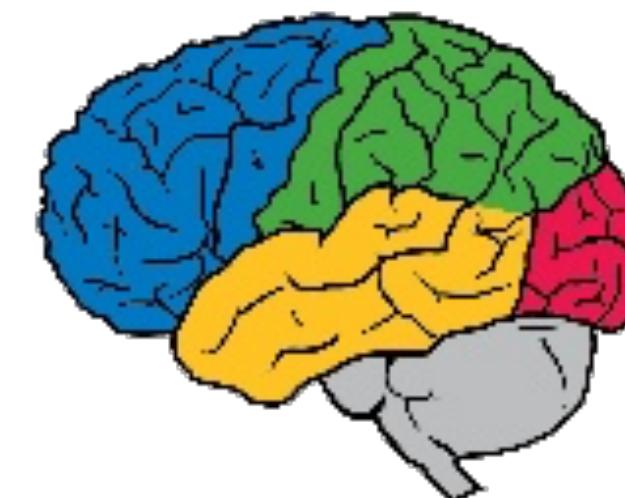


Meta Reinforcement Learning

Chelsea Finn



Why are humans so good at RL?



People have prior experience.

People have an existing representation of the world.

Can we learn a representation under which RL is fast?

Key idea: Explicitly optimize for such a representation

“Learn how to reinforcement learn”

Outline

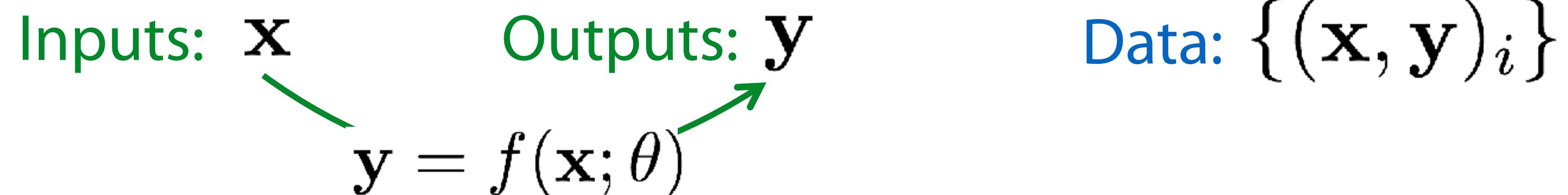
Meta-RL Problem Formulation & Examples

Method Classes: Recurrent Models, Gradient-Based Models

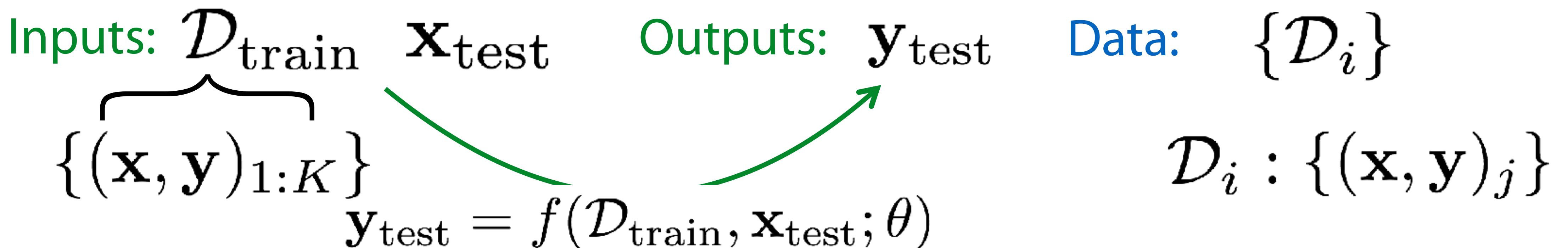
Challenges & Latest Developments

The Meta-Learning Problem

Supervised Learning:



Meta Supervised Learning:

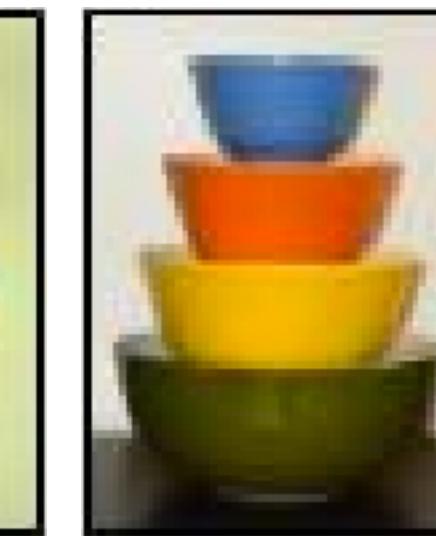


Why is this view useful?

Reduces the problem to the design & optimization of f .

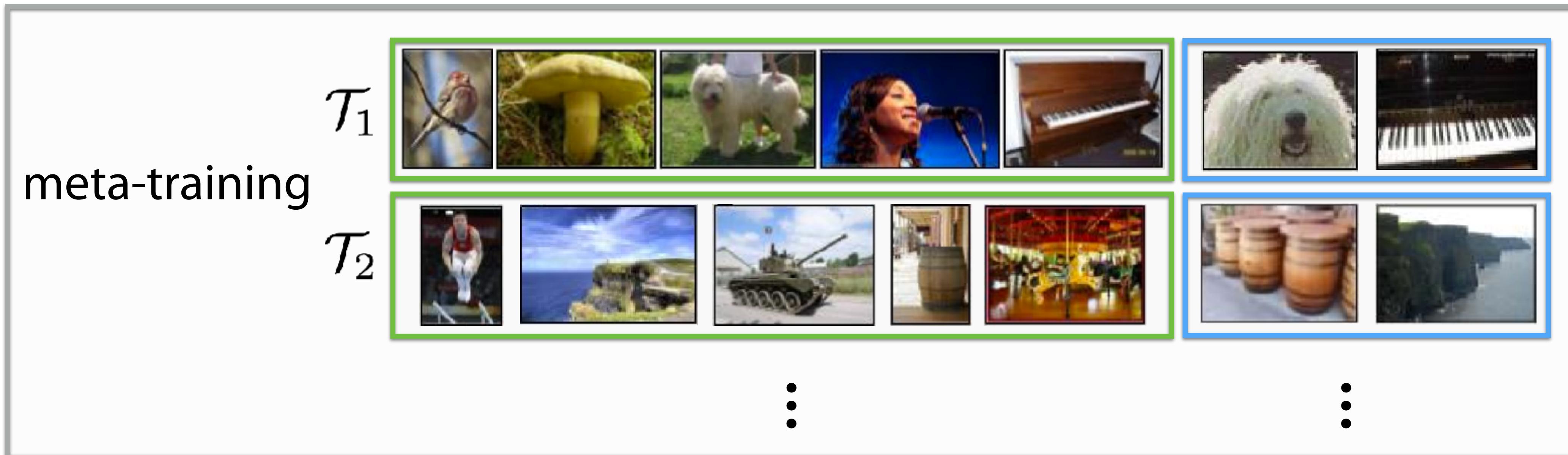
Example: Few-Shot Classification

Given 1 example of 5 classes:



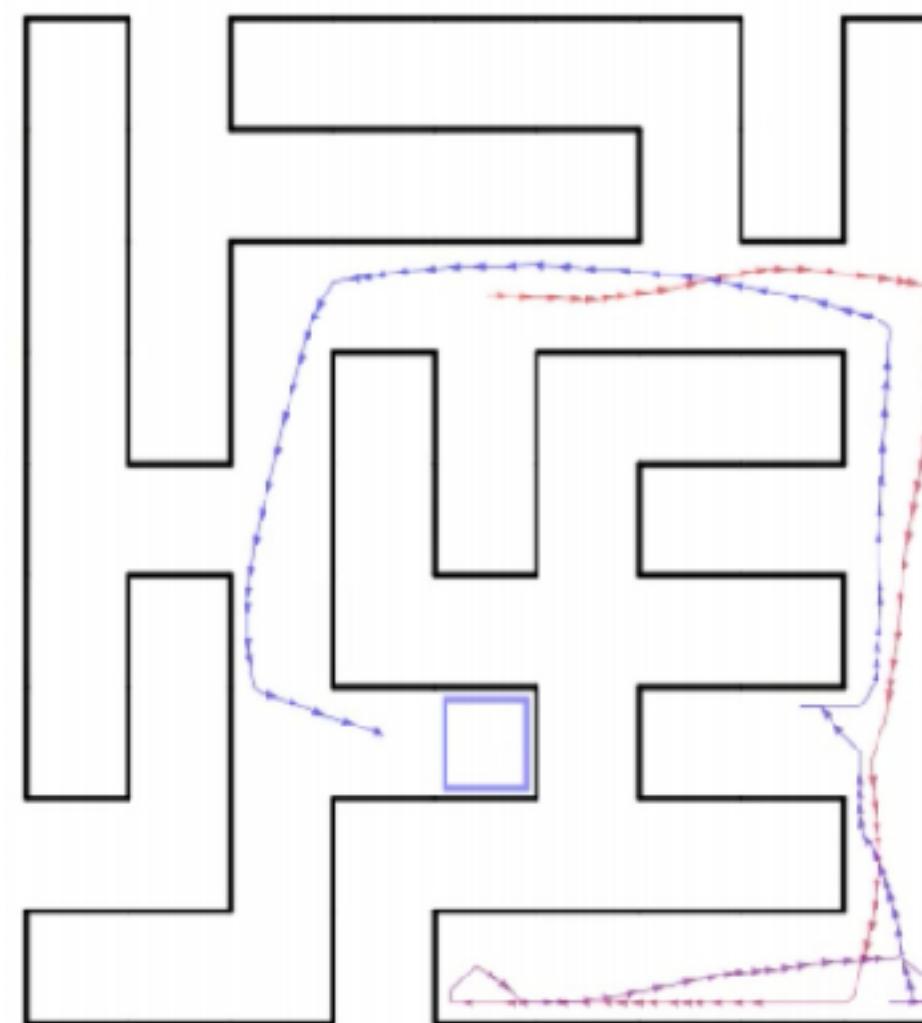
training data $\mathcal{D}_{\text{train}}$

Classify new examples
test set \mathbf{x}_{test}

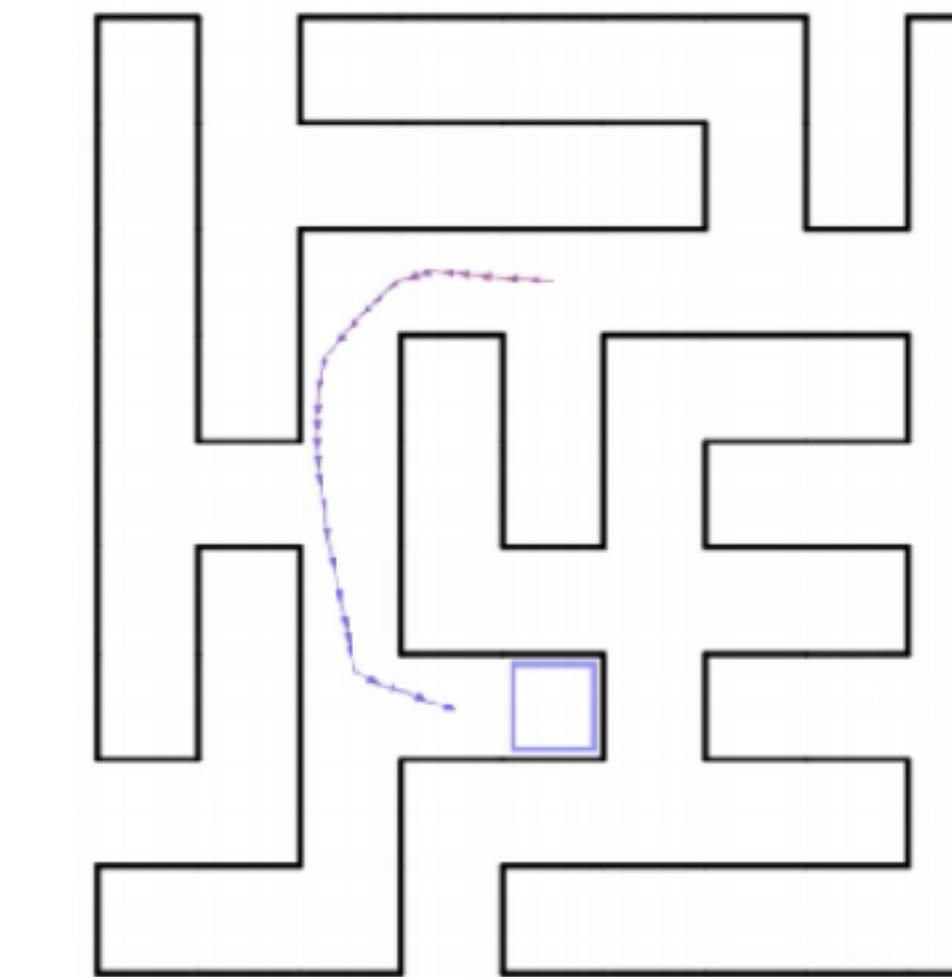


Meta-RL Example: Maze Navigation

Given a small amount of experience



Learn to solve the task



By learning how to learn many other tasks:

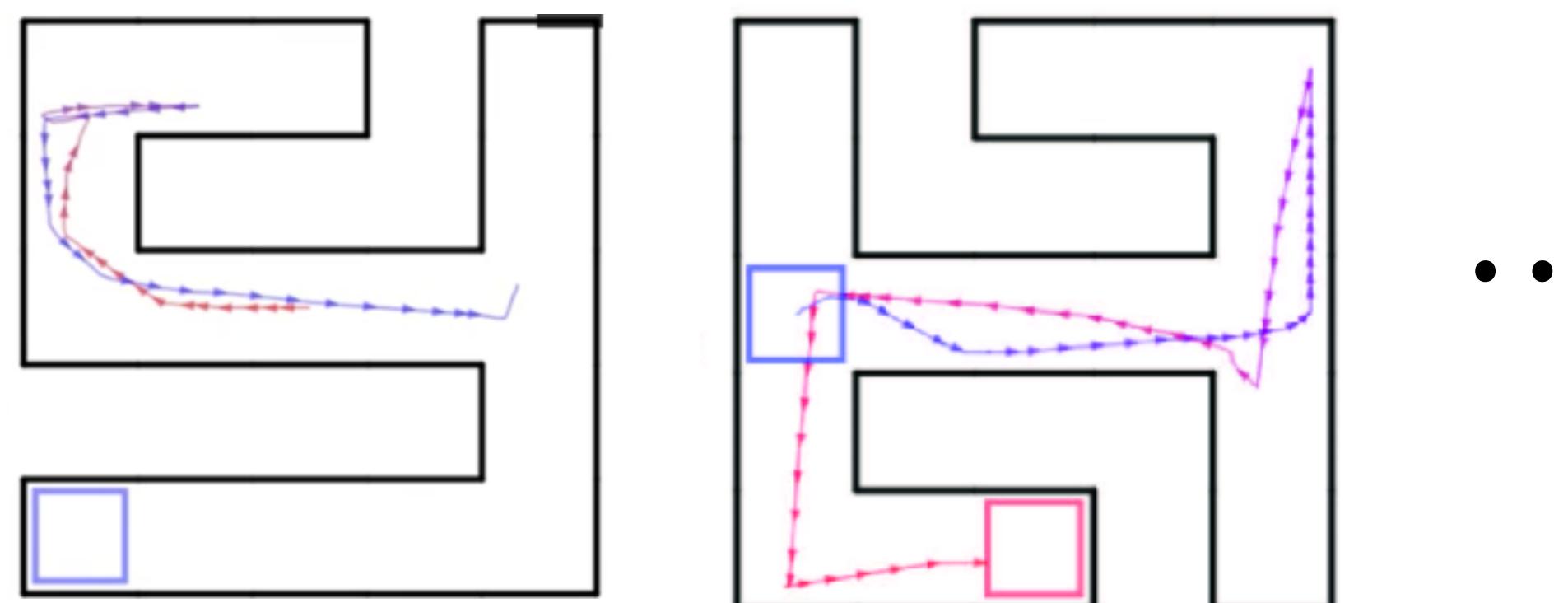
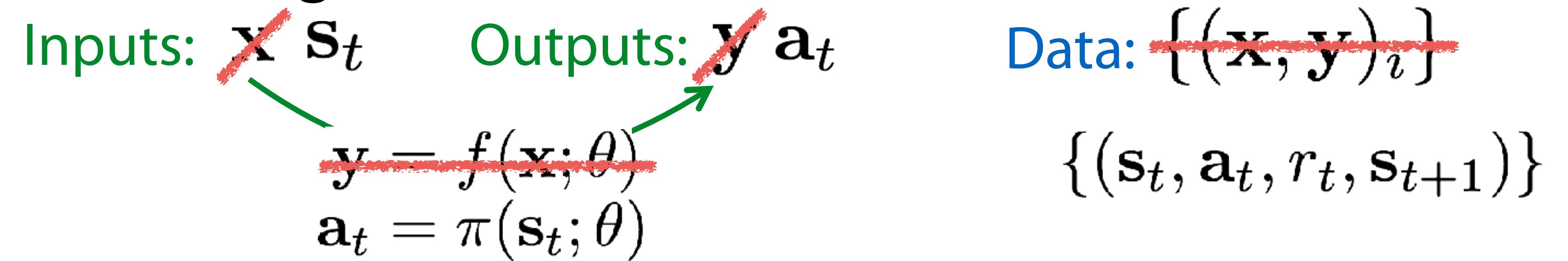


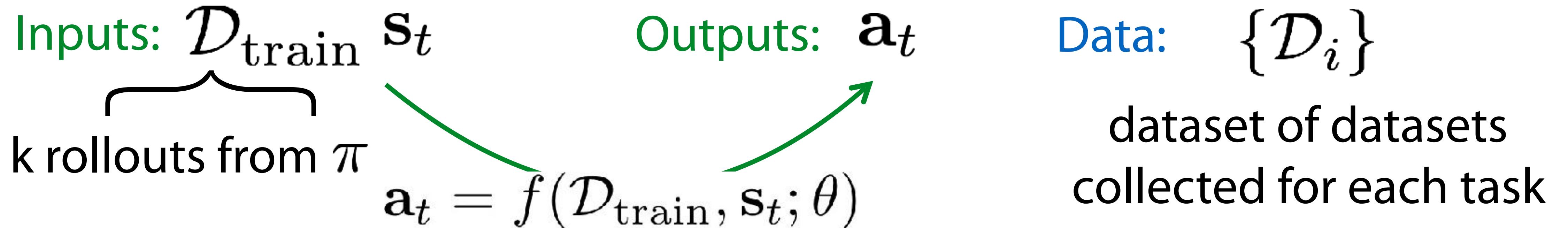
diagram adapted from Duan et al. '17

The Meta Reinforcement Learning Problem

Reinforcement Learning:



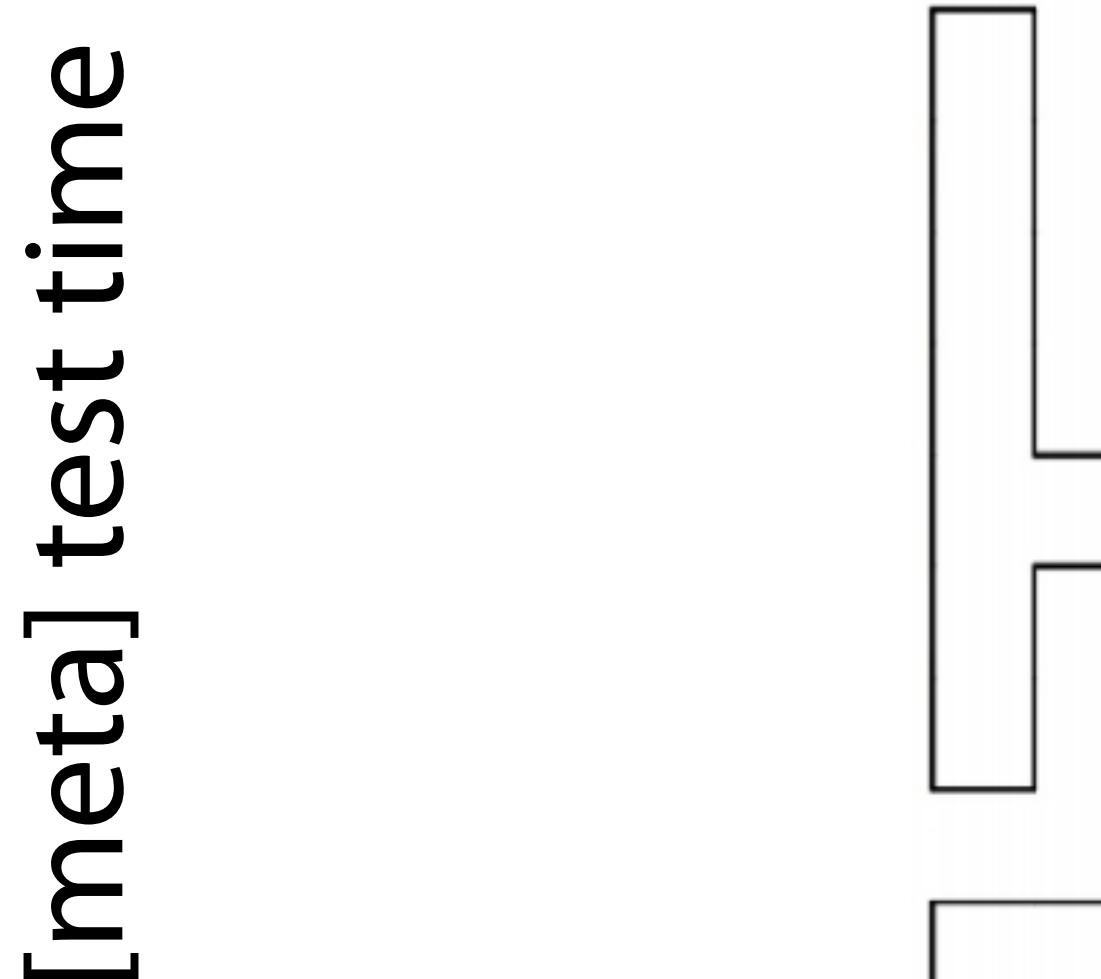
Meta Reinforcement Learning:



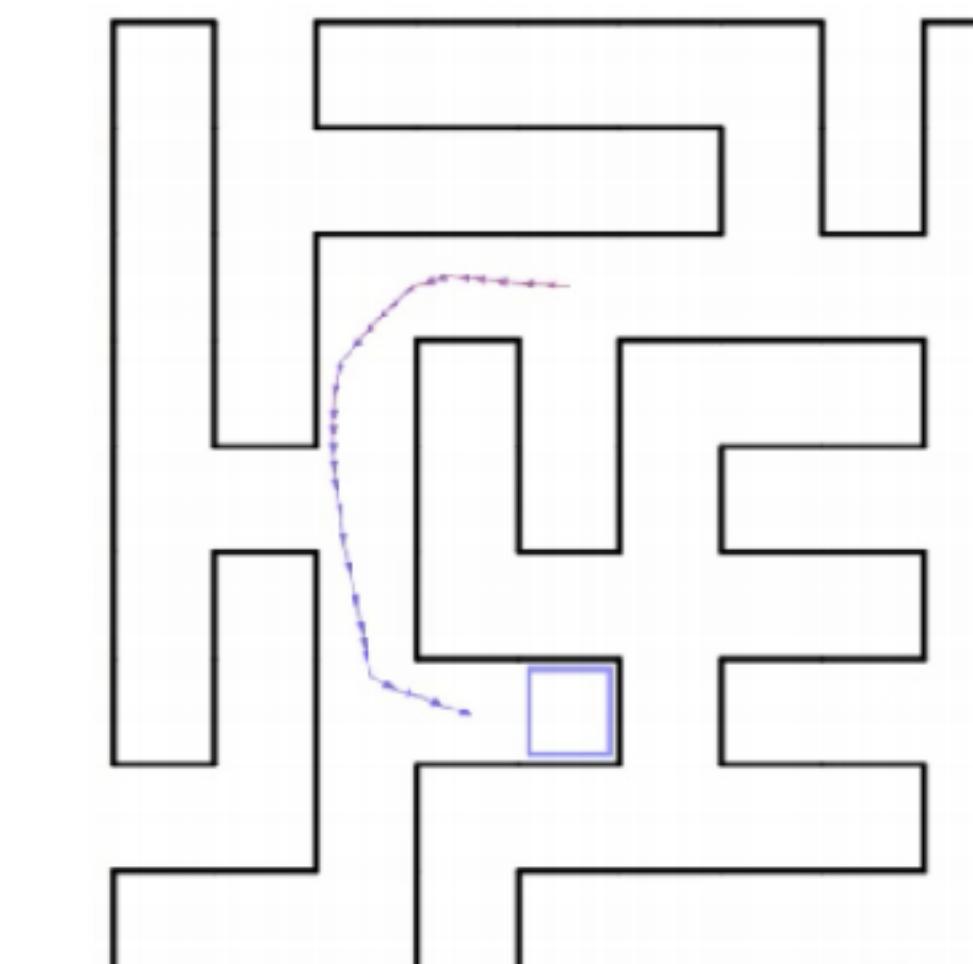
Design & optimization of f *and* collecting appropriate data
(learning to explore)

Meta-RL Example: Maze Navigation

Given a small amount of experience



Learn to solve the task



[meta] train time

By learning how to learn many other tasks:

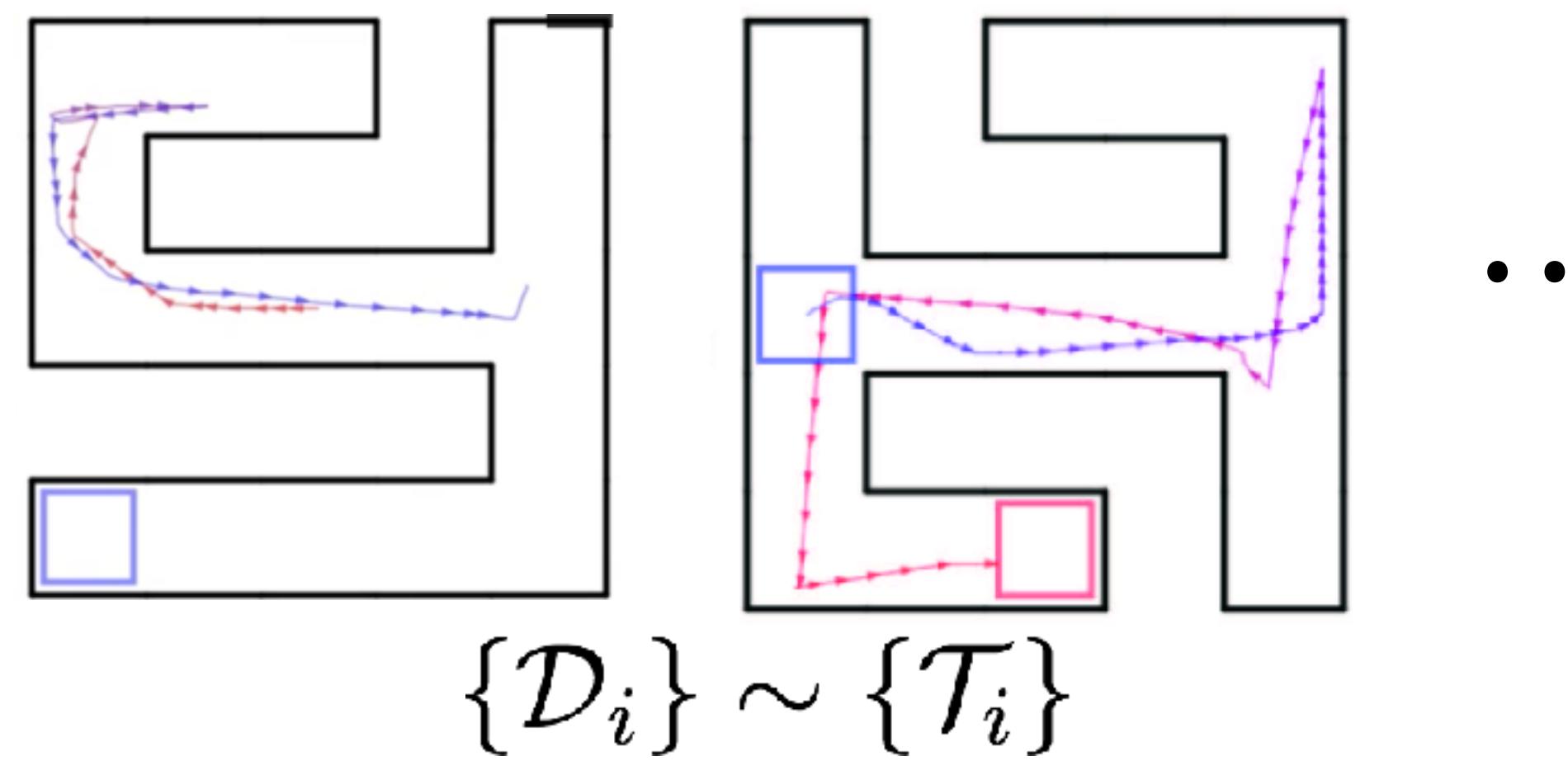
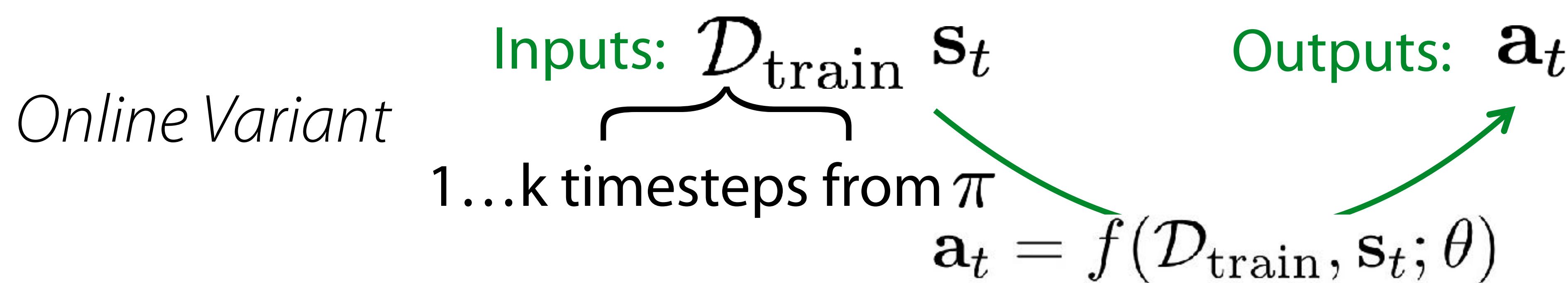
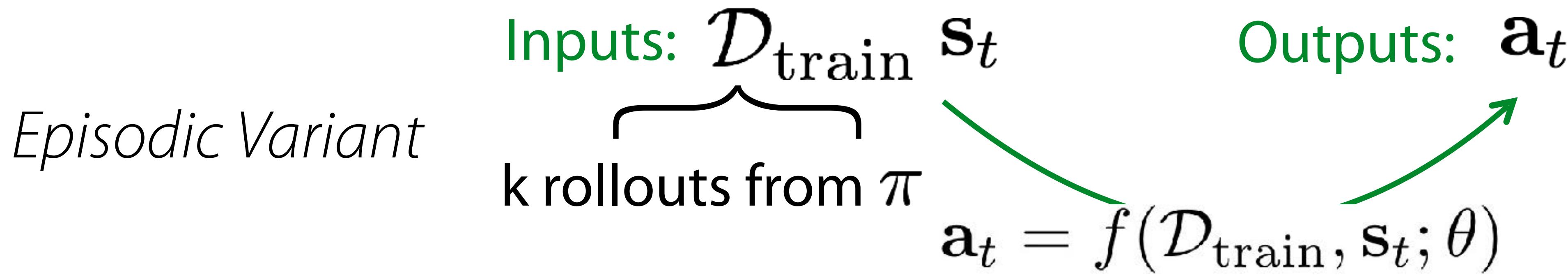


diagram adapted from Duan et al. '17

The Meta Reinforcement Learning Problem

Meta Reinforcement Learning:



Outline

Meta-RL Problem Formulation & Examples

Method Classes: Recurrent Models, Gradient-Based Models

Challenges & Latest Developments

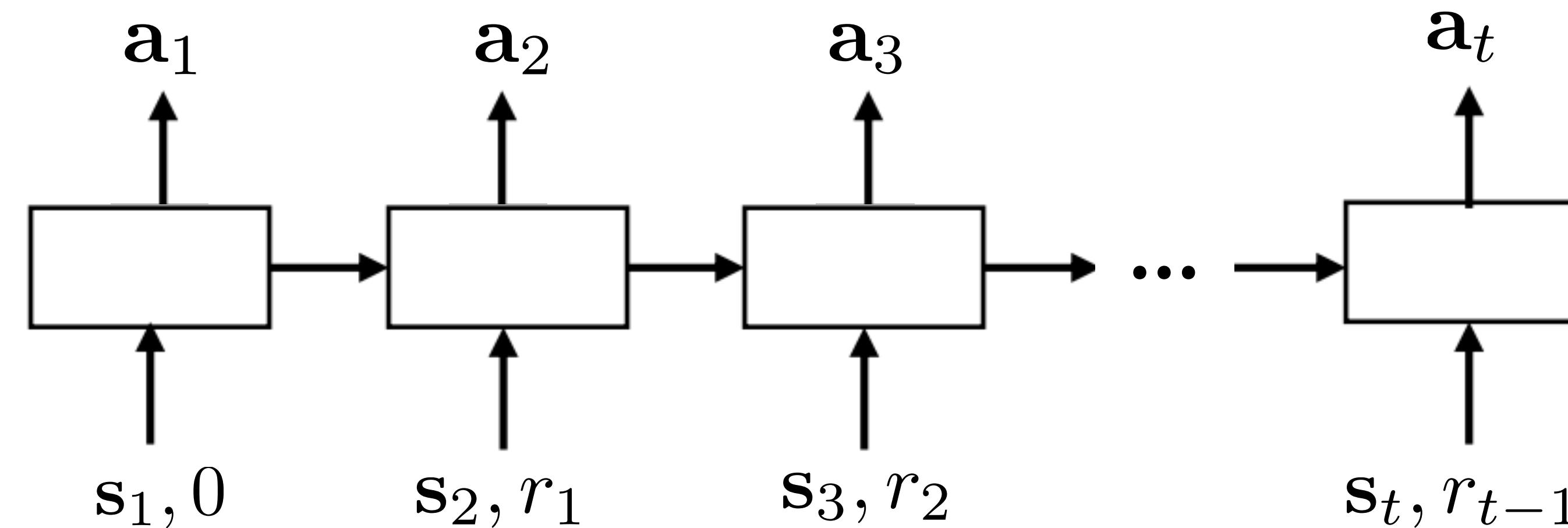
Design of f ?

$$\mathcal{D}_{\text{train}} \ s_t \longrightarrow a_t$$

Recurrent network
(LSTM, NTM, Conv)

$$a_t = f(\mathcal{D}_{\text{train}}, s_t; \theta)$$

Santoro et al. '16, Duan et al. '17, Wang et al. '17,
Mishra et al. '17, ...



Difference compared to using a recurrent policy?

Hidden state maintained **across episodes** within a task!

Case Study: Simple Neural Attentive Meta-Learner

Mishra et al, ICLR'18

A SIMPLE NEURAL ATTENTIVE META-LEARNER

Nikhil Mishra *†

UC Berkeley, Department of Electrical Engineering and Computer Science
Embodied Intelligence

{nmishra, rohaninejadm, c.xi, pabbeel}@berkeley.edu

Mostafa Rohaninejad*

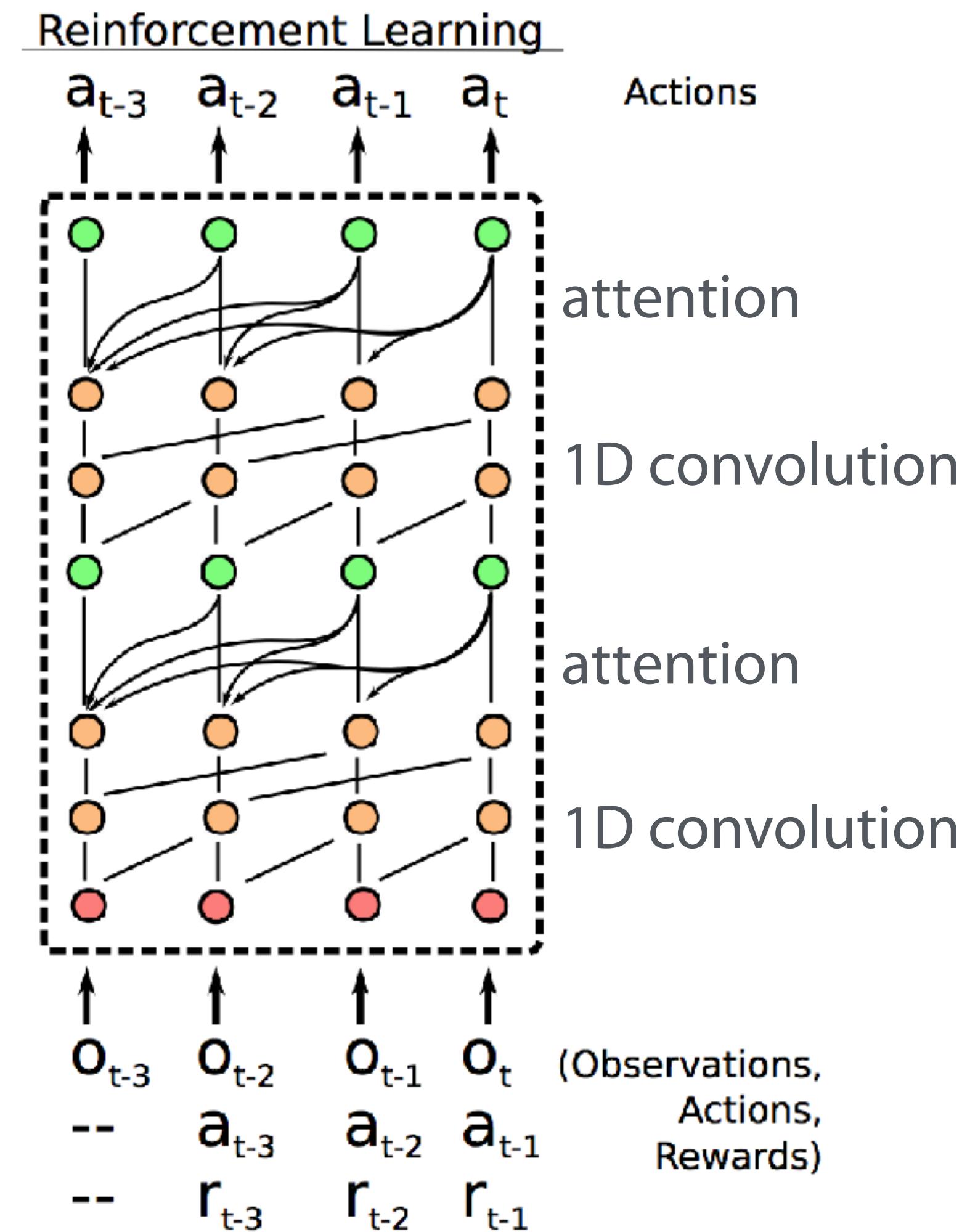
Xi Chen†

Pieter Abbeel†

Simple Neural Attentive Meta-Learner

Mishra et al, ICLR'18

interleave 1D convolutions and attention

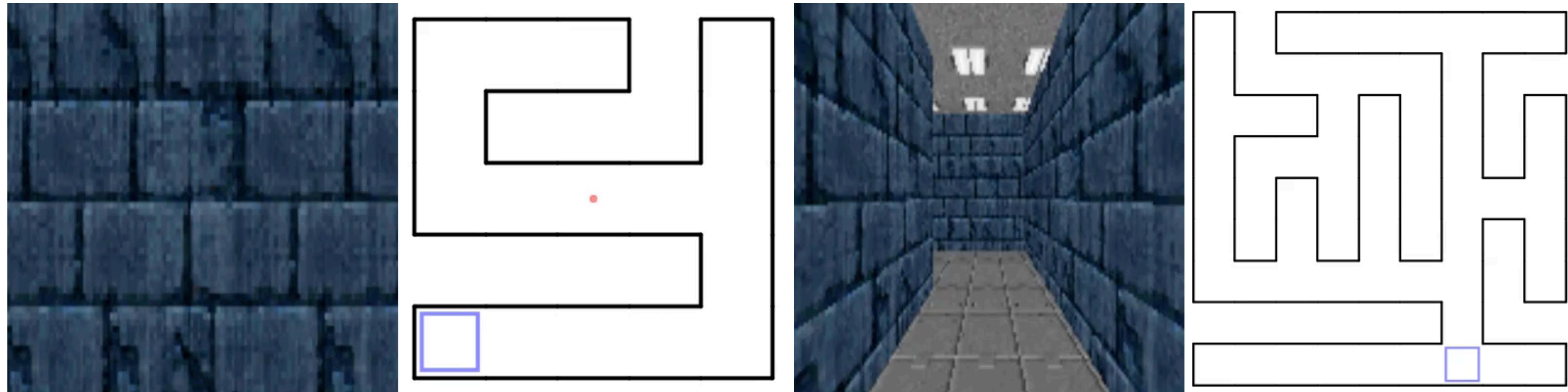


Simple Neural Attentive Meta-Learner

Mishra et al, ICLR'18

Experiment: Learning to visually navigate a maze

- train on 1000 small mazes
- test on held-out small mazes and large mazes



Simple Neural Attentive Meta-Learner

Mishra et al, ICLR'18

Experiment: Learning to visually navigate a maze

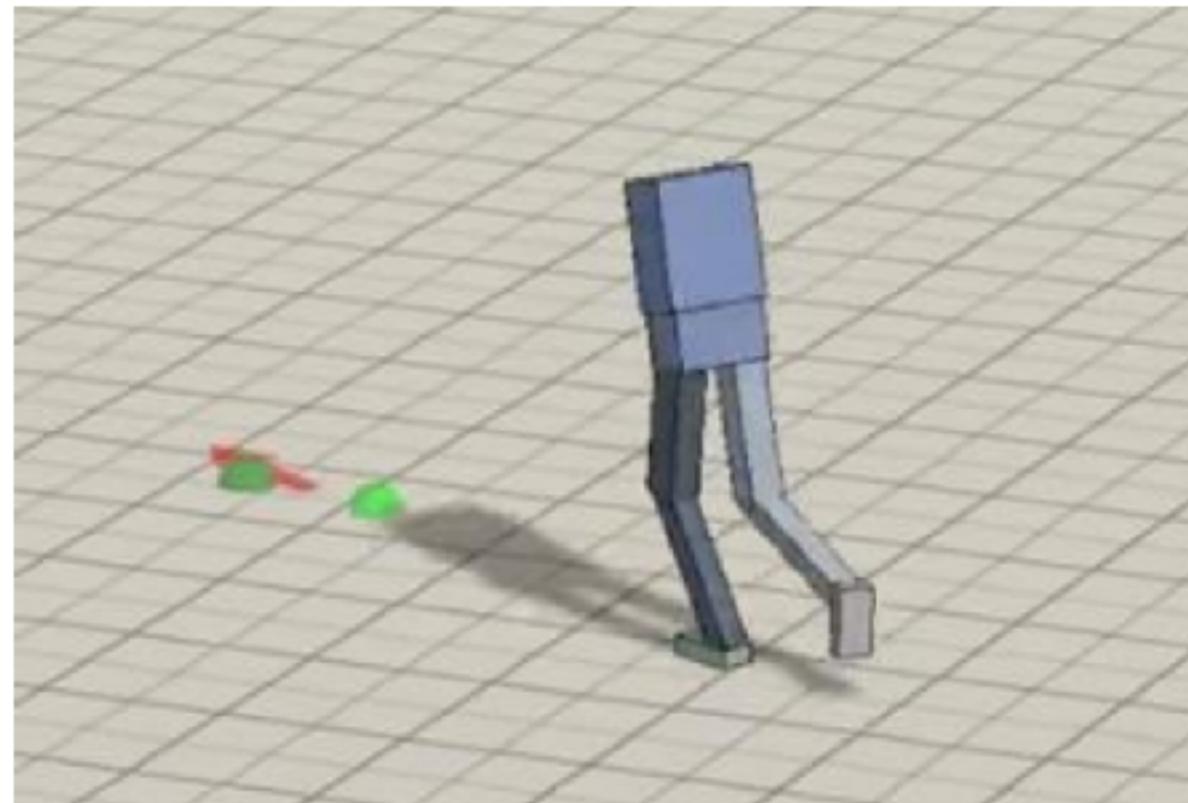
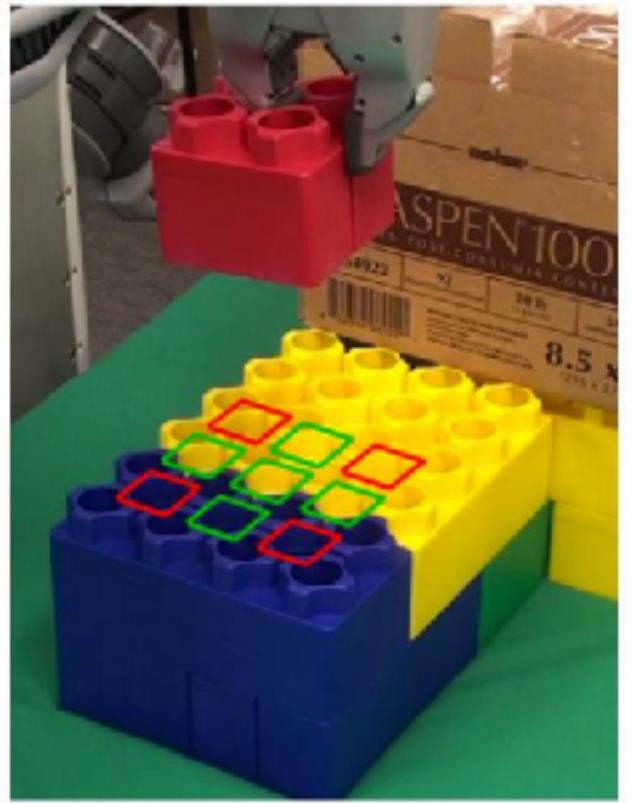
- train on 1000 small mazes
- test on held-out small mazes and large mazes

Method	Small Maze		Large Maze	
	Episode 1	Episode 2	Episode 1	Episode 2
Random	188.6 ± 3.5	187.7 ± 3.5	420.2 ± 1.2	420.8 ± 1.2
LSTM	52.4 ± 1.3	39.1 ± 0.9	180.1 ± 6.0	150.6 ± 5.9
SNAIL (ours)	50.3 ± 0.3	34.8 ± 0.2	140.5 ± 4.2	105.9 ± 2.4

Table 5: Average time to find the goal on each episode

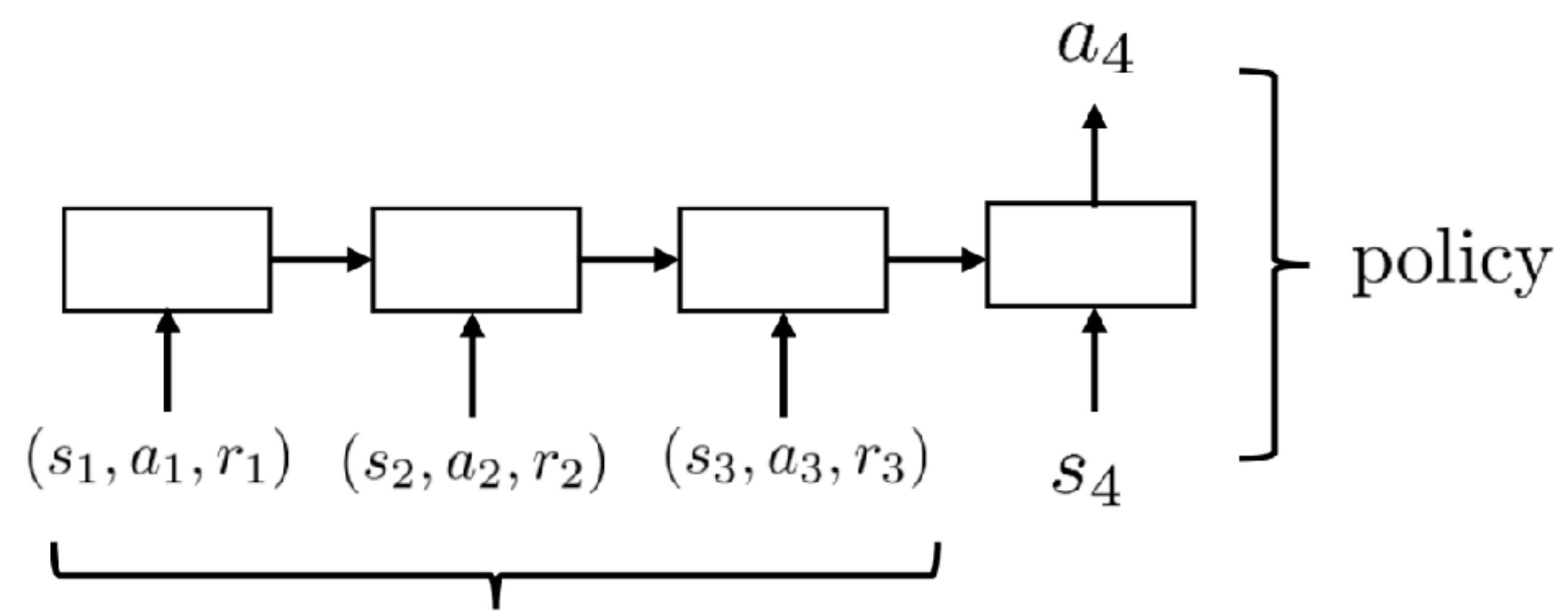
Digression: Connection to Contextual Policies

contextual policy: $\pi_\theta(\mathbf{a}|\mathbf{s}, \omega)$



ω : stack location

ω : walking direction



Contextual policy with experience as context.

What about goal-conditioned policies / value functions?

- rewards are a strict generalization of goals
- meta-RL objective is to *adapt* new tasks vs. *generalize* to new goals
(k-shot vs. 0-shot)

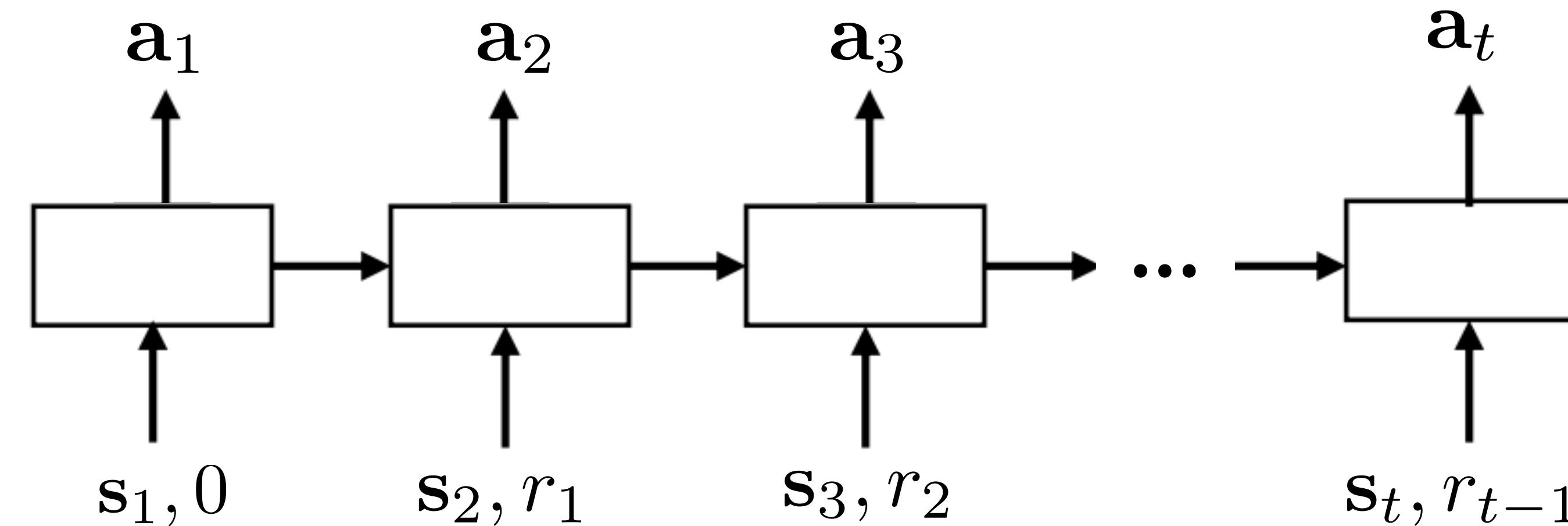
Design of f ?

$$\mathcal{D}_{\text{train}} \ x_{\text{test}} \xrightarrow{\text{green arrow}} y_{\text{test}}$$

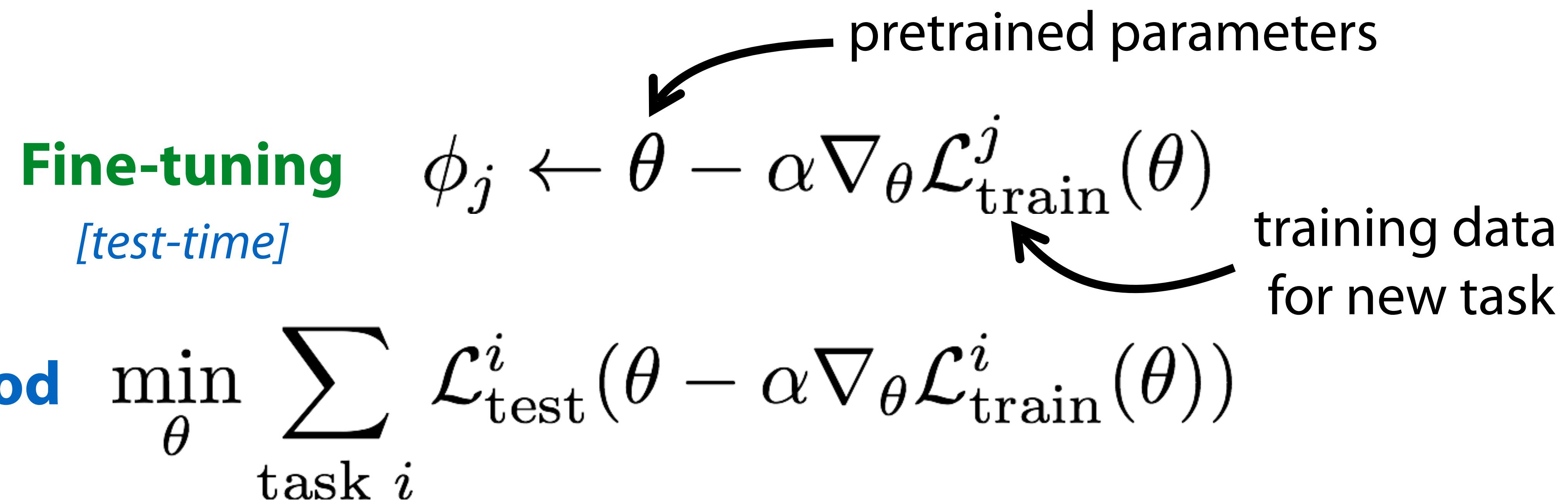
Recurrent network
(LSTM, NTM, Conv)

$$y_{\text{test}} = f(\mathcal{D}_{\text{train}}, x_{\text{test}}; \theta)$$

Santoro et al. '16, Duan et al. '17, Wang et al. '17,
Munkhdalai & Yu '17, Mishra et al. '17, ...



- + general & expressive
- + a variety of design choices in architecture
- complex model for complex task of learning
- impractical data requirements



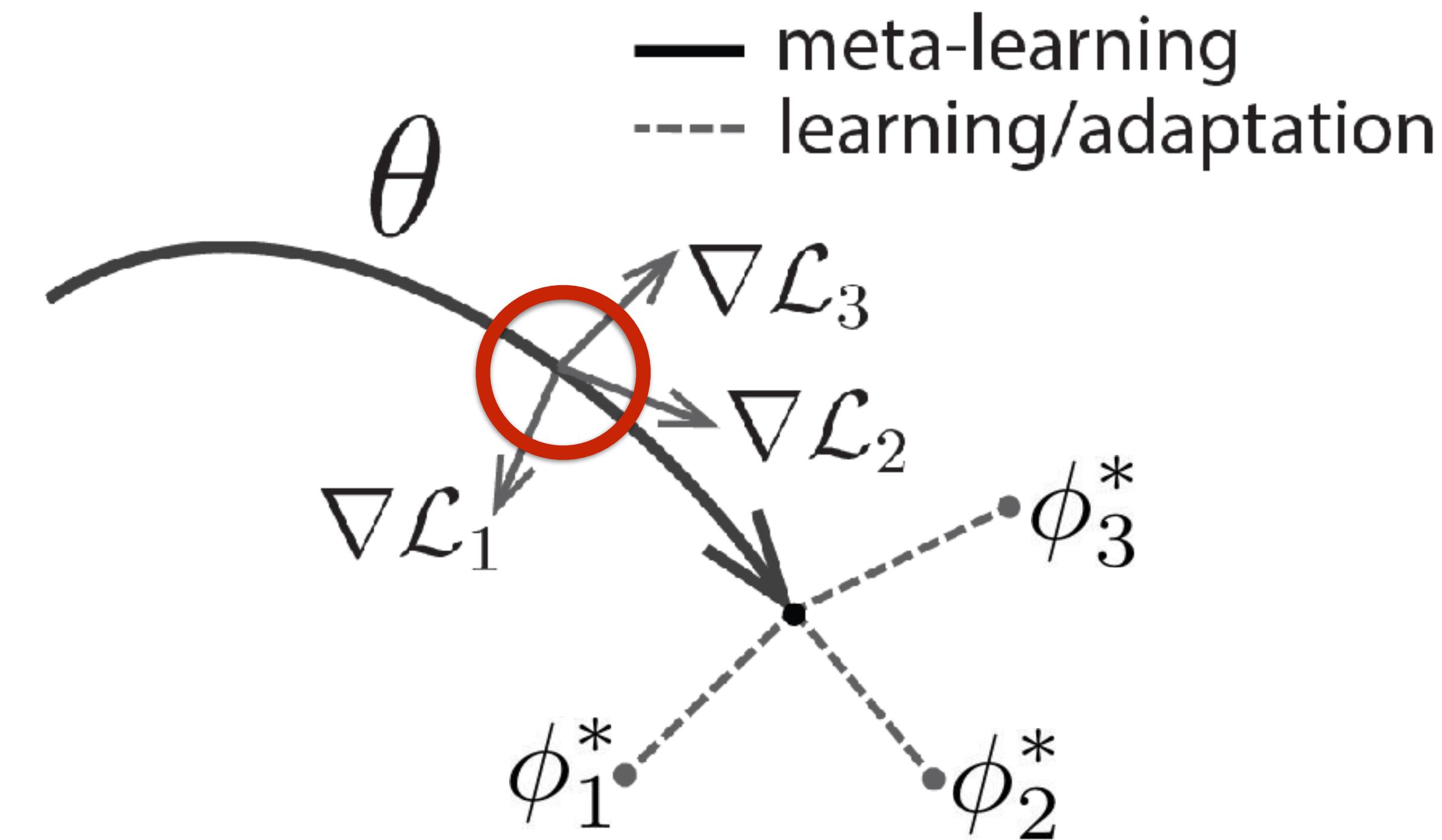
Key idea: Train over many tasks, to learn parameter vector θ that transfers

$$\min_{\theta} \sum_{\text{task } i} \mathcal{L}_{\text{test}}^i(\theta - \alpha \nabla_{\theta} \mathcal{L}_{\text{train}}^i(\theta))$$

θ

parameter vector
being meta-learned

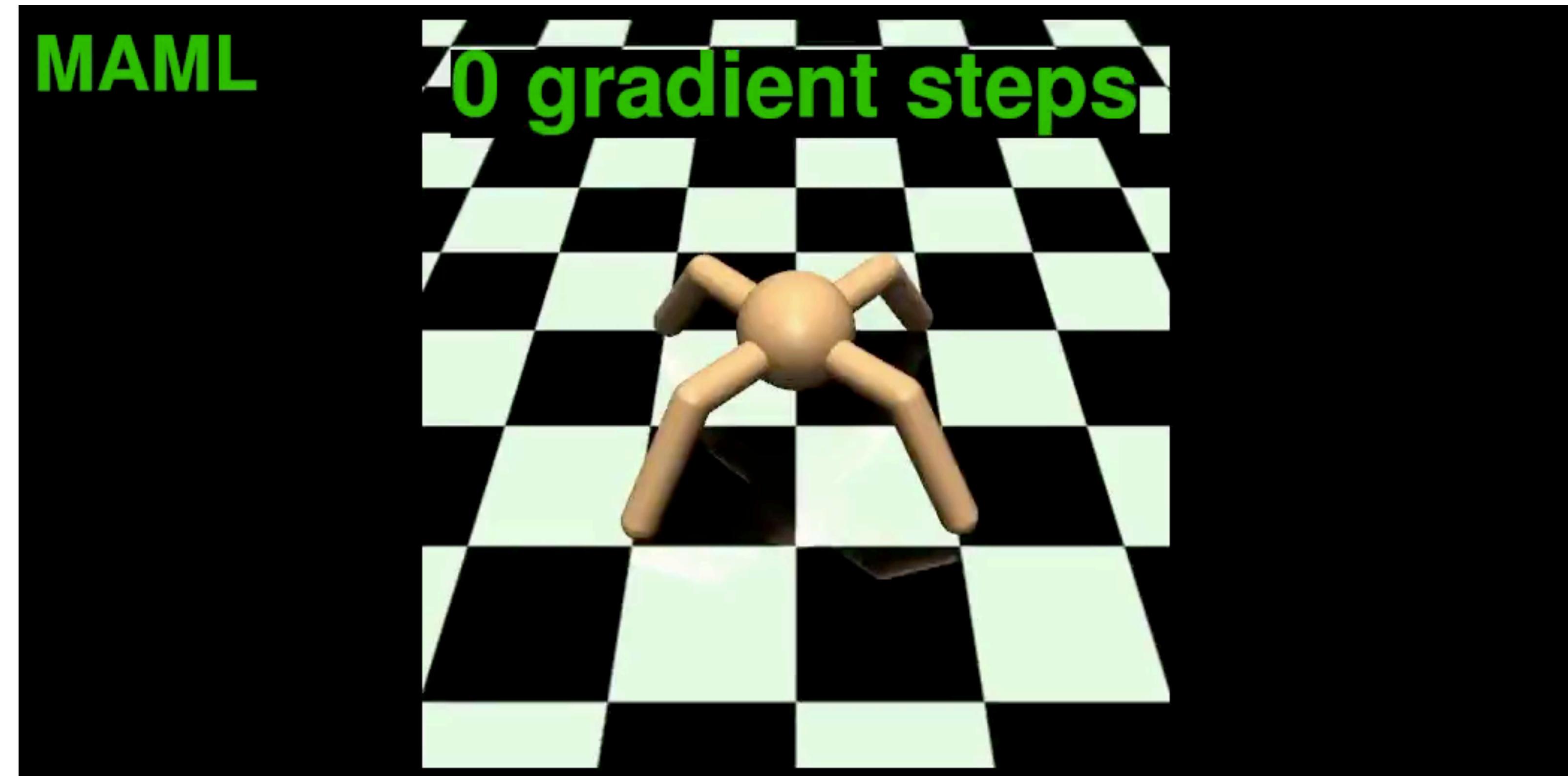
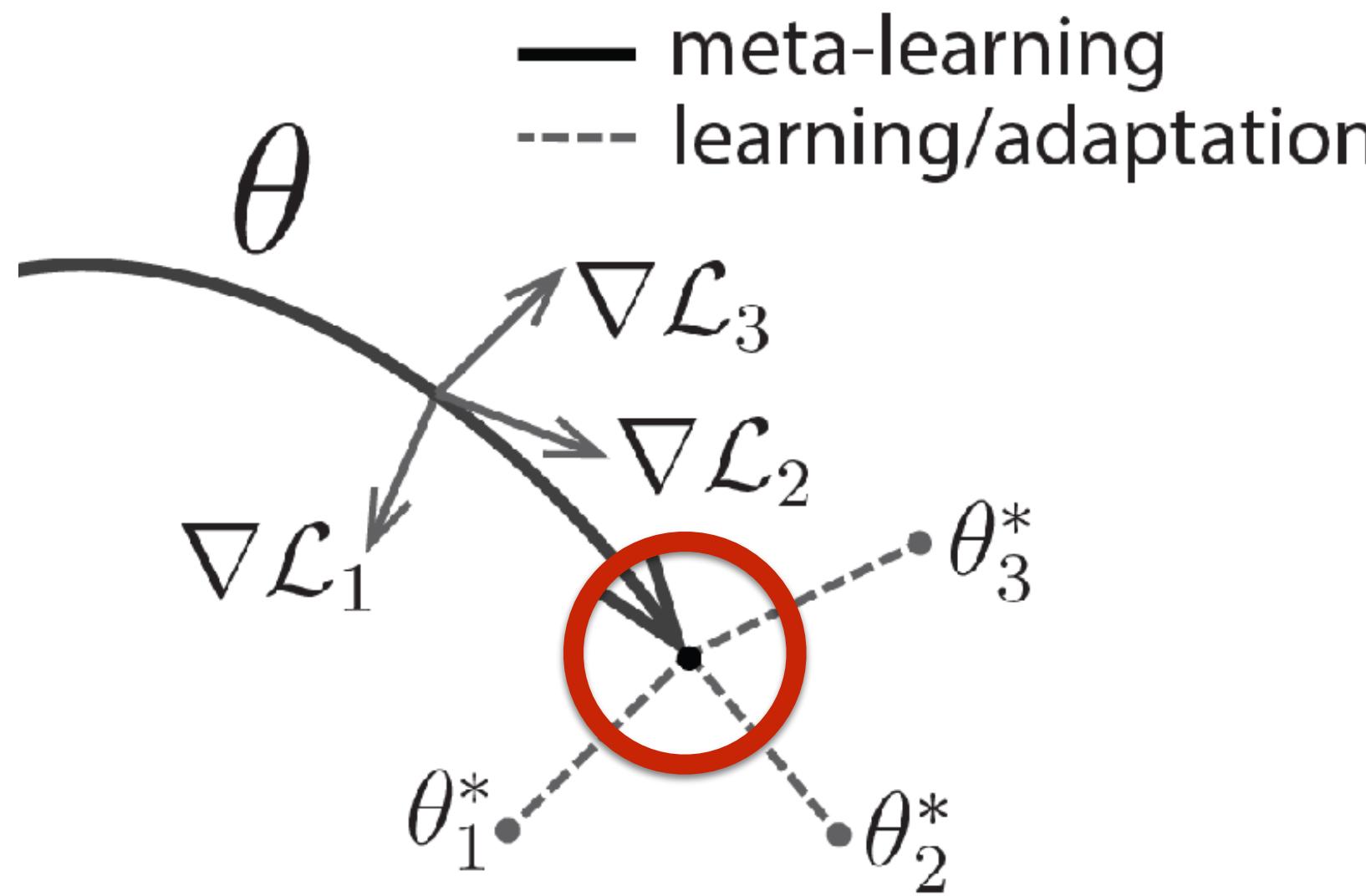
ϕ_i^* optimal parameter
vector for task i



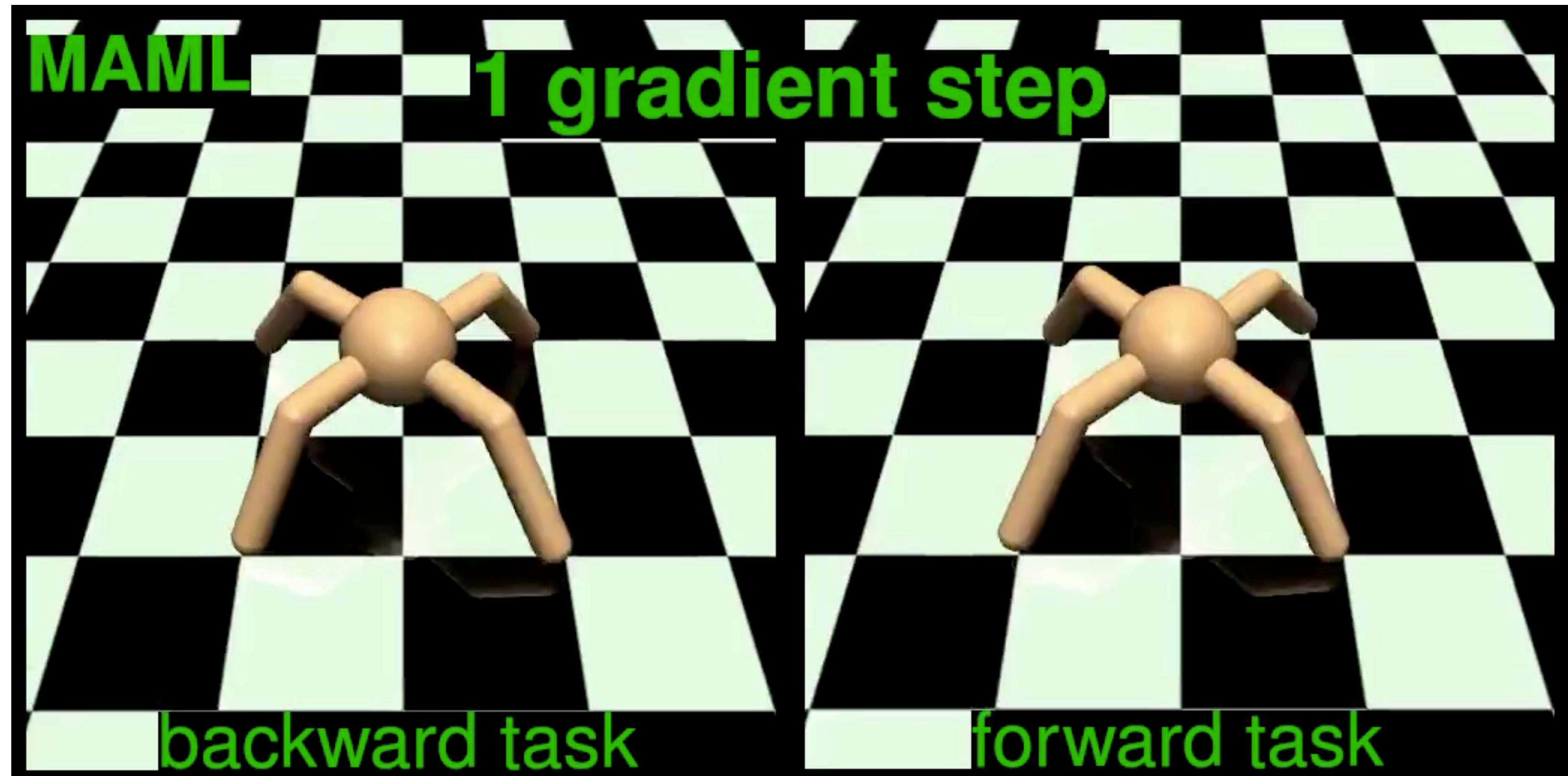
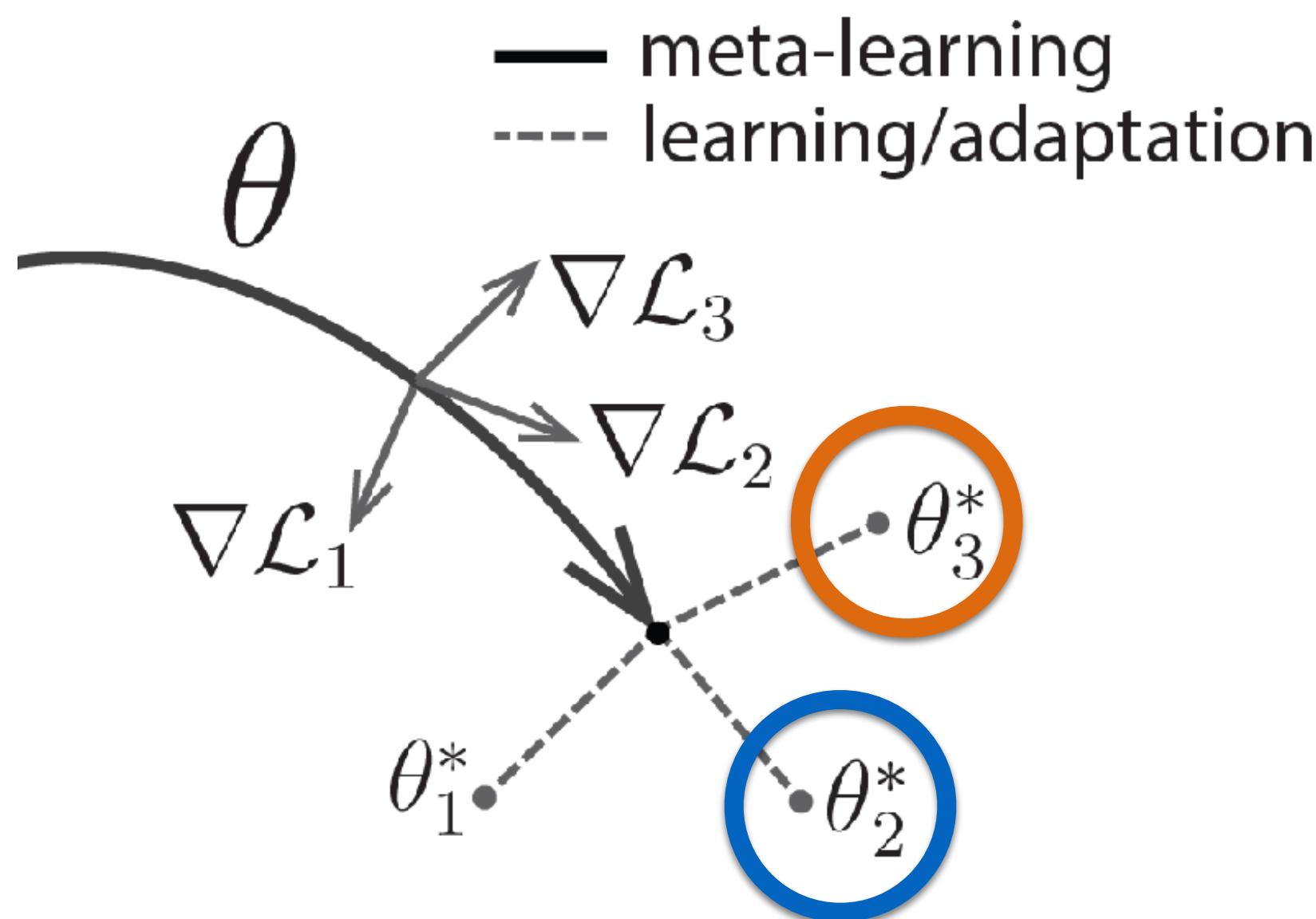
Model-Agnostic Meta-Learning

Can we learn a representation under which RL is fast?

Fast Adaptation in Reinforcement Learning



Fast Adaptation in Reinforcement Learning

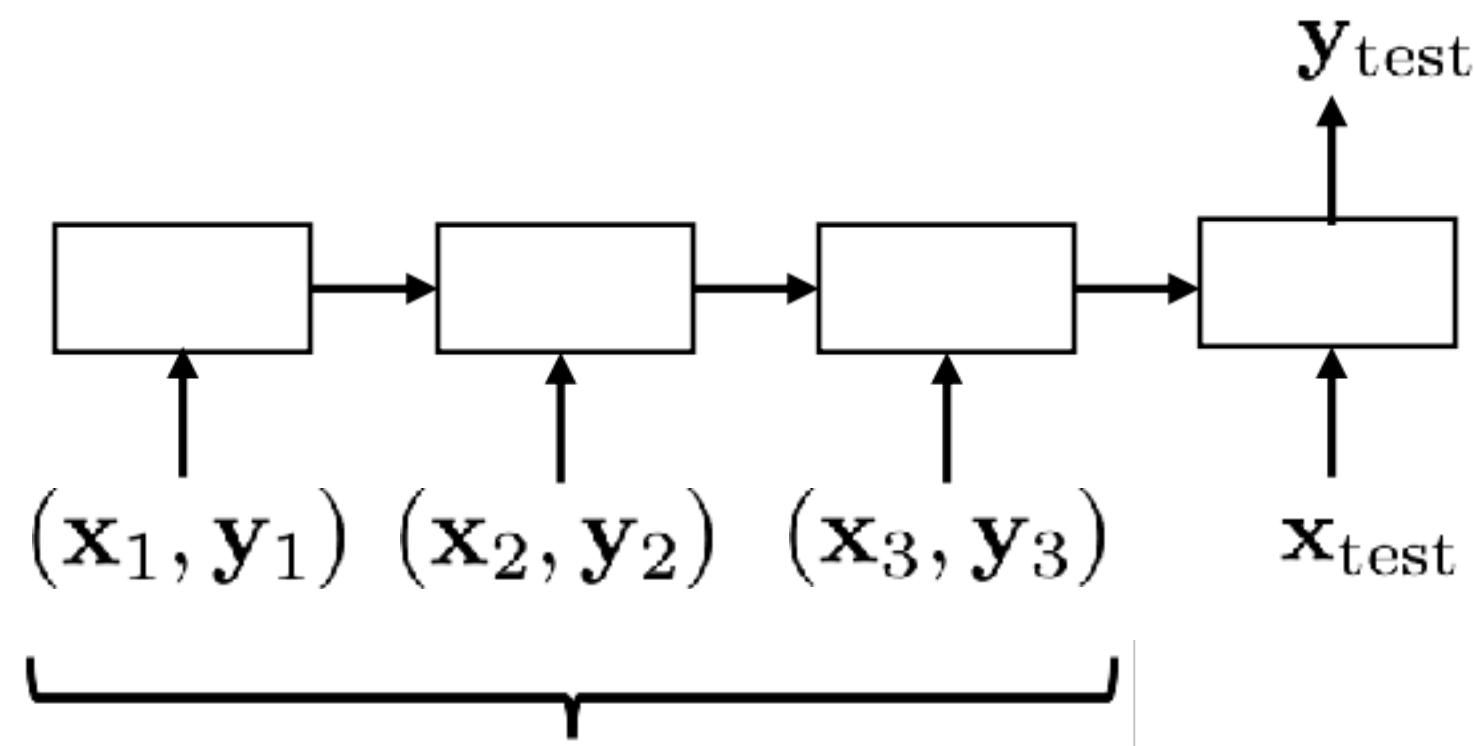


Design of f ?

$$\mathcal{D}_{\text{train}} \ x_{\text{test}} \xrightarrow{\quad} y_{\text{test}}$$

Recurrent network

$$y_{\text{test}} = f(\mathcal{D}_{\text{train}}, x_{\text{test}}; \theta)$$



network implements the
“learned learning procedure”

Does it converge?

- Sort of?

What does it converge to?

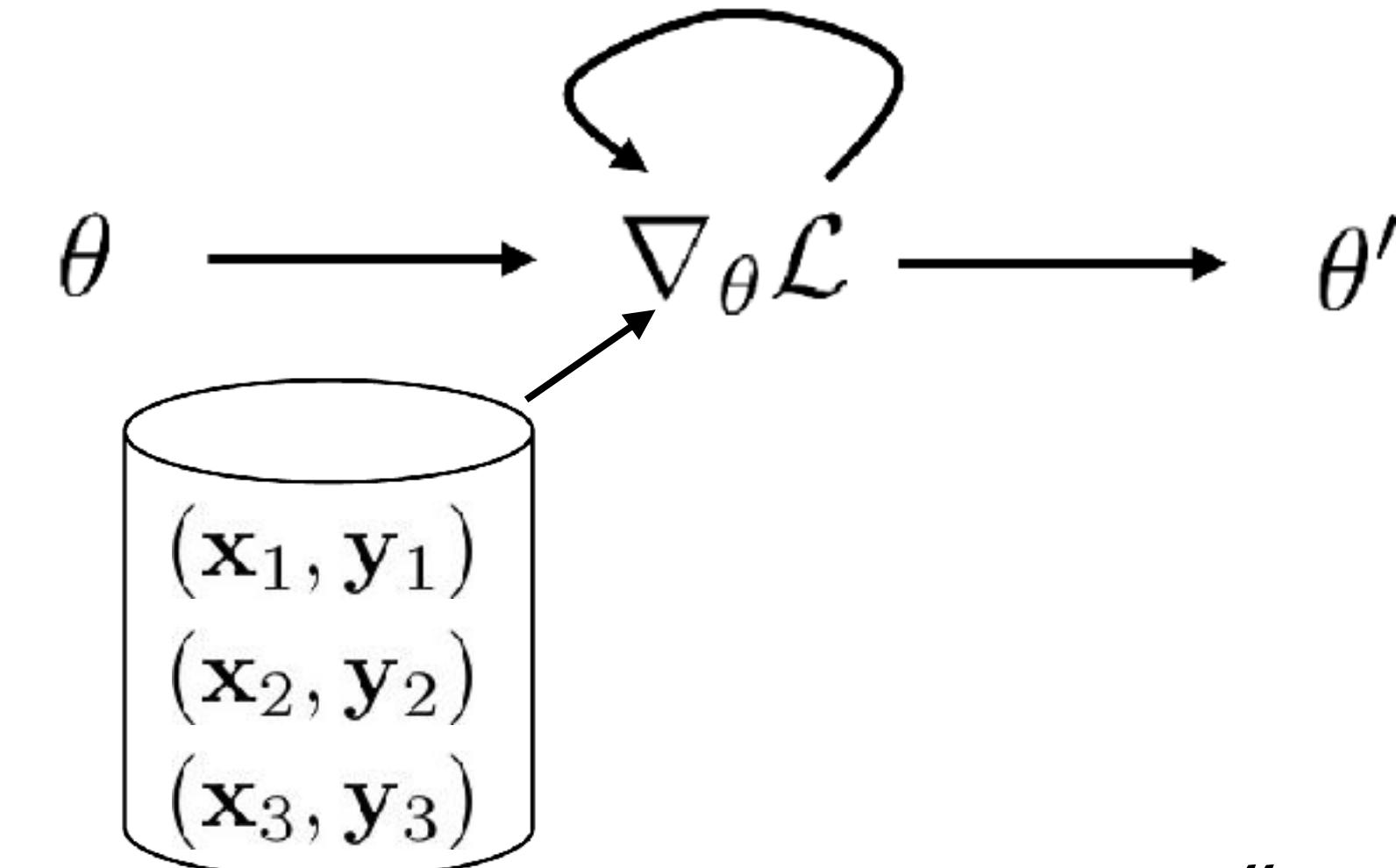
- Who knows...

What to do if not good enough?

- Nothing

MAML

$$y_{\text{test}} = f(x_{\text{test}}; \theta - \alpha \nabla_{\theta} \mathcal{L}(\mathcal{D}_{\text{train}}))$$



“consistency”

Does it converge?

- Yes (it's gradient descent...)

What does it converge to?

- A local optimum (it's gradient descent...)

What to do if not good enough?

- Keep taking gradient steps (it's gradient descent..)

Does this structure come at a cost?

Outline

Meta-RL Problem Formulation & Examples

Method Classes: Recurrent Models, Gradient-Based Models

Challenges & Latest Developments

Wishlist in a Meta-RL algorithm

	RL ² , SAIL	MAML-PG
Consistent	✗	✓
Expressive	✓	✗
Structured Exploration*	?	?
Efficient & off-policy	✗	✗

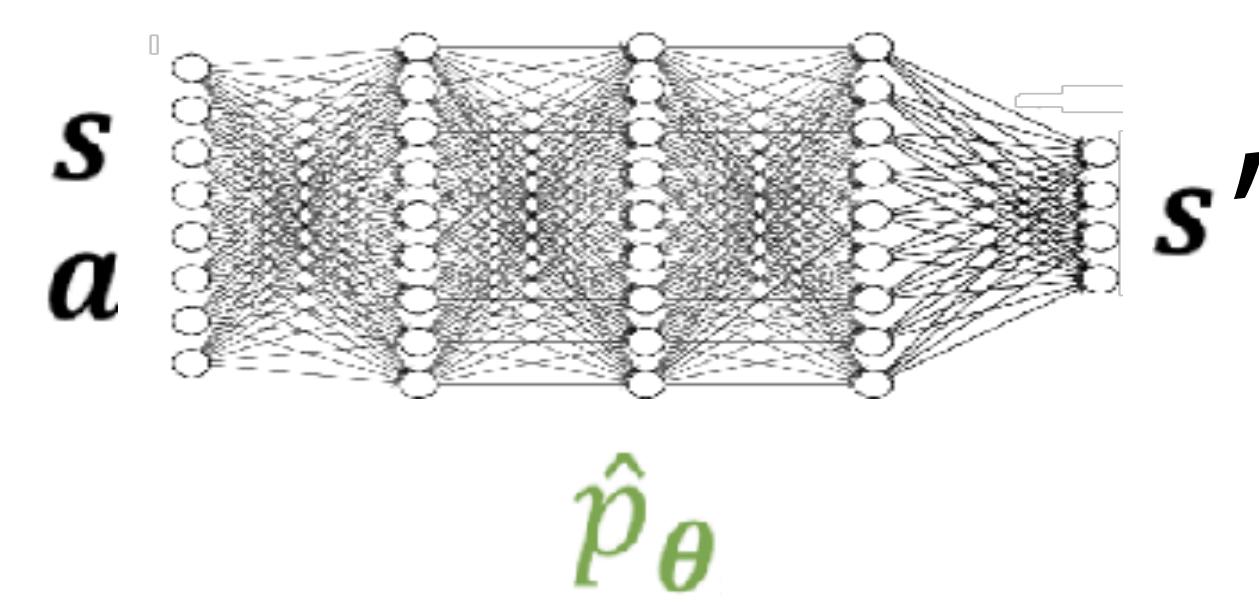
*Initial work in this direction: Gupta et al. NIPS '18, Stadie et al. NIPS'18

Model-based RL

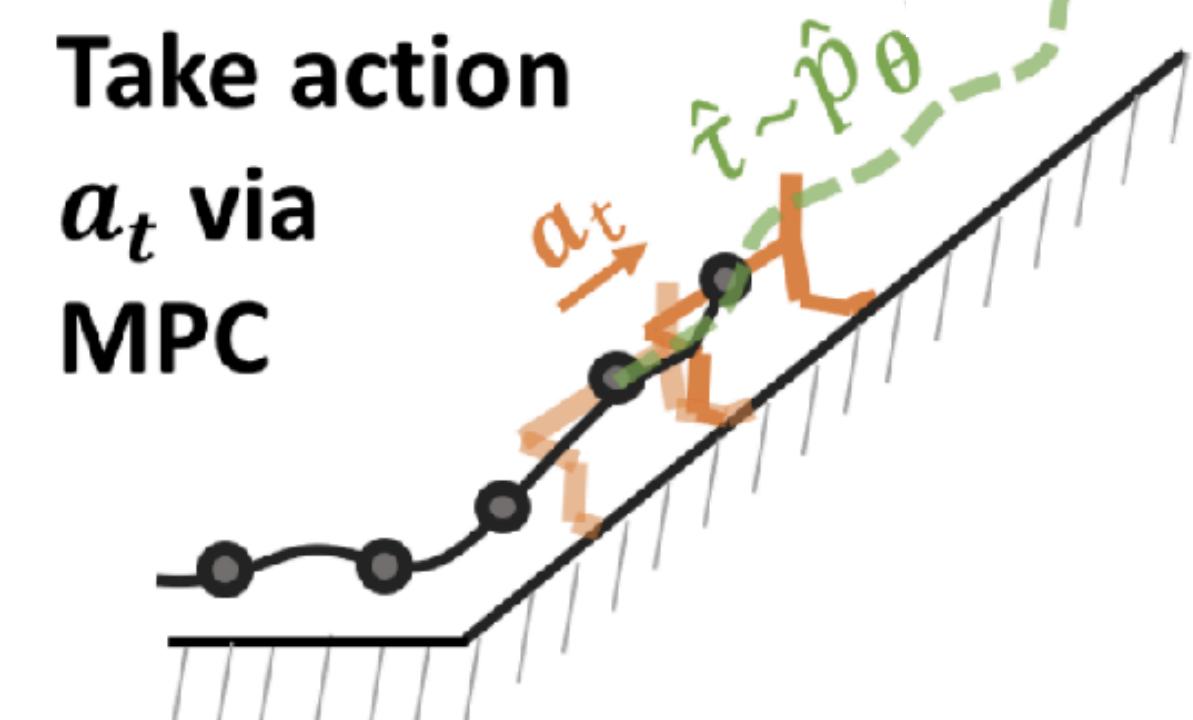
Collect data



Fit model



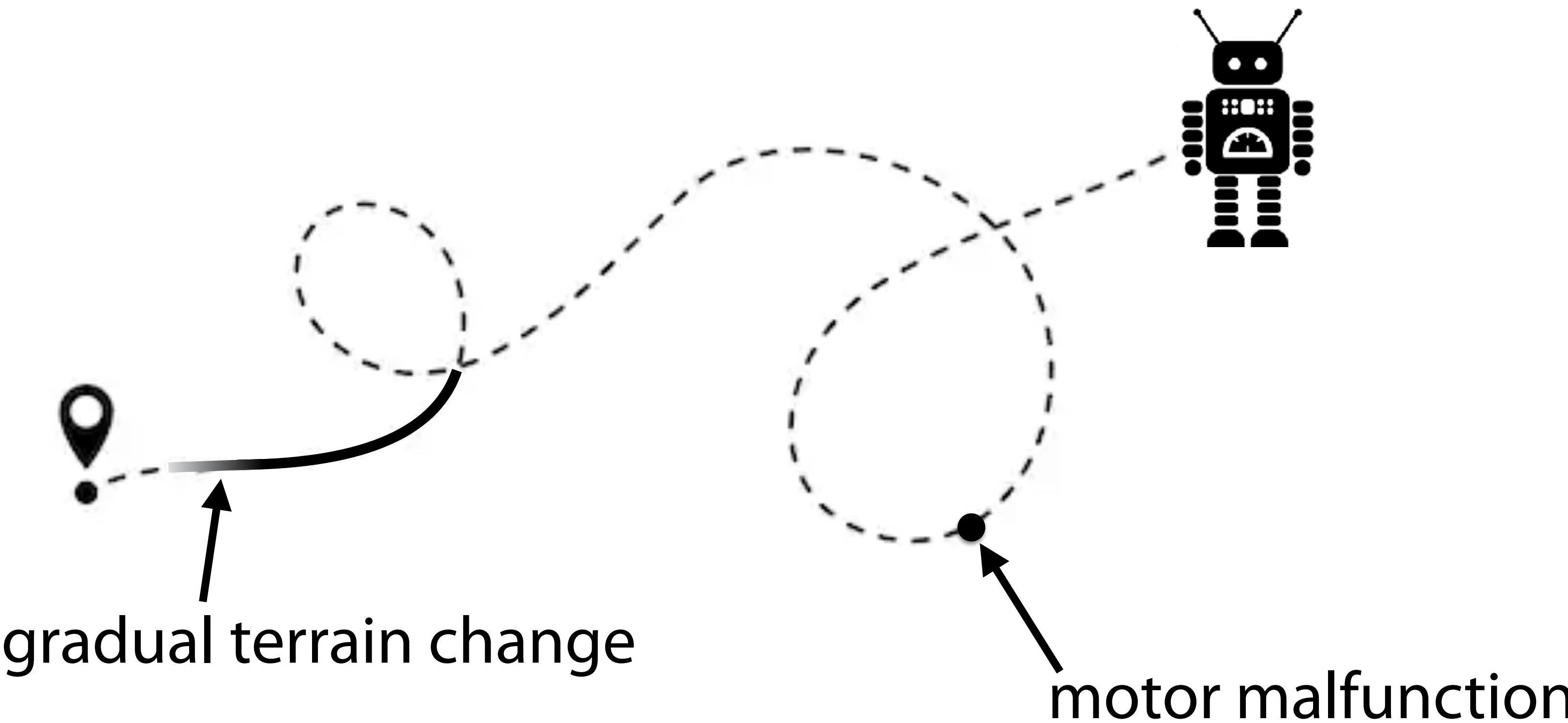
Plan using the model



Model-based RL already achieves zero-shot adaptation to new rewards.

Goal: learn to adapt model quickly to new environments
(online system ID)

Goal: learn to adapt model quickly to new environments



Anusha Nagabandi

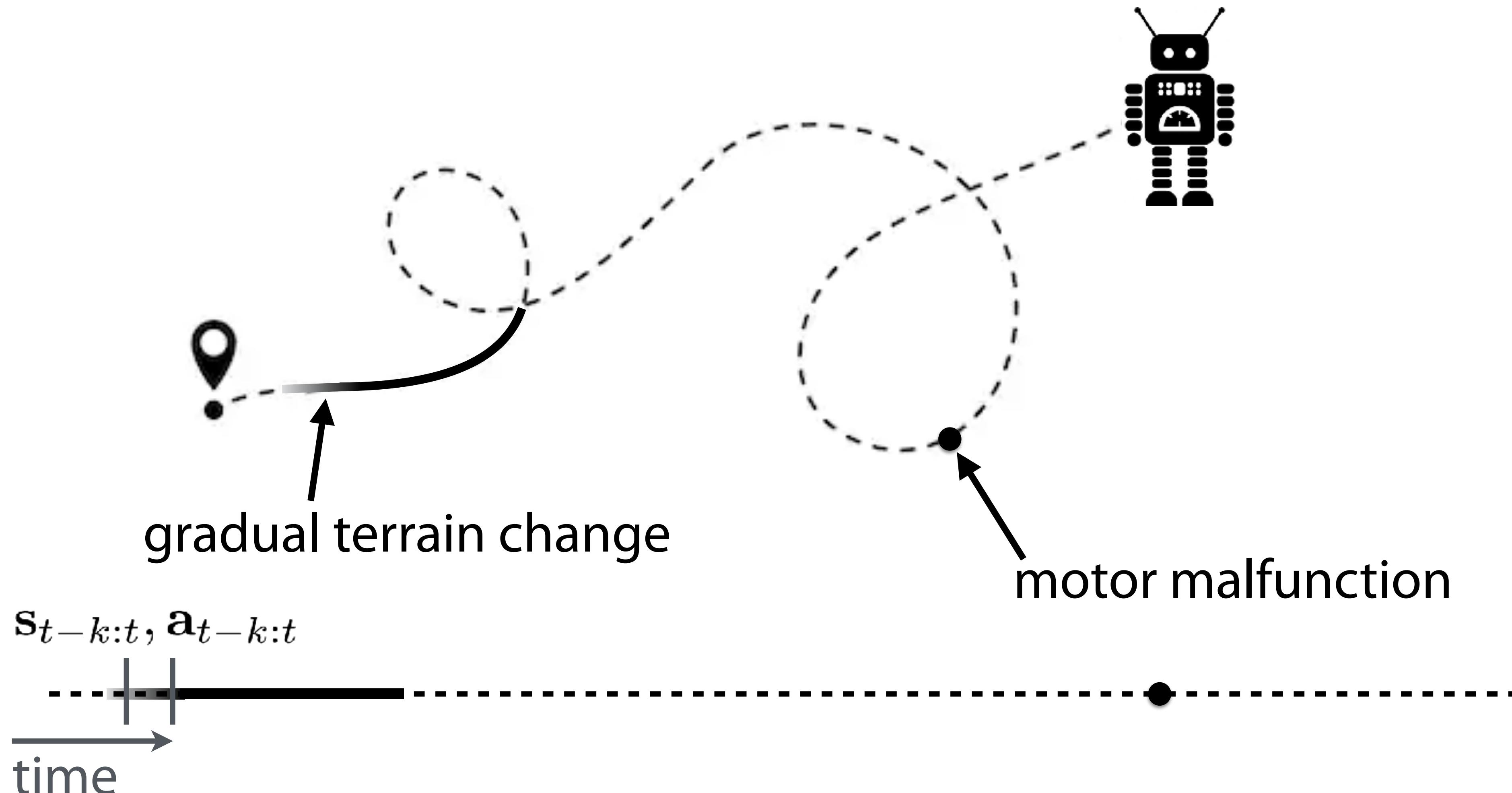


Ignasi Clavera



Simin Liu

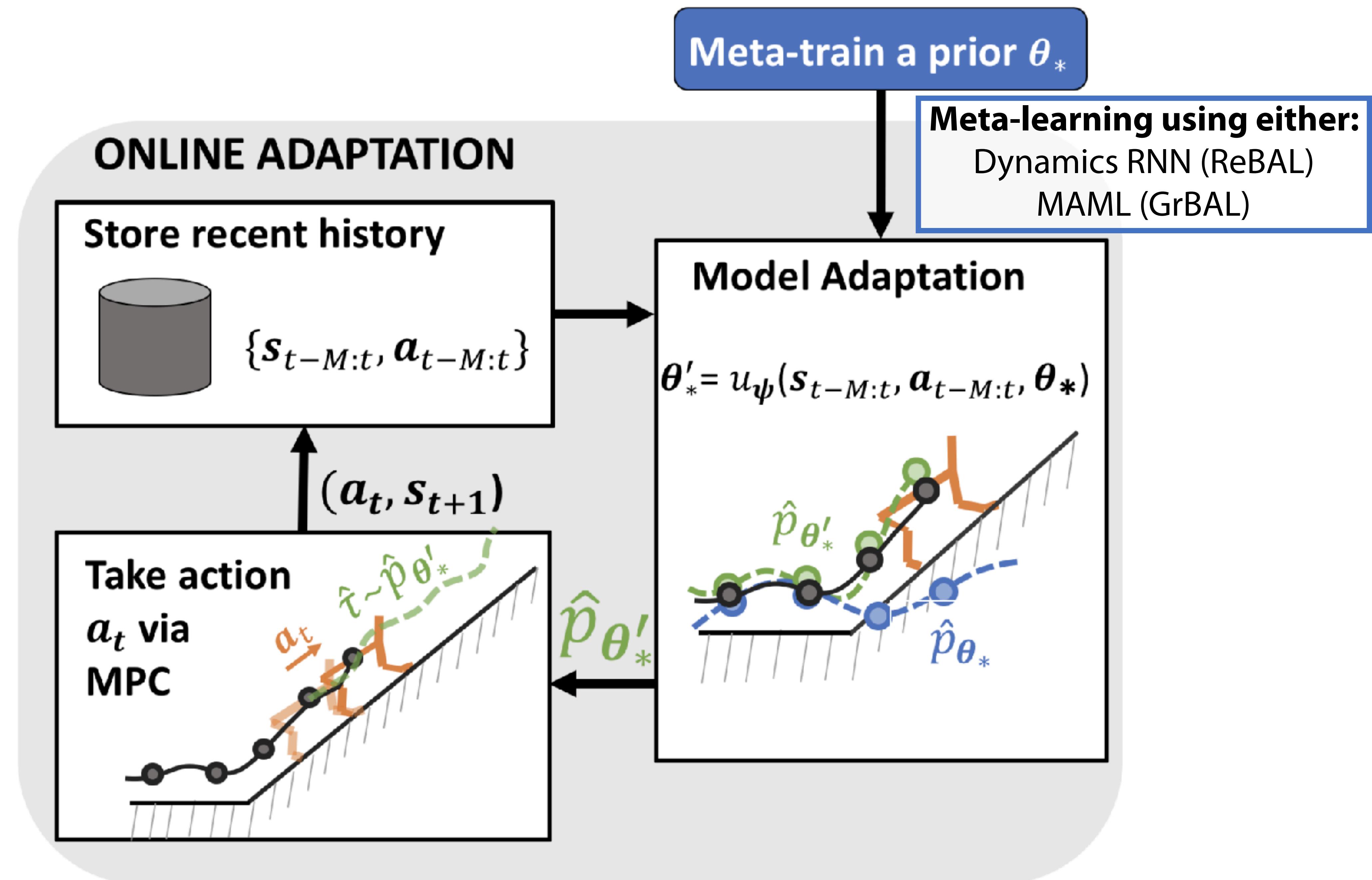
Goal: learn to adapt model quickly to new environments



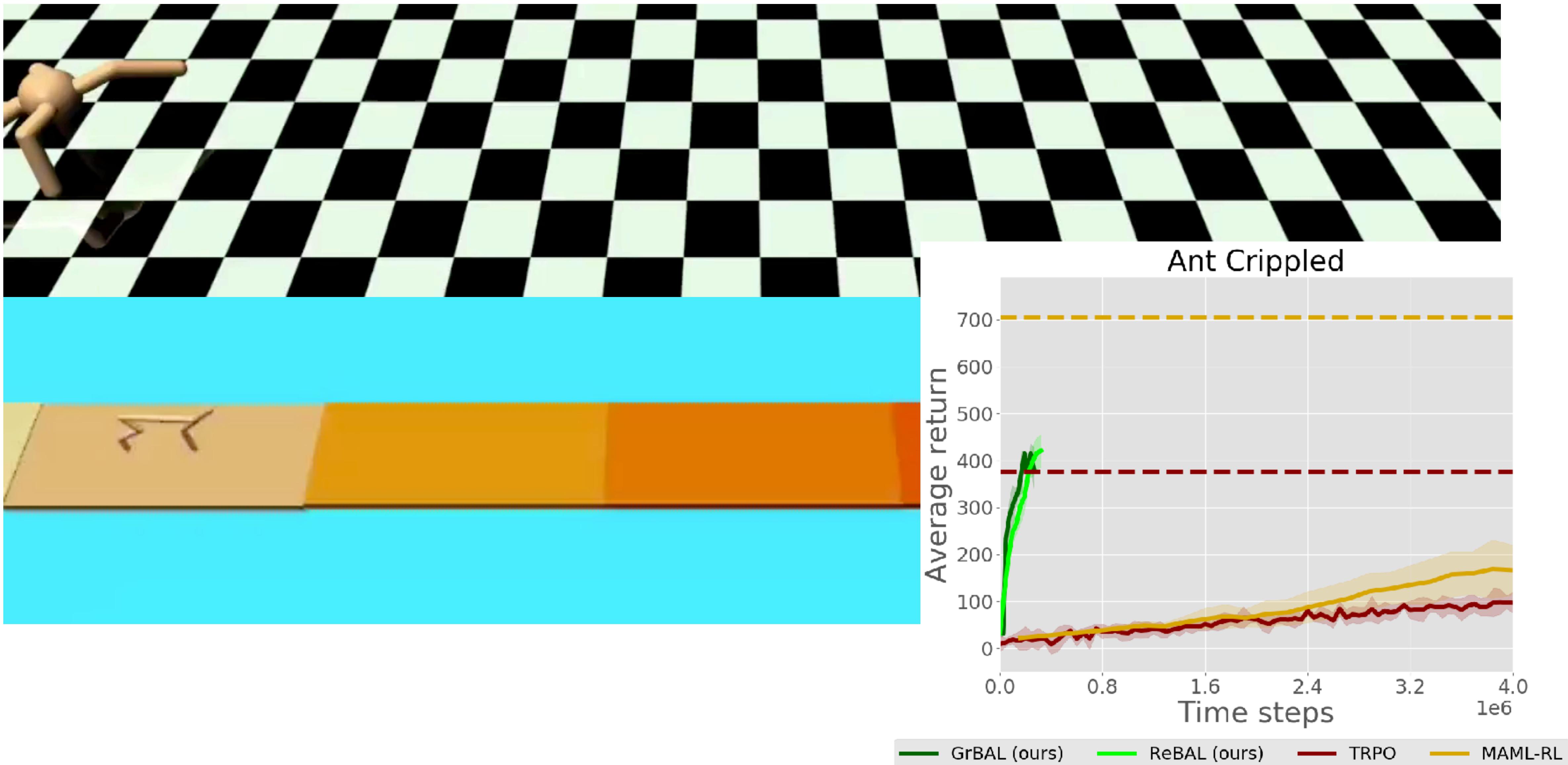
online adaptation = few-shot learning

tasks are temporal slices of experience

The task distribution is more-or-less free!



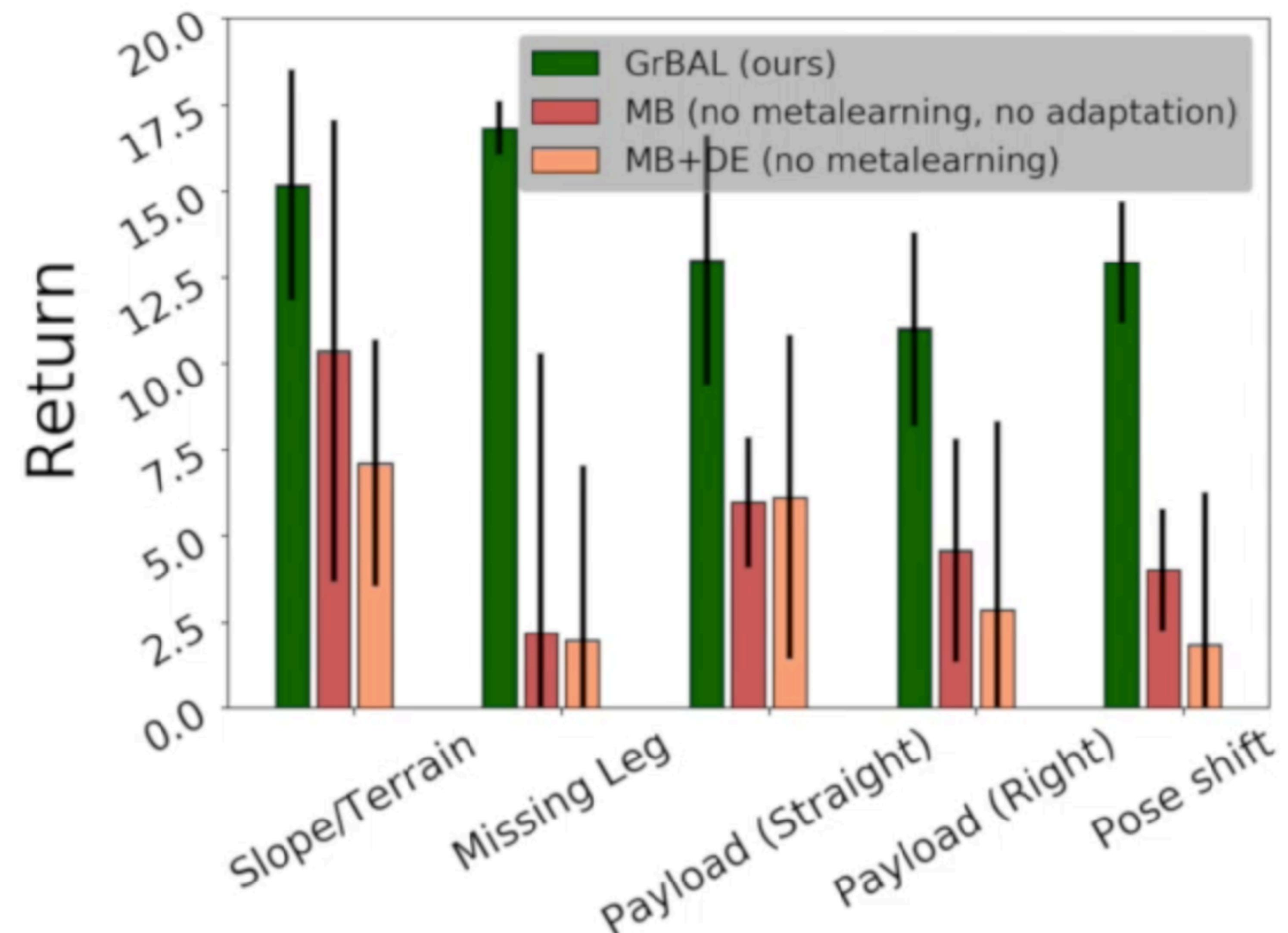
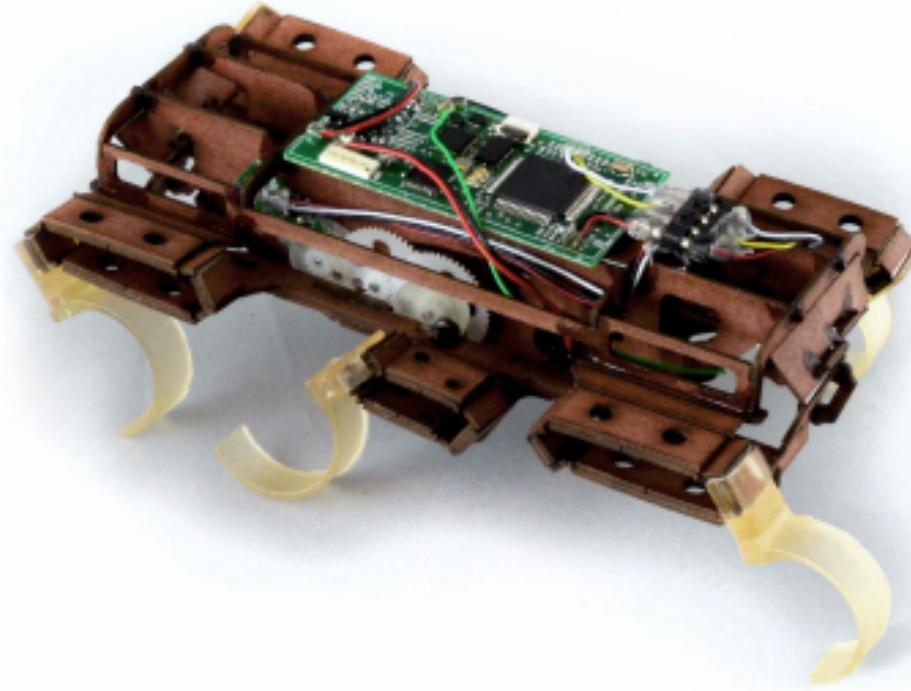
Dynamic Environments with Online Adaptation



Roach Robot

Meta-train on variable terrains

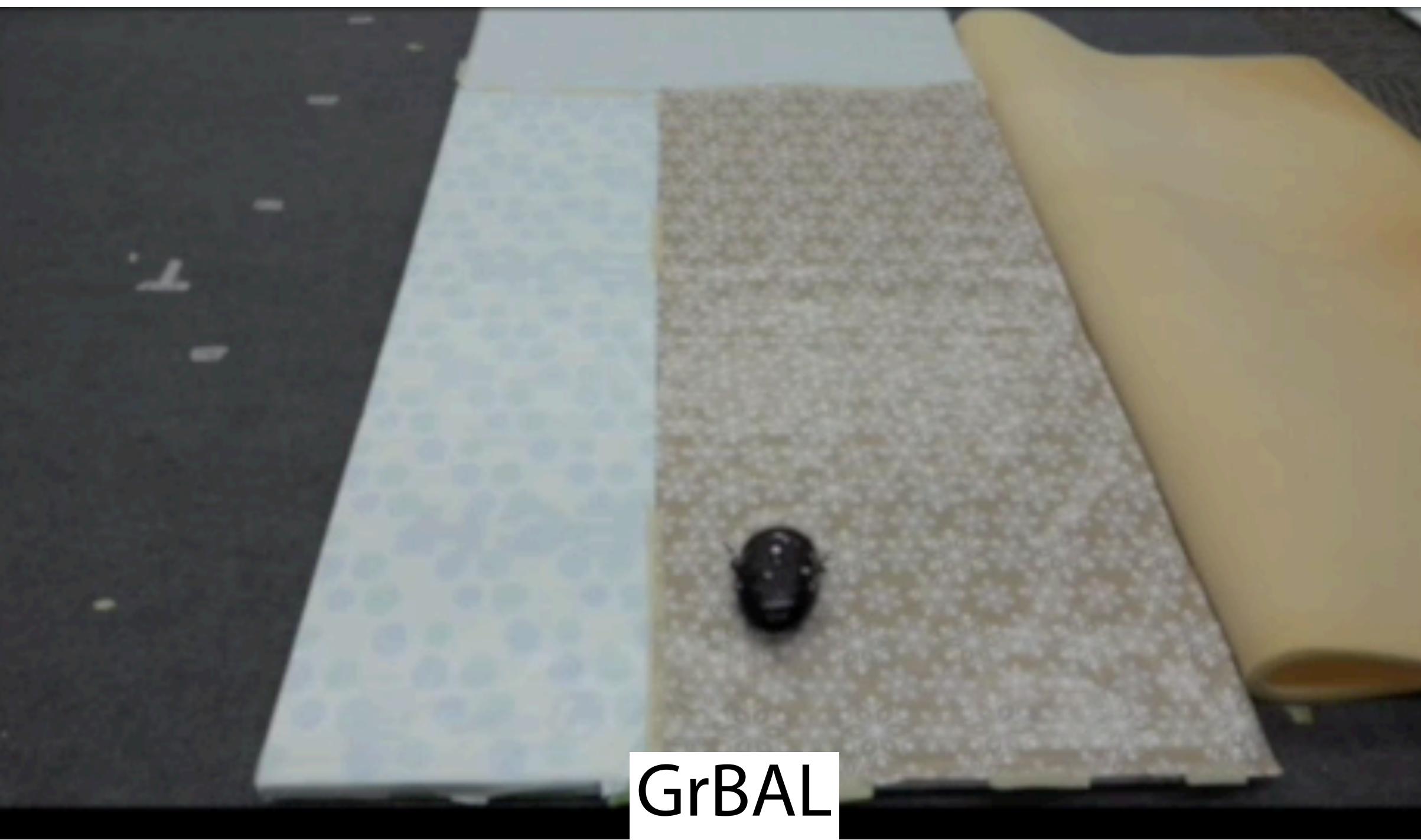
Meta-test with slope, missing leg, payload, calibration errors



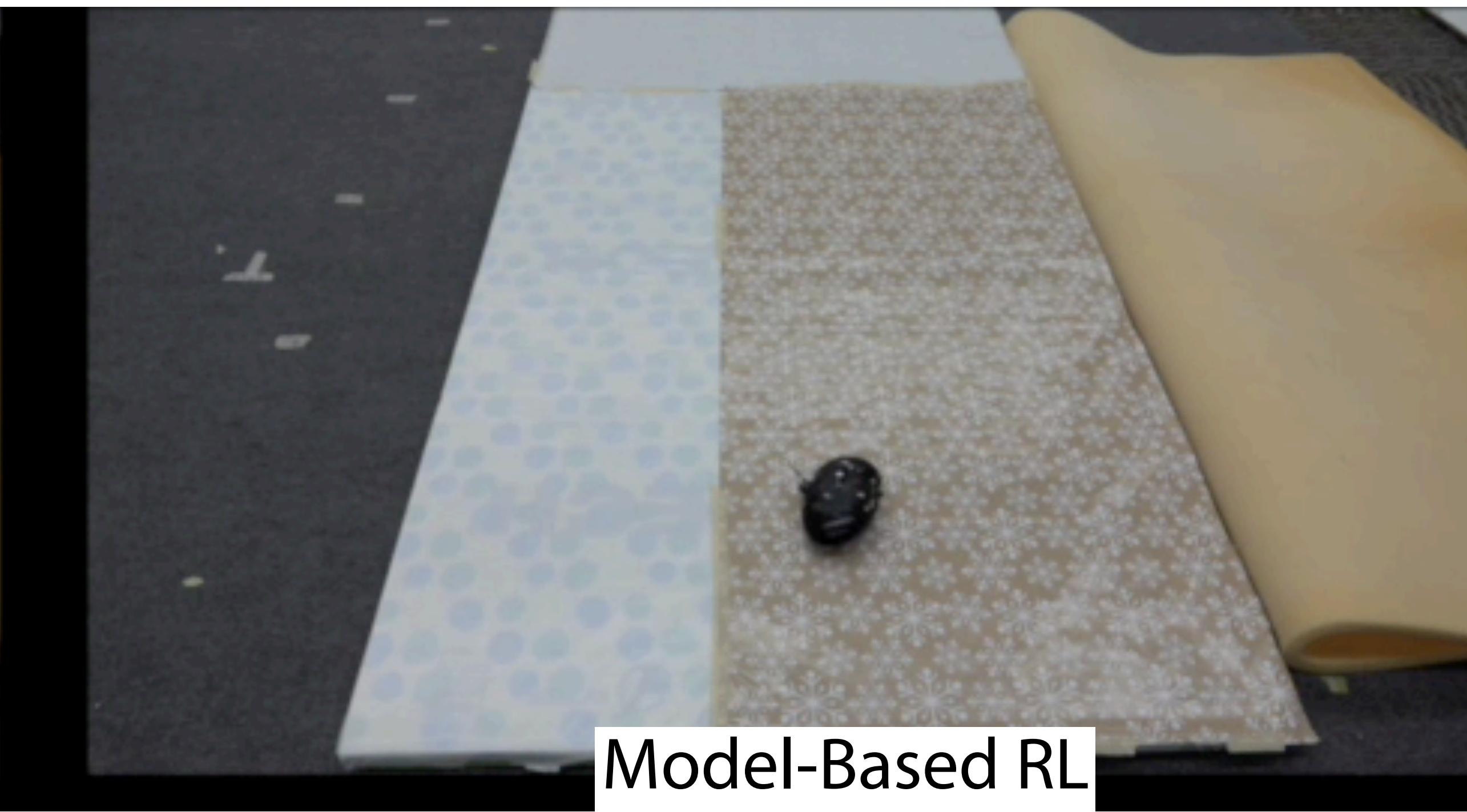
Roach Robot

Meta-train on variable terrains

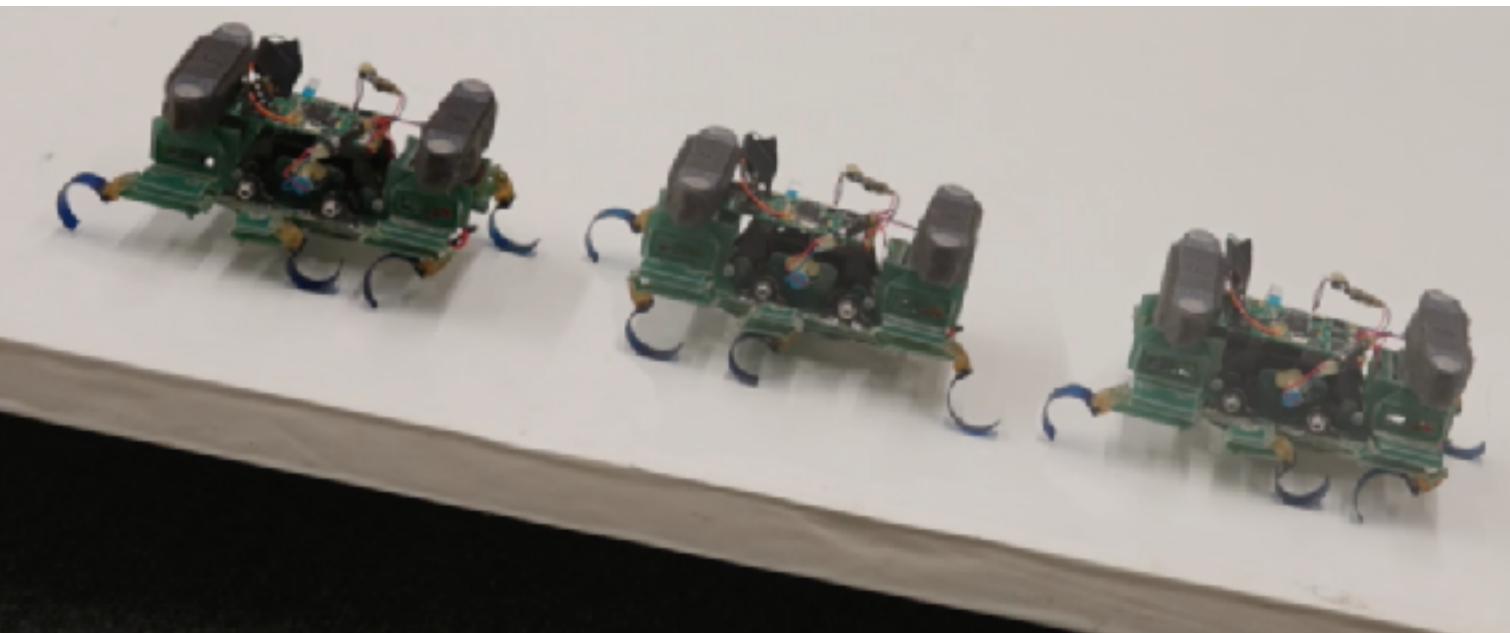
Meta-test with slope, missing leg, payload, calibration errors



GrBAL



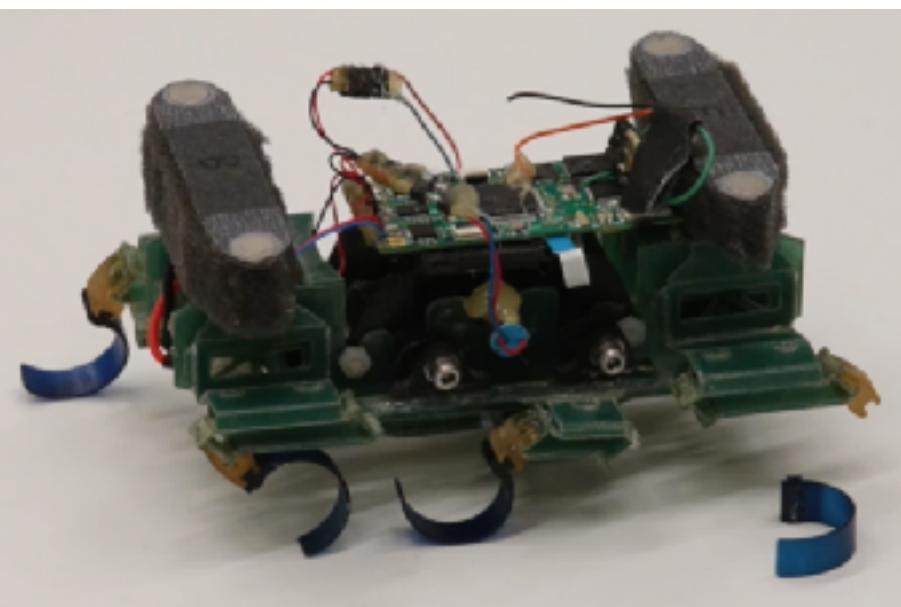
Model-Based RL
(no adaptation)



Roach Robot

Meta-train on variable terrains

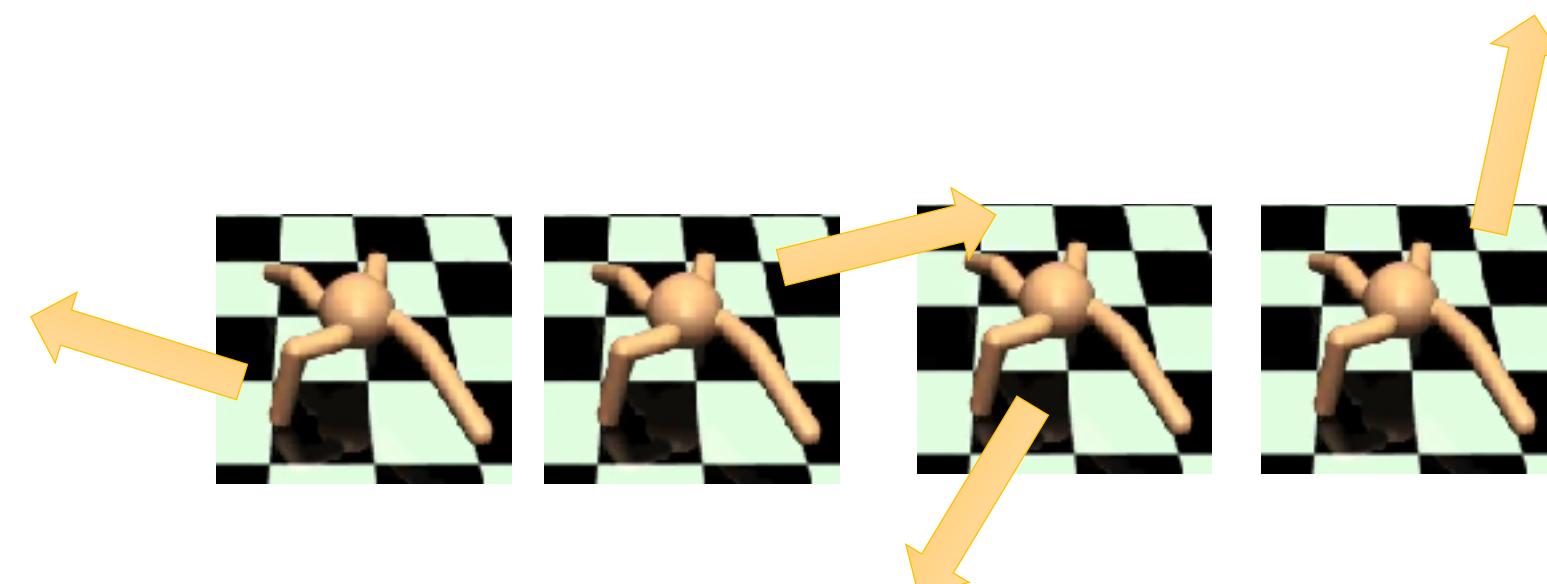
Meta-test with slope, missing leg, payload, calibration errors



Wishlist in a Meta-RL algorithm

	RL ² , SNAIL	MAML-PG	ReBAL, GrBAL	
Consistent	✗	✓	✓	
Expressive	✓	✗	?	
Structured Exploration	?	?	✗	Still <i>lots</i> of room for exploration & improvement!
Efficient & off-policy	✗	✗	✓	

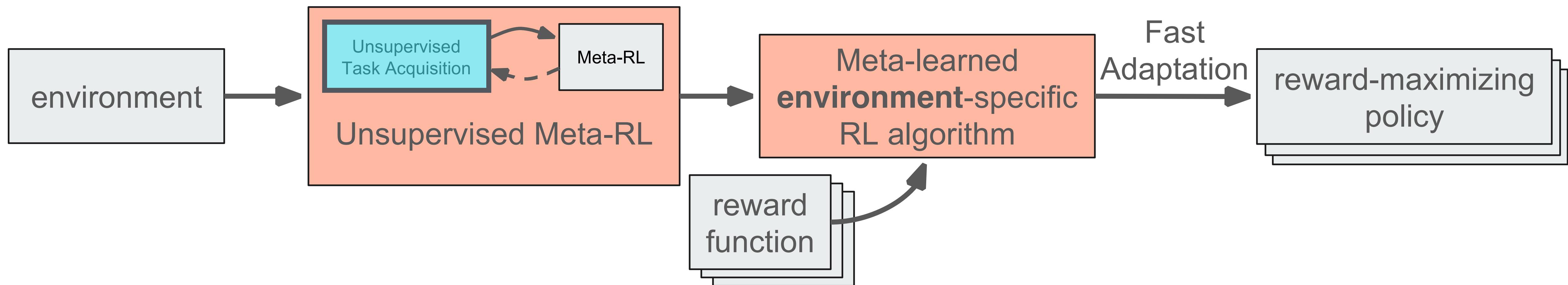
One more important challenge: where do the tasks come from?



So far: manual defined distribution of tasks.
(and corresponding rewards)

Unsupervised Meta-RL

General Recipe

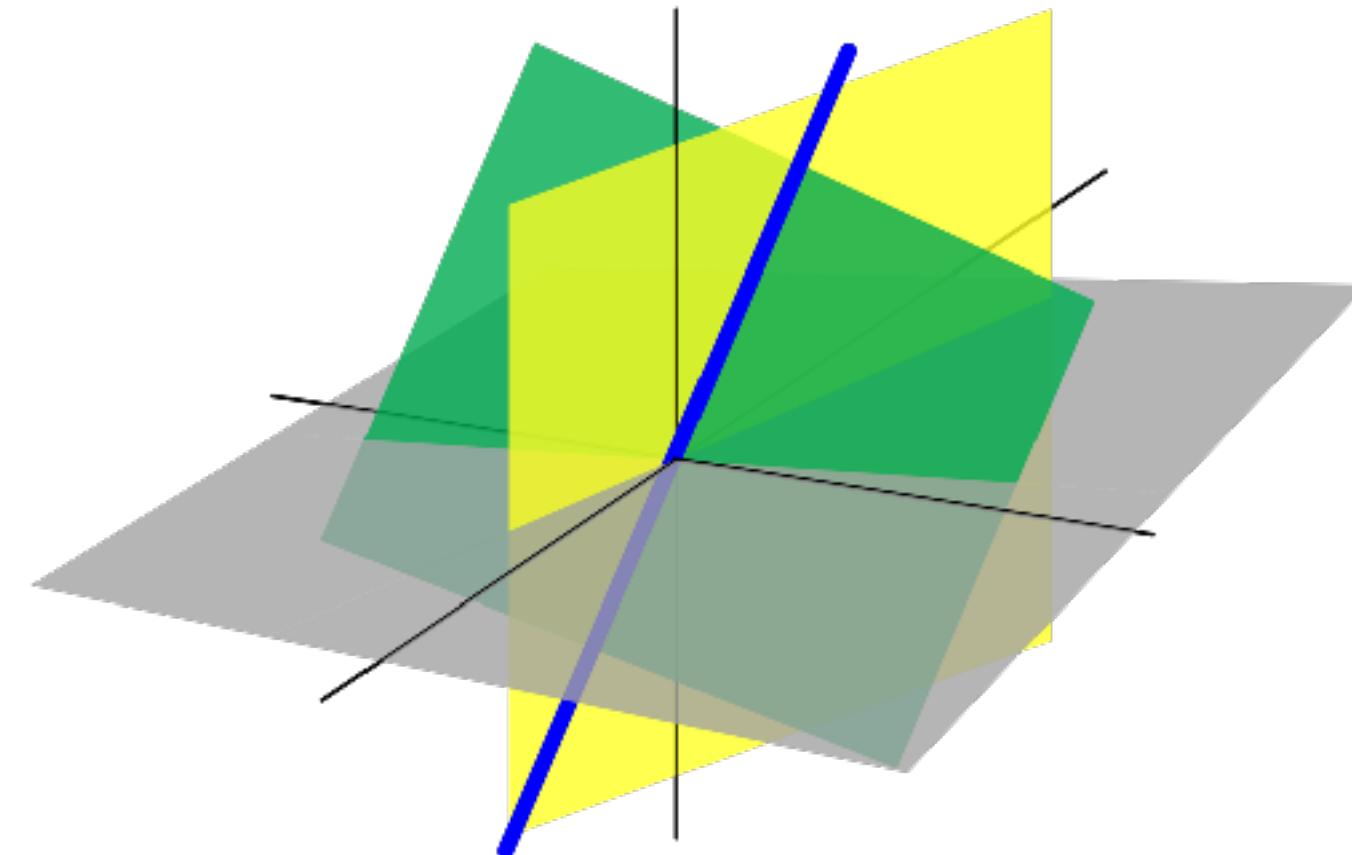
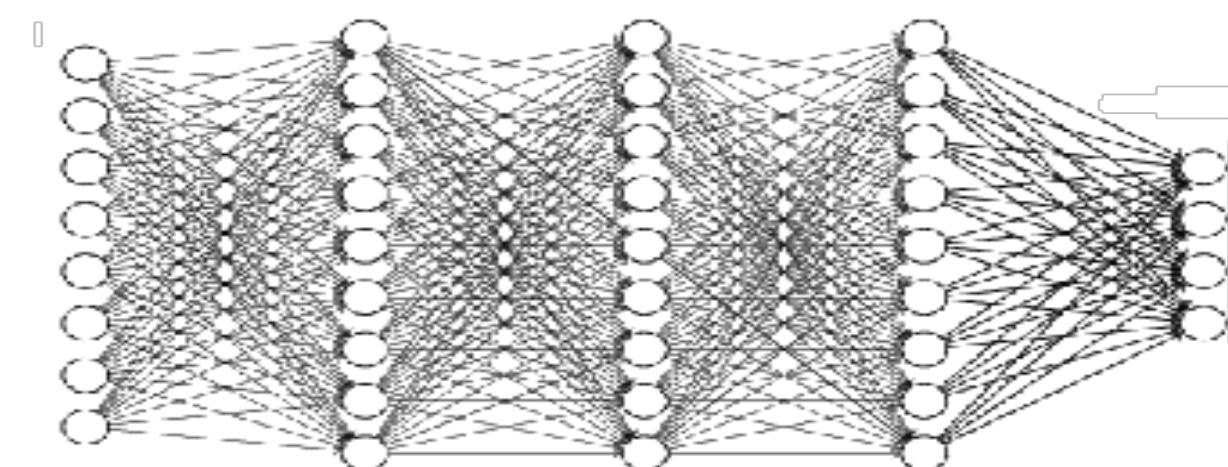


Random Task Proposals

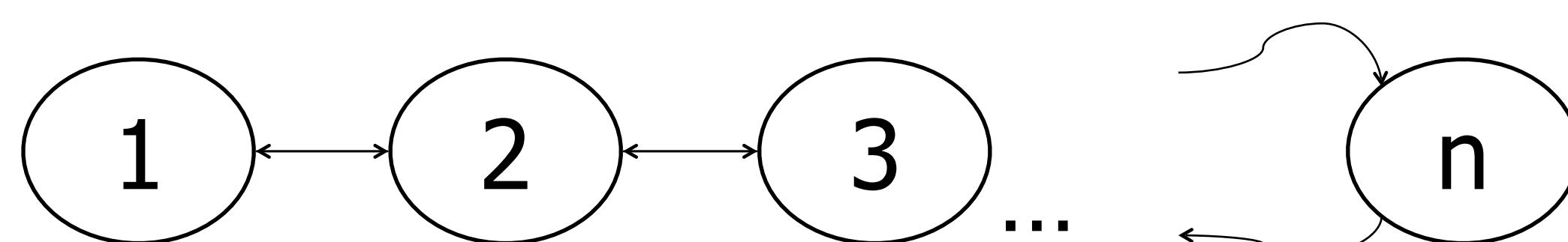
- Use randomly initialize discriminators for reward functions

$$R(s, z) = \log p_D(z|s)$$

D → randomly initialized network

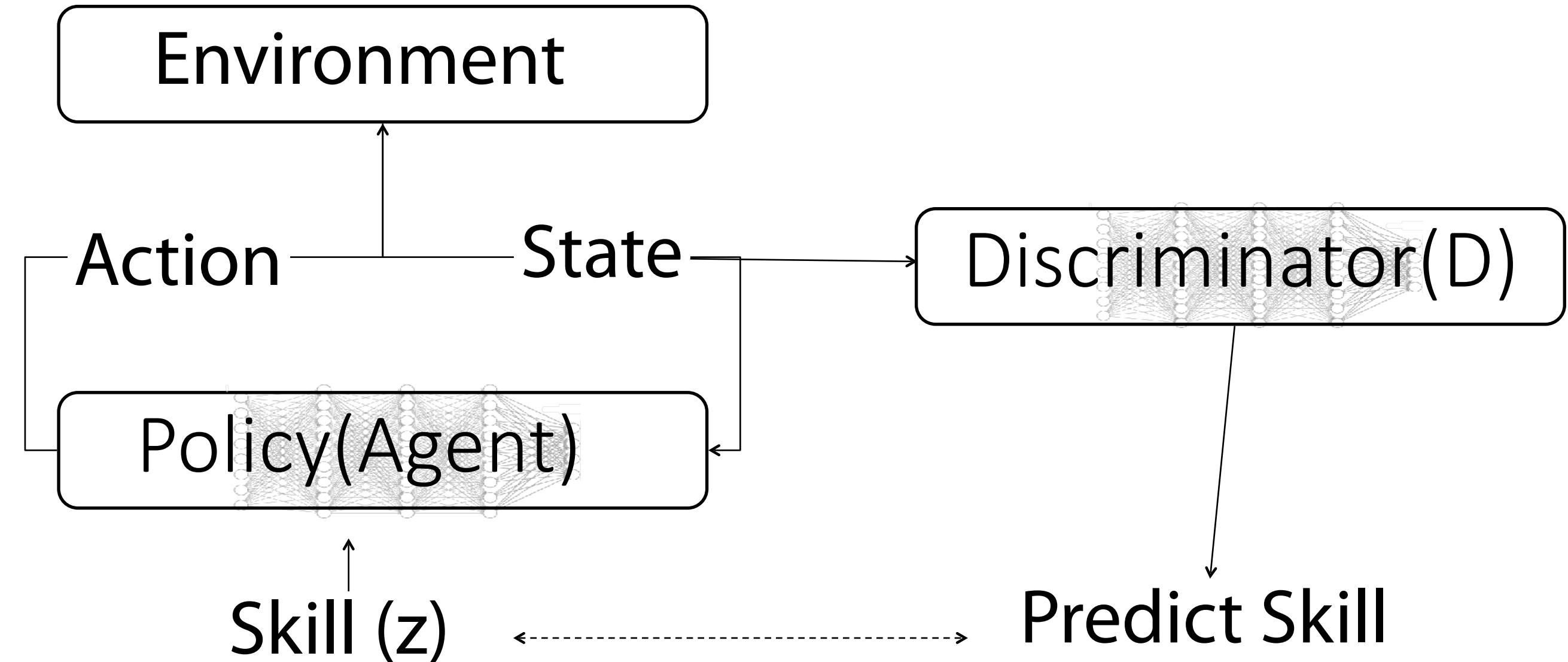


- Important: Random functions over state space, not random policies



Random policy – exponential
Random reward – polynomial

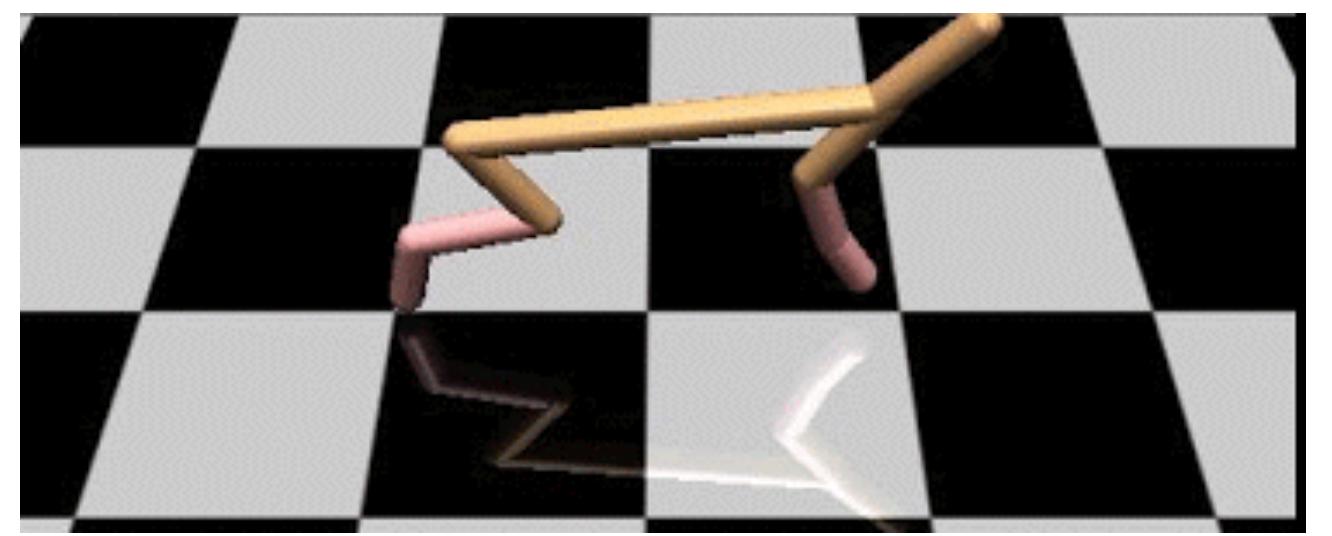
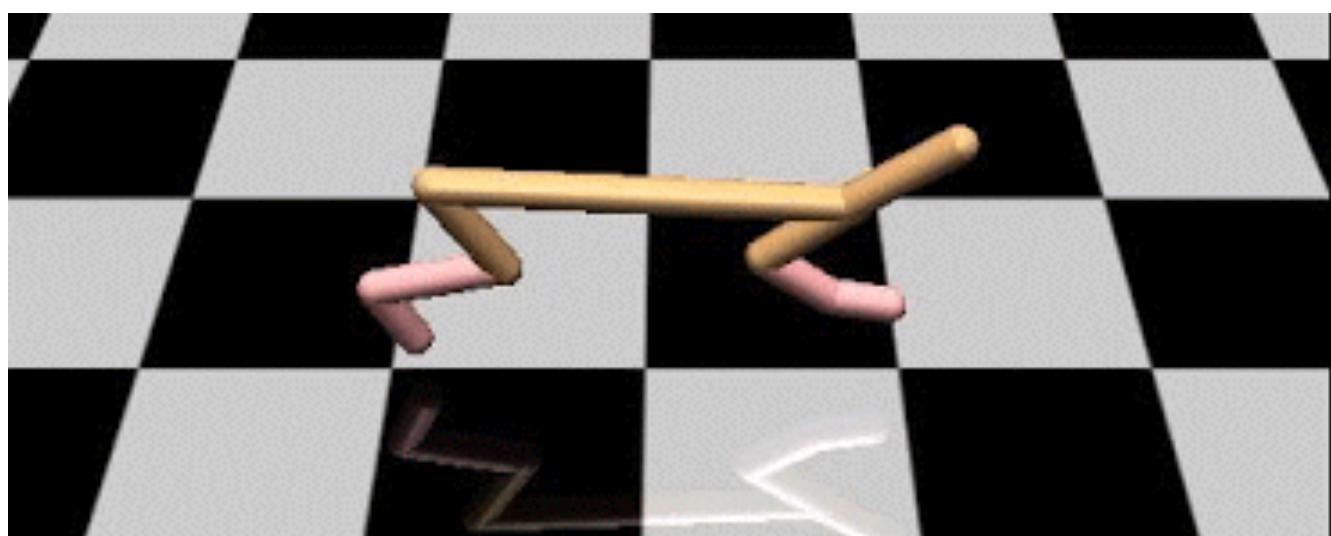
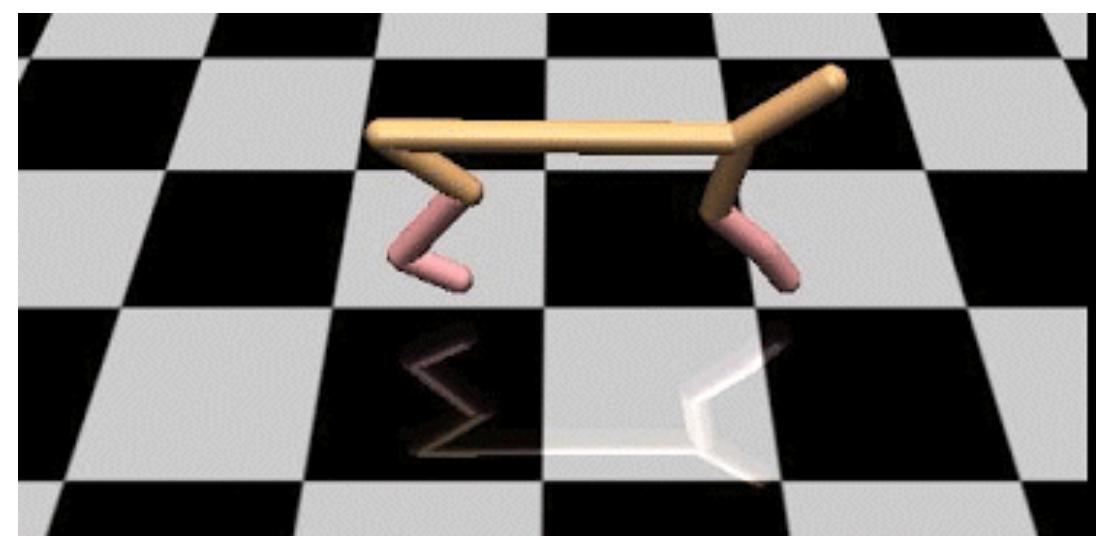
Diversity-Driven Proposals



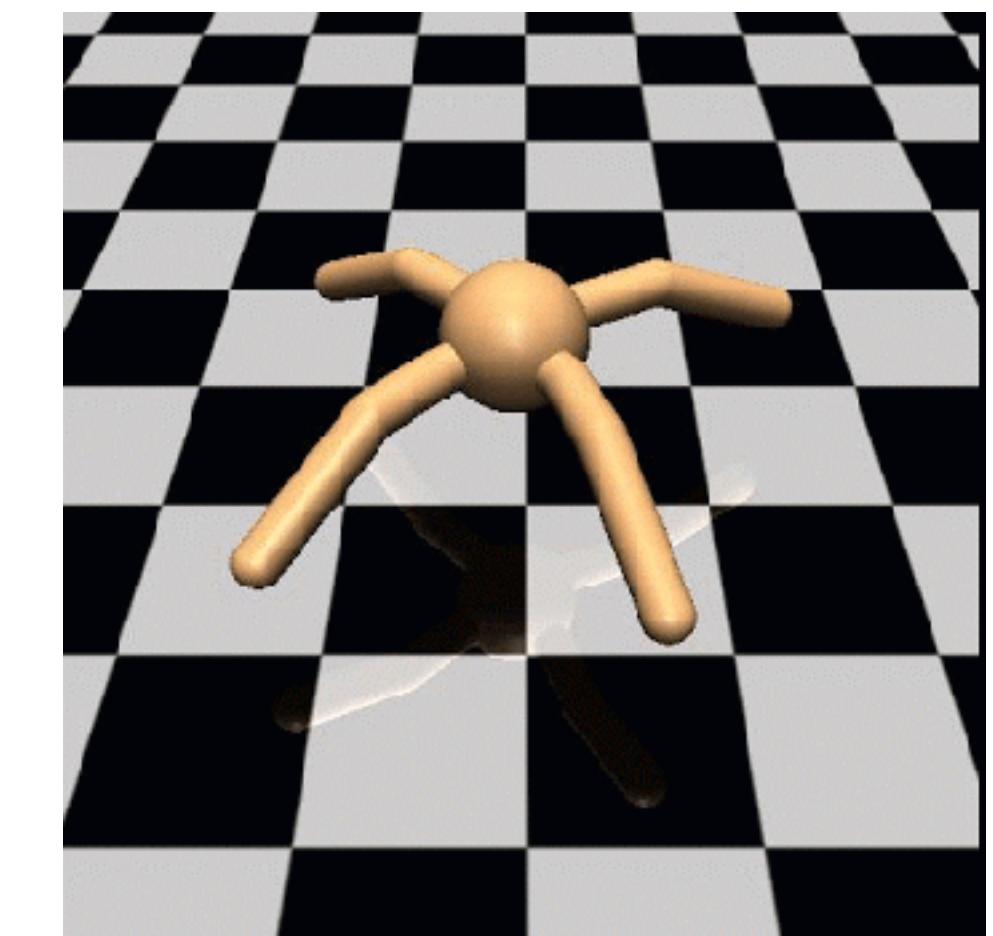
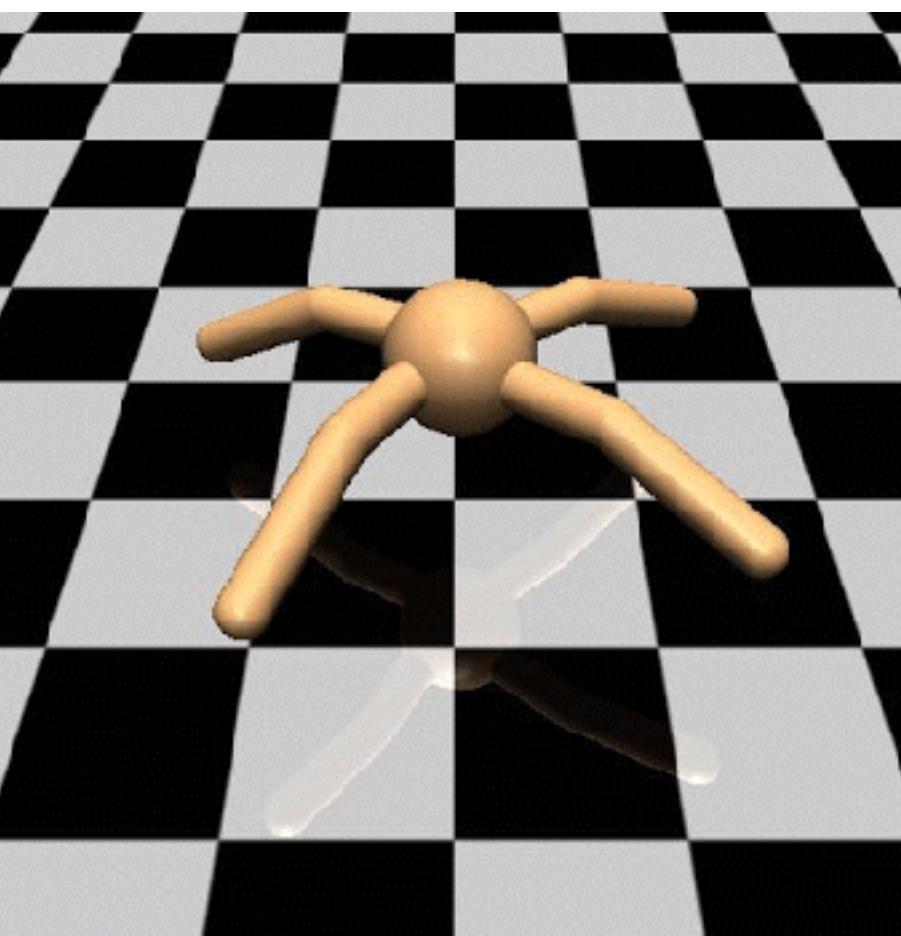
- Policy → visit states which are discriminable
- Discriminator → predict skill from state

Task Reward for UML: $R(s, z) = \log p_D(z|s)$

Examples of Acquired Tasks



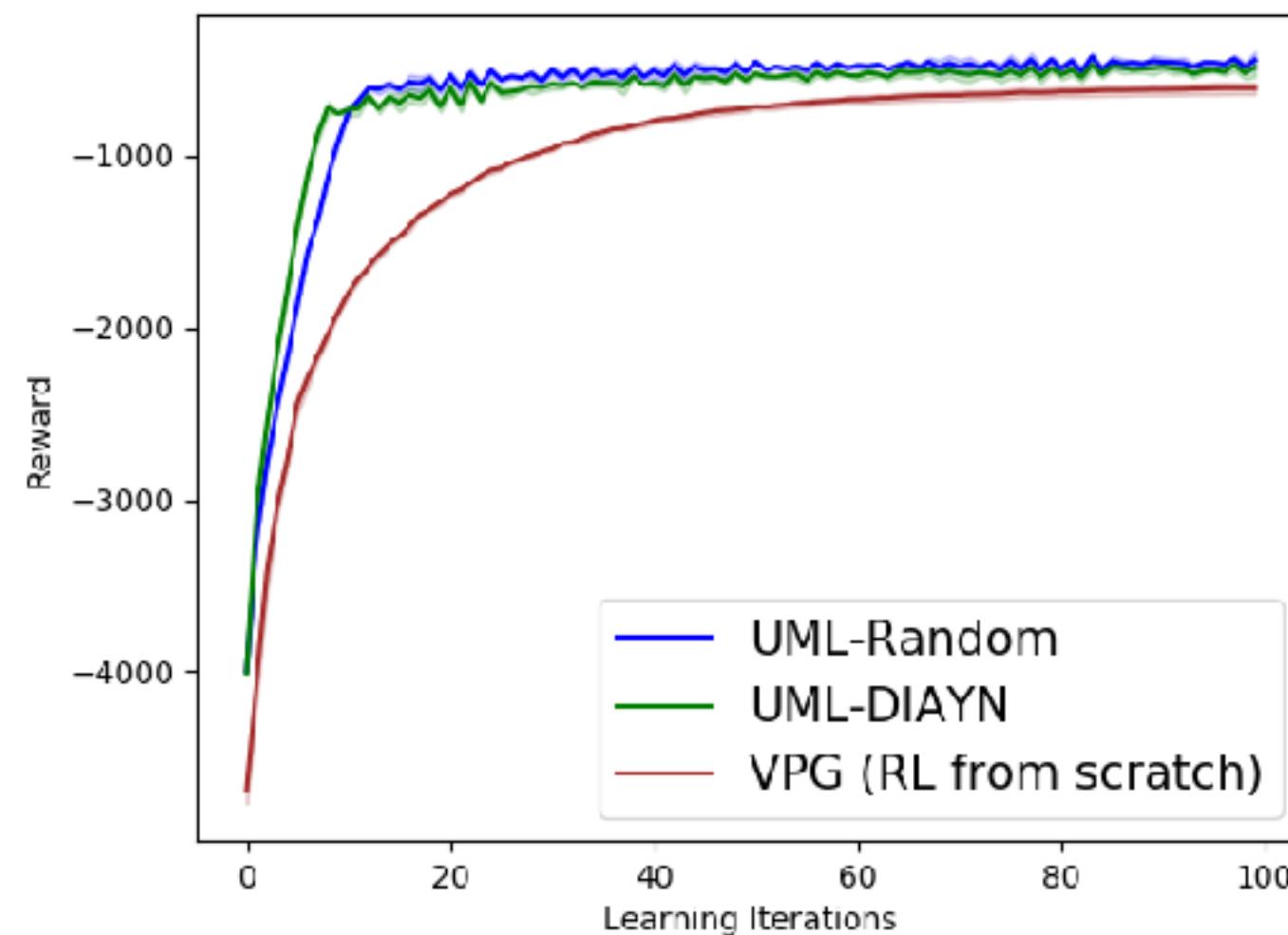
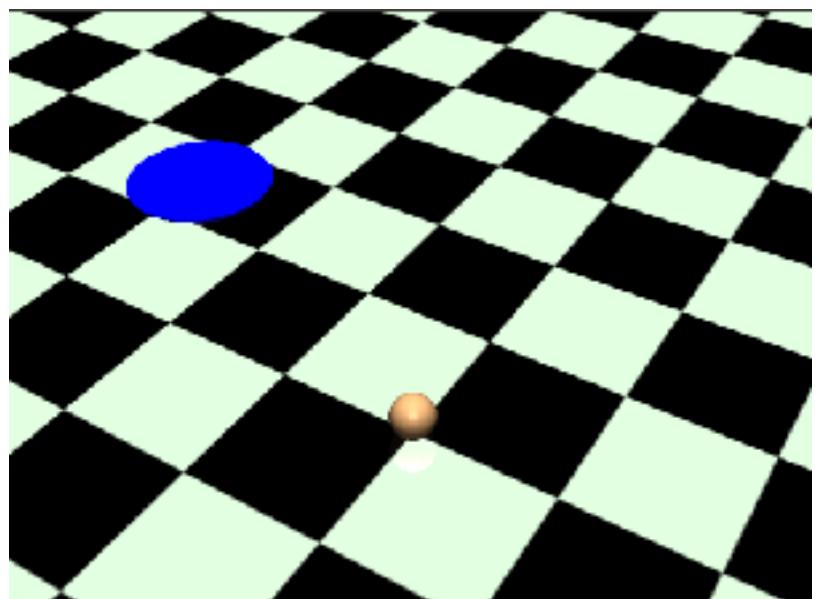
Cheetah



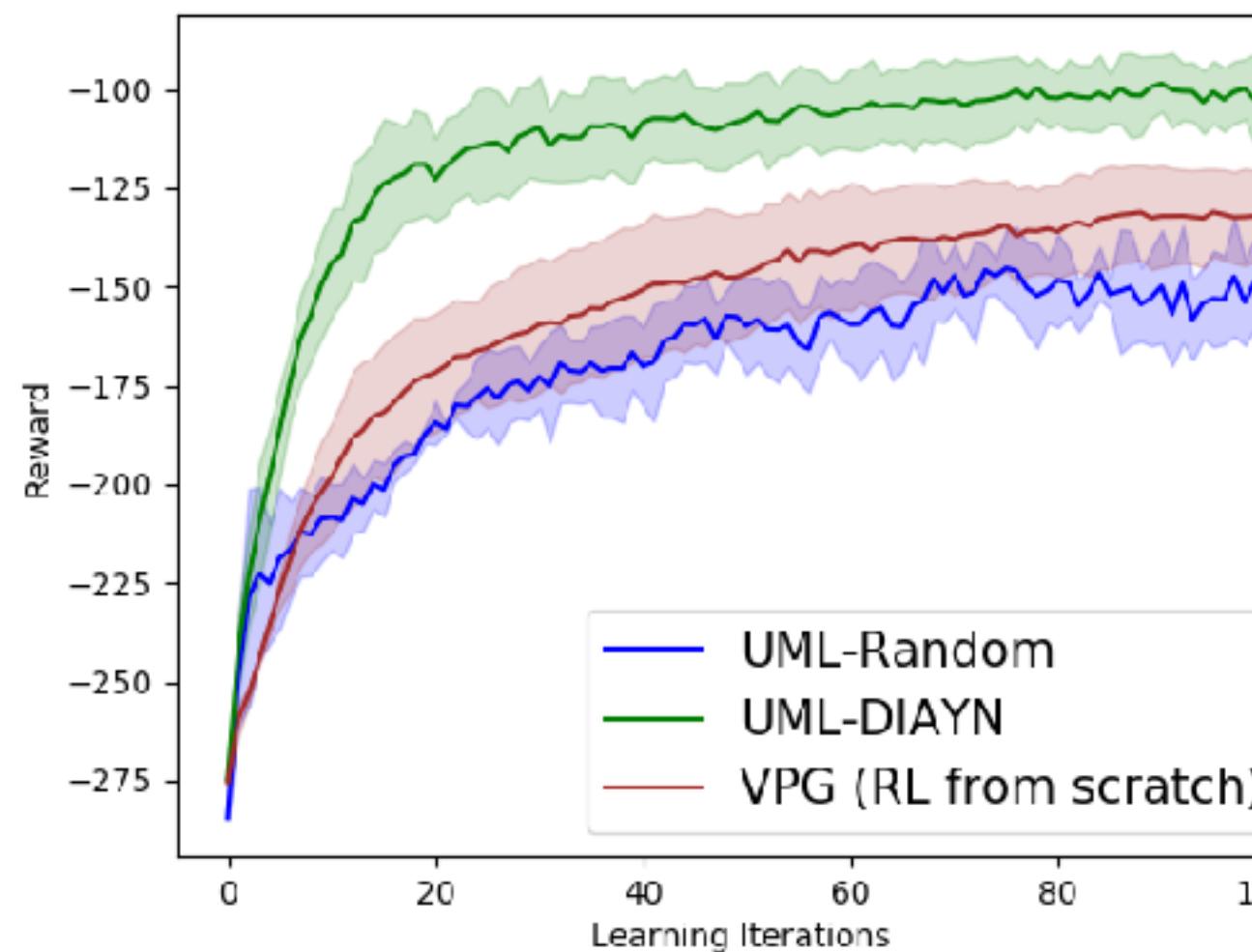
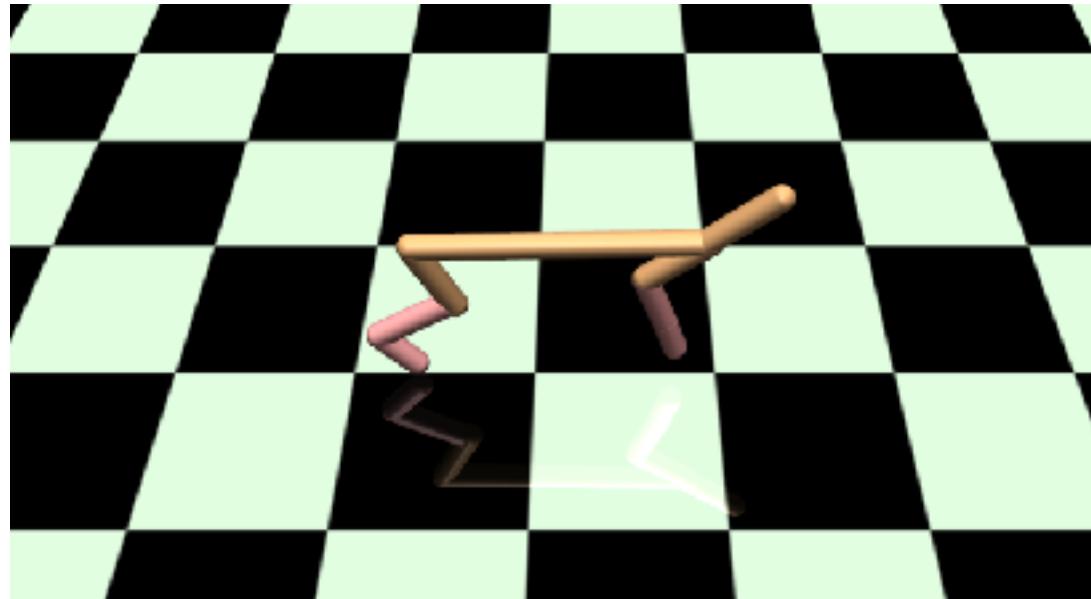
Ant

Does it work?

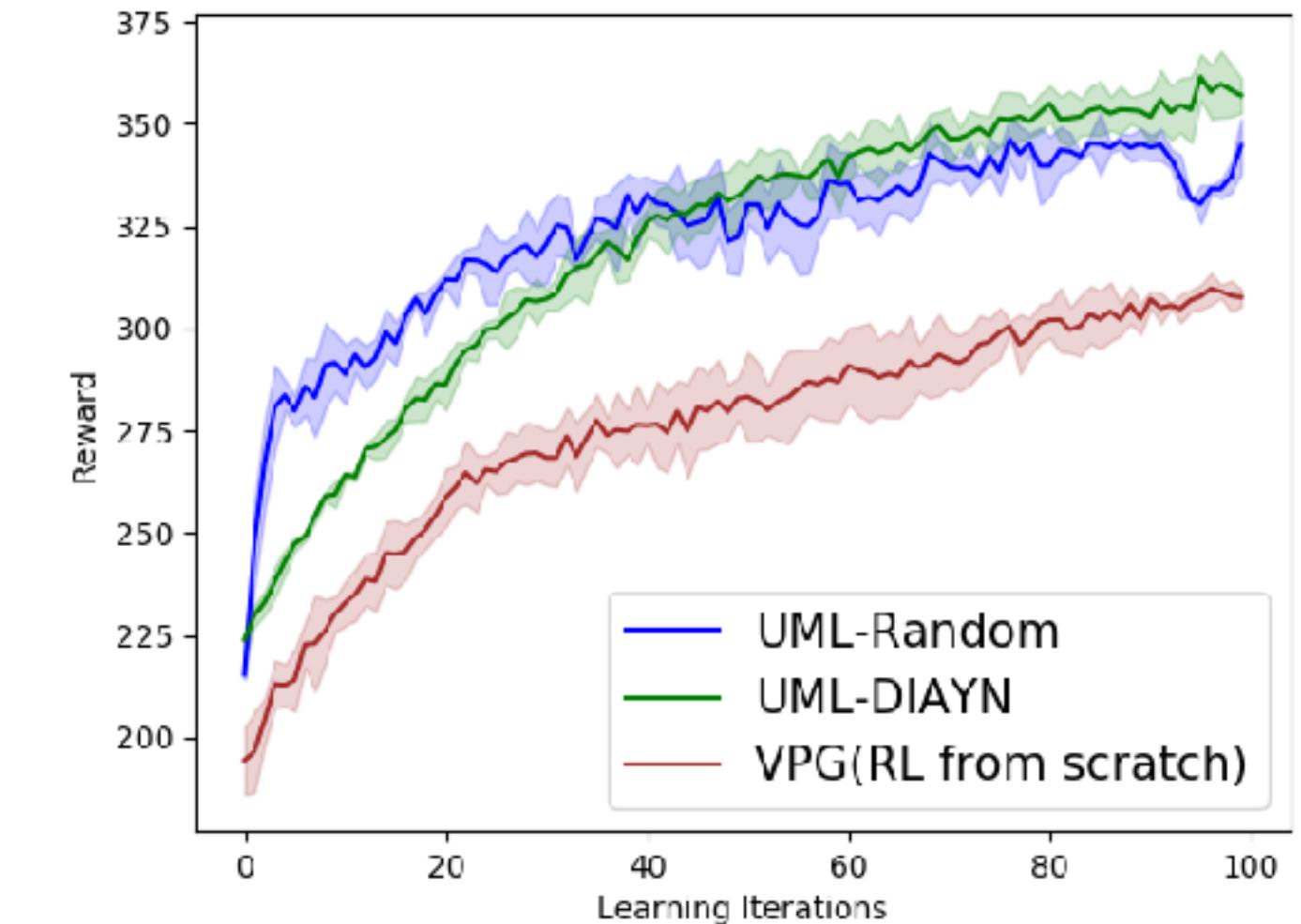
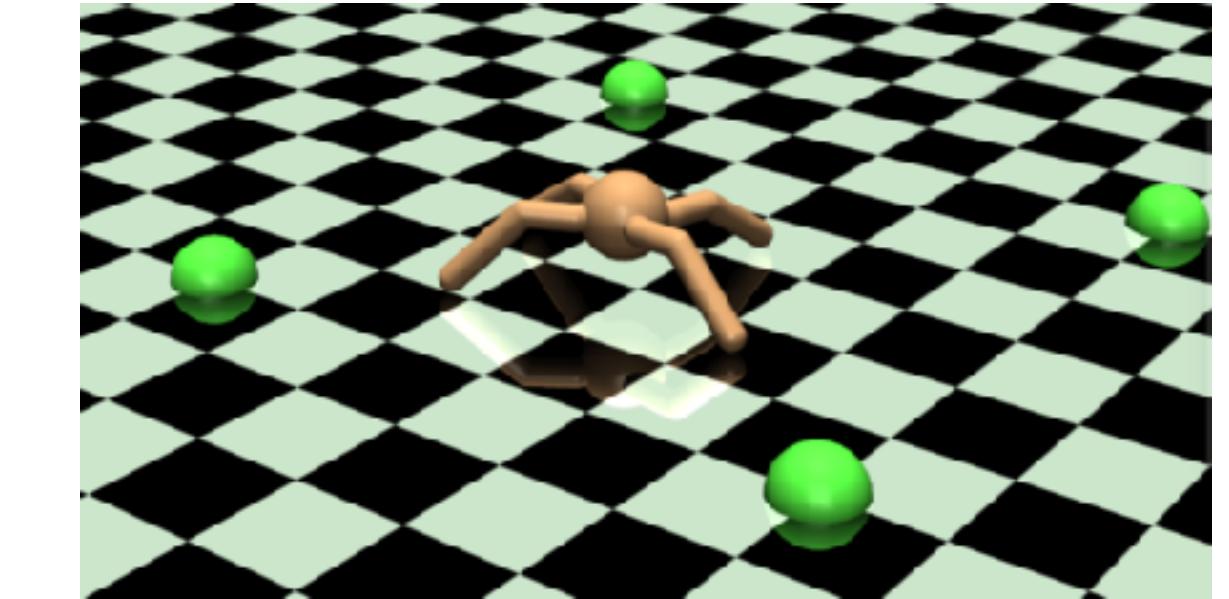
2D Navigation



Cheetah



Ant



Meta-test performance with rewards

Takeaway: Relatively simple mechanisms for proposing tasks work surprisingly well.

Key References

Recurrent Models: Yan Duan et al. '17 (RL^{^2}), Jane Wang et al. CogSci '17 (Learning to Reinforcement Learn), Nikhil Mishra et al. ICLR '18 (SNAIL)

Model-Agnostic Meta-Learning: Finn et al. ICML '17

Further Reading

Value-based: Sung et al. arXiv '17 (*Meta-critic networks*)

Exploration: Gupta et al. NIPS '18, Stadie et al. NIPS '18

Unsupervised: Gupta et al. arXiv '18 (*Unsupervised Meta-Learning for Reinforcement Learning*)

Model-based: Nagabandi*, Clavera* et al. arXiv '18

Evolutionary Strategies: Houthooft et al. NIPS '18

Cognitive Science: Wang et al. Nature Neuroscience '18 (*PFC as a Meta-RL System*)

Takeaway

Don't reinforcement learn from scratch.
Meta-RL provides an approach to do this.

Open Problems

- a meta-RL algorithm that is both **consistent** and **expressive**
- meta-RL algorithms that learn to **explore in a structured way**
- constructing **task distributions** for meta-RL in an automated way
- meta-RL algorithms that can **incorporate off-policy data** effectively

Collaborators

Sergey Levine



Pieter Abbeel



Anusha Nagabandi



Ignasi Clavera



Simin Liu



Ron Fearing



Abhishek Gupta



Ben Eysenbach



