# Guided Policy Search

Sergey Levine

# Learning on PR2



10x real time

iteration 1
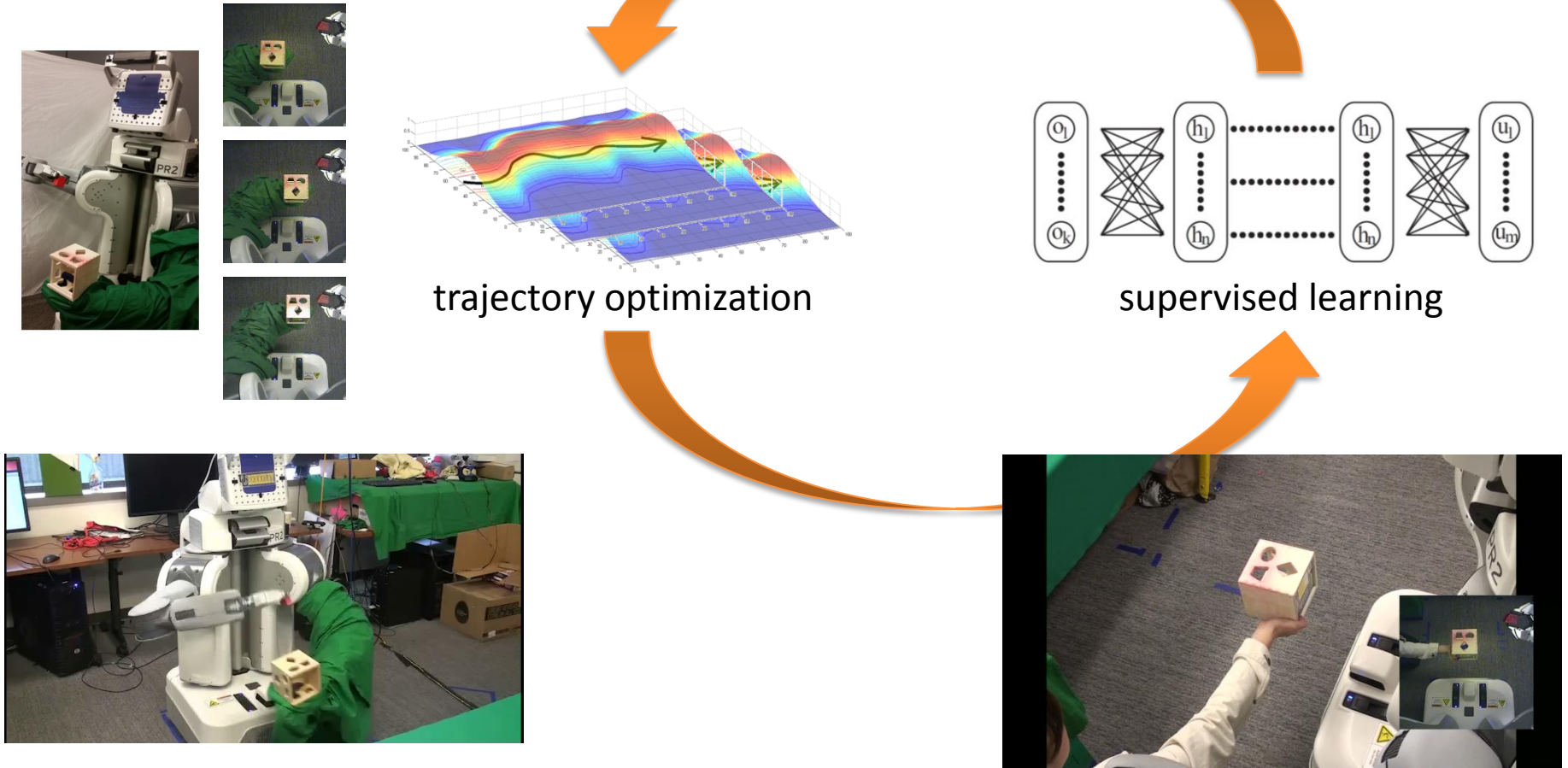
# Shape sorting cube



Learned Visuomotor Policy:  Shape sorting cube

# Visuomotor Policies



**Various Experiments**
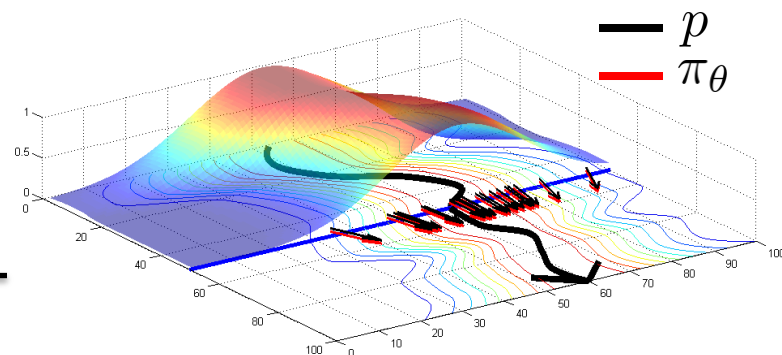Including the policy input

# Guided Policy Search



trajectory optimization

supervised learning

expectation under
current policy

$$\min_{\theta} \overbrace{E_{\pi_\theta}[c(\tau)]}$$

$$\min_{\theta, p(\tau)} E_p[c(\tau)] \quad \longleftarrow \text{trajectory distribution(s)}$$

$$s.t. \quad \pi_\theta(\mathbf{u}_t | \mathbf{o}(\mathbf{x}_t)) = p(\mathbf{u}_t | \mathbf{x}_t) \quad \forall t, \mathbf{x}_t, \mathbf{u}_t$$


$p$
$\pi_\theta$

Lagrange multiplier

$$\mathcal{L}(\theta, p, \lambda) = E_p[c(\tau)] + \sum_{t=1}^{T} \overbrace{\lambda_t} D_t(\pi_\theta, p)$$

optimize
$\mathcal{L}(\theta, p, \lambda)$
w.r.t. $p(\tau)$

optimize
$\mathcal{L}(\theta, p, \lambda)$
w.r.t. $\theta$

update $\lambda$ with
subgradient descent:

$$\lambda_t \leftarrow \lambda_t + \eta D_t(\pi_\theta, p)$$

# Supervised Learning Objective

$$\theta = \arg\min_{\theta} \sum_{t=1}^{T} E_{p(\mathbf{x_t})}[\rho_t D_{KL}(\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)\|p(\mathbf{u}_t|\mathbf{x}_t)) + \lambda_t^T E_{\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)}[\mathbf{u}_t]]$$

$$\pi_\theta(\mathbf{u}_t|\mathbf{x}_t) = \mathcal{N}(\mu_\pi(\mathbf{x}_t), \Sigma_\pi(\mathbf{x}_t))$$

$$p(\mathbf{u}_t|\mathbf{x}_t) = \mathcal{N}(\mu_p(\mathbf{x}_t), \Sigma_p(\mathbf{x}_t))$$

generate samples from $p(\mathbf{x}_t)$ by executing $p(\mathbf{u}_t|\mathbf{x}_t)$

$$\frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \rho_t D_{KL}(\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)\|p(\mathbf{u}_t|\mathbf{x}_t)) + \lambda_t^T \mu_\pi(\mathbf{x}_t)$$

$$\frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \frac{\rho_t}{2} \left[ (\mu_\pi(\mathbf{x}_t) - \mu_p(\mathbf{x}_t))^T \Sigma_p(\mathbf{x}_t)^{-1} (\mu_\pi(\mathbf{x}_t) - \mu_p(\mathbf{x_t})) + \text{tr}\left(\Sigma_\pi(\mathbf{x}_t)\Sigma_p(\mathbf{x}_t)^{-1}\right) + \log\frac{|\Sigma_p(\mathbf{x}_t)|}{|\Sigma_\pi(\mathbf{x}_t)|} \right] + \lambda_t^T \mu_\pi(\mathbf{x}_t)$$

$$\mu_\pi : \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \frac{\rho_t}{2} \|\mu_\pi(\mathbf{x}_t) - \mu^\star(\mathbf{x}_t)\|_{\Sigma_p(\mathbf{x}_t)^{-1}} \qquad \mu^\star(\mathbf{x_t}) = \mu_p(\mathbf{x}_t) - \Sigma_p(\mathbf{x}_t)\lambda_t$$

# Trajectory Optimization (<u>without</u> GPS)

Goal: optimize Gaussian trajectory distribution $p(\tau)$ w.r.t. $E_p[c(\tau)]$

Must optimize time-varying linear-Gaussian controller $p(\mathbf{u}_t | \mathbf{x}_t)$

Controller has form $p(\mathbf{u}_t | \mathbf{x}_t) = \mathcal{N}(\mathbf{K}_t \mathbf{x}_t, \Sigma_t)$

Use LQR to get $\mathbf{K}_t$, but what is $\Sigma_t$?

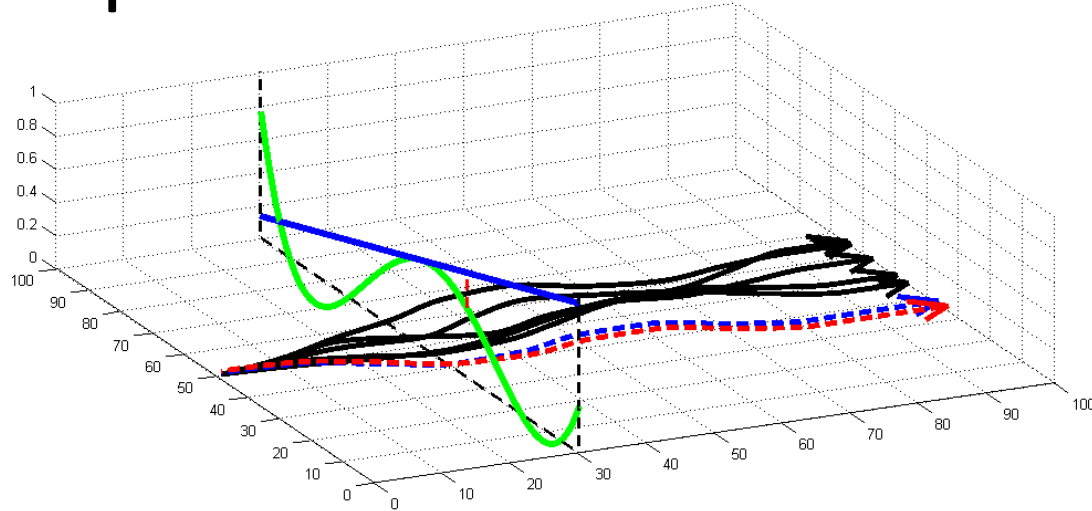If $\mathbf{x}_t$ is Markovian, $\Sigma_t = 0$ always (but this is boring...)

Let's instead optimize $E_p[c(\tau)] - \mathcal{H}(p) = \underbrace{\sum_{t=1}^{T} E_{p(\mathbf{x}_t, \mathbf{u}_t)}[c(\mathbf{x}_t, \mathbf{u}_t)] - \mathcal{H}(p(\mathbf{u}_t | \mathbf{x}_t))}_{E_p[c(\tau)] - \mathcal{H}(p)}$

(we'll see why soon...)

Maximum entropy solution is simply $p(\mathbf{u}_t | \mathbf{x}_t) = \mathcal{N}(\mathbf{K}_t \mathbf{x}_t, \underbrace{\mathbf{R}_t + \mathbf{B}_t^T \mathbf{P}_{t+1} \mathbf{B}}_{})$

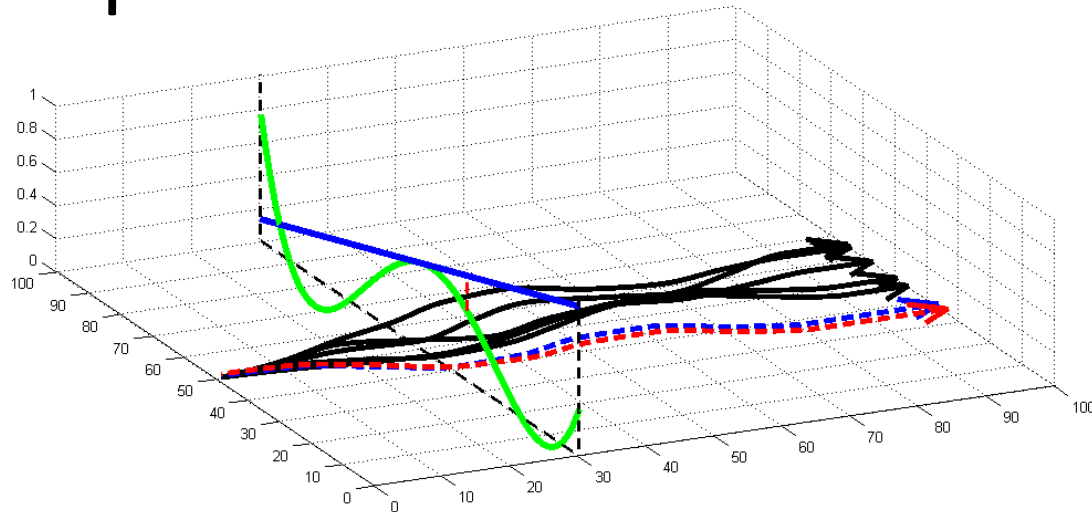LQR cost-to-go w.r.t. $\mathbf{u}_t$, sometimes written as $Q_{\mathbf{uut}}$

# Trajectory Optimization



1. Run time-varying policy $p(\mathbf{u}_t|\mathbf{x}_t)$ on robot $N$ times
2. Collect dataset $\mathcal{D} = \{\tau_i\}$ where $\tau_i = \{\mathbf{x}_{1i}, \mathbf{u}_{1i}, \ldots, \mathbf{x}_{Ti}, \mathbf{u}_{Ti}\}$
3. For each $t \in \{0, \ldots, T-1\}$, fit linear Gaussian $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$
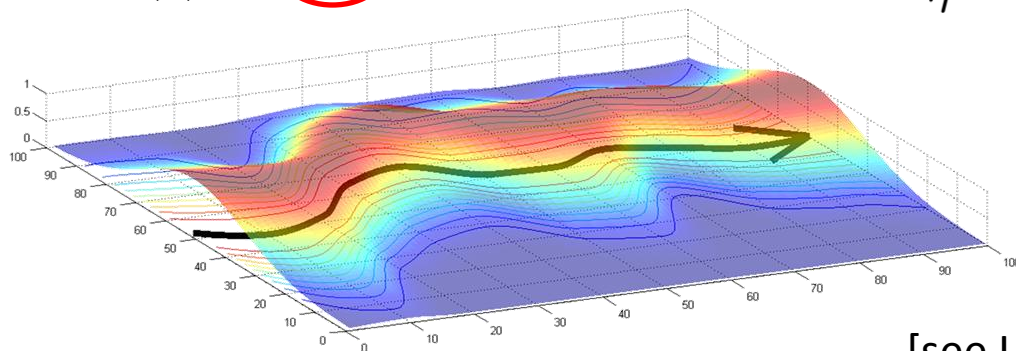4. Solve control problem to get new $p(\mathbf{u}_t|\mathbf{x}_t)$
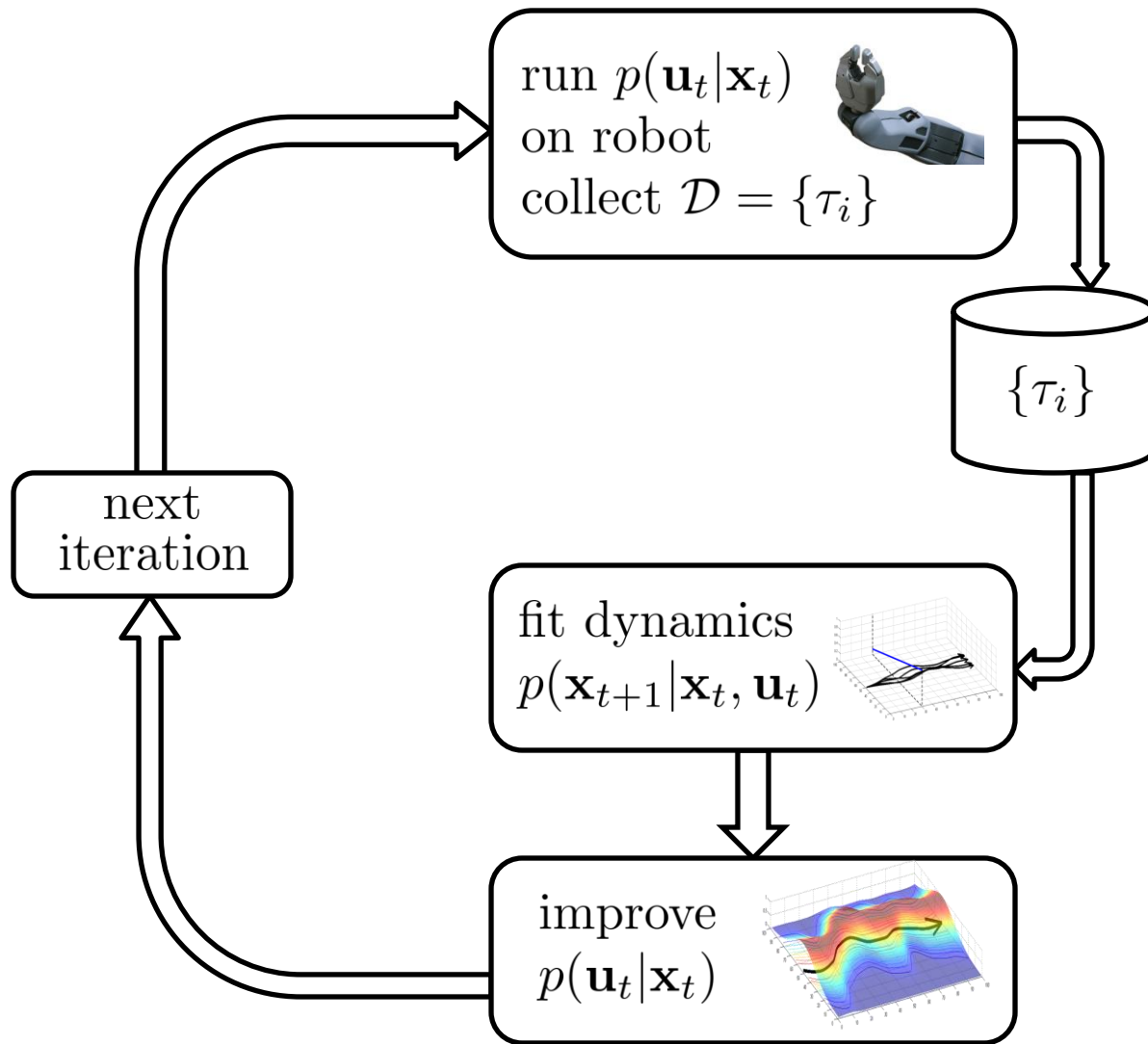
# Trajectory Optimization



$$\min_{p(\tau)} E_p[c(\tau)] \text{ s.t. } D_{KL}(p(\tau)\|\bar{p}(\tau)) \leq \epsilon$$

new     old

$$\frac{1}{\eta}\mathcal{L}(p,\eta) = E_p\left[\frac{1}{\eta}c(\tau) - \log \bar{p}(\tau)\right] - \mathcal{H}(p) - \epsilon$$

$$\min_{p(\tau)} E_p[\tilde{c}(\tau)] - \mathcal{H}(p) \qquad \tilde{c}(\tau) = \frac{1}{\eta}c(\tau) - \log \bar{p}(\tau)$$



[see Levine & Abbeel '14 for details]

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

$\{\tau_i\}$

next
iteration

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

improve
$p(\mathbf{u}_t|\mathbf{x}_t)$

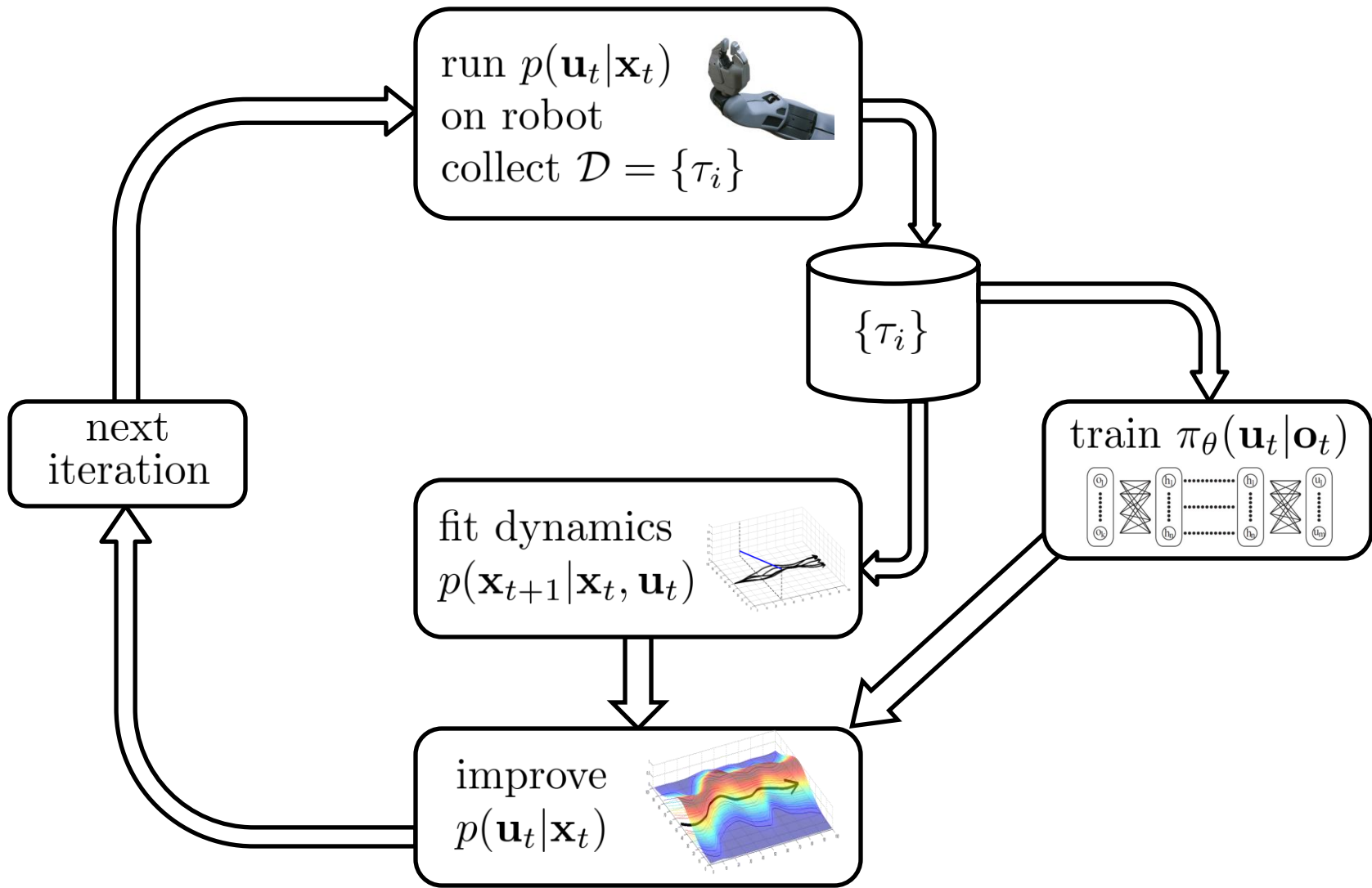[see L. et al. NIPS '14 for details]

# Trajectory Optimization (with GPS)

$$\min_{p(\tau)} \sum_{t=1}^{T} E_{p(\mathbf{x}_t, \mathbf{u}_t)} [c(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{u}_t^T \lambda_t + \rho_t D_{KL}(p(\mathbf{u}_t|\mathbf{x}_t) \| \pi_\theta(\mathbf{u}_t|\mathbf{x}_t)))]$$

$$\text{s.t. } D_{KL}(p(\tau) \| \bar{p}(\tau)) \le \epsilon$$

$$\mathcal{L}(p, \eta) = \sum_{t=1}^{T} E_{p(\mathbf{x}_t, \mathbf{u}_t)} [c(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{u}_t^T \lambda_t + \rho_t D_{KL}(p(\mathbf{u}_t|\mathbf{x}_t) \| \pi_\theta(\mathbf{u}_t|\mathbf{x}_t)) + \eta D_{KL}(p(\mathbf{u}_t|\mathbf{x}_t \| \bar{p}(\mathbf{u}_t, \mathbf{x}_t)))]$$

$$\mathcal{L}(p, \eta) = \sum_{t=1}^{T} E_{p(\mathbf{x}_t, \mathbf{u}_t)} [\underbrace{c(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{u}_t^T \lambda_t - \rho_t \log \pi_\theta(\mathbf{u}_t|\mathbf{x}_t) - \eta \bar{p}(\mathbf{u}_t, \mathbf{x}_t)}_{\tilde{c}(\mathbf{x}_t, \mathbf{u}_t)}] - (\rho_t + \eta) \mathcal{H}(p(\mathbf{u}_t|\mathbf{x}_t))$$
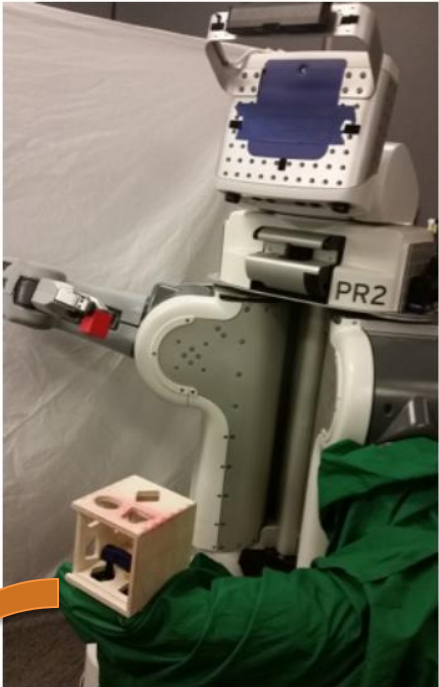
$$\mathcal{L}(p, \eta) = \sum_{t=1}^{T} E_{p(\mathbf{x}_t, \mathbf{u}_t)} [\tilde{c}(\mathbf{x}_t, \mathbf{u}_t)] - \nu_t \mathcal{H}(p(\mathbf{u}_t|\mathbf{x}_t))$$
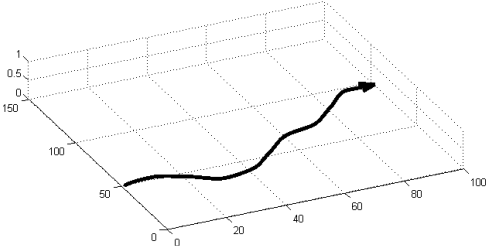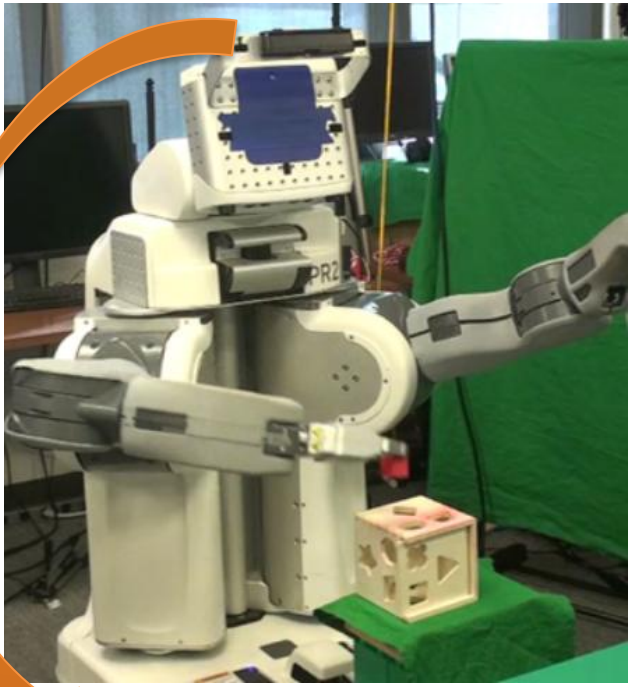
maximum entropy objective (like before)

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

$\{\tau_i\}$

next
iteration

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

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

improve
$p(\mathbf{u}_t|\mathbf{x}_t)$

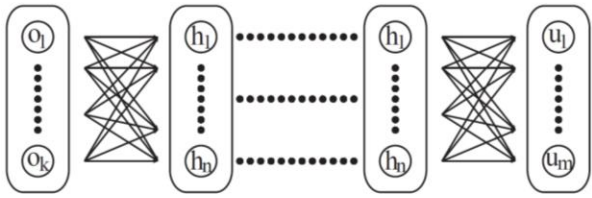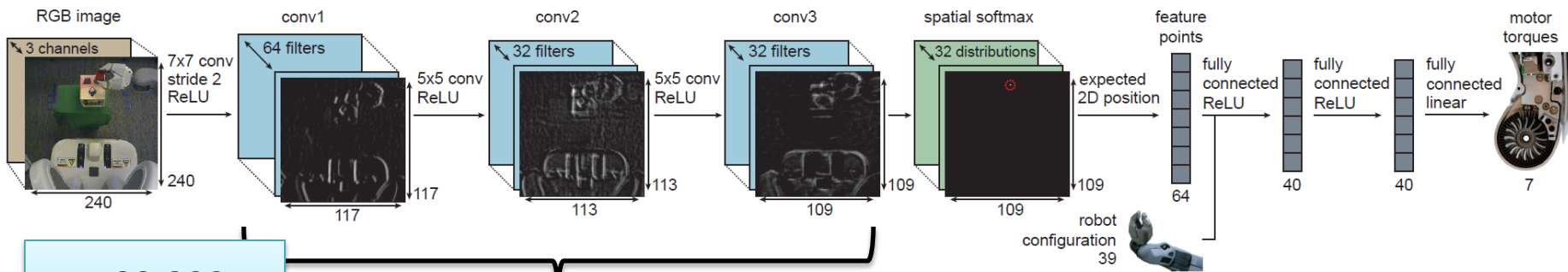[see L. et al. NIPS '14 for details]

# Instrumented Training

training time

test time



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

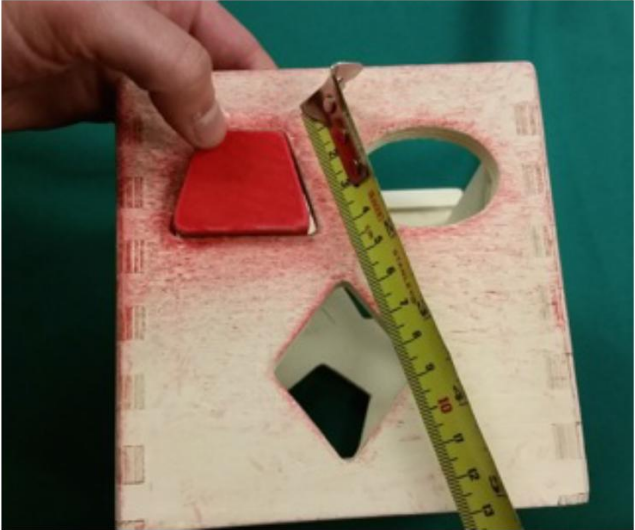$$\mathbf{o}_t \rightarrow \mathbf{u}_t$$

RGB image — 3 channels — 240 × 240

conv1 — 64 filters — 7x7 conv stride 2 ReLU — 117 × 117

conv2 — 32 filters — 5x5 conv ReLU — 113 × 113

conv3 — 32 filters — 5x5 conv ReLU — 109 × 109

spatial softmax — 32 distributions — 109 × 109 — expected 2D position

feature points — 64

fully connected ReLU — 40

fully connected ReLU — 40

fully connected linear — 7

motor torques

robot configuration — 39

~ 92,000 parameters

Chelsea Finn

$$\text{conv3 } a_{ij} \quad (\text{x}32)$$

$$\frac{e^{a_{ij}}}{\sum_{ij} e^{a_{ij}}}$$

| $(-2,2)$ | $(-1,2)$ | $(1,2)$ | $(2,2)$ |
| $(-2,1)$ | $(-1,1)$ | $(1,1)$ | $(2,1)$ |
| $(-2,-1)$ | $(-1,-1)$ | $(1,-1)$ | $(2,-1)$ |
| $(-2,-2)$ | $(-1,-2)$ | $(1,-2)$ | $(2,-2)$ |

$\otimes$

$\sum$

$(-0.2, 0.4)$ $\quad (\text{x}32)$

# Experimental Tasks

# Shape sorting cube



Learned Visuomotor Policy:  Shape sorting cube

# Hanger



**Learned Visuomotor Policy: Hanger Task**

# Hammer



Learned Visuomotor Policy: Hammer Task
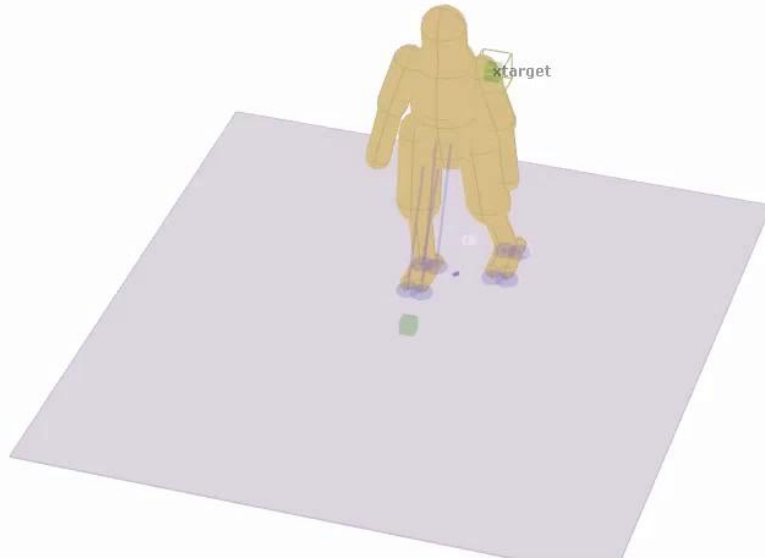
# Bottle

**Learned Visuomotor Policy: Bottle Task**

# Locomotion

Igor Mordatch

better trajectory optimization +
large scale simulation

veltarget_x

veltarget_y

xtarget

# Darwin Robot

better trajectory optimization +
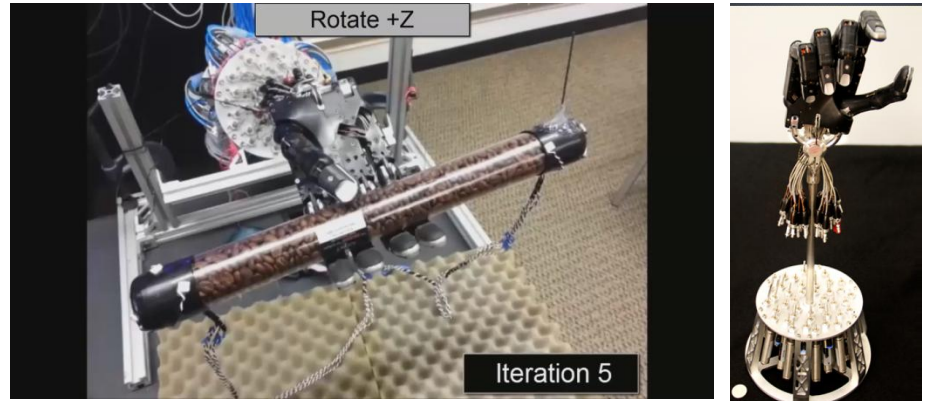large scale simulation +
adaptation to real world dynamics

Igor Mordatch



Mordatch, Mishra, Eppner, Abbeel

# Guided Policy Search Applications

## manipulation



with N. Wagener and P. Abbeel

## dexterous hands



with V. Kumar and E. Todorov

## locomotion
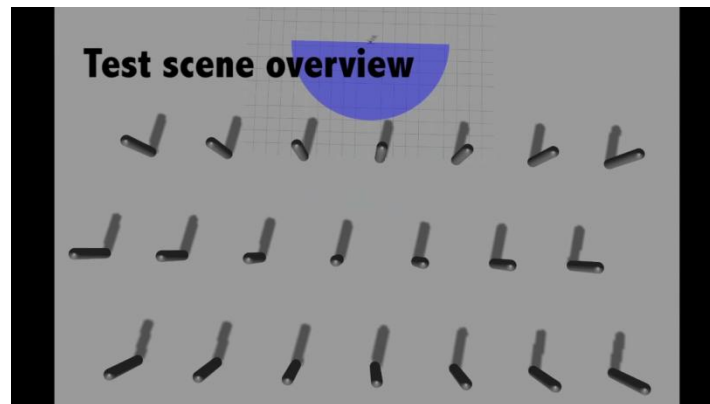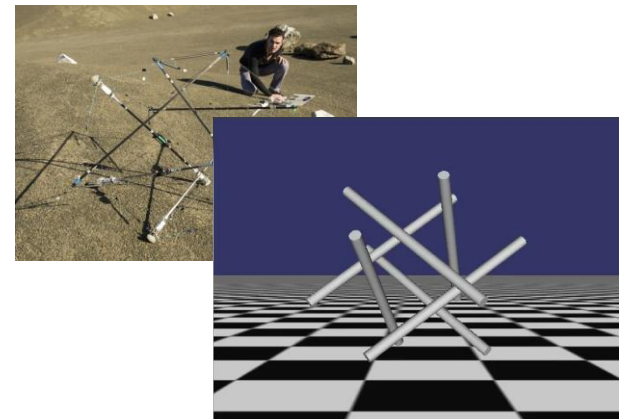


constrained GPS
300–400 N pushes

400N

with V. Koltun

## aerial vehicles



Test scene overview

with G. Kahn, T. Zhang, P. Abbeel

## tensegrity robot



with M. Zhang, K. Caluwaerts, P. Abbeel

# DAGGER

A simpler way to turn policy search into supervised learning

Requires a "stronger" teacher – must give optimal action $\mathbf{u}$ in any state $\mathbf{x}$

Typically used for imitation learning from a human expert

Initialize $\mathcal{D} \leftarrow \emptyset$.
Initialize $\hat{\pi}_1$ to any policy in $\Pi$.
**for** $i = 1$ **to** $N$ **do**

typically 0.0, except when i = 1, then 1.0

    Let $\pi_i = \beta_i \pi^* + (1 - \beta_i)\hat{\pi}_i$.
    Sample $T$-step trajectories using $\pi_i$.
    Get dataset $\mathcal{D}_i = \{(s, \pi^*(s))\}$ of visited states by $\pi_i$
    and actions given by expert.
    Aggregate datasets: $\mathcal{D} \leftarrow \mathcal{D} \bigcup \mathcal{D}_i$.
    Train classifier $\hat{\pi}_{i+1}$ on $\mathcal{D}$.
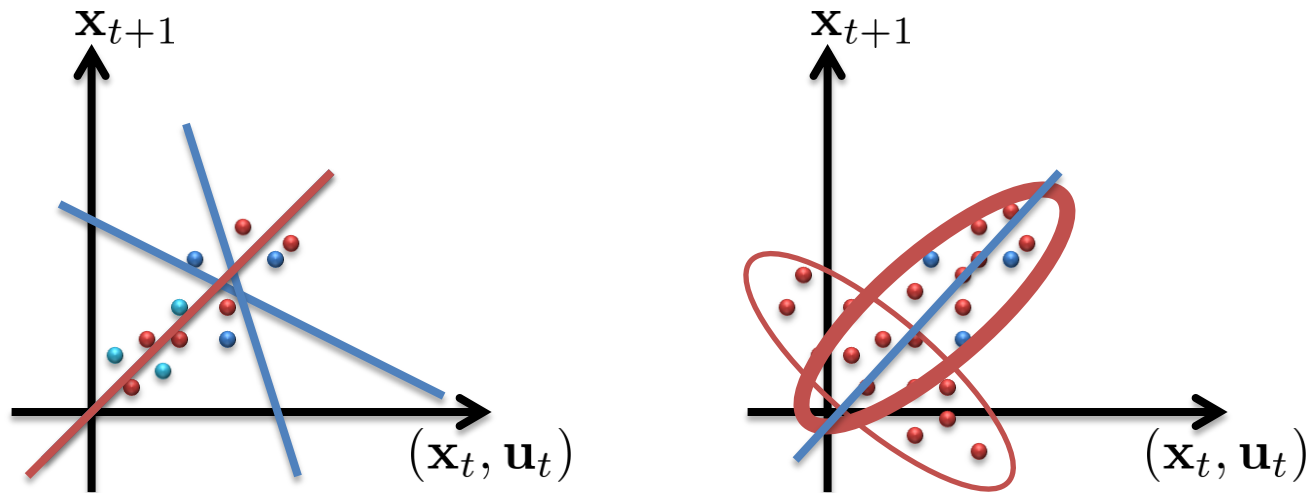**end for**
**Return** best $\hat{\pi}_i$ on validation.

# DAGGER Video

See http://videolectures.net/aistats2011_ross_reduction/

# Trajectory Optimization – Dynamics Fitting



1. Run time-varying policy $p(\mathbf{u}_t|\mathbf{x}_t)$ on robot $N$ times
2. Collect dataset $\mathcal{D} = \{\tau_i\}$ where $\tau_i = \{\mathbf{x}_{1i}, \mathbf{u}_{1i}, \ldots, \mathbf{x}_{Ti}, \mathbf{u}_{Ti}\}$
3. For each $t \in \{0, \ldots, T-1\}$, fit linear Gaussian $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$
4. Solve control problem to get new $p(\mathbf{u}_t|\mathbf{x}_t)$

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

prior

data

next iteration

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

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

update $\eta$

[see L. et al. NIPS '14 for details]
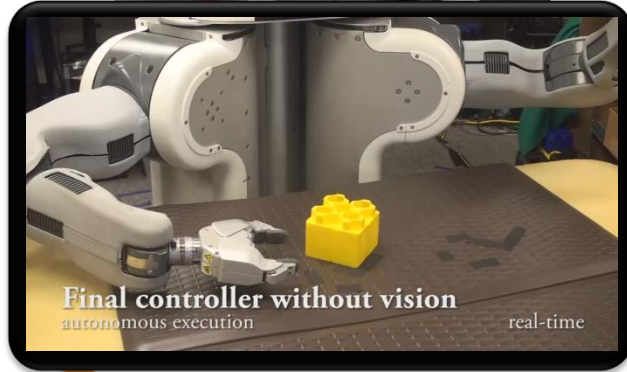
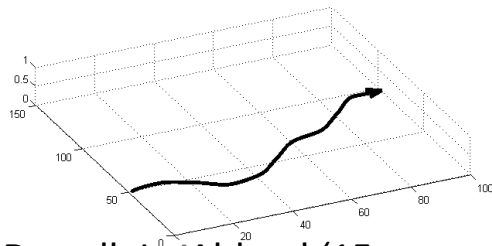# Learned Motion Skills

# More Visuomotor Experiments



real time
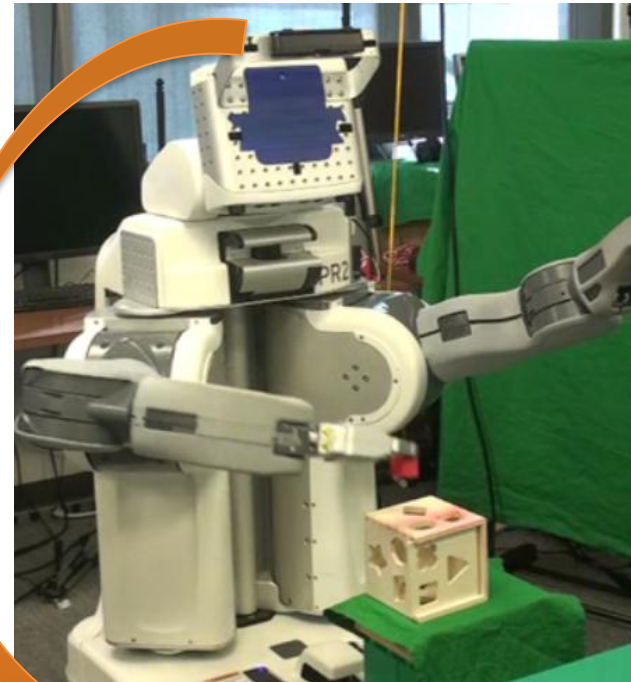
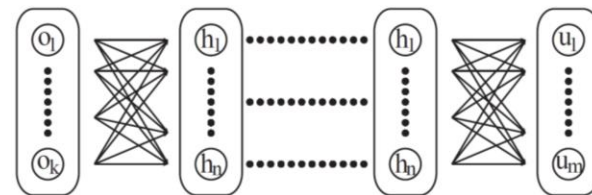autonomous execution

# Beyond Instrumented Training

training time

test time



Final controller using our method
autonomous execution                    real-time

Final controller without vision
autonomous execution                    real-time

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

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

Finn, Tan, Duan, Darrell, L., Abbeel '15

# Learning Visual State Spaces



$$\tilde{\mathbf{x}}_t = \begin{bmatrix} \mathbf{x}_t \\ \mathbf{f}_t \end{bmatrix}$$

1. Set target end-effector pose

Bag Transfer Task