

A Tour of Reinforcement Learning

The View from Continuous Control

Benjamin Recht
University of California, Berkeley



trustable, scalable, predictable

Control Theory



~~Reinforcement Learning~~ is the study of how to use past data to enhance the future manipulation of a dynamical system

Disciplinary Biases

AE/CE/EE/ME



continuous

model → *action*

IEEE Transactions

CS



discrete

data → *action*

Science Magazine

Disciplinary Biases

AE/CE/EE/ME

CS

Today's talk will try to unify these camps and point out how to merge their perspectives.

continuous

model → *action*

IEEE Transactions

discrete

data → *action*

Science Magazine



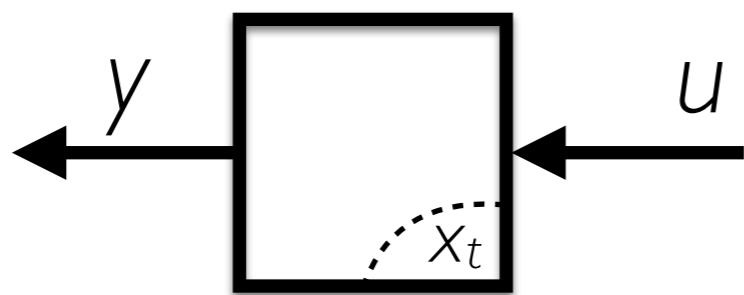
Main research challenge: What are the fundamental limits of learning systems that interact with the physical environment?

**How well must we understand a system
in order to control it?**

theoretical
foundations

- statistical learning theory
- robust control theory
- core optimization

Control theory is the study of dynamical systems with inputs



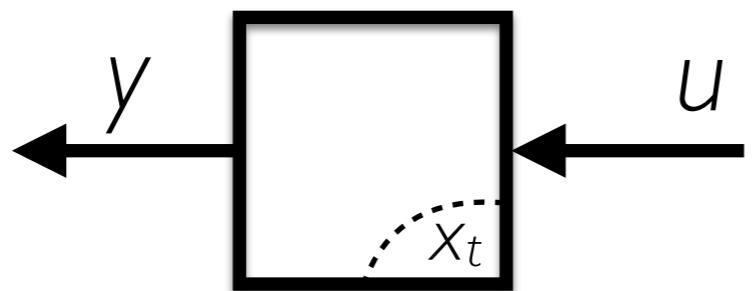
$$x_{t+1} = f(x_t, u_t)$$

x_t is called the *state*, and the dimension of the state is called the *degree*, d .

u_t is called the *input*, and the dimension is p .

Reinforcement

~~Learning Control theory~~ is the study of ^{discrete} dynamical systems with inputs



$$p(x_{t+1} \mid \text{past}) = p(x_{t+1} \mid x_t, u_t)$$

Markov Decision Process (MDP)

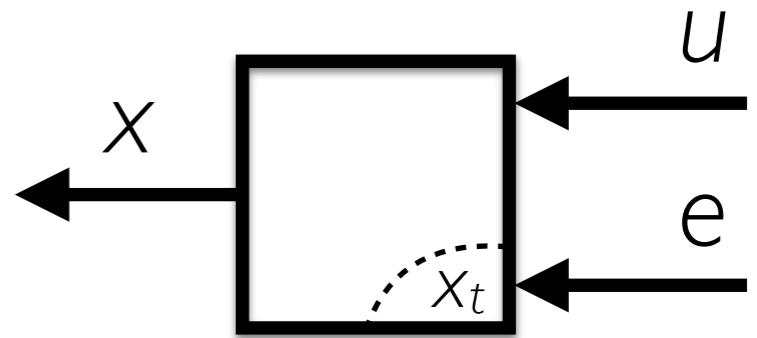
x_t is the *state*, and it takes values in $[d]$

u_t is called the *input*, and takes values in $[p]$.

Optimal control

$$\text{minimize} \quad \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right]$$

$$\text{s.t.} \quad \begin{aligned} x_{t+1} &= f_t(x_t, u_t, e_t) \\ u_t &= \pi_t(\tau_t) \end{aligned}$$



C_t is the cost. If you maximize, it's called a reward.

e_t is a noise process

f_t is the state-transition function

$\tau_t = (u_1, \dots, u_{t-1}, x_0, \dots, x_t)$ is an observed trajectory

$\pi_t(\tau_t)$ is the policy. This is the optimization decision variable.



Newton's
Laws

$$\begin{aligned} z_{t+1} &= z_t + v_t \\ v_{t+1} &= v_t + a_t \\ m a_t &= u_t \end{aligned}$$

minimize $\sum_{t=0}^T \mathbb{1}_{\{(x_t)_1\} > \epsilon}$

subject to $x_{t+1} = \begin{bmatrix} I & I \\ 0 & I \end{bmatrix} x_t + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} u_t$

$$x_t = \begin{bmatrix} z_t \\ v_t \end{bmatrix}$$



Newton's
Laws

$$\begin{aligned}z_{t+1} &= z_t + v_t \\v_{t+1} &= v_t + a_t \\m a_t &= u_t\end{aligned}$$

$$\text{minimize} \quad \sum_{t=0}^T (x_t)_1^2 + r u_t^2$$

$$\text{subject to} \quad x_{t+1} = \begin{bmatrix} I & | \\ 0 & | \end{bmatrix} x_t + \begin{bmatrix} 0 \\ |/m \end{bmatrix} u_t$$

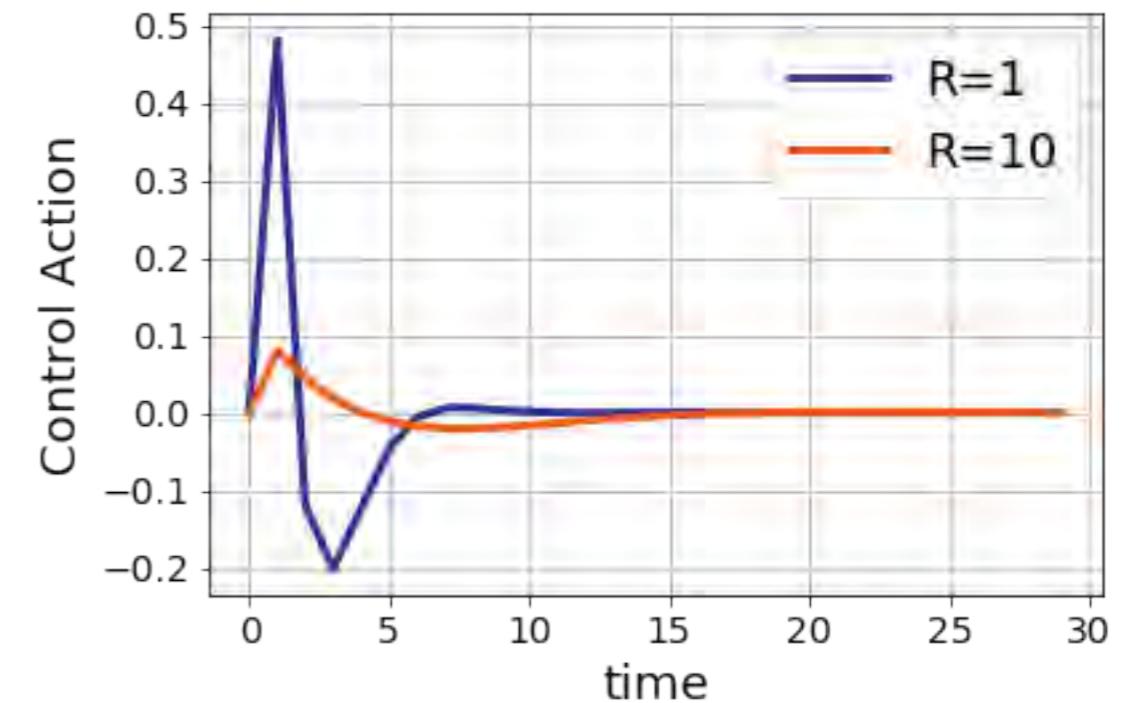
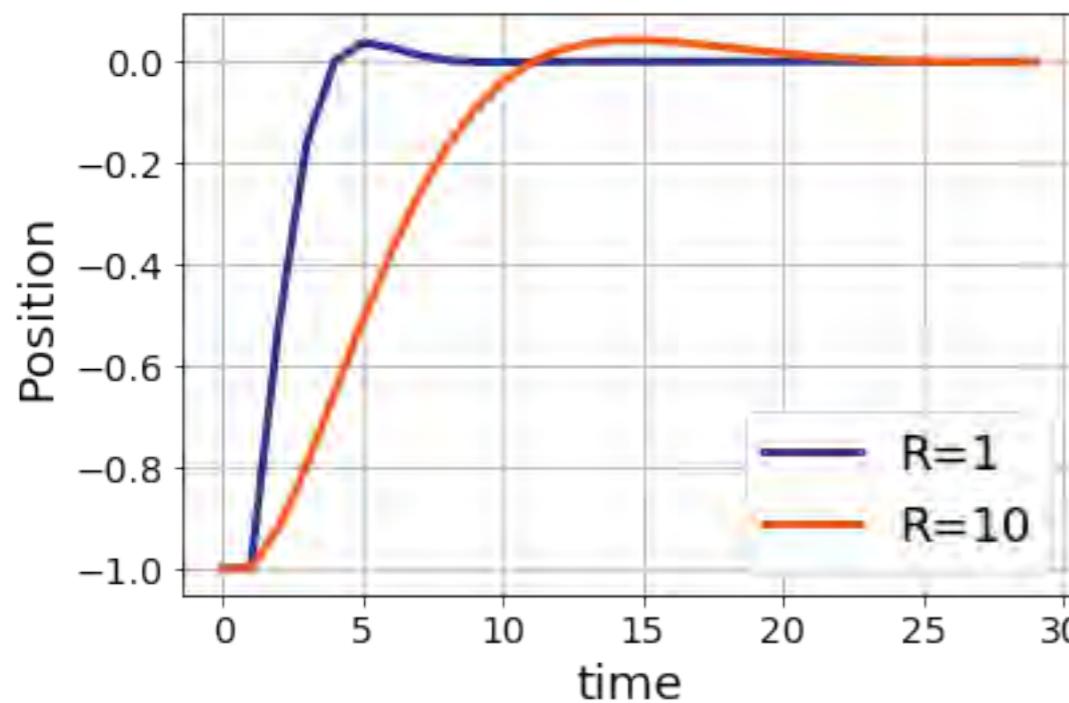
$$x_t = \begin{bmatrix} z_t \\ v_t \end{bmatrix}$$



minimize $\sum_{t=0}^T (x_t)_1^2 + ru_t^2$

subject to $x_{t+1} = \begin{bmatrix} I & I \\ 0 & I \end{bmatrix} x_t + \begin{bmatrix} 0 \\ I/m \end{bmatrix} u_t$

$$x_t = \begin{bmatrix} z_t \\ v_t \end{bmatrix}$$



“Simplest” Example: Linear Quadratic Regulator

$$\begin{aligned} \text{minimize} \quad & \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t \right] && \leftarrow \text{quadratic cost} \\ \text{s.t.} \quad & x_{t+1} = A x_t + B u_t + e_t && \leftarrow \text{linear dynamics} \end{aligned}$$



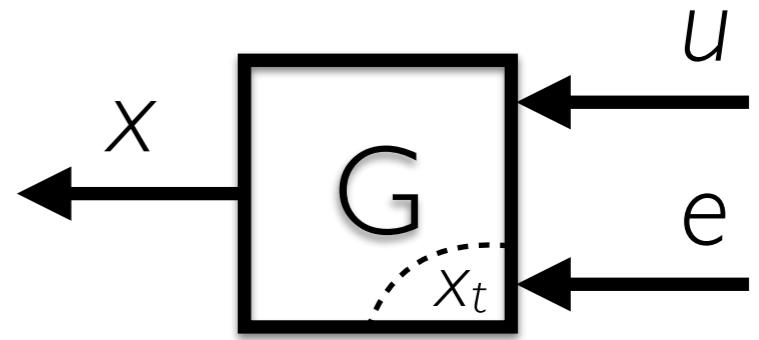
$$\text{minimize} \quad \sum_{t=0}^T (x_t)_1^2 + r u_t^2$$

$$\text{subject to} \quad x_{t+1} = \begin{bmatrix} I & I \\ 0 & I \end{bmatrix} x_t + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} u_t$$

$$x_t = \begin{bmatrix} z_t \\ v_t \end{bmatrix}$$

Optimal control

$$\begin{aligned} \text{minimize} \quad & \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ \text{s.t.} \quad & x_{t+1} = f_t(x_t, u_t, e_t) \\ & u_t = \pi_t(\tau_t) \end{aligned}$$



generic solutions with known dynamics:

Batch Optimization

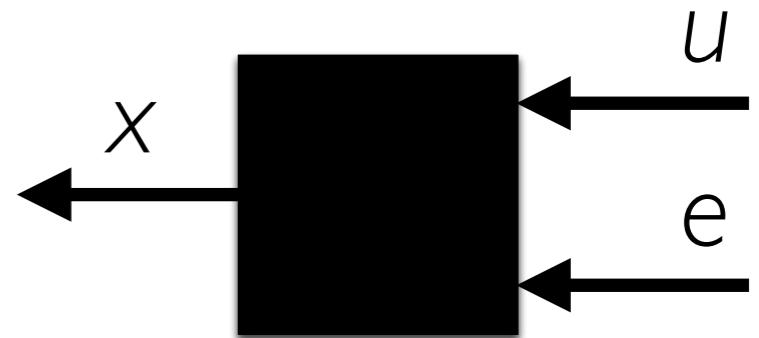
Dynamic Programming

Learning to control

$$\text{minimize} \quad \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right]$$

$$\text{s.t.} \quad x_{t+1} = f_t(x_t, u_t, e_t)$$

$$u_t = \pi_t(\tau_t)$$



C_t is the cost. If you maximize, it's called a reward.

e_t is a noise process

f_t is the state-transition function unknown!

$\tau_t = (u_1, \dots, u_{t-1}, x_0, \dots, x_t)$ is an observed trajectory

$\pi_t(\tau_t)$ is the policy. This is the optimization decision variable.

Major challenge: how to perform optimal control when the system is unknown?

Today: Reinvent RL attempting to answer this question

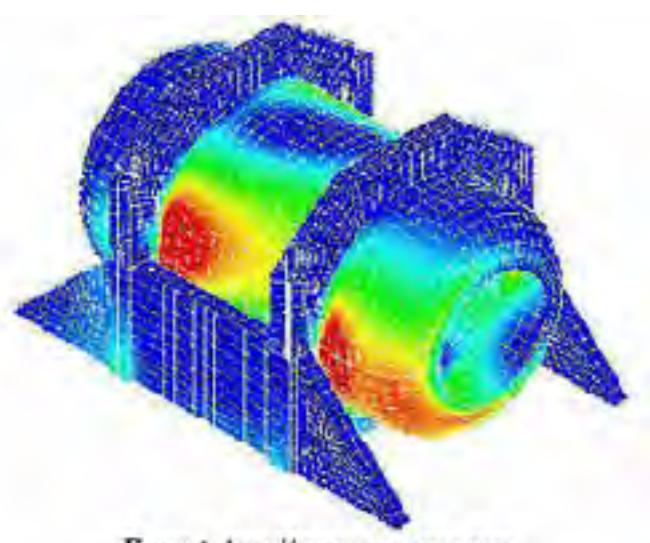
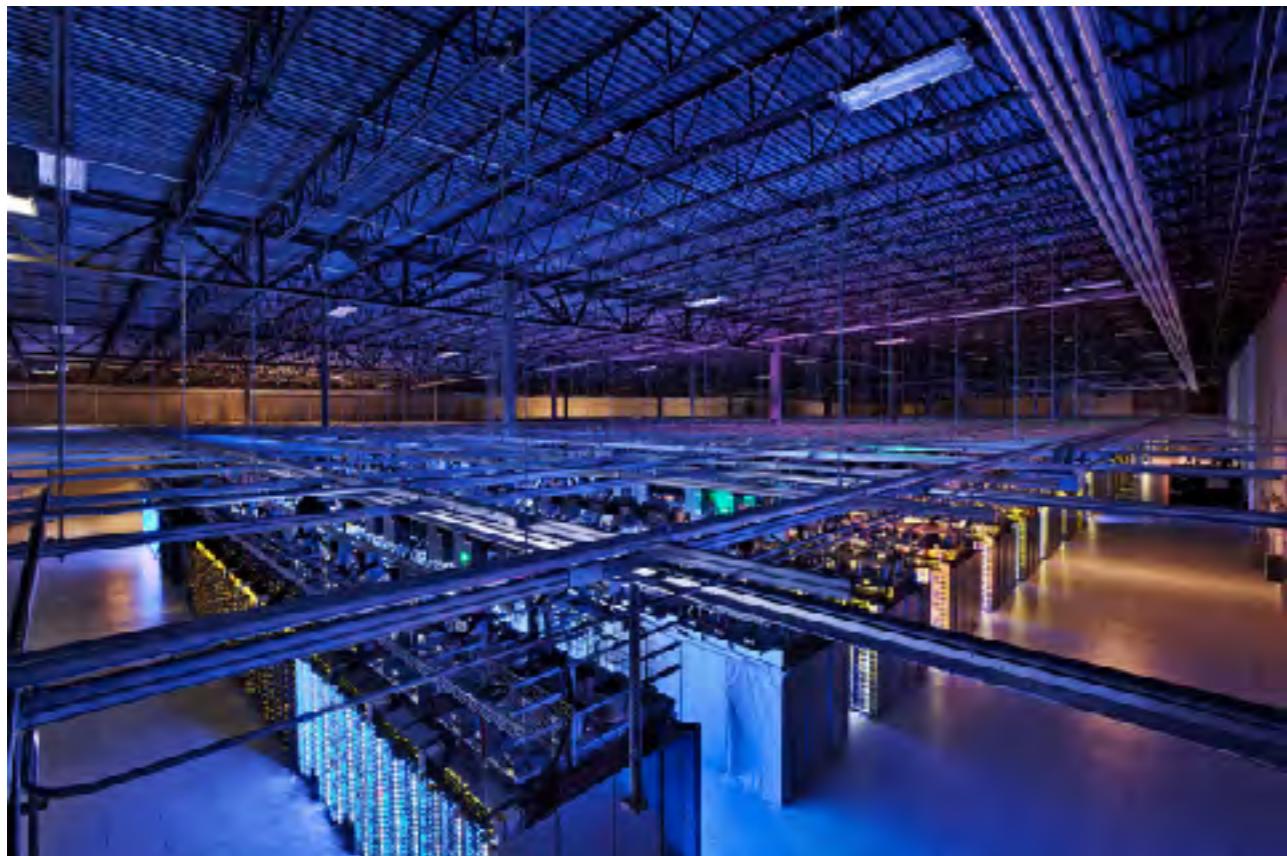
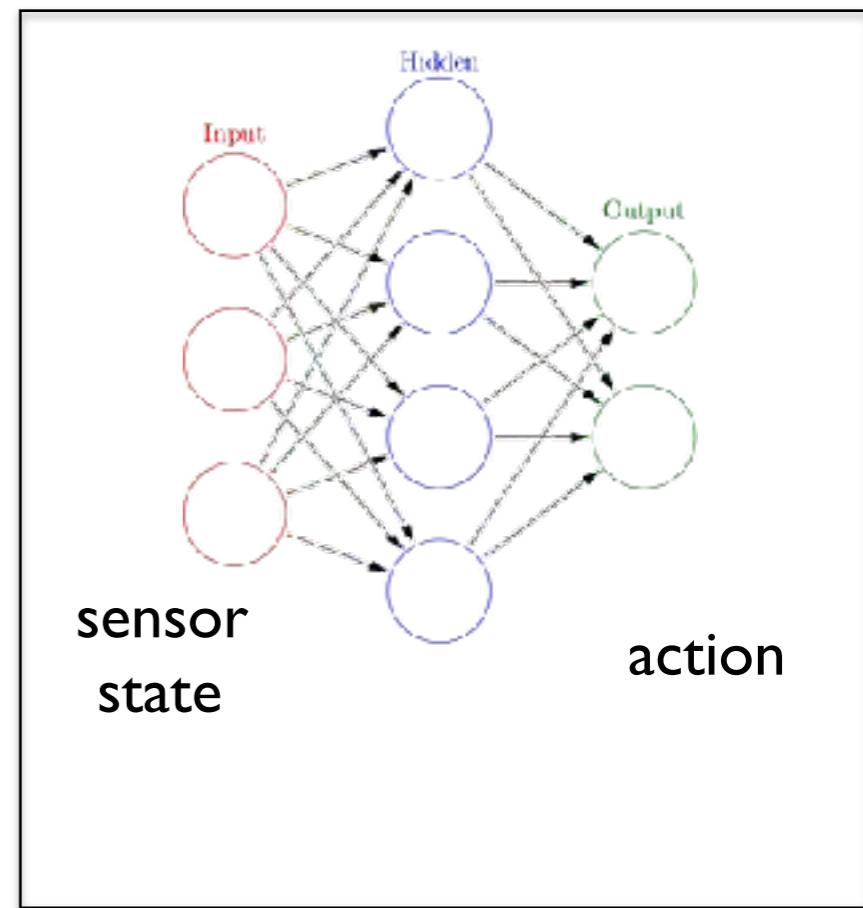
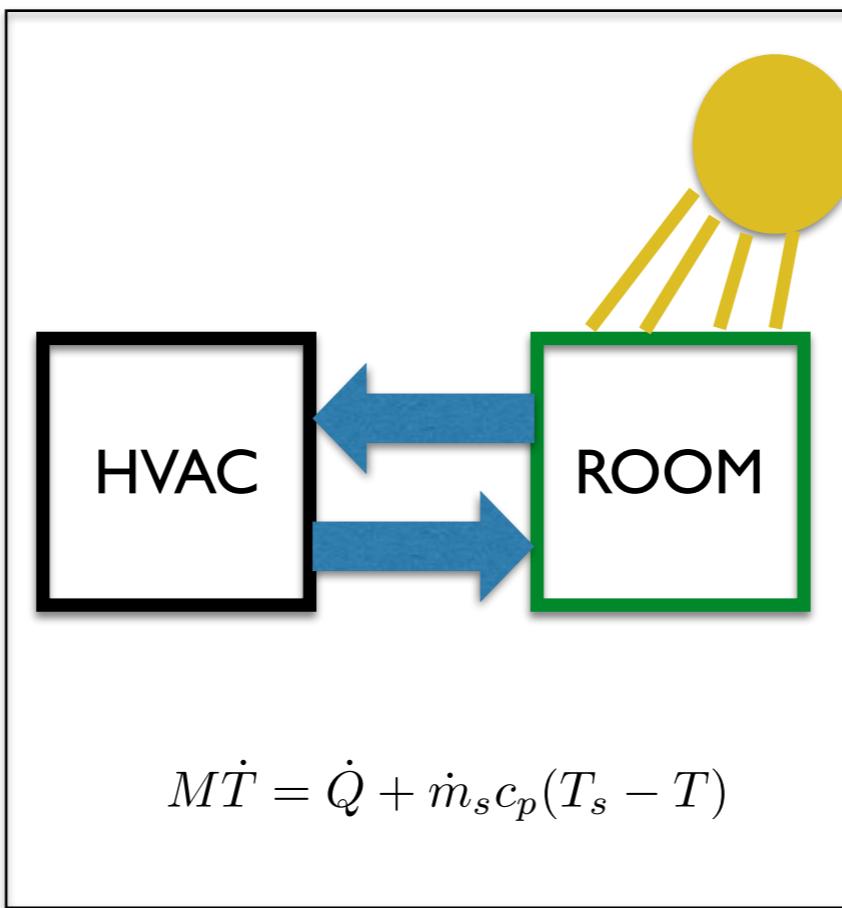
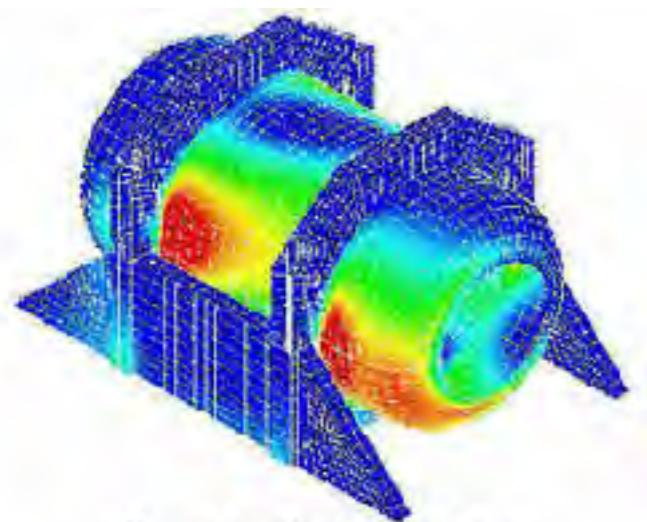


Figure 1 <http://www.noraneng.com>

$$\frac{\partial}{\partial t}(\rho \mathbf{u}) + \nabla \cdot (\rho \mathbf{u} \otimes \mathbf{u} + p \mathbf{I}) = \nabla \cdot \boldsymbol{\tau} + \rho \mathbf{g}$$

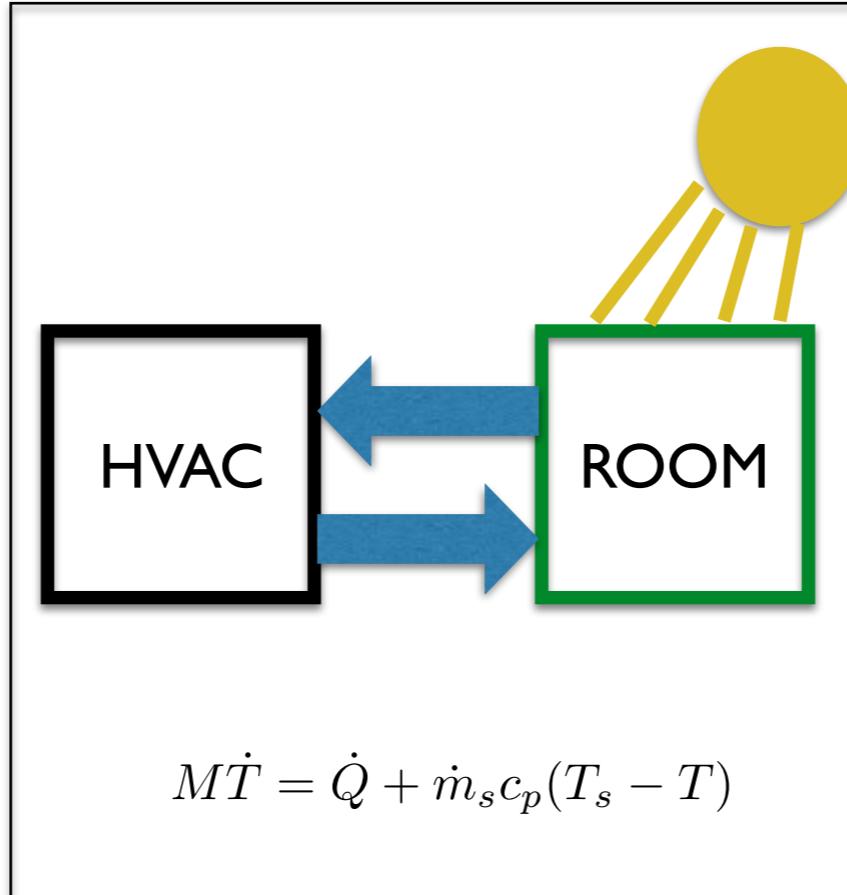


*Identify
everything*

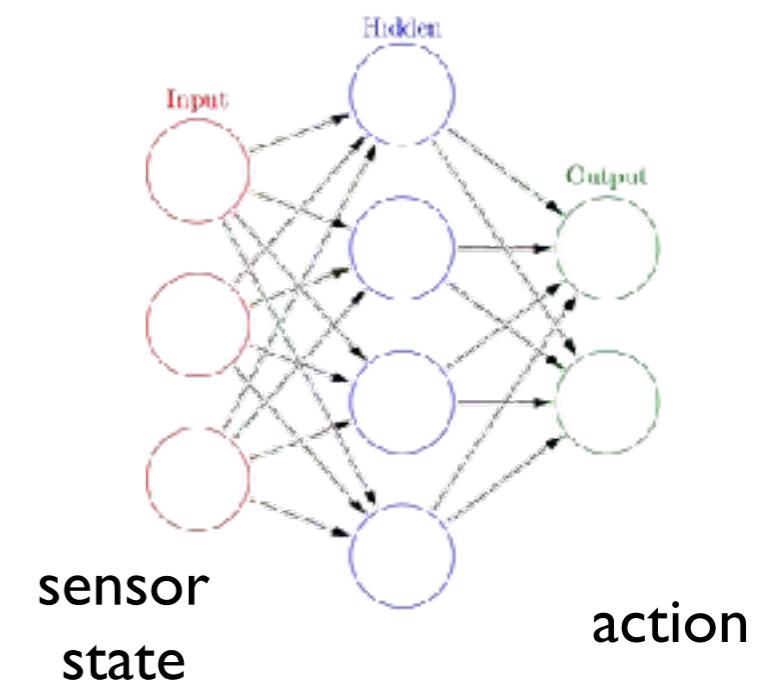


$$\frac{\partial}{\partial t}(\rho \mathbf{u}) + \nabla \cdot (\rho \mathbf{u} \otimes \mathbf{u} + p \mathbf{I}) = \nabla \cdot \boldsymbol{\tau} + \rho \mathbf{g}$$

*Identify a
coarse model*



*We don't need no
stinking models!*



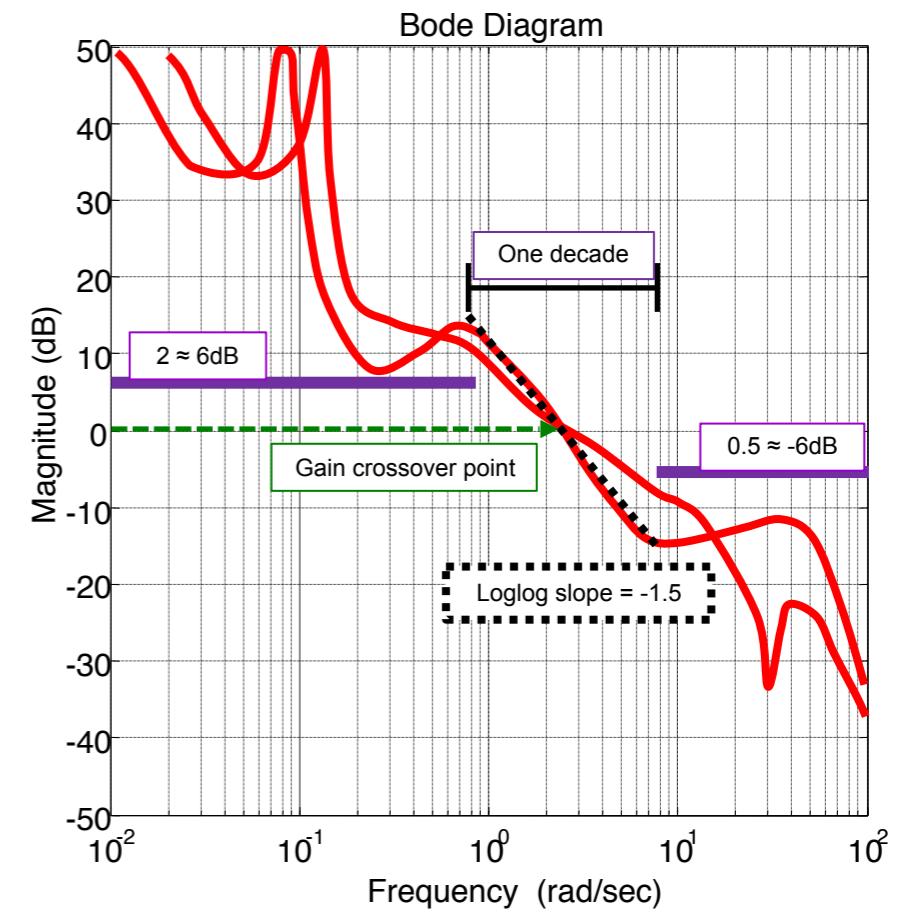
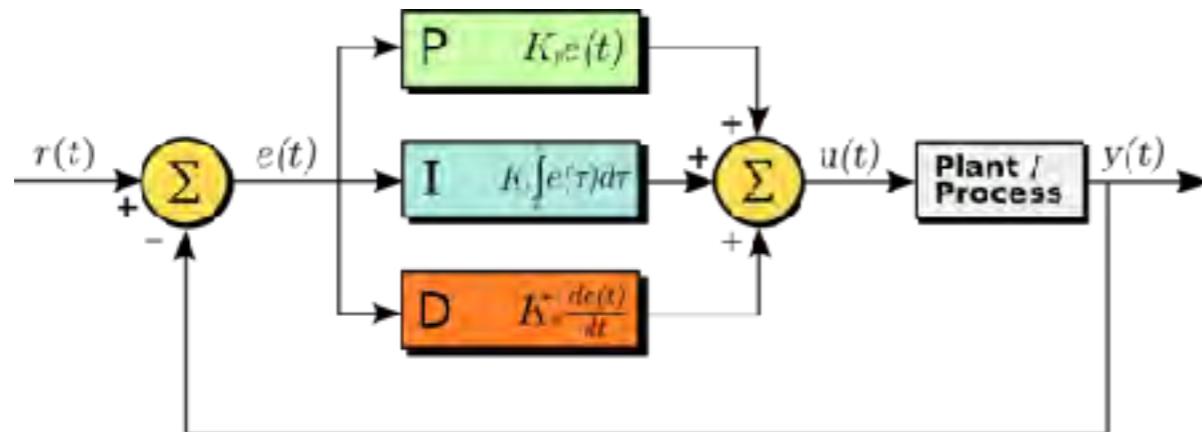
- PDE control
- High performance aerodynamics

- model predictive control

- reinforcement learning
- PID control?

We need robust fundamentals to
distinguish these approaches

But PID control works...



2 parameters suffice for 95% of all control applications.

How much needs to be modeled for more advanced control?

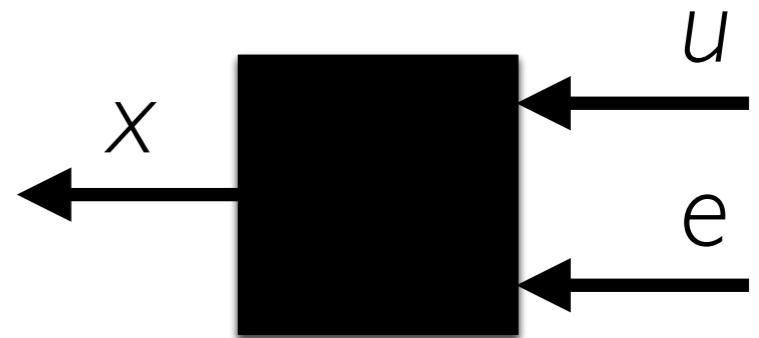
Can we learn to compensate for poor models, changing conditions?

Learning to control

$$\text{minimize} \quad \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right]$$

$$\text{s.t.} \quad x_{t+1} = f_t(x_t, u_t, e_t)$$

$$u_t = \pi_t(\tau_t)$$



C_t is the cost. If you maximize, it's called a reward.

e_t is a noise process

f_t is the state-transition function unknown!

$\tau_t = (u_1, \dots, u_{t-1}, x_0, \dots, x_t)$ is an observed trajectory

$\pi_t(\tau_t)$ is the policy. This is the optimization decision variable.

Major challenge: how to perform optimal control when the system is unknown?

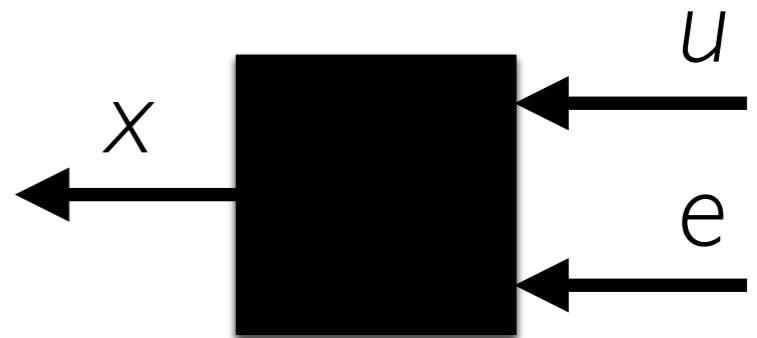
Today: Reinvent RL attempting to answer this question

Learning to control

$$\text{minimize} \quad \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right]$$

$$\text{s.t.} \quad x_{t+1} = f_t(x_t, u_t, e_t)$$

$$u_t = \pi_t(\tau_t)$$



Oracle: You can generate N trajectories of length T .

Challenge: Build a controller with smallest error
with fixed sampling budget ($N \times T$).

What is the optimal estimation/design scheme?

How many samples are needed for near optimal control?

The Linearization Principle

If a machine learning algorithm does crazy things when restricted to linear models, it's going to do crazy things on complex nonlinear models too.

Would you believe someone had a good SAT solver if it couldn't solve 2-SAT?

This has been a fruitful research direction:

- Recurrent neural networks (Hardt, Ma, R. 2016)
- Generalization and Margin in Neural Nets (Zhang et al 2017)
- Residual Networks (Hardt and Ma 2017)
- Bayesian Optimization (Jamieson et al 2017)
- Adaptive gradient methods (Wilson et al 2017)

“Simplest” Example: LQR

$$\begin{aligned} \text{minimize} \quad & \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t \right] \\ \text{s.t.} \quad & x_{t+1} = A x_t + B u_t + e_t \end{aligned}$$

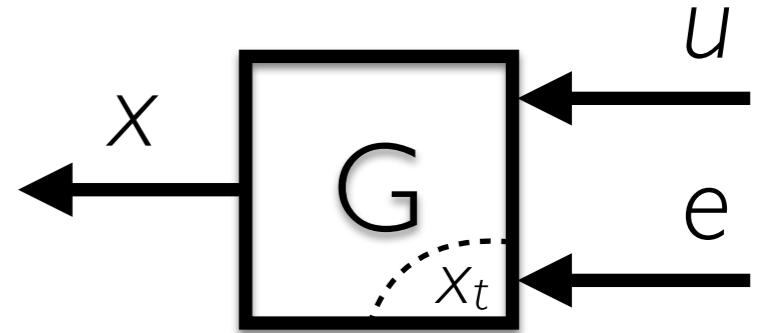


$$\text{minimize} \quad \sum_{t=0}^T (x_t)_1^2 + r u_t^2$$

$$\text{subject to} \quad x_{t+1} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} x_t + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} u_t$$

$$x_t = \begin{bmatrix} z_t \\ v_t \end{bmatrix}$$

RL Methods



$$\begin{aligned} & \text{minimize} && \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right] && \text{approximate dynamic programming} \\ & \text{s.t.} && x_{t+1} = f_t(x_t, u_t, e_t) && \text{model-based} \\ & && u_t = \pi_t(\tau_t) && \text{direct policy search} \end{aligned}$$

How to solve optimal control when the model f is unknown?

- Model-based: fit model from data
- Model-free
 - Approximate dynamic programming: estimate cost from data
 - Direct policy search: search for actions from data

$$\begin{aligned} & \text{minimize} && \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ & \text{s.t.} && x_{t+1} = f(x_t, u_t, e_t) \\ & && u_t = \pi_t(\tau_t) \end{aligned}$$

Model-based RL

Collect some simulation data. Should have

$$x_{t+1} \approx \varphi(x_t, u_t) + \nu_t$$

Fit dynamics with *supervised learning*:

$$\hat{\varphi} = \arg \min_{\varphi} \sum_{t=0}^{N-1} \|x_{t+1} - \varphi(x_t, u_t)\|^2$$

Solve approximate problem:

$$\begin{aligned} & \text{minimize} && \mathbb{E}_{\omega} \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ & \text{s.t.} && x_{t+1} = \varphi(x_t, u_t) + \omega_t \\ & && u_t = \pi(\tau_t) \end{aligned}$$

minimize

$$\mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right]$$

s.t.

$$\begin{aligned}x_{t+1} &= f_t(x_t, u_t, e_t) \\u_t &= \pi_t(\tau_t)\end{aligned}$$

~~Approximate~~ Dynamic Programming

Both the methods and analyses are complicated, but this is the core of classical RL.

Sadly, if you don't already know it, this probably won't make a ton of sense until the sixth time you see it...

Dynamic Programming

$$\text{minimize } \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) + C_f(x_{T+1}, u_{T+1}) \right] =: Q_1(x, u)$$

$$\begin{aligned} \text{s.t. } & x_{t+1} = f_t(x_t, u_t, e_t), \quad x_1 = x \\ & u_t = \pi_t(\tau_t), \quad u_1 = u \end{aligned}$$

“Q-function”

Terminal Q-function:

$$Q_{T+1}(x, u) = C_f(x, u)$$

Recursive formula (recurse backwards):

$$Q_k(x, u) = C_k(x, u) + \min_{u'} \mathbb{E}_e [Q_{k+1}(f_k(x, u, e), u')]$$

Optimal Policy:

$$\pi_k(\tau_k) = \arg \min_u Q_k(x_k, u)$$

“Simplest” Example: LQR

$$\begin{aligned} \text{minimize} \quad & \mathbb{E} \left[\sum_{t=1}^{T-1} x_t^* Q x_t + u_t^* R u_t + x_T^* P_T x_T \right] \\ \text{s.t.} \quad & x_{t+1} = Ax_t + Bu_t + e_t \end{aligned}$$

Dynamic Programming: $Q_T(x, u) = x^* P_T x$

$$\begin{aligned} Q_t(x, u) &= C_t(x, u) + \min_{u'} \mathbb{E}_e [Q_{t+1}(f_t(x, u, e), u')] \\ &= x^* Q x + u^* R u + (Ax + Bu)^* P_{t+1} (Ax + Bu) + c_t \end{aligned}$$

$$P_t = Q + A^* P_{t+1} A - A^* P_{t+1} B (R + B^* P_{t+1} B)^{-1} B^* P_{t+1} A$$

$$u_t = -(B^* P_{t+1} B + R)^{-1} B^* P_{t+1} A x_t =: K_t x_t$$

“Simplest” Example: LQR

$$\begin{aligned} \text{minimize} \quad & \lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t \right] \\ \text{s.t.} \quad & x_{t+1} = Ax_t + Bu_t + e_t \end{aligned}$$

When (A, B) known, optimal to build control $u_t = Kx_t$

$$u_t = -(B^*PB + R)^{-1}B^*PAx_t =: Kx_t$$

$$P = Q + A^*PA - A^*PB(R + B^*PB)^{-1}B^*PA$$

Discrete Algebraic Riccati Equation

- Dynamic programming has simple form because quadratics are miraculous.
- Solution is independent of noise variance.
- For finite time horizons, we could solve this with a variety of batch solvers.
- Note that the solution is only time invariant on the infinite time horizon.

minimize

$$\mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right]$$

s.t.

$$x_{t+1} = f_t(x_t, u_t, e_t)$$

$$u_t = \pi_t(\tau_t)$$

Approximate Dynamic Programming

Recursive formula:

$$Q_k(x, u) = C_k(x, u) + \mathbb{E}_e \left[\min_{u'} Q_{k+1}(f_k(x, u, e), u') \right]$$

Optimal Policy:

$$\pi_k(\tau_k) = \arg \min_u Q_k(x_k, u)$$

minimize
s.t.

$$\mathbb{E}_e \left[\sum_{t=1}^{\infty} \gamma^t C(x_t, u_t) \right]$$

$$x_{t+1} = f(x_t, u_t, e_t)$$

$$u_t = \pi(\tau_t)$$

discount factor

Approximate Dynamic Programming

Bellman Equation:

$$Q(x, u) = C(x, u) + \gamma \mathbb{E}_e \left[\min_{u'} Q(f(x, u, e), u') \right]$$

Optimal Policy:

$$\pi(x) = \arg \min_u Q(x, u)$$

Generate algorithms using the insight:

$$Q(x_k, u_k) \approx C(x_k, u_k) + \gamma \min_{u'} Q(x_{k+1}, u') + \nu_k$$

Q-learning:

stochastic approximation

$$Q_{\text{new}}(x_k, u_k) = (1 - \eta) Q_{\text{old}}(x_k, u_k) - \eta \left(C(x_k, u_k) + \gamma \min_{u'} Q_{\text{old}}(x_{k+1}, u') \right)$$

$$\begin{aligned} \text{minimize} \quad & \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ \text{s.t.} \quad & x_{t+1} = f_t(x_t, u_t, e_t) \\ & u_t = \pi_t(\tau_t) \end{aligned}$$

Direct Policy Search

Sampling to Search

$$\begin{aligned} \min_{z \in \mathbb{R}^d} \Phi(z) &= \min_{p(z)} \mathbb{E}_p[\Phi(z)] \\ &\leq \min_{\vartheta} \mathbb{E}_{p(z; \vartheta)}[\Phi(z)] =: J(\vartheta) \end{aligned}$$

- Search over probability distributions
- Use function approximations that might not capture optimal distribution
- Can build (incredibly high variance) stochastic gradient estimates by sampling:

$$\nabla J(\vartheta) = \mathbb{E}_{p(z; \vartheta)} [\Phi(z) \nabla_{\vartheta} \log p(z; \vartheta)]$$

Reinforce Algorithm

$$J(\vartheta) := \mathbb{E}_{p(z; \vartheta)} [\Phi(z)]$$

$$\begin{aligned}\nabla_{\vartheta} J(\vartheta) &= \int \Phi(z) \nabla_{\vartheta} p(z; \vartheta) dz \\ &= \int \Phi(z) \left(\frac{\nabla_{\vartheta} p(z; \vartheta)}{p(z; \vartheta)} \right) p(z; \vartheta) dz \\ &= \int (\Phi(z) \nabla_{\vartheta} \log p(z; \vartheta)) p(z; \vartheta) dz \\ &= \mathbb{E}_{p(z; \vartheta)} [\Phi(z) \nabla_{\vartheta} \log p(z; \vartheta)]\end{aligned}$$

Reinforce Algorithm

$$J(\vartheta) := \mathbb{E}_{p(z;\vartheta)} [\Phi(z)]$$

$$\nabla J(\vartheta) = \mathbb{E}_{p(z;\vartheta)} [\Phi(z) \nabla_{\vartheta} \log p(z; \vartheta)]$$

Sample

$$z_k \sim p(z; \vartheta_k)$$

Compute

$$G(z_k, \vartheta_k) = \Phi(z_k) \nabla_{\vartheta_k} \log p(z_k; \vartheta_k)$$

Update

$$\vartheta_{k+1} = \vartheta_k - \alpha_k G(z_k, \vartheta_k)$$

Reinforce Algorithm

$$J(\vartheta) := \mathbb{E}_{p(z; \vartheta)} [\Phi(z)]$$

Sample

$$z_k \sim p(z; \vartheta_k)$$

Compute

$$G(z_k, \vartheta_k) = \Phi(z_k) \nabla_{\vartheta_k} \log p(z_k; \vartheta_k)$$

Update

$$\vartheta_{k+1} = \vartheta_k - \alpha_k G(z_k, \vartheta_k)$$

Generic algorithm for solving discrete optimization:

$$z \in \{-1, 1\}^d$$

$$p(z; \vartheta) = \prod_{i=1}^d \frac{\exp(z_i \vartheta_i)}{\exp(-\vartheta_i) + \exp(\vartheta_i)}$$

$$\vartheta_{k+1} = \vartheta_k - \alpha_k \Phi(z_k) (z_k + \tanh(\vartheta_k))$$

Does this “solve” *any* discrete problem?

Random Search

$$\begin{aligned} \text{minimize } & \mathbb{E}_{e_t, \omega} \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ \text{s.t. } & x_{t+1} = f_t(x_t, u_t, e_t) \\ & u_t = \pi(\tau_t; \vartheta + \omega) \end{aligned}$$

Direct Policy Search

$$G(\omega, \vartheta) = \left(\sum_{t=1}^T C(x_t, u_t) \right) \nabla \log p(\omega)$$

parameter perturbation

$$G^{(m)}(\omega, \vartheta) = \frac{1}{m} \sum_{i=1}^m \frac{C(\vartheta + \sigma \omega_i) - C(\vartheta - \sigma \omega_i)}{2\sigma} \omega_i$$

$$\begin{aligned} C(\vartheta) &= \sum_{t=1}^T C(x_t, u_t) \\ \omega_i &\sim \mathcal{N}(0, I) \end{aligned}$$

random finite difference approximation to the gradient

aka... (μ, λ) -Evolution Strategies
 SPSA
 Bandit Convex Opt

Random Search for LQR

$$\begin{aligned} \text{minimize} \quad & \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t \right] \\ \text{s.t.} \quad & x_{t+1} = Ax_t + Bu_t + e_t \end{aligned}$$

“Greedy strategy”: Build control $u_t = Kx_t$

- Sample a random perturbation: $\nu \sim \mathcal{N}(0, \sigma^2 I)$
- Collect samples from control $u_t = (K + \nu)x_t : \tau = \{x_1, \dots, x_T\}$
- Compute cost: $J(\tau) = \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t$
- Update: $K \leftarrow K - \alpha_t J(\tau) \nu$

Policy Gradient

$$\begin{aligned} \text{minimize} \quad & \mathbb{E}_{e_t, u_t} \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ \text{s.t.} \quad & x_{t+1} = f_t(x_t, u_t, e_t) \\ & u_t \sim p(u|x_t; \vartheta) \end{aligned}$$

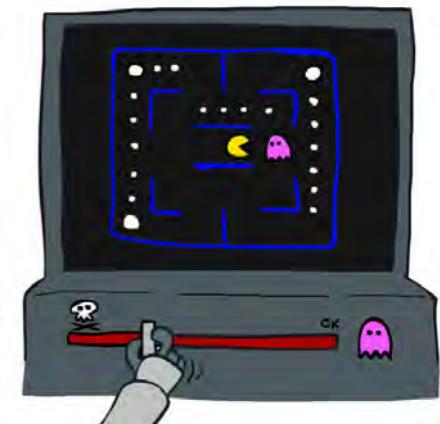
probabilistic policy

Direct Policy Search

$$G(\tau, \vartheta) = \left(\sum_{t=1}^T C(x_t, u_t) \right) \cdot \left(\sum_{t=0}^{T-1} \nabla_{\vartheta} \log p_{\vartheta}(u_t | x_t; \vartheta) \right)$$

Policy Gradient for LQR

$$\begin{aligned} & \text{minimize} && \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t \right] \\ & \text{s.t.} && x_{t+1} = Ax_t + Bu_t + e_t \end{aligned}$$



“Greedy strategy”: Build control $u_t = Kx_t$

- Sample a bunch of random vectors: $\nu_t \sim \mathcal{N}(0, \sigma^2 I)$
- Collect samples from control $u_t = Kx_t + \nu_t$: $\tau = \{x_1, \dots, x_T\}$
- Compute cost: $C(\tau) = \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t$
- Update: $K_{\text{new}} \leftarrow K_{\text{old}} - \alpha_t C(\tau) \sum_{t=0}^{T-1} \nu_t x_t^*$

policy gradient
only has access
to 0-th order
information!!!

$$\begin{aligned}
 & \text{minimize} && \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\
 & \text{s.t.} && x_{t+1} = f_t(x_t, u_t, e_t) \\
 & && u_t = \pi_t(\tau_t)
 \end{aligned}$$

Direct Policy Search

Policy Gradient

$$\begin{aligned}
 & \text{minimize} && \mathbb{E}_{e_t, u_t} \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\
 & \text{s.t.} && x_{t+1} = f_t(x_t, u_t, e_t) \\
 & && u_t \sim p(u|x_t; \vartheta)
 \end{aligned}$$

probabilistic policy

Random Search

$$\begin{aligned}
 & \text{minimize} && \mathbb{E}_{e_t, \omega} \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\
 & \text{s.t.} && x_{t+1} = f_t(x_t, u_t, e_t) \\
 & && u_t = \pi(\tau_t; \vartheta + \omega)
 \end{aligned}$$

parameter perturbation

- Reinforce applied to either problems does not depend on the dynamics.
- Both are Derivative-free algorithms!

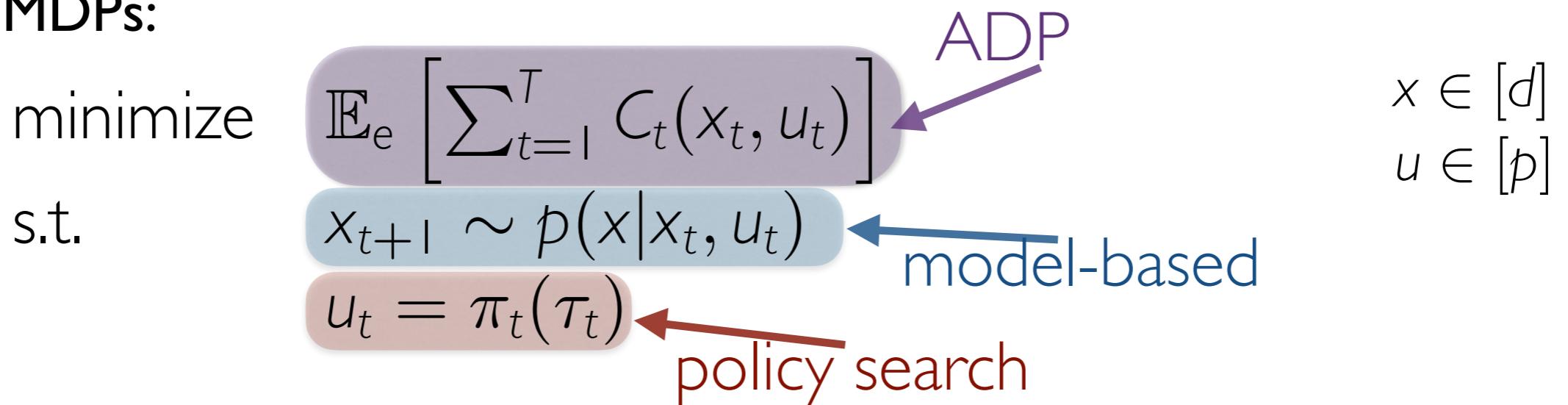
$$\begin{aligned} \text{minimize} \quad & \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ \text{s.t.} \quad & x_{t+1} = f_t(x_t, u_t, e_t) \\ & u_t = \pi_t(\tau_t) \end{aligned}$$

Direct Policy Search

- Reinforce is NOT Magic
 - What is the variance?
 - What is the approximation error?
 - Necessarily becomes derivative free as you are accessing the decision variable by sampling
- But it's certainly super easy!

Sample Complexity?

Discrete MDPs:



Algorithm Class	Samples per iteration	Parameters	“optimal” error after T
Model-based	1	d^2p	$\sqrt{\frac{d^2p}{T}}$
ADP	1	dp	$\sqrt{\frac{dp}{T}}$
Policy search	1	dp	$\sqrt{\frac{dp}{T}}$

Sample Complexity?

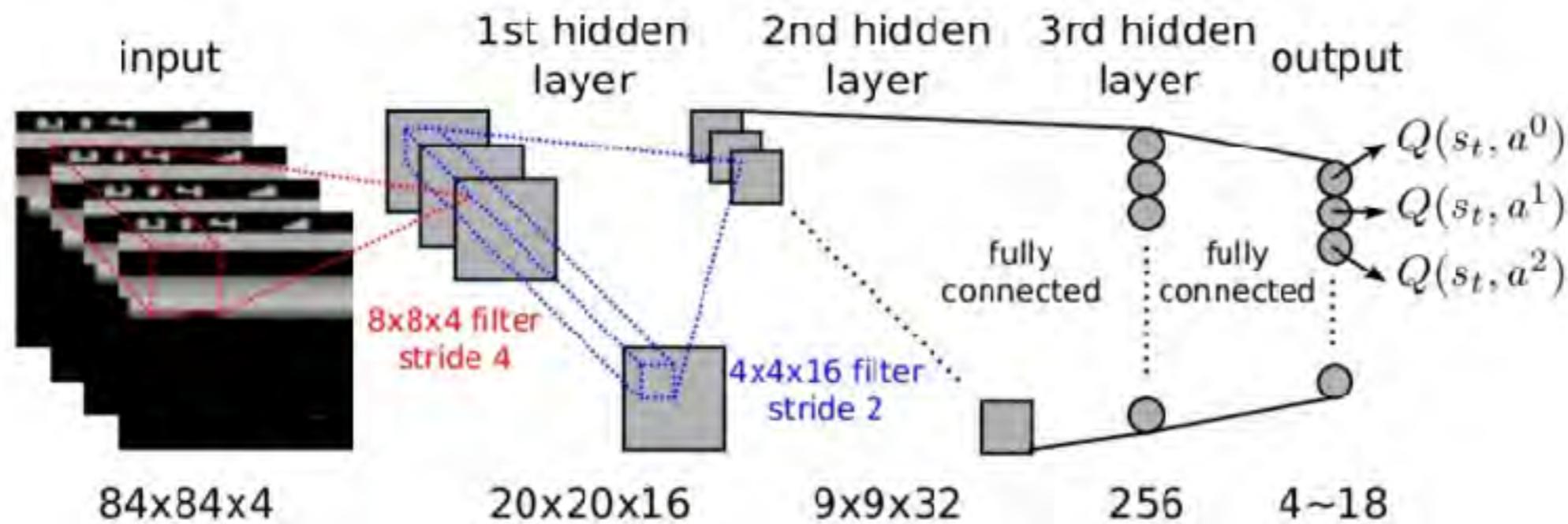
Continuous Control:

$$\begin{aligned}
 & \text{minimize} && \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\
 & \text{s.t.} && x_{t+1} = f_t(x_t, u_t, e_t) \\
 & && u_t = \pi_t(\tau_t)
 \end{aligned}$$

$x \in \mathbb{R}^d$
 $u \in \mathbb{R}^p$

Algorithm	Samples per iteration	LQR parameters	“optimal” error after T
Model-based	d	$d^2 + dp$	$C \sqrt{\frac{d+p}{T}}$
ADP	1	$\binom{d+p}{2}$	$C \frac{d+p}{\sqrt{T}}$
Policy search	1	dp	$C \sqrt{\frac{dp}{T}}$

Deep Reinforcement Learning



- Simply parameterize Q-function or policy as a deep net
- Note, ADP is tricky to analyze with function approximation
- Policy search is considerably more straightforward: make the log-prob a deep net.

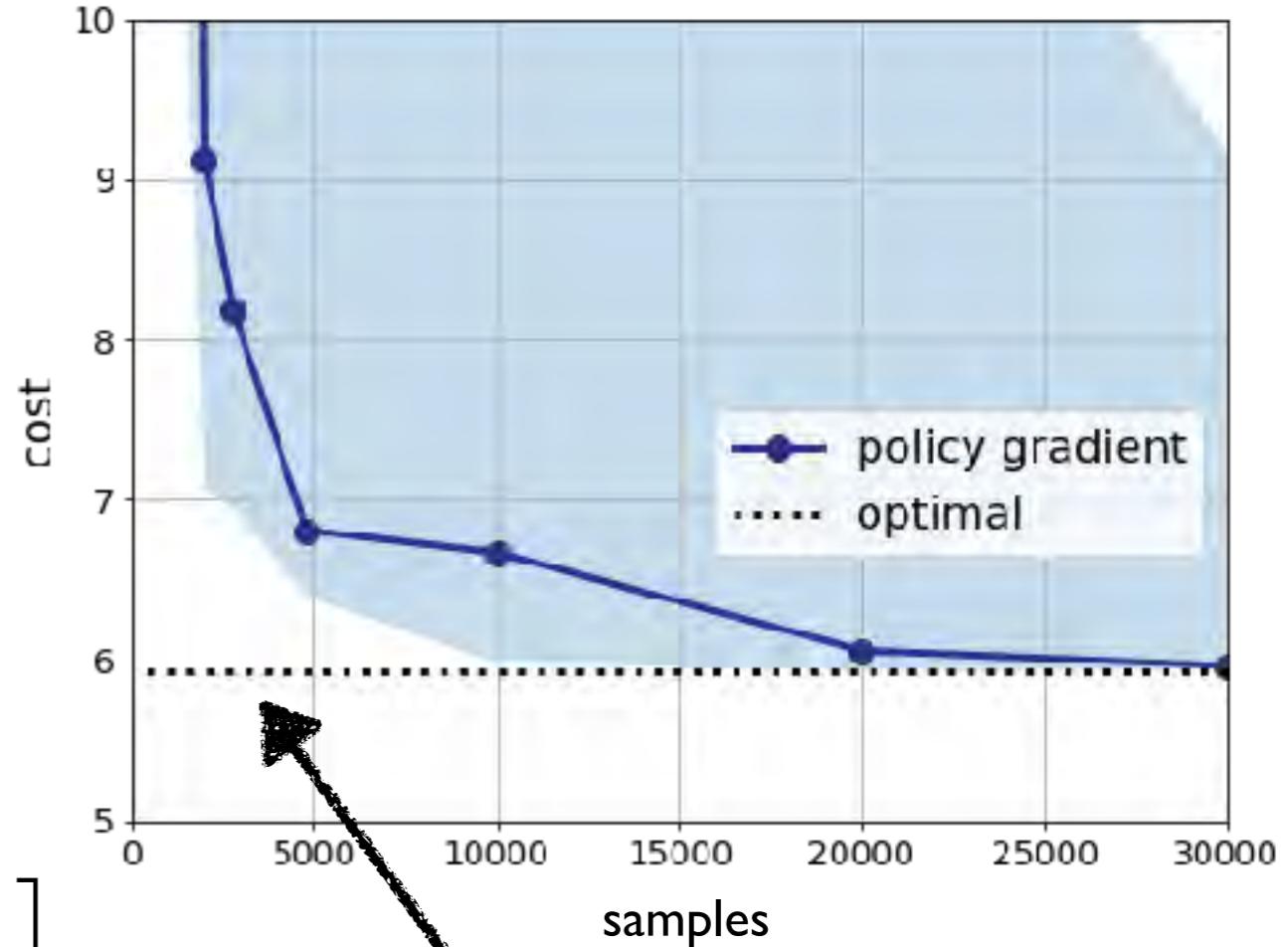
“Simplest” Example: LQR



minimize $\sum_{t=0}^T (x_t)_1^2 + ru_t^2$

subject to $x_{t+1} = \begin{bmatrix} I & I \\ 0 & I \end{bmatrix} x_t + \begin{bmatrix} 0 \\ I/m \end{bmatrix} u_t$

$$x_t = \begin{bmatrix} z_t \\ v_t \end{bmatrix}$$



Model-based and ADP
with 10 samples

Extraordinary Claims Require Extraordinary Evidence*



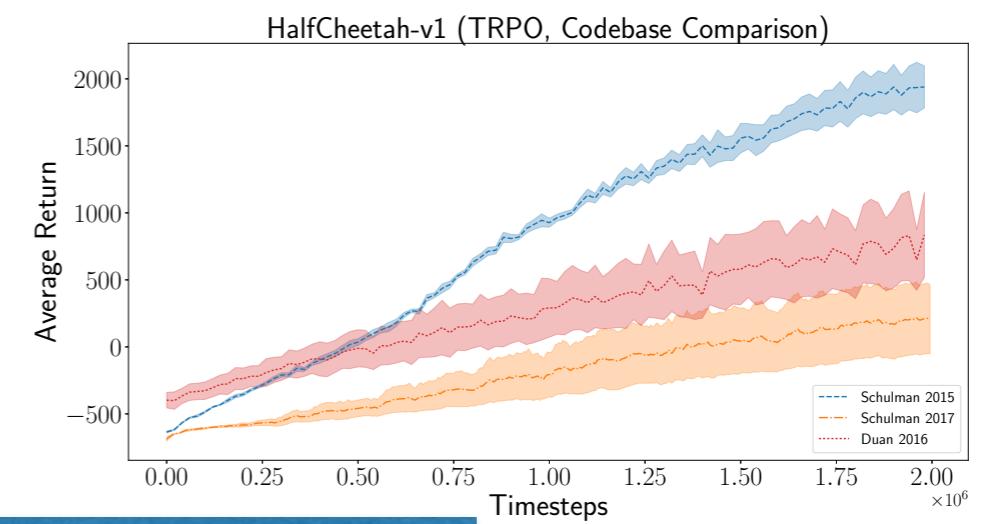
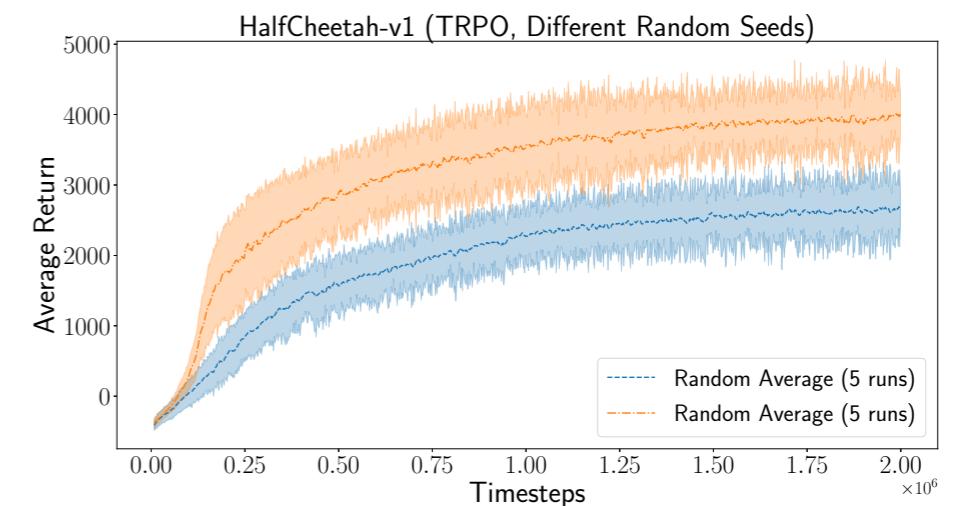
* only if your prior is correct

blog.openai.com/openai-baselines-dqn/

“Reinforcement learning results are tricky to reproduce: performance is very noisy, algorithms have many moving parts which allow for subtle bugs, and many papers don’t report all the required tricks.”

“RL algorithms are challenging to implement correctly; good results typically only come after fixing many seemingly-trivial bugs.”

[arxiv:1709.06560](https://arxiv.org/abs/1709.06560)



There has to be a better way!

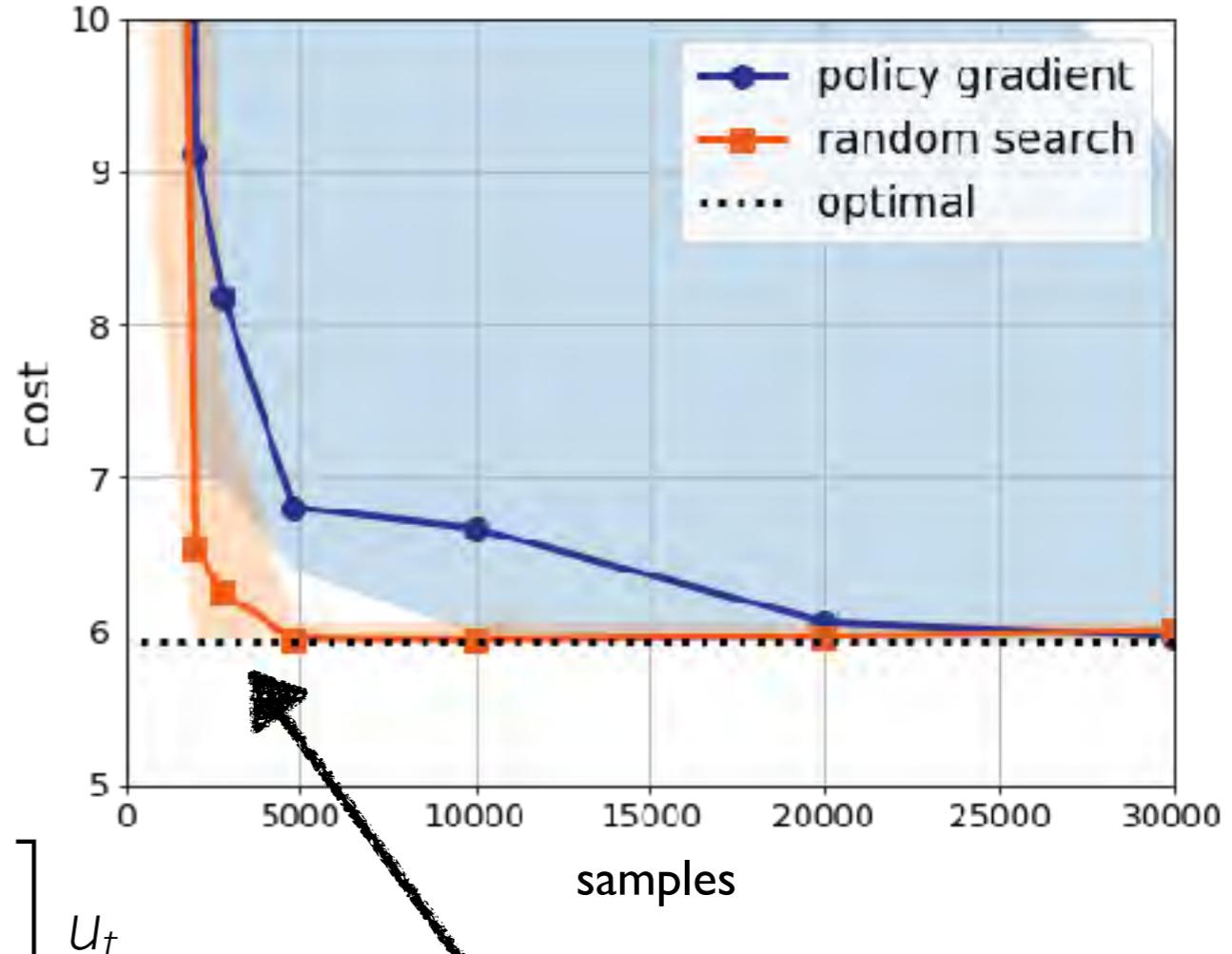
“Simplest” Example: LQR



minimize $\sum_{t=0}^T (x_t)_1^2 + ru_t^2$

subject to $x_{t+1} = \begin{bmatrix} I & I \\ 0 & I \end{bmatrix} x_t + \begin{bmatrix} 0 \\ I/m \end{bmatrix} u_t$

$$x_t = \begin{bmatrix} z_t \\ v_t \end{bmatrix}$$



Model-based and ADP
with 10 samples

Extraordinary Claims Require Extraordinary Evidence*



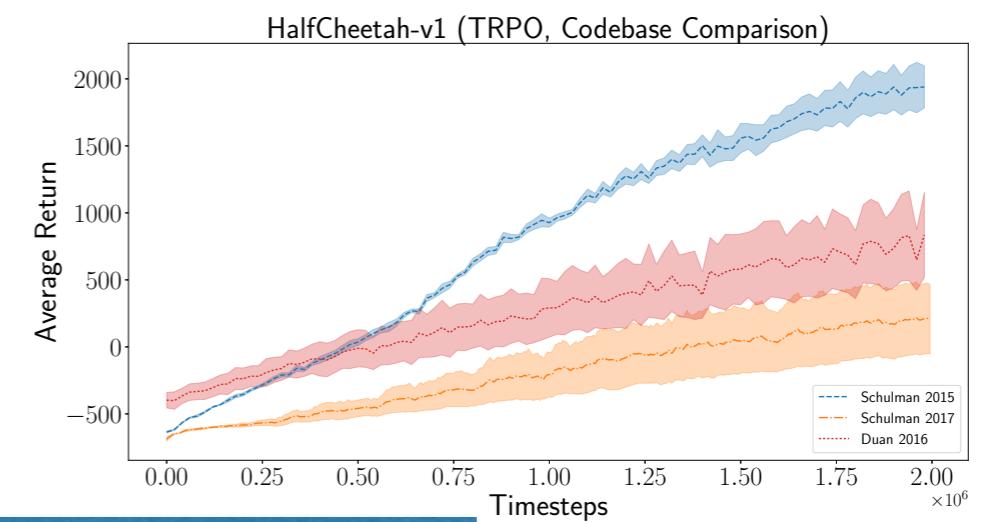
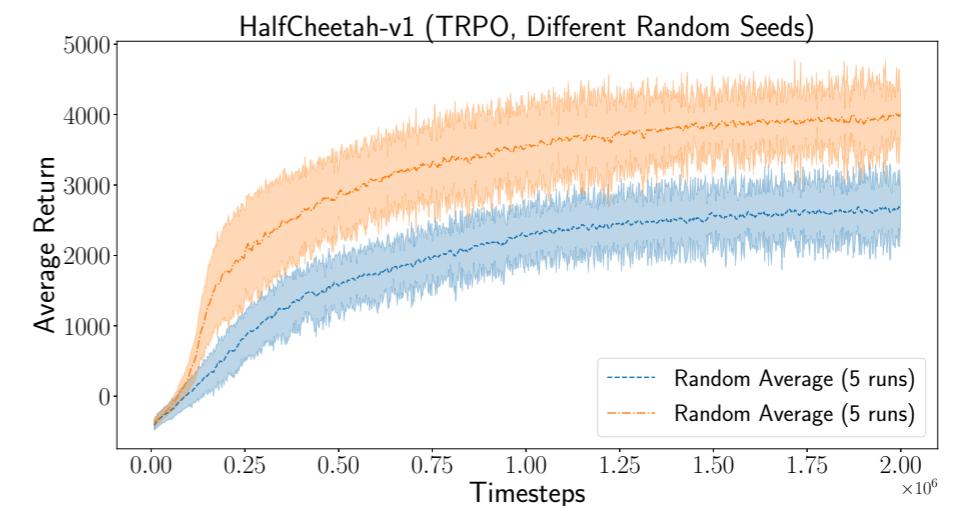
* only if your prior is correct

blog.openai.com/openai-baselines-dqn/

“Reinforcement learning results are tricky to reproduce: performance is very noisy, algorithms have many moving parts which allow for subtle bugs, and many papers don’t report all the required tricks.”

“RL algorithms are challenging to implement correctly; good results typically only come after fixing many seemingly-trivial bugs.”

[arxiv:1709.06560](https://arxiv.org/abs/1709.06560)



There has to be a better way!

$$\begin{aligned} & \text{minimize} && \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ & \text{s.t.} && x_{t+1} = f(x_t, u_t, e_t) \\ & && u_t = \pi_t(\tau_t) \end{aligned}$$

Model-based RL

Collect some simulation data. Should have

$$x_{t+1} \approx \varphi(x_t, u_t) + \nu_t$$

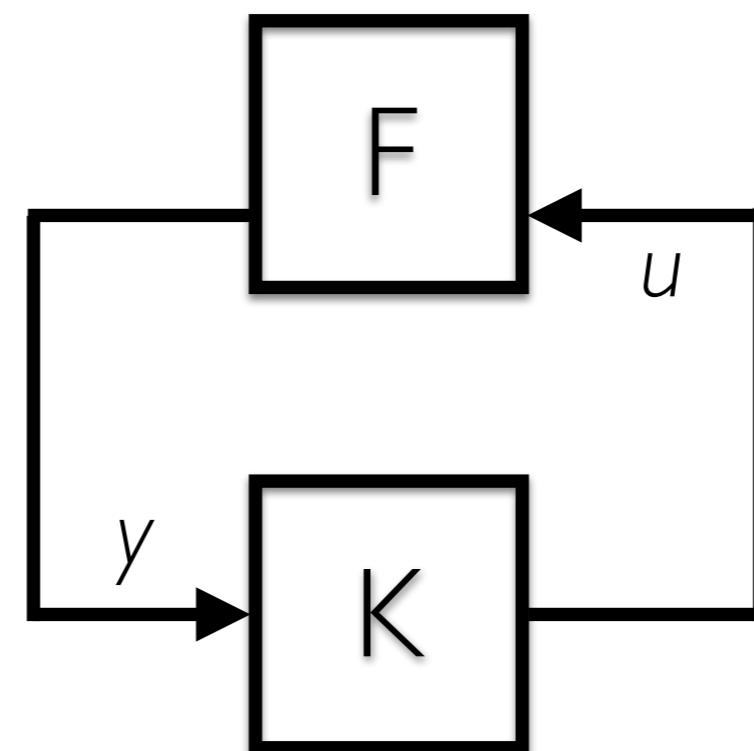
Fit dynamics with *supervised learning*:

$$\hat{\varphi} = \arg \min_{\varphi} \sum_{t=0}^{N-1} \|x_{t+1} - \varphi(x_t, u_t)\|^2$$

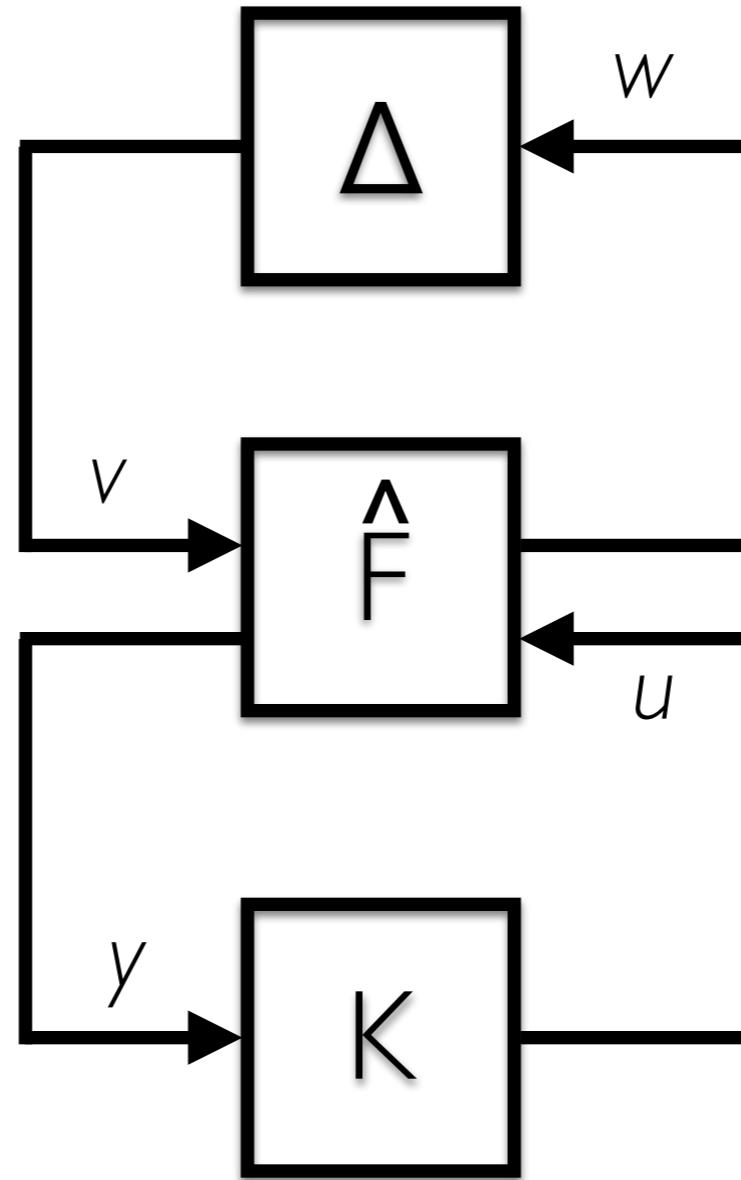
Solve approximate problem:

$$\begin{aligned} & \text{minimize} && \mathbb{E}_{\omega} \left[\sum_{t=1}^T C_t(x_t, u_t) \right] \\ & \text{s.t.} && x_{t+1} = \varphi(x_t, u_t) + \omega_t \\ & && u_t = \pi(\tau_t) \end{aligned}$$

Coarse-ID control



Coarse-ID control



High dimensional
stats bounds the
error

Coarse-grained
model is trivial
to fit

Design robust
control for
feedback loop

Coarse-ID control (static case)

$$\begin{array}{ll}\text{minimize}_u & x^* Q x \\ \text{subject to} & x = Bu + x_0\end{array}$$

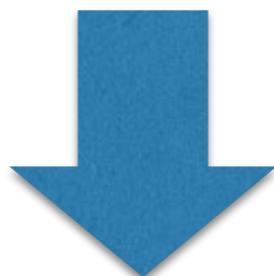
B unknown!

Collect data: $\{(x_i, u_i)\}$ $x_i = Bu_i + x_0 + e_i$

Estimate B : $\underset{B}{\text{minimize}} \quad \sum_{i=1}^N \|Bu_i + x_0 - x_i\|^2 \quad \rightarrow \hat{B}$

Guarantee: $\|B - \hat{B}\| \leq \epsilon$ with high probability

Note: $x = \hat{B}u + x_0 + \Delta_B u$



Robust optimization
problem:

$$\begin{array}{ll}\text{minimize}_u & \sup_{\|\Delta_B\| \leq \epsilon} \|Q^{1/2}(x - \Delta_B u)\| \\ \text{subject to} & x = \hat{B}u + x_0\end{array}$$

Coarse-ID control (static case)

$$\begin{array}{ll} \text{minimize}_u & x^* Q x \\ \text{subject to} & x = Bu + x_0 \end{array}$$

B unknown!

Collect data: $\{(x_i, u_i)\}$ $x_i = Bu_i + x_0 + e_i$

Estimate B : $\underset{B}{\text{minimize}} \quad \sum_{i=1}^N \|Bu_i + x_0 - x_i\|^2 \quad \rightarrow \hat{B}$

Guarantee: $\|B - \hat{B}\| \leq \epsilon$ with high probability

Solve robust
optimization problem:

$$\begin{array}{ll} \text{minimize}_u & \sup_{\|\Delta_B\| \leq \epsilon} \|Q^{1/2}(x - \Delta_B u)\| \\ \text{subject to} & x = \hat{B}u + x_0 \end{array}$$

Relaxation:
(Triangle inequality!)

$$\begin{array}{ll} \text{minimize}_u & \|Q^{1/2}x\| + \epsilon\lambda\|u\| \\ \text{subject to} & x = \hat{B}u + x_0 \end{array}$$

Coarse-ID control (static case)

$$\begin{array}{ll}\text{minimize}_u & x^* Q x \\ \text{subject to} & x = Bu + x_0\end{array}$$

B unknown!

Collect data: $\{(x_i, u_i)\}$ $x_i = Bu_i + x_0 + e_i$

Estimate B : $\underset{B}{\text{minimize}} \quad \sum_{i=1}^N \|Bu_i - x_i\|^2 \quad \rightarrow \hat{B}$

Guarantee: $\|B - \hat{B}\| \leq \epsilon$ with high probability

Relaxation:
(Triangle inequality!)

$$\begin{array}{ll}\text{minimize}_u & \|Q^{1/2}x\| + \epsilon\lambda\|u\| \\ \text{subject to} & x = \hat{B}u + x_0\end{array}$$

Generalization bound

$$\text{cost}(\hat{u}) \leq \text{cost}(u_\star) + 4\epsilon\lambda\|u_\star\|\|Q^{1/2}x_\star\| + 4\epsilon^2\lambda^2\|u_\star\|^2$$

Coarse-ID control (static case)

$$\begin{array}{ll} \text{minimize}_u & x^* Q x \\ \text{subject to} & x = Bu + x_0 \end{array}$$

B unknown!

Collect data: $\{(x_i, u_i)\}$ $x_i = Bu_i + x_0 + e_i$

Estimate B : $\underset{B}{\text{minimize}} \quad \sum_{i=1}^N \|Bu_i + x_0 - x_i\|^2 \quad \rightarrow \hat{B}$

Guarantee: $\|B - \hat{B}\| \leq \epsilon$ with high probability

Relaxation:
(Triangle inequality!)

$$\begin{array}{ll} \text{minimize}_u & \|Q^{1/2}x\| + \epsilon\lambda\|u\| \\ \text{subject to} & x = \hat{B}u + x_0 \end{array}$$

Generalization bound

$$\text{cost}(\hat{u}) = \text{cost}(u_\star) + \mathcal{O}(\epsilon)$$

Coarse-ID Control for LQR

$$\begin{aligned} \text{minimize} \quad & \lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t \right] \\ \text{s.t.} \quad & x_{t+1} = Ax_t + Bu_t + e_t \quad \text{Gaussian noise} \\ & \text{Assume stable } A \end{aligned}$$

Run an experiment for T steps with random input. Then

$$\text{minimize}_{(A,B)} \quad \sum_{i=1}^T \|x_{i+1} - Ax_i - Bu_i\|^2$$

If $T \geq \tilde{O} \left(\frac{\sigma^2(d+p)}{\lambda_{\min}(\Lambda_c)\epsilon^2} \right)$ where $\Lambda_c = A\Lambda_c A^* + BB^*$
controllability Gramian

then $\|A - \hat{A}\| \leq \epsilon$ and $\|B - \hat{B}\| \leq \epsilon$ w.h.p.

[Dean, Mania, Matni, R., Tu, 2017]
[Mania, R., Simchowitz, Tu, 2018]

Coarse-ID Control for LQR

$$\begin{aligned} \text{minimize}_u \quad & \sup_{\|\Delta_A\|_2 \leq \epsilon_A, \|\Delta_B\|_2 \leq \epsilon_B} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t \\ \text{s.t.} \quad & x_{t+1} = (\hat{A} + \Delta_A)x_t + (\hat{B} + \Delta_B)u_t \end{aligned}$$

Solving an SDP relaxation of this robust control problem yields

$$\frac{J(\hat{K}) - J_\star}{J_\star} \leq C \Gamma_{\text{cl}} \left(\lambda_{\min}(\Lambda_c)^{-1/2} + \|K_\star\|_2 \right) \sqrt{\frac{\sigma^2(d+p)}{T}} \quad \text{w.h.p.}$$

$$\Lambda_c = A \Lambda_c A^* + B B^*$$

controllability Gramian

$$\Gamma_{\text{cl}} := \|(zI - A - BK_\star)^{-1}\|_{\mathcal{H}_\infty}$$

closed loop gain

This also tells you when your cost is finite!
Extends to unstable A case as well.

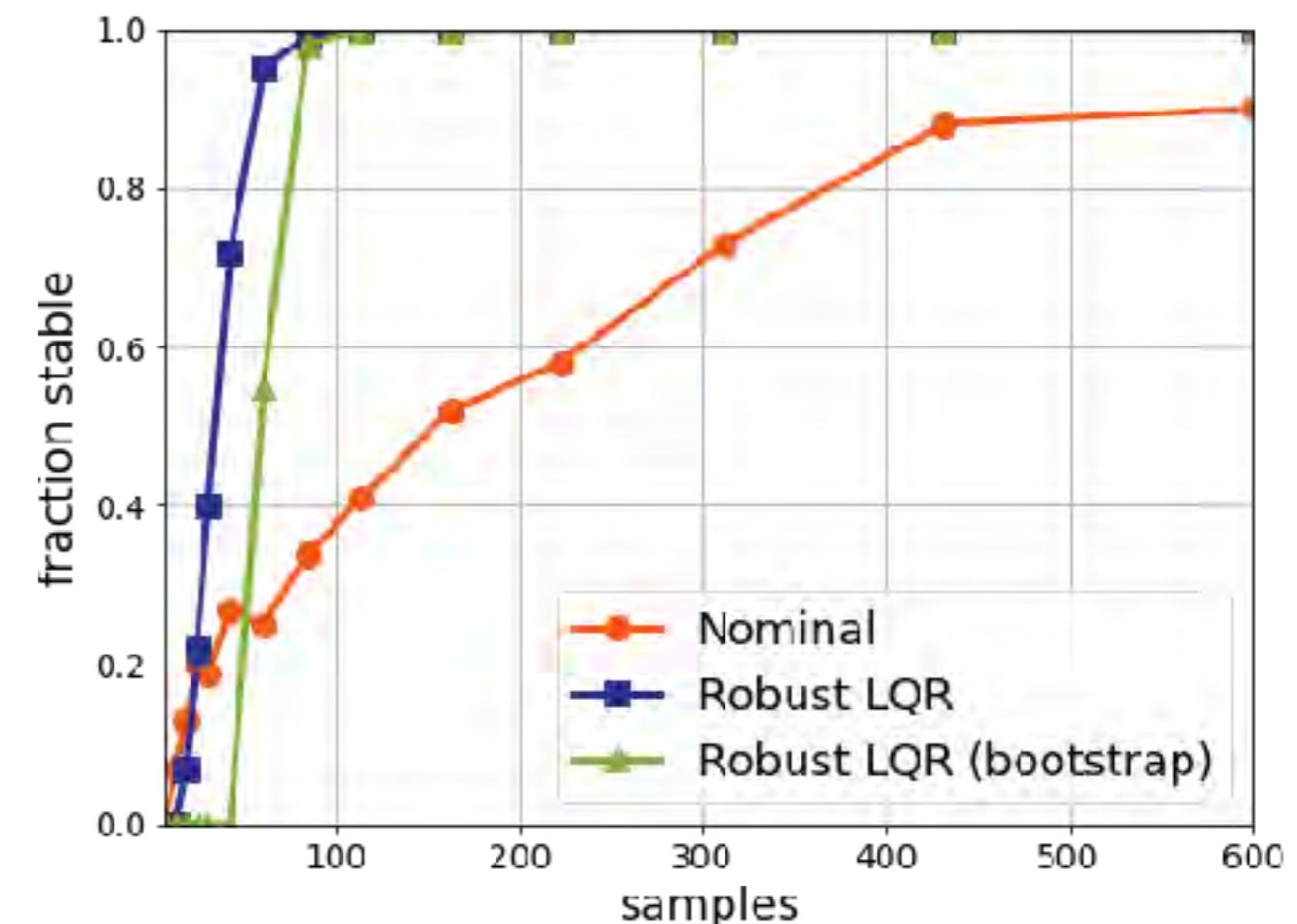
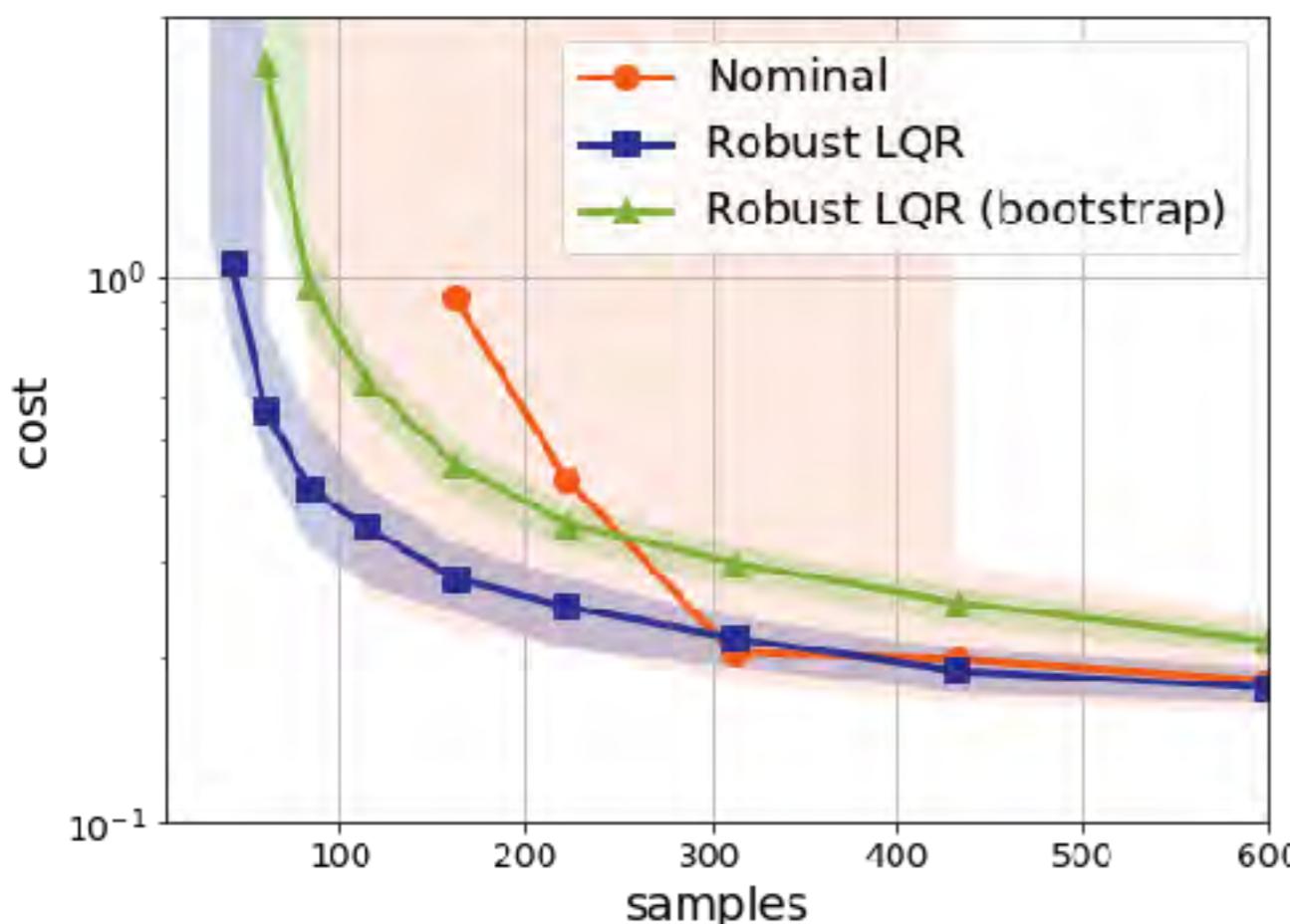


Why robust?



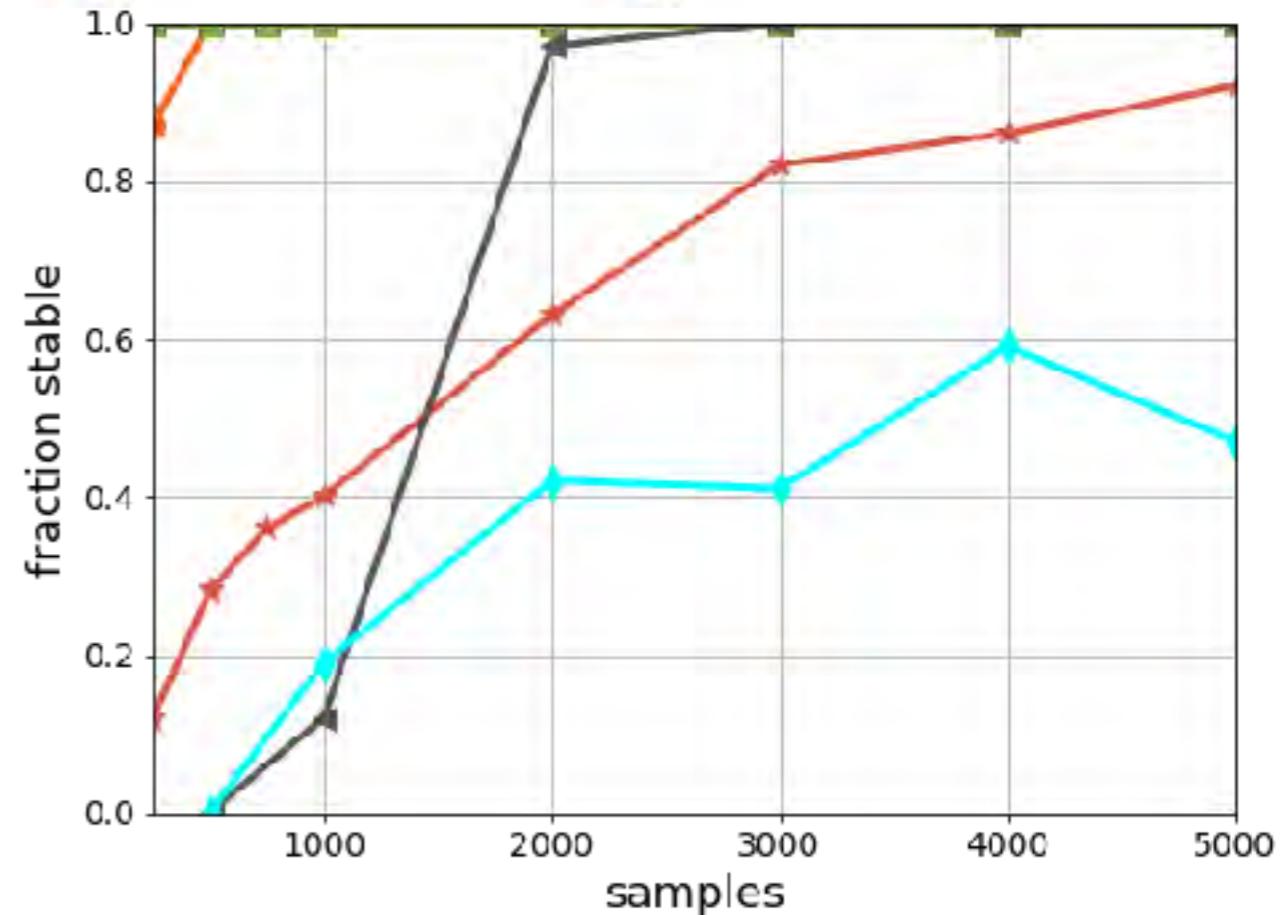
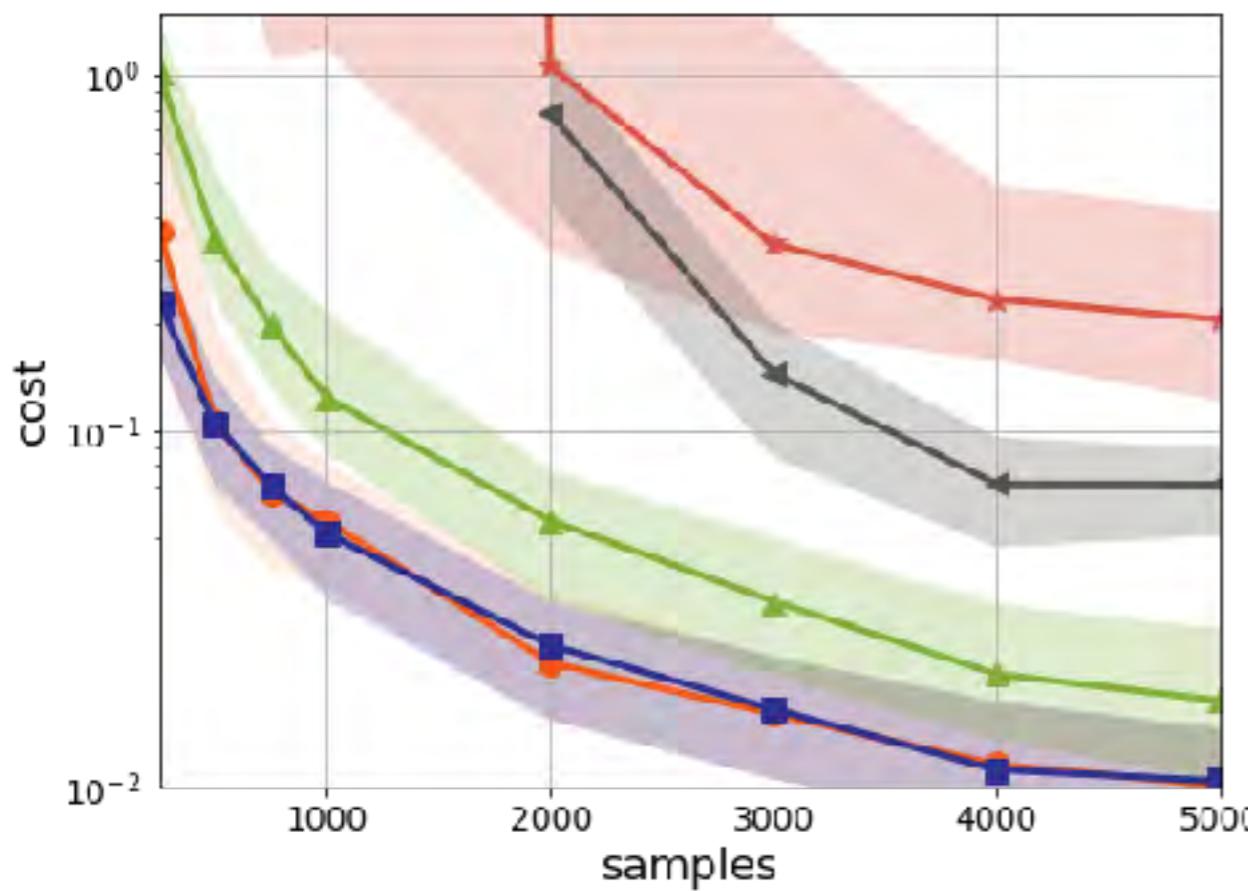
$$x_{t+1} = \begin{bmatrix} 1.01 & 0.01 & 0 \\ 0.01 & 1.01 & 0.01 \\ 0 & 0.01 & 1.01 \end{bmatrix} x_t + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} u_t + e_t$$

Slightly unstable system, system ID tends to think some nodes are stable



Least-squares estimate may yield unstable controller

Robust synthesis yields stable controller



Model-free
performs worse
than model-based

- Nominal
- Robust LQR
- ◆ Robust LQR (bootstrap)
- ★ LSPI
- ← Random Search
- ◆ Policy Gradient

Why has no one done this before?

- Coarse-ID control is the first non-asymptotic bound for this oracle model.
- Our guarantees for least-squares estimation required some heavy machinery.
 - Indeed, best bounds building on very recent papers.
- Our SDP relaxation uses brand new techniques in controller parameterization (*Systems Level Synthesis* by Matni et al.)
 - **Key insight:** Robustness makes analysis tractable
- *The Singularity has arrived!*
- Lots of work in the last years, to be highlighted in extended bib.

Even LQR is not simple!!!

$$\begin{aligned} \text{minimize} \quad & J := \sum_{t=1}^{\infty} x_t^T Q x_t + u_t^T R u_t \\ \text{s.t.} \quad & x_{t+1} = Ax_t + Bu_t + e_t \quad \text{Gaussian noise} \end{aligned}$$

$$\frac{J(\hat{K}) - J_*}{J_*} \leq C \Gamma_{\text{cl}} \left(\lambda_{\min}(\Lambda_c)^{-1/2} + \|K_*\|_2 \right) \sqrt{\frac{\sigma^2(d+p)\log(1/\delta)}{n}}$$

where $\Lambda_c = A\Lambda_c A^* + BB^*$
controllability Gramian

$\Gamma_{\text{cl}} := \|(zI - A - BK_*)^{-1}\|_{\mathcal{H}_\infty}$
closed loop gain



Hard to estimate
Control insensitive to mismatch



Easy to estimate
Control very sensitive to mismatch

50 papers on Cosma Shalizi's blog say otherwise!

Need to fix learning theory for time series.

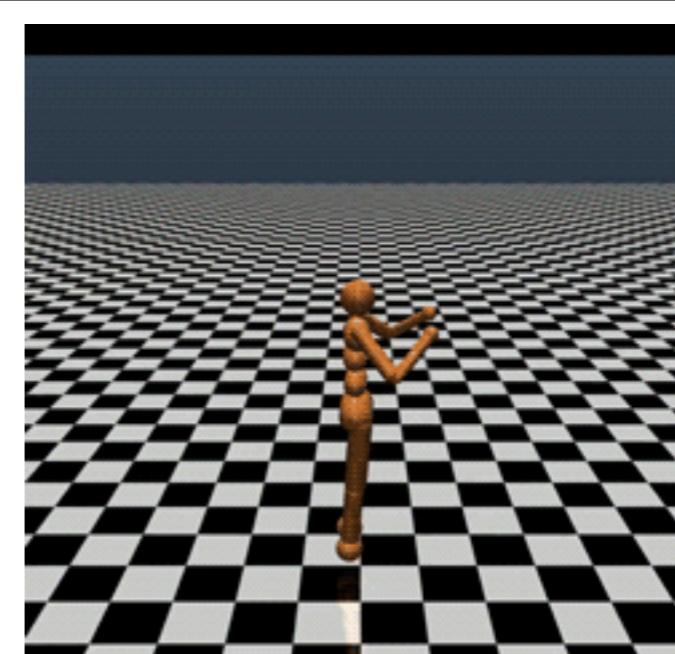
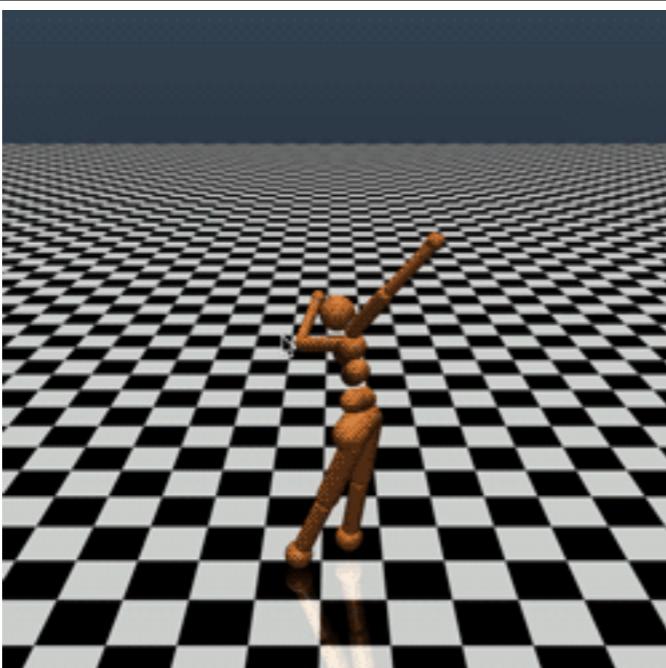
The Linearization Principle

If a machine learning algorithm does crazy things when restricted to linear models, it's going to do crazy things on complex nonlinear models too.

What happens when we return to nonlinear models?

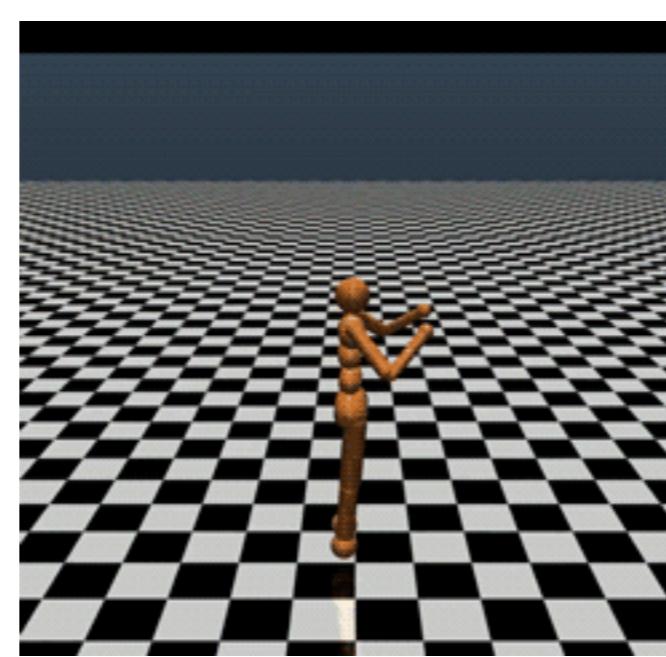
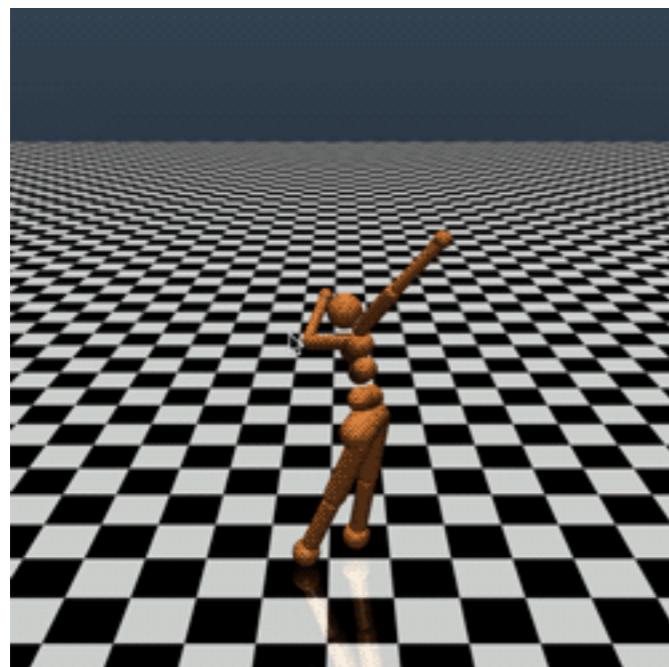
Random search of linear policies outperforms Deep Reinforcement Learning

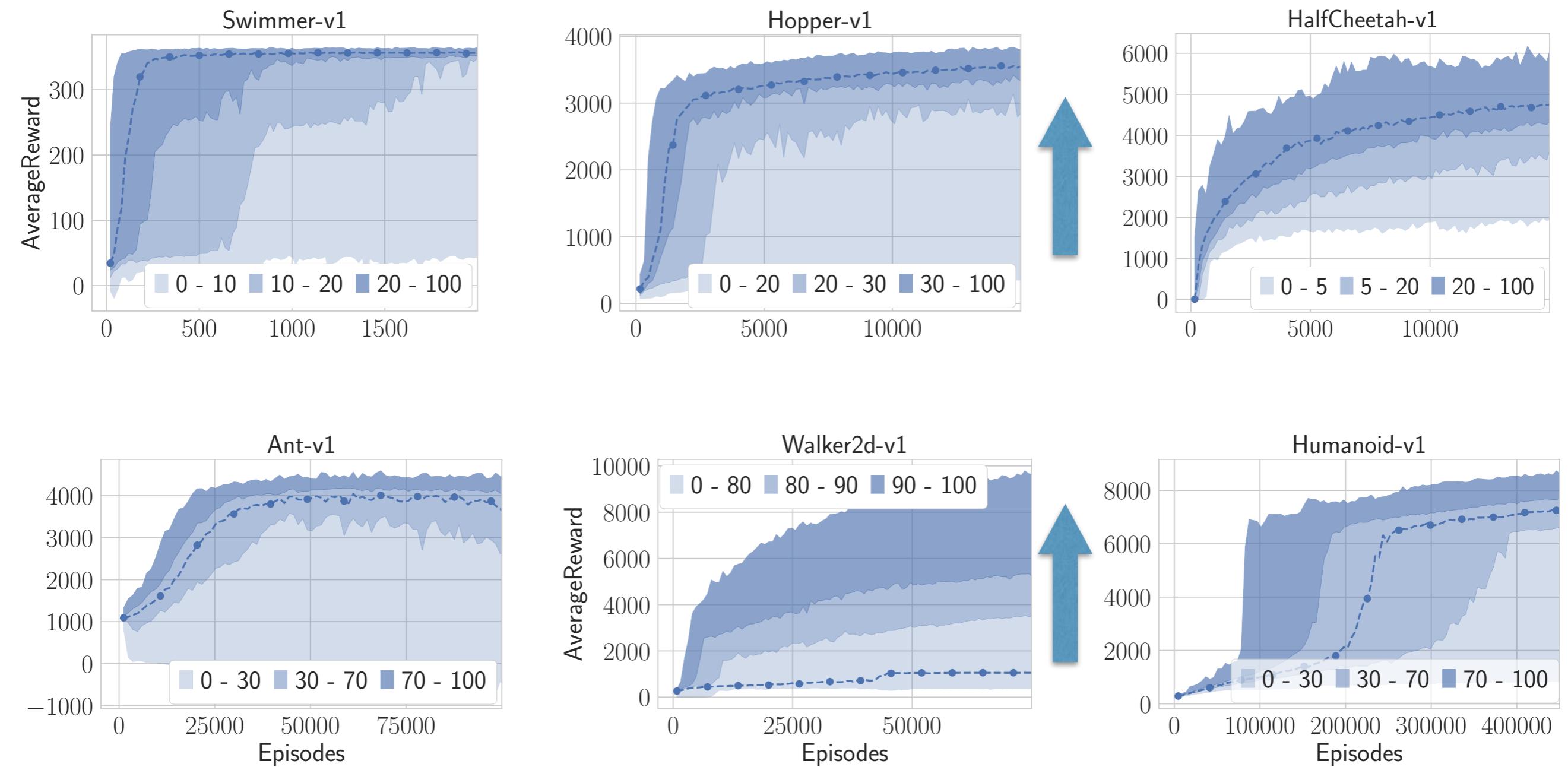
<i>Larger</i> is better		Maximum average reward		
Task	RS	NG-lin	NG-rbf	TRPO-nn
Swimmer-v1	365	366	365	131
Hopper-v1	3909	3651	3810	3668
HalfCheetah-v1	6722	4149	6620	4800
Walker	11389	5234	5867	5594
Ant	5146	4607	4816	5007
Humanoid	11600	6440	6849	6482



Larger is better

Task	RS	Maximum average reward		
		NG-lin	NG-rbf	TRPO-nn
Swimmer-v1	365	366	365	131
Hopper-v1	3909	3651	3810	3668
HalfCheetah-v1	6722	4149	6620	4800
Walker	11389	5234	5867	5594
Ant	5146	4607	4816	5007
Humanoid	11600	6440	6849	6482





Larger is better

Model Predictive Control

$$\text{minimize} \quad \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) + C_f(x_{T+1}) \right]$$

$$\text{s.t.} \quad x_{t+1} = f_t(x_t, u_t, e_t), \quad x_1 = x$$
$$u_t = \pi_t(\tau_t)$$

Optimal Policy: $\boxed{\pi(x) = \arg \min_u Q_1(x, u)}$

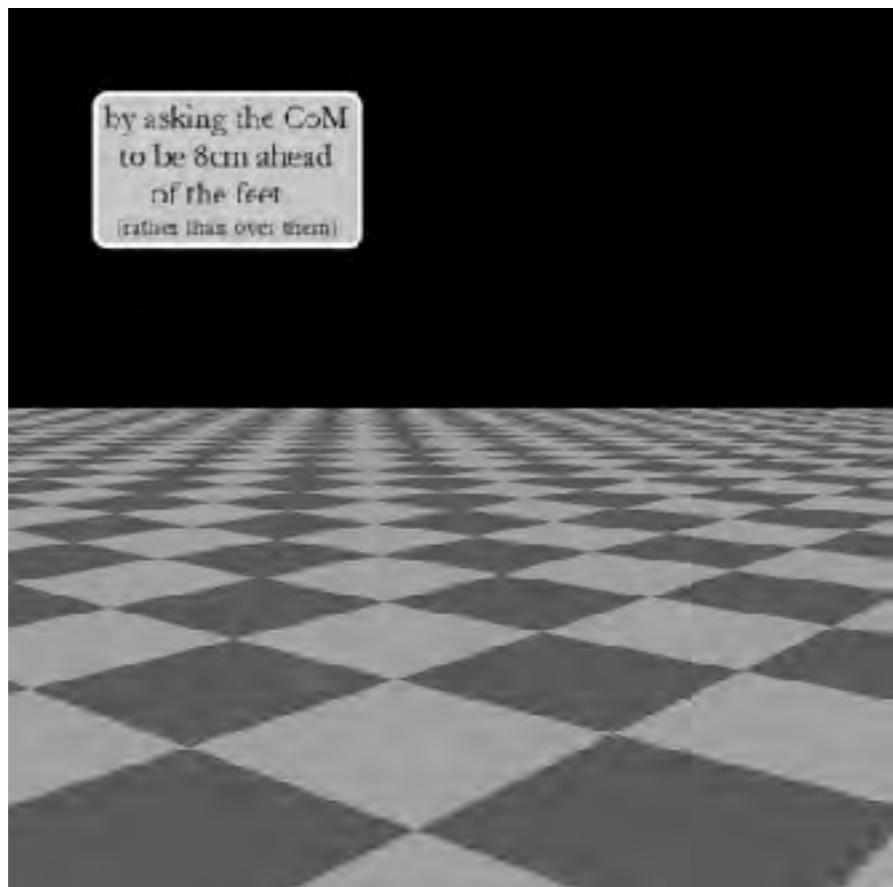
$$Q_1(x, u) = C_1(x, u) + \mathbb{E}_e \left[\min_{u'} Q_2(f_1(x, u, e), u') \right]$$

MPC: use the Q-function for *all* time steps

$$Q_1(x, u) = \sum_{t=1}^H C_t(x, u) + \mathbb{E}_e \left[\min_{u'} Q_{H+1}(f_H(x, u, e), u') \right]$$

MPC ethos: plan on short time horizons, use feedback to correct modeling error and disturbance.

Model Predictive Control



Videos from Todorov Lab

<https://homes.cs.washington.edu/~todorov/>

Learning in MPC

$$\begin{aligned} \text{minimize} \quad & \mathbb{E}_e \left[\sum_{t=1}^T C_t(x_t, u_t) + C_f(x_{T+1}) \right] \\ \text{s.t.} \quad & x_{t+1} = f_t(x_t, u_t, e_t), \quad x_1 = x \\ & u_t = \pi_t(\tau_t) \end{aligned}$$

Optimal Policy:

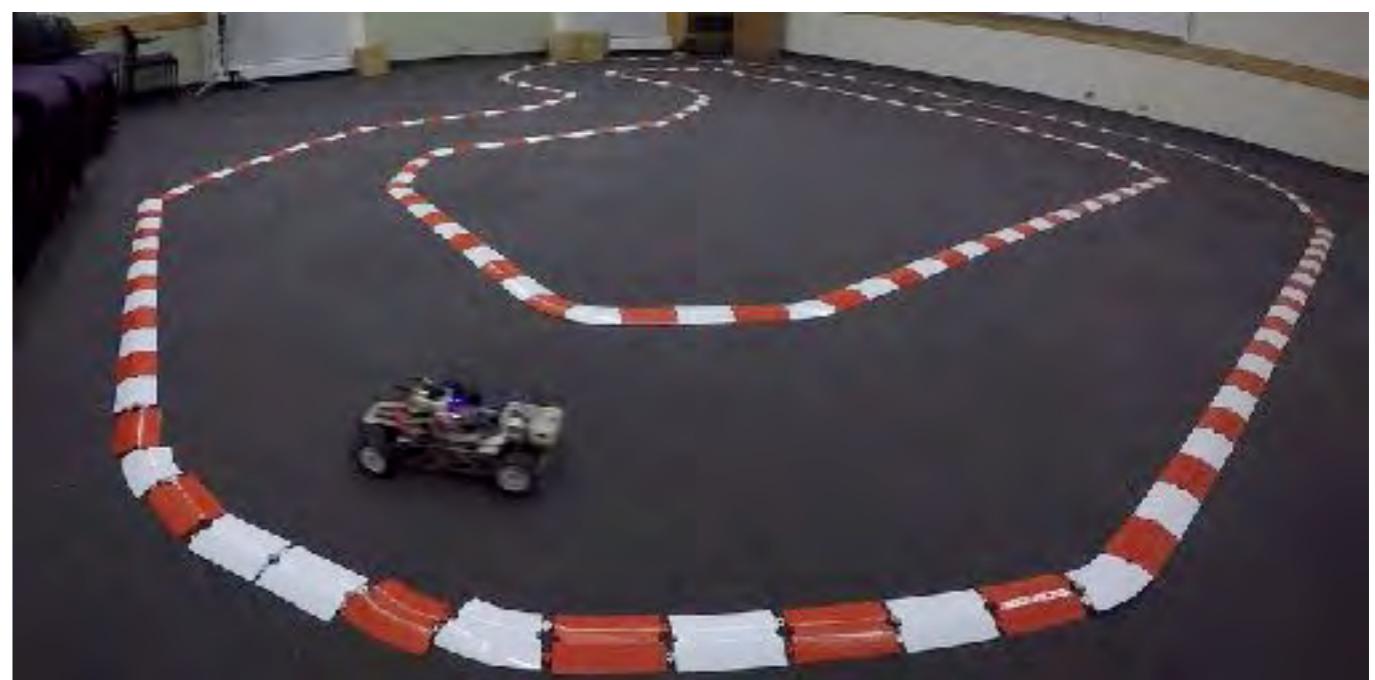
$$\pi(x) = \arg \min_u Q_1(x, u)$$

MPC: use the Q-function for *all* time steps

$$Q_1(x, u) = \sum_{t=1}^H C_t(x, u) + \mathbb{E}_e \left[\min_{u'} Q_{H+1}(f_H(x, u, e), u') \right]$$

Use past data to learn the terminal Q-function:

The value of a state is the minimum value seen for the remainder of the episode from that state.



So many things left to do...

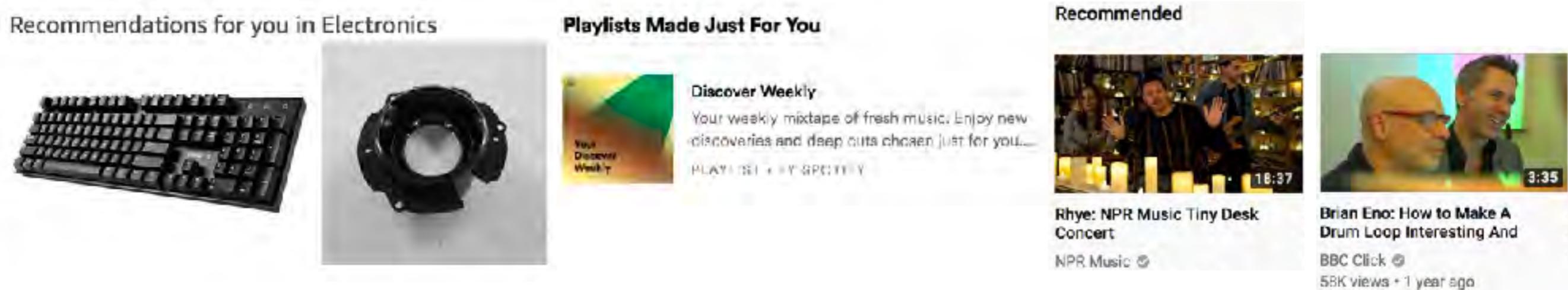
- Are the coarse-ID results optimal, even with respect to the parameters?
- Tight upper and lower sample complexities for LQR. (Is the optimal error scaling $T^{-1/2}$ or T^{-1} ?)
- Finite analysis of learning in MPC.
- Adaptive control.
- Iterative learning control.
- Nonlinear models, constraints, and improper learning.
- Safe exploration, earning about uncertain environments.
- Implementing in test-beds.

Actionable Intelligence

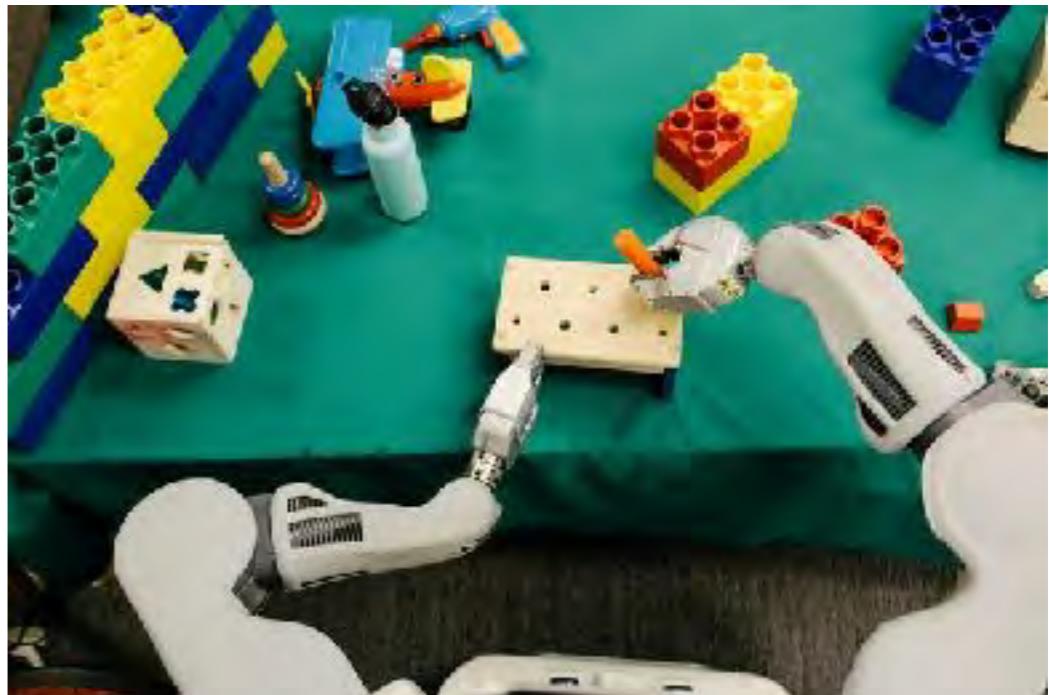
~~Control Theory~~

~~Reinforcement Learning~~ is the study of how to use past data to enhance the future manipulation of a dynamical system

Actionable Intelligence is the study of how to use past data to enhance the future manipulation of a dynamical system

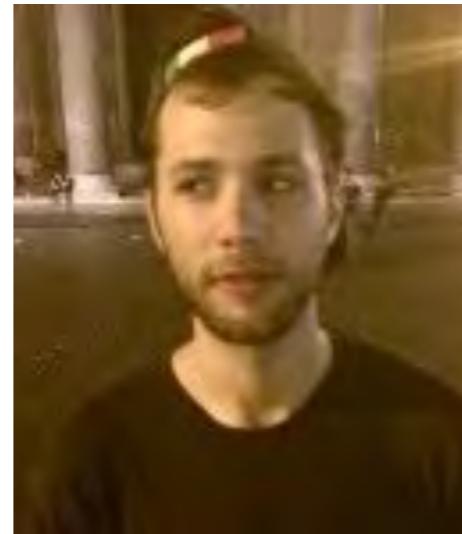


As soon as a machine learning system is unleashed in feedback with humans, that system is an actionable intelligence system, not a machine learning system.



Actionable Intelligence
trustable, scalable, predictable

Collaborators



Joint work with Sarah Dean, Aurelia Guy, Horia Mania, Nikolai Matni, Max Simchowitz, and Stephen Tu.

Recommended Texts

- D. Bertsekas. *Dynamic Programming and Optimal Control*. 4th edition, volumes 1 (2017) and 2 (2012). Athena Scientific.
- D. Bertsekas. and J.Tsitsiklis. *Neuro-dynamic Programming*. Athena Scientific, 1996.
- F. Borrelli, A. Bemporad, and M. Morari. *Predictive Control for Linear and Hybrid Systems*. Cambridge, 2017.
- B. Recht “A Tour of Reinforcement Learning: The View from Continuous Control.” [arXiv:1806.09460](https://arxiv.org/abs/1806.09460)

References from the Actionable Intelligence Lab

- [argmin.net](#)
- “On the Sample Complexity of the Linear Quadratic Regulator” S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu. [arXiv:1710.01688](#)
- “Non-asymptotic Analysis of Robust Control from Coarse-grained Identification.” S. Tu, R. Boczar, A. Packard, and B. Recht. [arXiv:1707.04791](#)
- “Least-squares Temporal Differencing for the Linear Quadratic Regulator” S. Tu and B. Recht. In submission to ICML 2018. [arXiv:1712.08642](#)
- “Learning without Mixing.” H. Mania, B. Recht, M. Simchowitz, and S. Tu. In submission to COLT 2018. [arXiv:1802.08334](#)
- “Simple random search provides a competitive approach to reinforcement learning.” H. Mania, A. Guy, and B. Recht. [arXiv:1803.07055](#)
- “Regret Bounds for Robust Adaptive Control of the Linear Quadratic Regulator.” S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu. [arXiv:1805.09388](#)

<https://people.eecs.berkeley.edu/~brecht/publications.html>

$$\begin{aligned}
 & \underset{u}{\text{minimize}} && \lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T x_t^* Q x_t + u_t^* R u_t \right] \\
 & \text{s.t.} && x_{t+1} = Ax_t + Bu_t + e_t
 \end{aligned}$$

Key to formulation:
 Write (x,u) as linear
 function of disturbance

$$\begin{bmatrix} x_t \\ u_t \end{bmatrix} = \sum_{k=1}^t \begin{bmatrix} \Phi_x[k] \\ \Phi_u[k] \end{bmatrix} e_{t-k}$$

$$\mathbb{E} [x_t^* Q x_t] = \sigma^2 \sum_{k=1}^t \text{Tr}(\Phi_x[k]^* Q \Phi_x[k])$$

$$\mathbb{E} [u_t^* R u_t] = \sigma^2 \sum_{k=1}^t \text{Tr}(\Phi_u[k]^* R \Phi_u[k])$$

$$\begin{aligned}
 & \underset{\Phi}{\text{minimize}} && \left\| \begin{bmatrix} Q^{\frac{1}{2}} & 0 \\ 0 & R^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \Phi_x \\ \Phi_u \end{bmatrix} \right\|_F^2 \\
 & \text{s.t.} && \Phi_x[t+1] = A\Phi_x[t] + B\Phi_u[t] \\
 & && \Phi_x[0] = I
 \end{aligned}$$

Key to formulation:
Write (x,u) as linear
function of disturbance

$$\begin{bmatrix} x_t \\ u_t \end{bmatrix} = \sum_{k=1}^t \begin{bmatrix} \Phi_x[k] \\ \Phi_u[k] \end{bmatrix} e_{t-k}$$

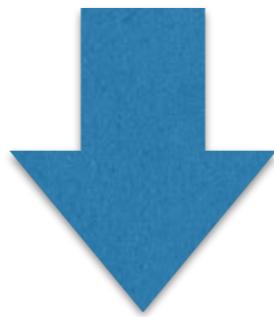
minimize
 Φ

$$\sup_{\|\Delta_A\|_2 \leq \epsilon_A, \|\Delta_B\|_2 \leq \epsilon_B} \left\| \begin{bmatrix} Q^{\frac{1}{2}} & 0 \\ 0 & R^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \Phi_x \\ \Phi_u \end{bmatrix} \right\|_F^2$$

s.t.

$$\Phi_x[t+1] = (\hat{A} + \Delta_A)\Phi_x[t] + (\hat{B} + \Delta_B)\Phi_u[t]$$

$$\Phi_x[0] = I$$



As in the static case, push
robustness into cost.

minimize
 Φ

s.t.

$$\Phi_x[t+1] = \hat{A}\Phi_x[t] + \hat{B}\Phi_u[t]$$

$$\Phi_x[0] = I$$

$$\sup_{\|\Delta_A\|_2 \leq \epsilon_A, \|\Delta_B\|_2 \leq \epsilon_B} \left\| \begin{bmatrix} Q^{\frac{1}{2}} & 0 \\ 0 & R^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \Phi_x \\ \Phi_u \end{bmatrix} (I + \Delta)^{-1} \right\|_F^2$$

SLS Formulation of Robust LQR

$$\begin{aligned} & \underset{\gamma \in [0, 1)}{\text{minimize}} \quad \frac{1}{1-\gamma} \quad \underset{\Phi_x, \Phi_u}{\text{min}} \quad \left\| \begin{bmatrix} Q^{\frac{1}{2}} & 0 \\ 0 & R^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \Phi_x \\ \Phi_u \end{bmatrix} \right\|_F \\ & \text{s.t.} \quad \Phi_x[t+1] = \hat{A}\Phi_x[t] + \hat{B}\Phi_u[t] \\ & \quad \Phi_x[0] = I \\ & \quad \sup_{\theta \in [0, 2\pi)} \left\| \sum_{k=1}^{T_{\max}} \begin{bmatrix} \frac{\epsilon_A}{\sqrt{\alpha}} \Phi_x[k] \\ \frac{\epsilon_B}{\sqrt{1-\alpha}} \Phi_u[k] \end{bmatrix} e^{ik\theta} \right\| \leq \gamma \end{aligned}$$

- Solvable by SDP for fixed γ
- Binary search over γ to find optimal solution

SLS Formulation of Robust LQR

$$\begin{aligned}
 & \text{minimize}_{X, Z, W, \gamma} \quad \frac{1}{(1-\gamma)^2} \{ \text{Trace}(QW_{11}) + \text{Trace}(RW_{22}) \} \\
 & \text{subject to} \quad \begin{bmatrix} X & X & Z^* \\ X & W_{11} & W_{12} \\ Z & W_{21} & W_{22} \end{bmatrix} \succeq 0 \\
 & \quad \begin{bmatrix} X - I & \hat{A}X + \hat{B}Z & 0 & 0 \\ (\hat{A}X + \hat{B}Z)^* & X & \epsilon_A X & \epsilon_B Z^* \\ 0 & \epsilon_A X & \alpha\gamma^2 I & 0 \\ 0 & \epsilon_B Z & 0 & (1 - \alpha)\gamma^2 I \end{bmatrix} \succeq 0.
 \end{aligned}$$

- Solvable by SDP for fixed γ
- Binary search over γ to find optimal solution
- Optimal controller is $K = -ZX^{-1}$