# Learning Dynamic Manipulation Skills under Unknown Dynamics with Guided Policy Search

## Sergey Levine          Pieter Abbeel
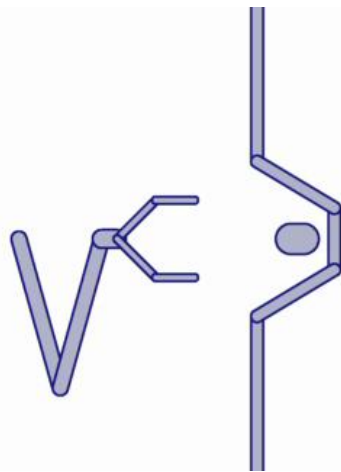
UC Berkeley                    UC Berkeley

Team TROOPER: Lockheed Martin, University of Pennsylvania, Rensselaer Polytechnic Institute
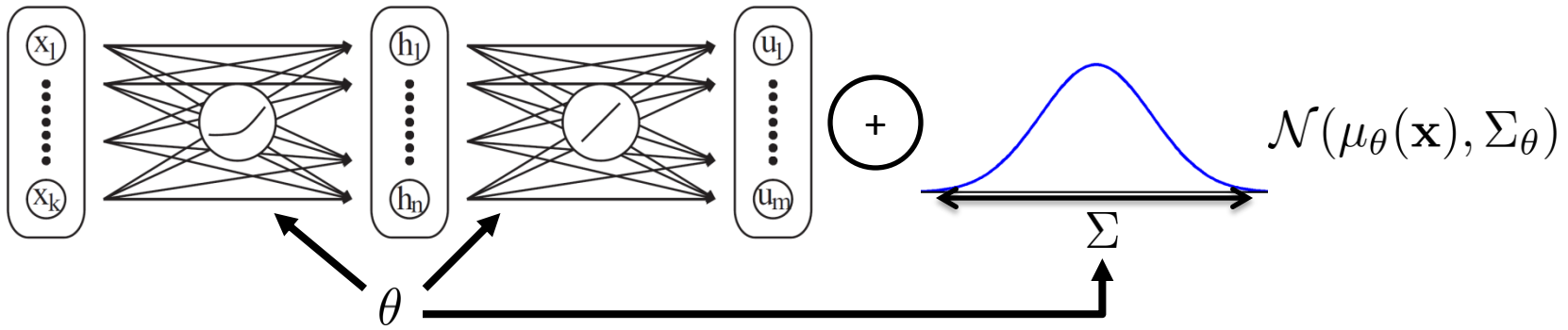


Philipp Krahenbuhl, Stanford University



2D Insertion
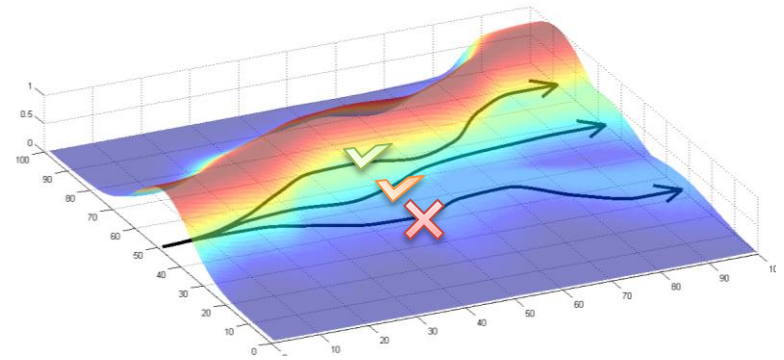unobserved slot positions
[neural network]

general-purpose neural network controller



$$\theta = \arg\min_\theta J(\theta)$$

$$J(\theta) = E_{\pi_\theta}\left[\sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

$$\pi_\theta(\mathbf{u}_t|\mathbf{x}_t) - \text{control policy}$$

$$\mathcal{N}(\mu_\theta(\mathbf{x}), \Sigma_\theta)$$

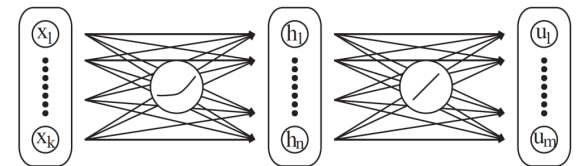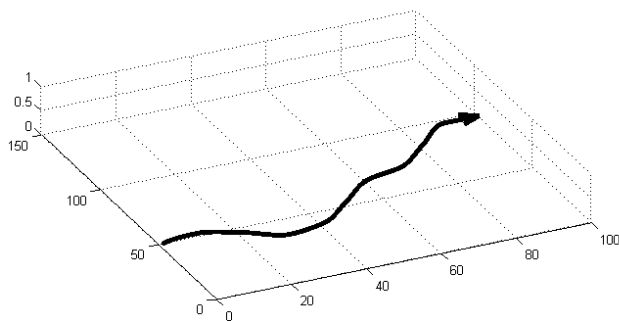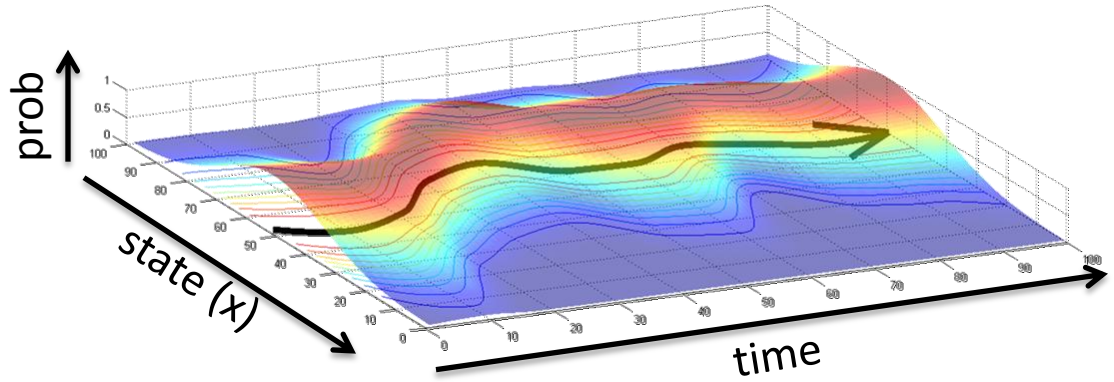| | | | |
|---|---|---|---|
| policy search (RL) | complex dynamics | complex policy | HARD |
| supervised learning | ~~complex dynamics~~ | complex policy | EASY |
| trajectory optimization | complex dynamics | ~~complex policy~~ | EASY |

trajectory optimization



supervised learning
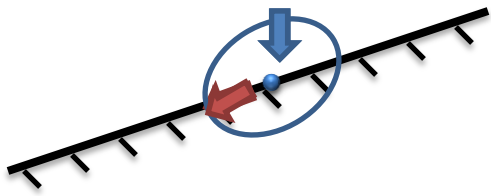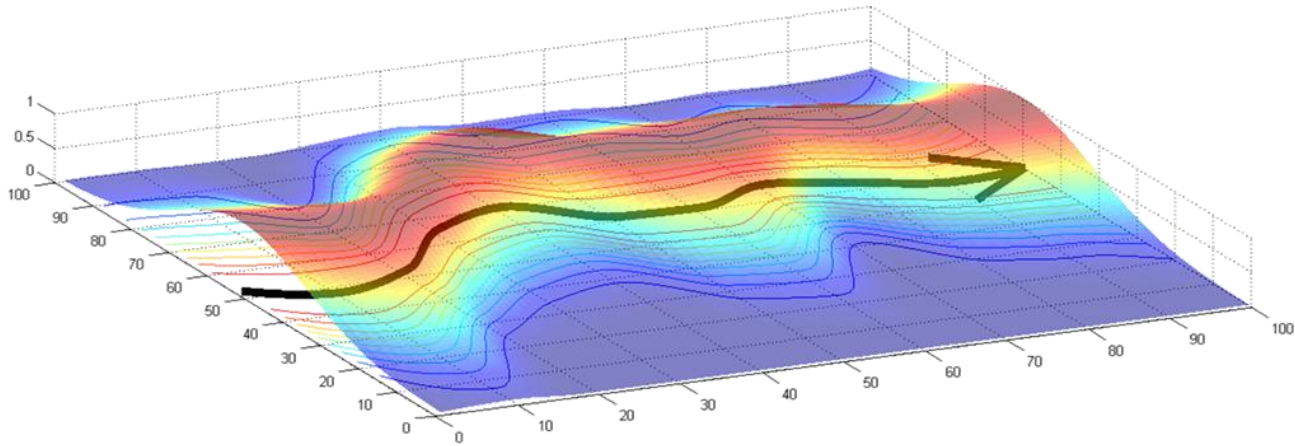
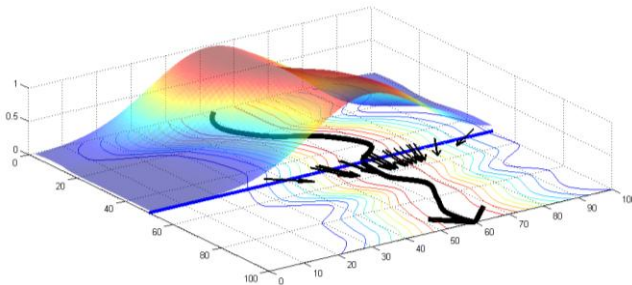# Trajectory Optimization



guided policy search

# Trajectory Optimization
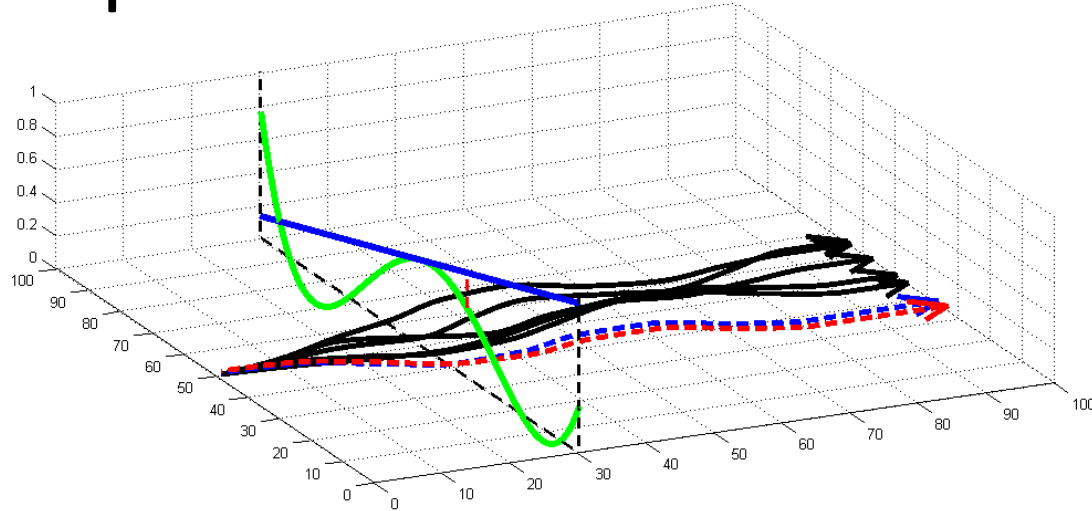


$$\min_{q(\tau)} E_q[c(\tau)] - \mathcal{H}(q)$$

approximate solution using iterative LQR
(similar to extended Kalman filter)

- locally linear dynamics
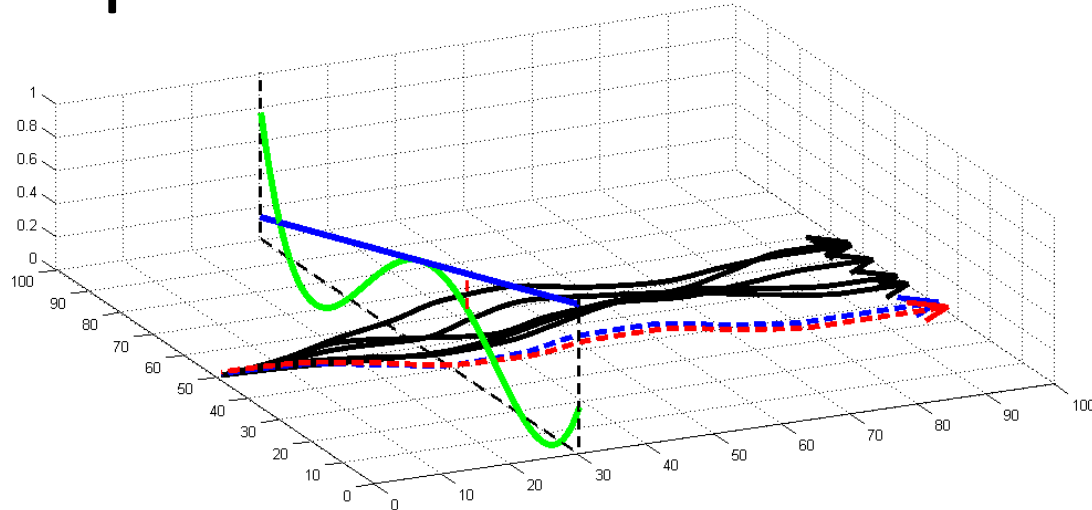- locally quadratic cost
- Gaussian distribution



$$q(\mathbf{u}_t|\mathbf{x}_t) = \mathcal{N}(\mathbf{k}_t + \mathbf{K}_t\mathbf{x}_t, \Sigma_t)$$

# Trajectory Optimization



1. Run time-varying policy $q(\mathbf{u}_t|\mathbf{x}_t)$ 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 $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$

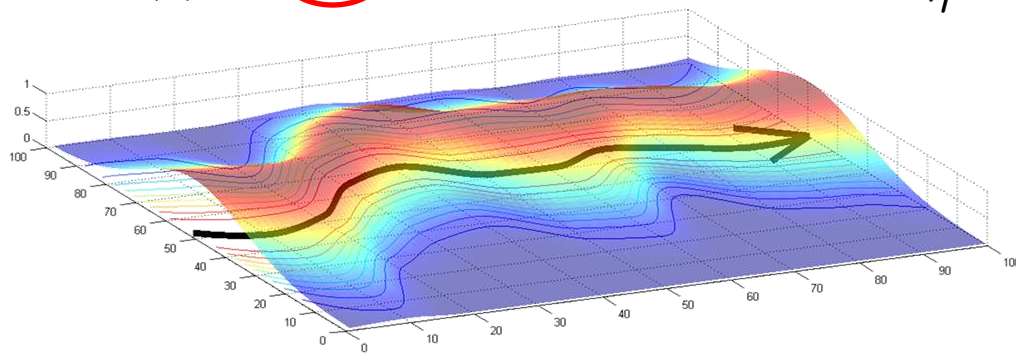4. Solve control problem to get new $q(\mathbf{u}_t|\mathbf{x}_t)$
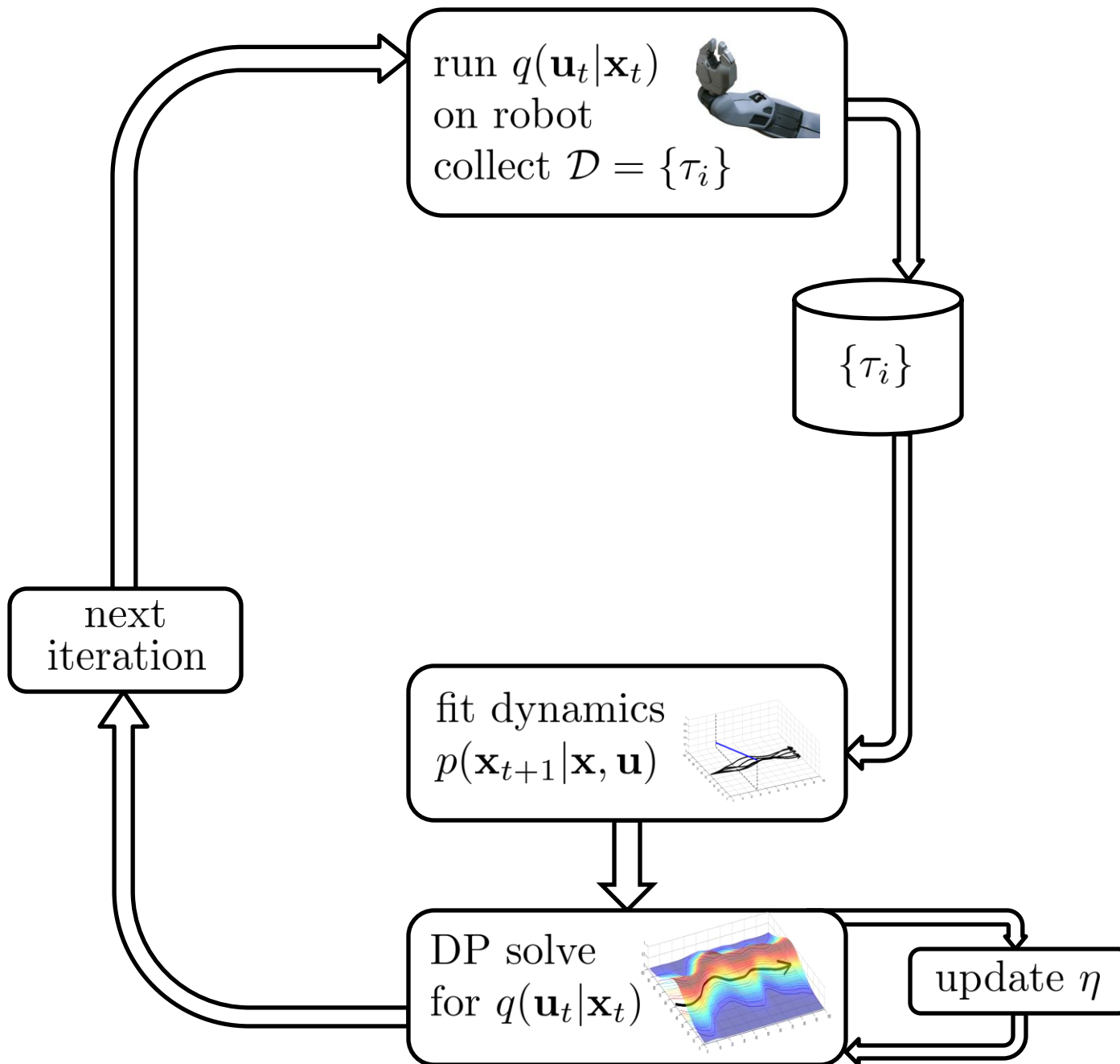
# Trajectory Optimization



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

new    old

$$\frac{1}{\eta}\mathcal{L}(q,\eta) = E_q\left[\frac{1}{\eta}c(\tau) - \log \bar{q}(\tau)\right] - \mathcal{H}(q) - \epsilon$$
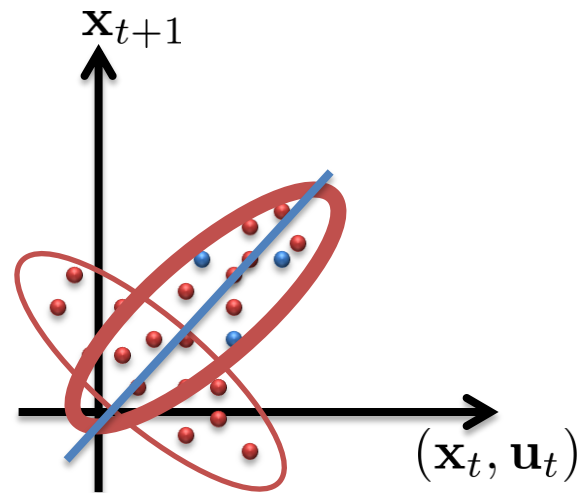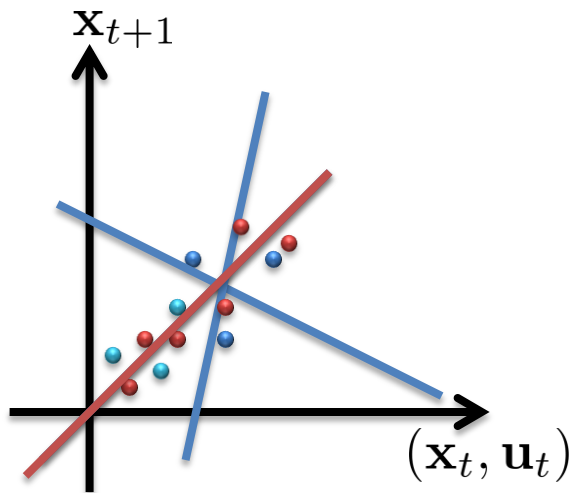
$$\min_{q(\tau)} E_q[c(\tau)] - \mathcal{H}(q) \qquad \tilde{c}(\tau) = \frac{1}{\eta}c(\tau) - \log \bar{q}(\tau)$$

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

$\{\tau_i\}$

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

next
iteration

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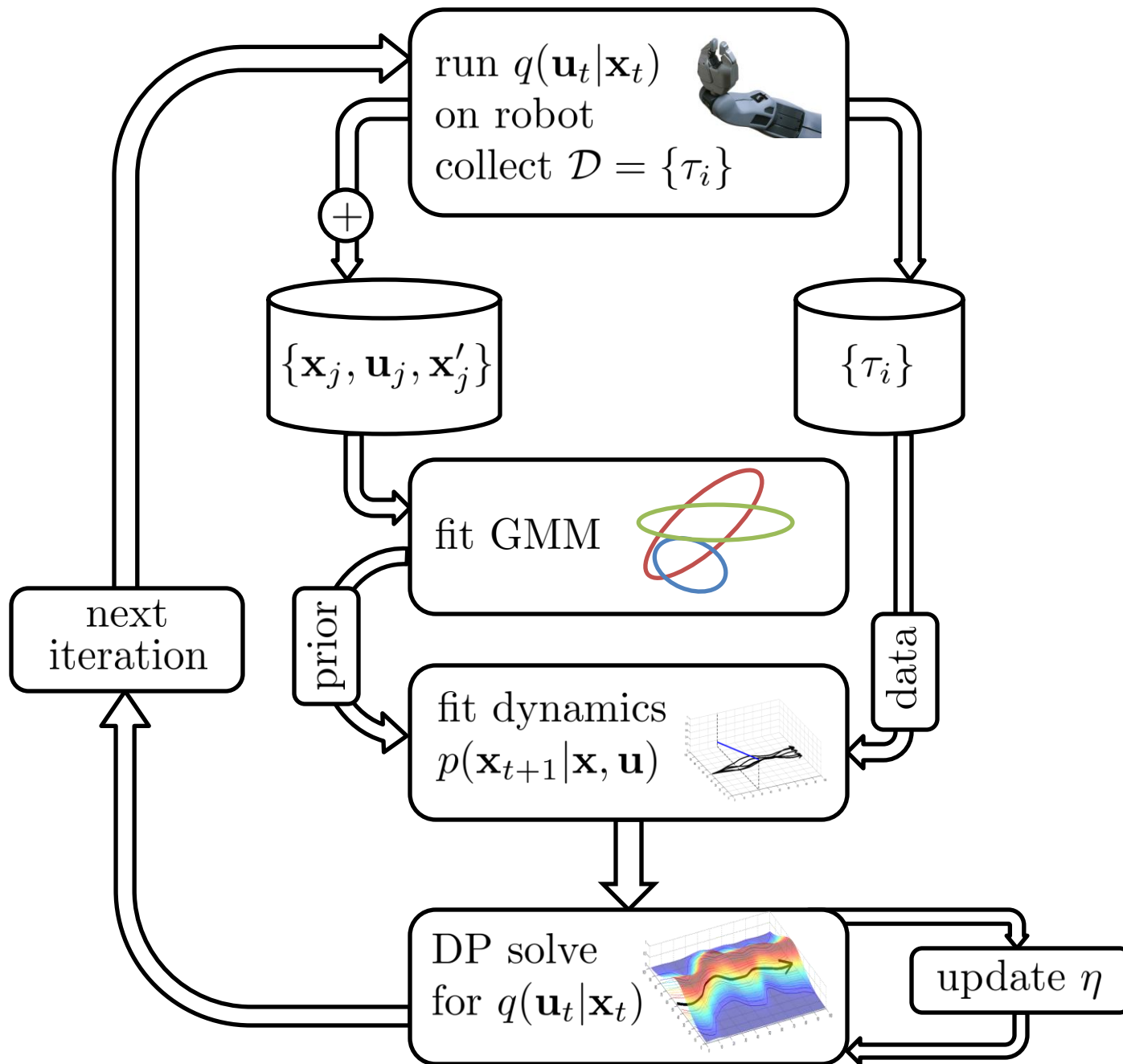
update $\eta$

# Trajectory Optimization



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

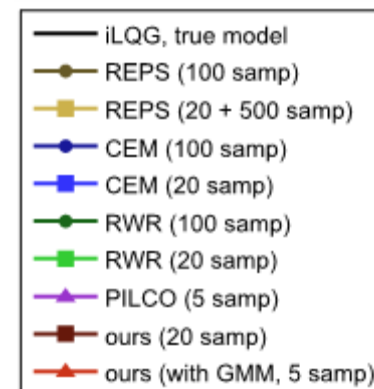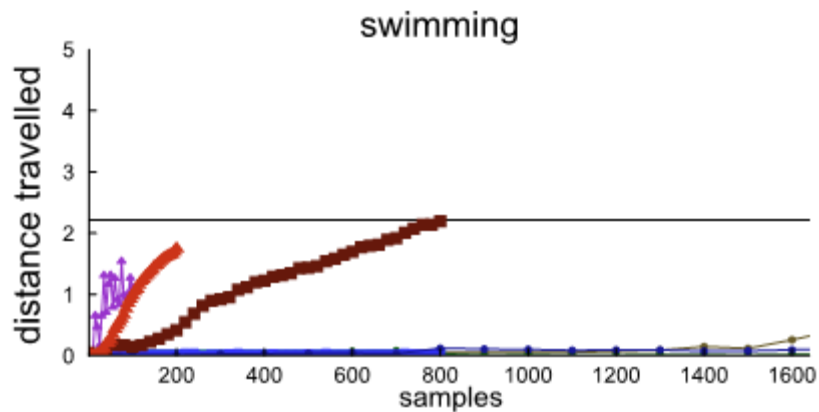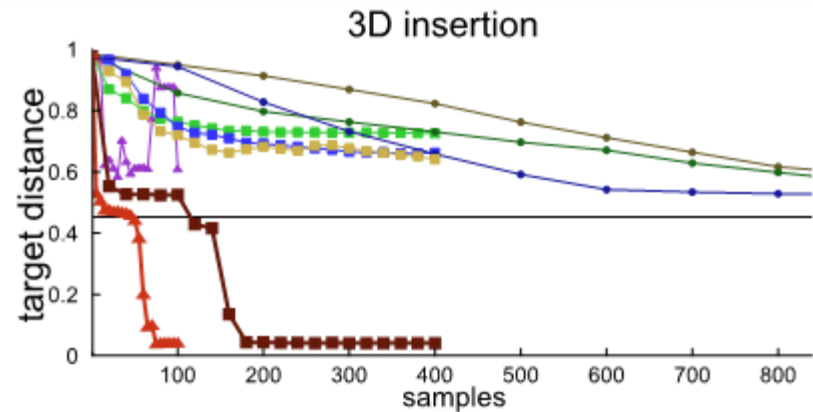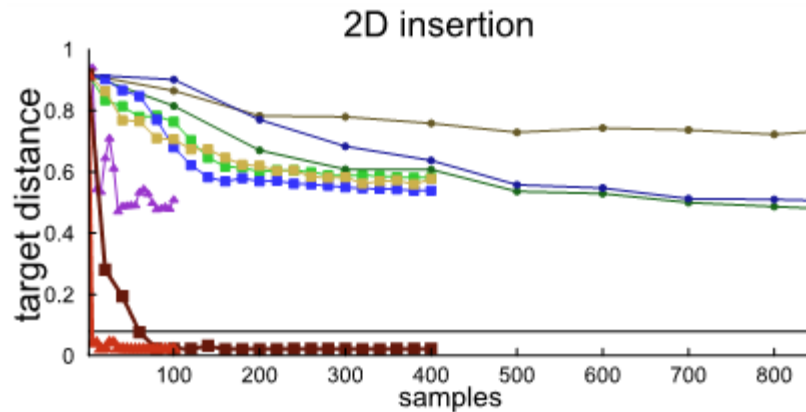# Trajectory Optimization

2D Insertion

optimized trajectory

linear Gaussian control

Swimming

optimized trajectory

[linear Gaussian]
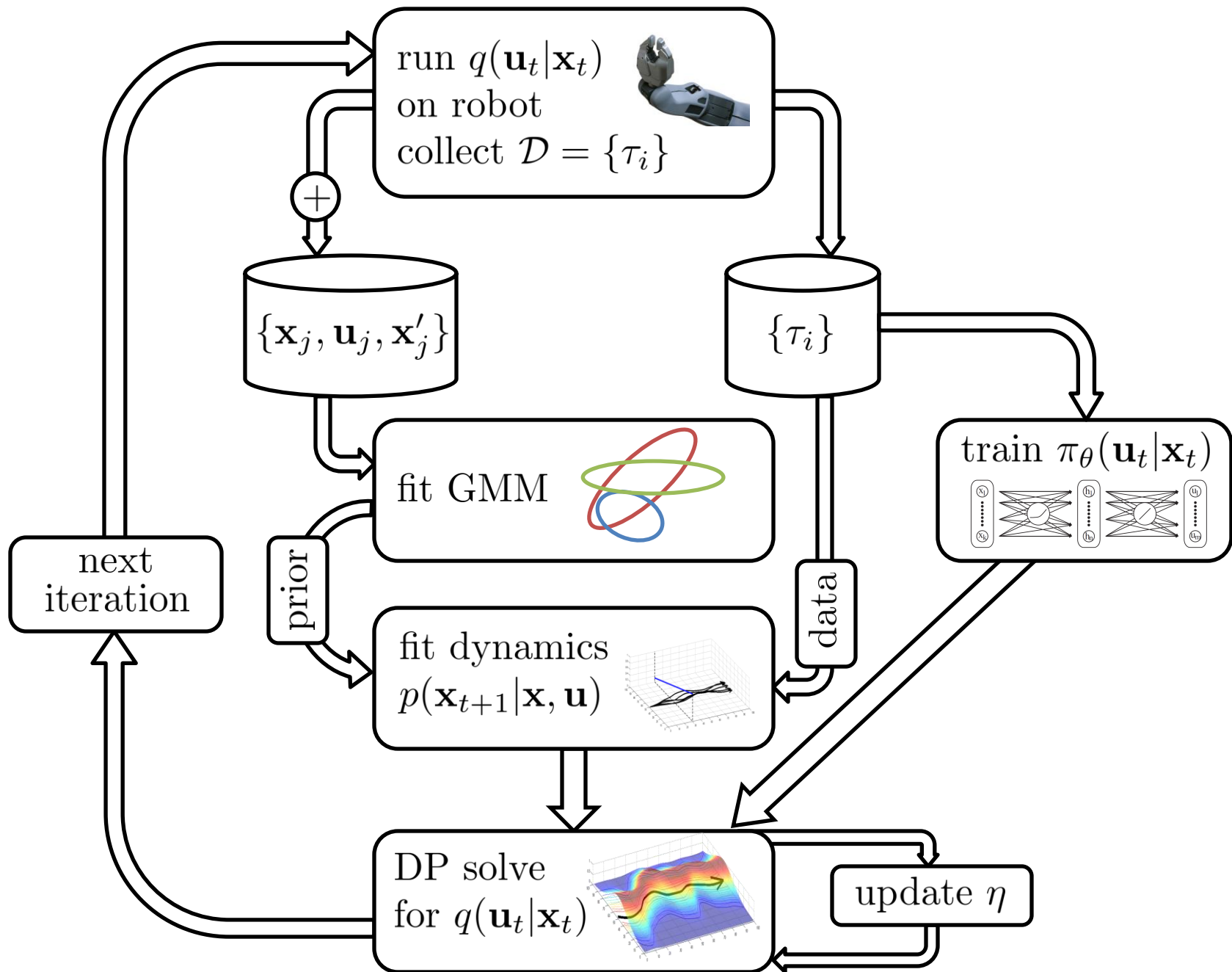
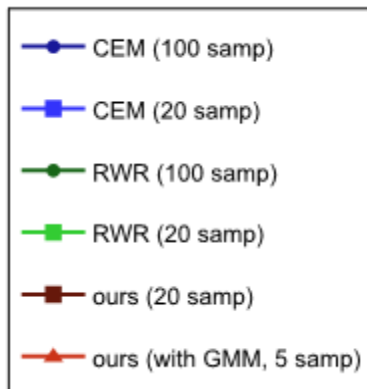# Trajectory Optimization

# Guided Policy Search



see Levine & Koltun, ICML 2014
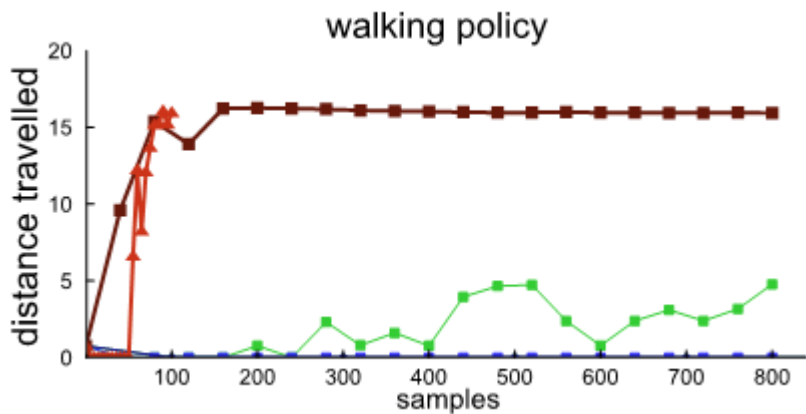
## 2D Insertion

unobserved slot positions

[neural network]

## Swimming

learned policy

[neural network]

## 2D insertion policy

## 3D insertion policy

## swimming policy

## walking policy

CEM (100 samp)

CEM (20 samp)

RWR (100 samp)

RWR (20 samp)

ours (20 samp)

ours (with GMM, 5 samp)

# Concluding Comments

- simple <u>linear</u> dynamics model

- fast, simple, standard LQR solver

- can handle contacts despite linear model

- fit very complex policies with guided policy search