

# RL Part 3.2: Probabilistic Policy Search

---

**Jan Peters**  
**Gerhard Neumann**



# What we have seen from the policy gradients

---

- Policy Search is a powerful and practical alternative to value function and model-based methods.
- Policy gradients have dominated this area for a long time and solidly working methods exist.
- They still need a lot of samples and **we need to tune the learning rate**
- Learning the **exploration rate** is still an open problem



# Outline of the Lecture

---

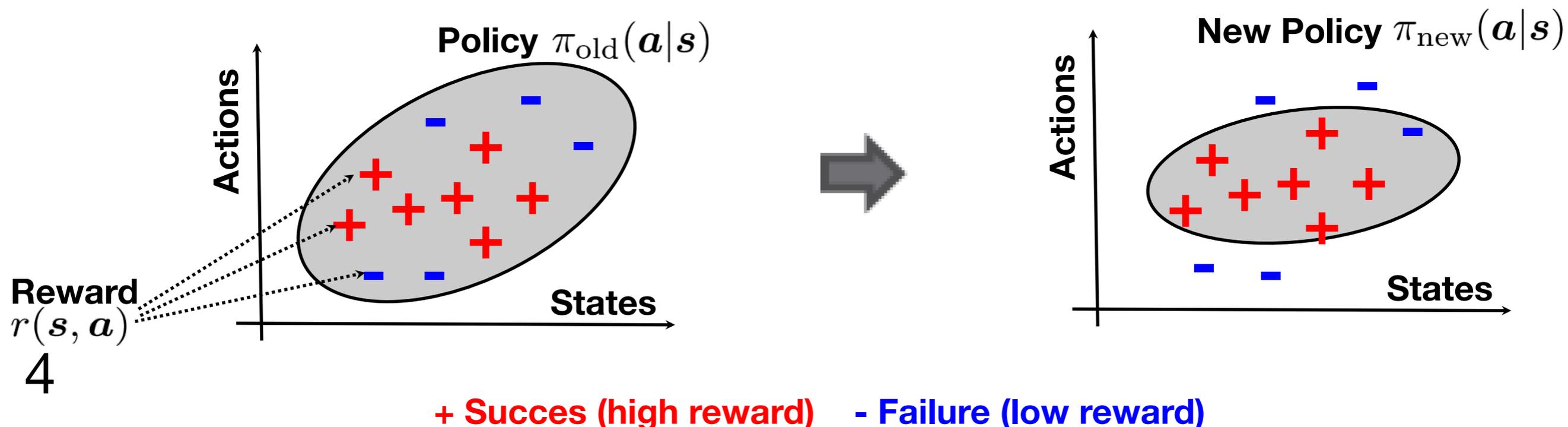
1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
3. Relative Entropy Policy Search (REPS)
4. REPS for Contextual Policy Search
5. Learning Versatile Solutions
6. Sequencing Movement Primitives

# Success Matching Principle



“When learning from a set of their own trials in iterated decision problems, humans attempt to match **not the best taken action** but the **reward-weighted frequency** of their actions and outcomes” (Arrow, 1958).

**Success-Matching:** Policy update learn to reproduce successful outcomes



# Episode-Based Success Matching



## Iterate:

**Sample and evaluate** parameters:

$$\boldsymbol{\theta}^{[i]} \sim \pi(\boldsymbol{\theta}; \boldsymbol{\omega}_k) \quad R^{[i]} = \sum_{t=1}^T r_t^{[i]}$$

**Compute „success probability“** for each sample

$$w^{[i]} = f(R^{[i]})$$

➔ transform reward in a **non-negative weight** (improper probability distribution)

**Compute „Success“ weighted** policy on the samples

$$p_k(\boldsymbol{\theta}^{[i]}) \propto w^{[i]} \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega}_k)$$

**Fit new parametric policy**  $\pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega}_{k+1})$  that best approximates  $p_k(\boldsymbol{\theta}^{[i]})$



# Episode-Based Success Matching

---

## 2 Open issues:

**How to fit the policy**  $\pi(\theta^{[i]}; \omega_{k+1})$  ?

**How to compute**  $w^{[i]} = f(R^{[i]})$  ?



# Outline of the Lecture

---

1. Introduction

**2. Policy Updates by Weighted Maximum Likelihood**

3. Relative Entropy Policy Search (REPS)

4. REPS for Contextual Policy Search

5. Conclusion

# Policy Fitting



**Problem:** We want to find a parametric distribution  $\pi(\boldsymbol{\theta}; \boldsymbol{\omega}_{k+1})$  that best fits the distribution  $p(\boldsymbol{\theta}^{[i]}) \propto w^{[i]} \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega}_k)$ ,

➔ **We can do that by minimizing:**

$$\boldsymbol{\omega}_{k+1} = \operatorname{argmin}_{\boldsymbol{\omega}} \operatorname{KL}(p(\boldsymbol{\theta}^{[i]}) || \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega}))$$

$$= \operatorname{argmin}_{\boldsymbol{\omega}} \int p(\boldsymbol{\theta}) \log \frac{p(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta}; \boldsymbol{\omega})} d\boldsymbol{\theta}$$

$$\approx \operatorname{argmax}_{\boldsymbol{\omega}} \sum_i \frac{p(\boldsymbol{\theta}^{[i]})}{\pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega}_k)} \log \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega})$$

We sampled from the old policy

$$\boldsymbol{\omega}_{k+1} = \operatorname{argmax}_{\boldsymbol{\omega}} \sum_i w^{[i]} \log \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega})$$

➔ The fitting of the policy is obtained by a **weighted maximum likelihood estimate**

➔ Closed form solutions exists, no learning rates

# Weighted Maximum Likelihood Solutions...



**For a Gaussian policy:**  $\pi(\boldsymbol{\theta}; \boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\boldsymbol{\mu} = \frac{\sum_i w^{[i]} \boldsymbol{\theta}^{[i]}}{\sum_i w^{[i]}}$$

Weighted mean

$$\boldsymbol{\Sigma} = \frac{\sum_i w^{[i]} (\boldsymbol{\theta}^{[i]} - \boldsymbol{\mu})(\boldsymbol{\theta}^{[i]} - \boldsymbol{\mu})^T}{\sum_i w^{[i]}}$$

Weighted covariance

**But more general:** Also for mixture models, GPs and so on...

# Difference to policy gradients



## Weighted Maximum Likelihood:

$$\omega_{k+1} = \operatorname{argmax}_{\omega} \sum_i w^{[i]} \log \pi(\theta^{[i]}; \omega)$$

**I.e.:** Set  $\nabla_{\omega} \sum_i w^{[i]} \log \pi(\theta^{[i]}; \omega) = \sum_i \nabla_{\omega} \log \pi(\theta^{[i]}; \omega) w^{[i]} = 0$

Solve in closed form for  $\omega$

## Policy Gradients:

$$\nabla_{\omega} J_{\omega} = \sum_{i=1}^N \nabla_{\omega} \log \pi(\theta_i; \omega_k) R_i, \quad \omega_{k+1} = \omega_k + \alpha \nabla_{\omega} J_{\omega}$$

# Computing the weights...



So **where are the weights**  $w^{[i]} = f(R^{[i]})$  coming from?

We need to transform the returns in an **improper probability distribution**

**Simple Way: Exponential transformation**  $w^{[i]} = \exp(\beta(R^{[i]} - \max R^{[i]}))$

$\beta$  ... temperature of the distribution

Often set by heuristics, e.g.:  $\beta = \frac{10}{\max R^{[i]} - \min R^{[i]}}$

**Can be justified from different view-points**

**EM-Algorithms: PoWER, Reward-Weighted Regression**

**Optimal Control: PI2**

# Expectation Maximization

---



## Some notes on Expectation-Maximization in this context

EM is a method for Max. Likelihood in the case of latent (unobserved) variables

**Observed variable:** Reward Event  $C = 1$  (we want to get reward)

**Unobserved variable:** Trajectory  $\tau$  (or parameters) that created the reward event

**E-Step:** Estimate new desired distribution

**M-Step:** Estimate new policy parameters from the weighted samples

**Step-based Policy Search Versions of EM:** PoWER (Kober 2008), Reward-Weighted Regression (Peters 2007)

# Exponential Transformation

---



## Some notes on the exponential transformation

In stochastic environments, we do not optimize the expected reward any more as...

$$\mathbb{E}_{p(\tau)} [\exp(R(\tau))] \neq \exp(\mathbb{E}_{p(\tau)} [R(\tau)])$$

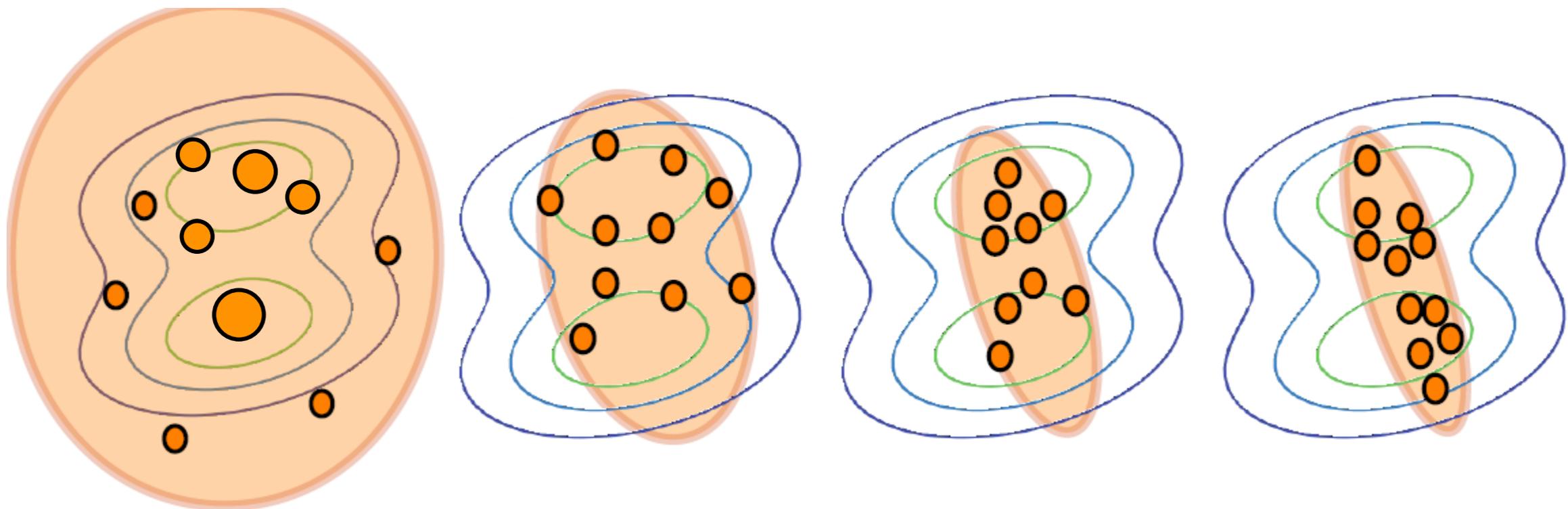
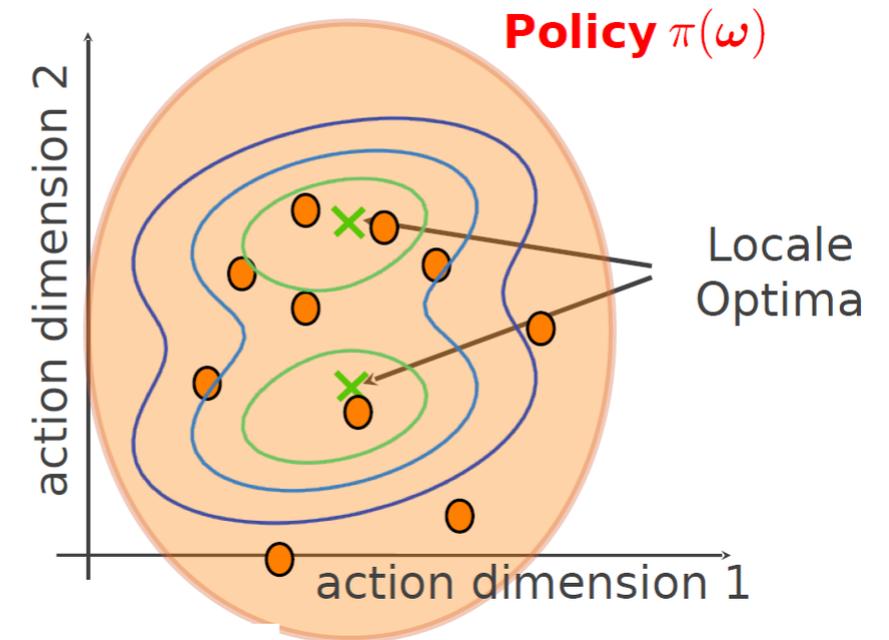
The objective gets „risk attracted“

For moderately stochastic environments it still works well

# Illustration on weighted ML



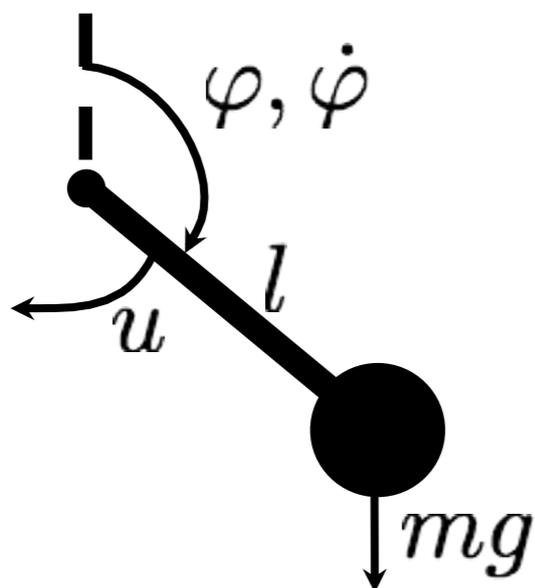
Example for a 2D parameter space:



# Underactuated Swing-Up



- swing heavy pendulum up



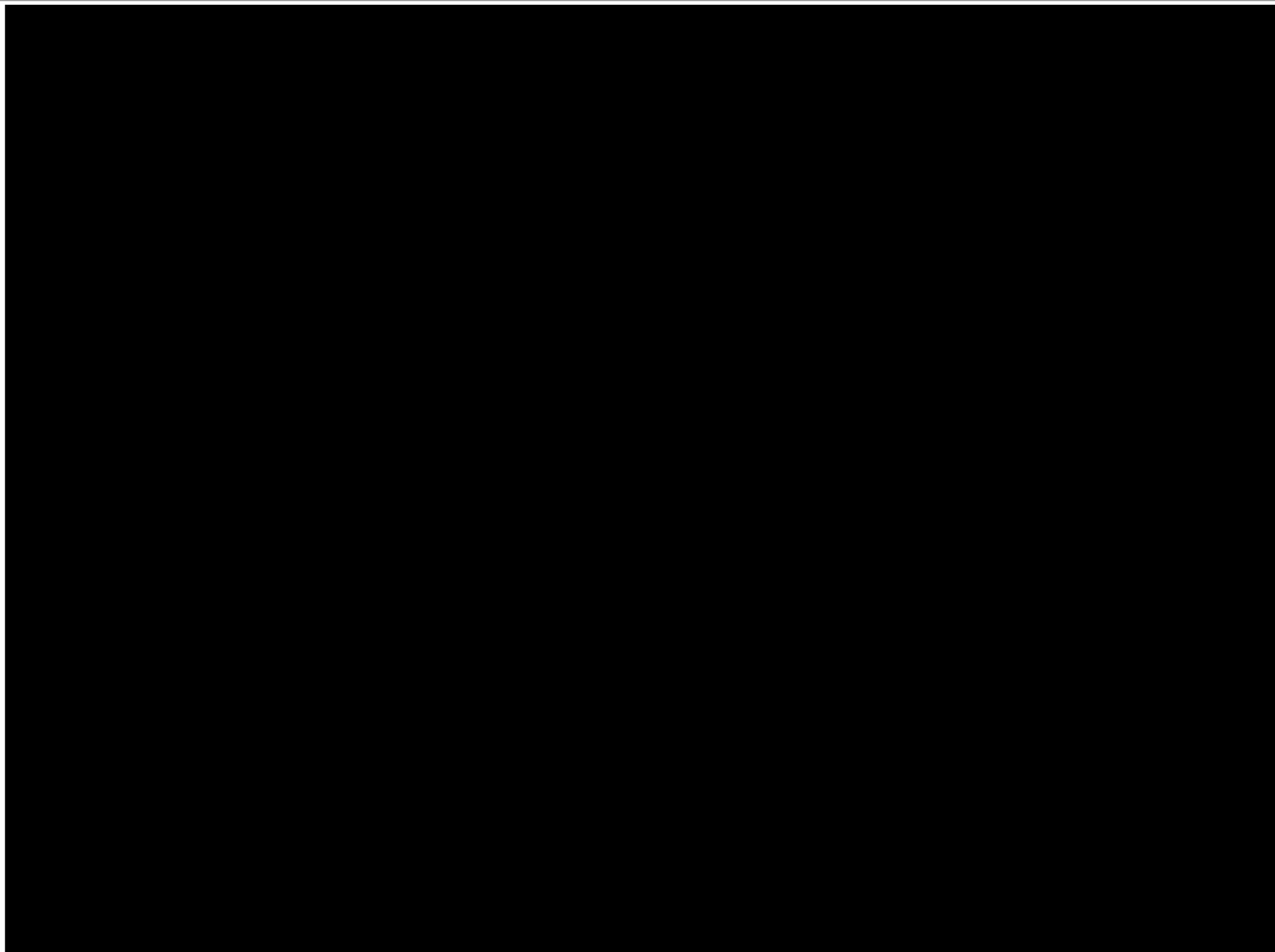
$$ml^2\ddot{\varphi} = -\mu\dot{\varphi} + mgl \sin \varphi + u$$
$$\varphi \in [-\pi, \pi]$$

- motor torques limited, Policy: DMPs

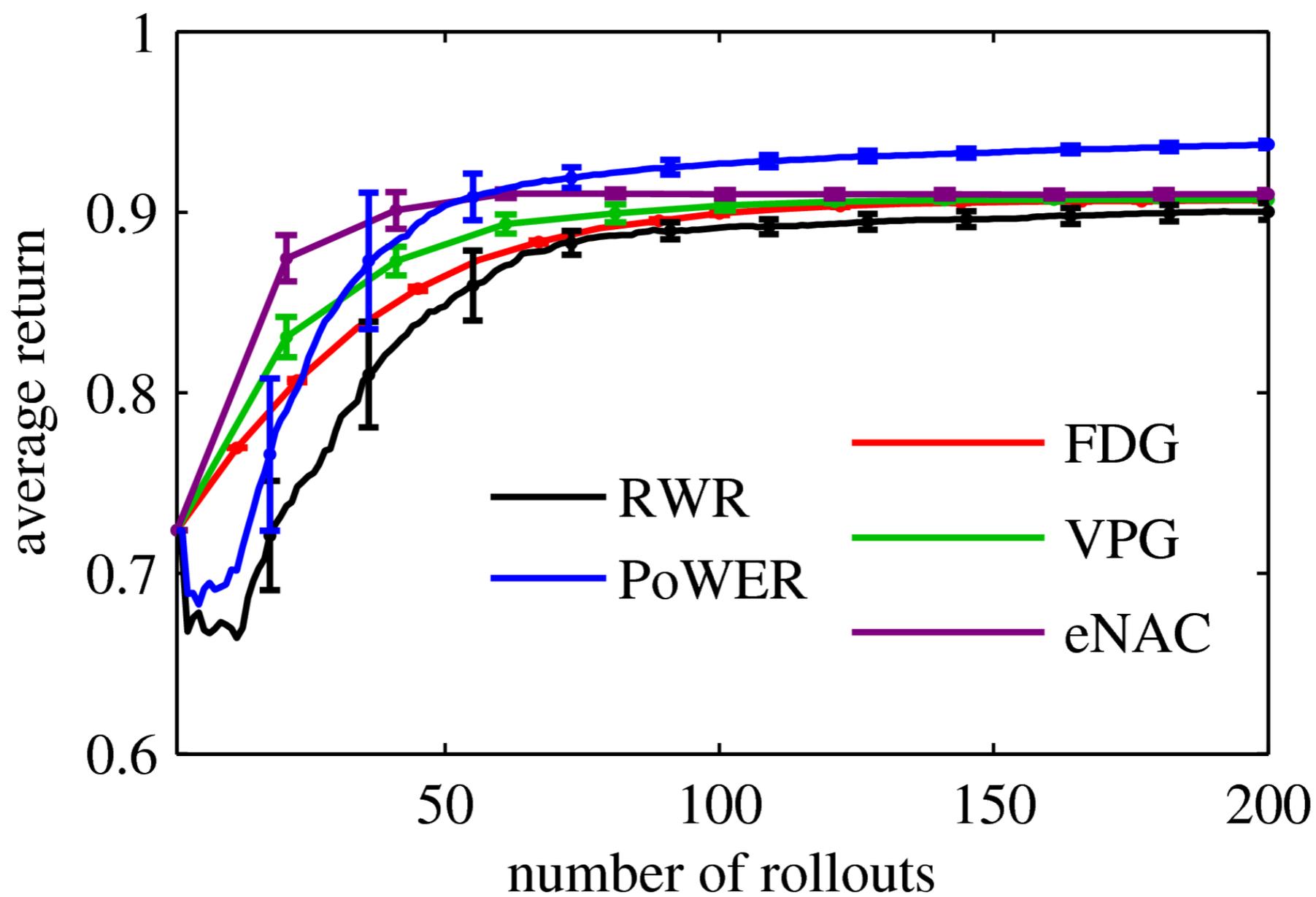
$$|u| \leq u_{max}$$

- reward function

$$r = \exp \left( -\alpha \left( \frac{\varphi}{\pi} \right)^2 - \beta \left( \frac{2}{\pi} \right)^2 \log \cos \left( \frac{\pi}{2} \frac{u}{u_{max}} \right) \right)$$

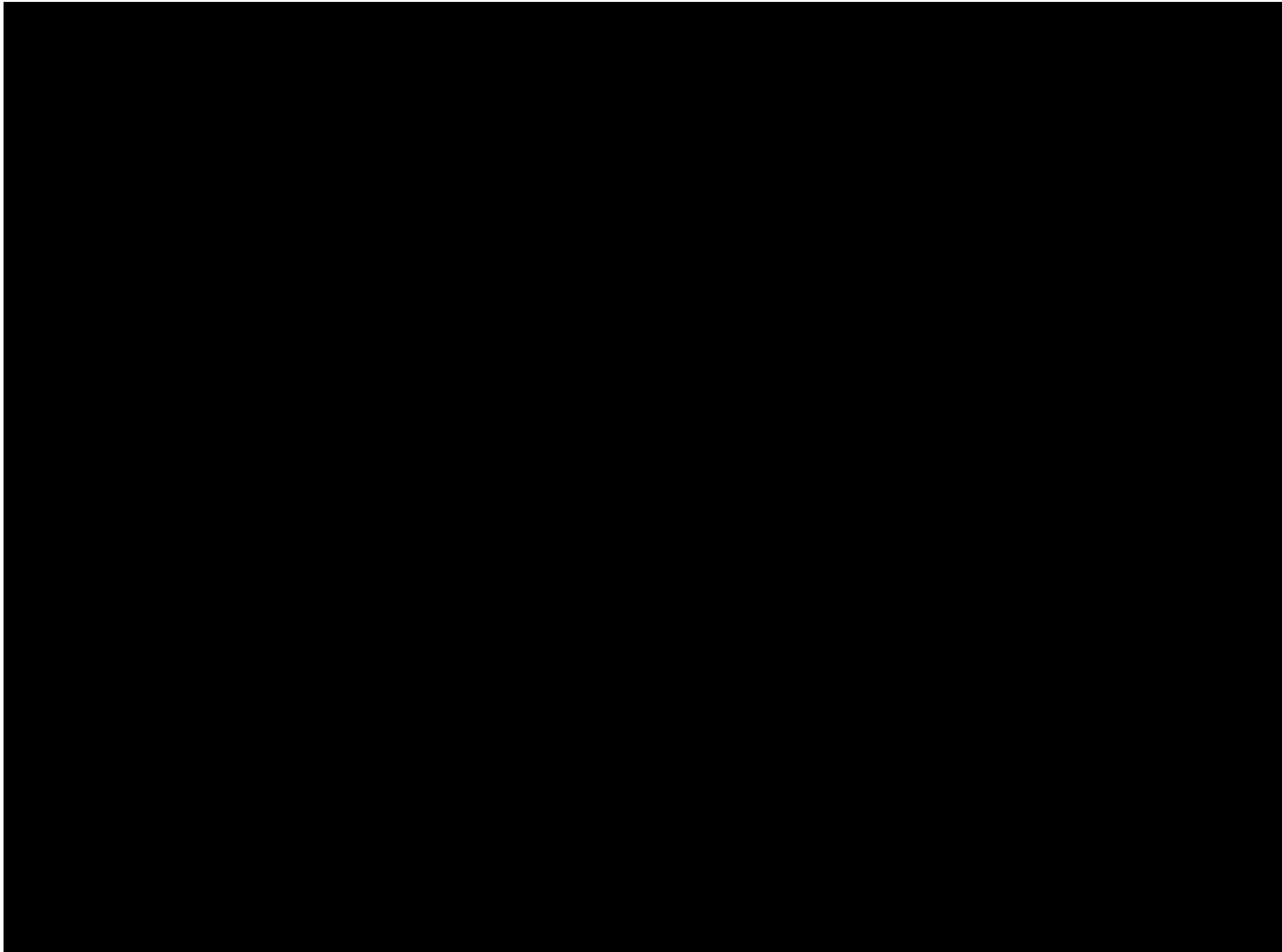


# Underactuated Swing-Up



# Ball in the Cup

---



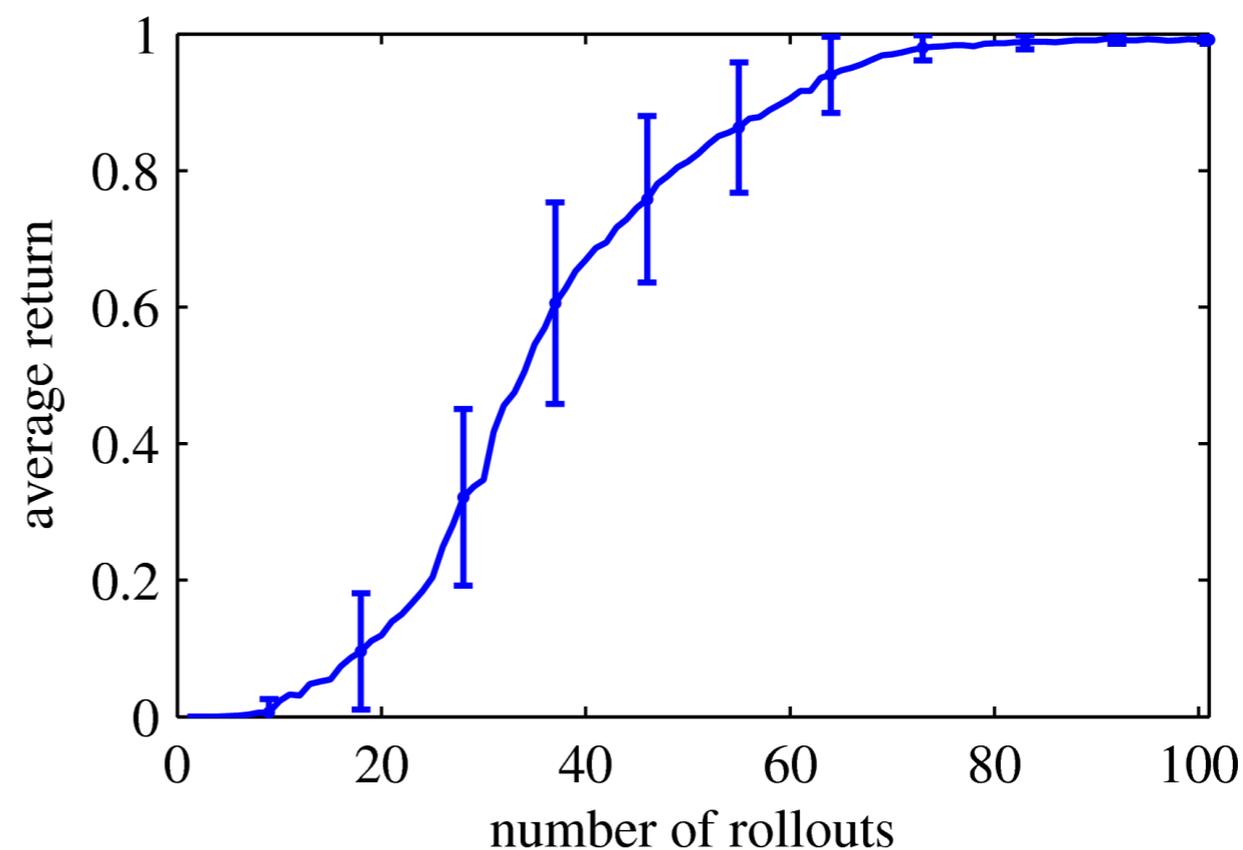
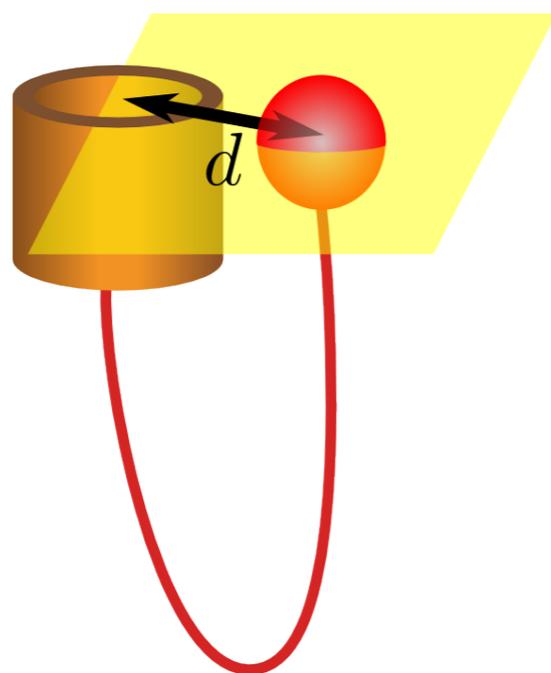
# Ball-in-a-Cup



## Reward function:

$$r_t = \begin{cases} \exp\left(-\alpha\left((x_c - x_b)^2 + (y_c - y_b)^2\right)\right) & \text{if } t = t_c \\ 0 & \text{if } t \neq t_c \end{cases}$$

## Policy: DMPs





# Outline of the Lecture

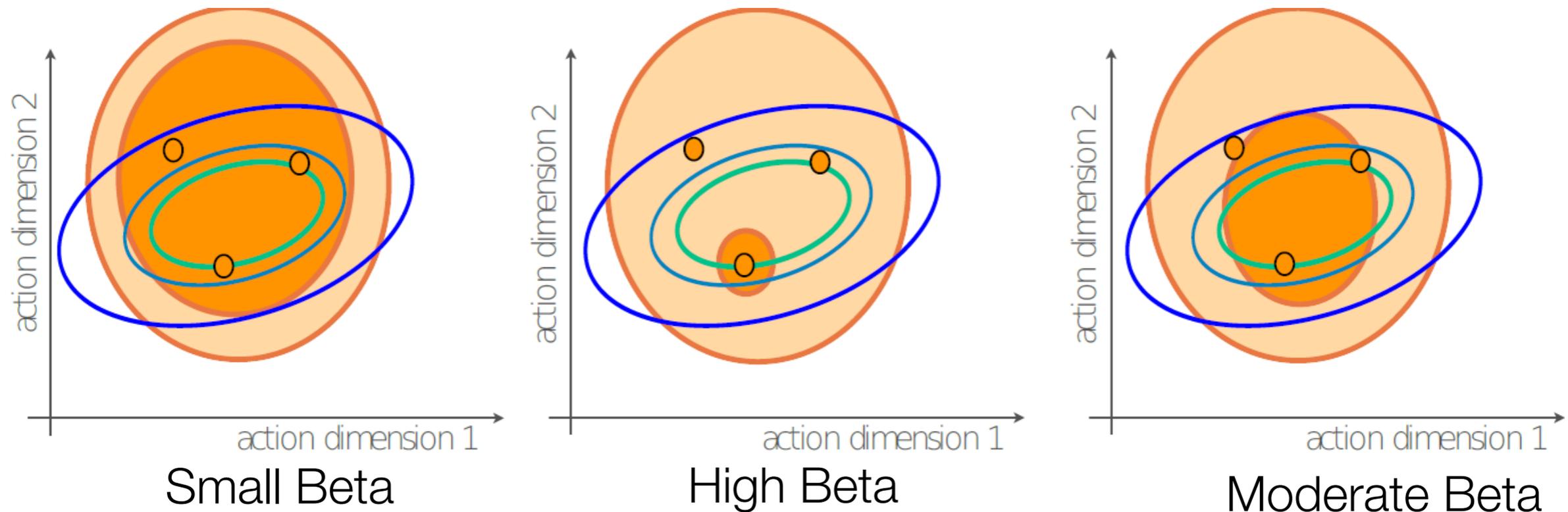
---

1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
- 3. Relative Entropy Policy Search (REPS)**
4. REPS for Contextual Policy Search
5. Learning Versatile Solutions
6. Sequencing Movement Primitives

# Policy Search: Choosing the step size



**What is a good desired distribution for the policy update?**



How can we choose the **exploration-exploitation tradeoff**?

**Again use a metric to control the step-size of the update**



# Relative Entropy Policy Search

## Relative entropy as metric between two policies

$$\text{KL}(\pi(\boldsymbol{\theta}) || q(\boldsymbol{\theta})) \leq \epsilon$$

We get the following optimization problem:

$$\max_{\pi} \sum_i \pi(\boldsymbol{\theta}^{[i]}) R(\boldsymbol{\theta}^{[i]}) \quad \text{Maximize Reward}$$

$$\text{s.t: } \sum_i \pi(\boldsymbol{\theta}^{[i]}) = 1 \quad \text{It's a distribution}$$

$$\text{KL}(\pi(\boldsymbol{\theta}) || q(\boldsymbol{\theta})) \leq \epsilon \quad \text{Stay close to the old policy } q(\boldsymbol{\theta})$$

**Policy Update is formulated as constraint optimization problem**

# Relative Entropy Policy Search



**We get the following optimization problem:**

$$\max_{\pi} \sum \pi(\boldsymbol{\theta}^{[i]}) R(\boldsymbol{\theta}^{[i]}) \quad \text{Maximize Reward}$$

$$\text{s.t: } \sum_i \pi(\boldsymbol{\theta}^{[i]}) = 1 \quad \text{It's a distribution}$$

$$\text{KL}(\pi(\boldsymbol{\theta}) || q(\boldsymbol{\theta})) \leq \epsilon \quad \text{Stay close to the old policy } q(\boldsymbol{\theta})$$

**Which has the following **analytic solution**:**

$$\pi(\boldsymbol{\theta}) \propto q(\boldsymbol{\theta}) \exp\left(\frac{\mathcal{R}_{\boldsymbol{\theta}}}{\eta}\right)$$

Thats exactly success matching with exponential transformation!

**Scalingfactor  $\eta$  :**

- **Automatically chosen from optimization** (Lagrange Multiplier)

- Specified by KL-bound  $\epsilon$



# Getting the Lagrangian multipliers

How to get  $\eta$  :

Solve **dual optimization** problem:

**Dual function:** 
$$h(\eta) = \eta\epsilon + \eta \log \int q(\boldsymbol{\theta}) \exp\left(\frac{\mathcal{R}\boldsymbol{\theta}}{\eta}\right) d\boldsymbol{\omega}$$

$$= \eta\epsilon + \eta \log \sum_{i=1}^N \frac{1}{N} \exp\left(\frac{\mathcal{R}^{[i]}}{\eta}\right)$$

**Minimize:**  $\eta^* = \operatorname{argmin}_{\eta} h(\eta) \quad \text{s.t: } \eta > 0$

**Log-sum-exp** softmax structure

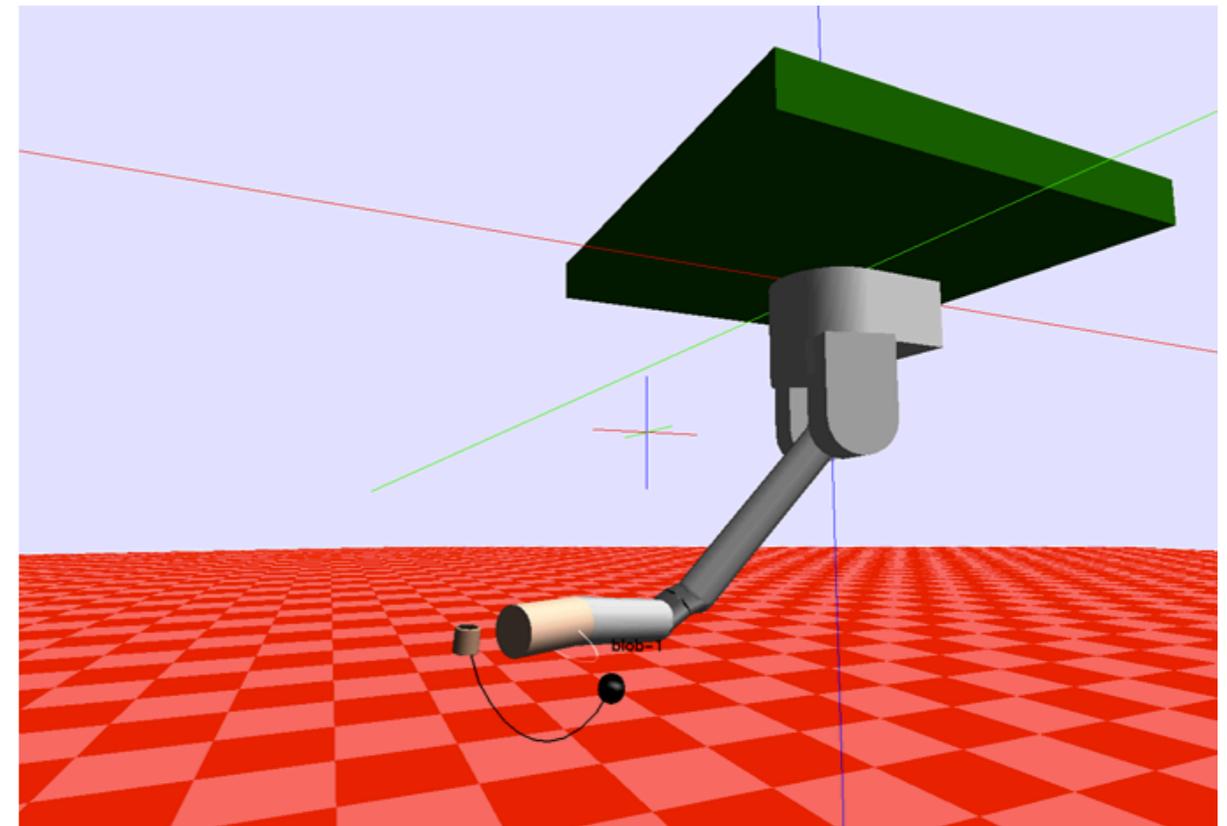
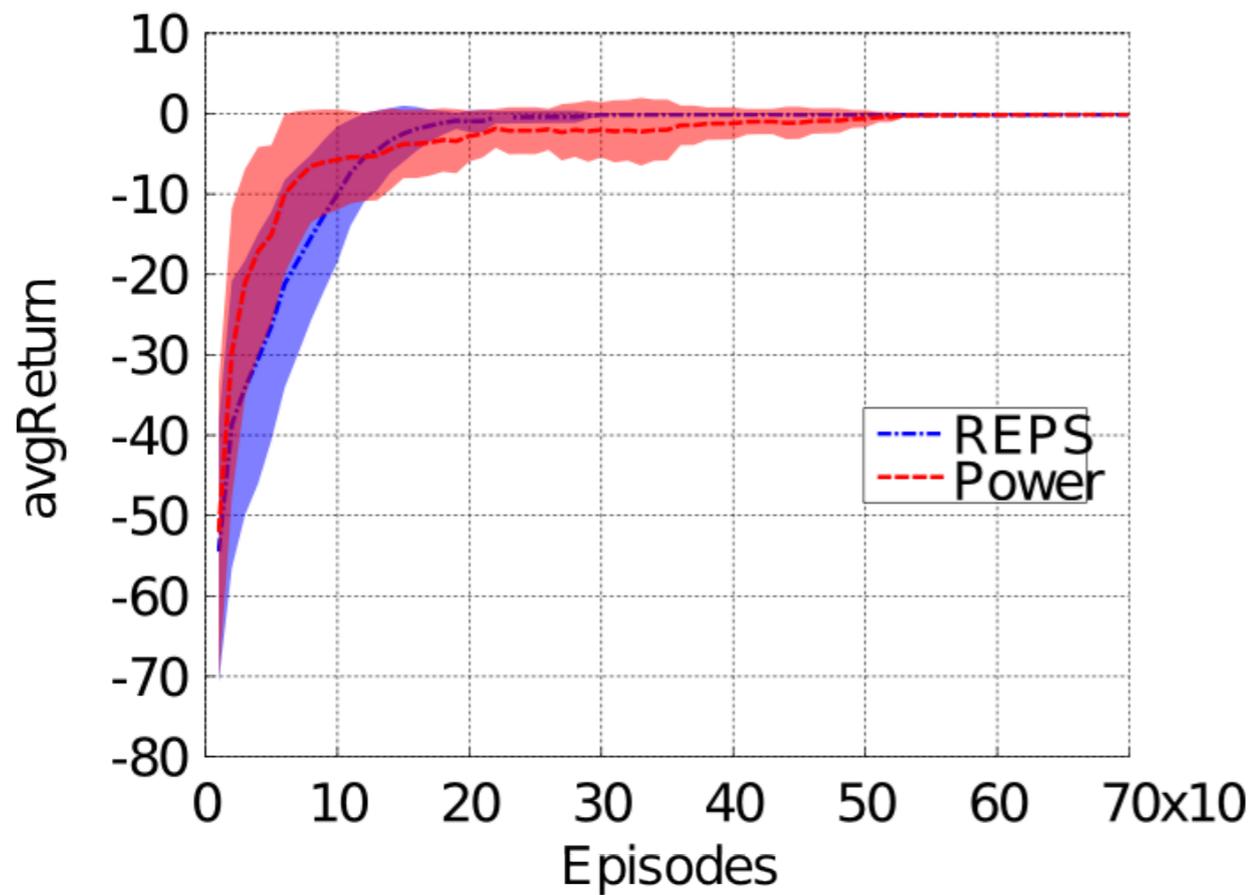
Optimized by standard optimization tools

(e.g. trust region algorithms)



# Results

## Comparison on simulated Ball In The Cup





# Outline of the Lecture

---

1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
3. Relative Entropy Policy Search (REPS)
4. **REPS for Contextual Policy Search**
5. Learning Versatile Solutions
6. Sequencing Movement Primitives

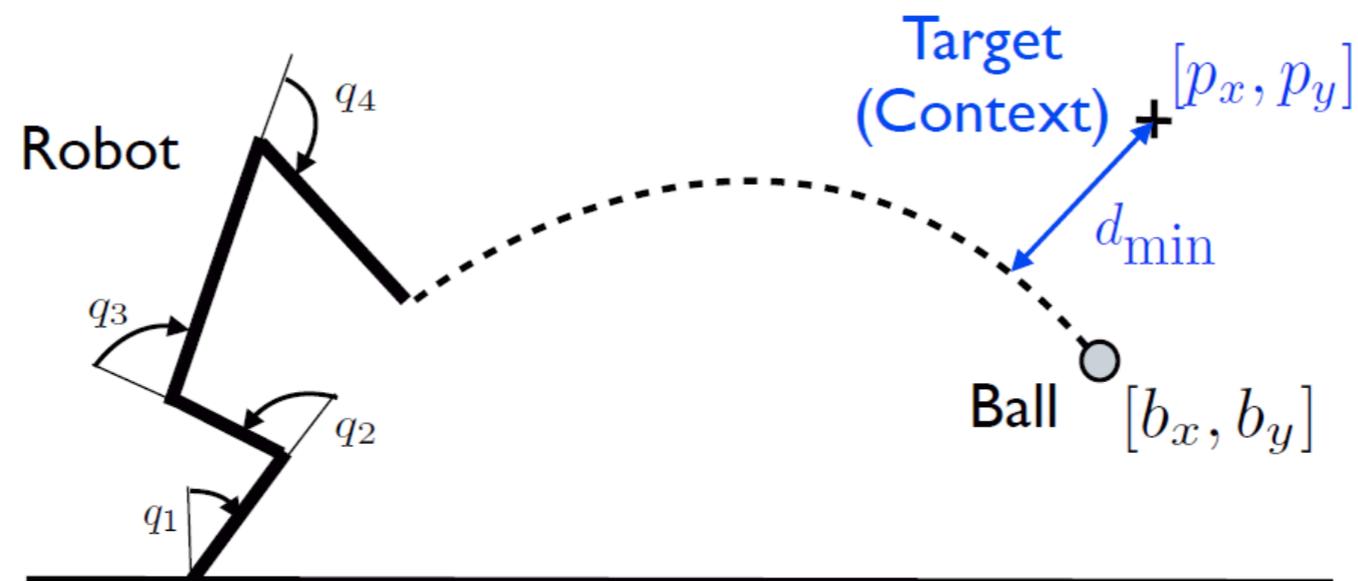


# Contextual Policy Search

## Contextual Policy Search

Context  $x$  describes objectives of the task (fixed before task execution)

E.g.: Target location to throw a ball



# Contextual Policy Search



## Contextual Policy Search

Context  $\mathbf{x}$  describes objectives of the task (fixed before task execution)

E.g.: Target location to throw a ball

We now want to learn an upper level policy  $\pi(\boldsymbol{\theta}|\mathbf{x};\boldsymbol{\omega})$  that adapts  $\boldsymbol{\theta}$  to the context

Data-set used for policy update

$$\mathcal{D}_{\text{episode}} = \{ \boldsymbol{\theta}^{[i]}, \mathbf{x}^{[i]}, R^{[i]} \}_{i=1\dots N}$$

**Goal:** maximize expected reward

$$J_{\pi} = \iint \mu_0(\mathbf{x}) \pi(\boldsymbol{\theta}|\mathbf{x}) \mathcal{R}_{\mathbf{x}\boldsymbol{\theta}} d\mathbf{x} d\boldsymbol{\theta}$$



# Contextual Policy Search

Optimize over the joint distribution:

$$p(\mathbf{x}, \boldsymbol{\theta}) = \mu(\mathbf{x})\pi(\boldsymbol{\theta}|\mathbf{x})$$

$$\max_p \sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) R(\mathbf{x}, \boldsymbol{\theta}) \quad \text{Maximize Reward}$$

$$\sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) = 1 \quad \text{It's a distribution}$$

$$\text{KL}(p(\mathbf{x}, \boldsymbol{\theta}) || q(\mathbf{x}, \boldsymbol{\theta})) \leq \epsilon \quad \text{Stay close to the data}$$

$$\forall \mathbf{x} \quad p(\mathbf{x}) = \sum_{\boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) = \mu_0(\mathbf{x}) \quad \text{Reproduce given context distribution}$$

## Problems:

- Context distribution can not be freely chosen by the algorithm
- Infinite amount of constraints
- For each context, we need many parameter vector samples

# Matching Feature Averages



$$\max_p \sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) R(\mathbf{x}, \boldsymbol{\theta})$$

Maximize Reward

$$\sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) = 1$$

It's a distribution

$$\text{KL}(\pi(\boldsymbol{\theta}|\mathbf{x})\mu(\mathbf{x})||q(\mathbf{x}, \boldsymbol{\theta})) \leq \epsilon$$

Stay close to the data

$$\sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) \phi(\mathbf{x}) = \hat{\phi}$$

**Reproduce given context  
feature averages**

Instead of matching the context distribution exactly,  
we can match only certain feature averages (moments) of the distribution

# Matching Feature Averages

---



$$\sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) \boldsymbol{\phi}(\mathbf{x}) = \hat{\boldsymbol{\phi}}$$

Reproduce given context  
feature averages

What does that mean? Example:

$$\boldsymbol{\phi}(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$

- ➔ Match first and second order moment
- ➔ Equivalent to matching mean and variance
- ➔ Exact for Gaussian distributions



# Matching Feature Averages

$$\max_p \sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) R(\mathbf{x}, \boldsymbol{\theta})$$

Maximize Reward

$$\sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) = 1$$

It's a distribution

$$\text{KL}(\pi(\boldsymbol{\theta}|\mathbf{x})\mu(\mathbf{x})||q(\mathbf{x}, \boldsymbol{\theta})) \leq \epsilon$$

Stay close to the data

$$\sum_{\mathbf{x}, \boldsymbol{\theta}} p(\mathbf{x}, \boldsymbol{\theta}) \phi(\mathbf{x}) = \hat{\phi}$$

Reproduce given context feature averages

Closed form solution:

$$\mu(\mathbf{x})\pi(\boldsymbol{\theta}|\mathbf{x}) \propto q(\mathbf{x}, \boldsymbol{\theta}) \exp\left(\frac{R_{\mathbf{x}\boldsymbol{\theta}} - V(\mathbf{x})}{\eta}\right)$$

We automatically get a baseline for the returns

$$V(\mathbf{x}) = \phi^T(\mathbf{x})\mathbf{v}$$

# Matching Feature Averages



$$\max_p \sum_i p(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]}) R(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]})$$

Maximize Reward

$$\sum_i p(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]}) = 1$$

It's a distribution

$$\text{KL}(\pi(\boldsymbol{\theta}|\mathbf{x})\mu(\mathbf{x})||q(\mathbf{x}, \boldsymbol{\theta})) \leq \epsilon$$

Stay close to the data

Reproduce given context  
feature averages

$$\text{e.g., } \phi(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$

Match Mean and Variance

$$\mu(\mathbf{x})\pi(\boldsymbol{\theta}|\mathbf{x}) \propto q(\mathbf{x}, \boldsymbol{\theta}) \exp\left(\frac{R_{\mathbf{x}\boldsymbol{\theta}} - V(\mathbf{x})}{\eta}\right)$$

# Matching Feature Averages



$$\max_p \sum_i p(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]}) R(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]})$$

Maximize Reward

$$\sum_i p(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]}) = 1$$

It's a distribution

$$\text{KL}(\pi(\boldsymbol{\theta}|\mathbf{x})\mu(\mathbf{x})||q(\mathbf{x}, \boldsymbol{\theta})) \leq \epsilon$$

Stay close to the data

$$\sum_{\mathbf{x}} p(\mathbf{x})\phi(\mathbf{x}) = \hat{\phi}$$

Reproduce given context  
feature averages

$$\text{e.g., } \phi(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$

Match Mean and Variance

$$\mu(\mathbf{x})\pi(\boldsymbol{\theta}|\mathbf{x}) \propto q(\mathbf{x}, \boldsymbol{\theta}) \exp\left(\frac{R_{\mathbf{x}\boldsymbol{\theta}} - V(\mathbf{x})}{\eta}\right)$$

# Getting the Lagrangian multipliers



## How to get $\eta, \mathbf{v}$

Solve **dual optimization** problem:

**Dual function:**

$$h(\eta, \mathbf{v}) = \eta\epsilon + \hat{\phi}^T \mathbf{v} + \eta \log \sum_i \frac{1}{N} \exp \left( \frac{\mathcal{R}^{[i]} - \phi^T(\mathbf{x}^{[i]})\mathbf{v}}{\eta} \right)$$

**Minimize:**  $[\eta^*, \mathbf{v}^*] = \operatorname{argmin}_{\eta} h(\eta, \mathbf{v}) \quad \text{s.t: } \eta > 0$

**Integral is over the context-parameter space**

We can use  $(\mathbf{x}^{[i]}, \boldsymbol{\theta}^{[i]})$  samples instead of many samples  $\boldsymbol{\theta}^{[ij]}$  per context  $\mathbf{x}^{[i]}$



# Contextual Policies with weighted ML

**Estimate parametric policy**  $\pi_{\omega}(\boldsymbol{\theta}|\boldsymbol{x})$  :

If  $\pi_{\omega}(\boldsymbol{\theta}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{K}\boldsymbol{x} + \boldsymbol{k}, \boldsymbol{\Sigma}_{\omega})$  is Gaussian:

$$\begin{bmatrix} \boldsymbol{k}^T \\ \mathbf{K}^T \end{bmatrix} = (\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{A}$$

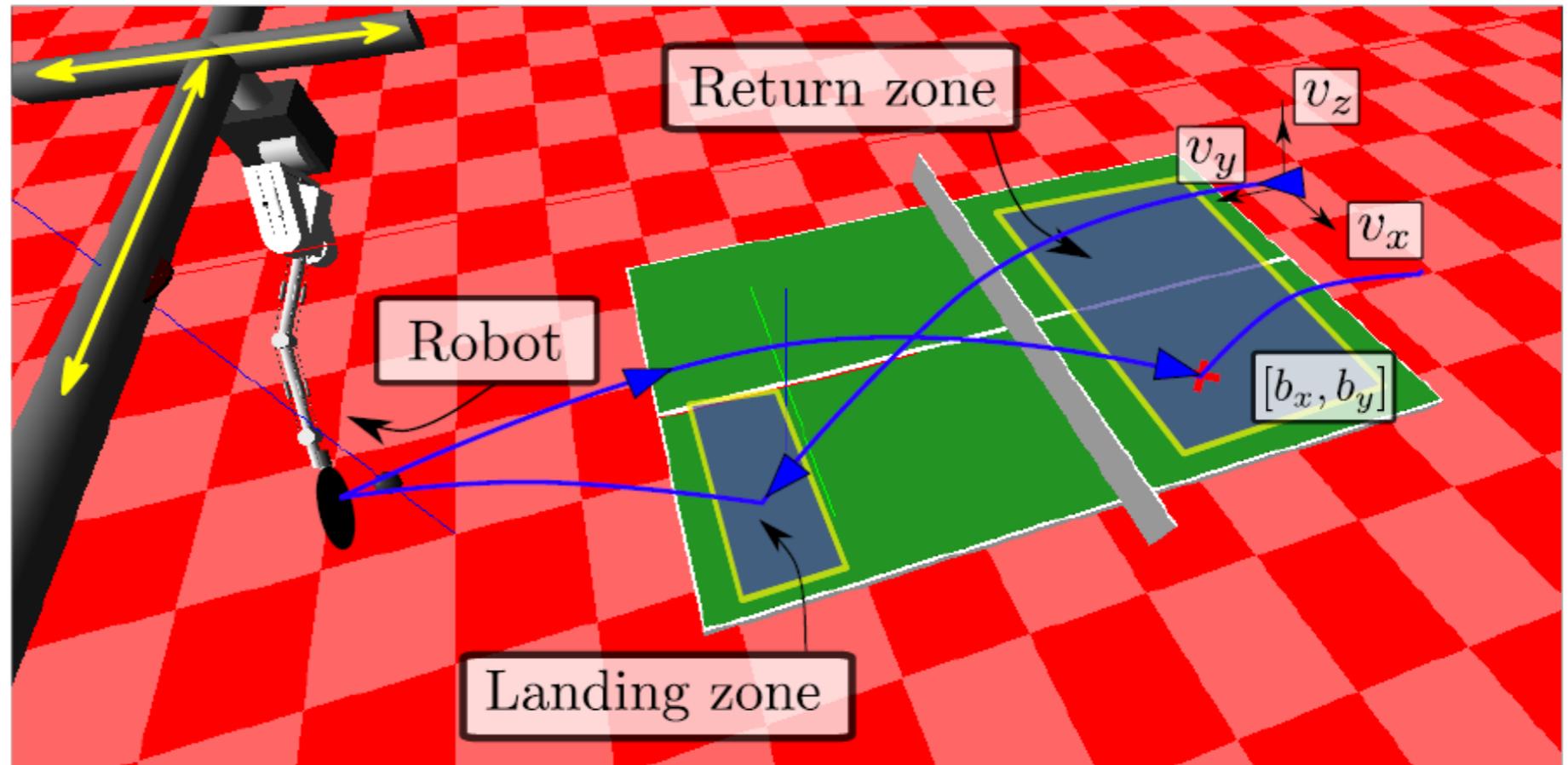
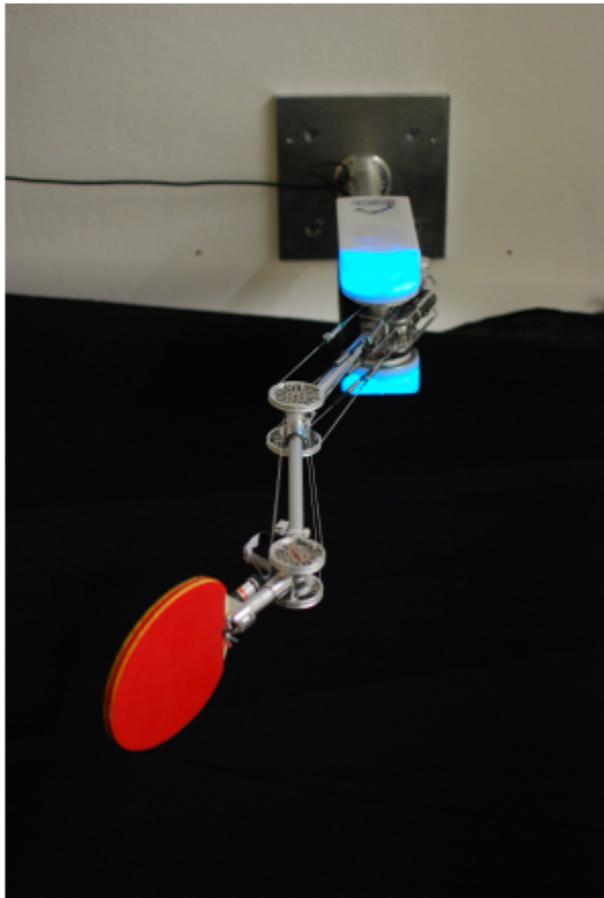
$$\boldsymbol{\mu}_i = \boldsymbol{k} + \mathbf{K} \boldsymbol{x}_i \quad \boldsymbol{\Sigma} = \frac{\sum_i p_i (\boldsymbol{\theta}^{[i]} - \boldsymbol{\mu}_i)(\boldsymbol{\theta}^{[i]} - \boldsymbol{\mu}_i)^T}{\sum_i p_i}$$

- X ... input data matrix (including 1 for the bias)
- D ... diagonal weighting matrix
- A .... Parameter matrix

Just standard **weighted linear regression**...

# Table tennis experiments

---



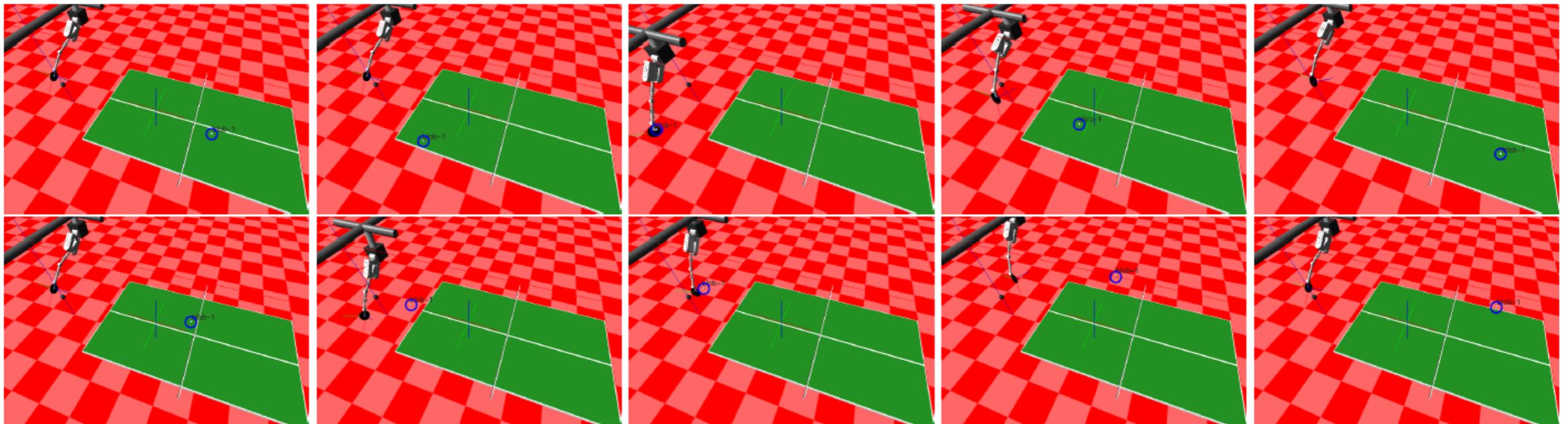
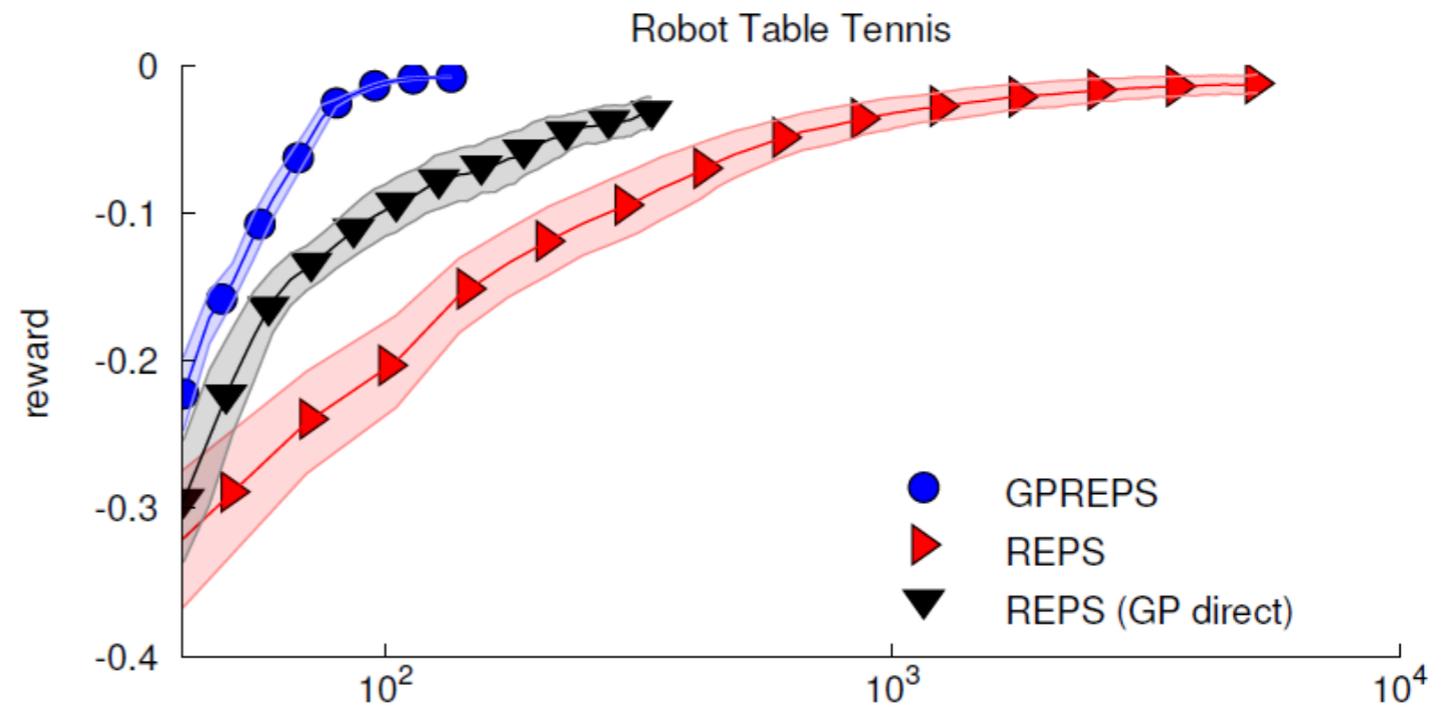
[Kupscik, Neumann et al, submitted, 2013]

# Table tennis experiments

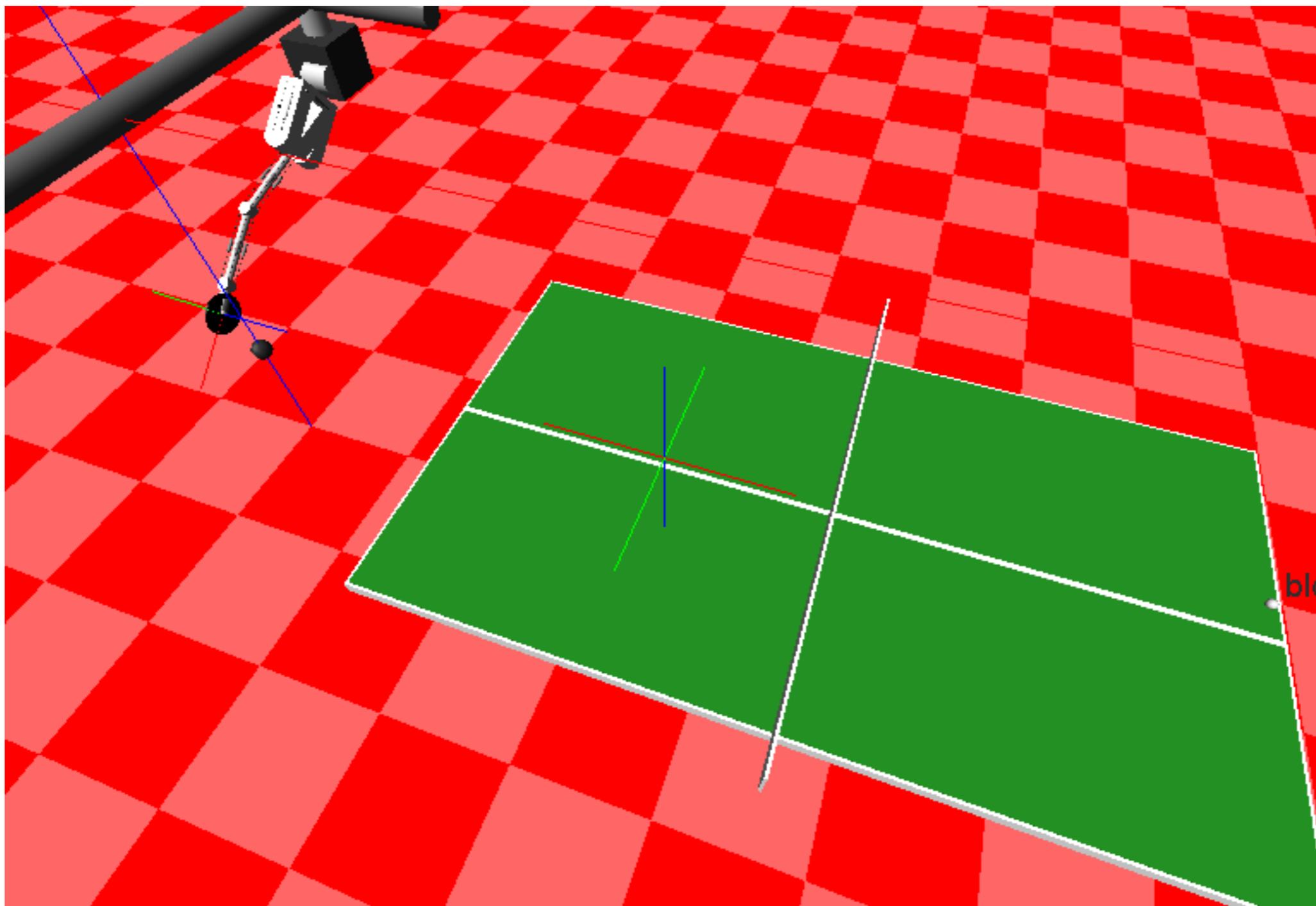


## REPS with learned forward models

- Complex behavior can be learned within 100 episodes



# Table tennis experiments





# Outline of the Lecture

---

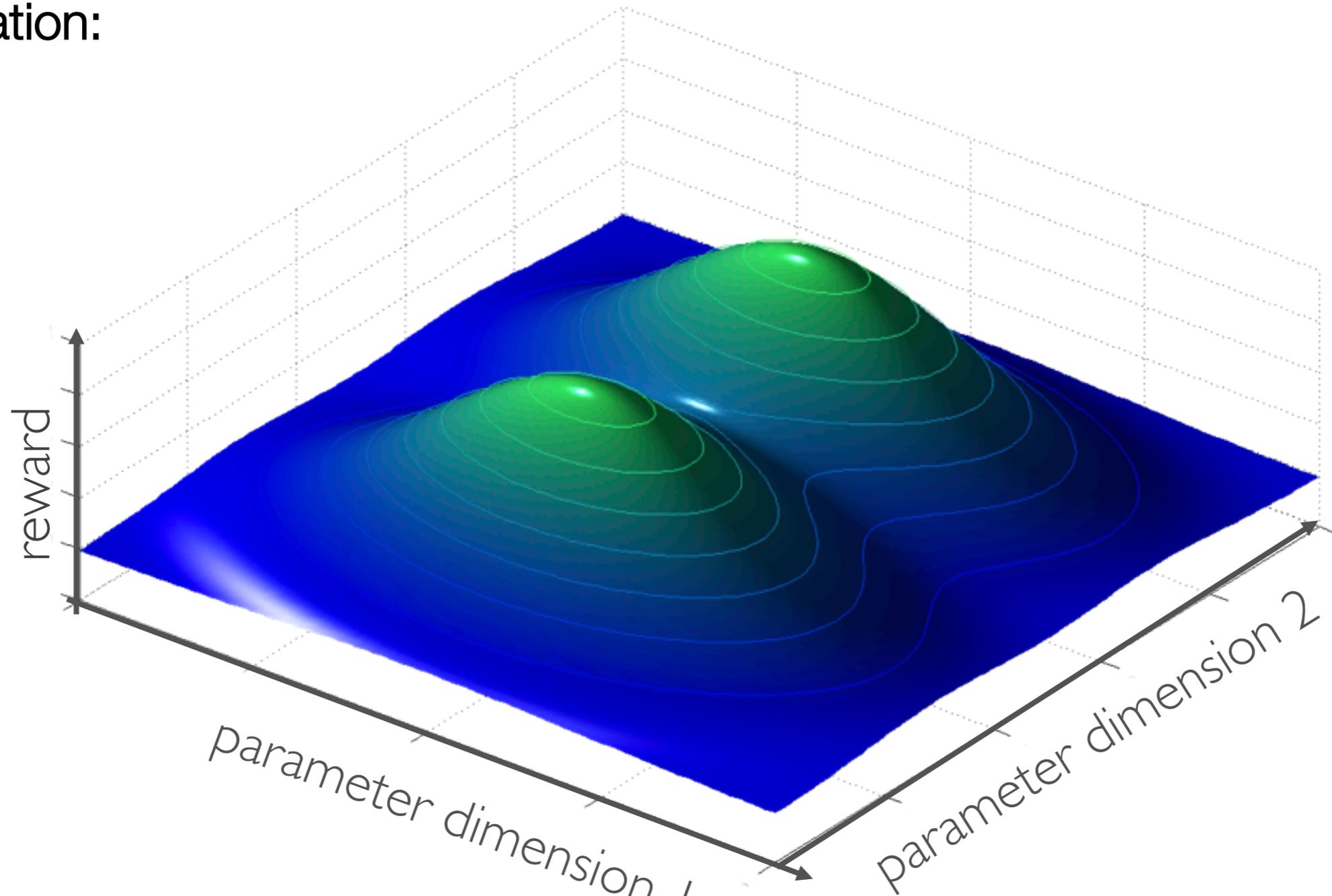
1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
3. Relative Entropy Policy Search (REPS)
4. REPS for Contextual Policy Search
5. **Learning Versatile Solutions**
6. Sequencing Movement Primitives

# Versatile Solutions: Illustration

Many motor-tasks have **multiple solutions**:

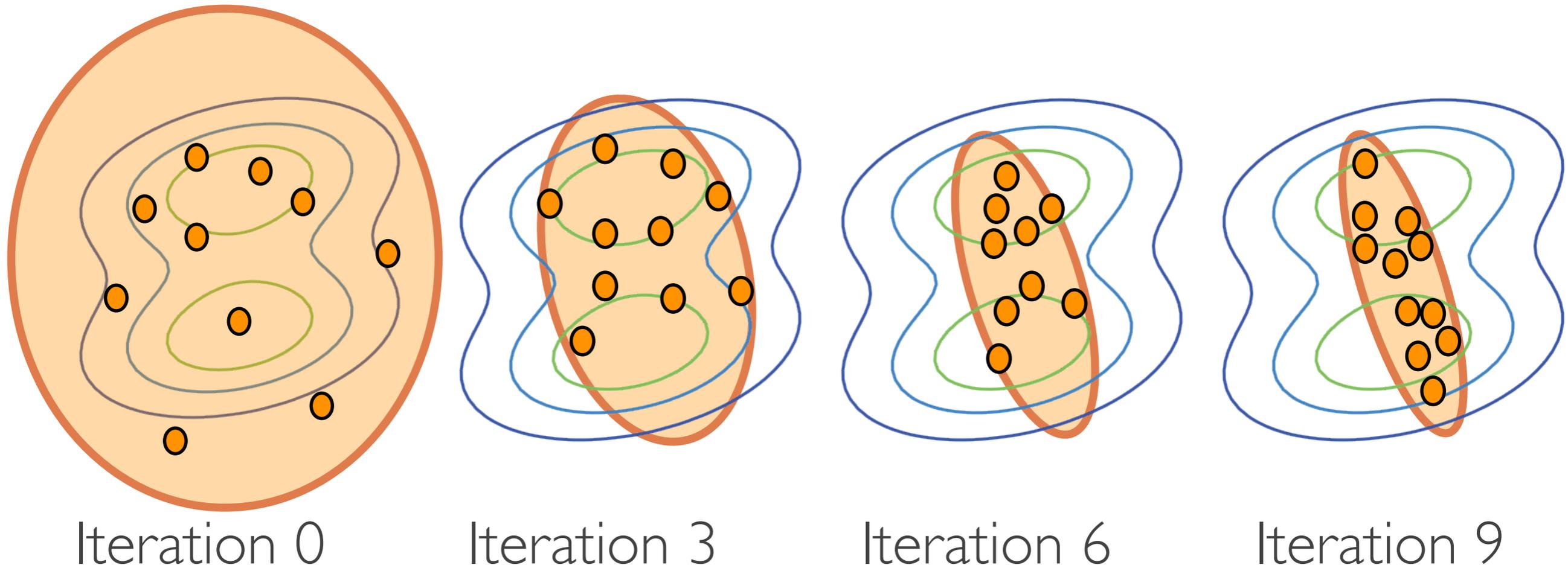
- ➔ More difficult policy search problem
- ➔ We want to find all these solutions

Illustration:



# Illustration

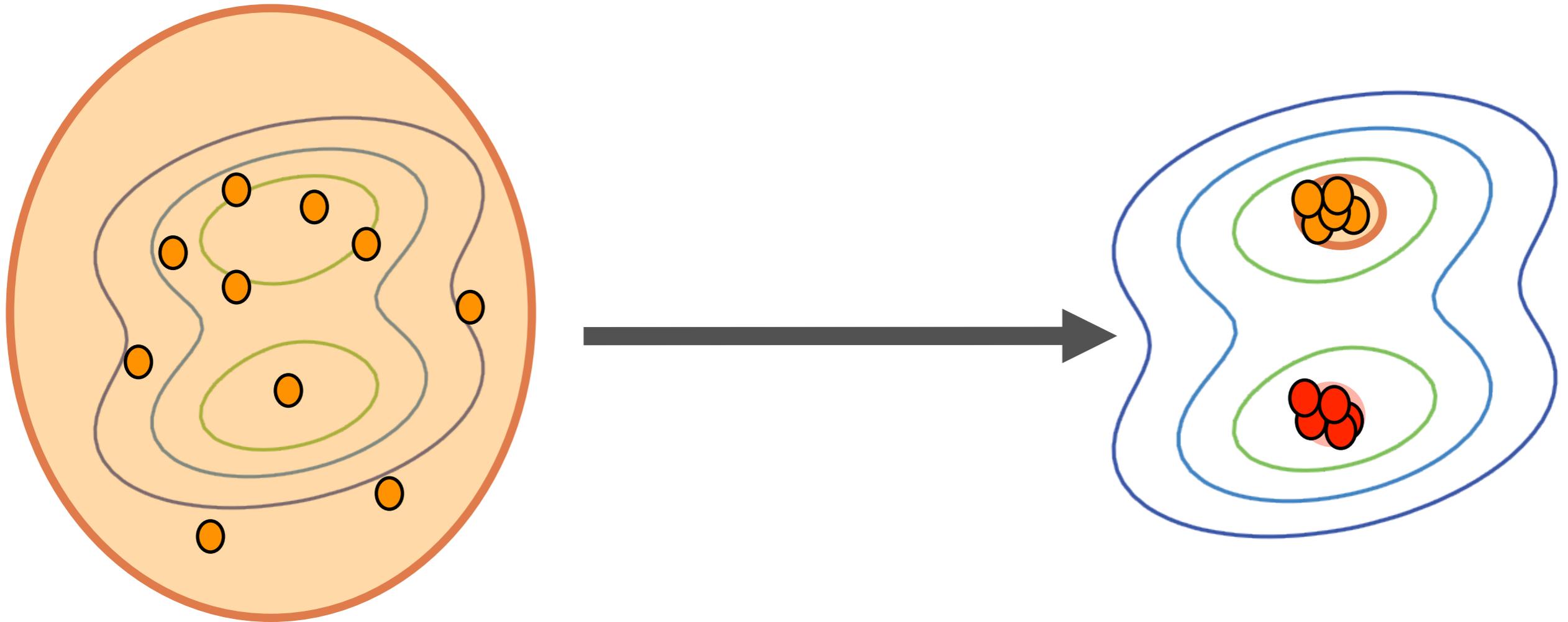
Current Formulation:



Policy averages over several modes.

# Illustration

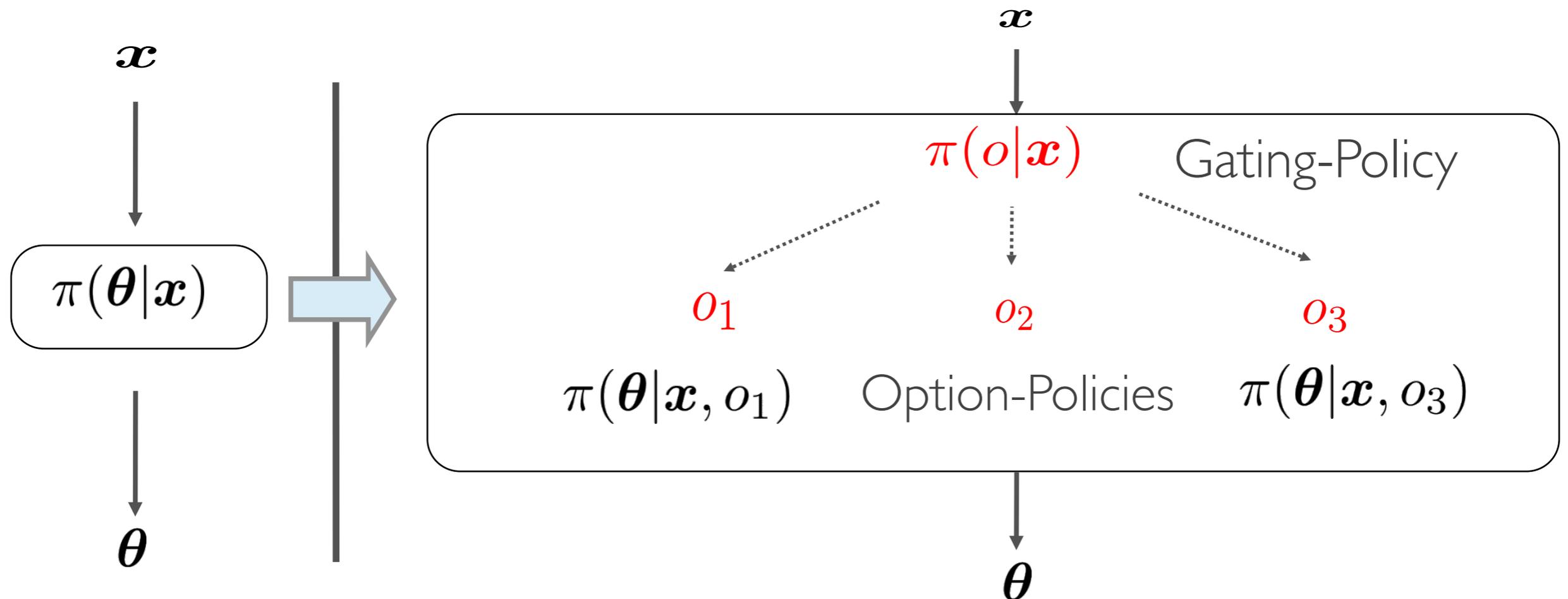
We want to find **both** solutions!



# Introduce Hierarchy

Upper-level policy as combination of options

- Selection of the option: **Gating-policy**
- Selection of the parameters: **Option-policy**



# “Naive” Hierarchical Approach

$$\max_{\pi, \mu} \sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega | \mathbf{x}, \mathbf{o}) \pi(\mathbf{o} | \mathbf{x}) R_{\mathbf{x}\omega} \quad \text{Maximize reward}$$

$$\sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega | \mathbf{x}, \mathbf{o}) \pi(\mathbf{o} | \mathbf{x}) = 1 \quad \text{Distribution}$$

$$\sum_{\mathbf{x}} \mu(\mathbf{x}) \phi(\mathbf{x}) = \hat{\phi} \quad \text{Reproduce Context-Features}$$

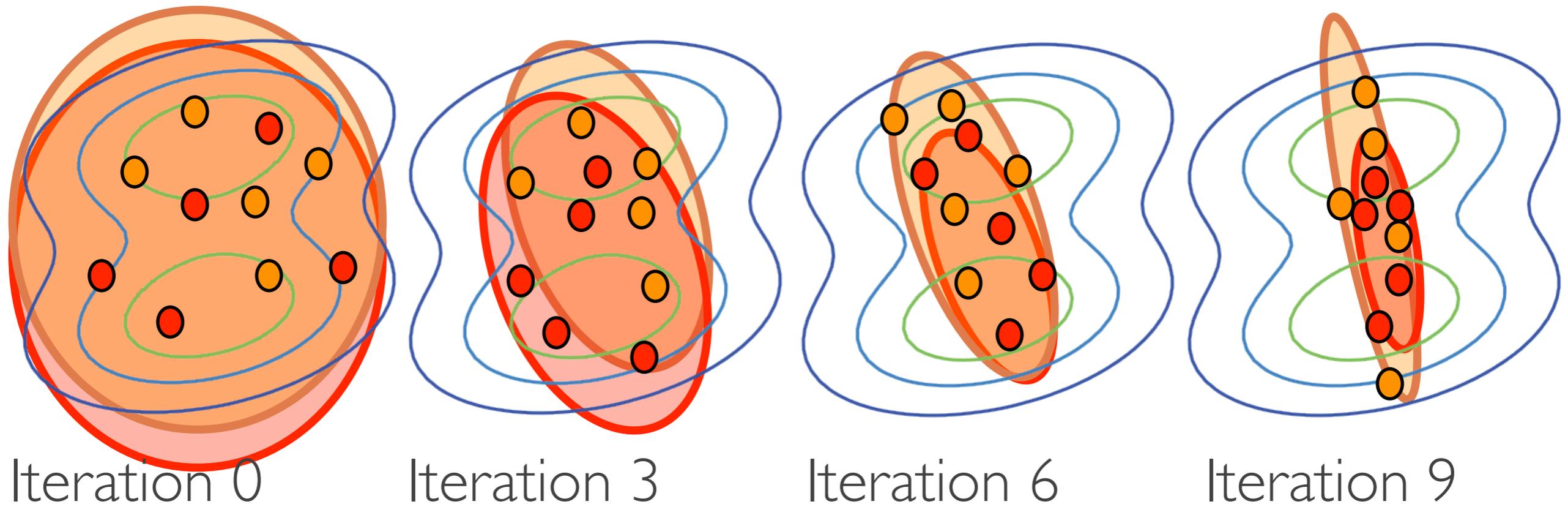
$$\epsilon \geq \sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega, \mathbf{o} | \mathbf{x}) \log \frac{\mu(\mathbf{x}) \pi(\omega, \mathbf{o} | \mathbf{x})}{q(\mathbf{x}, \omega, \mathbf{o})} \quad \text{Stay close to the “data”}$$

?

Versatile Solutions

# Illustration

“Naive” Approach:



Multiple Options, BUT no separation

# Learning versatile Options

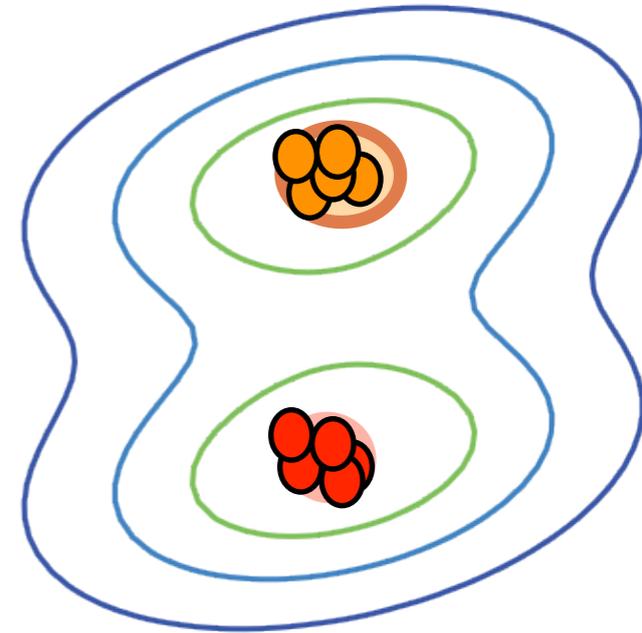
Options should represent distinct solutions.

➔ Limit the overlap of the options

High entropy of  $p(o|\mathbf{x}, \boldsymbol{\theta})$  ➔ high overlap

Limit the entropy ➔ less overlap

$$\kappa \geq \mathbb{E} \left[ \underbrace{- \sum_o p(o|\mathbf{x}, \boldsymbol{\theta}) \log p(o|\mathbf{x}, \boldsymbol{\theta})}_{\text{Entropy}} \right]$$



# Hierarchical REPS (HiREPS)

$$\max_{\pi, \mu} \sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega | \mathbf{x}, \mathbf{o}) \pi(\mathbf{o} | \mathbf{x}) R_{\mathbf{x}\omega}$$

Maximize reward

$$\sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega | \mathbf{x}, \mathbf{o}) \pi(\mathbf{o} | \mathbf{x}) = 1$$

Distribution

$$\sum_{\mathbf{x}} \mu(\mathbf{x}) \phi(\mathbf{x}) = \hat{\phi}$$

Reproduce Context-Features

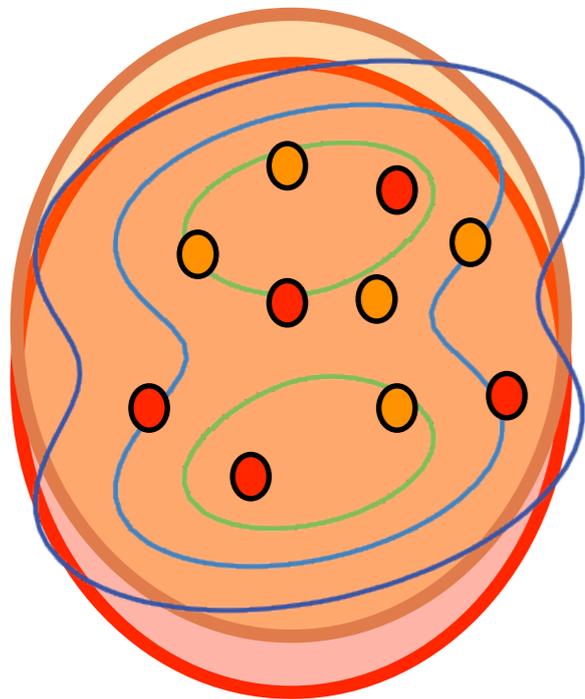
$$\epsilon \geq \sum_{\mathbf{x}, \omega, \mathbf{o}} \mu(\mathbf{x}) \pi(\omega, \mathbf{o} | \mathbf{x}) \log \frac{\mu(\mathbf{x}) \pi(\omega, \mathbf{o} | \mathbf{x})}{q(\mathbf{x}, \omega, \mathbf{o})}$$

Stay close to the “data”, no wild exploration

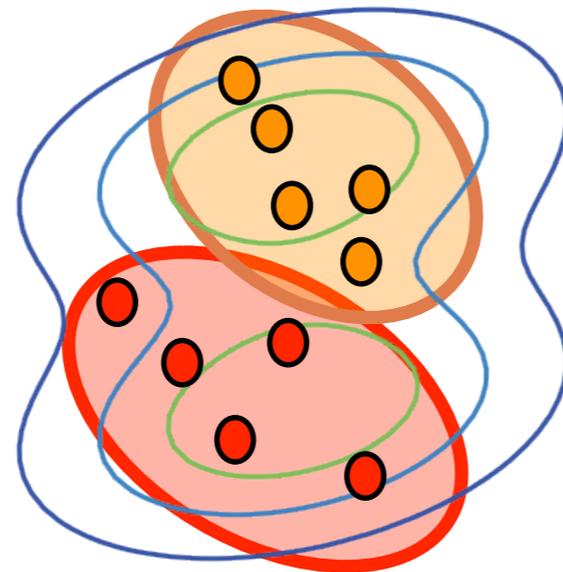
$$\kappa \geq \mathbb{E} [-p(\mathbf{o} | \mathbf{x}, \omega) \log p(\mathbf{o} | \mathbf{x}, \omega)]$$

Versatile Solutions

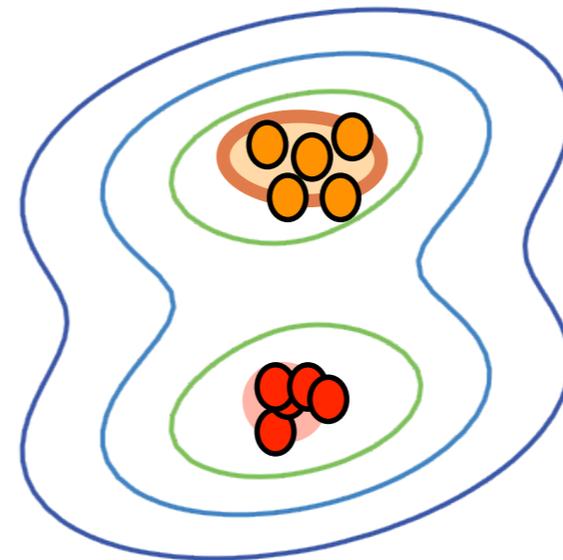
# HiREPS



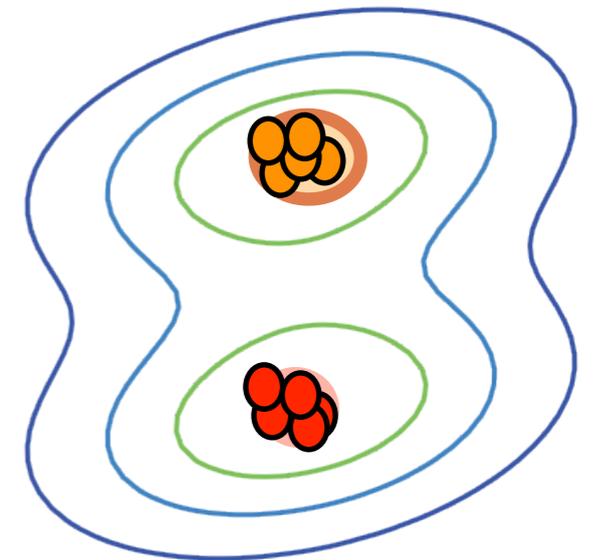
Iteration 0



Iteration 3



Iteration 6



Iteration 9

Learning of versatile, distinct solutions due to separation of options.

# Tetherball



**Pole**

**Barrett  
WAM**



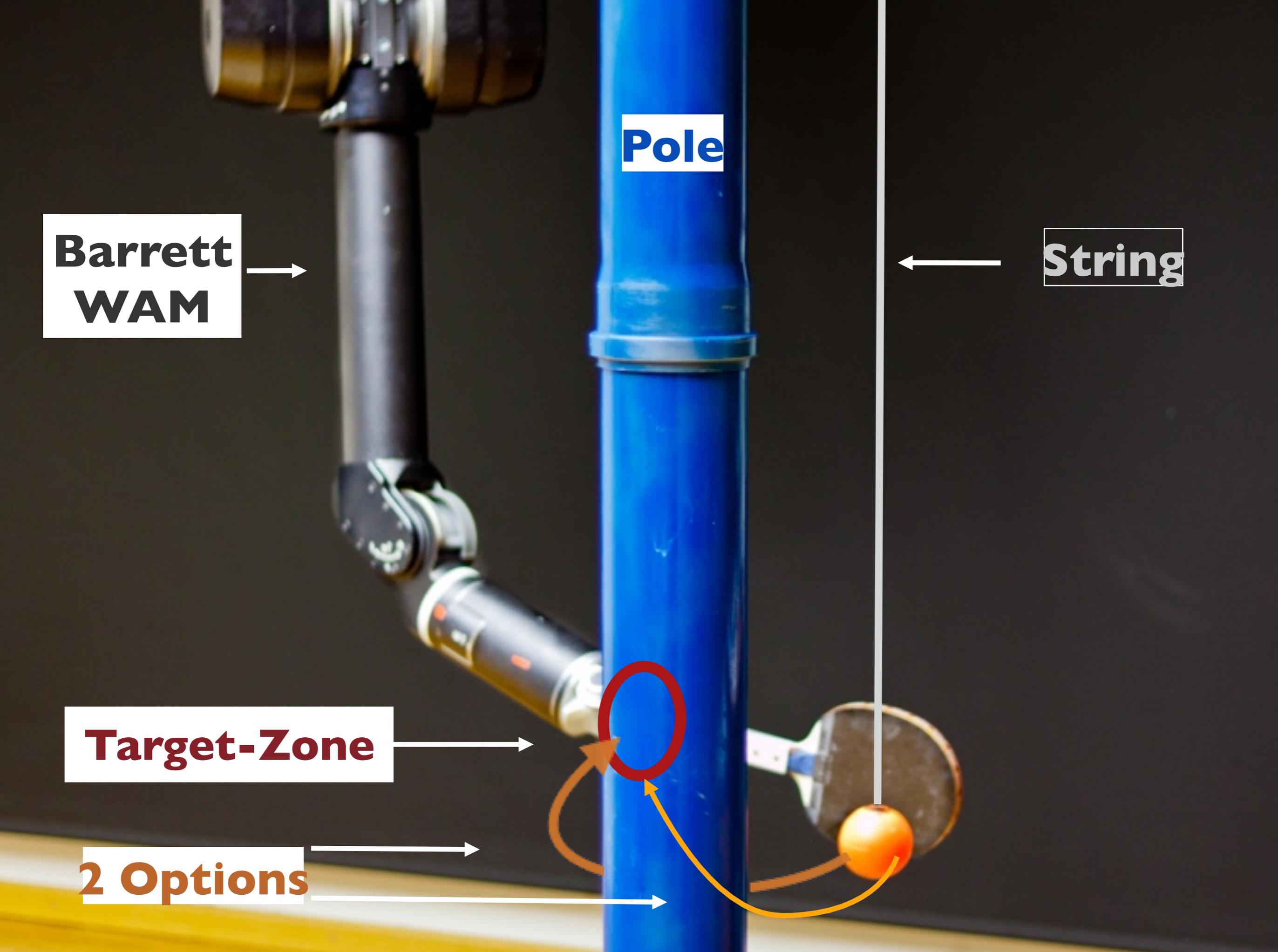
**String**



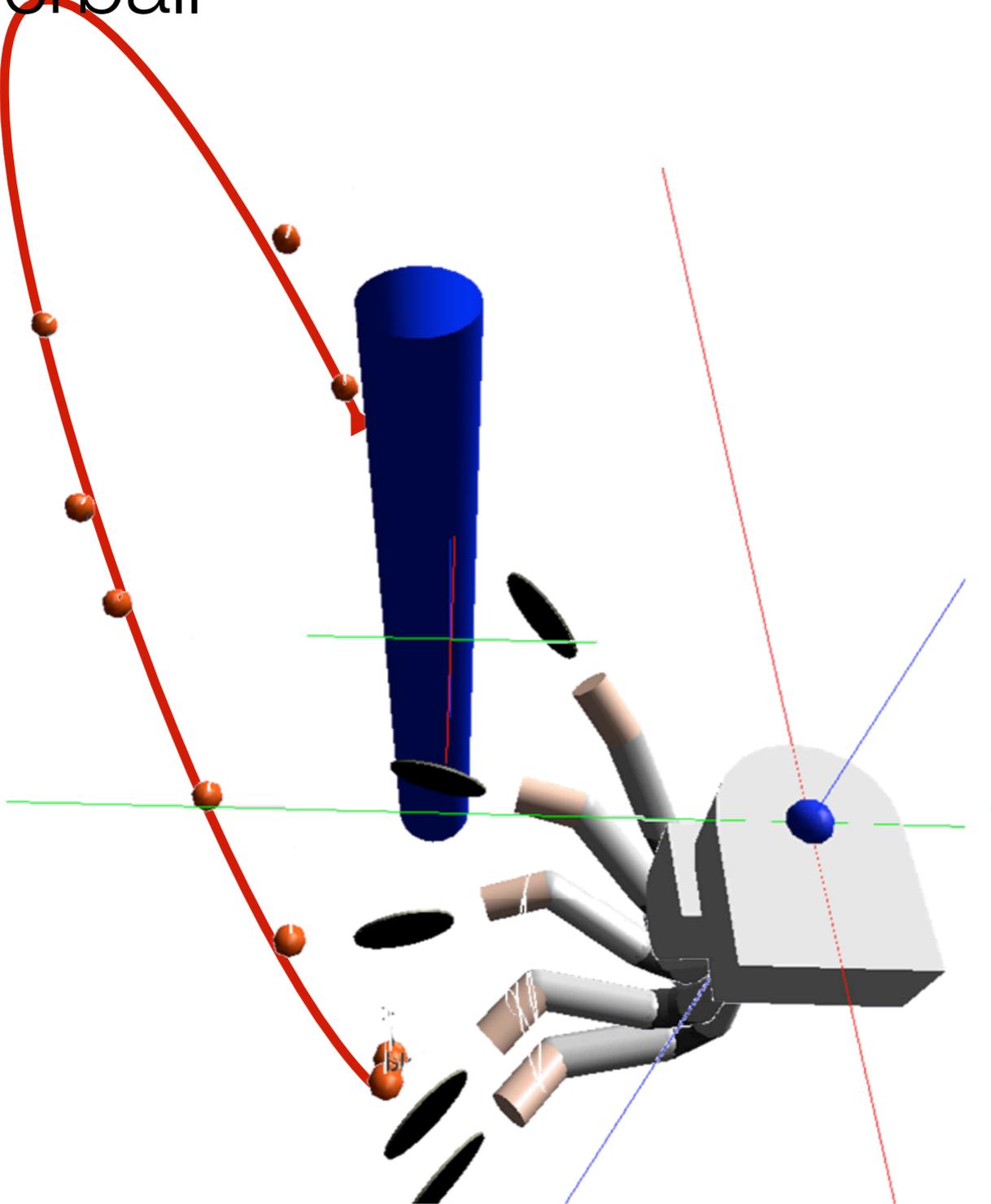
**Target-Zone**



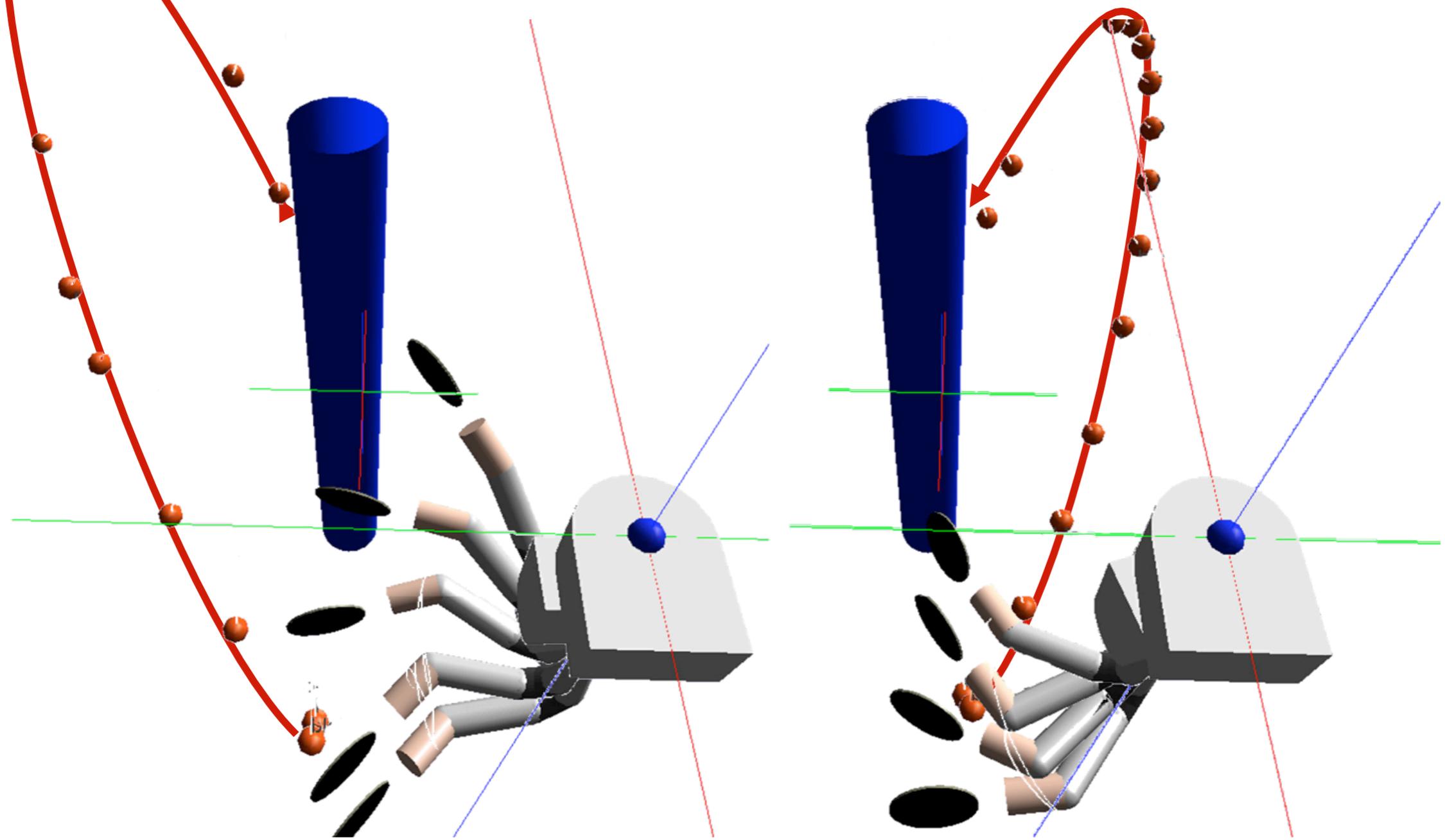
**2 Options**



# Tetherball

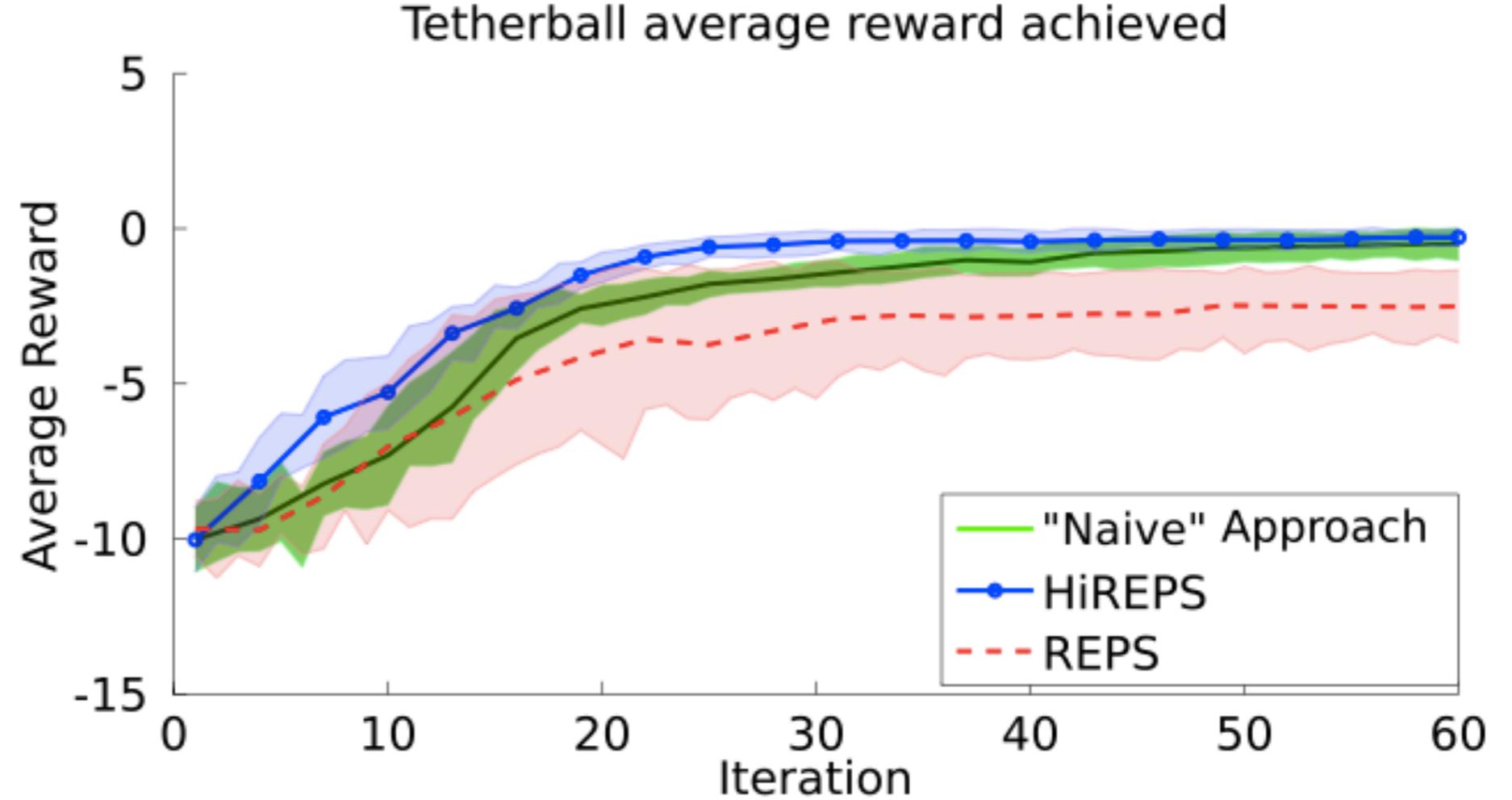


# Tetherball



HiREPS learns distinct solutions.

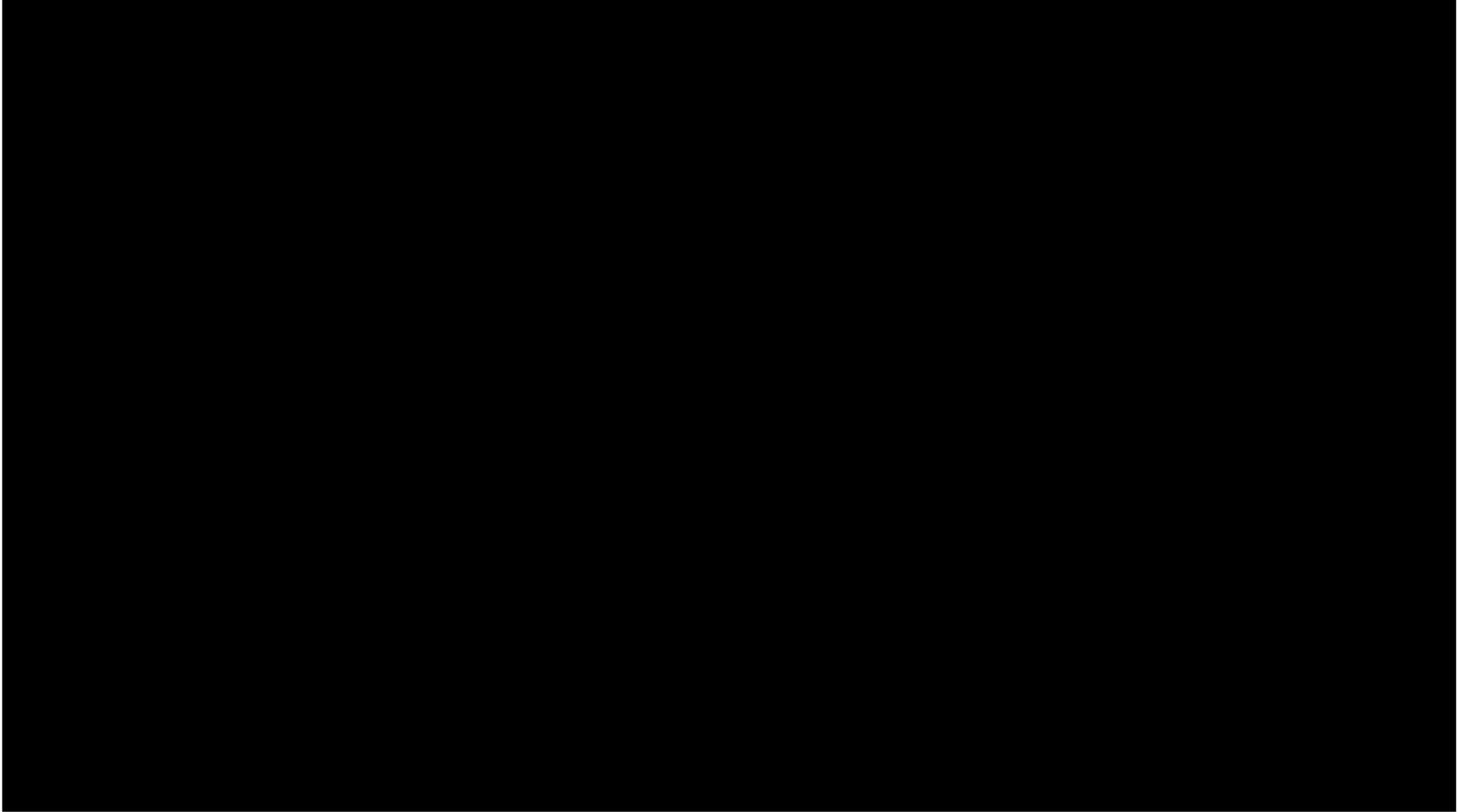
# Results:



Finds [several solutions](#)

Improved convergence, [no averaging](#) over different solutions

Video





# Outline of the Lecture

---

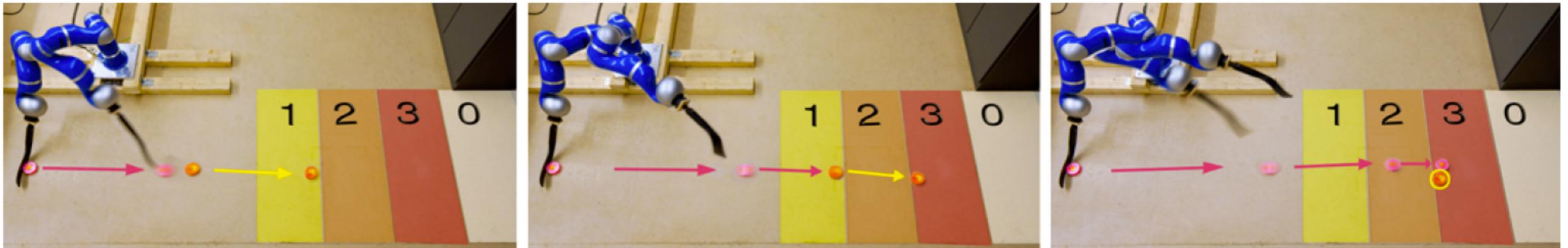
1. Introduction
2. Policy Updates by Weighted Maximum Likelihood
3. Relative Entropy Policy Search (REPS)
4. REPS for Contextual Policy Search
5. Learning Versatile Solutions
6. **Sequencing Movement Primitives**

# Sequencing of Building Blocks

Many motor tasks require a **sequence of elemental building blocks** to fulfill the task

- ➔ The context of **later building blocks** depends on the execution of previous ones
- ➔ We need to learn the **long-term effects** of the building blocks

**Sequential Robot-Hockey Task:** place target-puck in reward zone ,3' after three shoots





# Sequencing of Building Blocks

**Goal:** Sequence several building blocks  $k$  with parameters  $\theta_k$ ,

React to the outcome  $x_k$  of the previous action  $\theta_{k-1}$

**Introduce  $K$  decision steps**

**For each decision step,** learn individual upper-level policy

Maximize the **expected return over all decision steps**

$$J_\pi = \sum_{k=1}^K \iint \mu_k(\mathbf{x}) \pi_k(\boldsymbol{\theta}|\mathbf{x}) R_{\mathbf{x}\boldsymbol{\theta}}^k d\boldsymbol{\theta} d\mathbf{x}$$

**Context distributions:**  $\mu_k(\mathbf{s})$  is specified by the previous policies  $\pi_{l < k}(\boldsymbol{\theta}_l | \mathbf{x}_l)$

$$\mu_k(\mathbf{x}') = \iint \mu_{k-1}(\mathbf{x}) \pi_{k-1}(\boldsymbol{\theta}|\mathbf{x}) p(\mathbf{x}'|\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} d\boldsymbol{\theta}$$

How to compute the policy  $\pi_k(\boldsymbol{\theta}|\mathbf{x})$ ?

**Exploit:** Maximize reward

**Explore:** Stay close to old exploration policy  $q_k(\mathbf{x}, \boldsymbol{\theta})$

Estimate a distribution

Reproduce context distribution

$$\operatorname{argmax}_{p(\mathbf{x}, \boldsymbol{\theta})} \sum_k \sum_{\boldsymbol{\theta}, \mathbf{x}} p_k(\mathbf{x}, \boldsymbol{\theta}) R_{\mathbf{x}\boldsymbol{\theta}, k}$$

$$\text{s.t.: } \text{KL}(p_k(\mathbf{x}, \boldsymbol{\theta}) || q_k(\mathbf{x}, \boldsymbol{\theta})) \leq \epsilon, \quad \forall k$$

$$\sum_{\mathbf{x}, \boldsymbol{\theta}} p_k(\mathbf{x}, \boldsymbol{\theta}) = 1, \quad \forall k$$

$$p_k(\mathbf{x}') = \sum_{\mathbf{x}, \boldsymbol{\theta}} p_{k-1}(\mathbf{x}, \boldsymbol{\theta}) p(\mathbf{x}' | \mathbf{x}, \boldsymbol{\theta})$$

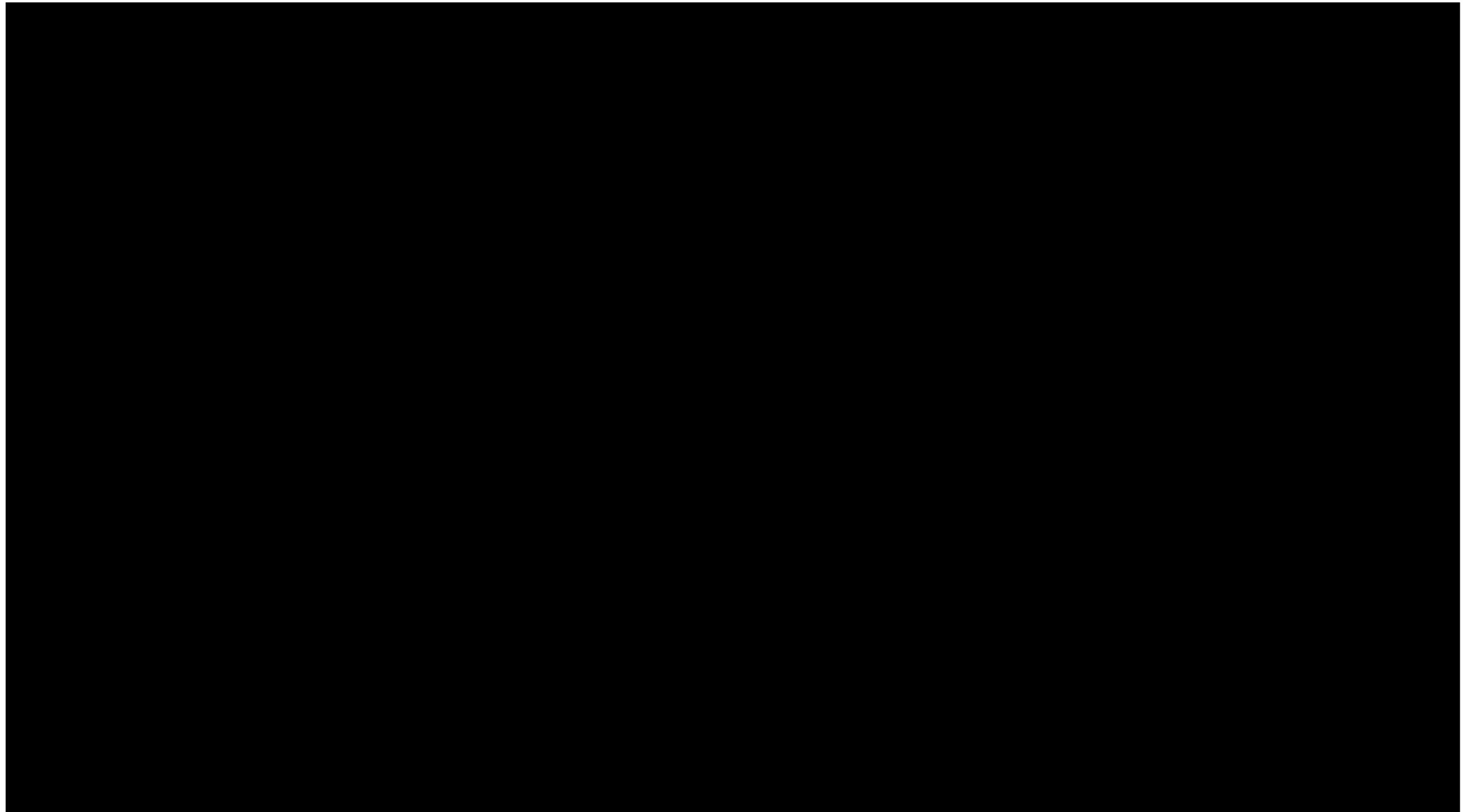
$$\text{Solution: } p_k(\mathbf{x}, \boldsymbol{\theta}) \propto q_k(\mathbf{x}, \boldsymbol{\theta}) \exp\left(\frac{R_{\mathbf{x}\boldsymbol{\theta}, k} + \mathbb{E}_{p(\mathbf{x}' | \mathbf{x}, \boldsymbol{\theta})} [V_{k+1}(\mathbf{x}')] - V_k(\mathbf{x})}{\eta_k}\right)$$

$\mathbb{E}_{p(\mathbf{x}' | \mathbf{x}, \boldsymbol{\theta})} [V_{k+1}(\mathbf{x}')] \dots$  Encodes long-term reward

# Video



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

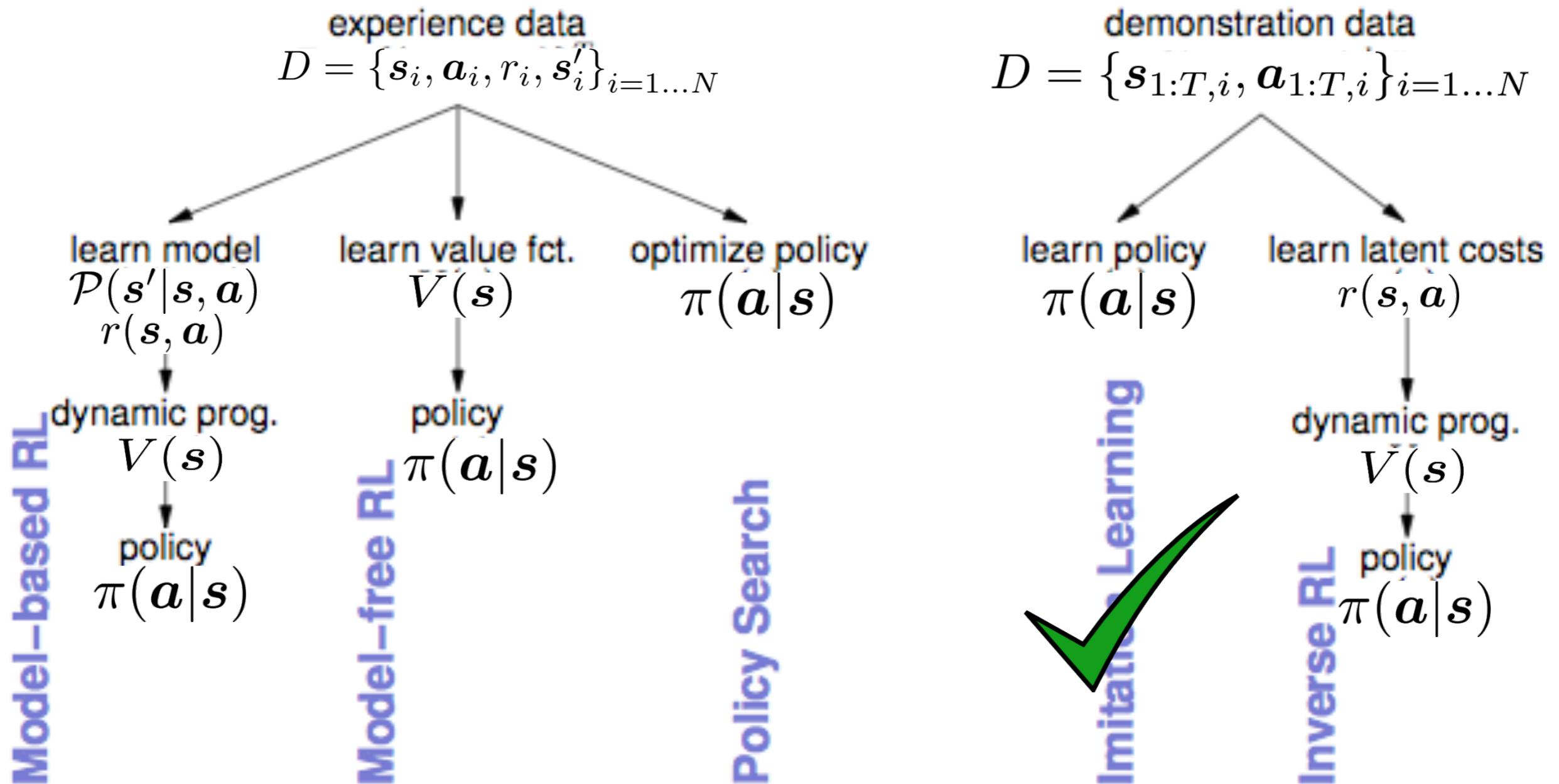


# Conclusion

## Probabilistic Policy Search Methods:

- ➔ Policy update reduces to **weighted maximum likelihood estimates** of the parameters
- ➔ Any type of **structured policy** can be used (e.g. mixture model)
- ➔ Weights are specified by **exponential transformation** of the returns
- ➔ REPS **optimizes the temperature** of this transformation to match a desired Kullback-Leibler divergence
- ➔ **Contextual policy search** can be used for multi-task learning

# Bigger Picture



# Wrap-Up: Model-Based

## **Model Complexity:** Very High

Learn forward model  $f : (\mathcal{R}^{|S|+|A|}) \rightarrow \mathcal{R}^{|S|}$

Need to be able to do dynamic programming (e.g. LQR)

**Small modelling error** can have a big effect on the policy

## **Scalability:** Poor (with some positive exceptions)

Learning high-dimensional (or discontinuous) models is very hard

## **Data-Efficiency:** Excellent

Use every transition to learn model

Model can be reused for different tasks

## **Other Limitations:**

Distance between two policies is hard to control

Huge computation times

# Wrap-Up: Value Based

## **Model Complexity:** OK

Learn Q-Function  $Q : (\mathcal{R}^{|S|+|A|}) \rightarrow \mathcal{R}$

**Small function approximation error** can have a big effect on the policy

## **Scalability:** Poor (with some positive exceptions)

Function approximation in high-dimensional state spaces is difficult

Policy is hard to obtain in high-dimensional action spaces

## **Data-Efficiency:** OK (online TD learning) to good (batch methods)

Batch: Reuse every transition

Online: Every transition is just used once

## **Other Limitations:**

Policy update is again unbounded, might lead to oscillations

# Wrap-Up: Step-Based Policy Search

**Model Complexity:** None (no approximation errors)

Need to evaluate reward to come  $Q_t^{[i]}$

**Scalability:** Good

Parametrized policies are a compact representation that allow learning also for high-D robots

Only works for a medium amount of parameters (a few hundred)

**Data-Efficiency:** Poor

Use every state action pair with reward to come

High variance in reward to come due to exploration in action space

**Other Limitations:**

Mainly used for learning single trajectories (e.g. DMPs)

# Wrap-Up: Episode-Based Policy Search

**Model Complexity:** None (no approximation errors)

Need to evaluate return for each trajectory  $R^{[i]}$

**Scalability:** Good

Parametrized policies

Only works for a small amount of parameters (around hundred)

**Data-Efficiency:** Poor

Each rollout is just one sample

High variance in returns in case of stochastic environments

**Other Limitations:**

Mainly used for learning single trajectories (e.g. DMPs)