

CS 188: Artificial Intelligence

Robotics *



Instructors: Pieter Abbeel & Anca Dragan --- University of California, Berkeley

These slides were created by Dan Klein, Pieter Abbeel and Anca Dragan for CS188 Intro to AI at UC Berkeley.

All CS188 materials are available at <http://ai.berkeley.edu>

PR1



[Wyrobek, Berger, Van der Loos, Salisbury, ICRA 2008]

Personal Robotics Hardware



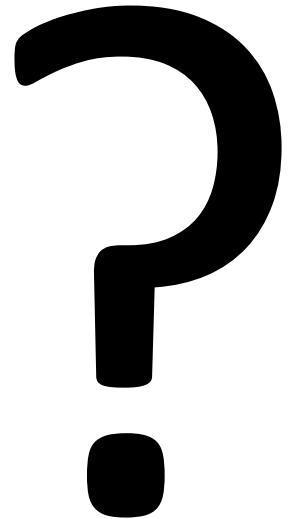
PR2
Willow Garage
\$400,000
2009



Baxter
Rethink Robotics
\$30,000
2013



Fetch
Fetch Robotics
~\$80,000
2015



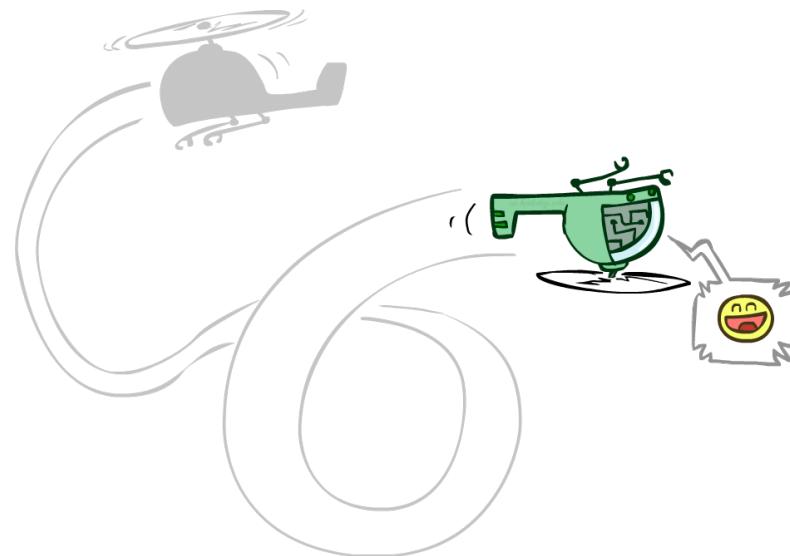
Robot Learning

- Apprenticeship learning
- Reinforcement learning

Robot Learning

- *Apprenticeship learning*
- Reinforcement learning

Robotic Helicopters



Motivating Example



- How do we execute a task like this?

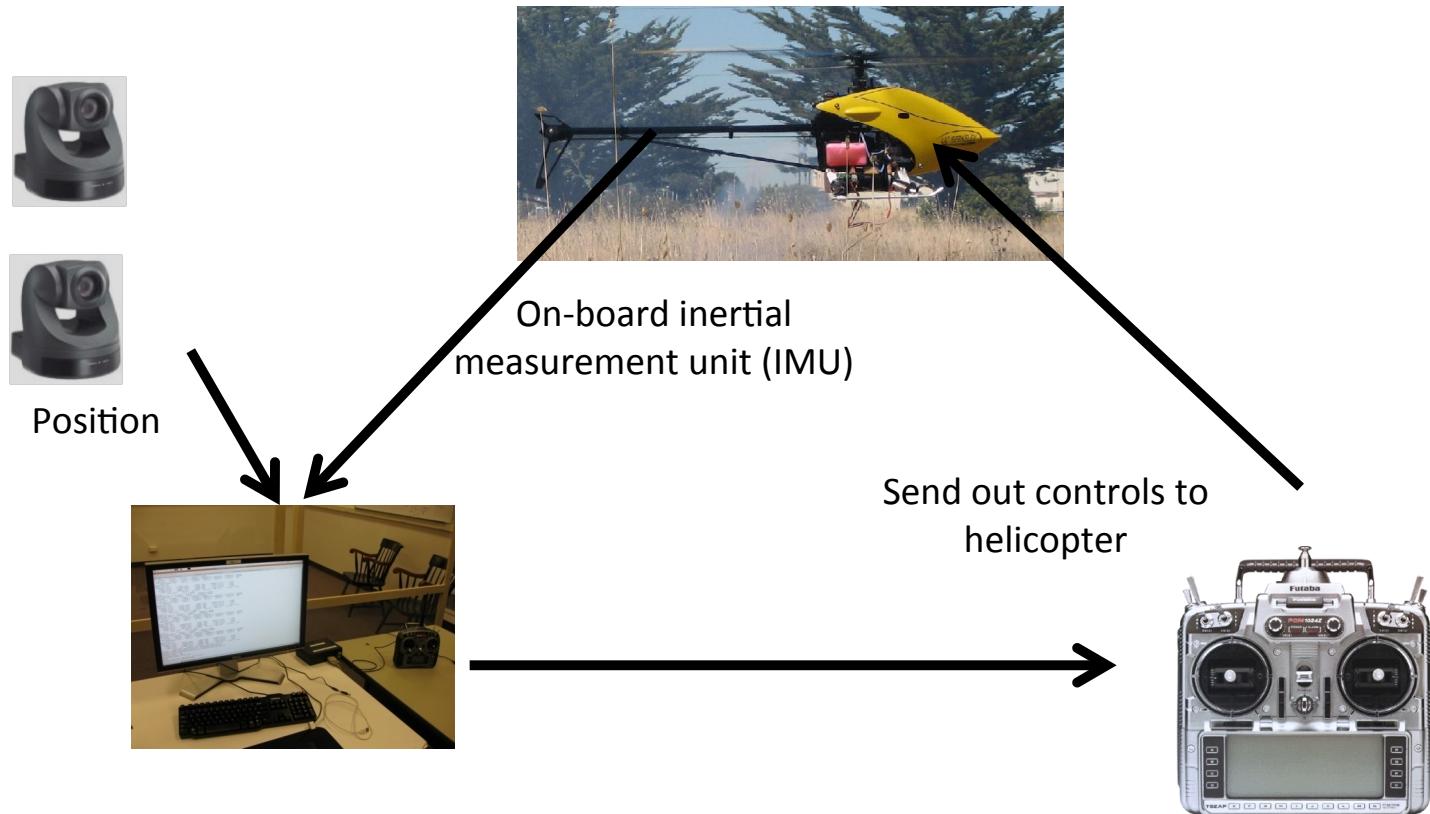
[VIDEO: tictoc_results.wmv]

Autonomous Helicopter Flight

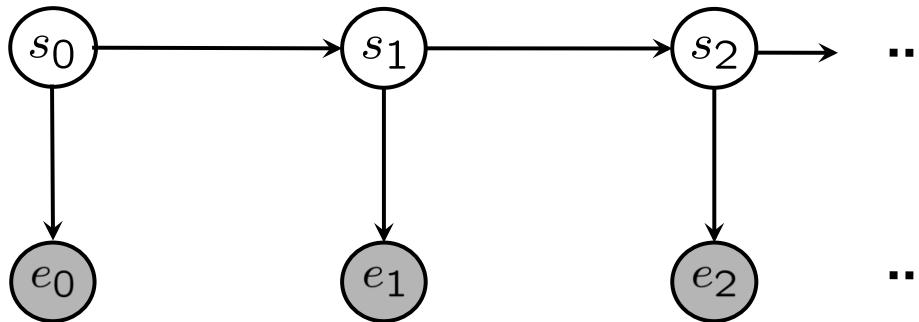


- Key challenges:
 - Track helicopter position and orientation during flight
 - Decide on control inputs to send to helicopter

Autonomous Helicopter Setup



HMM for Tracking the Helicopter



- State: $s = (x, y, z, \phi, \theta, \psi, \dot{x}, \dot{y}, \dot{z}, \dot{\phi}, \dot{\theta}, \dot{\psi})$
- Measurements: [observation update]
 - 3-D coordinates from vision, 3-axis magnetometer, 3-axis gyro, 3-axis accelerometer
- Transitions (dynamics): [time elapse update]
 - $s_{t+1} = f(s_t, a_t) + w_t$ f : encodes helicopter dynamics, w : noise

Helicopter MDP

- State: $s = (x, y, z, \phi, \theta, \psi, \dot{x}, \dot{y}, \dot{z}, \dot{\phi}, \dot{\theta}, \dot{\psi})$
- Actions (control inputs):
 - a_{lon} : Main rotor longitudinal cyclic pitch control (affects pitch rate)
 - a_{lat} : Main rotor latitudinal cyclic pitch control (affects roll rate)
 - a_{coll} : Main rotor collective pitch (affects main rotor thrust)
 - a_{rud} : Tail rotor collective pitch (affects tail rotor thrust)
- Transitions (dynamics):
 - $s_{t+1} = f(s_t, a_t) + w_t$
[f encodes helicopter dynamics]
[w is a probabilistic noise model]
- Can we solve the MDP yet?



Problem: What's the Reward?

- Reward for hovering:

$$\begin{aligned} R(s) = & -\alpha_x(x - x^*)^2 \\ & -\alpha_y(y - y^*)^2 \\ & -\alpha_z(z - z^*)^2 \\ & -\alpha_{\dot{x}}\dot{x}^2 \\ & -\alpha_{\dot{y}}\dot{y}^2 \\ & -\alpha_{\dot{z}}\dot{z}^2 \end{aligned}$$

[VIDEO: ihover-newencoding09.wmv]

Hover



[Ng et al, 2004]

Problem: What's the Reward?

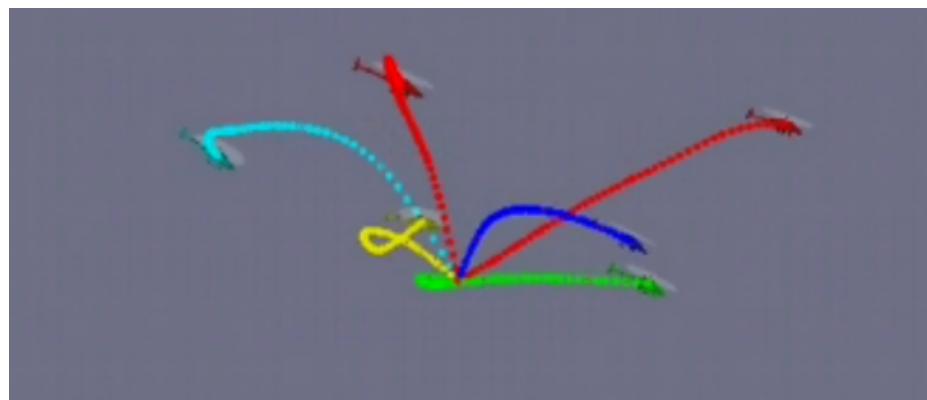
- Rewards for “Flip”?
 - Problem: what's the target trajectory?
 - Just write it down by hand?

[VIDEO: 20061204---bad.wmv]

Flips (?)



Helicopter Apprenticeship?

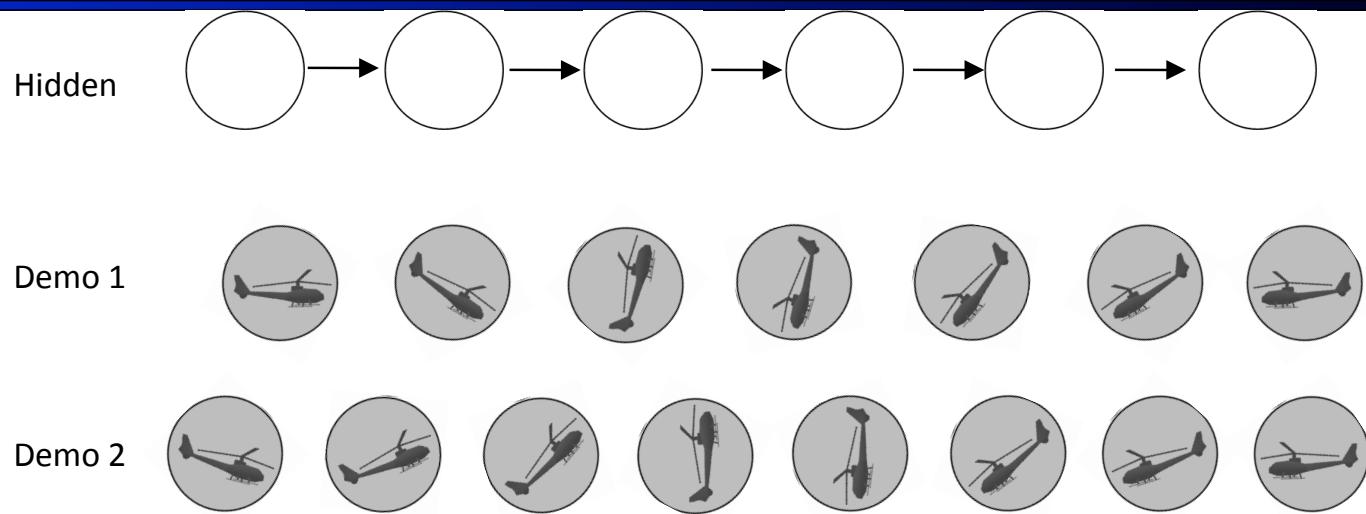


[VIDEO: airshow_unaligned.wmv]

Demonstrations



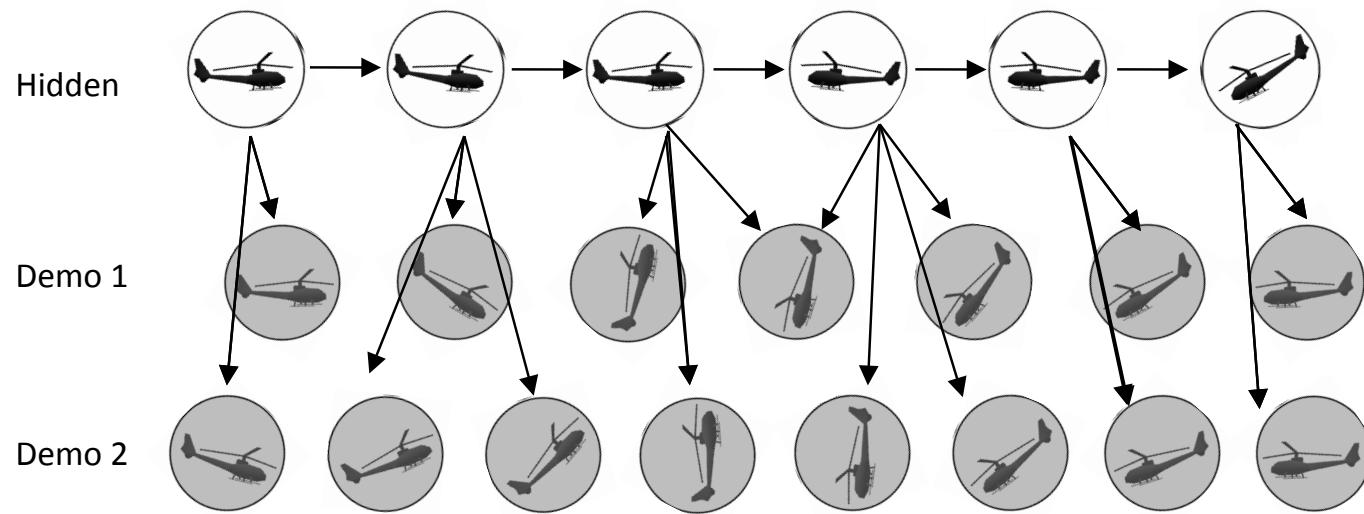
Learning a Trajectory



- **HMM-like generative model**
 - Dynamics model used as HMM transition model
 - Demos are observations of hidden trajectory
- **Problem: how do we align observations to hidden trajectory?**

Abbeel, Coates, Ng, IJRR 2010

Probabilistic Alignment using a Bayes' Net

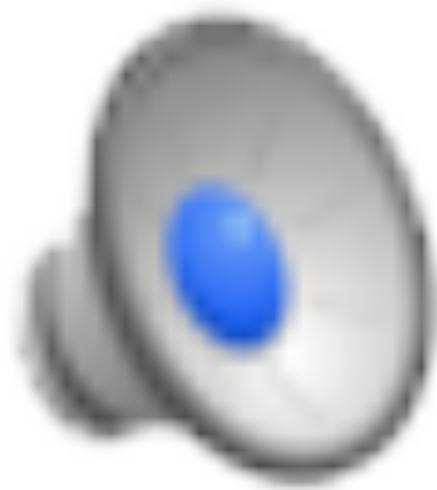


- 
 - Dynamic Time Warping
(Needleman&Wunsch 1970, Sakoe&Chiba, 1978)
 - Extended Kalman filter / smoother

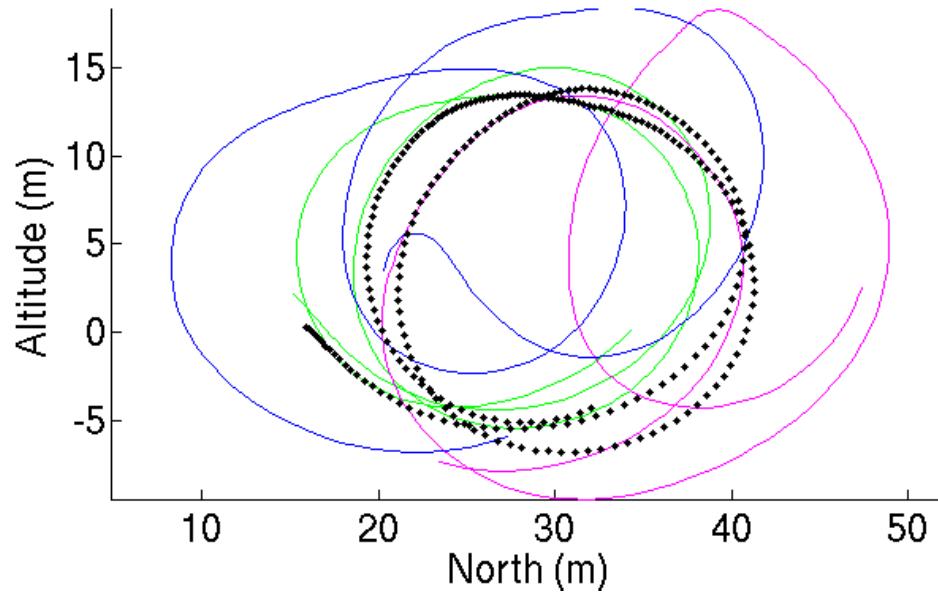
Abbeel, Coates, Ng, IJRR 2010

[VIDEO: airshow_unaligned.wmv]

Aligned Demonstrations



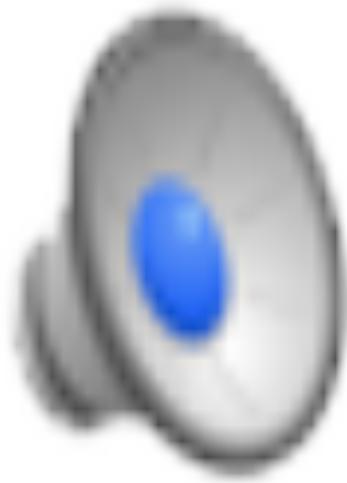
Alignment of Samples



- Result: inferred sequence is much cleaner!

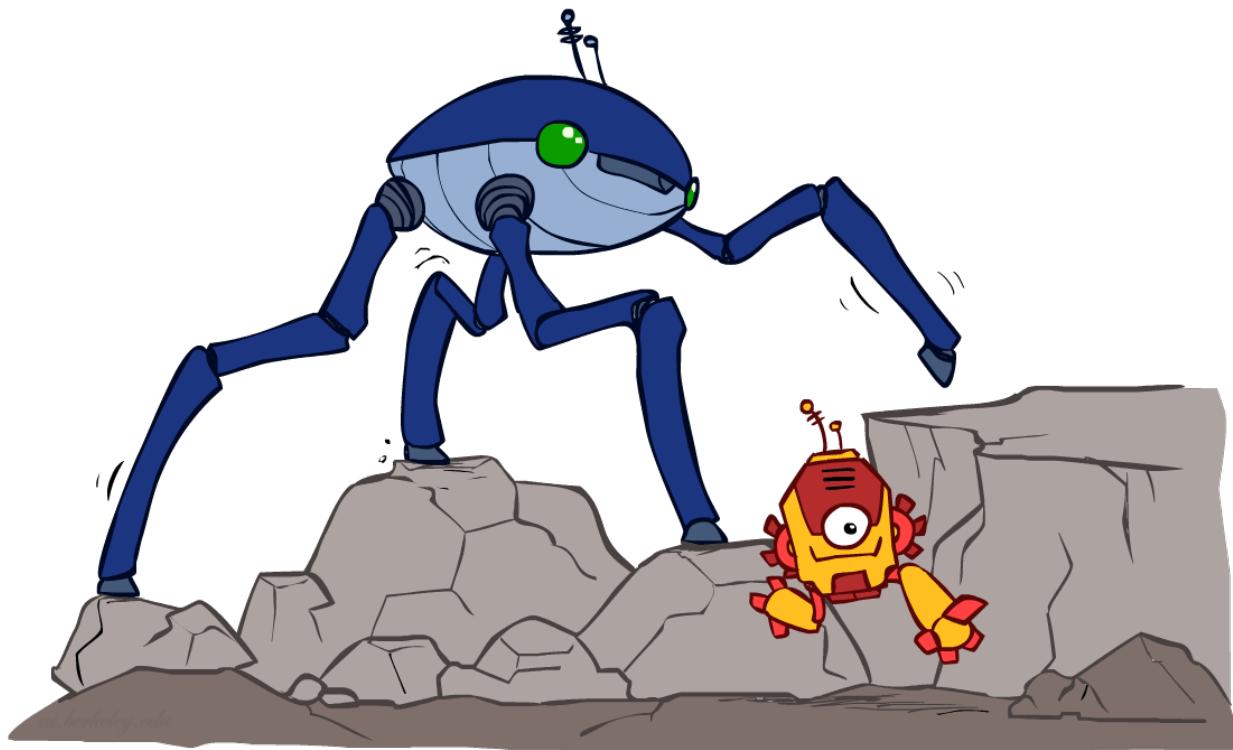
[VIDEO: airshow_trimmed.wmv]

Final Behavior



[Abbeel, Coates, Quigley, Ng, 2010]

Legged Locomotion



Quadruped



- Low-level control problem: moving a foot into a new location
→ search with successor function ~ moving the motors
- High-level control problem: where should we place the feet?
 - Reward function $R(x) = w \cdot f(s)$ [25 features]

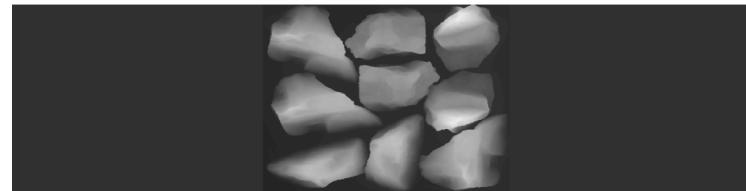
[Kolter, Abbeel & Ng, 2008]

Experimental setup

- Demonstrate path across the “training terrain”



- Run apprenticeship to learn the reward function
- Receive “testing terrain”---height map.



- Find the optimal policy with respect to the *learned reward function* for crossing the testing terrain.

[Kolter, Abbeel & Ng, 2008]

[VIDEO: quad initial.wmv]

Without learning



[VIDEO: quad initial.wmv]

With learned reward function

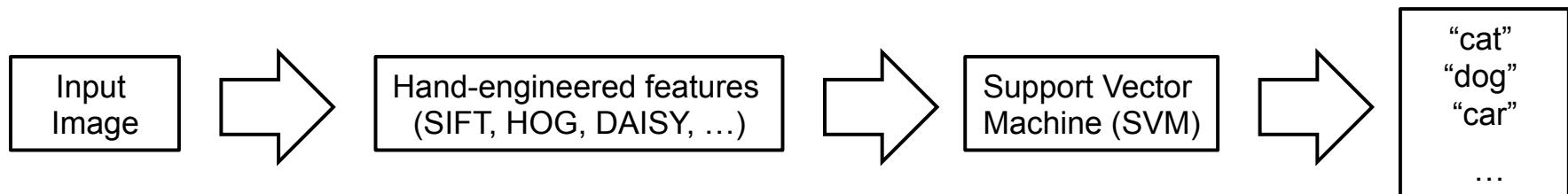


Robot Learning

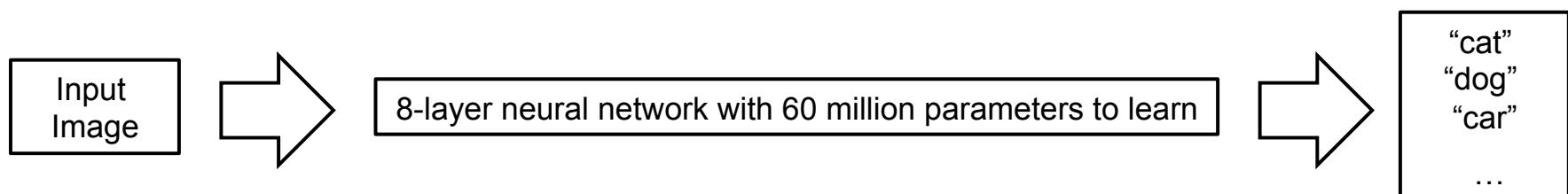
- Apprenticeship learning
- *Reinforcement learning*

Object Detection in Computer Vision

- State-of-the-art object detection until 2012:

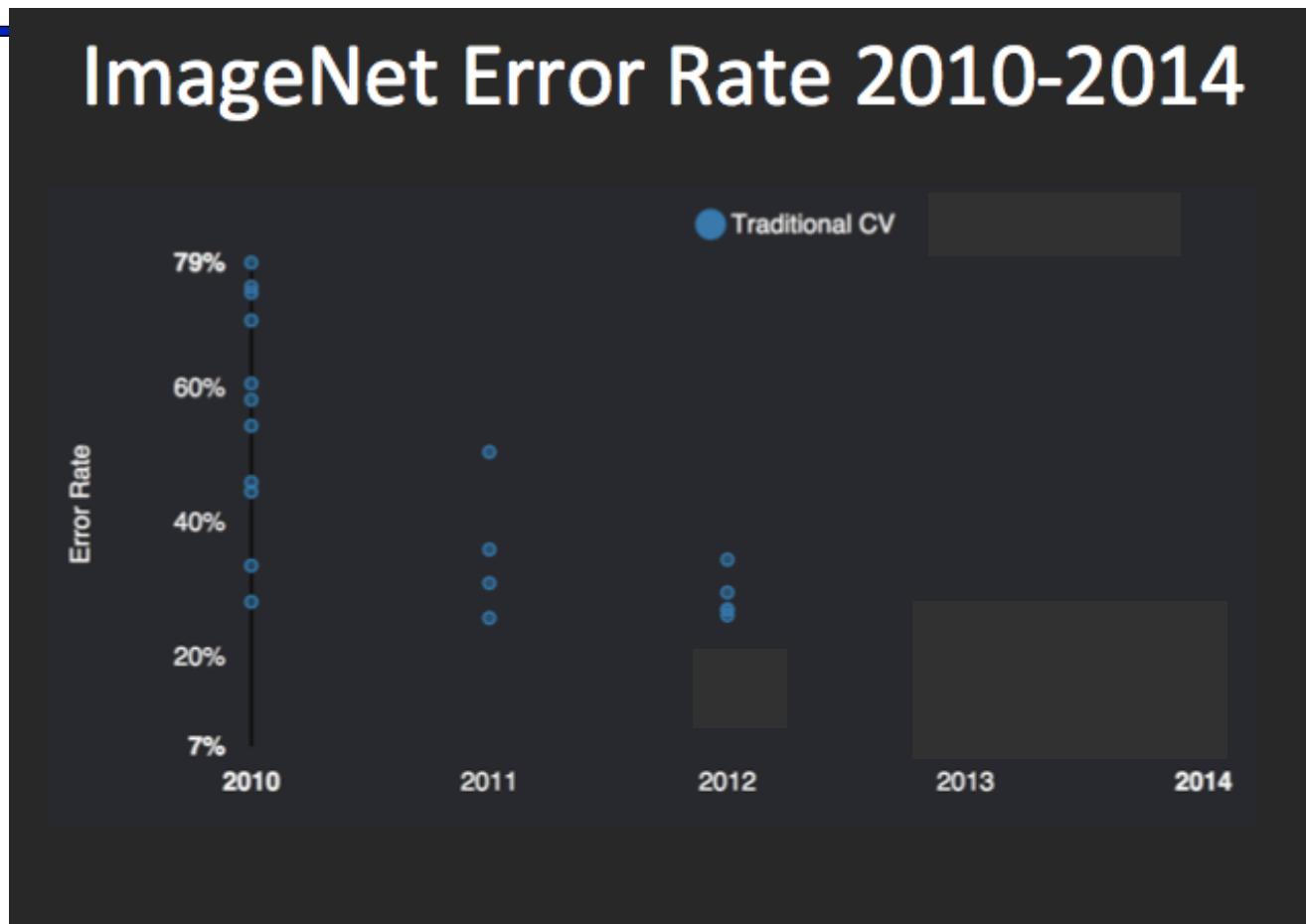


- Deep Supervised Learning (Krizhevsky, Sutskever, Hinton 2012; also LeCun, Bengio, Ng, Darrell, ...):

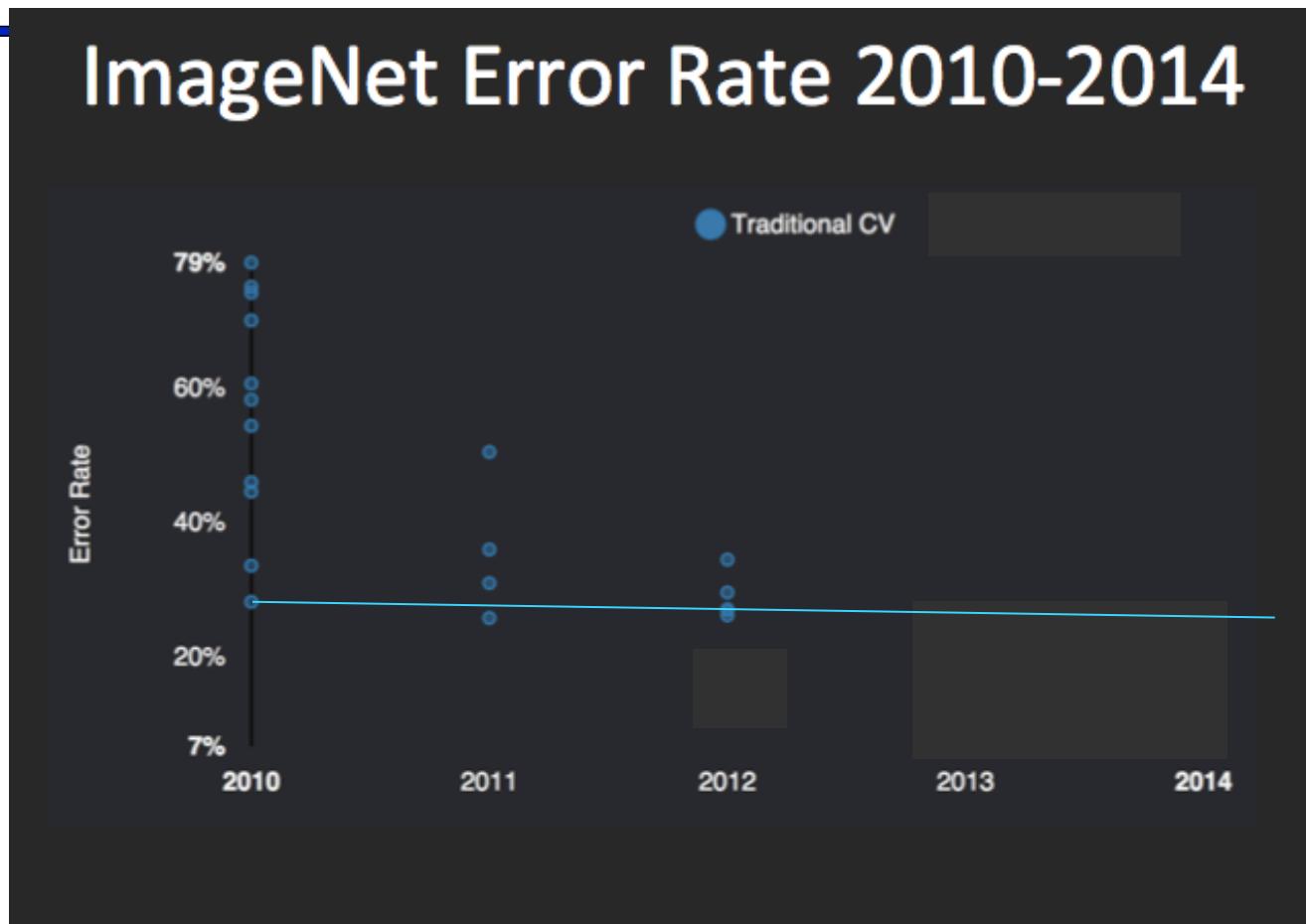


- ~1.2 million training images from ImageNet [Deng, Dong, Socher, Li, Li, Fei-Fei, 2009]

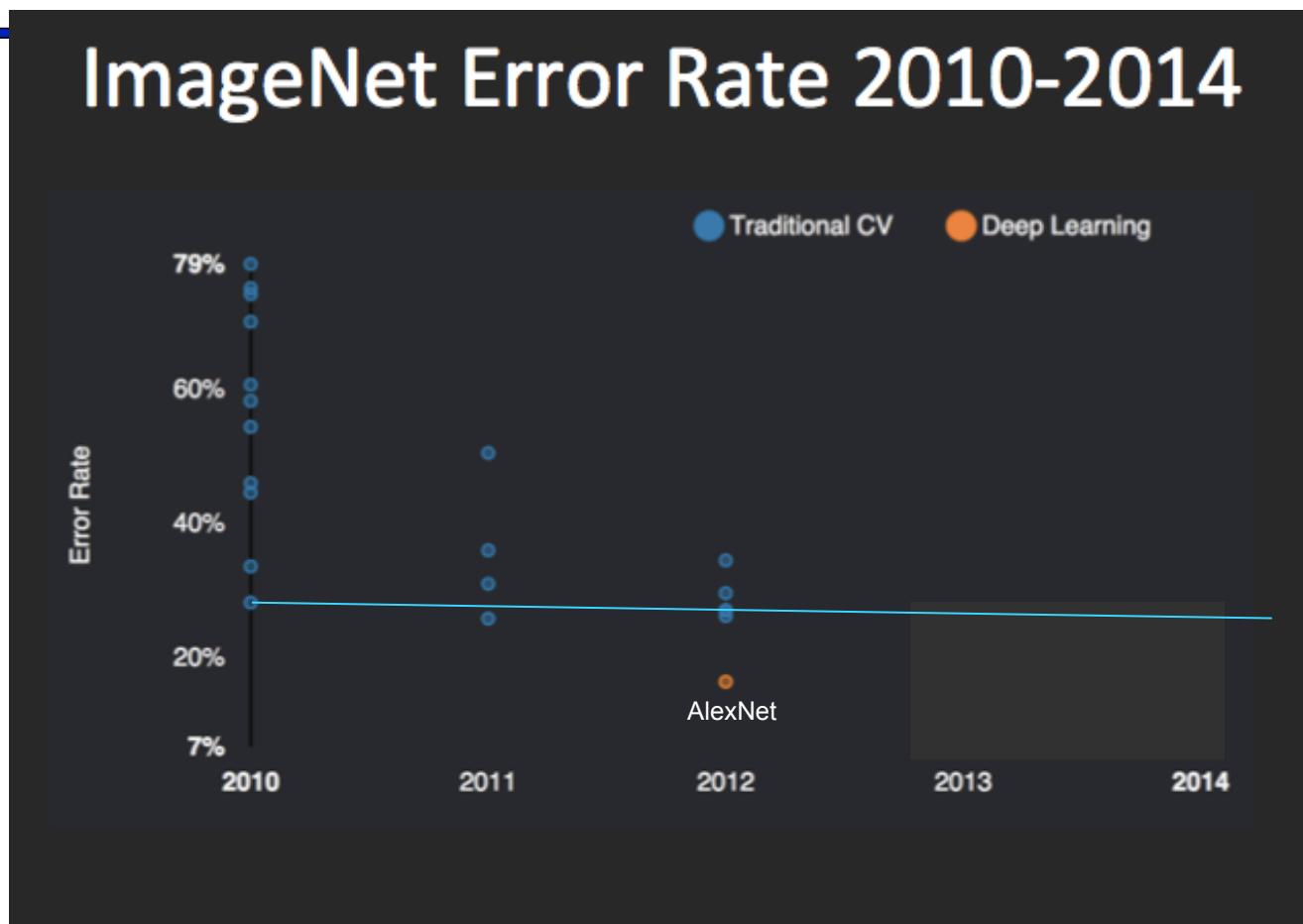
Performance



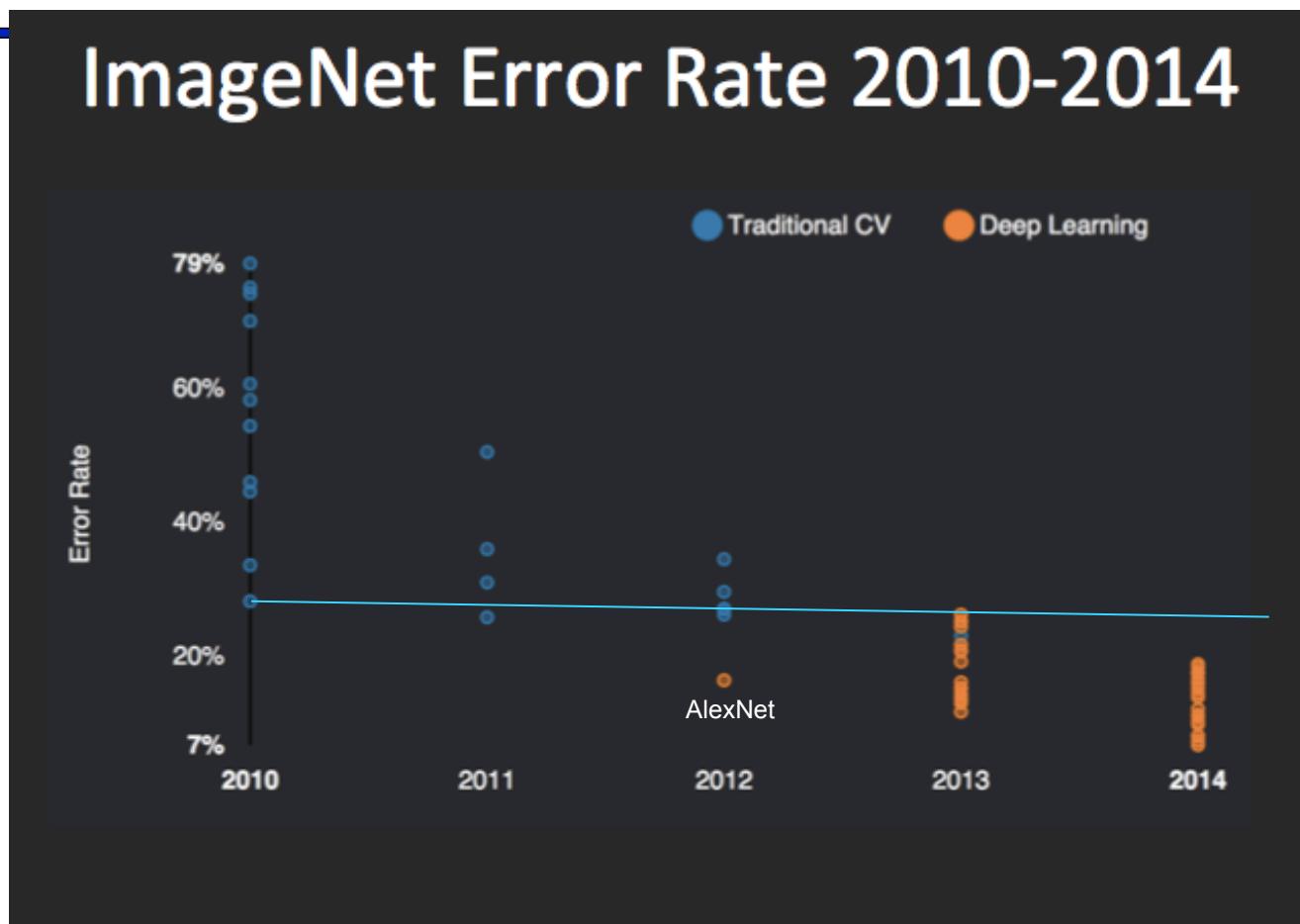
Performance



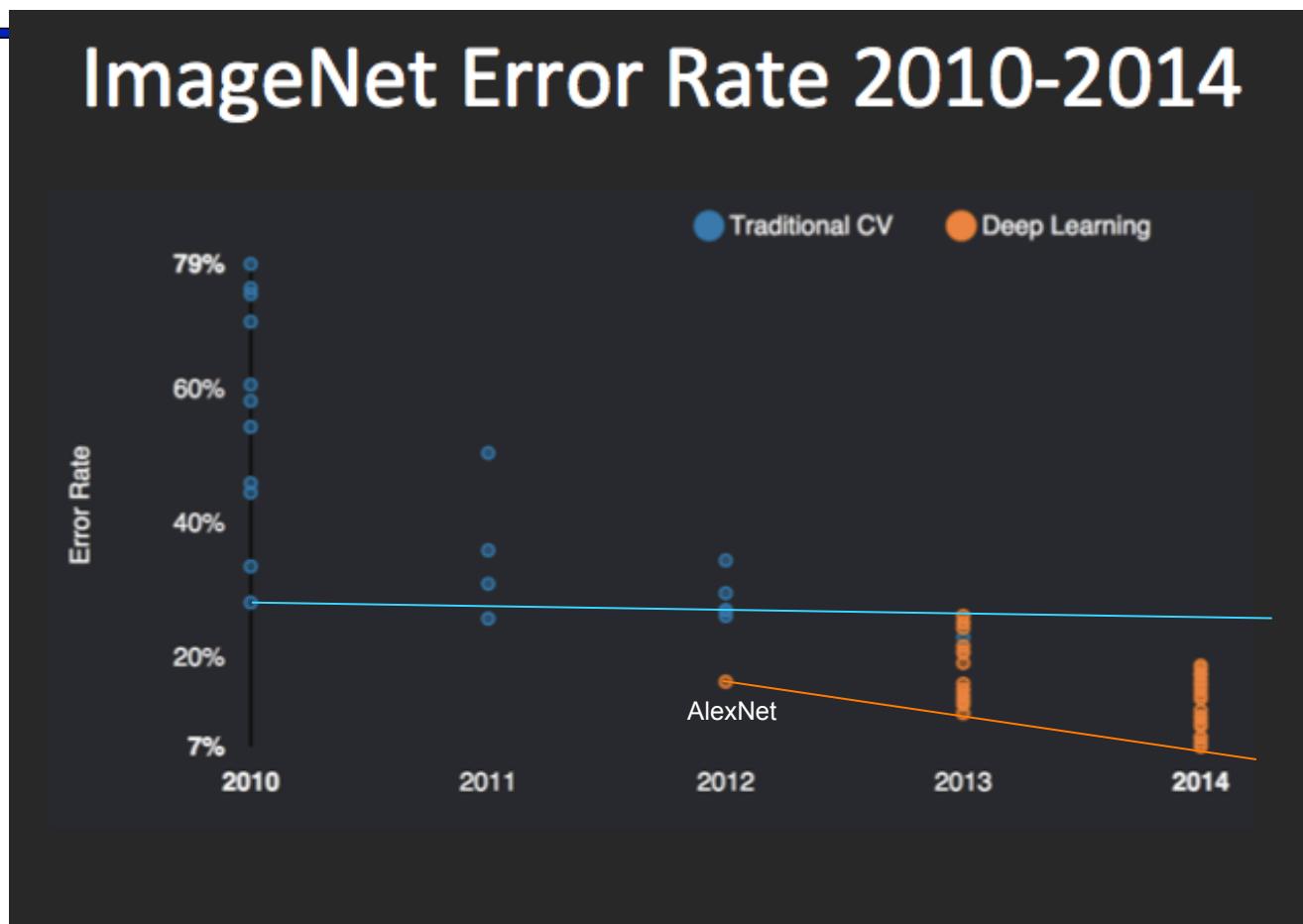
Performance



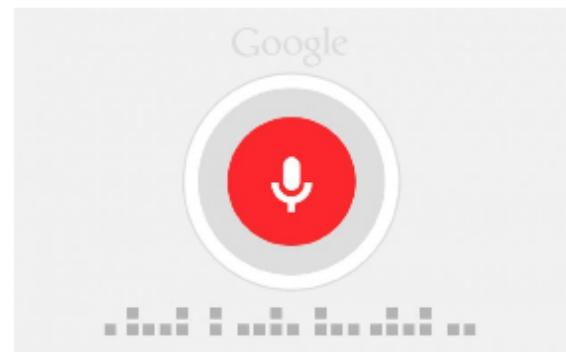
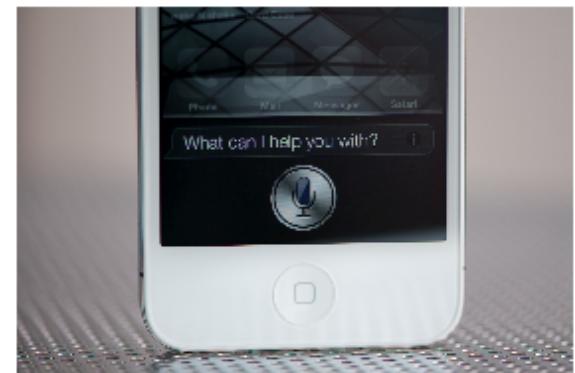
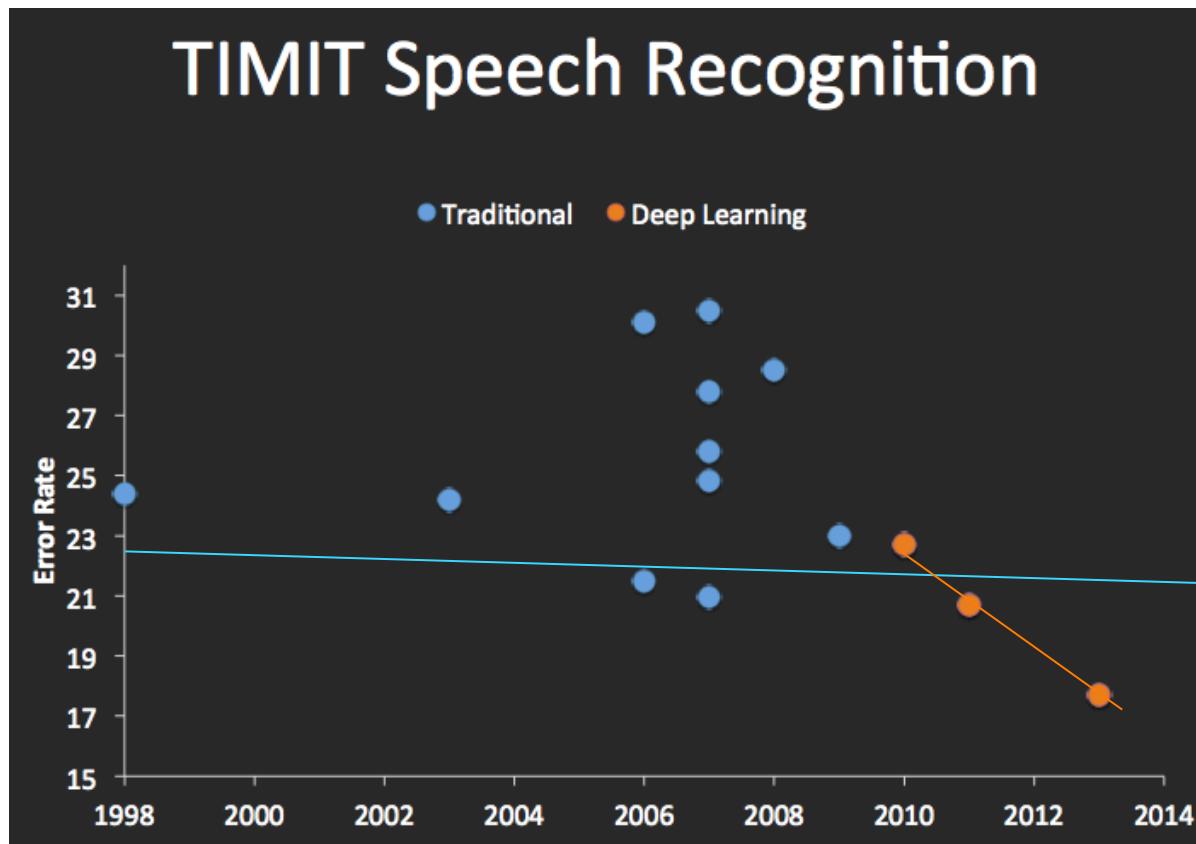
Performance



Performance



Speech Recognition



graph credit Matt Zeiler, Clarifai

History

- 1958 Rosenblatt proposed perceptrons
- 1980 Neocognitron (Fukushima, 1980)
- 1982 Hopfield network, SOM (Kohonen, 1982), Neural PCA (Oja, 1982)
- 1985 Boltzmann machines (Ackley et al., 1985)
- 1986 Multilayer perceptrons and backpropagation (Rumelhart et al., 1986)
- 1988 RBF networks (Broomhead&Lowe, 1988)
- 1989 Autoencoders (Baldi&Hornik, 1989), Convolutional network (LeCun, 1989)
- 1992 Sigmoid belief network (Neal, 1992)
- 1993 Sparse coding (Field, 1993) (Olshausen, 1996)
- 2000s Sparse, Probabilistic, and Energy models (Hinton, Bengio, LeCun, Ng)

Is deep learning 3, 30, or 60 years old?

based on history by K. Cho



Rosenblatt's Perceptron

What's Changed

- Data

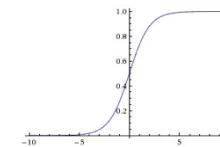
- 1.2M training examples
- 2048 (different crops)
- 90 (PCA re-colorings)

- Compute power

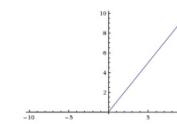
- Two NVIDIA GTX 580 GPUs
- 5-6 days of training time

- Nonlinearity

Sigmoid



→ ReLU



- Regularization

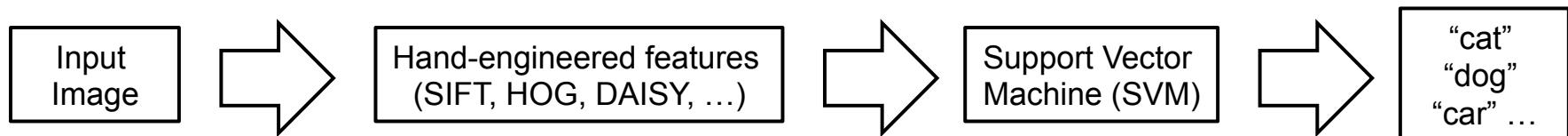
- Drop-out

- Exploration of model structure

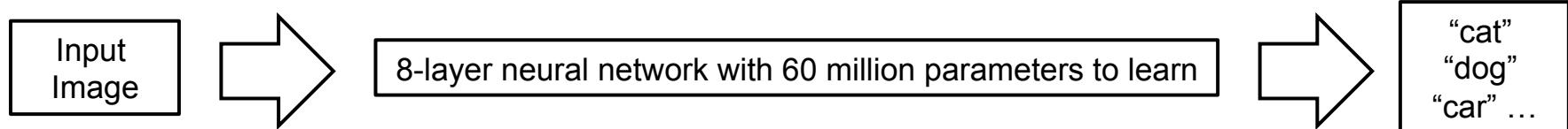
- Optimization know-how

Object Detection in Computer Vision

- State-of-the-art object detection until 2012:



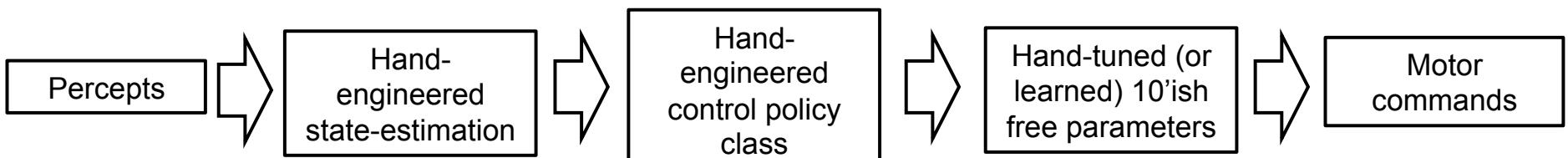
- Deep Supervised Learning (Krizhevsky, Sutskever, Hinton 2012; also LeCun, Bengio, Ng, Darrell, ...):



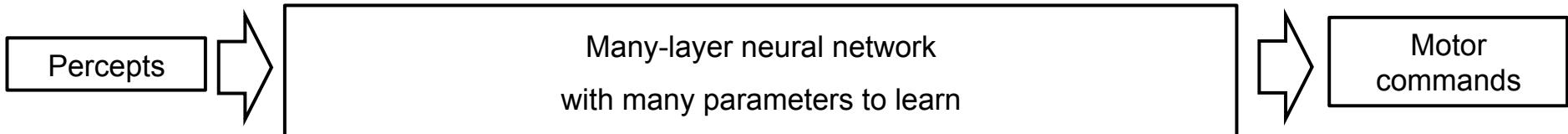
- ~1.2 million training images from ImageNet [Deng, Dong, Socher, Li, Li, Fei-Fei, 2009]

Robotics

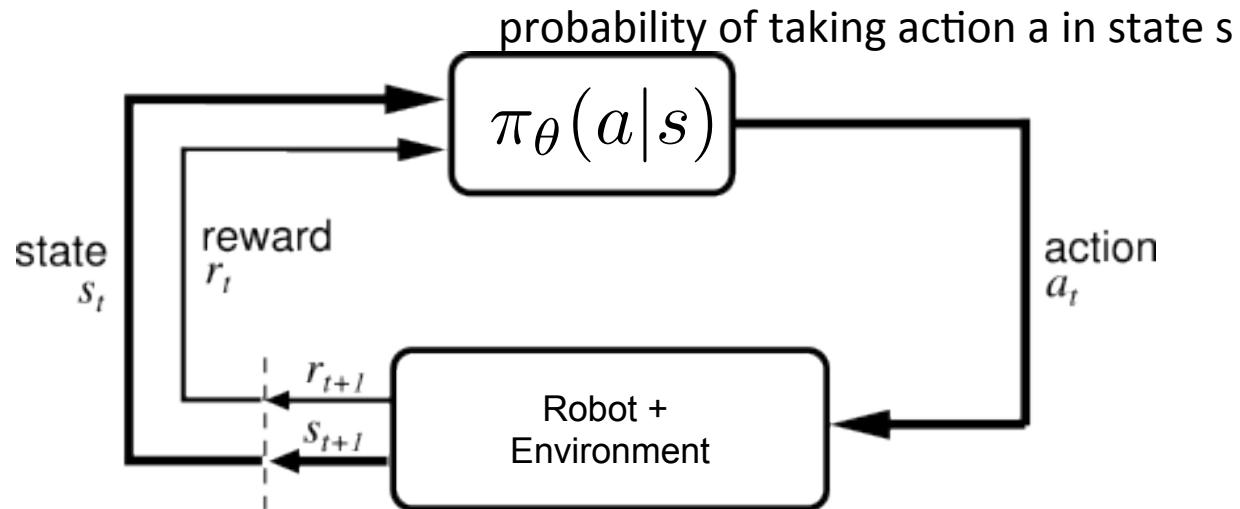
- **Current state-of-the-art robotics**



- **Deep reinforcement learning**



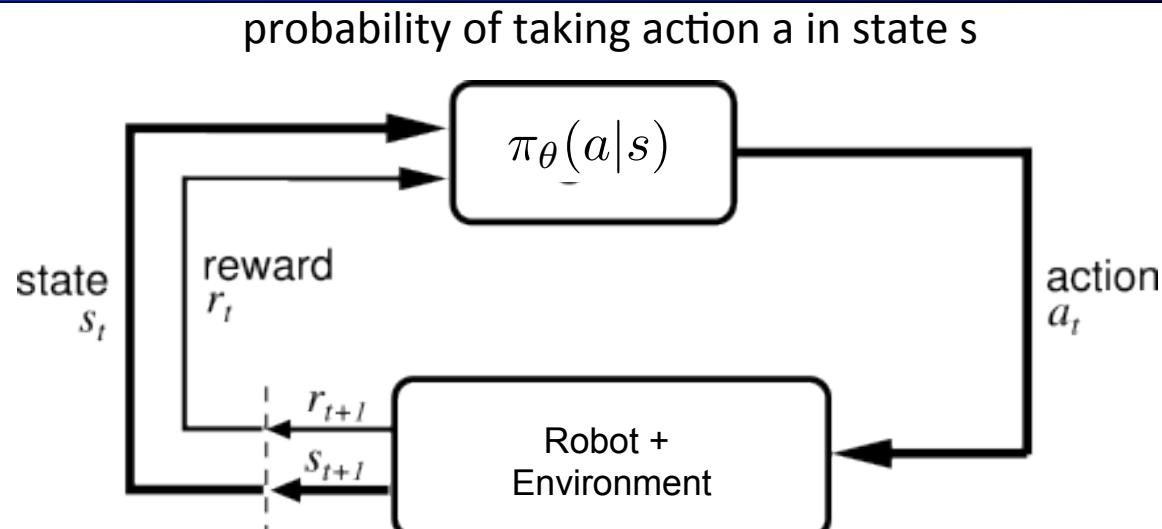
Reinforcement Learning (RL)



- Robotics
- Marketing / Advertising
- Dialogue
- Optimizing operations / logistics
- Queue management
- ...

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

Reinforcement Learning



- Goal:

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

- Additional challenges:
 - Stability
 - Credit assignment
 - Exploration

From Pixels to Actions?



Pong



Enduro

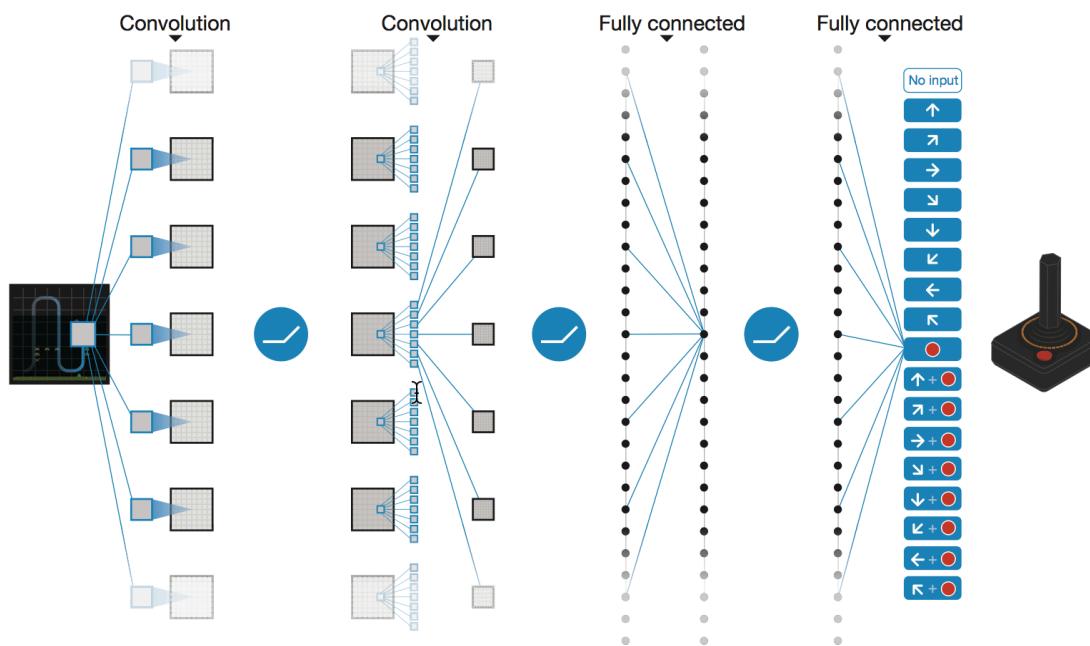


Beamrider



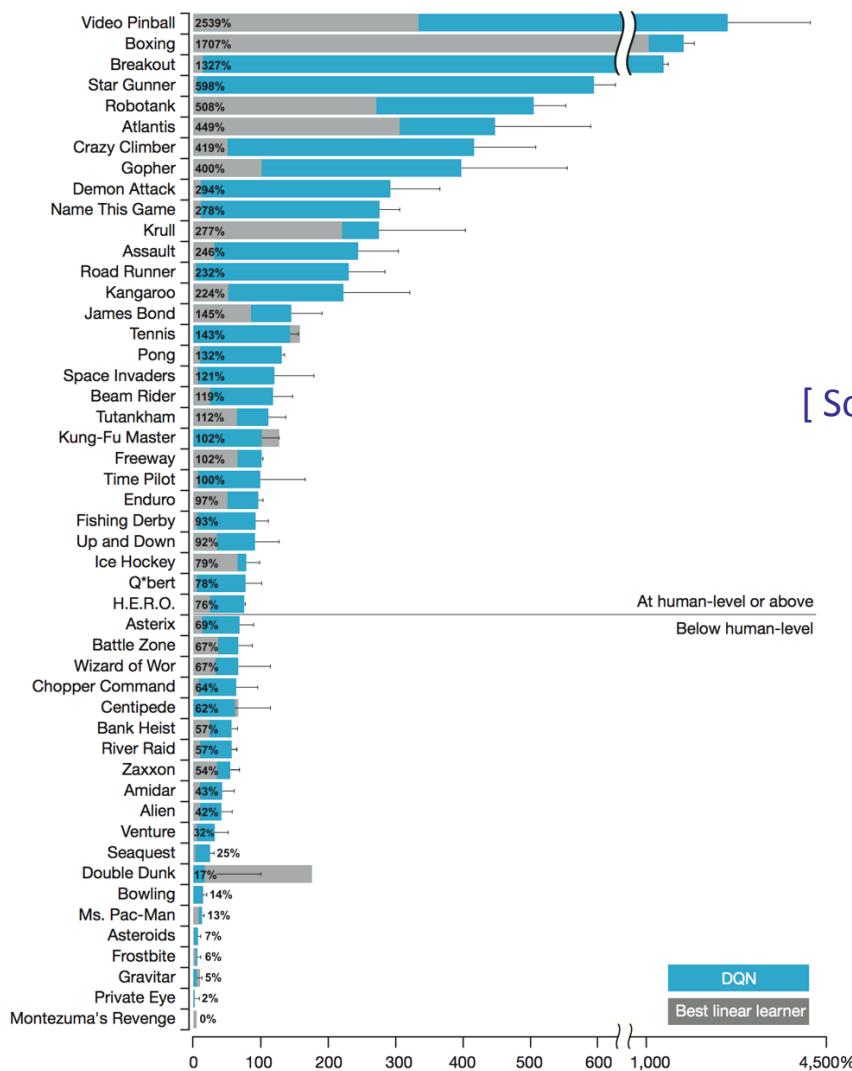
Q*bert

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)]



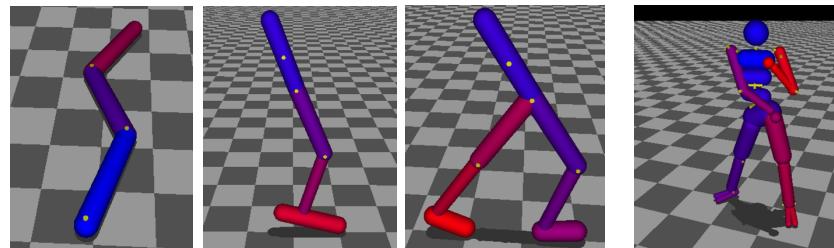
[Source: Mnih et al., Nature 2015 (DeepMind)]

Deep Q-Network (DQN)

- Approach:
 - Q-learning with e-greedy and deep network as function approximator
- Key idea 1: stabilizing Q-learning
 - Mini-batches of size 32 (vs. single sample updates)
 - Q-values used to compute temporal difference only updated every 10,000 updates
- Key idea 2: lots of data / compute
 - trained for a total of 50 million frames (*=38 days of game experience*) and use a replay memory of one million most recent frames

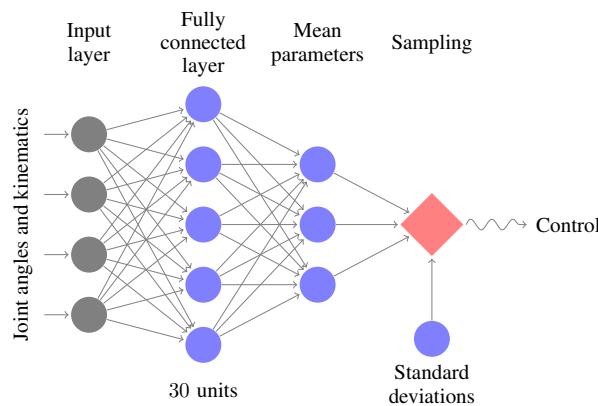
How About Continuous Control, e.g., Locomotion?

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

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

- 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

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

$$\max_{\delta\theta} \quad \hat{g}^\top \delta\theta$$

$$\text{s.t. } KL(P(\tau; \theta) || P(\tau; \theta + \delta\theta)) \leq \varepsilon$$

- Trust Region:
 - Sampled evaluation of gradient
 - Gradient only locally a good approximation
 - Change in policy changes state-action visitation frequencies

Stability

[Schulman, Levine, Moritz, Jordan, Abbeel, ICML 2015]

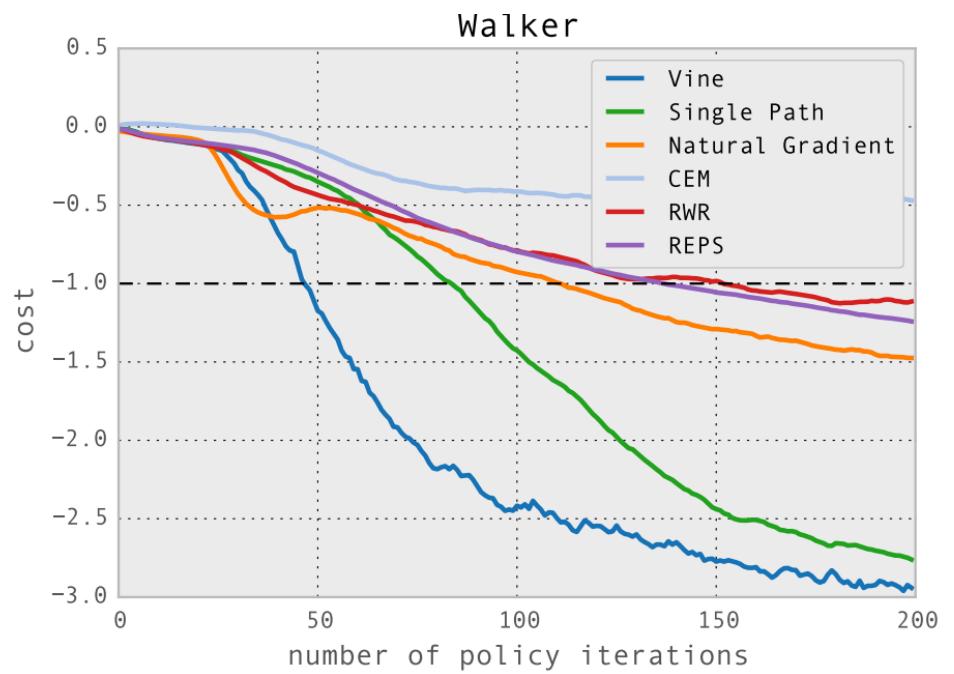
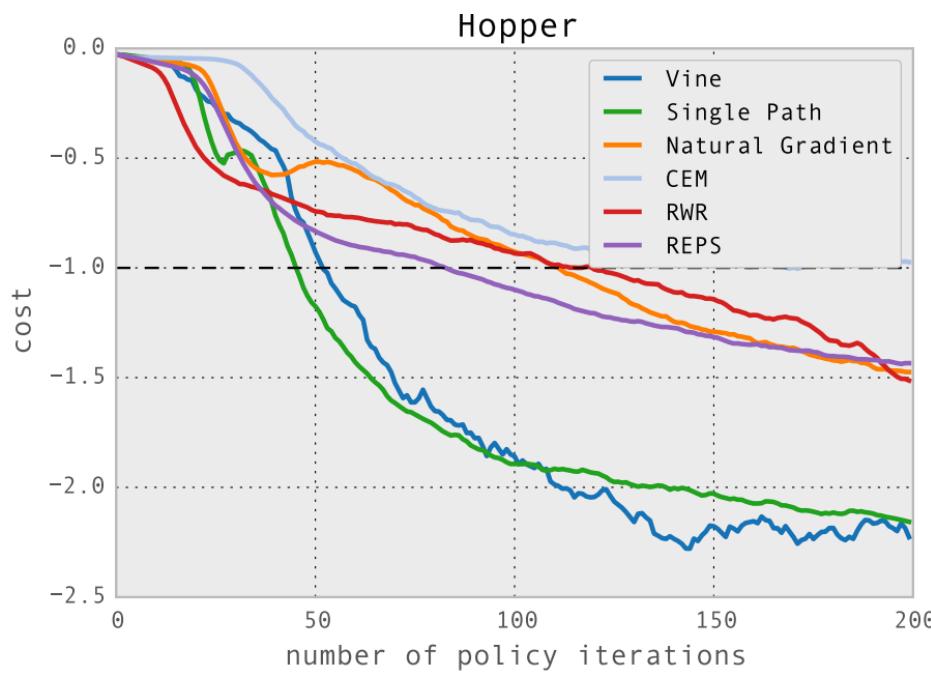
Experiments in Locomotion

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

The following gaits were obtained

[Schulman, Levine, A.]

Learning Curves -- Comparison

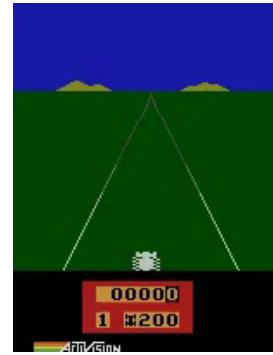


Atari Games

- 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]



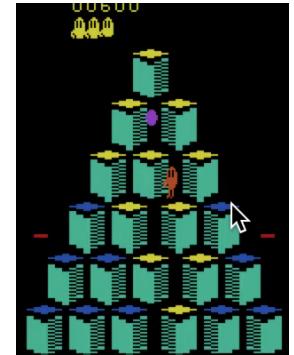
Pong



Enduro



Beamrider



Q*bert

Revisiting the Gradient Estimate

Trust Region
Policy
Optimization
(TRPO)

$$\max_{\delta\theta} \hat{g}^\top \delta\theta$$

$$\text{s.t. } KL(P(\tau; \theta) || P(\tau; \theta + \delta\theta)) \leq \varepsilon$$

$$\hat{g} = \sum_i \sum_{t=0}^H \nabla_\theta \log \pi_\theta(a_t | s_t) \left(\sum_{k=t}^H R(s_k) - \hat{V}_\phi(s_t) \right)$$

Generalized
Advantage
Estimation (GAE)

Key ideas:
Credit Assignment

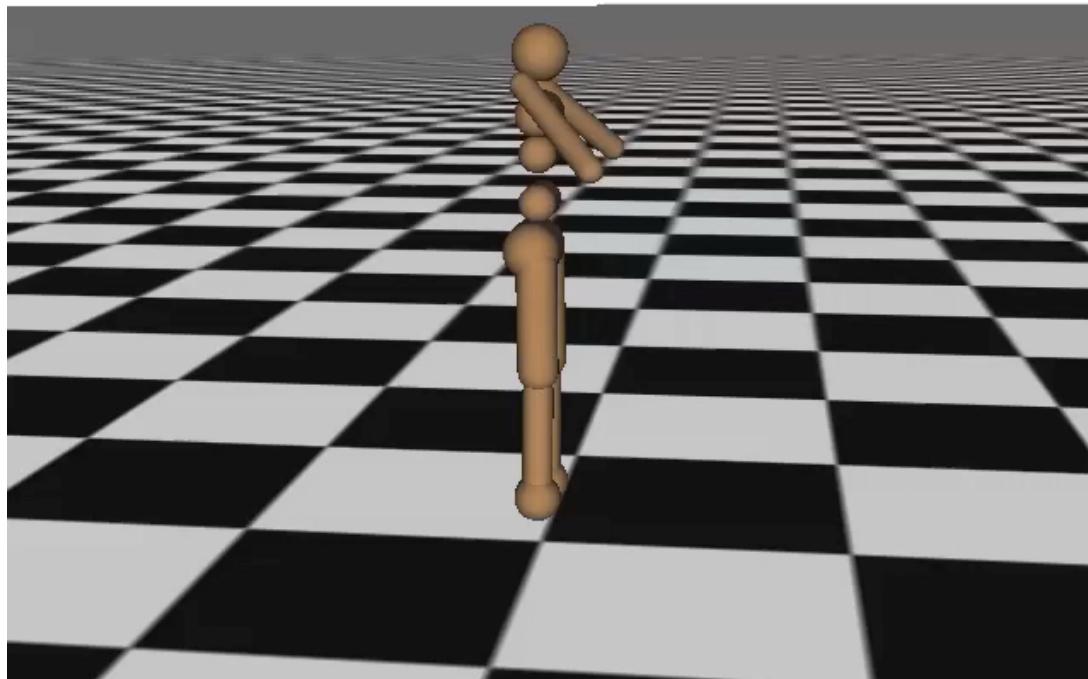
$$\hat{Q}_\phi(s_t)$$

- Deep neural net
- Trust-region optimization

[Schulman, Moritz, Levine, Jordan, Abbeel, 2015 (under review ICLR)]

Learning Locomotion

Iteration 0



[Schulman, Moritz, Levine, Jordan, Abbeel, 2015]

In Contrast: Darpa Robotics Challenge

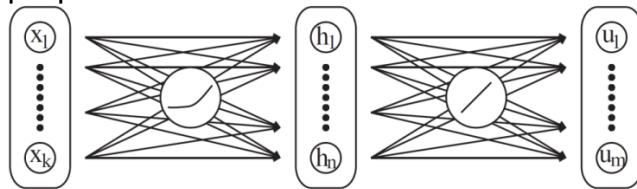


How About Real Robotic Visuo-Motor Skills?



Guided Policy Search

general-purpose neural network controller



policy search (RL)

complex dynamics

complex policy

HARD

supervised learning

complex dynamics

complex policy

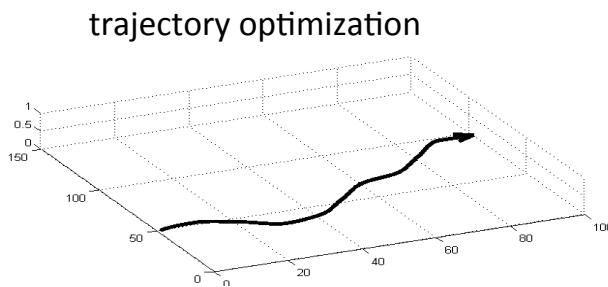
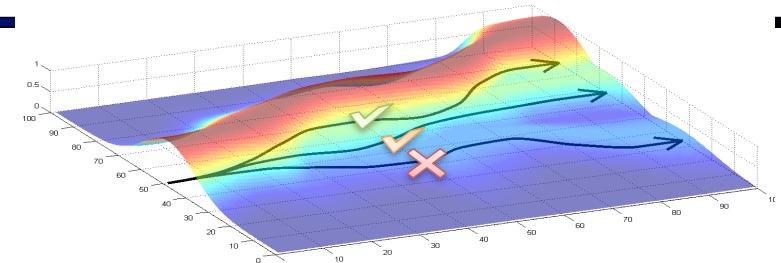
EASY

trajectory optimization

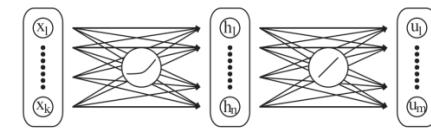
complex dynamics

complex policy

EASY



Supervised learning



Guided Policy Search



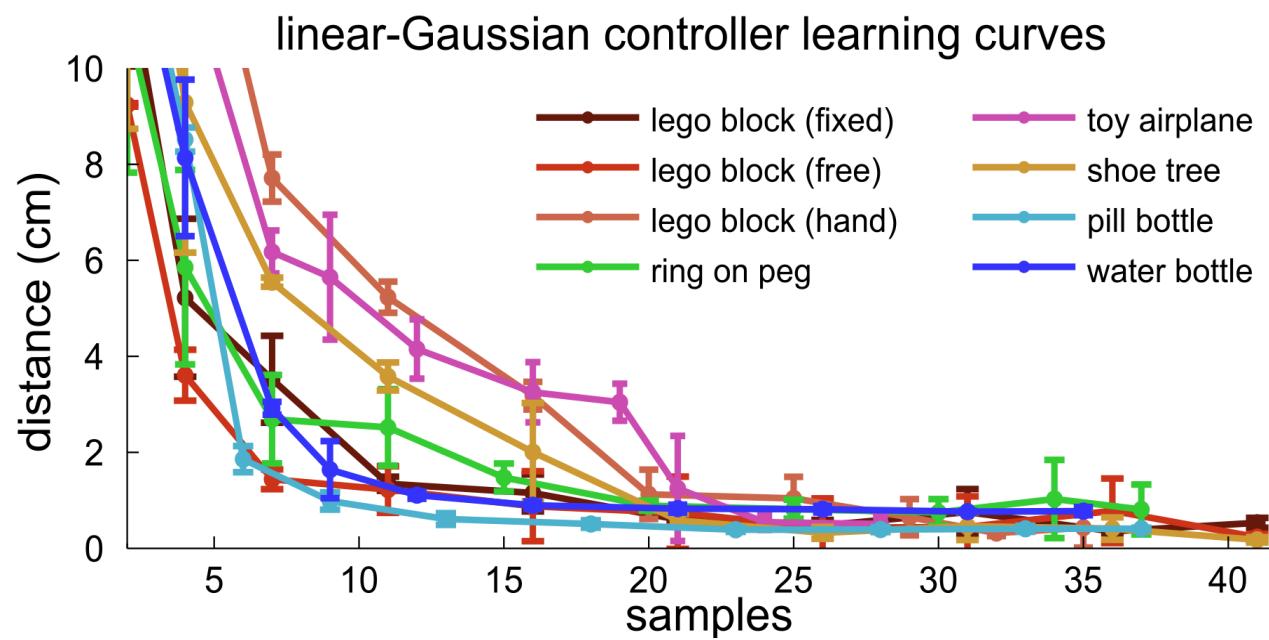
- Issue with just optimal trajectories
 - Don't learn how to steer back onto optimal trajectories
→ Distribution over near-optimal trajectories
- Issue with two-phase pipeline
 - Representational mismatch trajectory distribution vs. neural net
→ Joint optimization

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

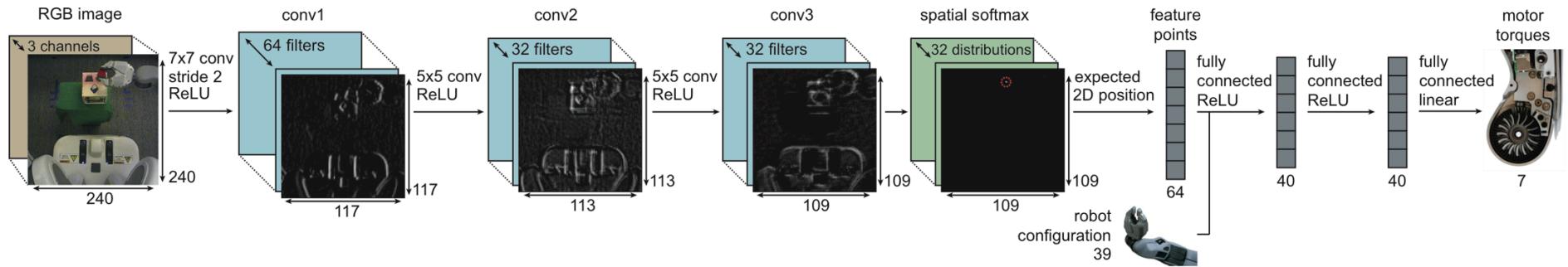
Block Stacking – Learning the Controller for a Single Instance



Linear-Gaussian Controller Learning Curves

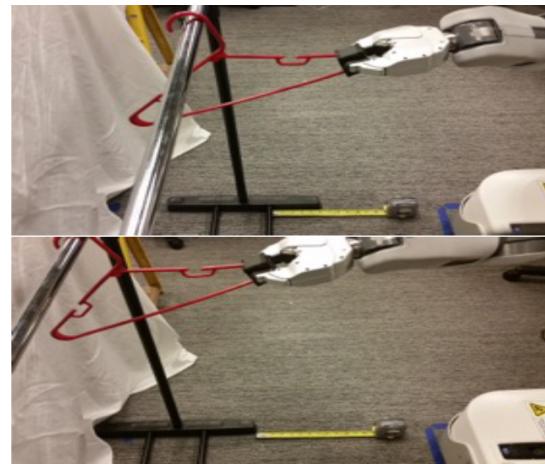
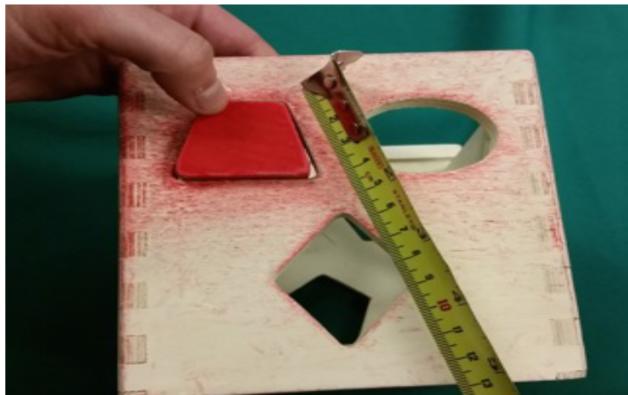


Architecture (92,000 parameters)



[Levine*, Finn*, Darrell, Abbeel, 2015, TR at: rll.berkeley.edu/deeplearningrobotics]

Experimental Tasks



Learning

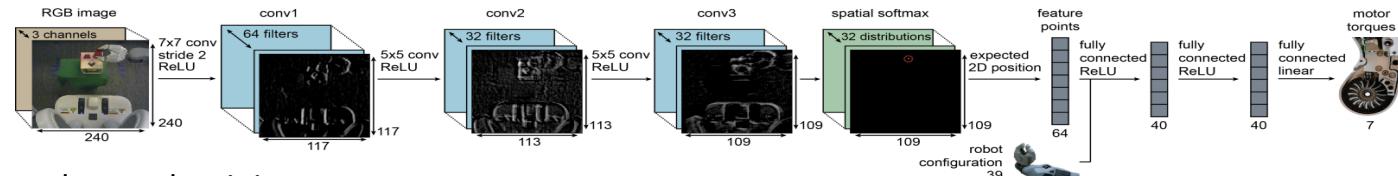


Learned Skills

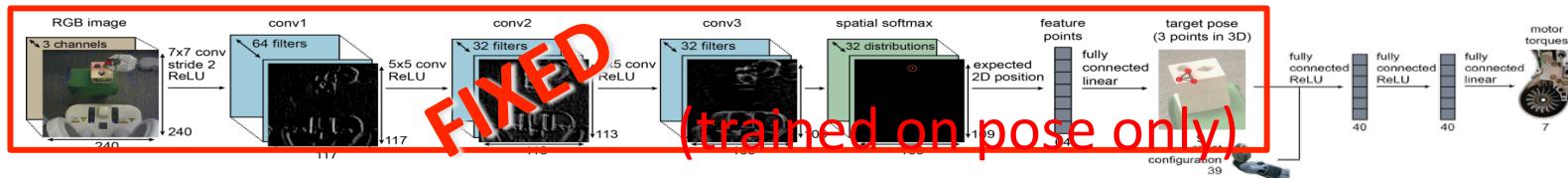


[Levine*, Finn*, Darrell, Abbeel, 2015, TR at: rll.berkeley.edu/deeplearningrobotics]

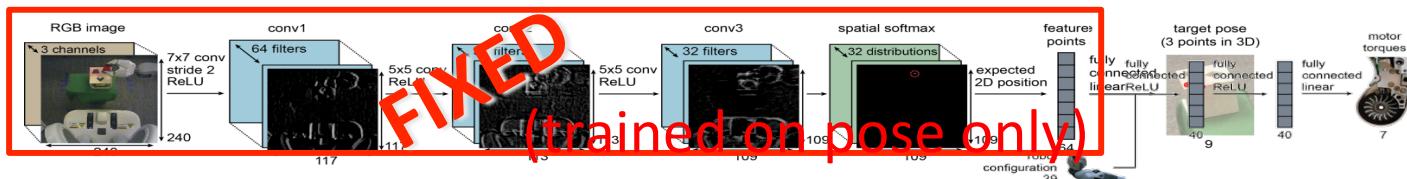
Comparisons



end-to-end training



pose prediction

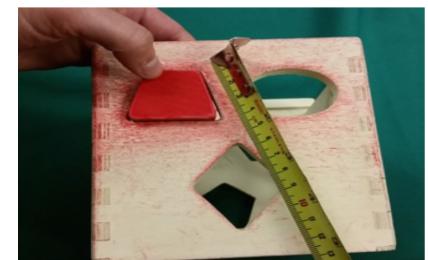
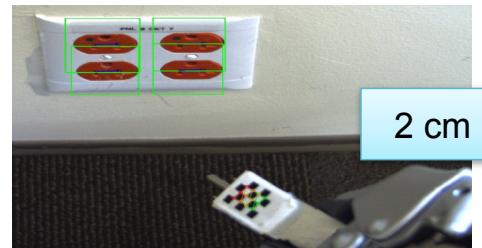


pose features

Comparisons

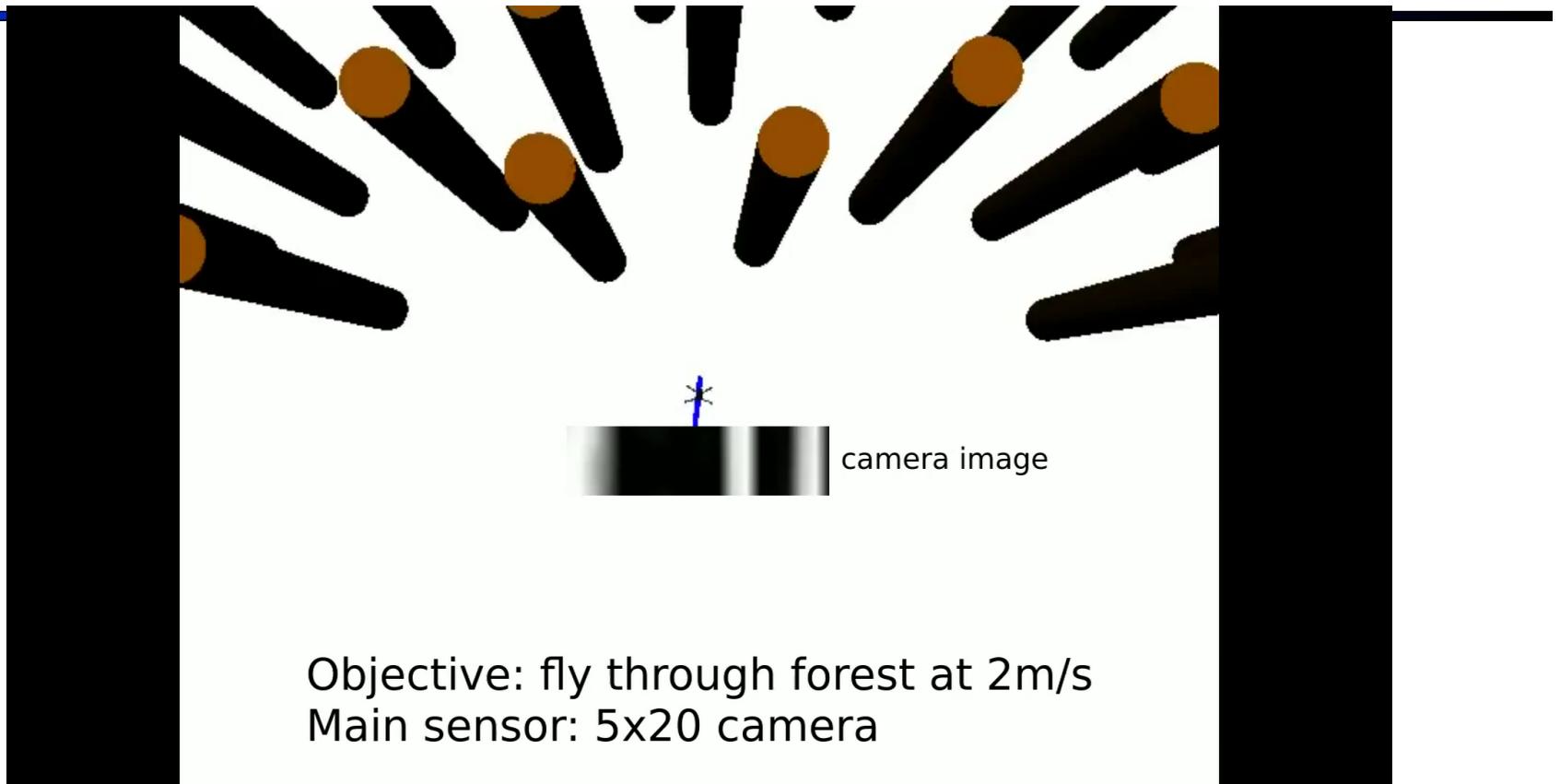
| | |
|---------------------|--------------|
| coat hanger | success rate |
| pose prediction | 55.6% |
| pose features | 88.9% |
| end-to-end training | 100% |
| shape sorting cube | success rate |
| pose prediction | 0% |
| pose features | 70.4% |
| end-to-end training | 96.3% |
| toy claw hammer | success rate |
| pose prediction | 8.9% |
| pose features | 62.2% |
| end-to-end training | 91.1% |
| bottle cap | success rate |
| pose prediction | n/a |
| pose features | 55.6% |
| end-to-end training | 88.9% |

| network architecture | test error (cm) |
|--|--------------------|
| softmax + feature points (ours) | 1.30 ± 0.73 |
| softmax + fully connected layer | 2.59 ± 1.19 |
| fully connected layer | 4.75 ± 2.29 |
| max-pooling + fully connected | 3.71 ± 1.73 |



Meeussen et al. (Willow Garage)

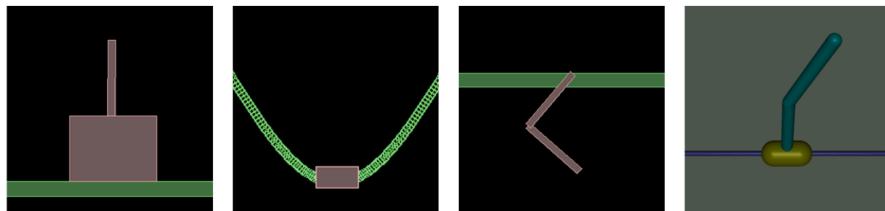
Experiments: Learned Neural Network Policy



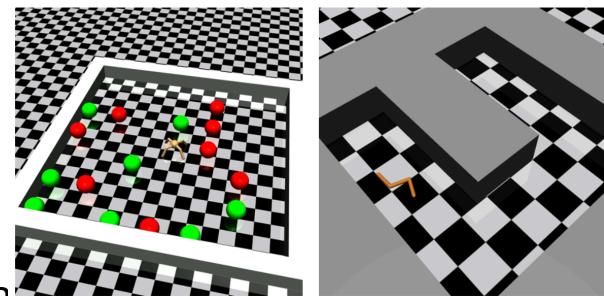
[Khan, Zhang, Levine, Abbeel 2016]

Deep RL Benchmarking -- Tasks

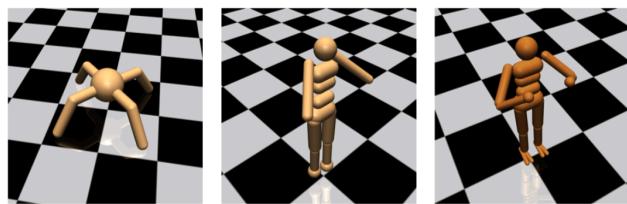
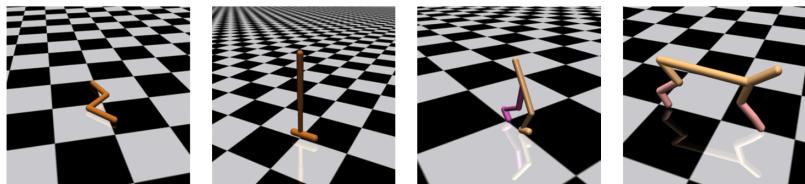
1. Basic tasks



3. Hierarchical



2. Locomotion



4. Partially Observable

sensing, delayed action, sysID

5. Driving...

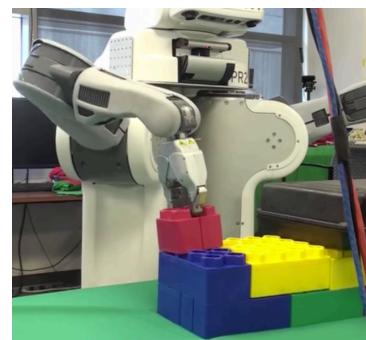
Deep RL Benchmarking -- Algorithms

- 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)
- ...

Started benchmarking (on arXiv soon...)

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.

Frontiers: Applications

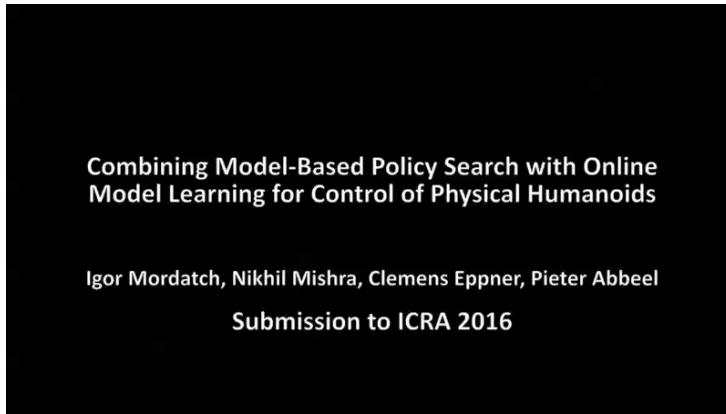


Frontiers

- Shared and transfer learning



- Transfer simulation -> real world



[Mordatch, Mishra, Eppner, Abbeel ICRA 2016]

- Memory

- Estimation

- Unsupervised + RL

Next Time: Natural Language Processing!
