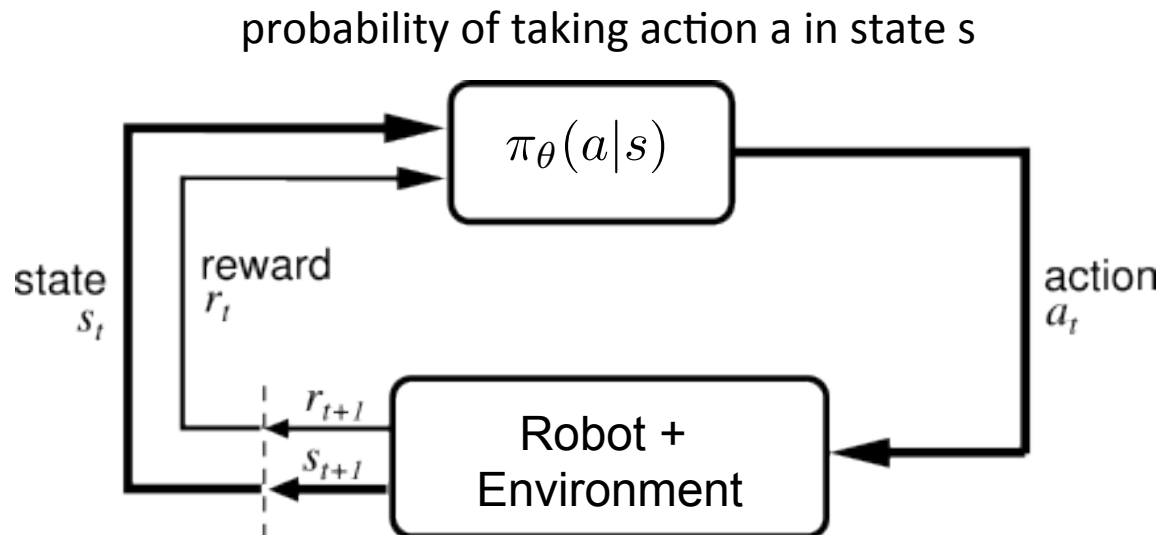# Deep Reinforcement Learning for Robotics
## Pieter Abbeel

## UC Berkeley / OpenAI / gradescope.com [code: ICML2016]

# Deep Reinforcement Learning (RL)
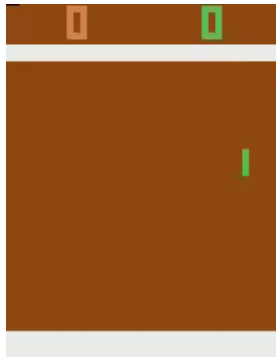
probability of taking action a in state s



$$\pi_\theta(a|s)$$

state $s_t$

reward $r_t$

$r_{t+1}$

$s_{t+1}$

Robot + Environment

action $a_t$

- Goal:

$$\max_\theta \quad \mathrm{E}[\sum_{t=0}^{H} R(s_t)|\pi_\theta]$$

- **Additional challenges:**

  - **Stability**

  - **Credit assignment**

  - **Exploration**

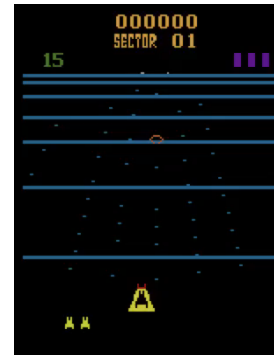Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# From Pixels to Actions?



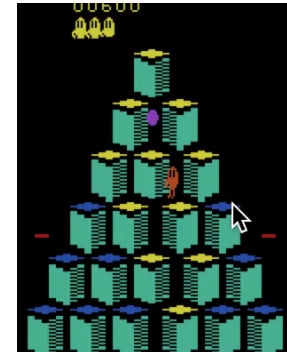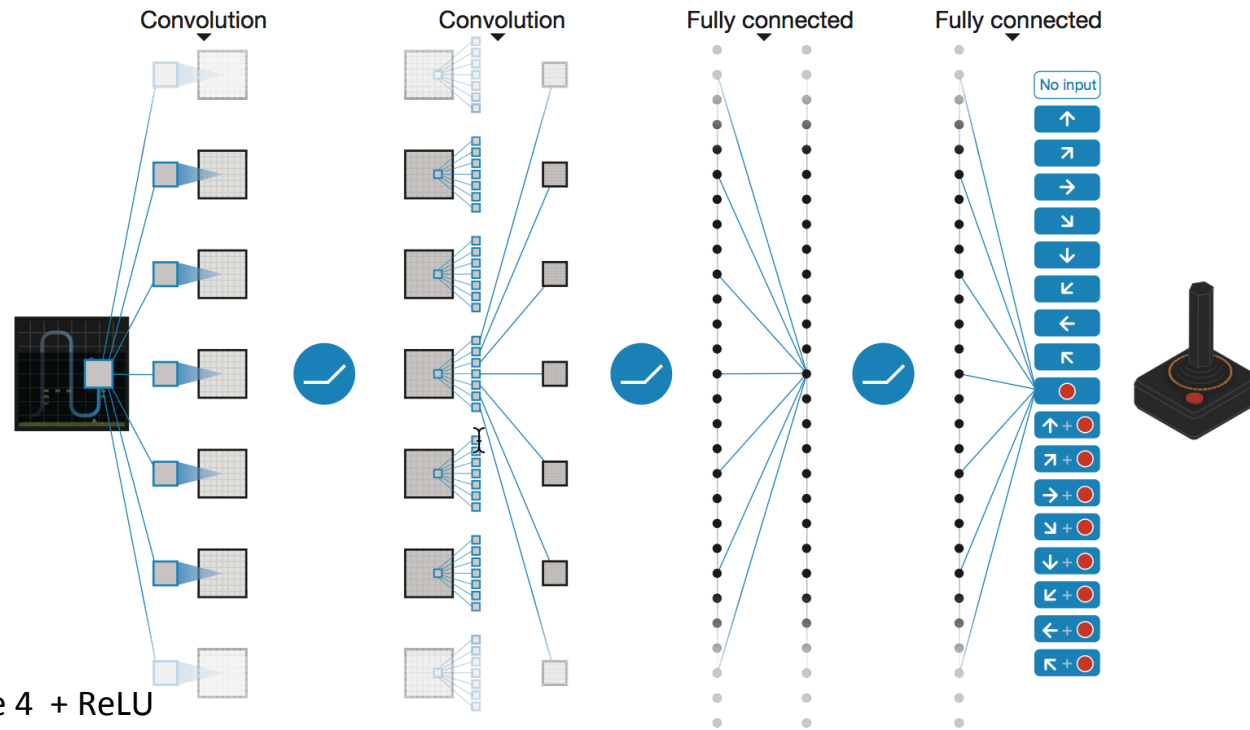Pong    Enduro    Beamrider    Q*bert
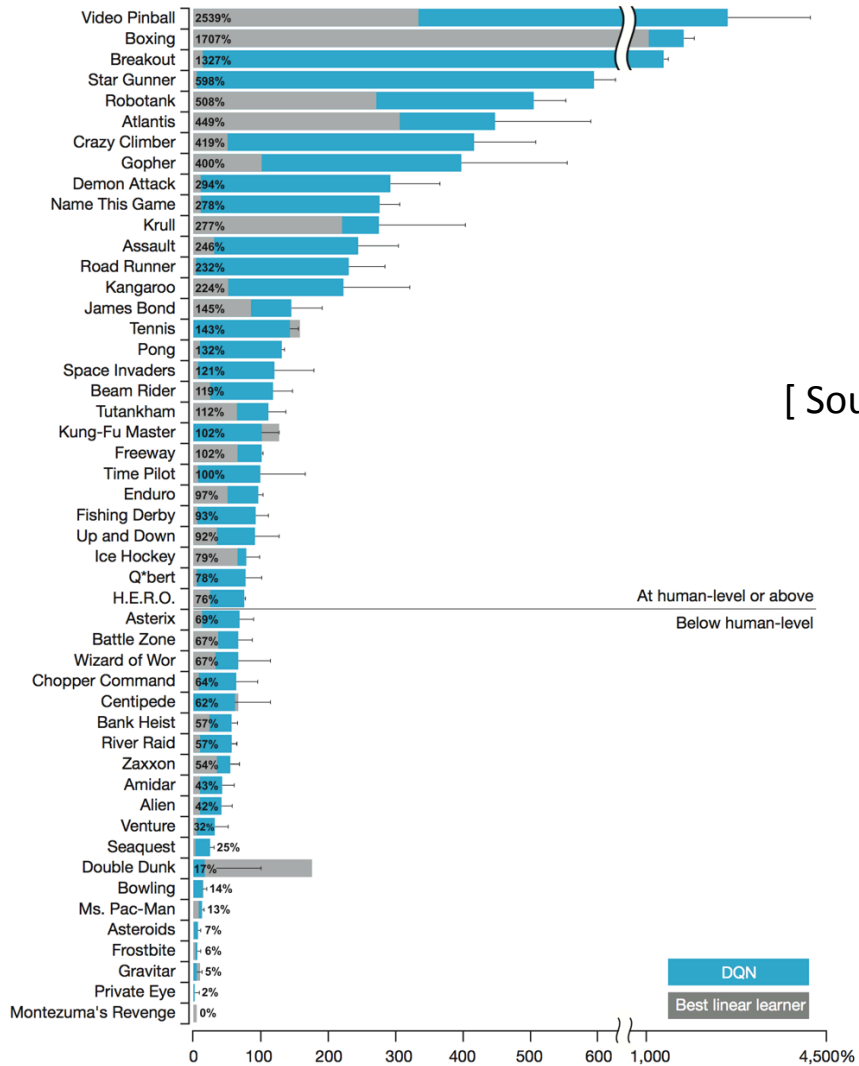
# Deep Q-Network (DQN): From Pixels to Joystick Commands



32 8x8 filters with stride 4  + ReLU
64 4x4 filters with stride 2  + ReLU
64 3x3 filters with stride 1  + ReLU
fully connected 512 units + ReLU
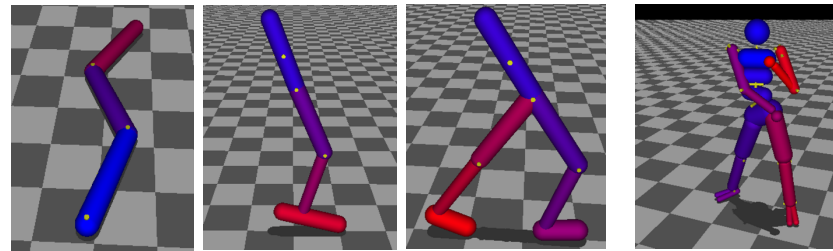fully connected output units, one per action

[Source: Mnih et al., Nature 2015 (DeepMind) ]

Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

Video Pinball 2539%
Boxing 1707%
Breakout 1327%
Star Gunner 598%
Robotank 508%
Atlantis 449%
Crazy Climber 419%
Gopher 400%
Demon Attack 294%
Name This Game 278%
Krull 277%
Assault 246%
Road Runner 232%
Kangaroo 224%
James Bond 145%
Tennis 143%
Pong 132%
Space Invaders 121%
Beam Rider 119%
Tutankham 112%
Kung-Fu Master 102%
Freeway 102%
Time Pilot 100%
Enduro 97%
Fishing Derby 93%
Up and Down 92%
Ice Hockey 79%
Q*bert 78%
H.E.R.O. 76%
Asterix 69%
Battle Zone 67%
Wizard of Wor 67%
Chopper Command 64%
Centipede 62%
Bank Heist 57%
River Raid 57%
Zaxxon 54%
Amidar 43%
Alien 42%
Venture 32%
Seaquest 25%
Double Dunk 17%
Bowling 14%
Ms. Pac-Man 13%
Asteroids 7%
Frostbite 6%
Gravitar 5%
Private Eye 2%
Montezuma's Revenge 0%

At human-level or above
Below human-level

DQN
Best linear learner

0   100   200   300   400   500   600   1,000   4,500%

[ Source: Mnih et al., Nature 2015 (DeepMind) ]

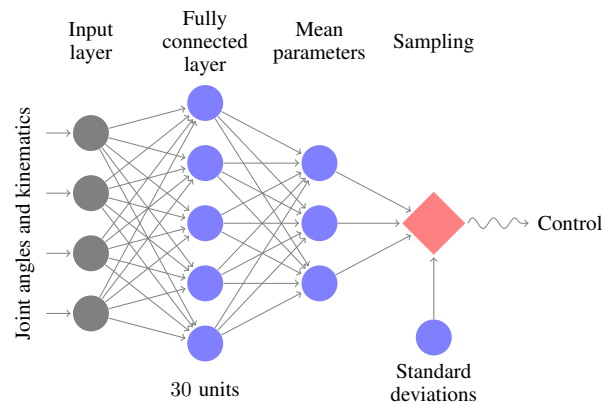Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# How About Continuous Control, e.g., Locomotion?

Robot models in physics simulator
(MuJoCo, from Emo Todorov)



Input: joint angles and velocities
Output: joint torques

Neural network architecture:

# Challenges with Q-Learning

- How to score every possible action?

- How to ensure monotonic progress?

# Policy Optimization

$$\max_{\theta} \quad \mathrm{E}[\sum_{t=0}^{H} R(s_t)|\pi_\theta]$$

- Often simpler to represent good policies than good value functions

- True objective of expected cost is optimized (vs. a surrogate like Bellman error)

- Existing work: (natural) policy gradients

  - Challenges:  good, large step directions

# Trust Region Policy Optimization

$$\max_{\theta} \quad \mathrm{E}[\sum_{t=0}^{H} R(s_t)|\pi_\theta]$$

$$\max_{\delta\theta} \quad \hat{L}(\theta + \delta\theta)$$

$$\text{s.t.} \quad \mathrm{KL}\left(P(\tau;\theta)||P(\tau;\theta+\delta\theta)\right) \leq \varepsilon$$
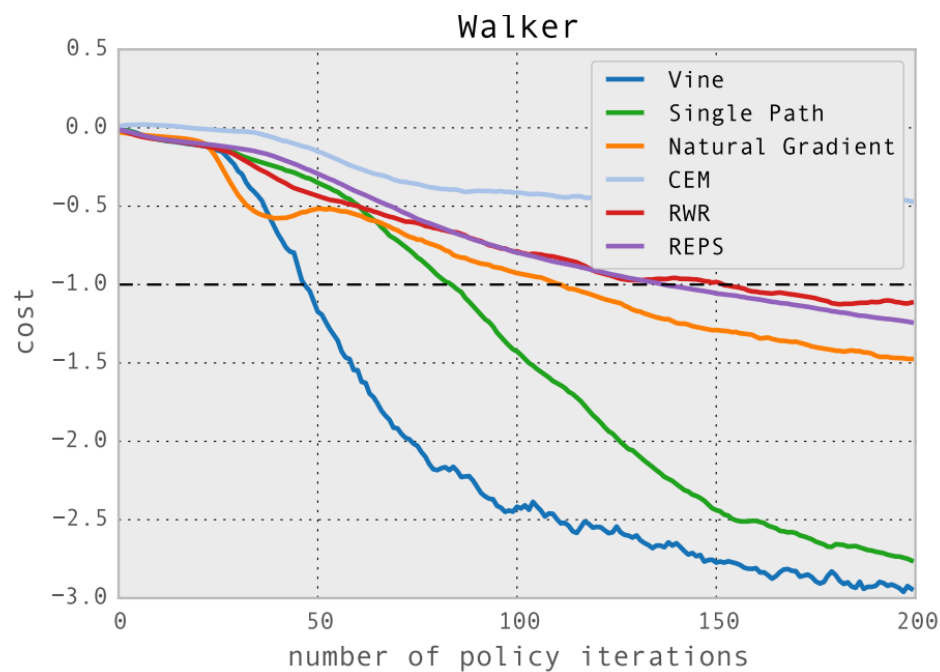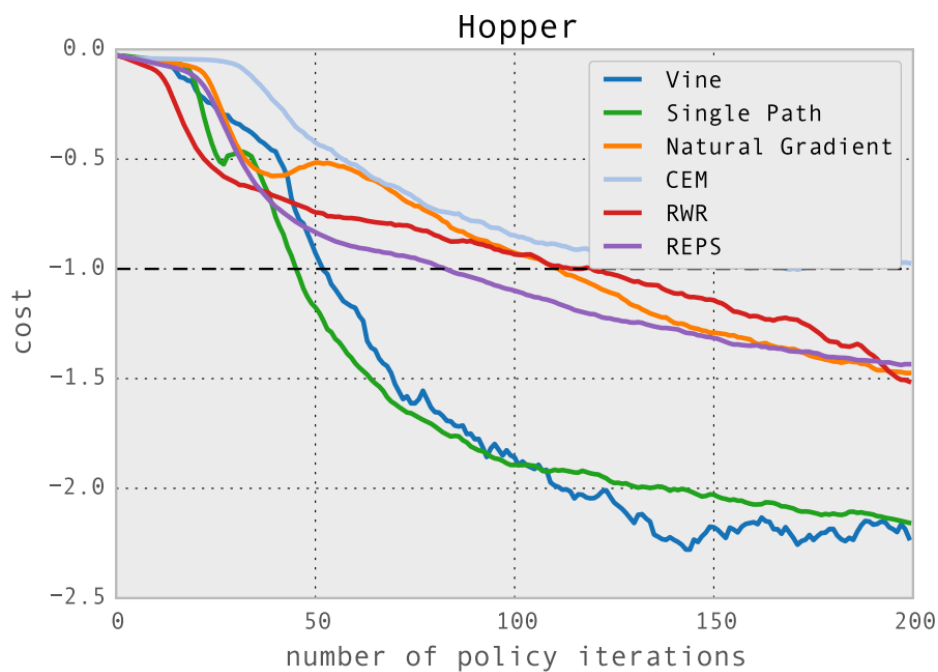
- $\hat{L}$ : Surrogate Objective

- $\mathrm{KL}$ : Trust region

# Experiments in Locomotion



Our algorithm was tested on
three locomotion problems
in a physics simulator

The following gaits were obtained

# Learning Curves -- Comparison

# Generalized Advantage Estimation (GAE)

Objective:

$$\max_{\theta} \quad \mathrm{E}[\sum_{t=0}^{H} R(s_t)|\pi_\theta]$$

Gradient:
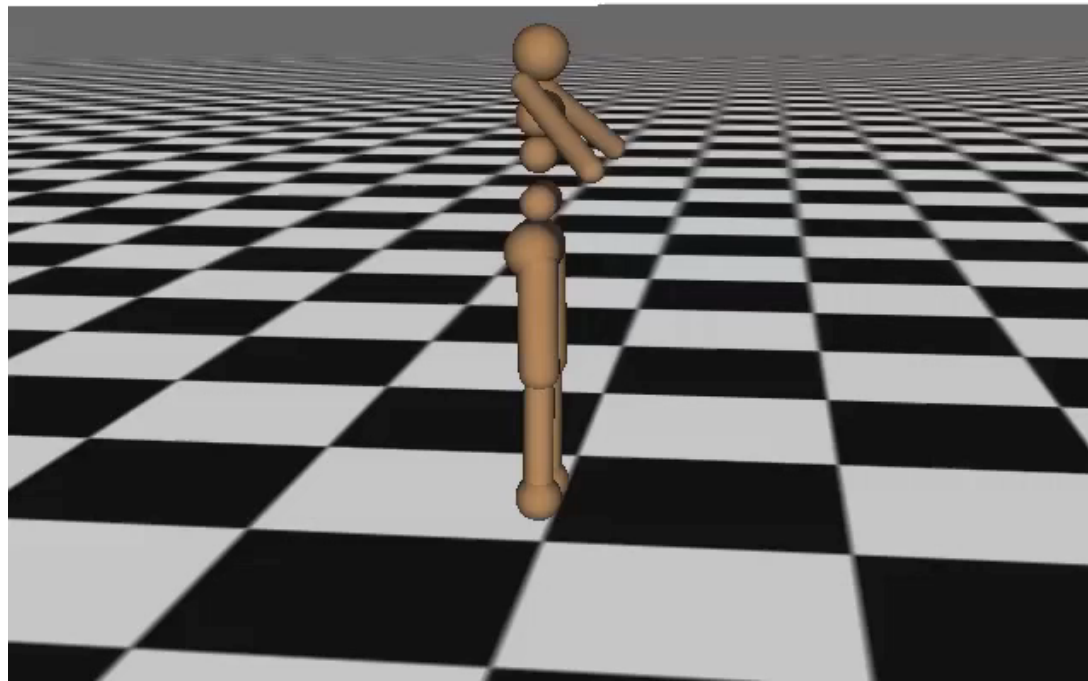
$$\mathrm{E}[\sum_{t=0}^{H} \nabla_\theta \log \pi_\theta(a_t|s_t)\left(\sum_{k=t}^{H} R(s_k) - V(s_t)\right)]$$

single sample estimate of advantage

- Generalized Advantage Estimation

  - Exponential interpolation between actor-critic and Monte Carlo estimates

  - Trust region approach to (high-dimensional) value function estimation

[Schulman, Moritz, Levine, Jordan, Abbeel, ICLR 2016]

# Learning Locomotion through Trust Region Policy Optimization (TRPO)

Iteration 0



[Schulman, Moritz, Levine, Jordan, Abbeel, ICLR 2016]

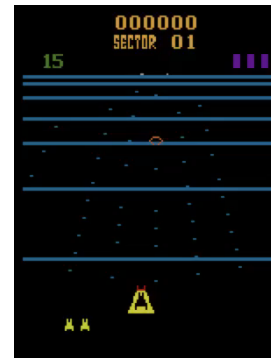Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# Atari Games

- Deep Q-Network (DQN) [Mnih et al, 2013/2015]

- Dagger with Monte Carlo Tree Search [Xiao-Xiao et al, 2014]

- Trust Region Policy Optimization [Schulman, Levine, Moritz, Jordan, Abbeel, 2015]
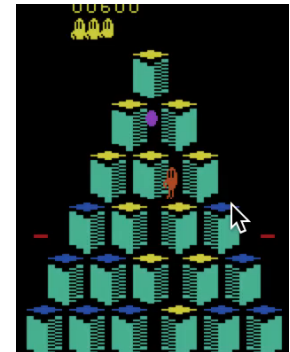
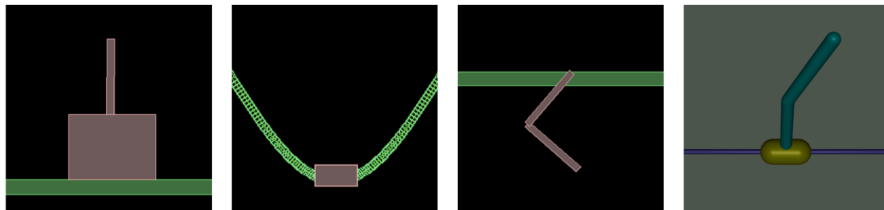- A3C [Mnih et al., 2016]



Pong          Enduro          Beamrider          Q*bert

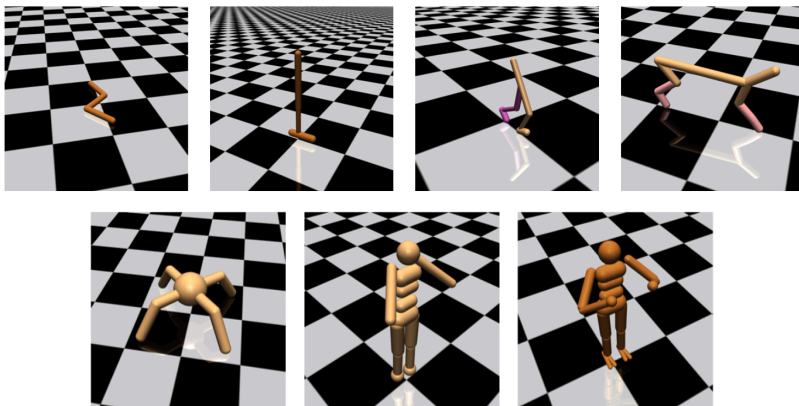# Deep RL Benchmarking

- Tasks

- Algorithms

- Experimental setup

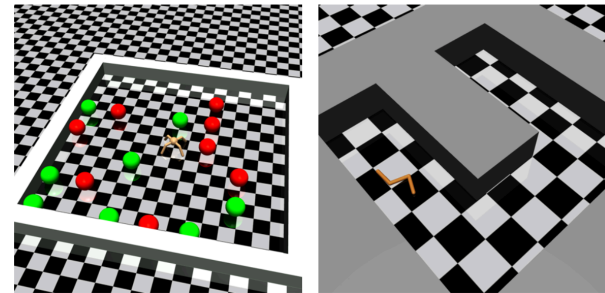Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# Deep RL Benchmarking -- Tasks

## 1. Basic tasks



## 2. Locomotion



## 3. Hierarchical



## 4. Partially observable

sensing, delayed action, sysID

## 5. Driving…

Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# Deep RL Benchmarking -- Algorithms

- Reinforce

- Truncated Natural Policy Gradient

- Reward-Weighted Regression (RWR)

- Relative Entropy Policy Search (REPS)

- Trust-Region Policy Optimization (TRPO)

- Cross-Entropy Method (CEM)

- Covariance Matrix Adaptation Evolution Strategy (CMA-ES)
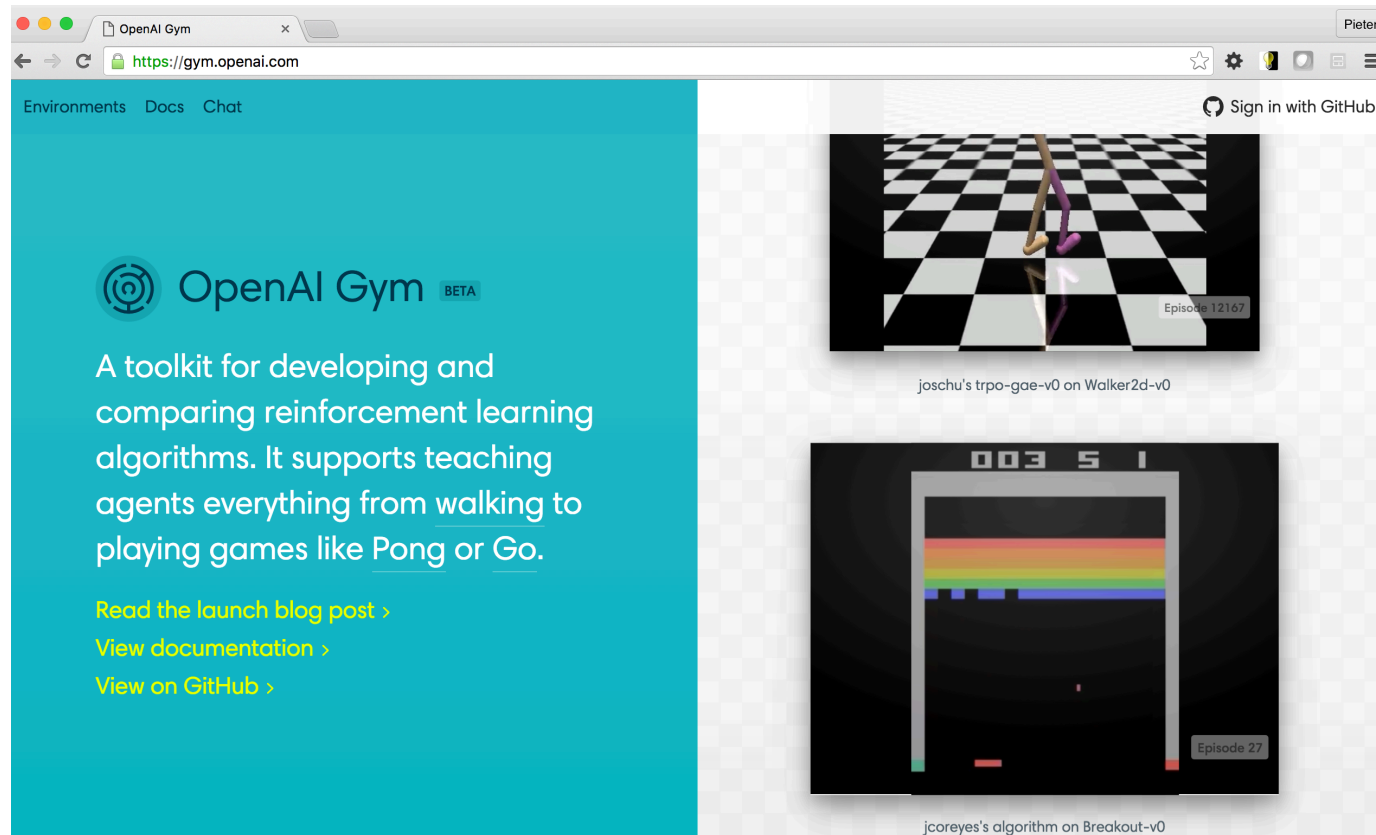
- Deep Deterministic Policy Gradients (DDPG)

- …

Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# Benchmarking [Duan et al, ICML 2016]

*Table 1.* Performance of the implemented algorithms in terms of average return over all training iterations for five different random seeds (same across all algorithms). The results of the best-performing algorithm on each task are highlighted in boldface. In the tasks column, the partially observable variants of the tasks are annotated as follows: LS stands for limited sensors, NO for noisy observations and delayed actions, and SI for system identifications. The notation N/A denotes that an algorithm has failed on the task at hand, e.g., CMA-ES leading to out-of-memory errors in the Full Humanoid task.

| Task | Random | VPG | TNPG | RWR | REPS | TRPO | CEM | CMA-ES |
|---|---|---|---|---|---|---|---|---|
| Cart-Pole Balancing | $77.1 \pm 0.0$ | $4693.7 \pm 14.0$ | $3986.4 \pm 748.9$ | $4861.5 \pm 12.3$ | $565.6 \pm 137.6$ | **$4869.8 \pm 37.6$** | $4815.4 \pm 4.8$ | $2440.4 \pm 568.3$ |
| Inverted Pendulum* | $-153.4 \pm 0.2$ | $13.4 \pm 18.0$ | $209.7 \pm 55.5$ | $84.7 \pm 13.8$ | $-113.3 \pm 4.6$ | **$247.2 \pm 76.1$** | $38.2 \pm 25.7$ | $-40.1 \pm 5.7$ |
| Mountain Car | $-415.4 \pm 0.0$ | $-67.1 \pm 1.0$ | $-66.5 \pm 4.5$ | $-79.4 \pm 1.1$ | $-275.6 \pm 166.3$ | **$-61.7 \pm 0.9$** | $-66.0 \pm 2.4$ | $-85.0 \pm 7.7$ |
| Acrobot | $-1904.5 \pm 1.0$ | $-508.1 \pm 91.0$ | $-395.8 \pm 121.2$ | $-352.7 \pm 35.9$ | $-1001.5 \pm 10.8$ | **$-326.0 \pm 24.4$** | $-436.8 \pm 14.7$ | $-785.6 \pm 13.1$ |
| Double Inverted Pendulum* | $149.7 \pm 0.1$ | $4116.5 \pm 65.2$ | **$4455.4 \pm 37.6$** | $3614.8 \pm 368.1$ | $446.7 \pm 114.8$ | $4412.4 \pm 50.4$ | $2566.2 \pm 178.9$ | $1576.1 \pm 51.3$ |
| Swimmer* | $-1.7 \pm 0.1$ | $92.3 \pm 0.1$ | **$96.0 \pm 0.2$** | $60.7 \pm 5.5$ | $3.8 \pm 3.3$ | $96.0 \pm 0.2$ | $68.8 \pm 2.4$ | $64.9 \pm 1.4$ |
| Hopper | $8.4 \pm 0.0$ | $714.0 \pm 29.3$ | $1155.1 \pm 57.9$ | $553.2 \pm 71.0$ | $86.7 \pm 17.6$ | **$1183.3 \pm 150.0$** | $63.1 \pm 7.8$ | $20.3 \pm 14.3$ |
| 2D Walker | $-1.7 \pm 0.0$ | $506.5 \pm 78.8$ | **$1382.6 \pm 108.2$** | $136.0 \pm 15.9$ | $-37.0 \pm 38.1$ | $1353.8 \pm 85.0$ | $84.5 \pm 19.2$ | $77.1 \pm 24.3$ |
| Half-Cheetah | $-90.8 \pm 0.3$ | $1183.1 \pm 69.2$ | $1729.5 \pm 184.6$ | $376.1 \pm 28.2$ | $34.5 \pm 38.0$ | **$1914.0 \pm 120.1$** | $330.4 \pm 274.8$ | $441.3 \pm 107.6$ |
| Ant* | $13.4 \pm 0.7$ | $548.3 \pm 55.5$ | $706.0 \pm 127.7$ | $37.6 \pm 3.1$ | $39.0 \pm 9.8$ | **$730.2 \pm 61.3$** | $49.2 \pm 5.9$ | $17.8 \pm 15.5$ |
| Simple Humanoid | $41.5 \pm 0.2$ | $87.6 \pm 2.3$ | $276.6 \pm 31.7$ | $56.7 \pm 3.8$ | $36.8 \pm 4.0$ | **$493.9 \pm 75.5$** | $98.8 \pm 11.5$ | $115.9 \pm 2.7$ |
| Full Humanoid | $13.2 \pm 0.1$ | $214.4 \pm 35.0$ | $282.9 \pm 19.8$ | $42.3 \pm 3.4$ | $45.5 \pm 4.1$ | **$455.9 \pm 22.8$** | $104.0 \pm 14.5$ | N/A $\pm$ N/A |
| Cart-Pole Balancing (LS)* | $77.1 \pm 0.0$ | $420.9 \pm 265.5$ | $945.1 \pm 27.8$ | $68.9 \pm 1.5$ | $898.1 \pm 22.1$ | **$960.2 \pm 46.0$** | $227.0 \pm 223.0$ | $68.0 \pm 1.6$ |
| Inverted Pendulum (LS) | $-122.1 \pm 0.1$ | $-13.4 \pm 3.2$ | $0.7 \pm 6.1$ | $-107.4 \pm 0.2$ | $-87.2 \pm 8.0$ | **$4.5 \pm 4.1$** | $-81.2 \pm 33.2$ | $-62.4 \pm 3.4$ |
| Mountain Car (LS) | $-83.0 \pm 0.0$ | $-81.2 \pm 0.6$ | $-65.7 \pm 9.0$ | $-81.7 \pm 0.1$ | $-82.6 \pm 0.4$ | **$-64.2 \pm 9.5$** | $-68.9 \pm 1.3$ | $-73.2 \pm 0.6$ |
| Acrobot (LS)* | $-393.2 \pm 0.0$ | $-128.9 \pm 11.6$ | $-84.6 \pm 2.9$ | $-235.9 \pm 5.3$ | $-379.5 \pm 1.4$ | **$-83.3 \pm 9.9$** | $-149.5 \pm 15.3$ | $-159.9 \pm 7.5$ |
| Cart-Pole Balancing (NO)* | $101.4 \pm 0.1$ | $616.0 \pm 210.8$ | **$916.3 \pm 23.0$** | $93.8 \pm 1.2$ | $99.6 \pm 7.2$ | $606.2 \pm 122.2$ | $181.4 \pm 32.1$ | $104.4 \pm 16.0$ |
| Inverted Pendulum (NO) | $-122.2 \pm 0.1$ | $6.5 \pm 1.1$ | **$11.5 \pm 0.5$** | $-110.0 \pm 1.4$ | $-119.3 \pm 4.2$ | $10.4 \pm 2.2$ | $-55.6 \pm 16.7$ | $-80.3 \pm 2.8$ |
| Mountain Car (NO) | $-83.0 \pm 0.0$ | $-74.7 \pm 7.8$ | $-64.5 \pm 8.6$ | $-81.7 \pm 0.1$ | $-82.9 \pm 0.1$ | **$-60.2 \pm 2.0$** | $-67.4 \pm 1.4$ | $-73.5 \pm 0.5$ |
| Acrobot (NO)* | $-393.5 \pm 0.0$ | $-186.7 \pm 31.3$ | $-164.5 \pm 13.4$ | $-233.1 \pm 0.4$ | $-258.5 \pm 14.0$ | **$-149.6 \pm 8.6$** | $-213.4 \pm 6.3$ | $-236.6 \pm 6.2$ |
| Cart-Pole Balancing (SI)* | $76.3 \pm 0.1$ | $431.7 \pm 274.1$ | **$980.5 \pm 7.3$** | $69.0 \pm 2.8$ | $702.4 \pm 196.4$ | $980.3 \pm 5.1$ | $746.6 \pm 93.2$ | $71.6 \pm 2.9$ |
| Inverted Pendulum (SI) | $-121.8 \pm 0.2$ | $-5.3 \pm 5.6$ | **$14.8 \pm 1.7$** | $-108.7 \pm 4.7$ | $-92.8 \pm 23.9$ | $14.1 \pm 0.9$ | $-51.8 \pm 10.6$ | $-63.1 \pm 4.8$ |
| Mountain Car (SI) | $-82.7 \pm 0.0$ | $-63.9 \pm 0.2$ | $-61.8 \pm 0.4$ | $-81.4 \pm 0.1$ | $-80.7 \pm 2.3$ | **$-61.6 \pm 0.4$** | $-63.9 \pm 1.0$ | $-66.9 \pm 0.6$ |
| Acrobot (SI)* | $-387.8 \pm 1.0$ | $-169.1 \pm 32.3$ | **$-156.6 \pm 38.9$** | $-233.2 \pm 2.6$ | $-216.1 \pm 7.7$ | $-170.9 \pm 40.3$ | $-250.2 \pm 13.7$ | $-245.0 \pm 5.5$ |
| Swimmer + Gathering | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ |
| Ant + Gathering | $-5.8 \pm 5.0$ | **$-0.1 \pm 0.1$** | $-0.4 \pm 0.1$ | $-5.5 \pm 0.5$ | $-6.7 \pm 0.7$ | $-0.4 \pm 0.0$ | $-4.7 \pm 0.7$ | N/A $\pm$ N/A |
| Swimmer + Maze | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ |
| Ant + Maze | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | N/A $\pm$ N/A |

# rllab
## [Duan et al, ICML 2016]



Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# Open AI Gym

# Curiosity-driven Exploration

$$r'(s_t, a_t, s_{t+1}) = r(s_t, a_t) + \eta D_{\mathrm{KL}}[p(\theta|\xi_t, a_t, s_{t+1}) \| p(\theta|\xi_t)]$$
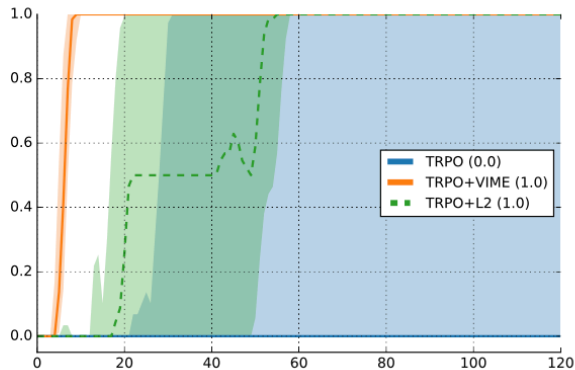
Building on:
    Curiosity: Schmidhuber , 1991; Sun, Gomez,
        Schmidhuber, 2011; Schmidhuber, 2010
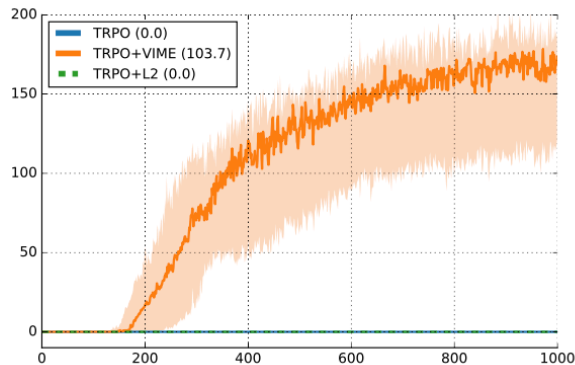    Bayesian neural nets: Blundell, Cornebise, Kavukcuoglu,
        Wierstra, 2015

# Curiosity-driven Exploration

(a) MountainCar

(b) CartPoleSwingup

(c) HalfCheetah

# Curiosity-driven Exploration

[Houthooft, Chen, Duan, Schulman, Turck, Abbeel, 2016]



Swimmer + Food Collection

# Curiosity-driven Exploration

[Houthooft, Chen, Duan, Schulman, Turck, Abbeel, 2016]



**TRPO**

**TRPO + VIME**

Swimmer + Food Collection

Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope
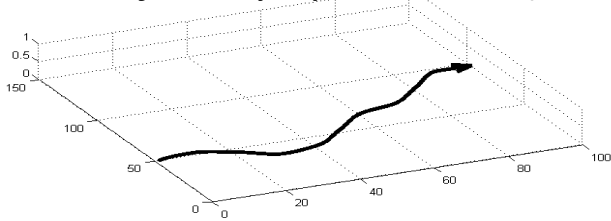
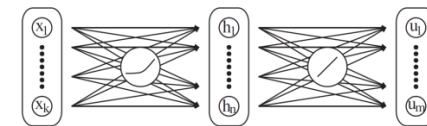# How About Real Robotic Visuo-Motor Skills?



Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# Guided Policy Search

Model-Based RL (through trajectory optimization)



**Supervised learning**
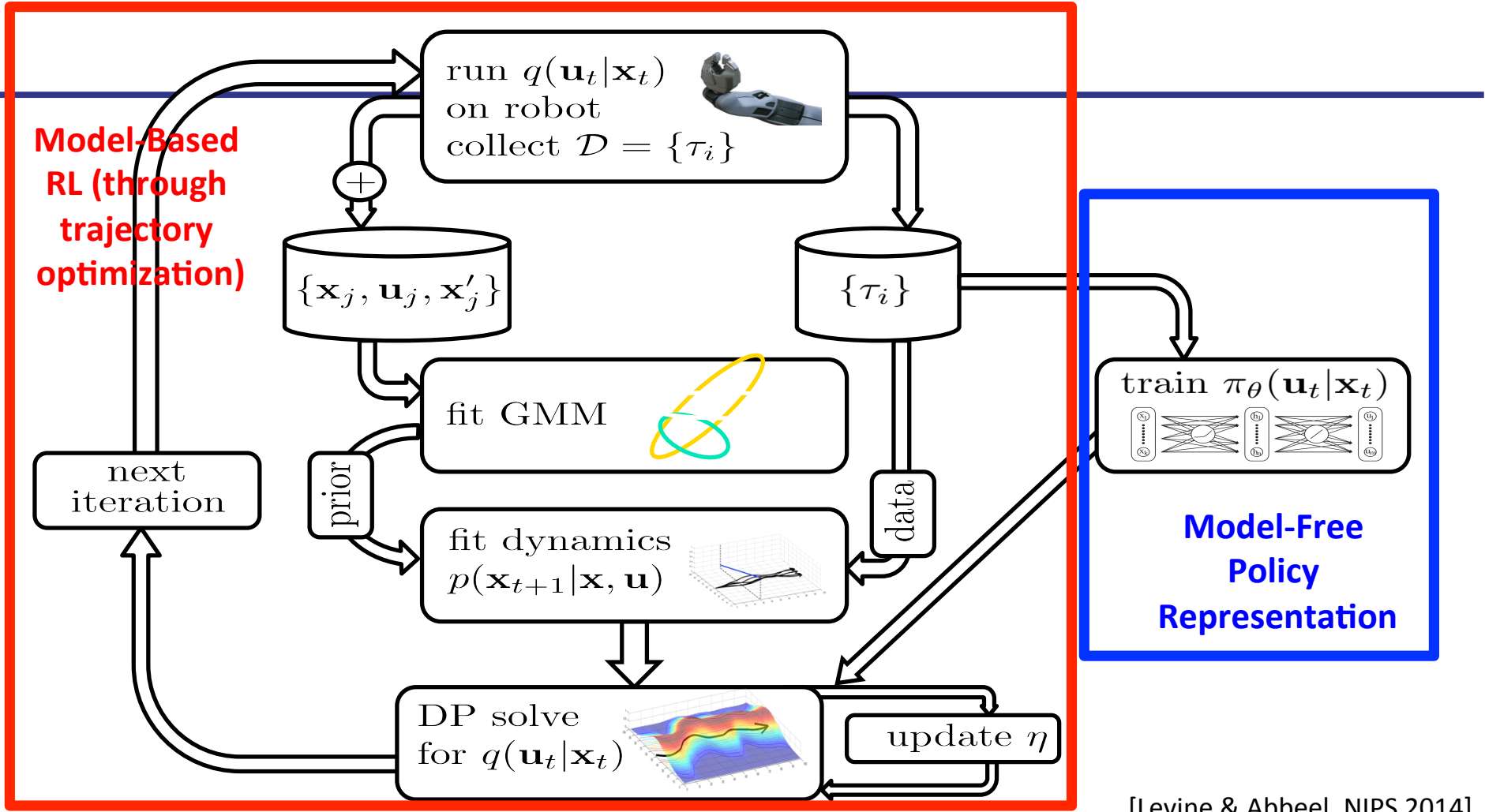
General Neural Net Policy



- Issue with two-phase pipeline

  - Representational mismatch trajectory distribution vs. neural net

→ Joint optimization

$$\max_{\{\pi^{(i)}\},\theta} \sum_i \mathrm{E}[\sum_{t=0}^{H} R(s_t^{(i)}) \mid \pi^{(i)}] - \lambda \sum_i \|\pi^{(i)} - \pi_\theta\|$$

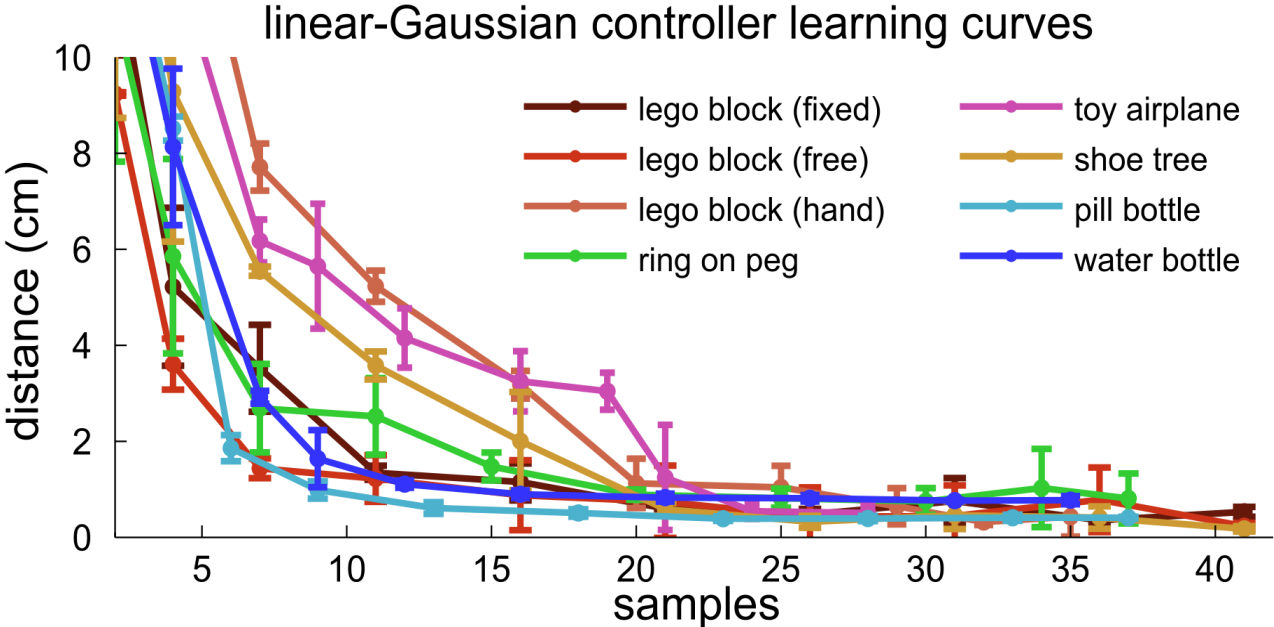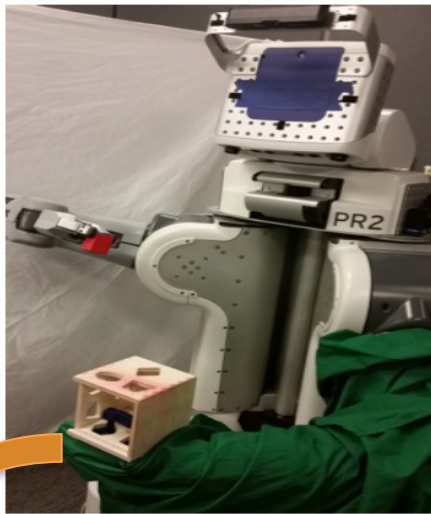**Model-Based RL (through trajectory optimization)**

run $q(\mathbf{u}_t|\mathbf{x}_t)$ on robot collect $\mathcal{D} = \{\tau_i\}$

$\{\mathbf{x}_j, \mathbf{u}_j, \mathbf{x}'_j\}$

$\{\tau_i\}$

fit GMM

next iteration

prior

data

fit dynamics $p(\mathbf{x}_{t+1}|\mathbf{x}, \mathbf{u})$

DP solve for $q(\mathbf{u}_t|\mathbf{x}_t)$

update $\eta$

[Levine & Abbeel, NIPS 2014]

**Model-Based RL (through trajectory optimization)**

run $q(\mathbf{u}_t|\mathbf{x}_t)$ on robot collect $\mathcal{D} = \{\tau_i\}$

$\{\mathbf{x}_j, \mathbf{u}_j, \mathbf{x}'_j\}$

$\{\tau_i\}$

fit GMM

next iteration

prior

fit dynamics $p(\mathbf{x}_{t+1}|\mathbf{x}, \mathbf{u})$

data

train $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$

**Model-Free Policy Representation**

DP solve for $q(\mathbf{u}_t|\mathbf{x}_t)$

update $\eta$

[Levine & Abbeel, NIPS 2014]

# Linear-Gaussian Controller Learning Curves



linear-Gaussian controller learning curves

# Instrumented Training



training time

test time

$$\mathbf{x}_t \rightarrow \mathbf{u}_t$$

$$\mathbf{o}_t \rightarrow \mathbf{u}_t$$

Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope
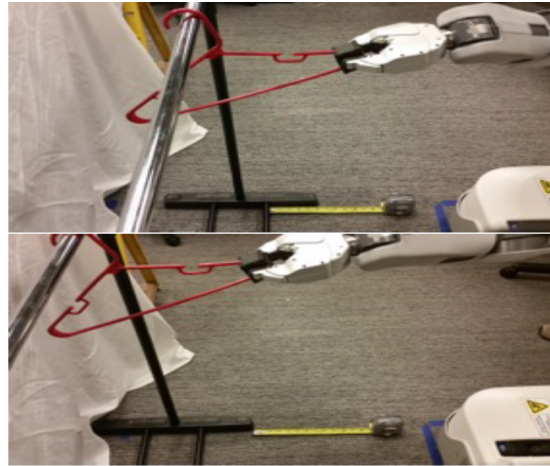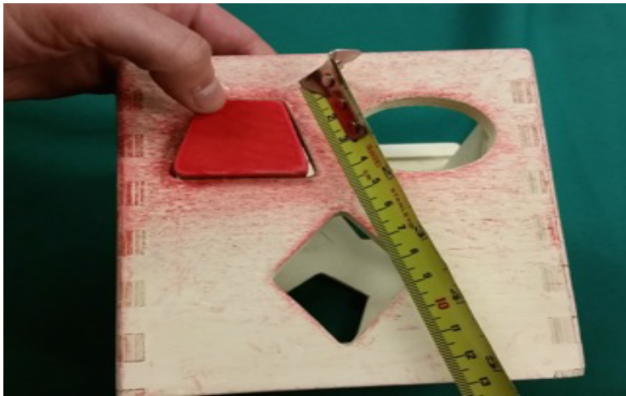
# $\pi_\theta$  Deep Spatial Neural Net Architecture



(92,000 parameters)

[Levine*, Finn*, Darrell, Abbeel, JMLR 2016]

# Experimental Tasks



[Levine*, Finn*, Darrell, Abbeel, JMLR 2016]

Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# Learning



[Levine*, Finn*, Darrell, Abbeel, JMLR 2016]

# Visuomotor Learning Directly in Visual Space



[Finn, Tan, Duan, Darrell, Levine,  Abbeel, ICRA 2016]

Related work: Embed to Control [Wattenberg,
Springenberg, Boedecker, Riedmiller, 2015]

Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

# Visuomotor Learning Directly in Visual Space



[Finn, Tan, Duan, Darrell, Levine, Abbeel, 2015]

# Visuomotor Cost Function Learning

- Learning from goal image can be great

- But:

  - Often other objects in environment --- don't actually expect to perfectly match example goal image

  - Goal image might not reveal much about how to get there

# Visuomotor Cost Function Learning

→ infer cost function from demonstrations

Challenges:
- underdefined problem
- difficult to evaluate learned cost
- large perceptual input spaces

# Prior Approaches | Desiderata

### repeatedly solve MDP



Abbeel & Ng '04
Ziebart et al. '08
Ratliff et al. '09

avoid (repeatedly) solving the MDP

### use known dynamics



Todorov '06
Levine et al. '12
Dragan et al. '12

handle unknown dynamics

### use hand-designed features



Boularias et al. '11
Kalakrishnan et al. '13
Doerr et al. '15

*learn features with flexible, nonlinear cost parametrization*

+ sample efficiency

# Guided Cost Learning



initial distribution $q_0$

human demonstrations

generate policy samples from q

Update cost using samples & demos

update q w.r.t. cost

**policy** q

**cost** c

# Guided Cost Learning



initial distribution $q_0$

human demonstrations

generate policy samples from q

**generator**

Update cost using samples & demos

**discriminator**

**update** q w.r.t. cost
**policy** q (partially optimize)

**cost** c

update cost in inner loop of policy optimization

# Guided Cost Learning



initial distribution $q_0$

human demonstrations

generate policy samples from q

generator

Update cost using samples & demos

discriminator

**update** q w.r.t. cost
**policy** q (partially optimize)

**cost** c

**Ho et al.,** ICML '16, arXiv '16
**Kim & Bengio,** arXiv '16
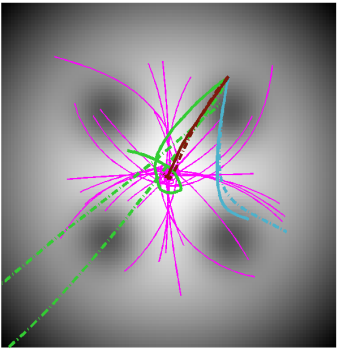
# Experiments

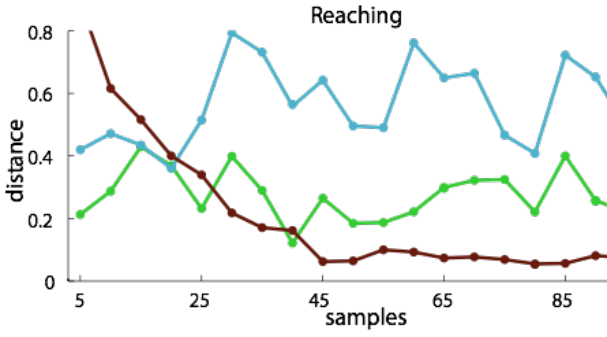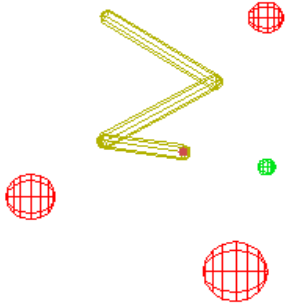## 2D navigation



## 2D reaching



## peg insertion



high-dimensional continuous states & actions

direct torque control

complex contact dynamics

# Experiments
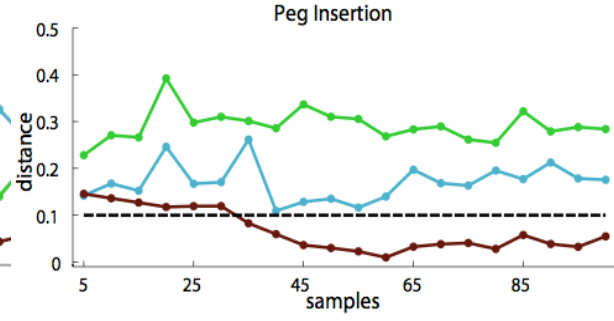
## 2D navigation



## 2D reaching



## peg insertion
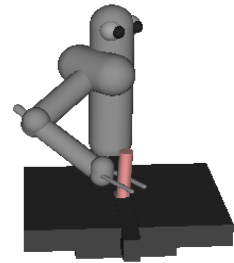




2D Navigation



Reaching
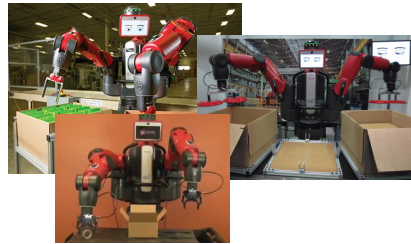


Peg Insertion
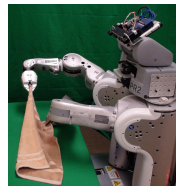
# Frontiers

- Shared and transfer learning



- Memory

  - Estimation

  - Temporal hierarchy / goal setting

- Applications



Pieter Abbeel -- UC Berkeley / OpenAI / Gradescope

Thank you