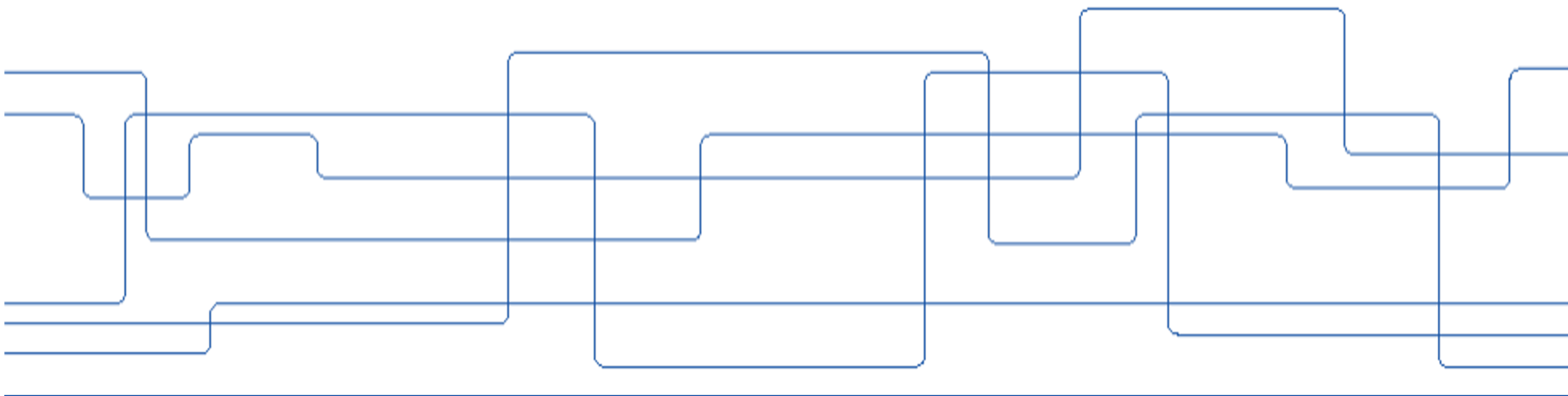


FDD3359 Reinforcement Learning Course

Meta Reinforcement Learning

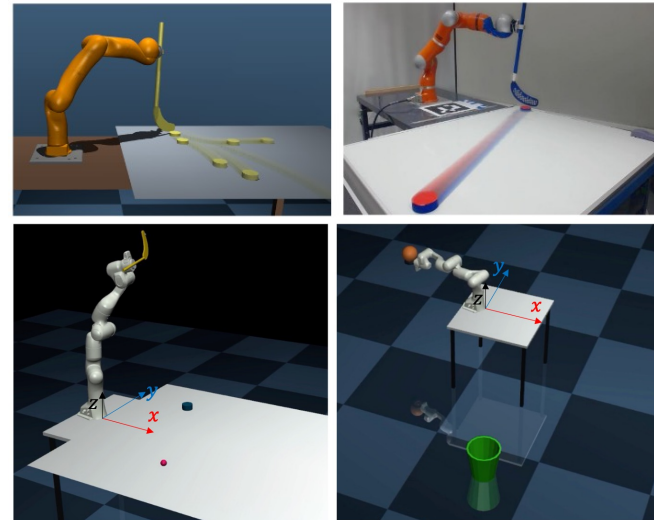
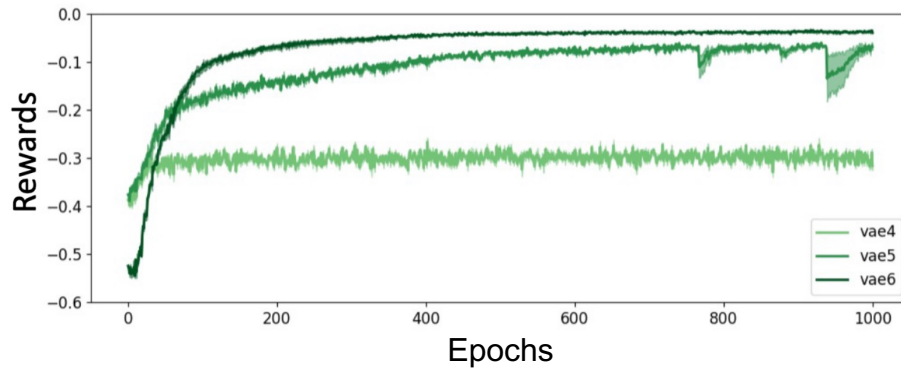
Ali Ghadirzadeh



Why Meta Reinforcement Learning?

Reinforcement Learning

- trains a policy from scratch
- is usually very data-inefficient
- does not usually transfer well

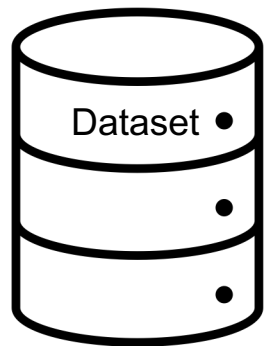


Leveraging prior knowledge to speed up solving new tasks

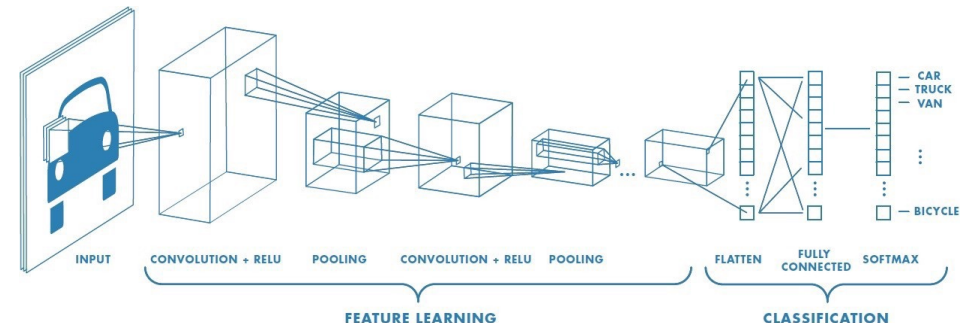
Meta Learning Background

Few-shot Supervised Learning

Standard Machine Learning



- A dataset
- A big neural network
- A loss function
- Some GPUs
- Stochastic gradient descent



Few-shot Learning

Training set consists of 5 classes



Classify novel test inputs



The Meta-Learning Problem
& Black-Box Meta-Learning
Prof. Chelsea Finn



Meta Learning Background

Few-shot Supervised Learning

Some contents are borrowed from
“The Meta-Learning Problem & Black-Box Meta-Learning”
by Prof. Chelsea Finn

Transfer learning via fine-tuning

$$\phi \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}^{\text{tr}})$$

Parameters pre-trained on \mathcal{D}_a

training data for new task \mathcal{T}_b
(typically for many gradient steps)



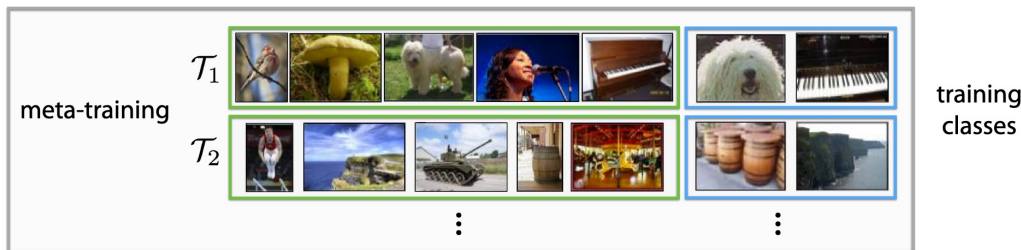
Meta Learning Background

Some contents are borrowed from
 “The Meta-Learning Problem & Black-Box Meta-Learning”
 by Prof. Chelsea Finn

Supervised Learning

$$\phi_i = f_{\theta}(\mathcal{D}_i^{\text{tr}})$$

$$\mathbf{y}^{\text{ts}} = g_{\phi_i}(\mathbf{x}^{\text{ts}})$$



Given 1 example of 5 classes:

Classify new examples

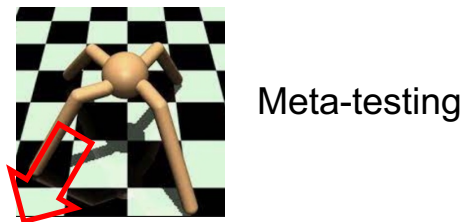
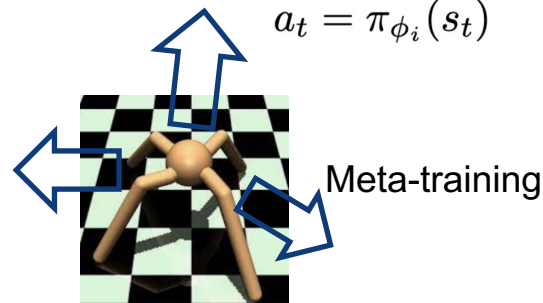


$$\max_{\theta} \sum_{\mathcal{T}_i} \mathcal{L}(f_{\theta}(\mathcal{D}_i^{\text{tr}}), \mathcal{D}_i^{\text{test}})$$

Reinforcement Learning

$$\phi_i = f_{\theta}(\{s_k, a_k, s_{k+1}, r_k\}_{k=0:T})$$

$$a_t = \pi_{\phi_i}(s_t)$$

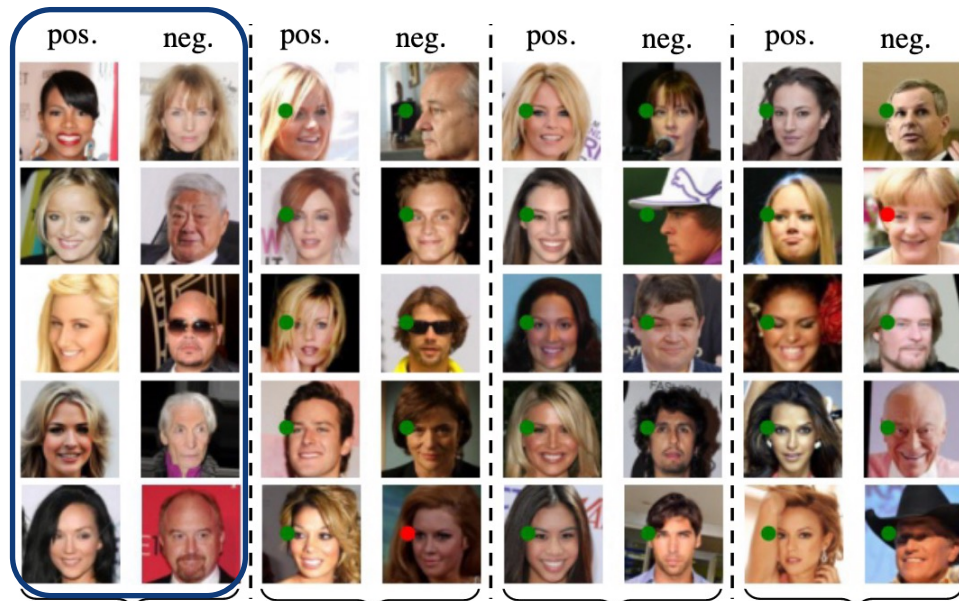


$$\max_{\theta} \sum_{\mathcal{M}_i} \mathbb{E}_{\pi_{\phi_i}} [R(\tau)]$$

Meta Learning Background

Few-shot Supervised Learning

Training data



Few-shot examples

	Task 1	Task 2	Task 3
• makeup	✗	✓	✓
• young	✓	✗	✓
• mouth open	✓	✓	✗

- Probabilistic Formulation
- Active Learning



Meta Reinforcement Learning

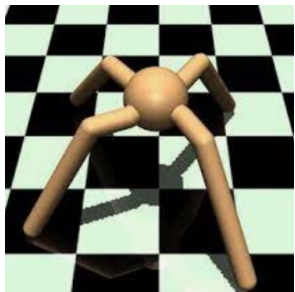
Meta-learning
(Outer loop optimization)

Exploration
Policy



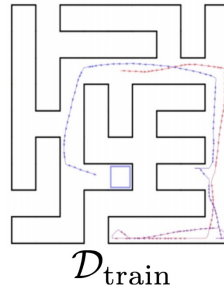
Task Policy

Adaptation
(Inner loop optimization)

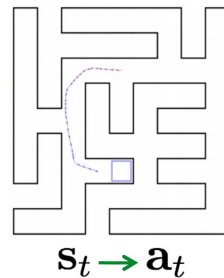


[meta] test time

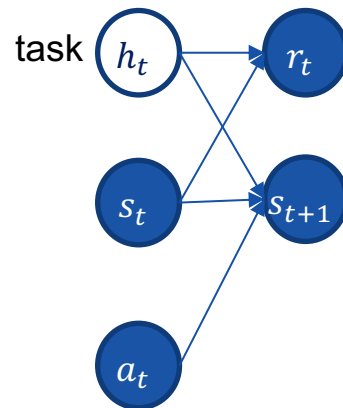
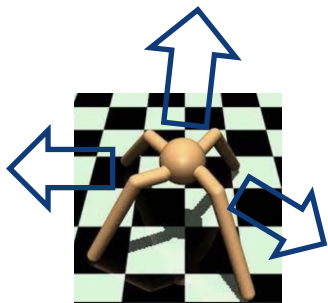
Given a small amount of experience



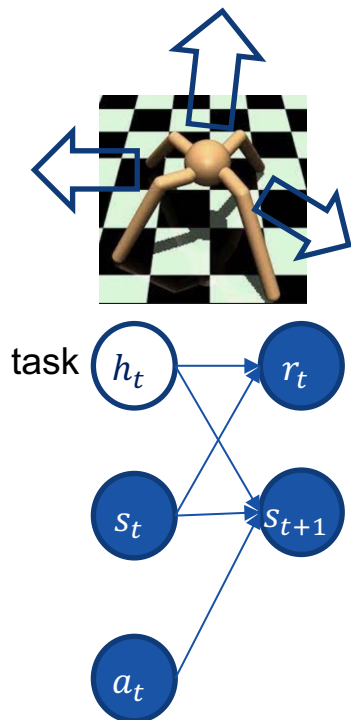
Learn to solve the task



POMDP View



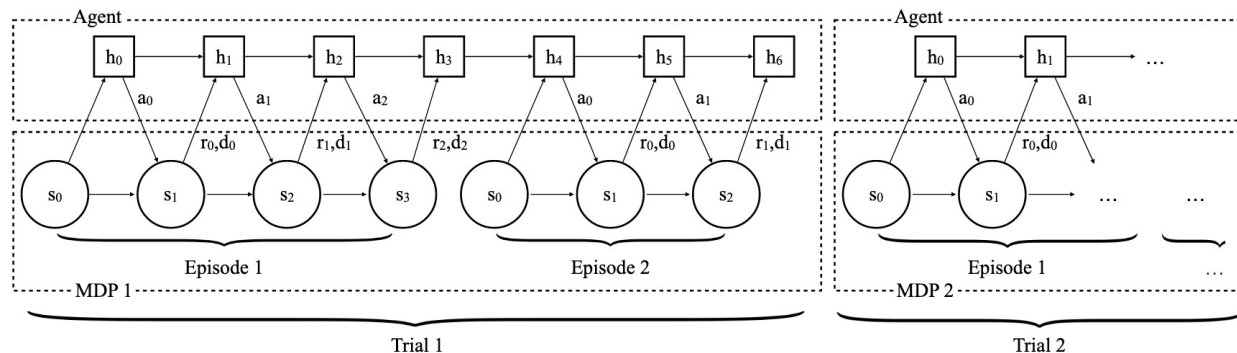
POMDP View



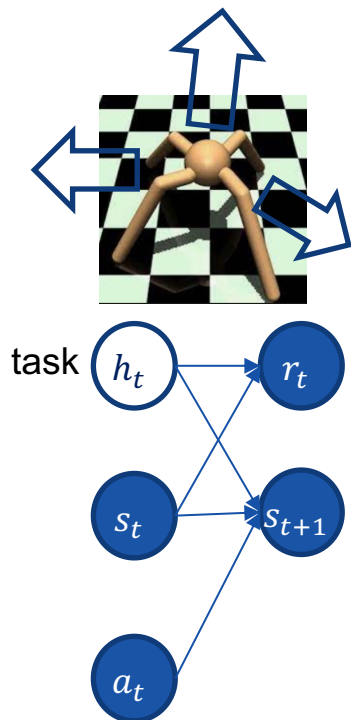
Recurrent Policy Training

RL²: FAST REINFORCEMENT LEARNING VIA SLOW REINFORCEMENT LEARNING

Yan Duan^{†‡}, John Schulman^{†‡}, Xi Chen^{†‡}, Peter L. Bartlett[†], Ilya Sutskever[‡], Pieter Abbeel^{†‡}



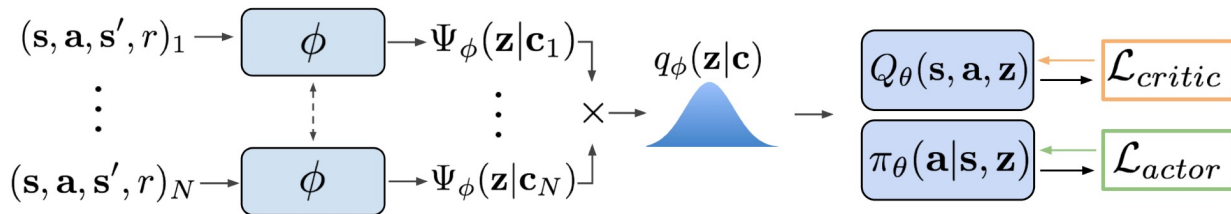
POMDP View



Latent-Variable Policy Training

Efficient Off-Policy Meta-Reinforcement Learning via Probabilistic Context Variables

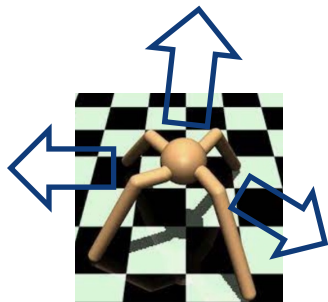
Kate Rakelly^{1*} Aurick Zhou^{1*} Deirdre Quillen¹ Chelsea Finn¹ Sergey Levine¹





Optimization-Based Meta-Learning

Some contents are borrowed from
“Optimization-Based Meta-Learning”
by Prof. Chelsea Finn



Fine-tuning

$$\phi \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}^{\text{tr}})$$

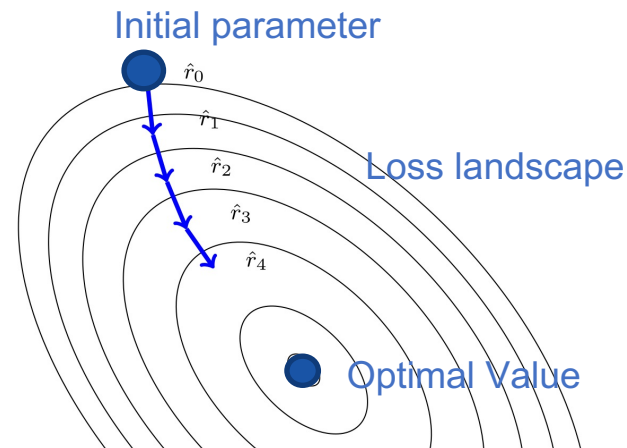
pre-trained parameters

training data for new task

(typically for many gradient steps)

How to improve learning efficiency?

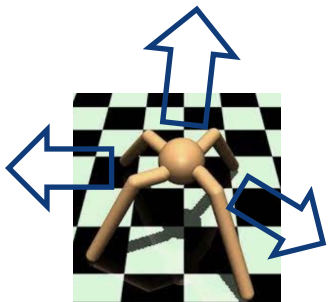
- Parameter Initialization θ
- Learning rate α
- Loss function \mathcal{L}
- Training data \mathcal{D}^{tr}





Optimization-Based Meta-Learning

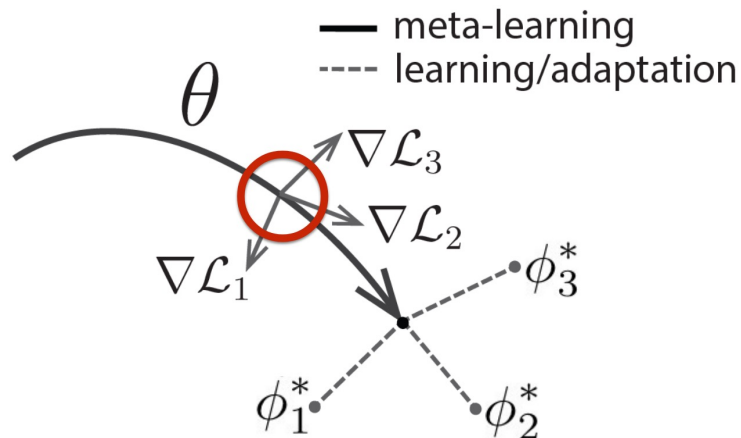
Some contents are borrowed from
“Optimization-Based Meta-Learning”
by Prof. Chelsea Finn



Meta-learning $\min_{\theta} \sum_{\text{task } i} \mathcal{L}(\theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_i^{\text{tr}}), \mathcal{D}_i^{\text{ts}})$

θ parameter vector
being meta-learned

ϕ_i^* optimal parameter
vector for task i



Model-Agnostic Meta-Learning



Optimization-Based Meta-Learning

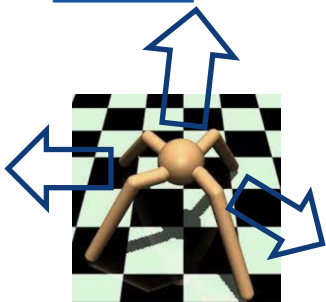
PROMP: PROXIMAL META-POLICY SEARCH

Jonas Rothfuss*
UC Berkeley, KIT
jonas.rothfuss@gmail.com

Dennis Lee*, Ignasi Clavera*
UC Berkeley
{dennis188, iclavera}@berkeley.edu

Tamim Asfour
Karlsruhe Inst. of Technology (KIT)
asfour@kit.edu

Pieter Abbeel
UC Berkeley, Covariant.ai
pabbeel@cs.berkeley.edu

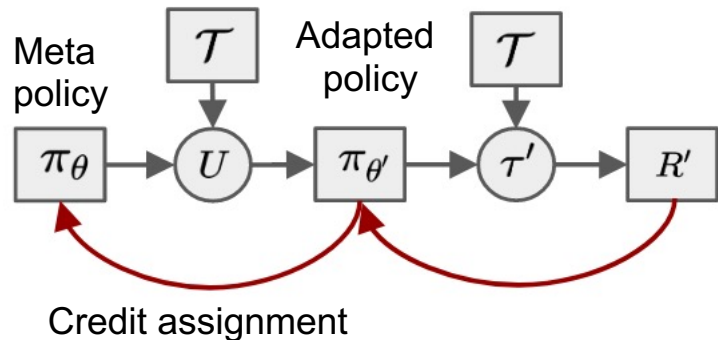


RL Objective

$$J = \mathbb{E}_{\tau \sim \pi_{\theta}} [\tau]$$

Policy Gradient

$$\nabla_{\theta} J_i(\theta) = E_{\tau \sim \pi_{\theta}, \tau_i} \left[\left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_t r_i(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

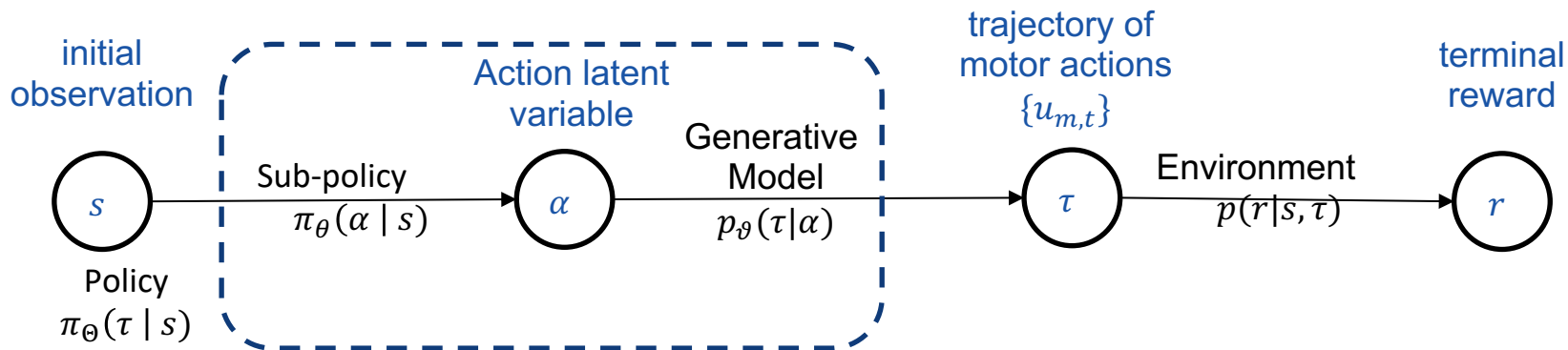


$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \rho(\tau)} \left[\mathbb{E}_{\substack{\tau \sim P_{\tau}(\tau|\theta) \\ \tau' \sim P_{\tau}(\tau'|\theta')}} \left[\nabla_{\theta} J_{\text{post}}(\tau, \tau') + \nabla_{\theta} J_{\text{pre}}(\tau, \tau') \right] \right]$$

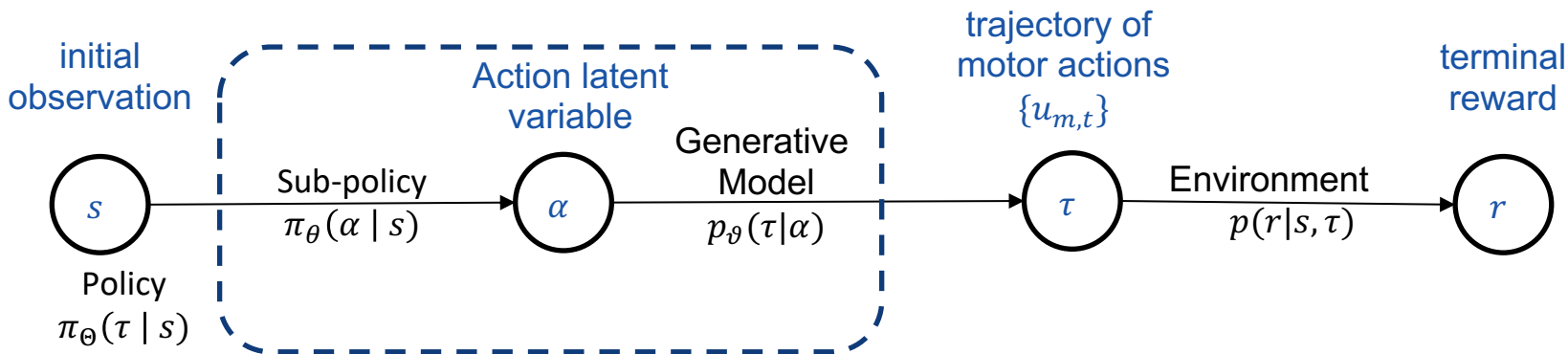
$$\nabla_{\theta} J_{\text{post}}(\tau, \tau') = \underbrace{\nabla_{\theta'} \log \pi_{\theta'}(\tau') R(\tau')}_{\nabla_{\theta'} J_{\text{outer}}} \underbrace{(I + \alpha R(\tau) \nabla_{\theta}^2 \log \pi_{\theta}(\tau))}_{\text{transformation from } \theta' \text{ to } \theta}$$

$$\nabla_{\theta} J_{\text{pre}}^I(\tau, \tau') = \alpha \nabla_{\theta} \log \pi_{\theta}(\tau) \left(\underbrace{(\nabla_{\theta} \log \pi_{\theta}(\tau) R(\tau))^{\top}}_{\nabla_{\theta} J_{\text{inner}}} \underbrace{(\nabla_{\theta'} \log \pi_{\theta'}(\tau') R(\tau'))}_{\nabla_{\theta'} J_{\text{outer}}} \right)$$

RL with Generative Models

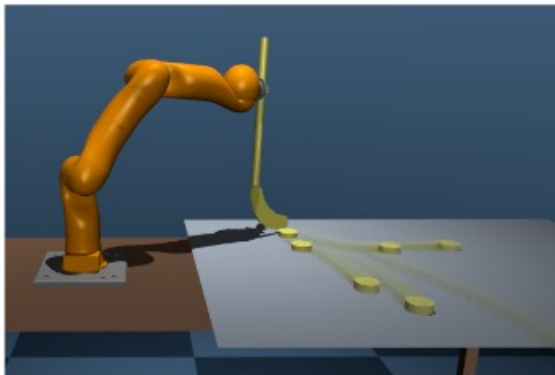
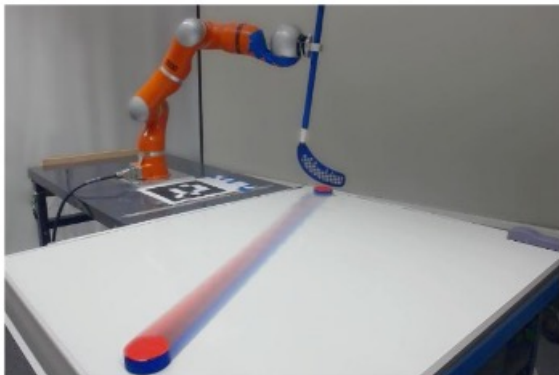


Sim-to-real Transfer Learning



Options:

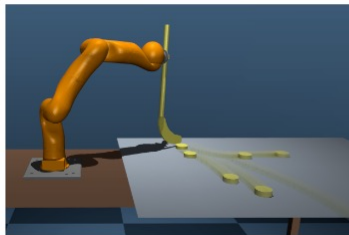
- Learn a general model
- Learn a learning agent



Arndt et al., ICRA20
Meta Reinforcement Learning for
Sim-to-real Domain Adaptation.

Learning a dynamic task efficiently?

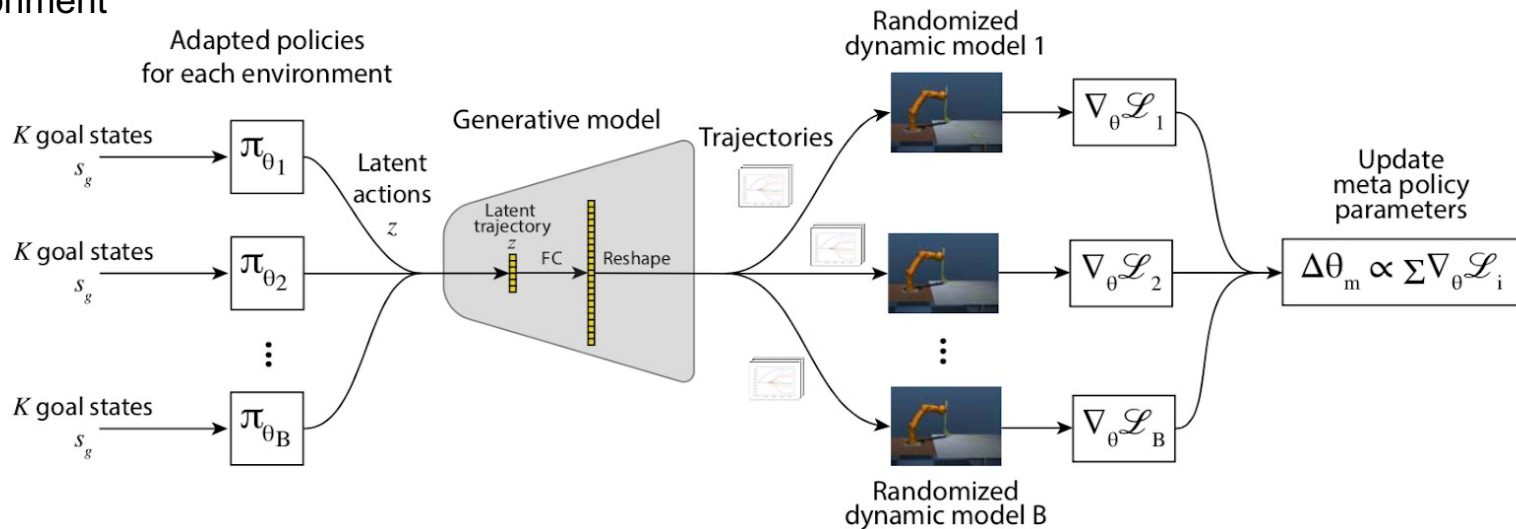
Sim-to-real Transfer Learning



Sim environment

TABLE I: Randomized parameters

Parameter	Minimum	Maximum
x linear friction (μ_x)	0.15	0.95
y linear friction (μ_y)	$0.7\mu_x$	$1.3\mu_x$
Torsional friction (μ_τ)	0.001	0.05
Rotational friction x (μ_{rx})	0.01	0.3
Rotational friction y (μ_{ry})	0.01	0.3
Puck mass	50g	500g
Initial puck position	$\epsilon \sim \mathcal{N}(0, 0.02)$	



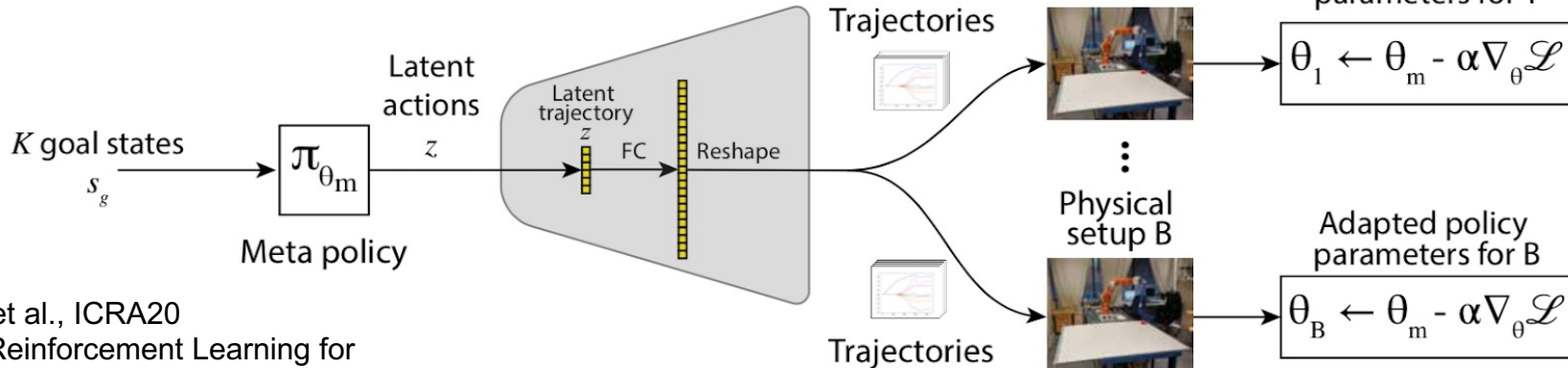
Sim-to-real Transfer Learning



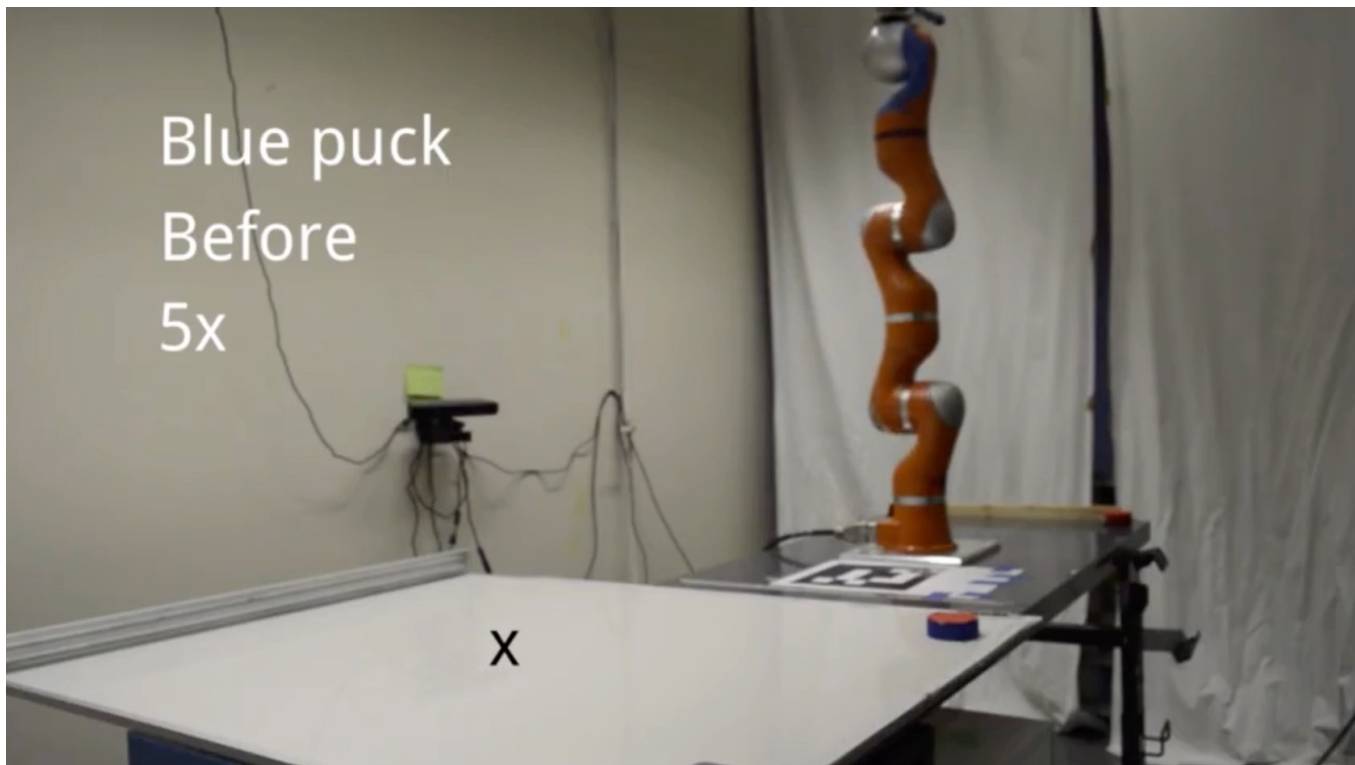
Real environment



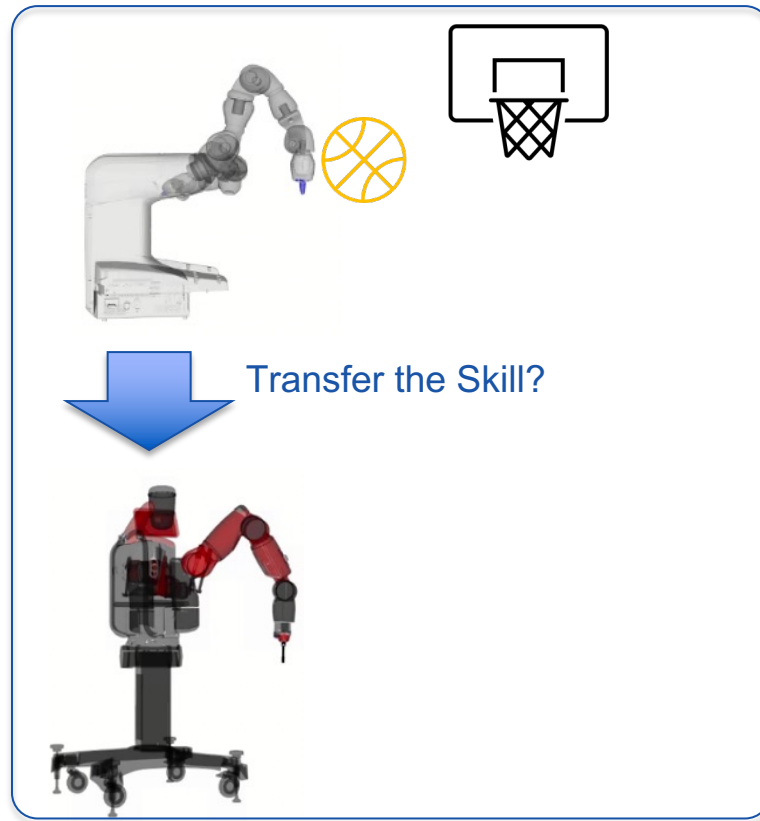
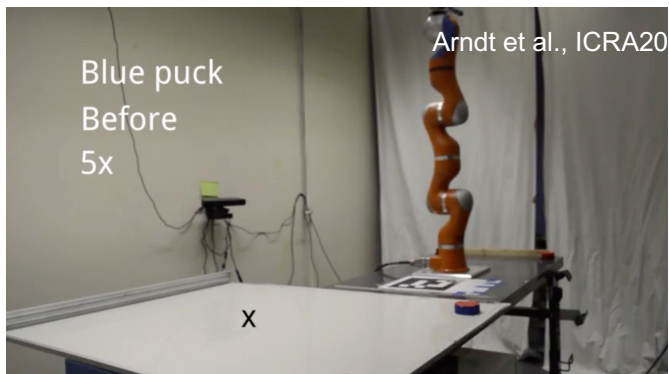
Generative model



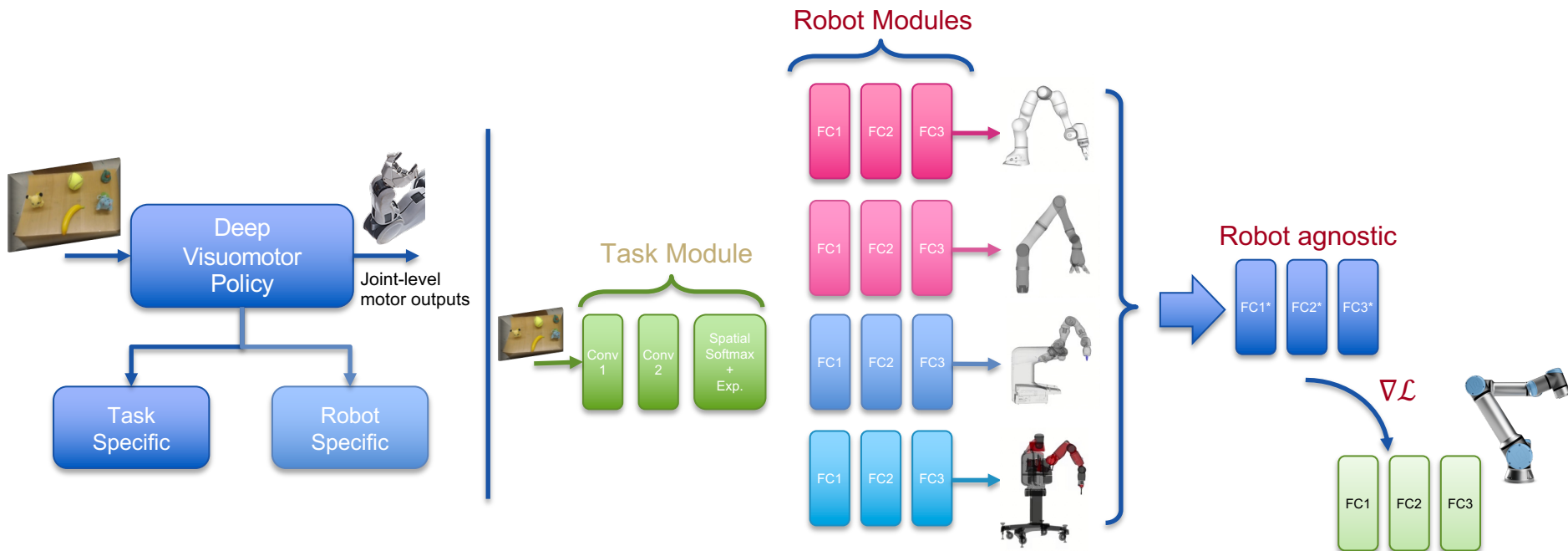
Sim-to-real Transfer Learning



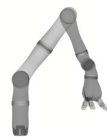
Multi-robot RL



Few-shot adaptation using meta-learning



Multi-robot learning



100 rand. robots
based on Kinova



100 rand. robots
based on YuMi



100 rand. robots
based on Franka



100 rand. robots
based on Baxter

In total 400 different robots in simulation

Initial Training on One Robot



Perception HL Policy Generative Model

- Adversarial feature training for generalizable robotic visuomotor control, Chen et al., ICRA20
- Data-efficient visuomotor policy training using reinforcement learning and generative models, Ghadirzadeh et al., arXiv20

Meta-learning on all 400 robots

- Collect sequence of motor outputs for each robot using, e.g., RRT* planner
- Obtain initialization over generative model parameters by meta-learning

motor outputs
 $\tau = \{u_{1:m,t+1:t+T}\}$

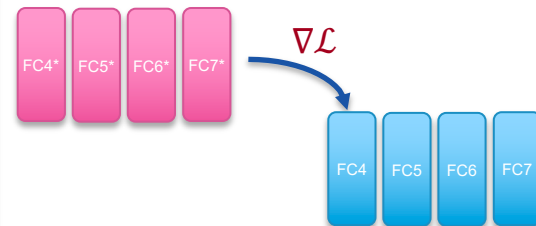


Auxiliary Generative Model
Enc. (Meta-Model)

Meta-testing (adapt to novel robots)



- Provide few demonstration for the novel robot
- Adapt the generative model using few gradient descents



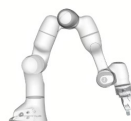
Model-Agnostic Meta-Learning (MAML)



100 rand. robots
based on Kinova



100 rand. robots
based on YuMi



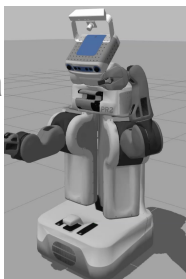
100 rand. robots
based on Franka



100 rand. robots
based on Baxter

For every **robot** i in dataset

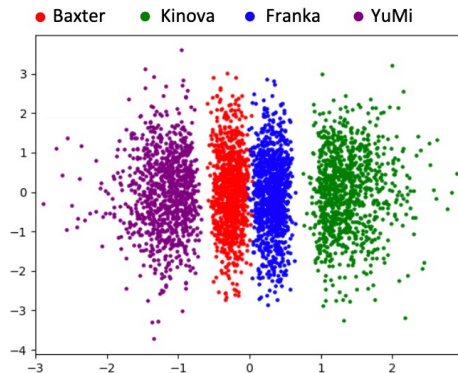
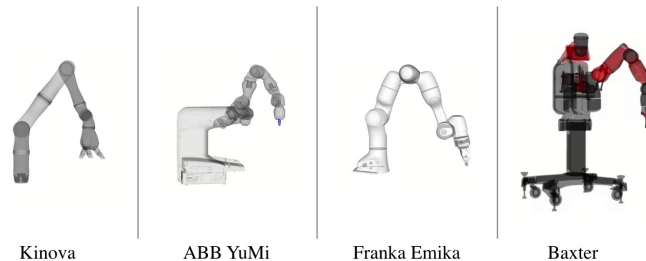
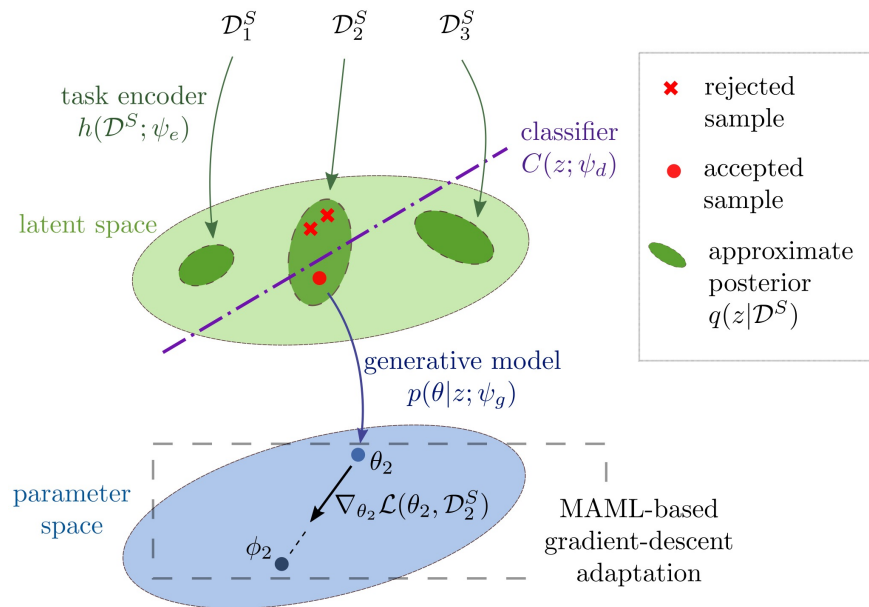
- Collect sequences of motor data
- Split the data and add to
 - Support dataset $\mathcal{D}_{\mathcal{T}_i}^{tr}$
 - Query dataset $\mathcal{D}_{\mathcal{T}_i}^{ts}$



Auxiliary Generative Model
Enc. (Meta-Model)

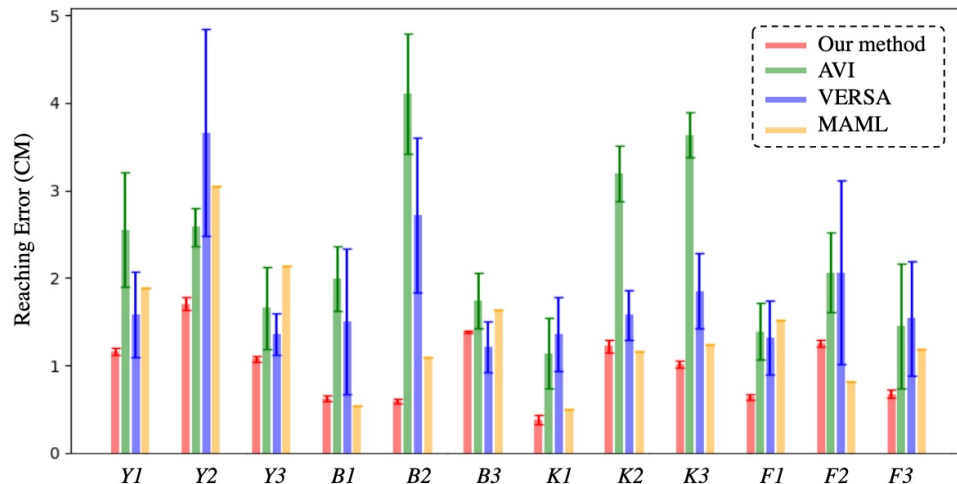
$$\min_{\phi} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}(\phi - \lambda \nabla_{\phi} \mathcal{L}(\phi, \mathcal{D}_{\mathcal{T}_i}^{tr}), \mathcal{D}_{\mathcal{T}_i}^{ts})$$

Conditional and Probabilistic MAML



Experimental Results

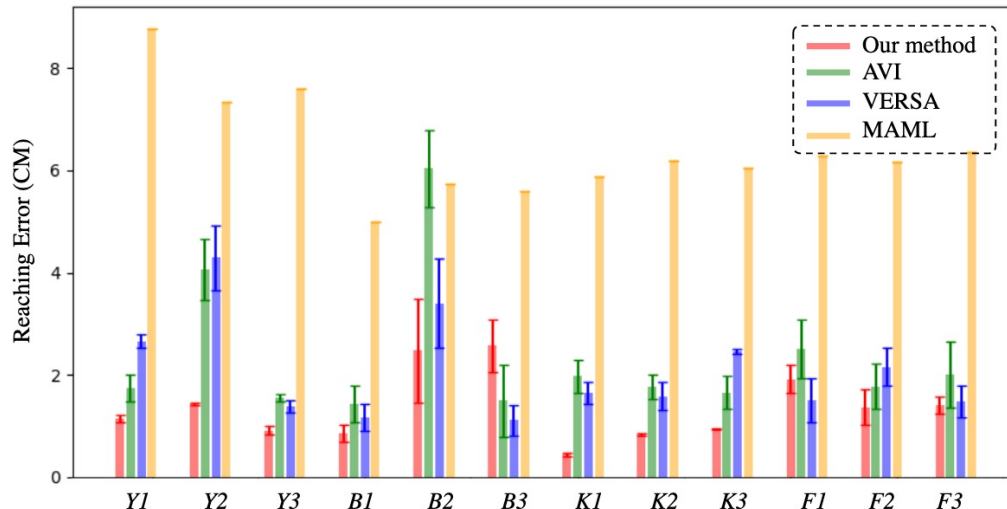
Experiments 1: meta-training and -testing on one platform (100 robots)



The average reaching error (in cm) of the adapted policies.
Y stands for YuMi, B for Baxter, K for Kinova, and F for Franka.

Experimental Results

Experiments 2: meta-training and -testing on all platforms (400 robots)



The average reaching error (in cm) of the adapted policies.
Y stands for YuMi, B for Baxter, K for Kinova, and F for Franka.