

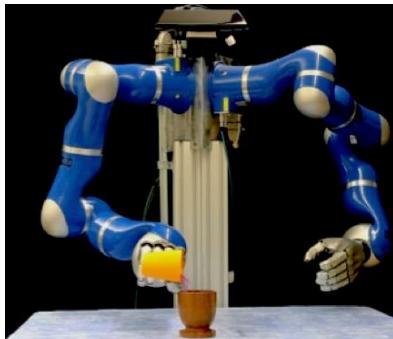
# *(REINFORCEMENT) LEARNING IN ROBOTICS*

Gerhard Neumann



# Long term vision

## Future: Autonomous robots



Household



Health / Elderly Care



Assist as Co-Workers



Farming



Hazardous Environments



Traffic

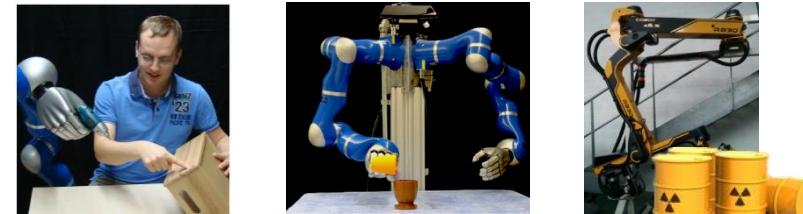
### Challenges:

- ✗ Complex skills
- ✗ Data-efficient self improvement
- ✗ Cooperation
- ✗ Perception

# How can we accomplish these challenges?

## Machine Learning

## Data-Driven Algorithms



## Autonomous Robots

## Domain Knowledge

Data-driven movement primitives

✗ Complex Skills

Reinforcement Learning

✗ Self-improvement

Representation Learning

✗ Cooperation

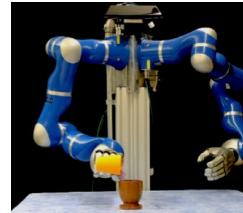
Multi-Agent Learning

✗ Perception

# How can we accomplish these challenges?

## Machine Learning

## Data-Driven Algorithms



## Autonomous Robots

## Domain Knowledge

Data-driven movement primitives

✗ Complex Skills

Reinforcement Learning

✗ Self-improvement

Multi-Agent Learning

✗ Cooperation

Representation Learning

✗ Perception

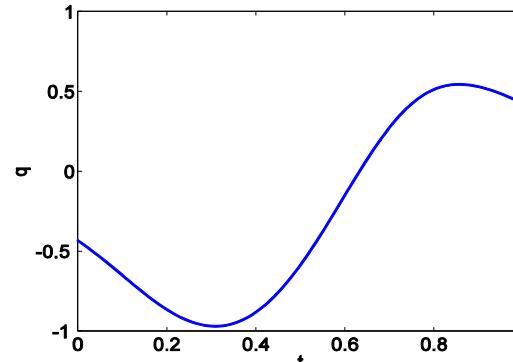
# Data-Driven Movement Representations

**Goal:** data-driven, compact, flexible

**Parametrized trajectories:**

$$\tau^* = q_{1:T} = f(\theta)$$

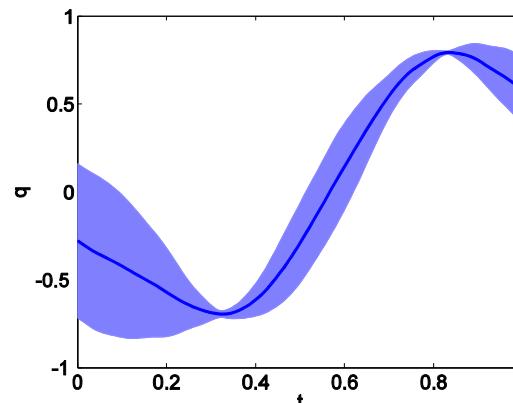
- Mean movement
- Followed by trajectory tracking controllers
- Example: Dynamic Movement Primitives (DMPs) [Ijspeert 2002]



**Parametrized trajectory distribution:**

$$\tau \sim p(\tau) \Leftrightarrow \theta \sim p(\theta)$$

- Family of movements
- **Gaussian:** Mean and Variance
- **Probabilistic Movement Primitives (ProMPs)** [Paraschos 2013]



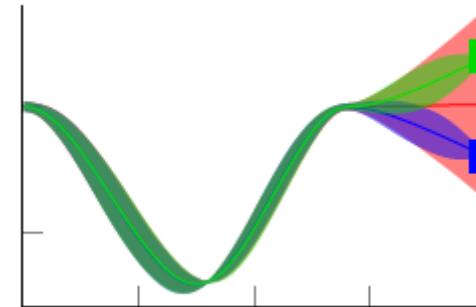
A. Paraschos, ..., G. Neumann, *Probabilistic Movement Primitives*, NIPS, 2013

# Adaptation and Variable Stiffness

Adapt final/intermediate position of the movement

- Conditioning/Bayes Theorem

$$p(\theta | \mathbf{q}_t = \mathbf{q}_t^*) = \frac{p(\mathbf{q}_t^* | \theta) p(\theta)}{p(\mathbf{q}_t^*)}$$

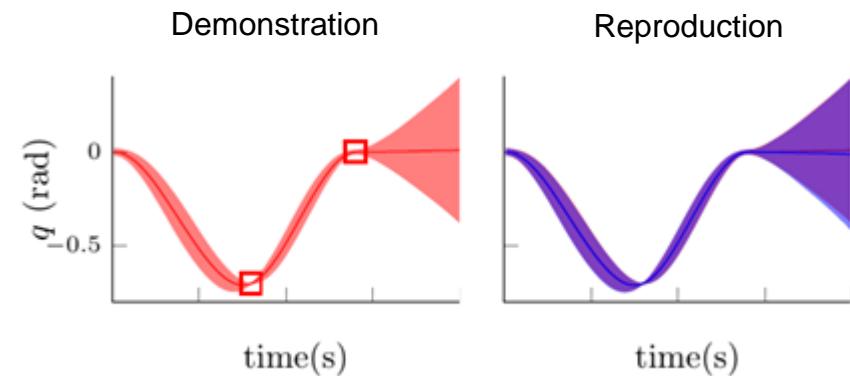


## Variable stiffness

- Obtain variable stiffness controller

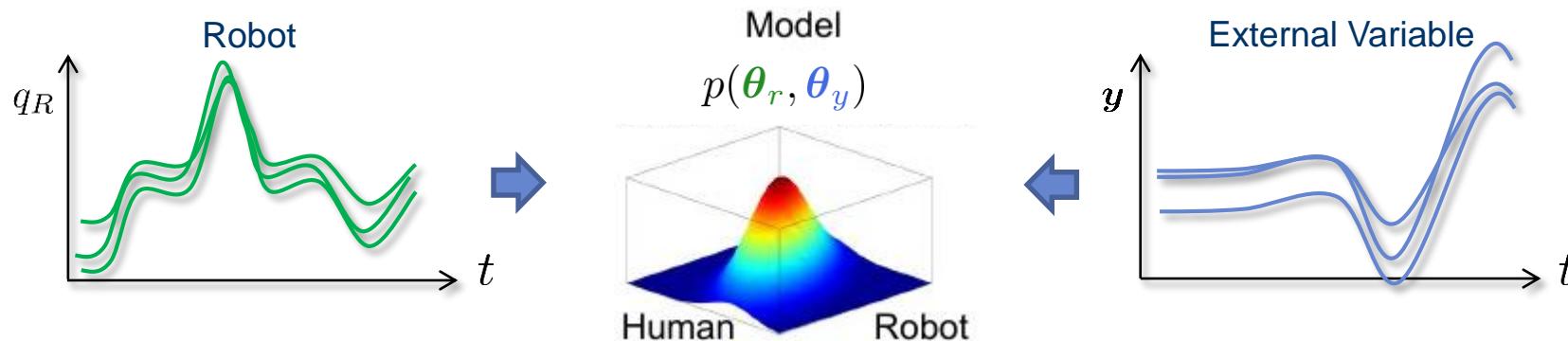
$$\mathbf{u}_t = \mathbf{K}_t \mathbf{x}_t + \mathbf{k}_t$$

- Matches moment of trajectory distribution
- High stiffness only if needed
- Requirement for safe interaction

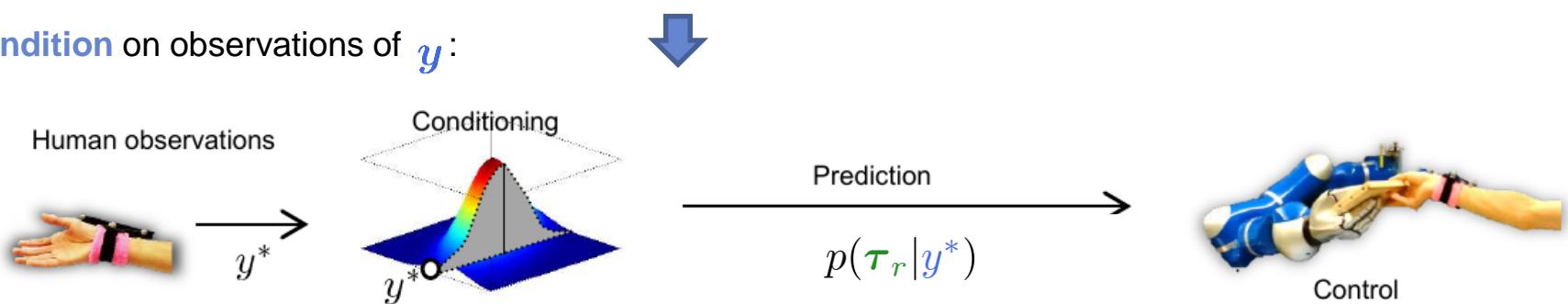


# Movement coupling

- E.g. couple with movement of interaction partner  $y \dots$
- Learn **joint distribution** of trajectory  $\tau_y$  and robot trajectory  $\tau_r$  from demonstrations



- **Condition** on observations of  $y^*$ :

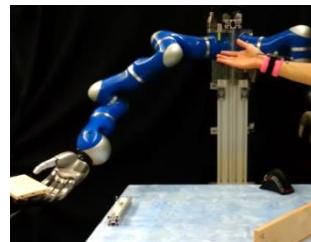


# Robot Assistant

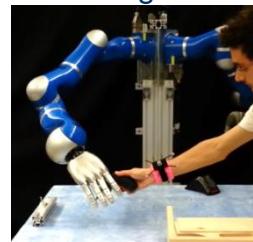
## Box assembly

- Learn interaction patterns by kinesthetic teach in

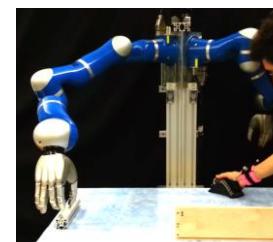
Plate handover



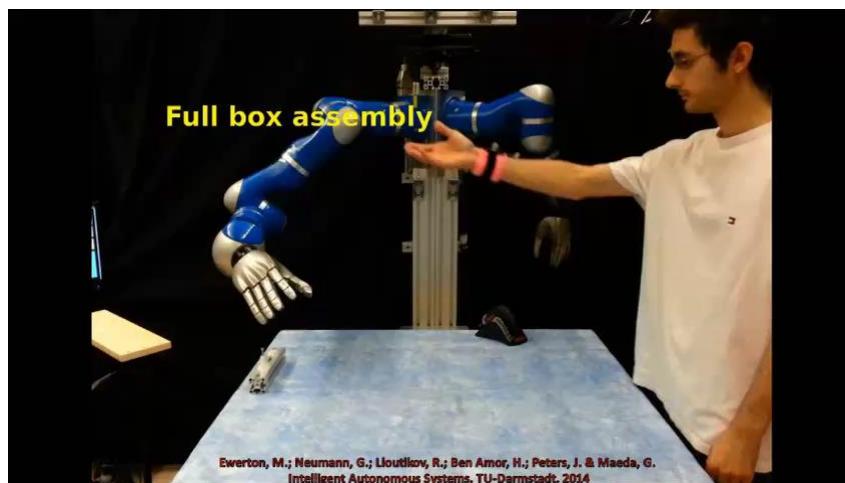
Holding tool



Screw handover



- Couple robot movement with human



## Flexible and Complex Skills:

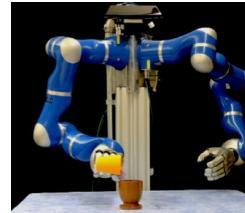
- Variable stiffness (for safe interaction)
- Modulation
- Coupling

Maeda, G.; Neumann, G.; et al. *Probabilistic Movement Primitives, for Coordination of Multiple Human-Robot Collaborative Tasks*,  
Autonomous Robots (AURO)

# How can we accomplish these challenges?

## Machine Learning

## Data-Driven Algorithms



## Autonomous Robots

## Domain Knowledge

Data-driven movement primitives

✗ Flexible Skills

Reinforcement Learning

✗ Self-improvement

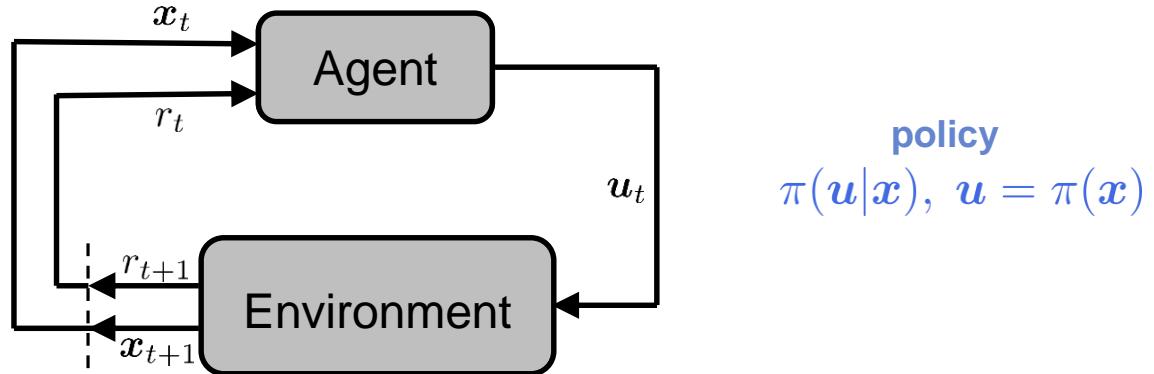
Multi-Agent Learning

✗ Cooperation

Representation Learning

✗ Perception

# Reinforcement Learning (RL)



Learning: Adapting the policy  $\pi(u|x)$  of the agent

**Objective:** Find policy that maximizes expected return

$$J_\pi = \mathbb{E} \left[ \sum_{t=0}^T r_t \middle| u_t \sim \pi(\cdot|x_t) \right]$$

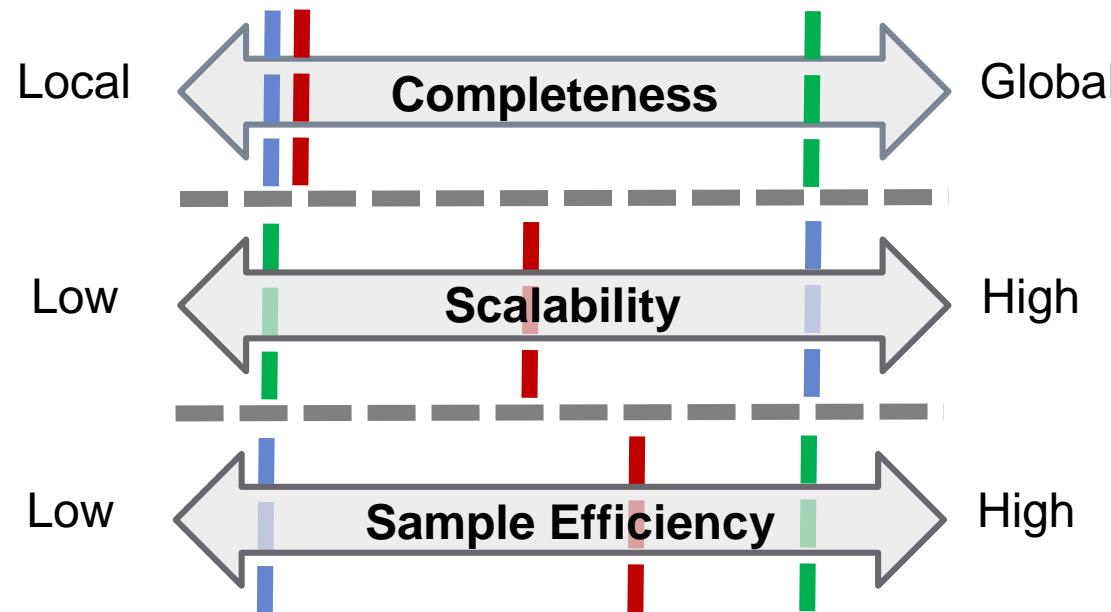
Learning with movement primitives

**Objective:** Find parameters  $\theta$  that maximize the expected return

$$R(\theta) = \mathbb{E} \left[ \sum_{t=0}^T r_t \middle| u_t = \pi_\theta(x_t) \right]$$

# Coarse Categorization of RL

- Deep Reinforcement Learning
- Bayesian Optimization
- Direct Policy Search / Stochastic Search



# Direct Policy Search

## Direct Policy Search by Stochastic Search:

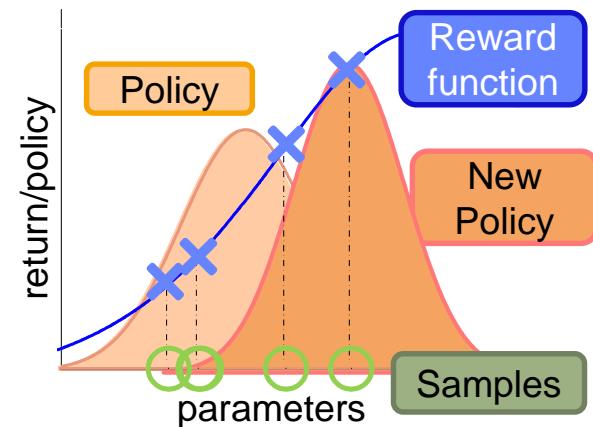
$$\pi^* = \arg \max_{\pi} \int \pi(\theta) R(\theta) d\theta$$

Parameters of Movement Primitive

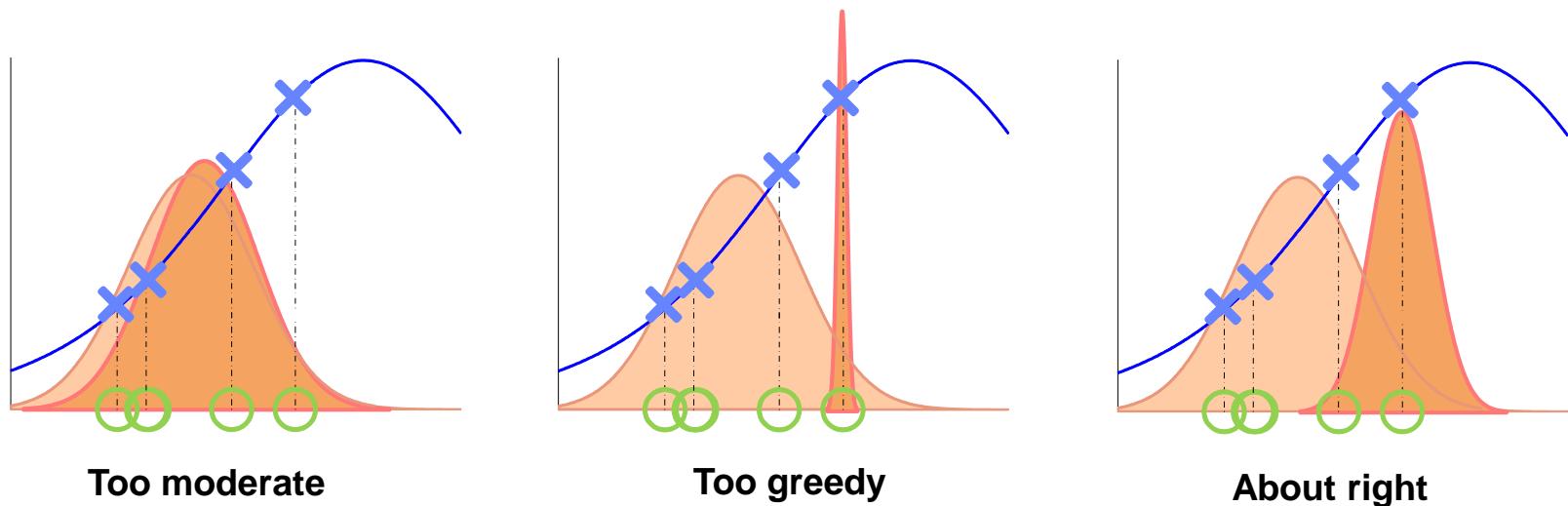
- Find distribution  $\pi(\theta)$  with maximum expected return  $R(\theta)$

## Three basic steps:

- **Explore:** Generate trajectories  $\tau^{[i]}$  following the policy  $\pi_k$
- **Evaluate:** Assess quality of trajectory or actions
- **Update:** Compute new policy  $\pi_{k+1}$



# How to find a good policy update?



**Control exploration-exploitation tradeoff:** immediate vs. long-term performance

# Information-Theoretic Policy Update

**Information-theoretic policy update:** incorporate information from new samples

1. Maximize return

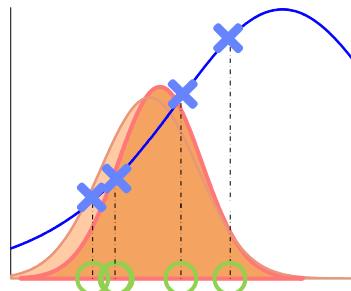
$$\arg \max_{\pi} \int \pi(\theta) R(\theta) d\theta$$

2. Bound information gain [Peters 2011]

$$\text{s.t. } \text{KL}(\pi || \pi_{\text{old}}) \leq \epsilon \quad \text{Control Greediness}$$

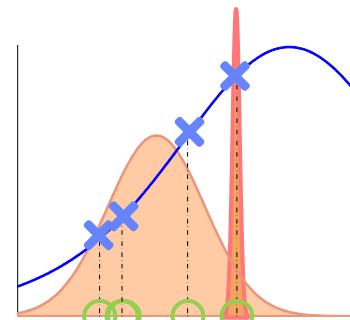
3. Bound entropy loss [Abdolmaleki 2015]

$$\underbrace{H(\pi_{\text{old}}) - H(\pi)}_{\text{loss in entropy}} \leq \gamma \quad \text{Control Entropy}$$



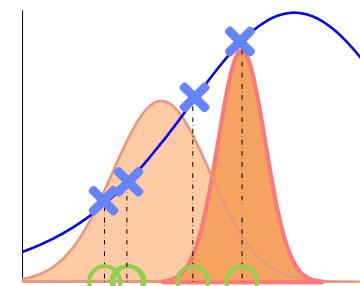
Too moderate

$$\epsilon \ll 1$$



Too greedy

$$\epsilon \gg 1 \quad \gamma \gg 1$$

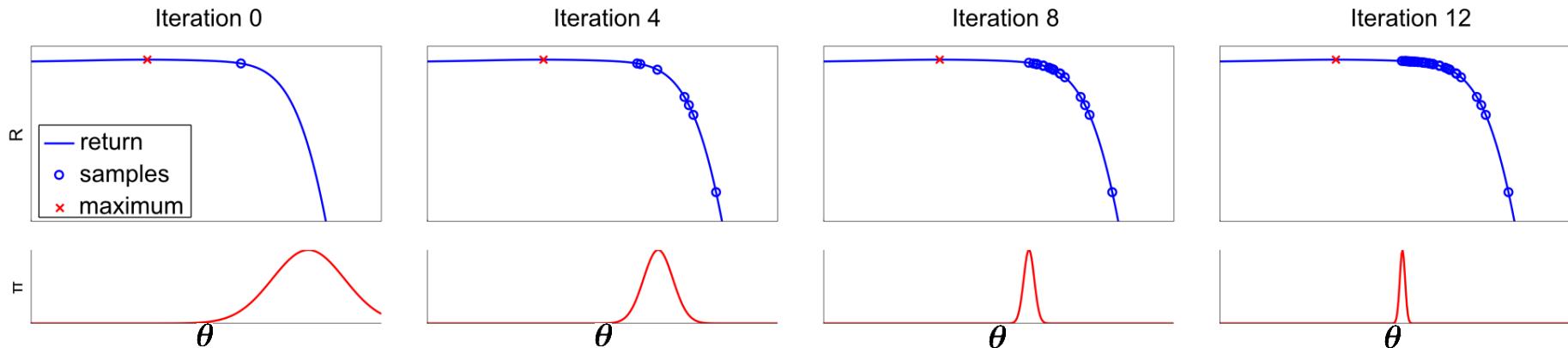


About right

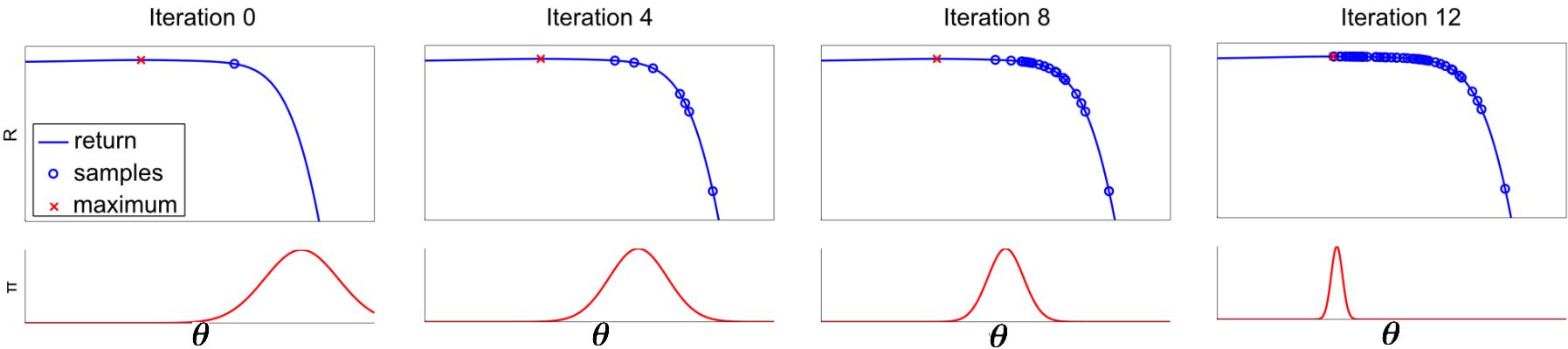
$$\epsilon \gg 1 \quad \gamma \ll 1$$

# Illustration: Distribution Update

No entropy loss bound



With bounded entropy loss



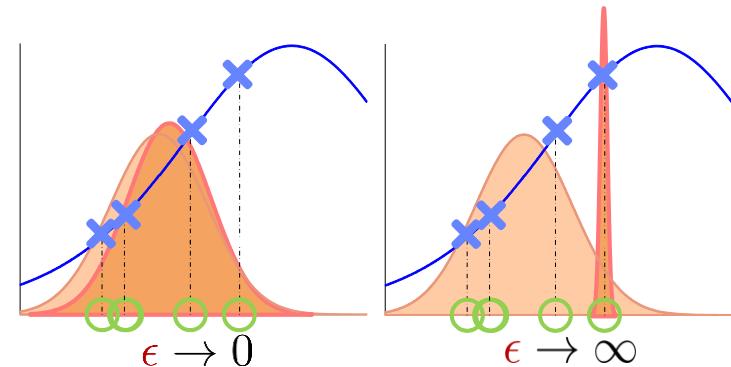
# Solution for Search Distribution

**Solution for unconstrained distribution:**  $\pi(\theta) \propto \pi_{\text{old}}(\theta)^{\frac{\eta}{\eta+\omega}} \exp\left(\frac{R(\theta)}{\eta + \omega}\right)$

- $\eta$  ... Lagrangian multiplier for:  $\text{KL}(\pi || \pi_{\text{old}}) \leq \epsilon$

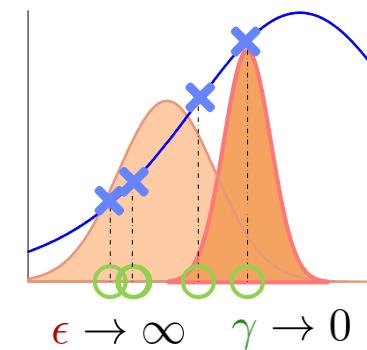
$$\epsilon \rightarrow 0 \quad \Rightarrow \quad \eta \rightarrow \infty \quad \Rightarrow \quad \pi \rightarrow \pi_{\text{old}}$$

$$\epsilon \rightarrow \infty \quad \Rightarrow \quad \eta \rightarrow 0 \quad \Rightarrow \quad \pi \rightarrow \text{greedy}$$



- $\omega$  ... Lagrangian multiplier for:  $H(\pi_{\text{old}}) - H(\pi) \leq \gamma$

$$\gamma \rightarrow 0 \quad \Rightarrow \quad \omega \gg 0 \quad \Rightarrow \quad \pi \rightarrow \text{more uniform}$$



**Gaussianity needs to be „enforced“ !**

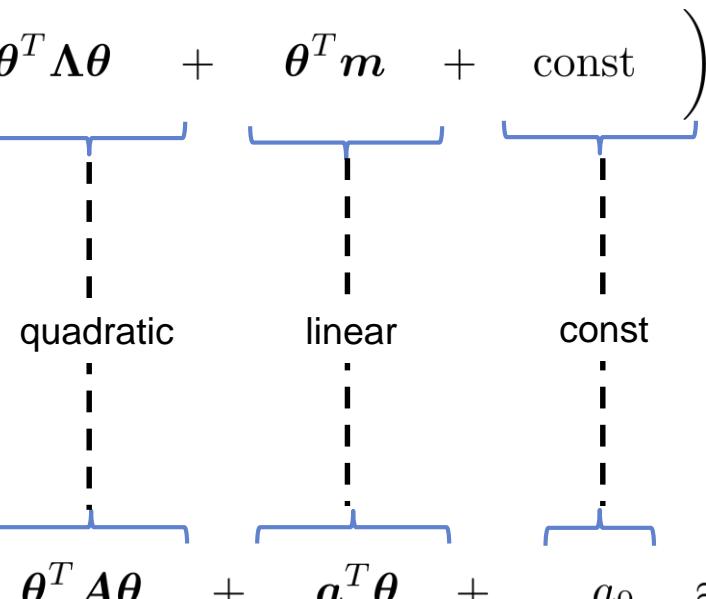
- Fit **new policy** on samples (REPS, [Daniel2012, Kupcsik2014, Neumann2014])
- Fit **return function** using compatible function approximation (MORE, [Abdolmaleki2015])

A. Abdolmaleki, ..., G. Neumann, *Model-Based Relative Entropy Stochastic Search*, NIPS 2015

# Fit Return Function

Use compatible function approximation:

- Gaussian distribution:  $\mathcal{N}[\theta | \mathbf{m}, \Lambda] \propto \exp\left(-\frac{1}{2}\theta^T \Lambda \theta + \theta^T \mathbf{m} + \text{const}\right)$
- Gaussian in canonical form (log linear)
- Parameters:** Precision  $\Lambda$  and linear part  $\mathbf{m}$



Match functional form:

$$\tilde{R}(\theta) = \theta^T \mathbf{A} \theta + \mathbf{a}^T \theta + a_0 \approx R(\theta)$$

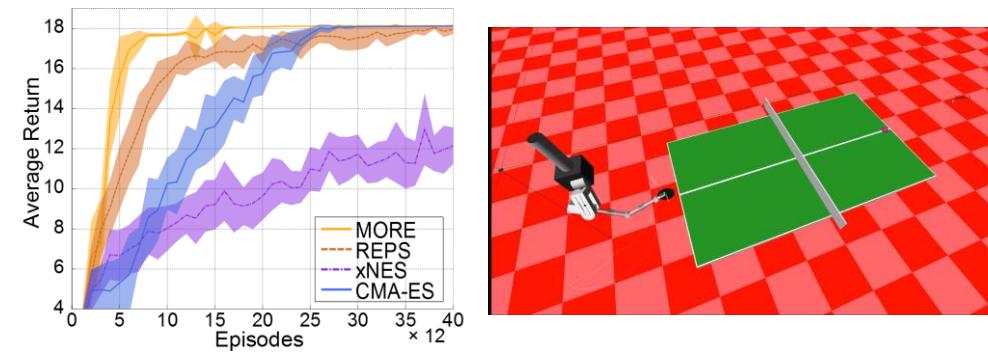
- Quadratic in  $\theta$ , but linear in parameters:  $\mathbf{w} = \{\mathbf{A}, \mathbf{a}, a_0\}$
- $\mathbf{w}$  obtained by **linear regression** on current set of samples
- Fits **local curvature** of the objective function

# MORE Algorithm

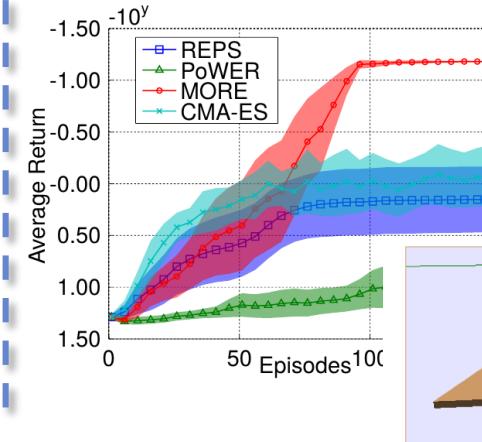
Model-Based Relative Entropy Stochastic Search (MORE) : [Abdolmaleki 2015]

1. Evaluation: Fit local surrogate  $\tilde{R}(\theta) \approx \theta^T A\theta + a^T \theta + a_0$
2. Update:  $\pi(\theta) \propto \pi_{\text{old}}(\theta)^{\frac{\eta}{\eta+\omega}} \exp\left(\frac{\tilde{R}(\theta)}{\eta+\omega}\right) \Rightarrow \pi(\theta) = \mathcal{N}(\theta|\mu^*, \Sigma^*)$

## Table-Tennis



## Beer-Pong



A. Abdolmaleki, ..., G. Neumann, Model-Based Relative Entropy Stochastic Search, NIPS 2015

# Adaptation of Skills

## Adapt parameters to context $s$

- Continuous valued vector
- Characterizes environment / objectives

## Learn contextual distribution:

$$\pi(\theta|s) = \mathcal{N}(\theta|M\phi(s), \Sigma)$$

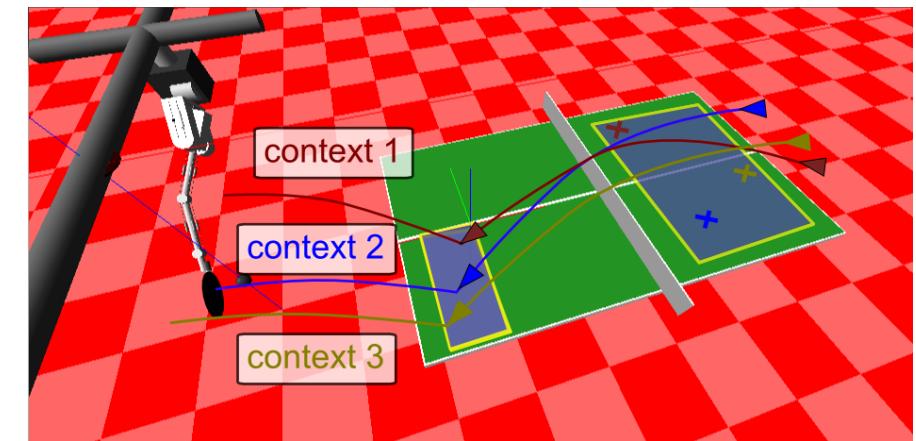
- Parameters are adapted to context

## Leads to equivalent formulation:

$$\arg \max_{\pi} \mathbb{E}_{p(s)} \left[ \int \pi(\theta|s) R(s, \theta) d\theta \right]$$

$$\text{s.t.: } \mathbb{E}_{p(s)} [\text{KL}(\pi(\cdot|s) || \pi_{\text{old}}(\cdot|s))] \leq \epsilon$$

$$H(\pi_{\text{old}}) - H(\pi) \leq \gamma$$

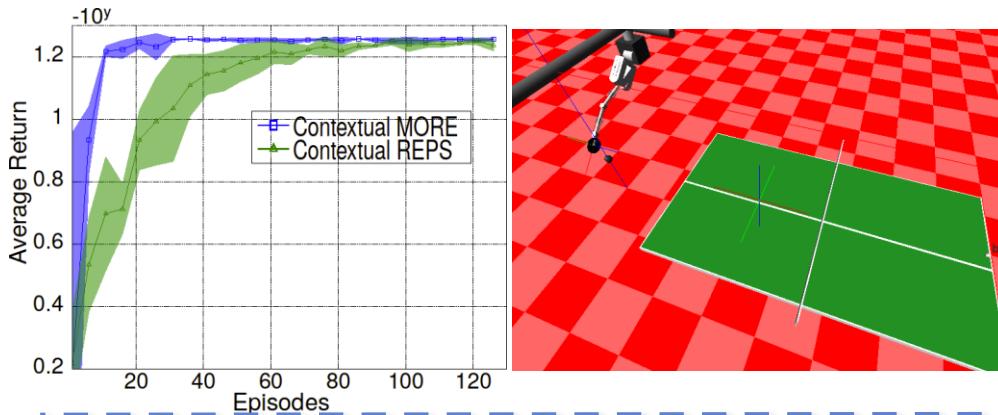


Abdolmaleki, ..., Neumann, Model-Based Relative Entropy Stochastic Search, NIPS 2015

# Adaptation of Skills

## Table-Tennis:

- Initial ball configuration
- Desired return impact point



## Robot-Hockey:

- Position of target puck
- Desired displacement

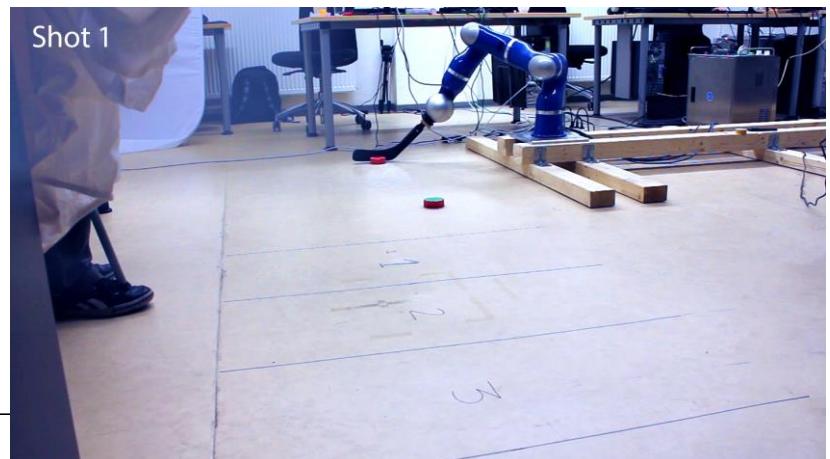


## Also works for high-D context:

- Raw images of ball positions
- Has to reach shown ball position
- Regularization necessary



V. Tangkaratt et al., "Policy search with high-dimensional context variables", in *AAAI Conference on Artificial Intelligence (AAAI)*, 2017.

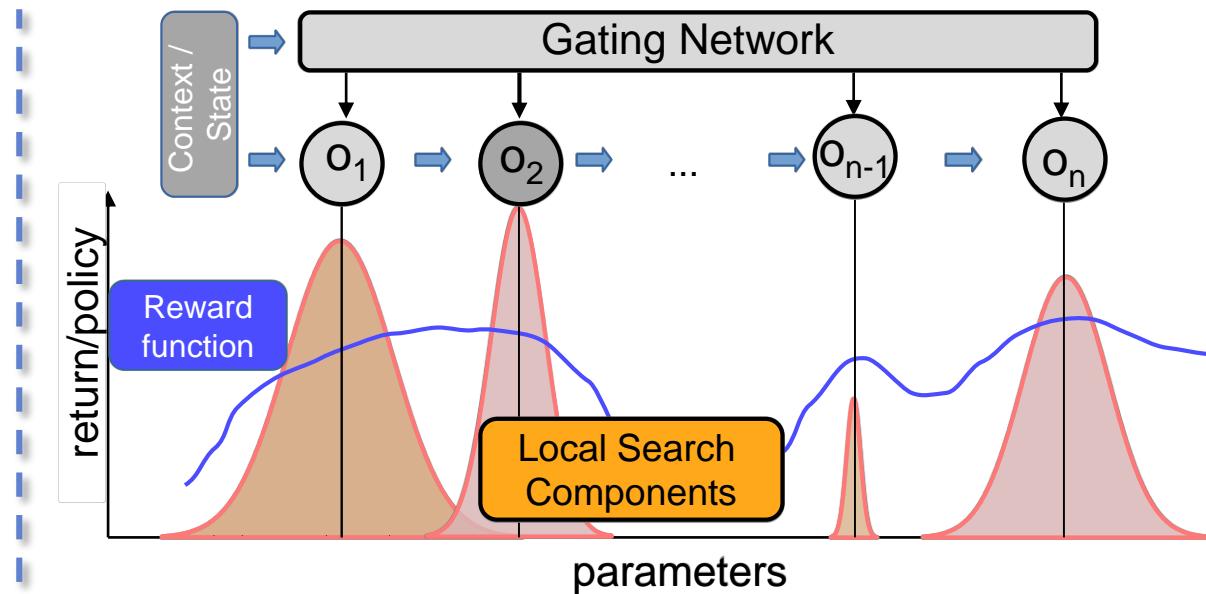


# Versatile Policies for Multi-Modal Solutions

## 1. Hierarchical Policy:

- **Gating Policy:**  $\pi(o|s)$ 
  - Select skill  $o$  given context
- **Local Sub-policies:**  $\pi(\theta|s, o)$ 
  - Select parameters given context

$$\pi(\theta|s) = \sum_o \pi(o|s)\pi(\theta|s, o)$$



## 2. Maximum Entropy RL

[Hannula et al, NIPS 2017]

$$J = \sum_o \pi(o) \int \pi(\theta|o) R(\theta) d\theta - \underbrace{\sum_o \pi(o) \int \pi(\theta|o) \log \pi(\theta)}_{H(\pi)}$$

→ **intractable**

# Variational Lower Bound

The objective  $J$  can be decomposed into:

$$J = \underbrace{U(\pi, \tilde{p}(o|\theta))}_{\text{Lower Bound}} + \text{KL}(p(o|\theta) \parallel \underbrace{\tilde{p}(o|\theta)}_{\text{Var. Dist.}})$$

Variational Distribution

- E-step: Tighten lower bound

$$\tilde{p}(o|\theta) = p(o|\theta) \Rightarrow \text{KL}(p(o|\theta) \parallel \tilde{p}(o|\theta)) = 0$$

- M-Step: Maximize lower bound

$$U(\pi, \tilde{p}(o|\theta)) = \sum_o \pi(o) \int \pi(\theta|o) (R(\theta) + \log \tilde{p}(o|\theta)) d\theta + H(\pi(\theta|o)) + H(\pi(o))$$

Yields decomposable max-ent RL problems with modified reward:

- **Gating:**

$$r(o) = \mathbb{E}_{\pi(\theta|o)} [R(\theta) + \log \tilde{p}(o|\theta)] + H(\pi(\theta|o))$$

- **Parameter policy:**

$$r_o(\theta) = R(\theta) + \log \tilde{p}(o|\theta)$$

# Application to Variational Inference

**Test on “simple application” first: can we fit multi-modal target distributions?**

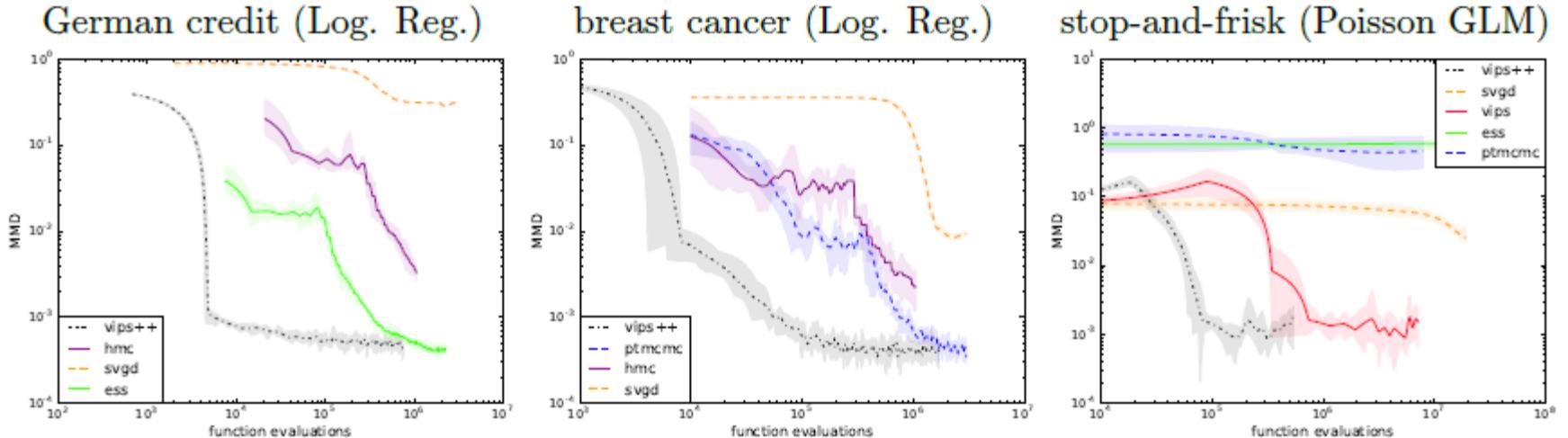
- Instead of maximizing the reward, we can minimize the I-projection

$$\arg \min_p \text{KL}(p||p^*) = \int p(\mathbf{x}) \log p^*(\mathbf{x}) + H(p)$$

- Same information geometric-methods can be used [Arenz & Neumann, ICML 2018]

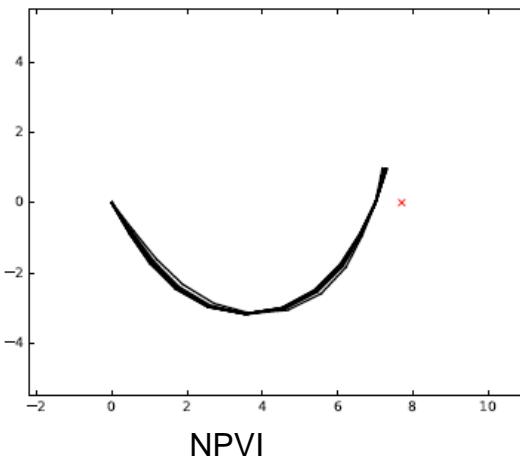
✓ Efficient

✓ Can find multiple modes

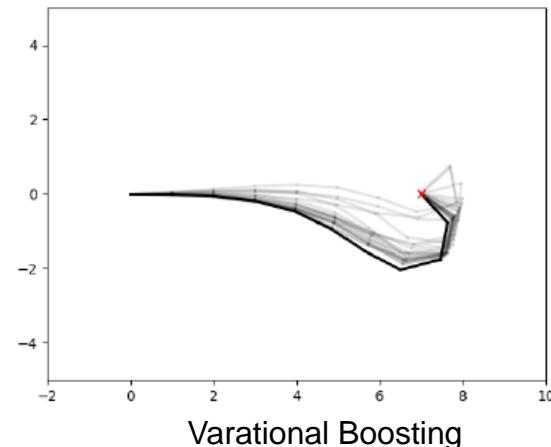


# Robotic Inspired Tasks

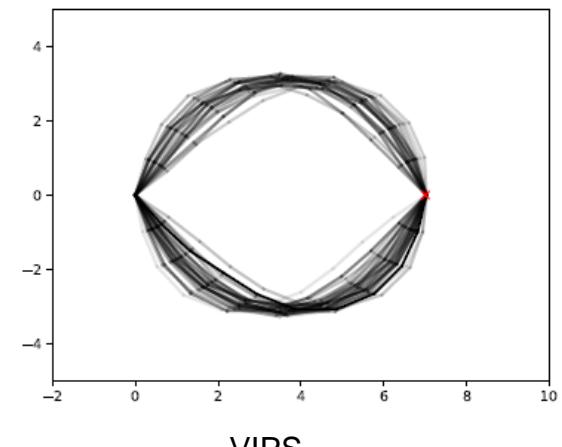
Planar 10-link robot arm:



NPVI

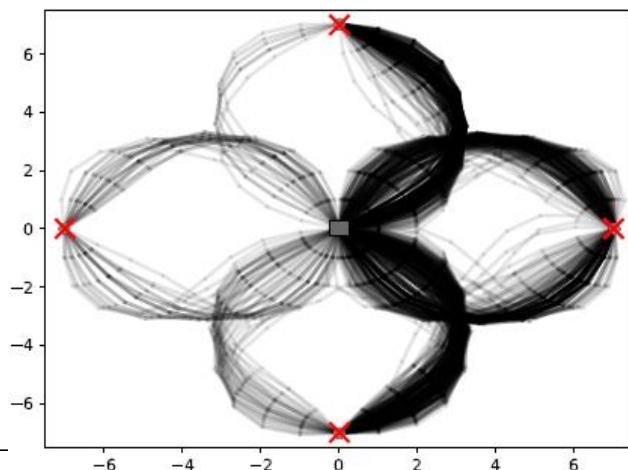


Variational Boosting



VIPS

4 Goals:



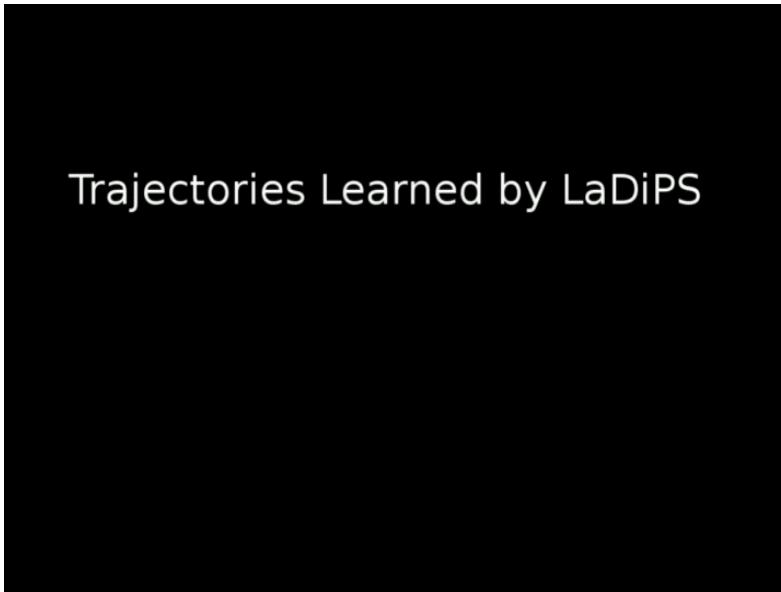
- ✓ Can find close to all solutions
- ✓ Policy search beats classical VI  
(faster, more sample efficient,  
much better approximations)

# Learning multiple skills

Many tasks can be performed in multiple ways

## Table-Tennis:

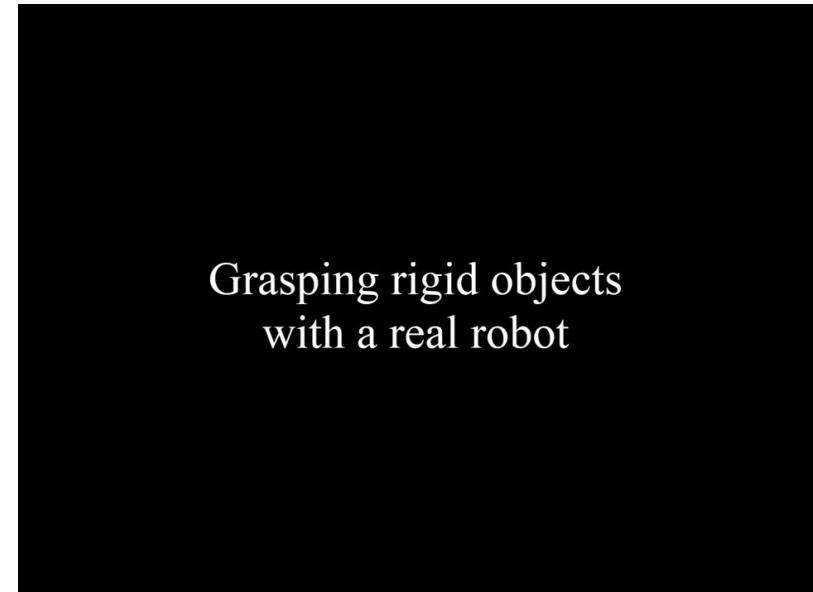
- **Context:** Initial ball configuration
- **Skill Library:** Forehand and Backhand
- **Skill Parametrization:** Movement primitives



Trajectories Learned by LaDiPS

## Multiple Grasp Types:

- **Context:** Local point-cloud features
- **Skill Library:** Grasp type
- **Skill Parametrization:** Grasp configuration



Grasping rigid objects  
with a real robot

F. End, R. Akroud, **G. Neumann**, *Layered Policy Search for Learning Hierarchical Skills*, ICRA 2017

T., Osa, J. Peters, **G. Neumann (2016)**: Experiments with Hierarchical Reinforcement Learning for Learning Multiple Grasp Types, ISER 2016

# Policy Search for Grasping

**Parameters encode:**

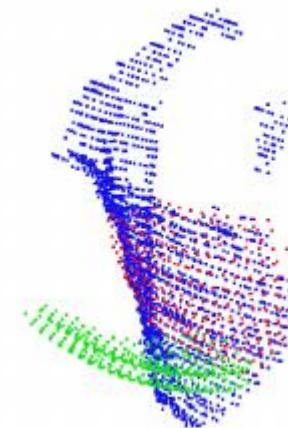
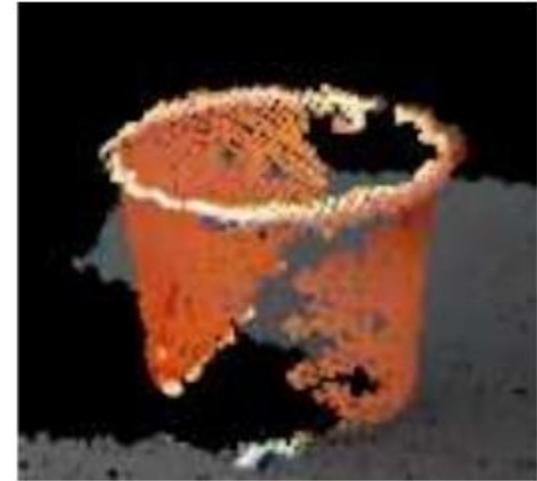
- Pre-Grasp configuration
- Grasp approach direction

**Context: 3D point cloud ?**

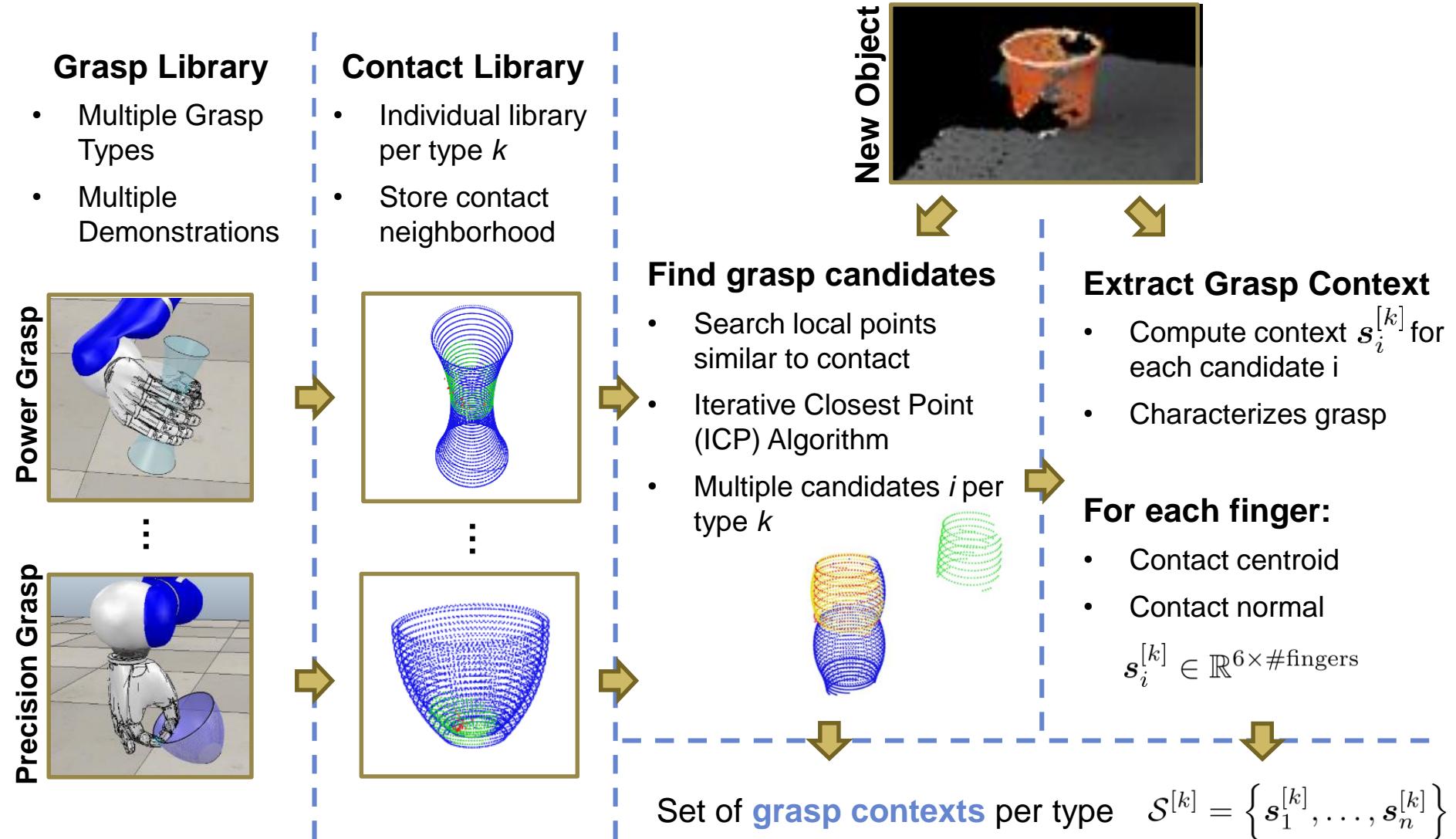
- Very high dimensional input space
- Incomplete
- Unsuitable for context

**However:**

- Only **local features** at contact are important for a specific grasp
- **Grasp candidate:** contact region per finger
- Each grasp candidate has **individual context vector**



# Context Generation



# Hierarchical Policy Search

Set of grasp candidates  $\mathcal{S}^{[k]} = \{s_1^{[k]}, \dots, s_n^{[k]}\}$

Upper Level

## Grasp Selection:

- Learn reward model per type:  $R_k(s) \approx \mathcal{GP}(0, k(s, s'))$
- Evaluate mean and variance of candidates:  $[R_i^{[k]}, \sigma_i^{[k]}] = R_k(s_i^{[k]})$
- Select with Upper Confidence Bounds (UCB):  $[k^*, i^*] = \arg \max_{[k,i]} R_i^{[k]} + \lambda \sigma_i^{[k]}$

Lower Level

## Grasp Adaptation: Type 1

- Learn kernelized contextual policy
- Use 12 states as prototypes

$$\pi_1(\theta|s) = \mathcal{N} \left( \sum_j \alpha_{1,j} k(s, s_j), \Sigma_1 \right)$$

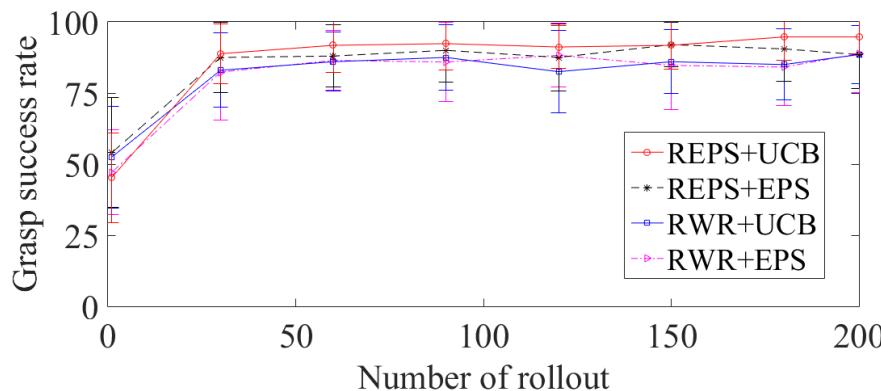
## Grasp Adaptation: Type k

- Learn kernelized contextual policy

$$\pi_k(\theta|s) = \mathcal{N} \left( \sum_j \alpha_{k,j} k(s, s_j), \Sigma_k \right)$$

Bayesian Optimization  
Contextual Policy Search

# Performance in simulation

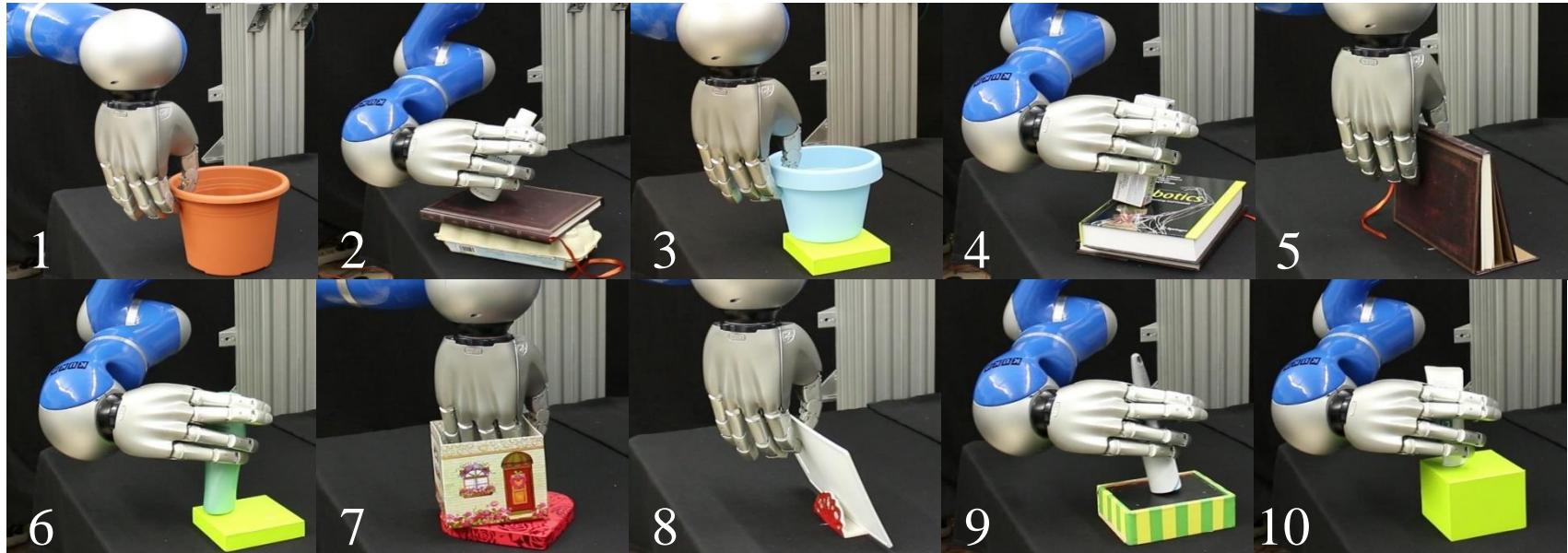


Learning phase  
-beginning-

- **Reward:**
  - Force closure
  - Success of lifting
- 12 initial demonstrations per grasp type
- 50% successrate before learning
- 92.5 % successrate after learning

$$s = [c_{\text{thumb}}, n_{\text{thumb}}, c_{\text{index}}, n_{\text{index}}]$$

# Testing on a real robot system

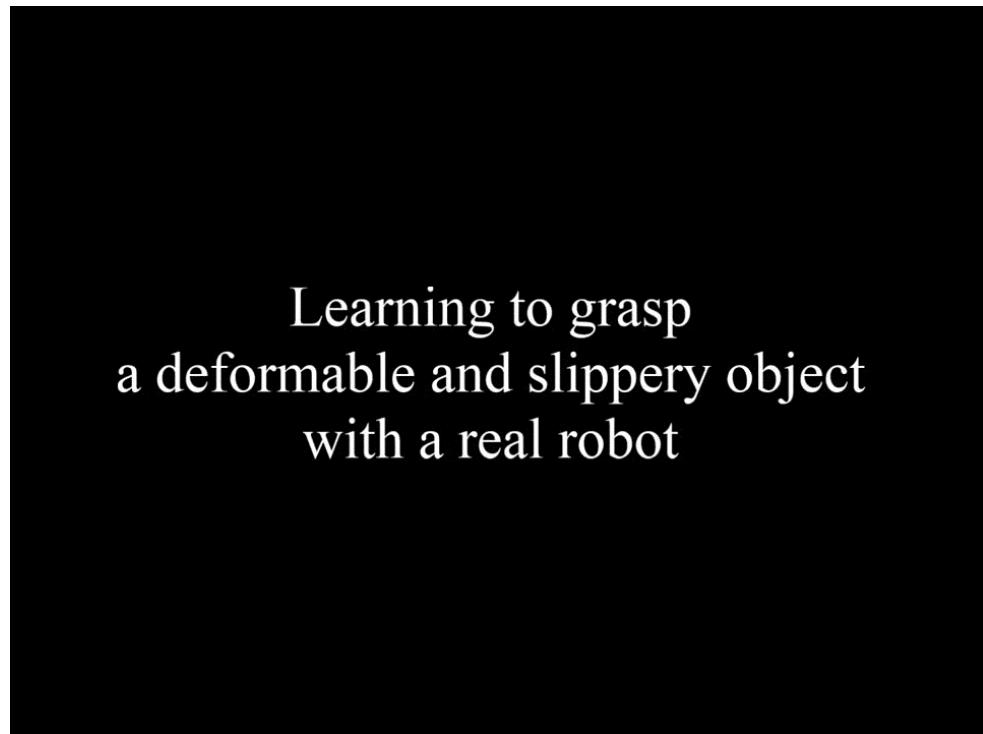
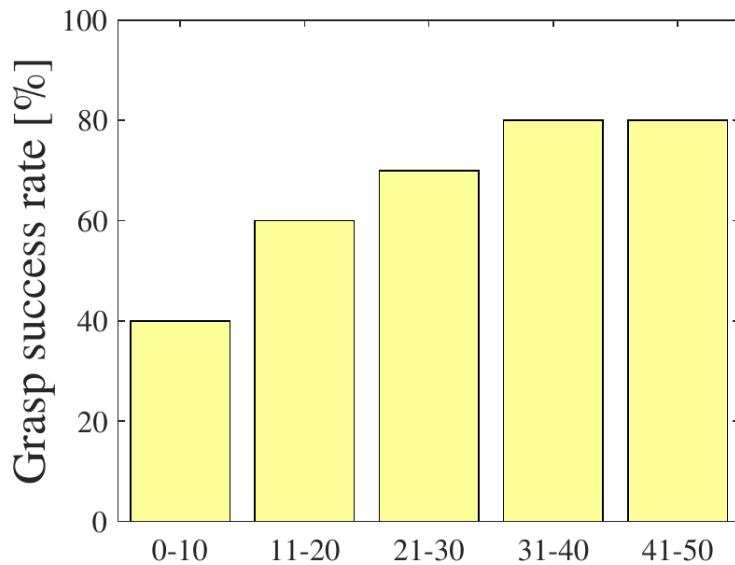


Object No.	1	2	3	4	5	6	7	8	9	10	Total
Success rate	5/5	5/5	4/5	3/5	5/5	5/5	4/5	5/5	5/5	5/5	92.0%

# Learning on a real robot system

**Many objects are deformable/slippery**

- Can we also learn grasp policies in such case?



# Application to Deep RL

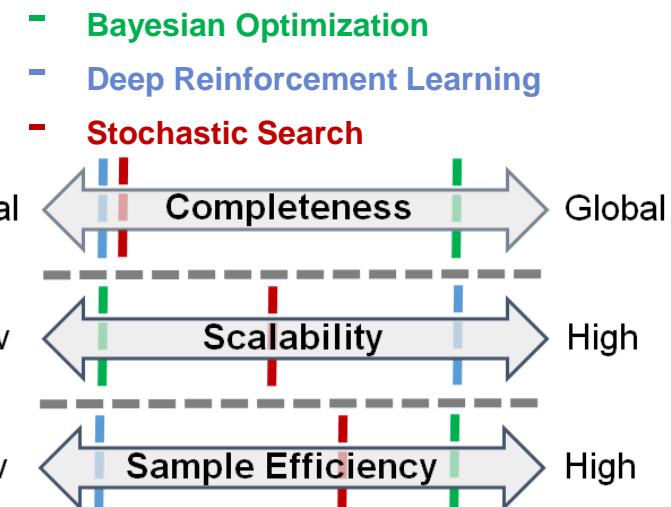
**Deep Policies:**  $\pi(a|s) \propto \exp(\phi_\beta(s, a)^T \theta)$

- $\phi_\beta(s, a)$  ... neural network with params  $\beta$
- $\theta$  ... log-linear parameters

**Objective:**

$$\pi_{\text{new}}^* = \operatorname{argmax}_\pi \mathbb{E}_{\mu_{\pi_{\text{old}}}(s)} \left[ \int \pi(a|s) A^{\pi_{\text{old}}}(s, a) da \right]$$

- where  $A^{\pi_{\text{old}}}(s, a)$  is the advantage function of the old policy



Some SOTA algorithms use similar trust regions (but use approximations to solve it)

- TRPO [Schulman2015], PPO [Schulman2016], MPO [Abdolmaleki2018]

Some SOTA algorithms use similar entropy regularization (not a bound)

- SAC [Harnoja2017], AC3

**Can we unify this using our framework?**

# Compatible Policy Search

**Compatible policy search (CPOS):**

$$\pi_{\text{new}}^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\mu_{\pi_{\text{old}}}(\mathbf{s})} \left[ \int \pi(\mathbf{a}|\mathbf{s}) A^{\pi_{\text{old}}}(\mathbf{s}, \mathbf{a}) d\mathbf{a} \right]$$

$$\text{s.t.: } \mathbb{E}_{\mu_{\pi_{\text{old}}}(\mathbf{s})} [\text{KL}(\pi(\cdot|\mathbf{s}) || \pi_{\text{old}}(\cdot|\mathbf{s}))] \leq \epsilon$$

$$H(\pi_{\text{old}}) - H(\pi) \leq \gamma$$

**Solution:**

$$\pi(\mathbf{a}|\mathbf{s}) \propto \pi_{\text{old}}(\mathbf{a}|\mathbf{s})^{\frac{\eta}{\eta+\omega}} \exp\left(\frac{A(\mathbf{s}, \mathbf{a})}{\eta + \omega}\right)$$

- Not a direct solution as it is hard to sample from this policy
- Yet, we can exploit special (compatible) structure of  $A(\mathbf{s}, \mathbf{a})$

# Compatible Policy Search

Solution (for log linear case):

$$\pi(a|s) \propto \pi_{\text{old}}(a|s)^{\frac{\eta}{\eta+\omega}} \exp\left(\frac{A(s, a)}{\eta + \omega}\right)$$



1. Move  $\pi_{\text{old}}$  into exp:

$$\pi(a|s) \propto \exp\left(\frac{\eta \log \pi_{\text{old}}(a|s) + A(s, a)}{\eta + \omega}\right)$$



2. Log Linear Model:

$$\pi(a|s) \propto \exp\left(\frac{\eta \phi_{\beta}(s, a)^T \theta_{\text{old}} + A(s, a)}{\eta + \omega}\right)$$



3. Compatible approximation:

$$\pi(a|s) \propto \exp\left(\phi_{\beta}(s, a)^T \frac{\eta \theta_{\text{old}} + w_A}{\eta + \omega}\right)$$

Log-Linear Models:

$$\log \pi_{\text{old}}(a|s) = \phi_{\beta}(s, a)^T \theta_{\text{old}} - \log Z$$

Compatible Function Approximation

$$A(s, a) \approx \phi_{\beta}(s, a)^T w_A$$

- For **Gaussians**: all **quadratic**, **linear** and **constant** terms

Compatible Policy Update:

$$\theta_{\text{new}} = \frac{\eta \theta_{\text{old}} + w_A}{\eta + \omega}$$

- Can also be obtained for non-linear term using Taylor expansion
- Equivalent to natural gradients if omega = 0 (no entropy bound)

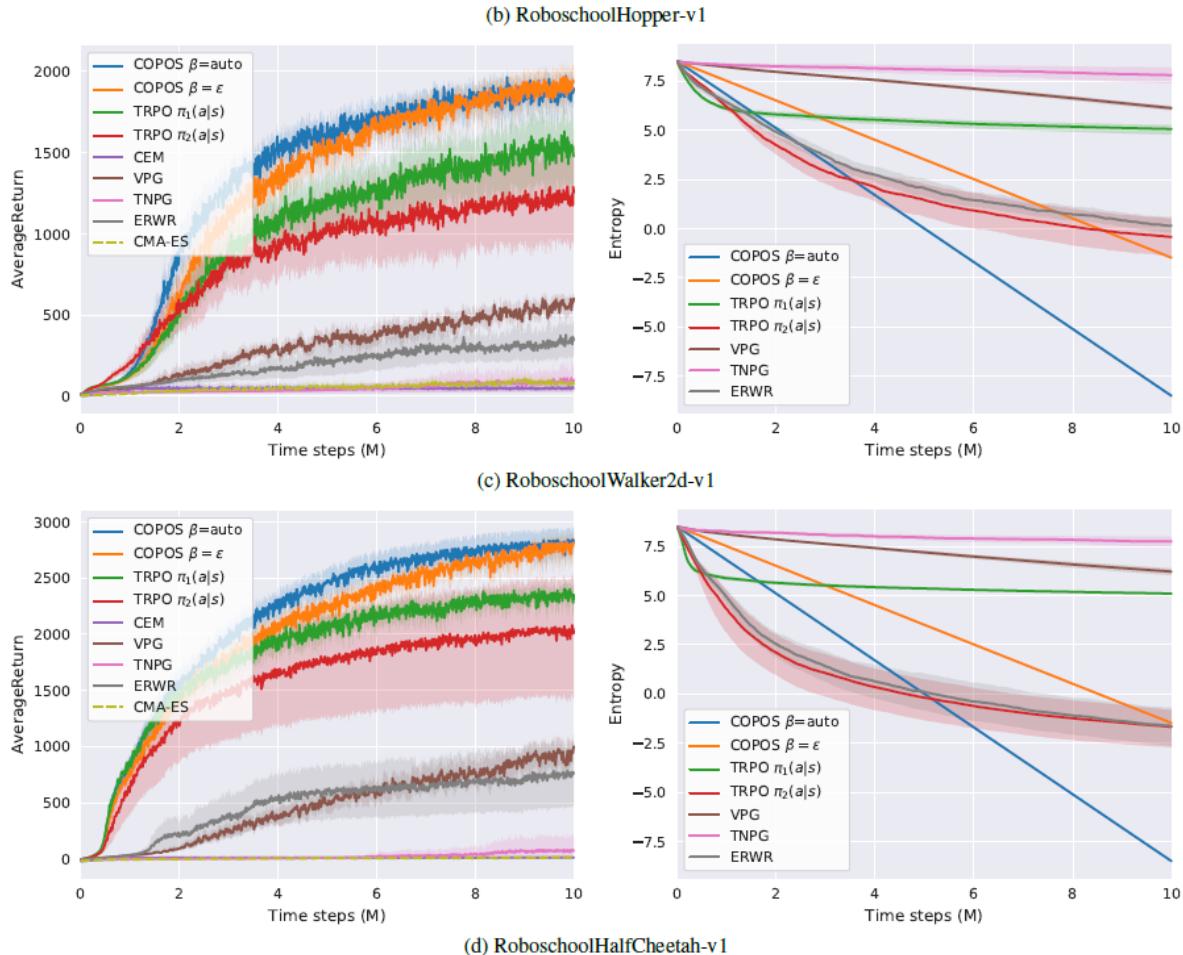
# Compatible Policy Search

## COPOS:

- Can compute KL and entropy constraints in closed form
- Entropy can be directly controlled
- Leads to improved exploration performance

## But:

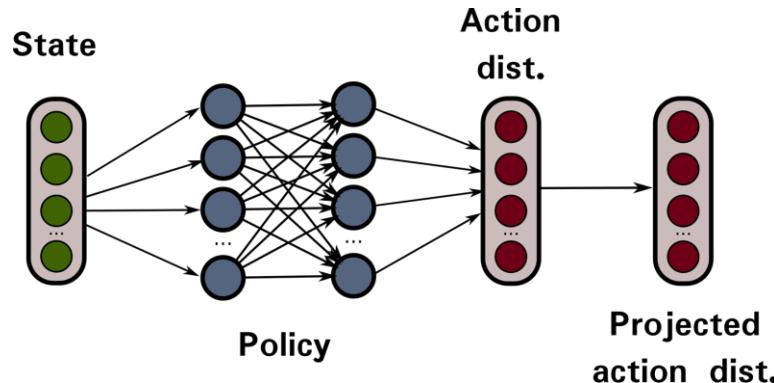
- Still uses approximations for non-linear part of NN!
- Compatible approximation quite complex to implement



# Projections in Policy search

Can we "control" the entropy of the policy in a simpler way?

- Projections hard constraining Shannon **entropy** of **Gaussian** or **soft-max** policies

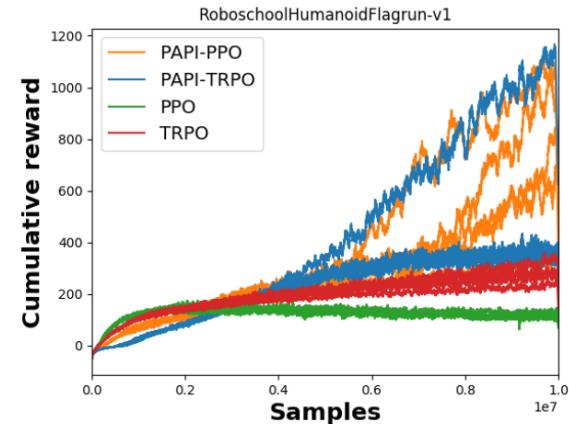
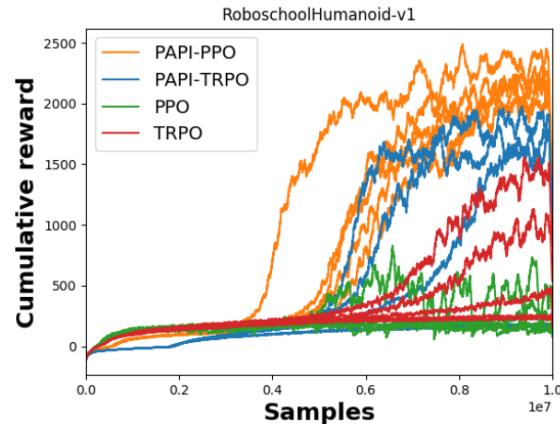


- Project solution of RL algorithm on the solution space given by the constraint
- Can be used on top of any other RL algorithm (TRPO, PPO,...)
- Projections that outperform other **KL**-constrained optimizers used in deep RL

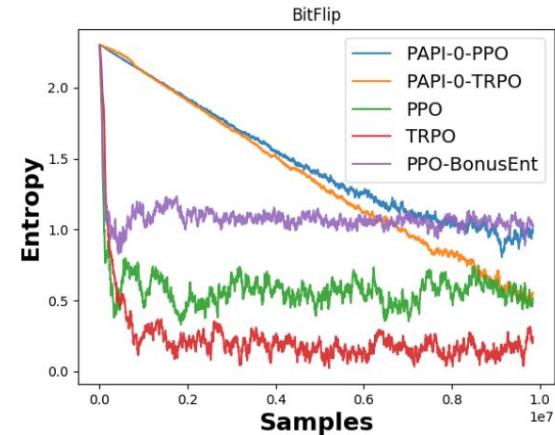
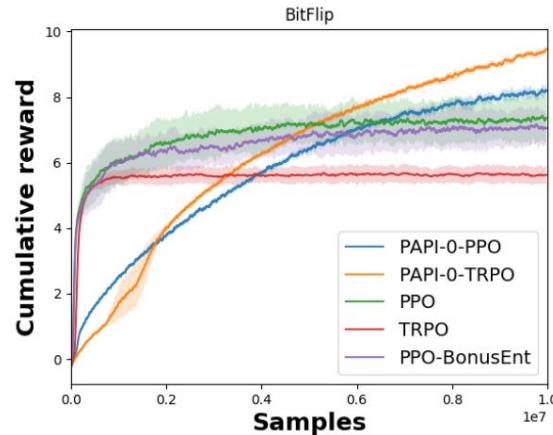
Riad Akroud, Joni Pajarinen, Gerhard Neumann, Jan Peters, "Projections for Approximate Policy Iteration Algorithms"

# Results for PAPI

Mojoco:



Complex exploration  
problem (bit flip):



# Conclusion: Motor Skill Learning

---

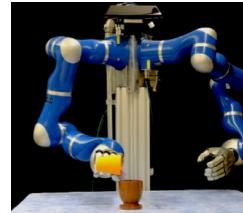
## Information-Geometric Motor Skill Learning:

- Efficient [local search](#) with Gaussian distributions
  - Entropy Loss Regularization
  - Compatible Function Approximation
- Skill [generalization](#) using contextual search
- Local point cloud features for grasp adaptation
- Hierarchical policy for grasp selection
- **Extends to deep RL:** Controlling entropy is key for good exploration!

# How can we accomplish these challenges?

## Machine Learning

## Data-Driven Algorithms



## Autonomous Robots

## Domain Knowledge

Data-driven movement primitives

✗ Flexible Skills

Reinforcement Learning

✗ Self-improvement

Multi-Agent Learning

✗ Cooperation

Representation Learning

✗ Perception

# Learning in Robot Swarms

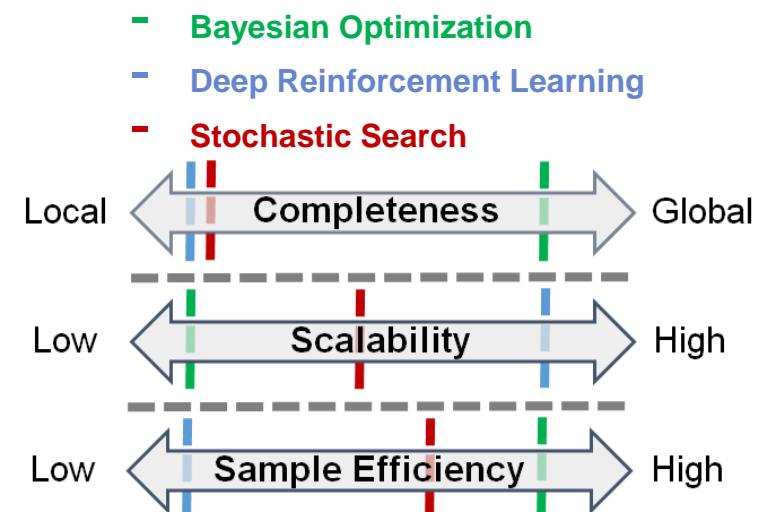
## Swarm Systems:

- Large number of heterogeneous agents
- Permutation invariance (no order of agents)
- Invariant to exact number of agents



## Deep Reinforcement Learning:

- Trust Region Policy Optimization (TRPO):  
 $s_i$
  - Use transitions from all agents to estimate gradient / Policies are shared across agents
- ✓ Can deal with complex observations
- ✗ Can not deal with varying number of agents
- ✗ Can not deal with permutation invariance



# Learning in Robot Swarms

## Swarm Systems:

- Large number of heterogeneous agents
- Learn decentralized control policies
- Each agent is described by state  $s_i$



**Local Observation Model:** Agent can observe (relative) state of their neighbors

$$o^{i,j} = f(s_i, s_j) \quad \text{E.g., distance, bearing to neighbours}$$

**Local Communication Model:** Agents can communicate further information about the interaction graph

$$o^{i,j} = f(s_i, s_j, \mathcal{G}) \quad \text{E.g., number of neighbors of agent j}$$

**Observation Set:**

$$O^i = \{o^{i,j} | j \in \mathcal{N}(i)\}$$

# Communication Models

**Local Observation Model:** Agent can observe (relative) state of their neighbors

$$o^{i,j} = f(s_i, s_j) \quad \text{E.g., distance, bearing to neighbours}$$

**Local Communication Model:** Agents can communicate further information about the interaction graph

$$o^{i,j} = f(s_i, s_j, \mathcal{G}) \quad \text{E.g., number of neighbors of agent j}$$

- Observation Set:**
- $$O^i = \{o^{i,j} | j \in \mathcal{N}(i)\}$$
- Similar properties than samples of distribution

---

**Mean Embedding of distribution**  $p_i(o|s) :$      $\mu_o^i = \frac{1}{|O^i|} \sum_{o^{i,j} \in O^i} \varphi(o^{i,j})$

**Feature vector:**  $\varphi(o)$

- Defined by kernels:               $k(o_i, o_j) = \langle \varphi(o_i), \varphi(o_j) \rangle$

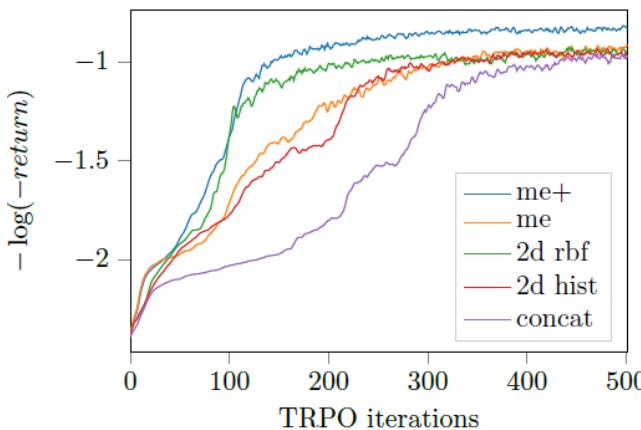
[Gebhard & Neumann, ICRA 2018]

- Defined by neural networks:     $\varphi(o) = \phi_\theta(o)$

[Huettenrauch & Neumann, submitted to JMLR]

**Used as input to neural network policy**

# Experiments: Pursuit Evasion

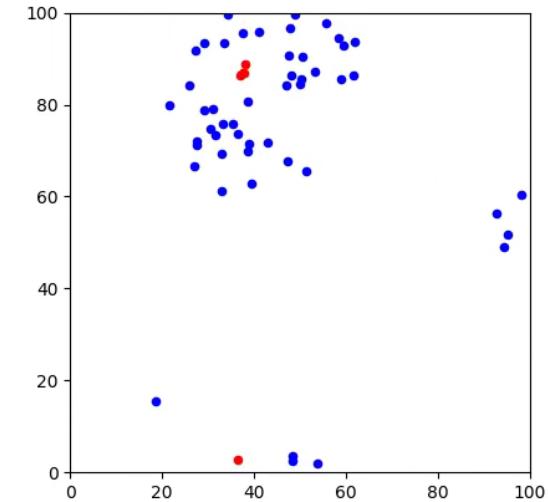


10 agents

## Conclusion:

- New approach for “set inputs”
- Distribution-based representation of agents
- Scales to a large number of agents
- More potential applications

50 agents / 5 evader



# How can we accomplish these challenges?

## Machine Learning

## Data-Driven Algorithms



## Autonomous Robots

## Domain Knowledge

Data-driven movement primitives

✗ Flexible Skills

Reinforcement Learning

✗ Self-improvement

Multi-Agent Learning

✗ Cooperation

Representation Learning

✗ Perception

# Recurrent Kalman Networks

**Goal:** State estimation for dynamical systems:

- Filtering
- Prediction



First step for decision making under uncertainty



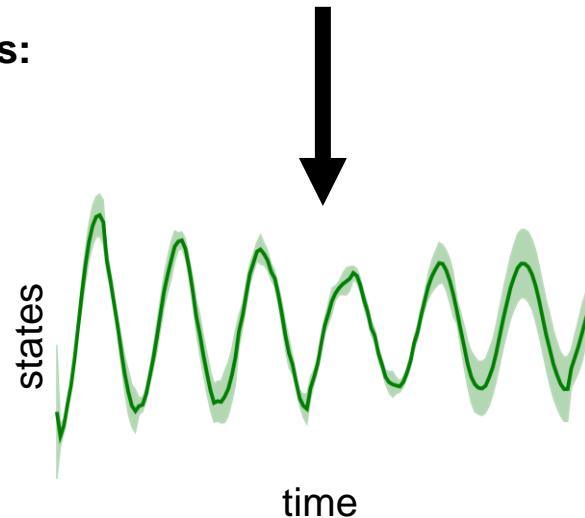
**Challenges:**

- High dimensional observations
- Partially observable
- Nonlinear dynamics
- Uncertainty



**(Deep Learning) Solutions:**

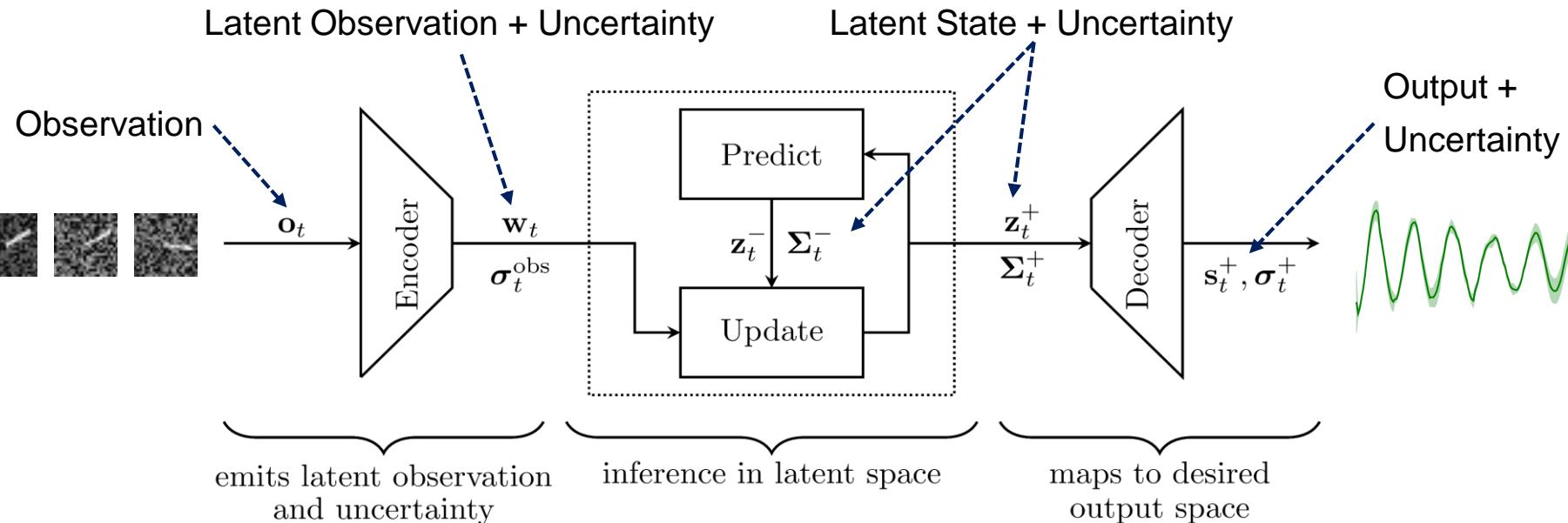
- CNNs
- RNNs
- Variational Inference (approximate)



**How can we propagate uncertainty through RNNs without approximations?**

- Recurrent cell based on Kalman filter

# Overview



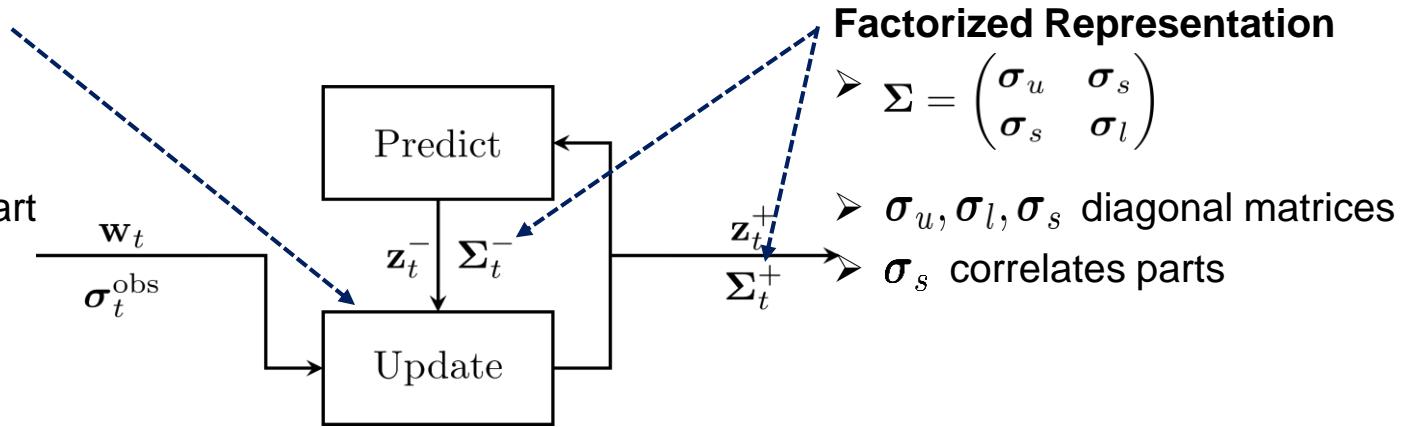
## Make backpropagation through Kalman filter feasible?

- Locally linear transition models, even for highly nonlinear systems
- High dimensional latent spaces
- Factorized state representation to avoid expensive and unstable matrix inversions

# Factorized State Representation

## Observation Model

- $\mathbf{H} = (\mathbf{I}_m \quad \mathbf{0}_{m \times m})$
- Splits latent state
  - 1. Observable part
  - 2. Memory part



## Results in simplified Kalman Update

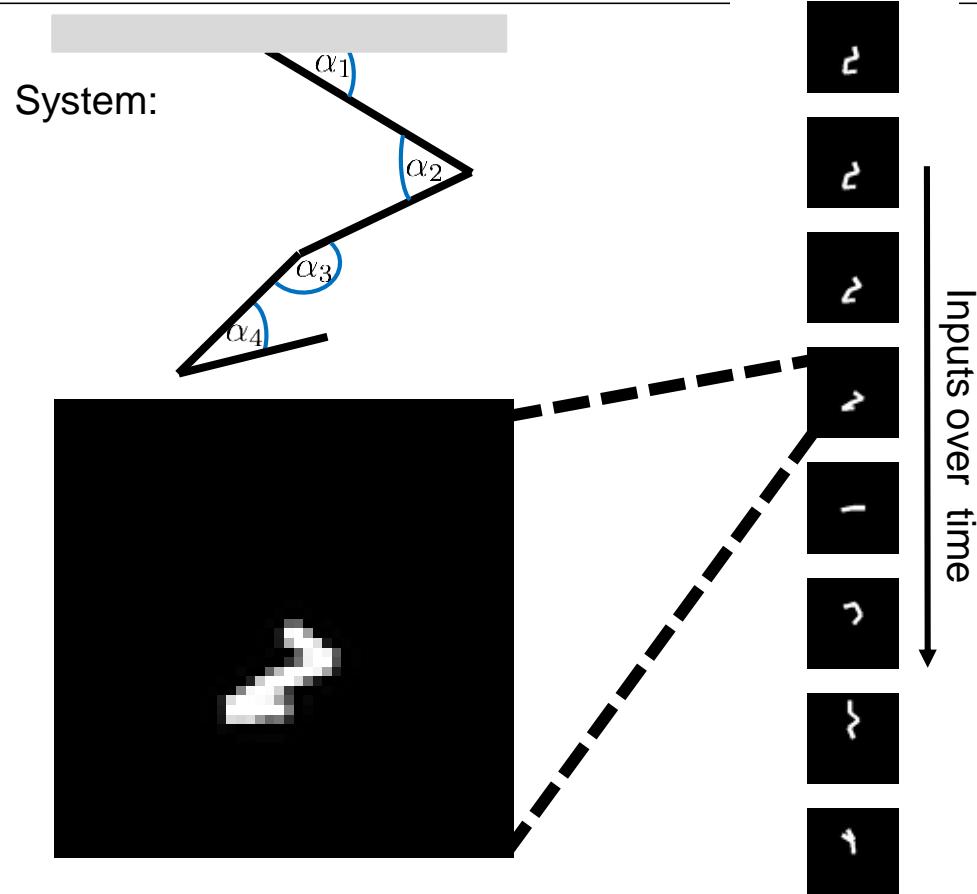
- No matrix inversion
  - Instead only pointwise operations
  - Assumptions not restrictive since latent space is learned
- } Makes inference and back-propagation feasible

# Quad Link Pendulum

- State (4 joint angles + velocity)
- Highly nonlinear dynamics
- Links occlude each other
- Estimate joint angles of all 4 links

	RKN	LSTM	GRU
Log Likelihood	14.534	11.960	10.346
RMSE	0.103	0.118	0.121

- Significantly better uncertainty estimate (higher log-likelihood)
- Better prediction (smaller RMSE)



# Summary & Conclusion

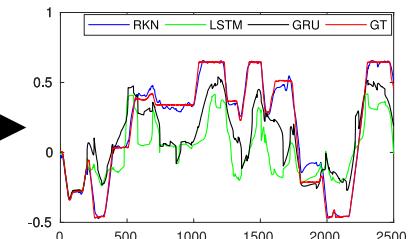
## Recurrent Kalman Networks...

- ... scale to real world systems
- ... allow direct state estimation from images
- ... use uncertainty in a principled manner to handle noise
- ... can be trained end-to-end without approximations

... use as basic building block for decision making under uncertainty in the future

## Additional Experiments

- Pendulum
- Image Imputation
- KITTI-Dataset for visual odometry
- Prediction for real pneumatic joint



- Comparison to recent approaches
- Code available

# Conclusion

---

- Information-theoretic methods are powerful tools to formalize the RL problem
- **Direct Policy Search:**
  - Data-efficient but limited complexity (feedback, perception, ...)
- **Deep RL:**
  - Data hungry but works with complex policies
  - “Controlling” exploration is key for both approaches
  - Still we need to close the “gap” and find the best of both worlds
  - Using “algorithmic priors” in deep learning is a promising solution
    - Represent distributions: Mean embeddings
    - Represent uncertainty: Kalman Filtering

# We are hiring!

## New Robot Manipulation group at the University of Tuebingen:

- Funded by Bosch (BCAI). Competitive Bosch salary

### Topics:

- Robot reinforcement learning
- Imitation learning
- Information theoretic methods
- Grasping and manipulation
- Robot vision
- Human-robot interaction

### Open positions :

- 6 PhD positions
- 2 Research Scientist (permanent positions)
- Many more open positions (PhD + research scientist) @  
BCAI headquarter in Renningen



EBERHARD KARLS  
**UNIVERSITÄT  
TÜBINGEN**

