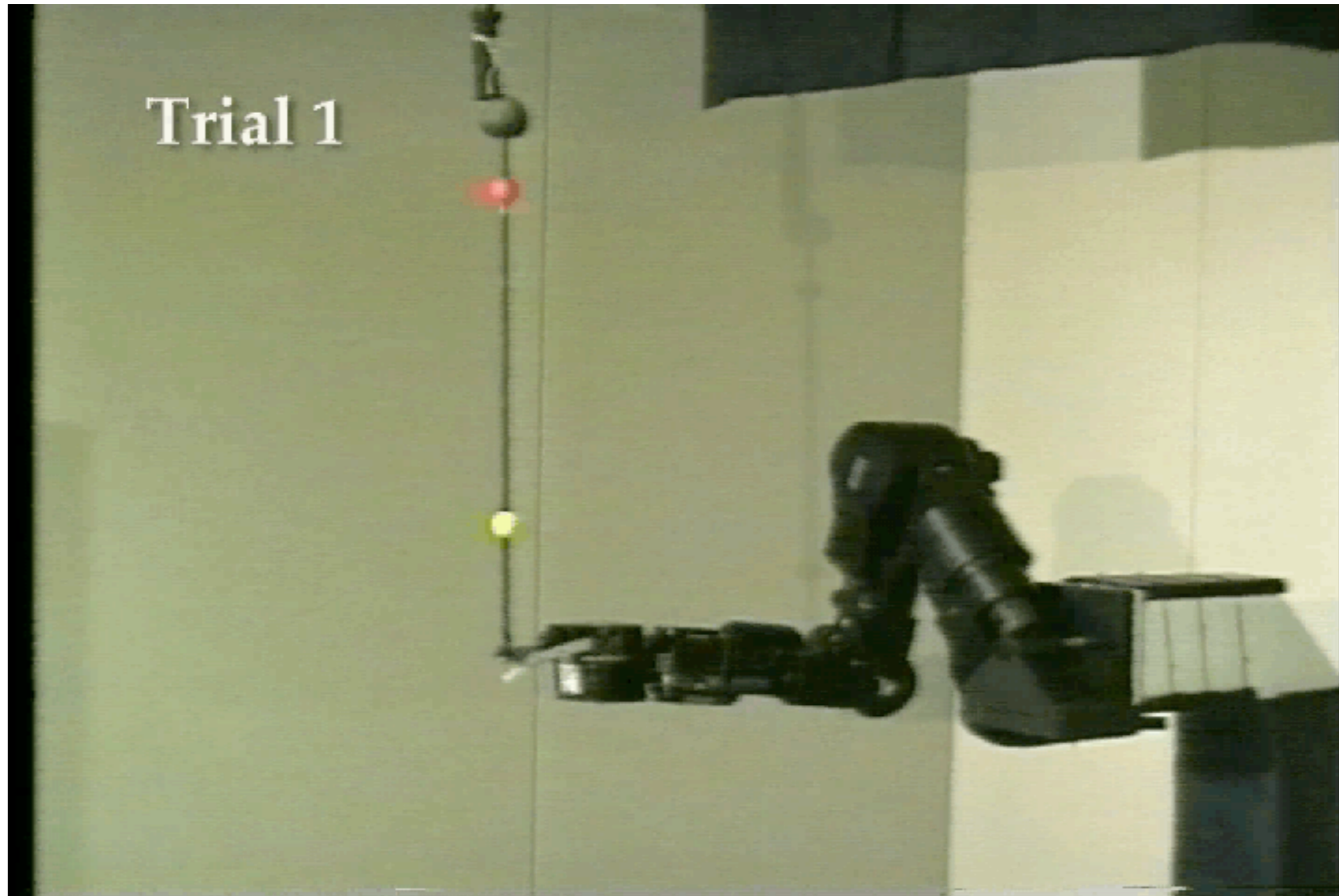
23

# How do we get the Optimal Policy

1. Forward Propagation: Run Simulator to Obtain Linearizations

2. Backward Solution: Compute Optimal Control Law

3. If not converged, go to 1.

# Model Learning with subsequent Policy Optimization

Atkeson & Schaal, 1996; Schaal, 1997

# Model Learning with subsequent Policy Optimization



Atkeson & Schaal, 1996; Schaal, 1997

26

# Outline of the Lecture

1. Introduction to Optimal Control

2. Solving Linear-Quadratic Optimal Control Problems

3. Optimal Control with Learned Models
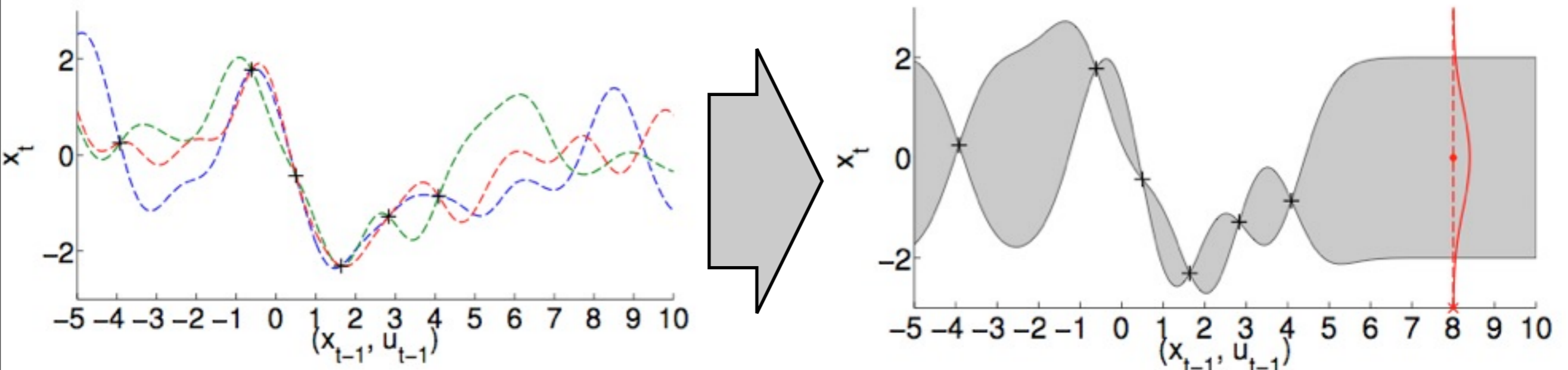
▶ 4. Hot story: Marc Deisenroth's PILCO Approach

5. Final Remarks

# Probabilistic Forward Models

Marc
Deisenroth

- Many forward models explain measured data.

- Choosing a bad model destroys will cause an optimization bias.

- Can we average ensure robustness towards bad approximations?

➡ Yes! Even in a Bayesian way with Gaussian Process Regression!

Deisenroth, Fox, Rasmussen, R:SS 2011

# Basic Idea

Marc
Deisenroth

- With a GP, you can compute the distributions over all future states based on all forward models weighted by their likelihood.
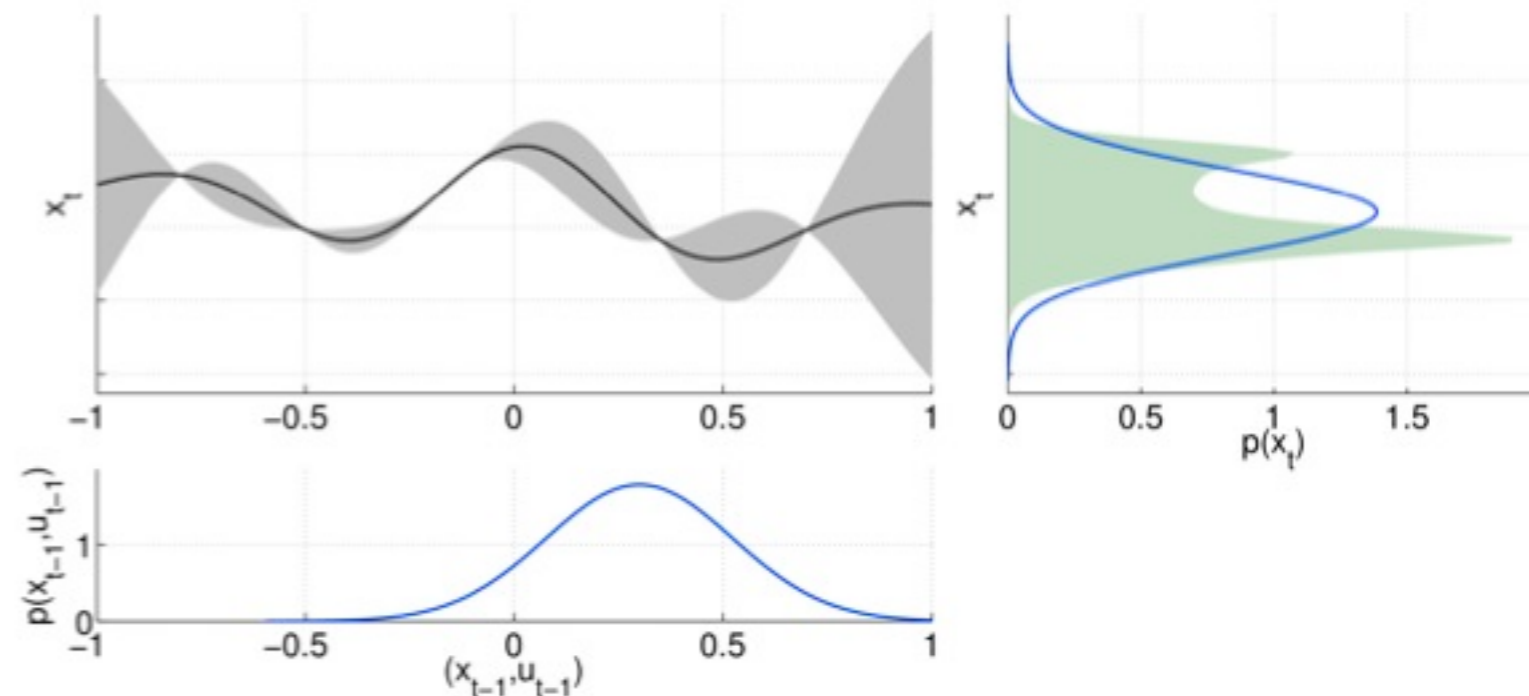
$$p(\mathbf{x}_1), \ldots, p(\mathbf{x}_T)$$

- The propagation is still approximate and done by moment matching

- From the distribution:

  ‣ Expected Return: $J^\pi(\boldsymbol{\theta})$

  ‣ Gradient: $\mathrm{d}J^\pi(\boldsymbol{\theta})/\mathrm{d}\boldsymbol{\theta}$

➡ We can do policy updates!

**29**

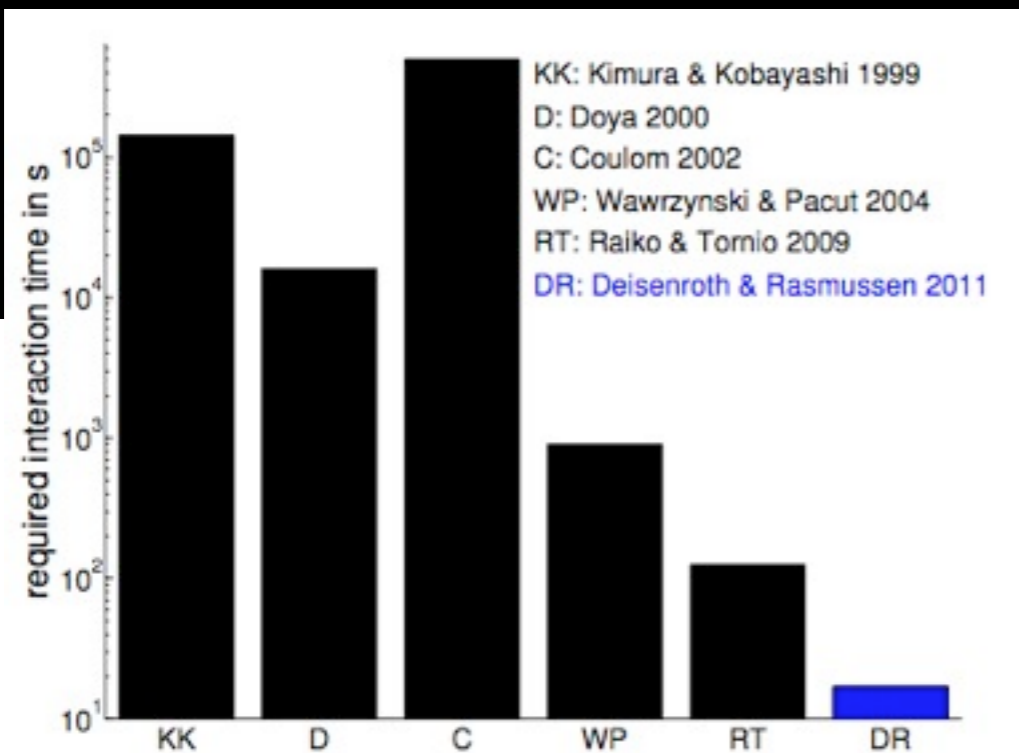Deisenroth, Fox, Rasmussen, R:SS 2011

# Applications

Marc
Deisenroth



Marc Peter Deisenroth, Carl Edward Rasmussen, Dieter Fox

Learning to Control a Low-Cost Robotic Manipulator
using Data-Efficient Reinforcement Learning

R:SS 2011

KK: Kimura & Kobayashi 1999
D: Doya 2000
C: Coulom 2002
WP: Wawrzynski & Pacut 2004
RT: Raiko & Tornio 2009
DR: Deisenroth & Rasmussen 2011

**30**

Deisenroth, Fox, Rasmussen, R:SS 2011

# Outline of the Lecture

1. Introduction to Optimal Control

2. Solving Linear-Quadratic Optimal Control Problems

3. Optimal Control with Learned Models

4. Hot story: Marc Deisenroth's PILCO Approach

▶ 5. Final Remarks

# Conclusions

- You have learned about optimal control today!

- Only two cased are solvable: linear & discrete!
  - Linear scales but does not generalize.
  - Discrete generalizes but does not scale.

- Using Learned Models, you can compute at least optimal "policy tubes".

- If you have many many tubes, in good regions, you have a policy.

- We will continue with Value Function and Policy Search Methods.

32

# Further Reading

- C. G. Atkeson (1994), Using Local Trajectory Optimizers to Speed Up Global Optimization in Dynamic Programming, Proceedings, Neural Information Processing Systems, Denver, Colorado, December, 1993, In: Neural Information Processing Systems 6, J. D. Cowan, G. Tesauro, and J. Alspector, eds. Morgan Kaufmann, 1994.

- Schaal, S. (1997). "Learning from demonstration". In: M.C. Mozer, M. Jordan, & T. Petsche (eds.), Advances in Neural Information Processing Systems 9, pp.1040-1046. Cambridge, MA: MIT Press

- Marc P. Deisenroth, Carl E. Rasmussen, Dieter Fox (2011). Learning to Control a Low-Cost Robotic Manipulator Using Data-Efficient Reinforcement Learning, Robotics: Science & Systems (RSS 2011)

33