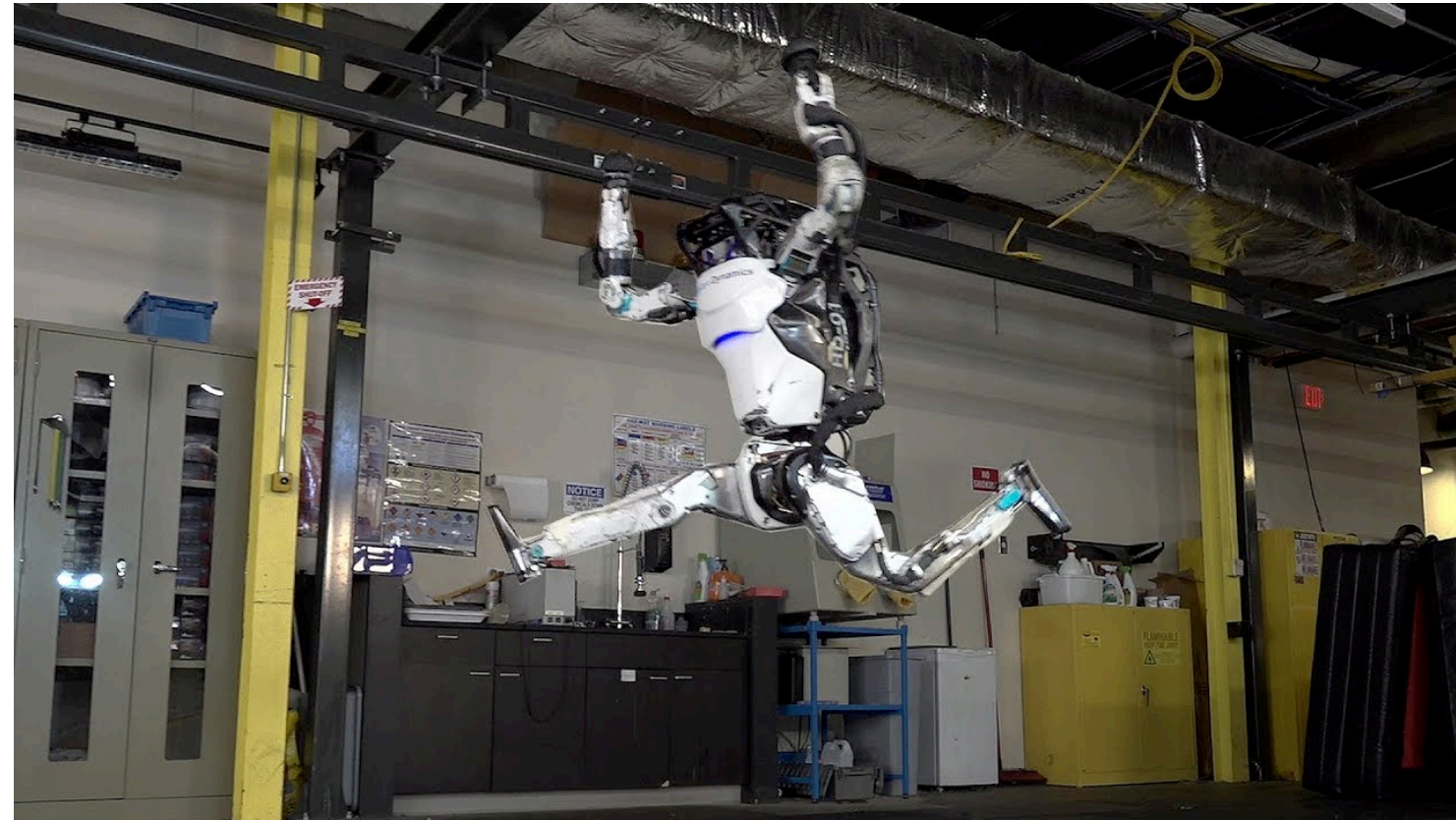


PWM: Policy Learning with Large World Models

Ignat Georgiev, Varun Giridhar, Nicklas Hansen, Animesh Garg

Motivation



Agile locomotion



Dexterous Manipulation

Common failure: contact-rich tasks!



Whole-Body Loco-Manipulation

The current state of robot learning

- **Reinforcement Learning (RL)**

- Formulated as MDP reward maximization
- Notoriously sample-inefficient
- Mostly simulation based
- Can do any single task given good simulation and data

- $$\max_{\theta} \mathbb{E}_{\substack{s_1 \sim \rho \\ a_h \sim \pi(\cdot | s_h)}} \left[\sum_{h=1}^H r(s_h, a_h) \right]$$

- **Behavior cloning (BC)**

- Formulated as supervised learning
- Very capable given optimal data
- Impressive multitasking
- Currently mostly simple tasks
- Lives and dies by its data

- $$\min_{\theta} \mathbb{E}_{\hat{a}_{h:h+k} \sim \pi(\cdot | s_h)} \|\hat{a}_{h:h+k} - a_{h:h+k}\|_2^2$$

Taking a page from deep learning

- Large models
- Large data
- Efficient optimization - SGD

The current state of robot learning

- Large models ✓
- Large data ✓
- Efficient optimization - SGD ✓
- **Behavior cloning (BC)**
 - Formulated as supervised learning
 - Very capable given optimal data
 - Impressive multitasking
 - Currently mostly simple tasks
 - Lives and dies by its data
 - $\min_{\theta} \mathbb{E}_{\hat{a}_{h:h+k} \sim \pi(\cdot | s_h)} \|\hat{a}_{h:h+k} - a_{h:h+k}\|_2^2$

The current state of robot learning

- **Reinforcement Learning (RL)**

- Formulated as MDP reward maximization
- Notoriously sample-inefficient
- Mostly simulation based
- Can do any single task given good simulation and data

- Large models 

- Large data 

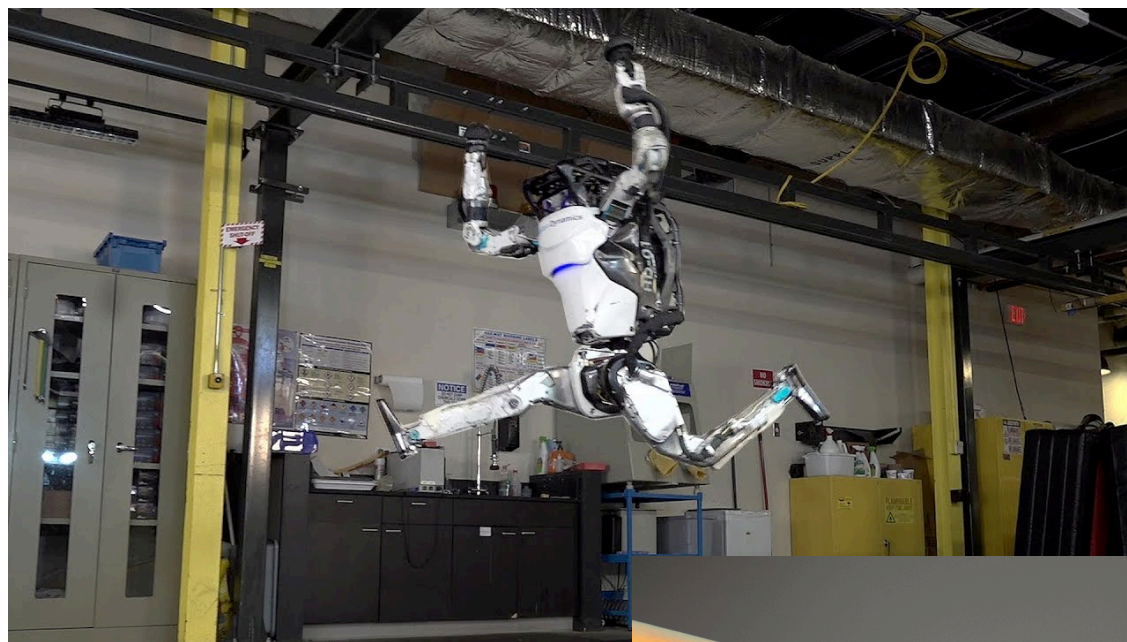
- Efficient optimization - SGD 

Most common: policy gradients (ZoG)

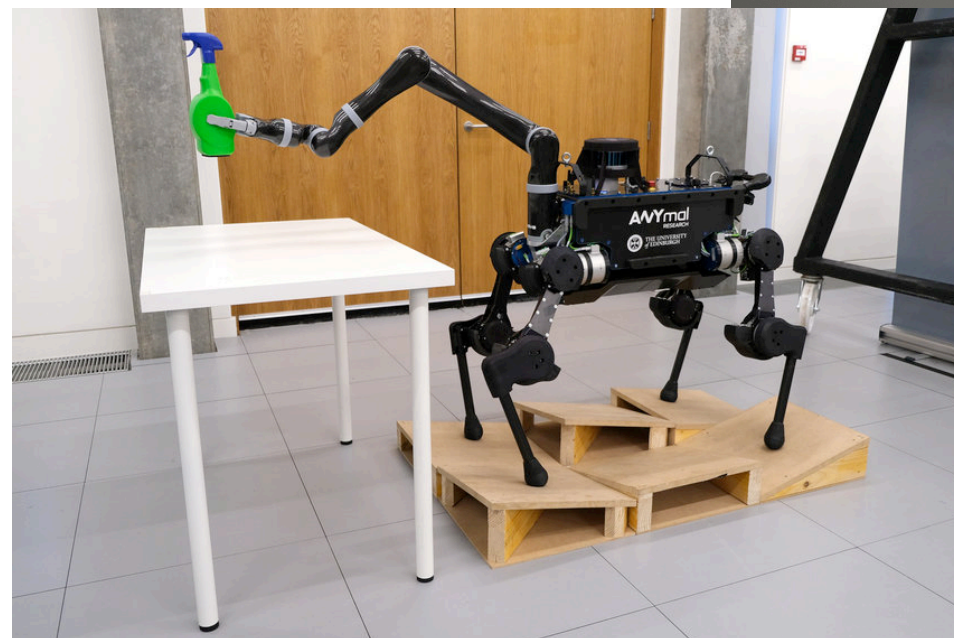
$$\nabla_{\theta}^{[0]} J(\theta) := \mathbb{E}_{a_h \sim \pi_{\theta}(\cdot | s_h)} \left[\left(\sum_{h=1}^H r(s_h, a_h) \right) \left(\sum_{h=1}^H \nabla_{\theta} \log \pi_{\theta}(a_h | s_h) \right) \right]$$

What are the current possibilities

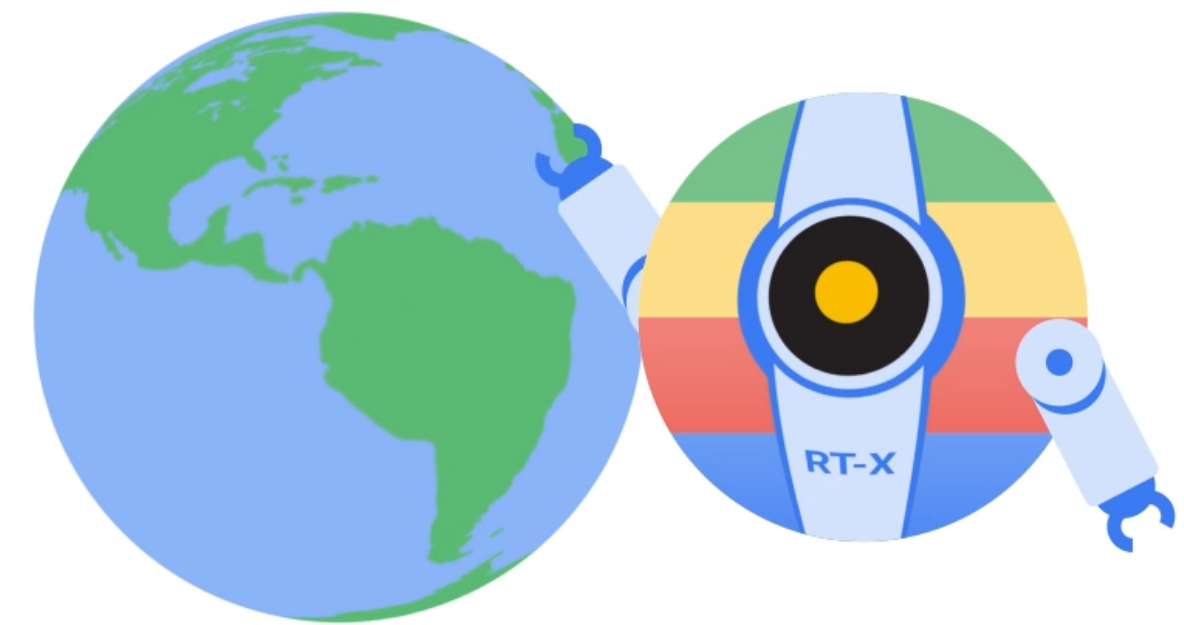
- RL



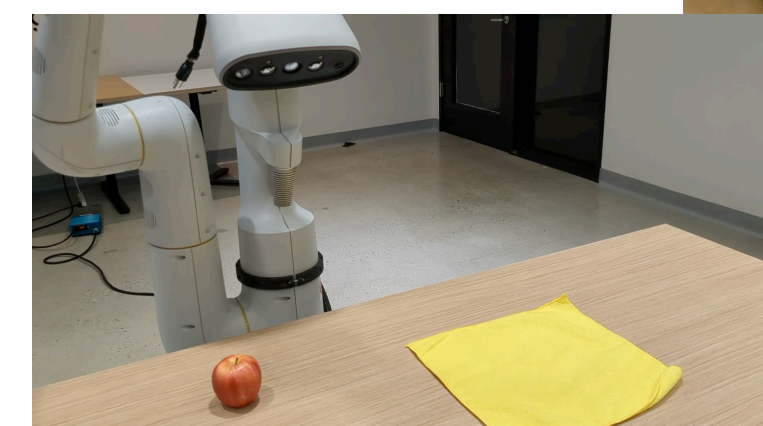
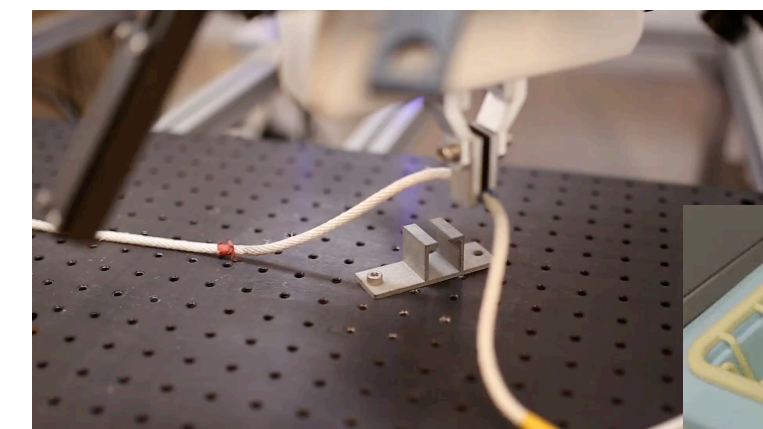
Maybe SGD is the answer?



- BC

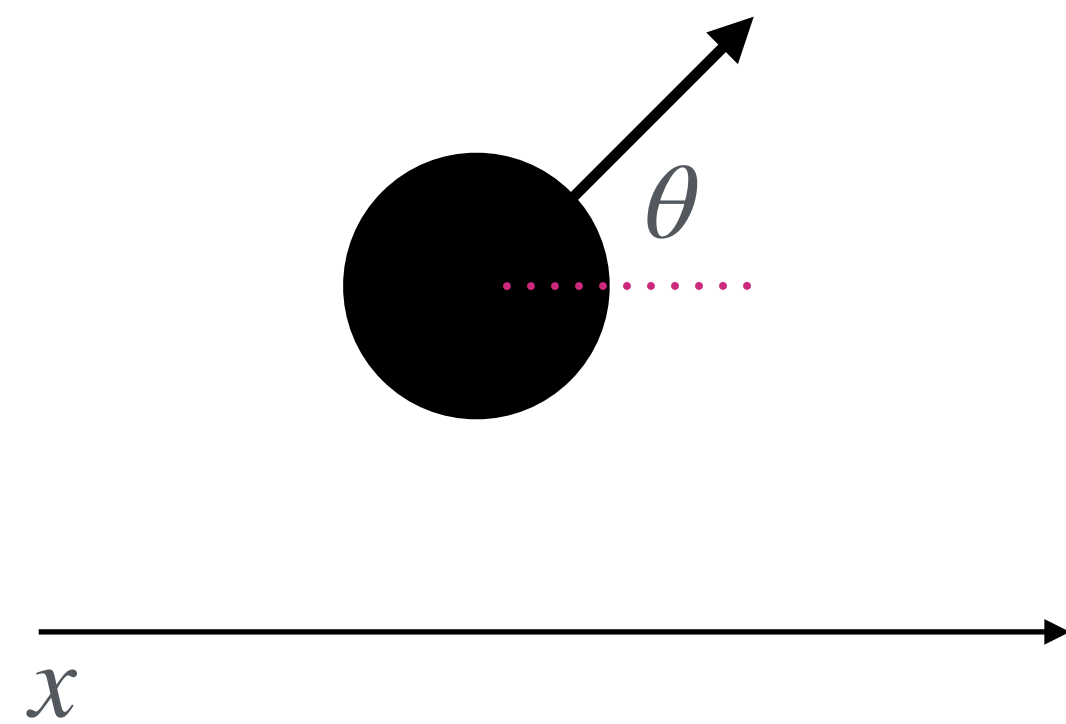


RT-X trained models are trained on 160k tasks



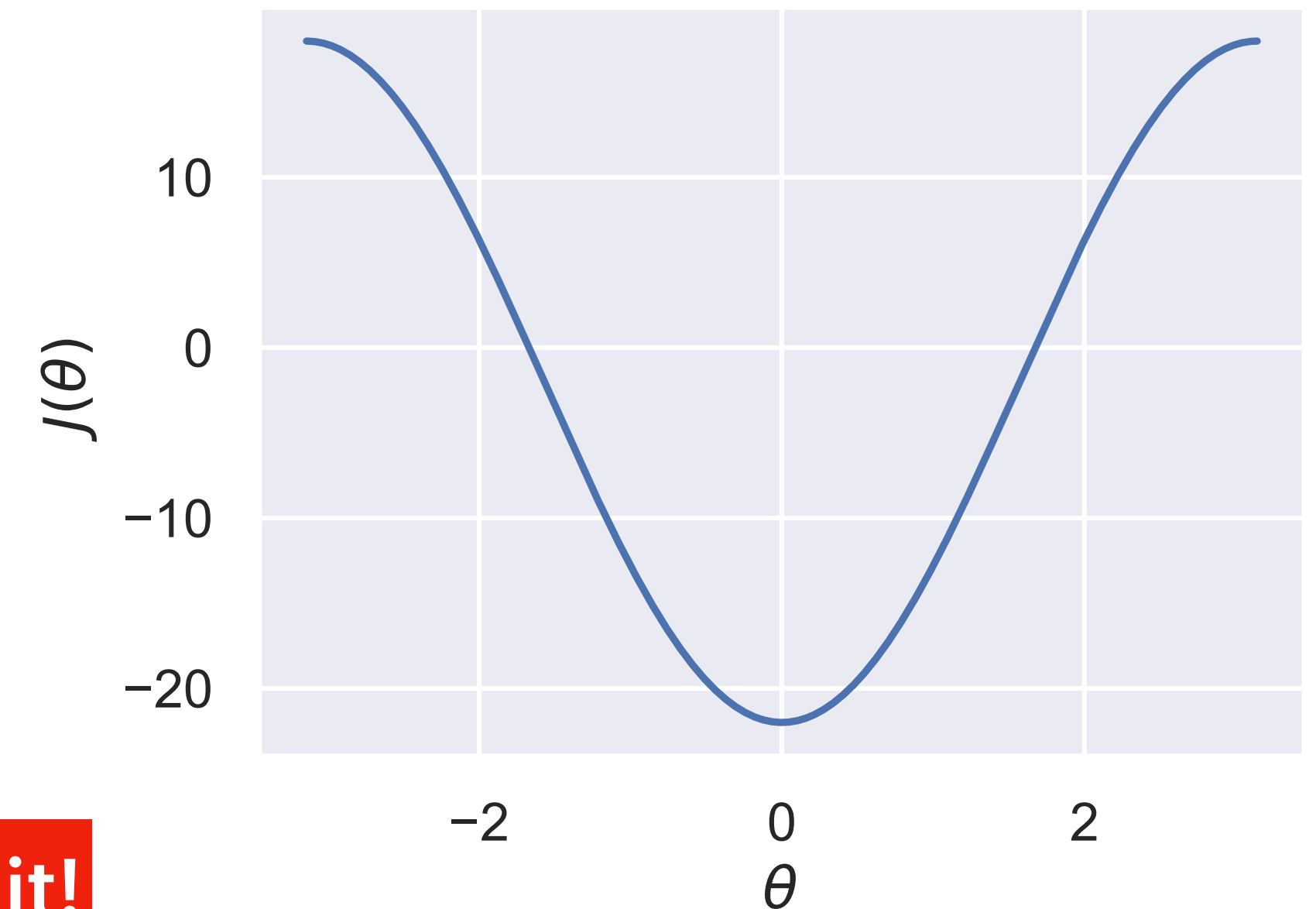
Why not SGD? Optimizing through contact

- What is so difficult?
- Let's optimize a simple task of a point mass thrown in free space for $t=2s$



$$J(\theta) = x_t = x_0 + v_x t + \frac{1}{2} a_x t^2$$

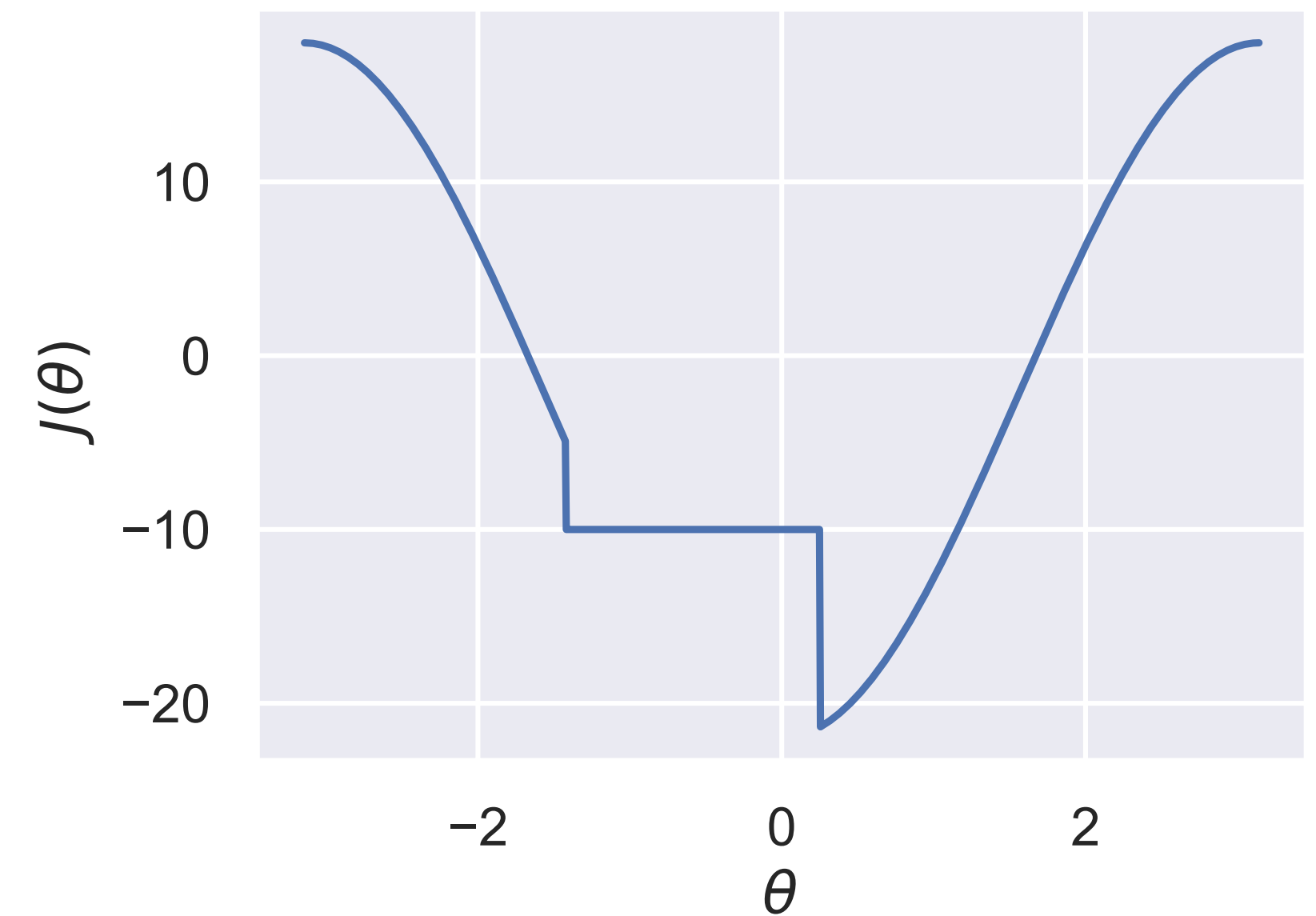
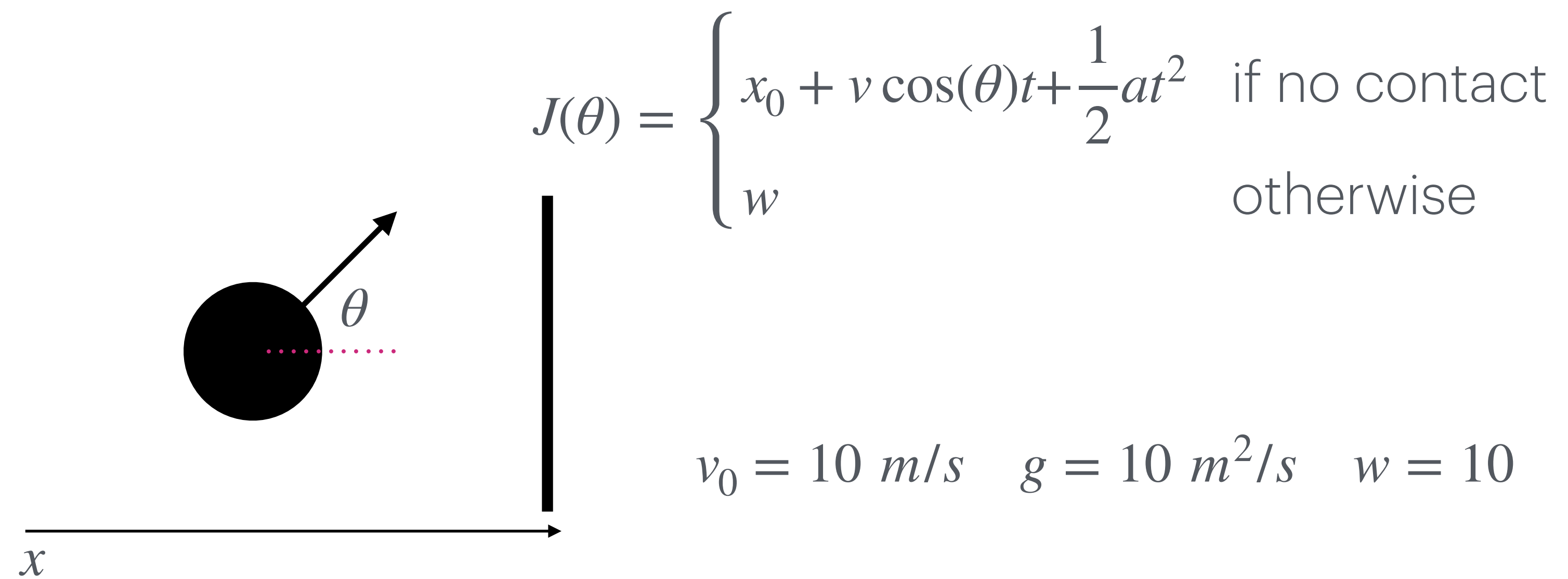
$$v_0 = 10 \text{ m/s} \quad g = 10 \text{ m}^2/\text{s}$$



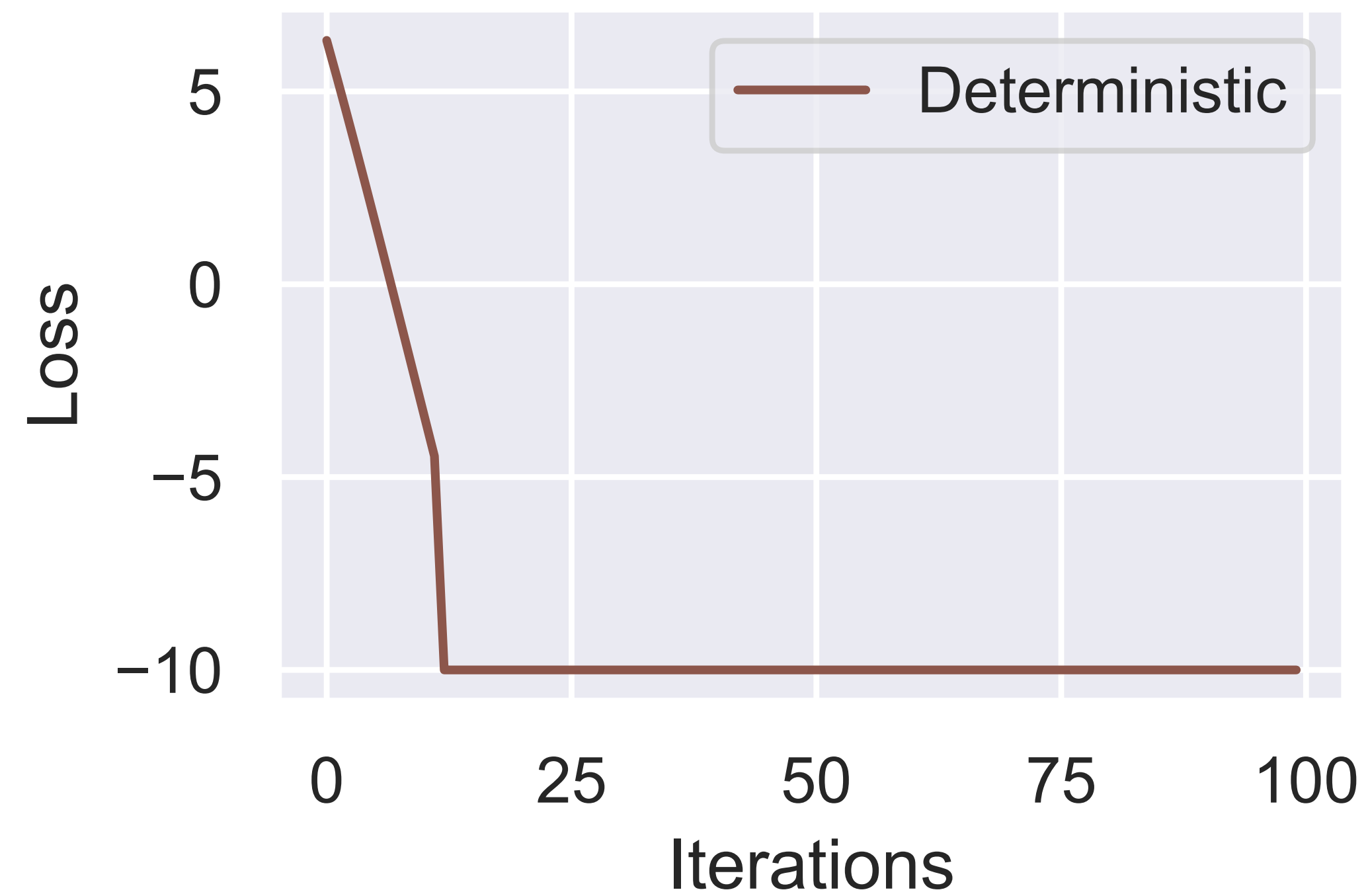
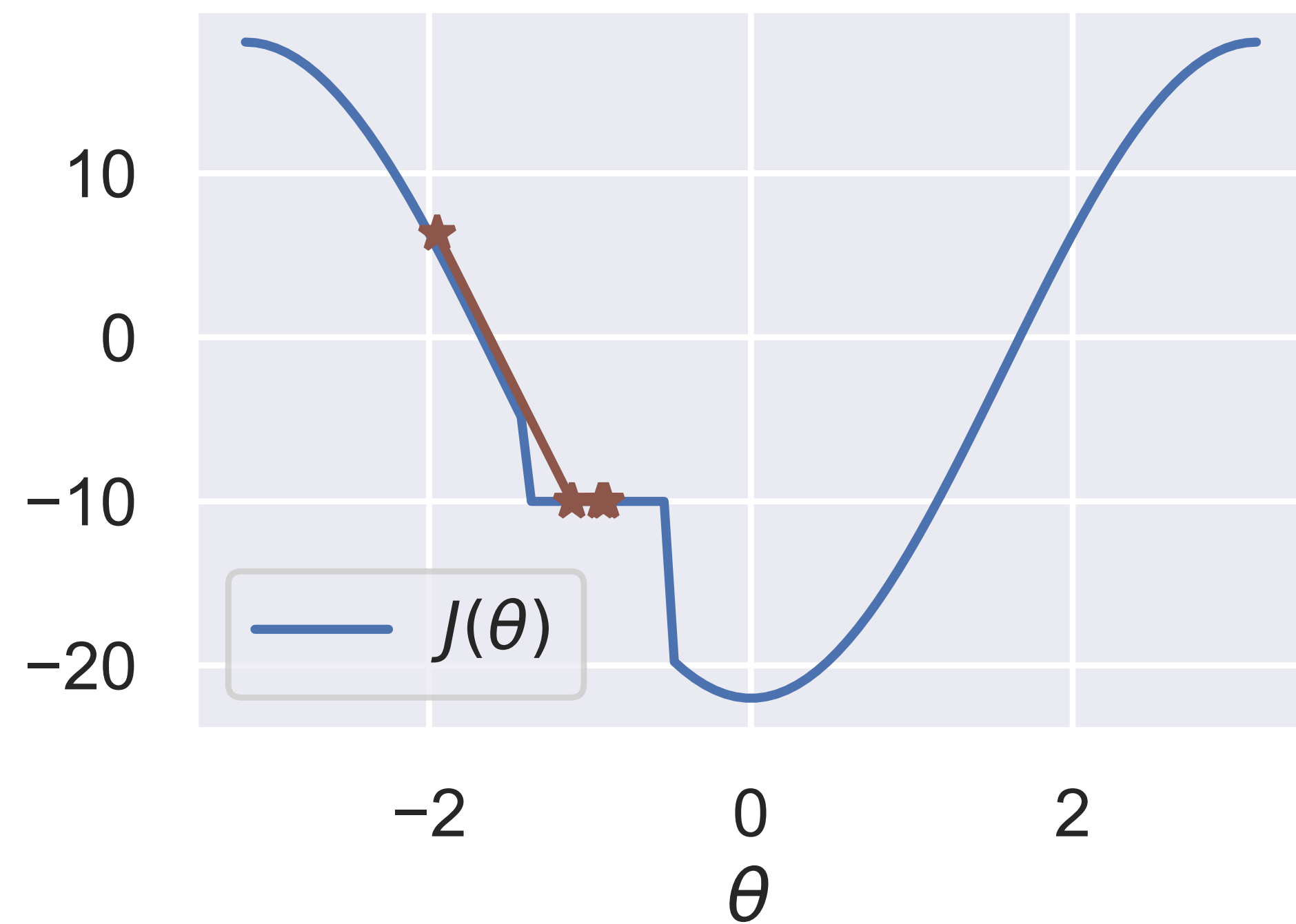
We can just throw GD at it!

Optimizing through contact

- Now add a wall! Assume that the ball sticks to it.

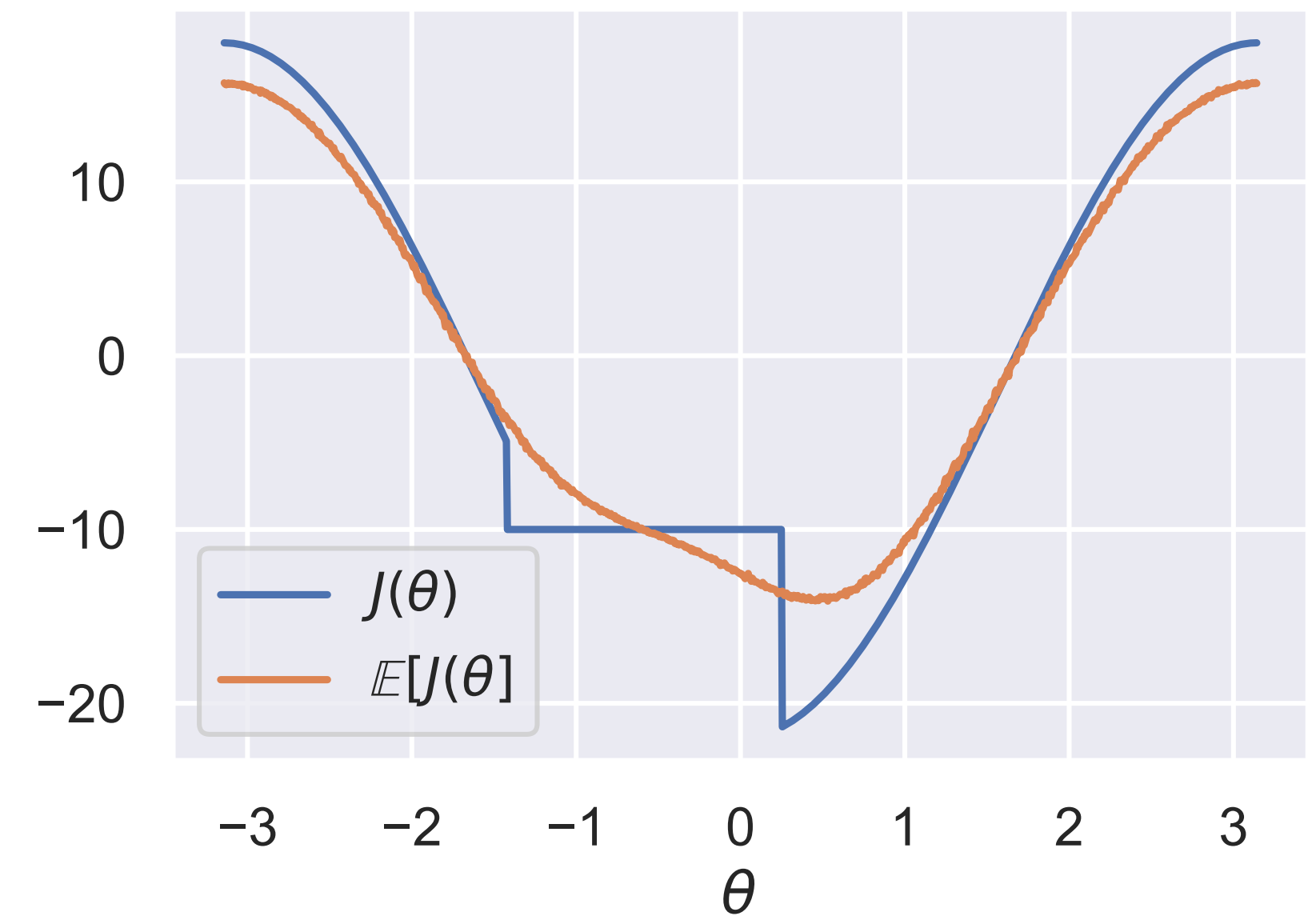
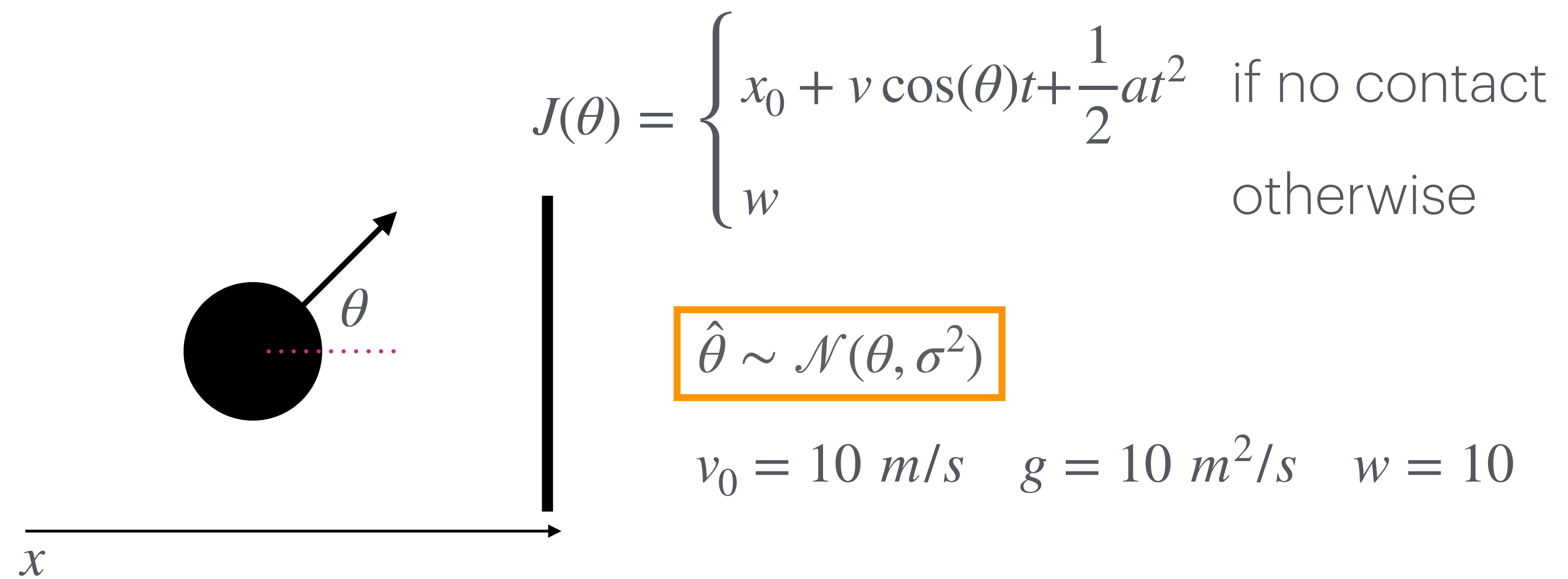


Optimizing through contact

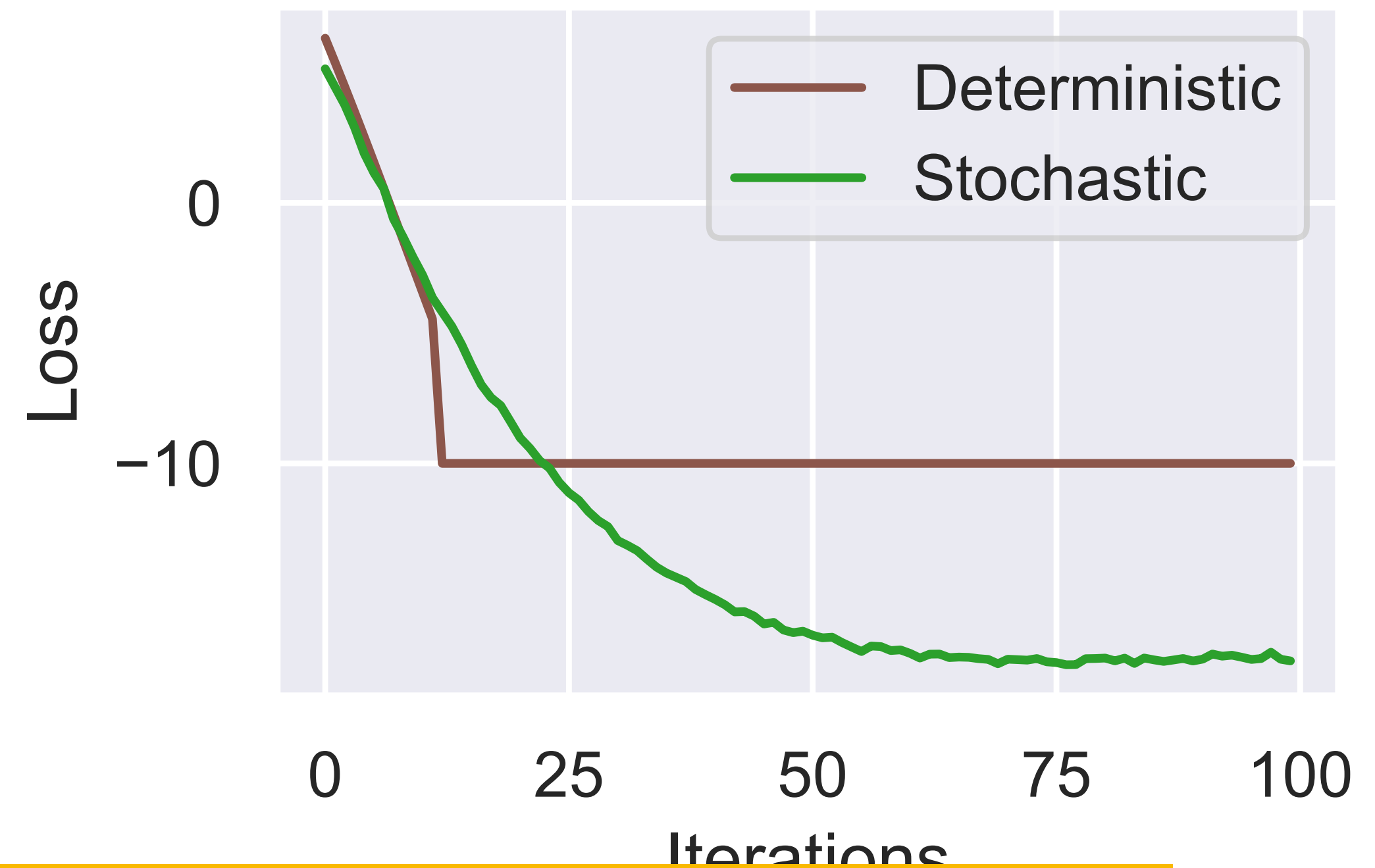
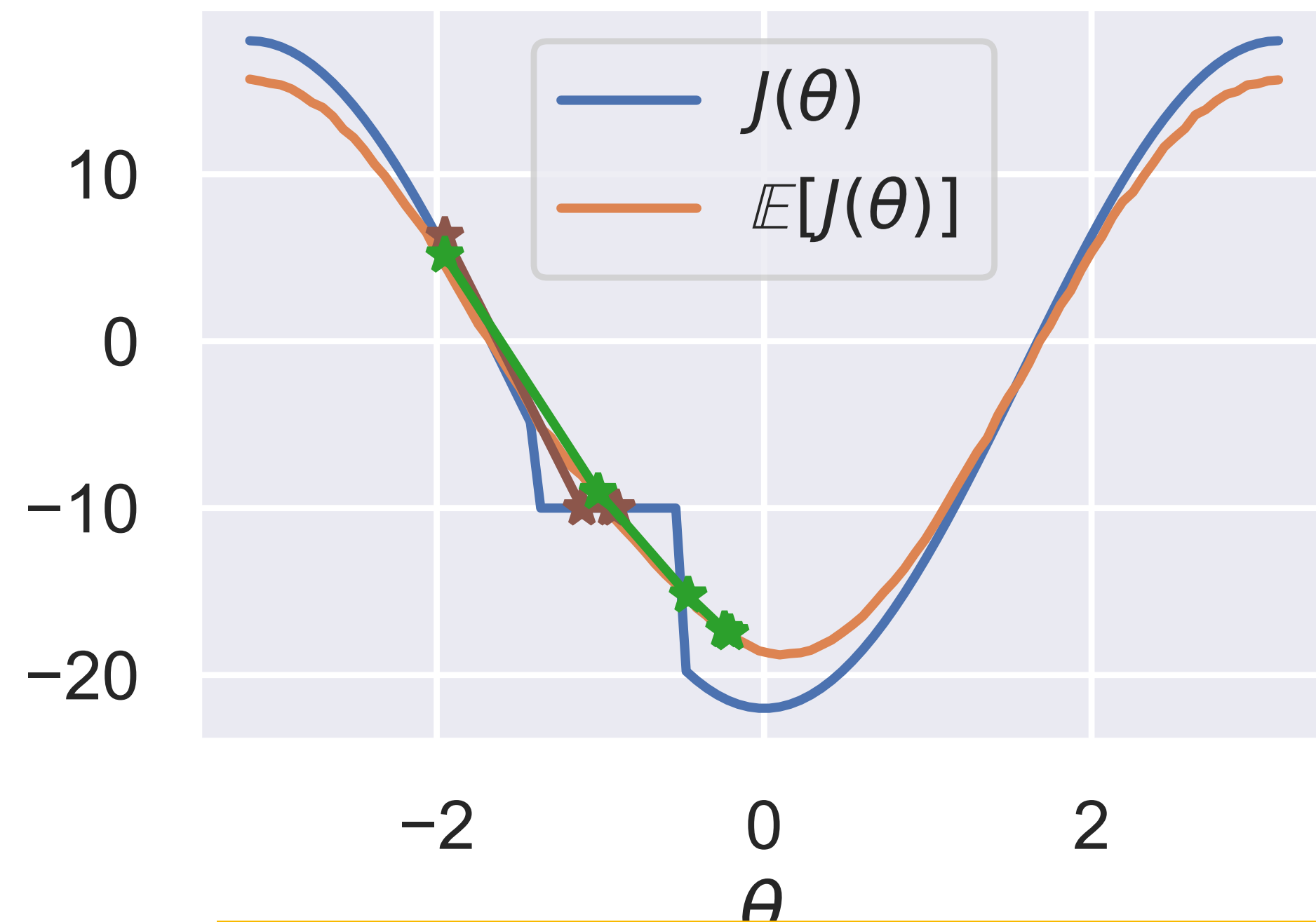


Optimizing through contact

- Make input (policy) noisy



Optimizing through contact



Takeaway: Stochastic optimization helps us solve non-smooth tasks

Takeaway: ZoG have been successful in RL as they can learn through contact

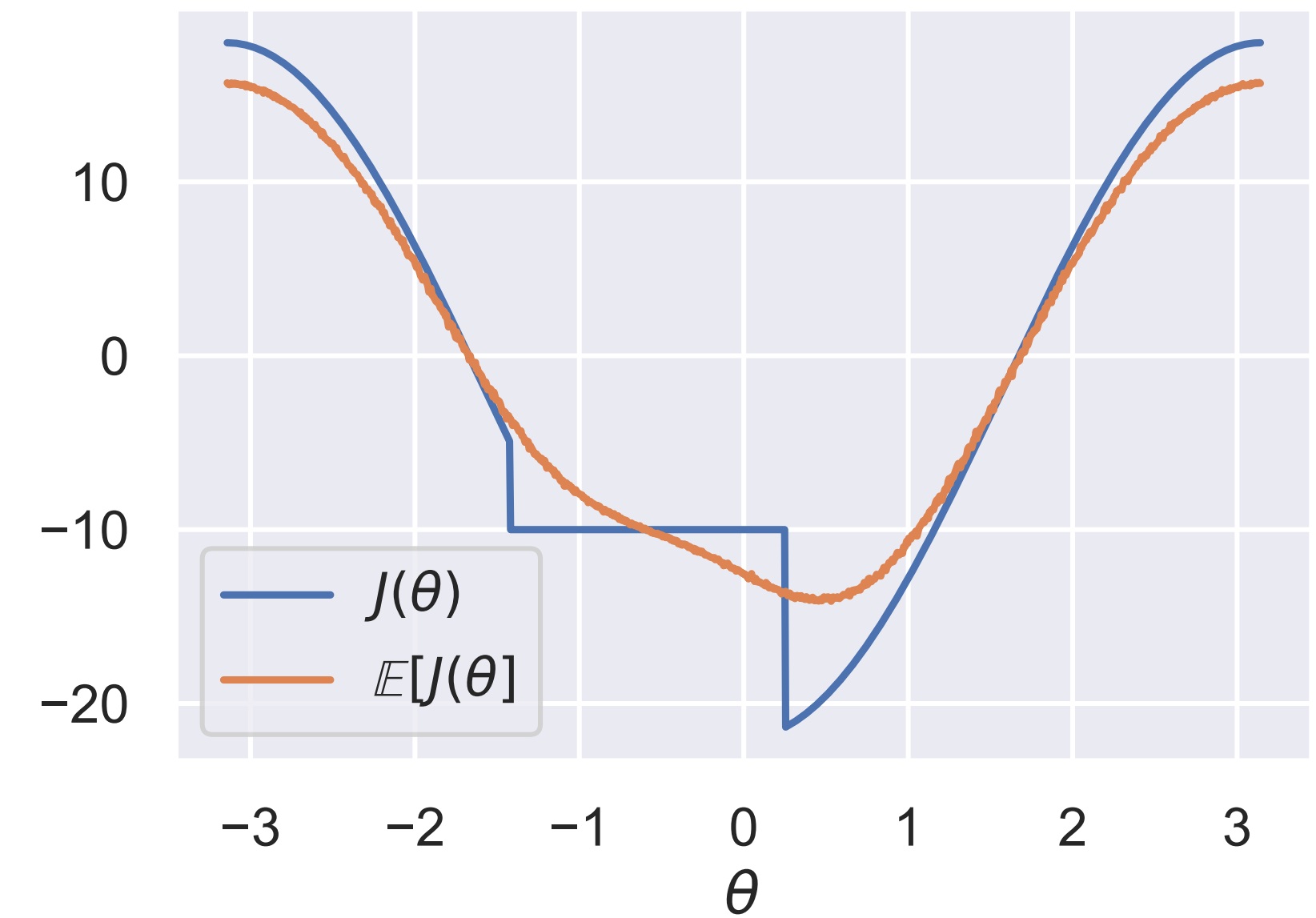
What can we do about it?

- ZoG RL by its formulation smooths out landscapes and enables optimization through contact.
- But it is SLOW and not scalable
- What else can we do?
 - Domain randomization (makes robust not good)
 - Initial state randomization
 -

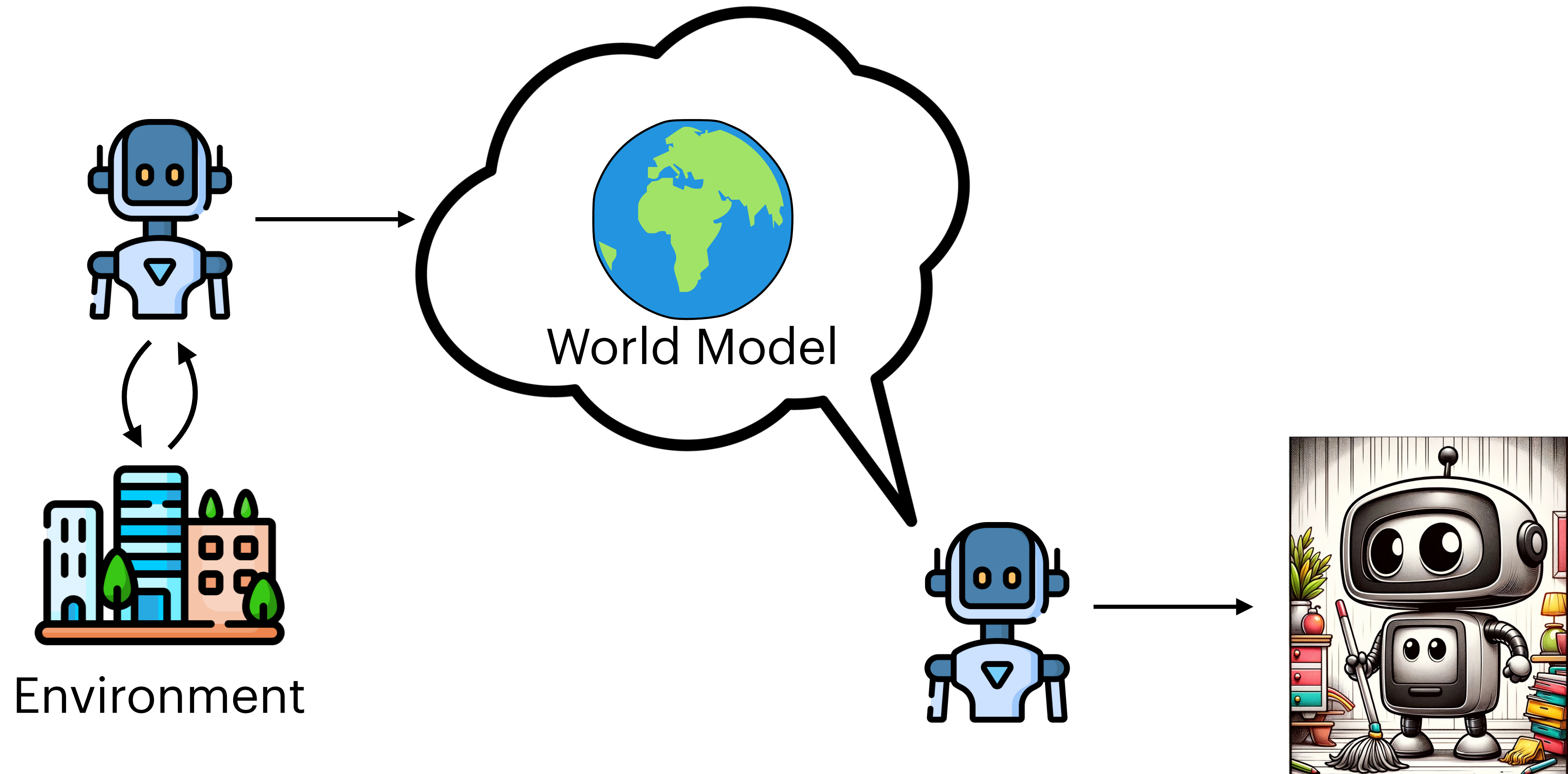
**All of these work but they
make data inefficient
approaches even more data
inefficient!**

Surrogate problem formulation

- What if instead of solving the real problem, we optimize a surrogate problem directly?
- Instead of $f(\theta)$, solve for $F_\phi(\theta)$
- We need these models to
 - Have minimal optimality gap $\|\hat{\theta} - \theta^*\|$
 - Be smoother than the original
$$\|\nabla F_\phi(\theta)\| \leq \|\nabla f(\theta)\|$$



World Model Framework



World models are smooth surrogates

- When regularized correctly, world models can act as smooth surrogates

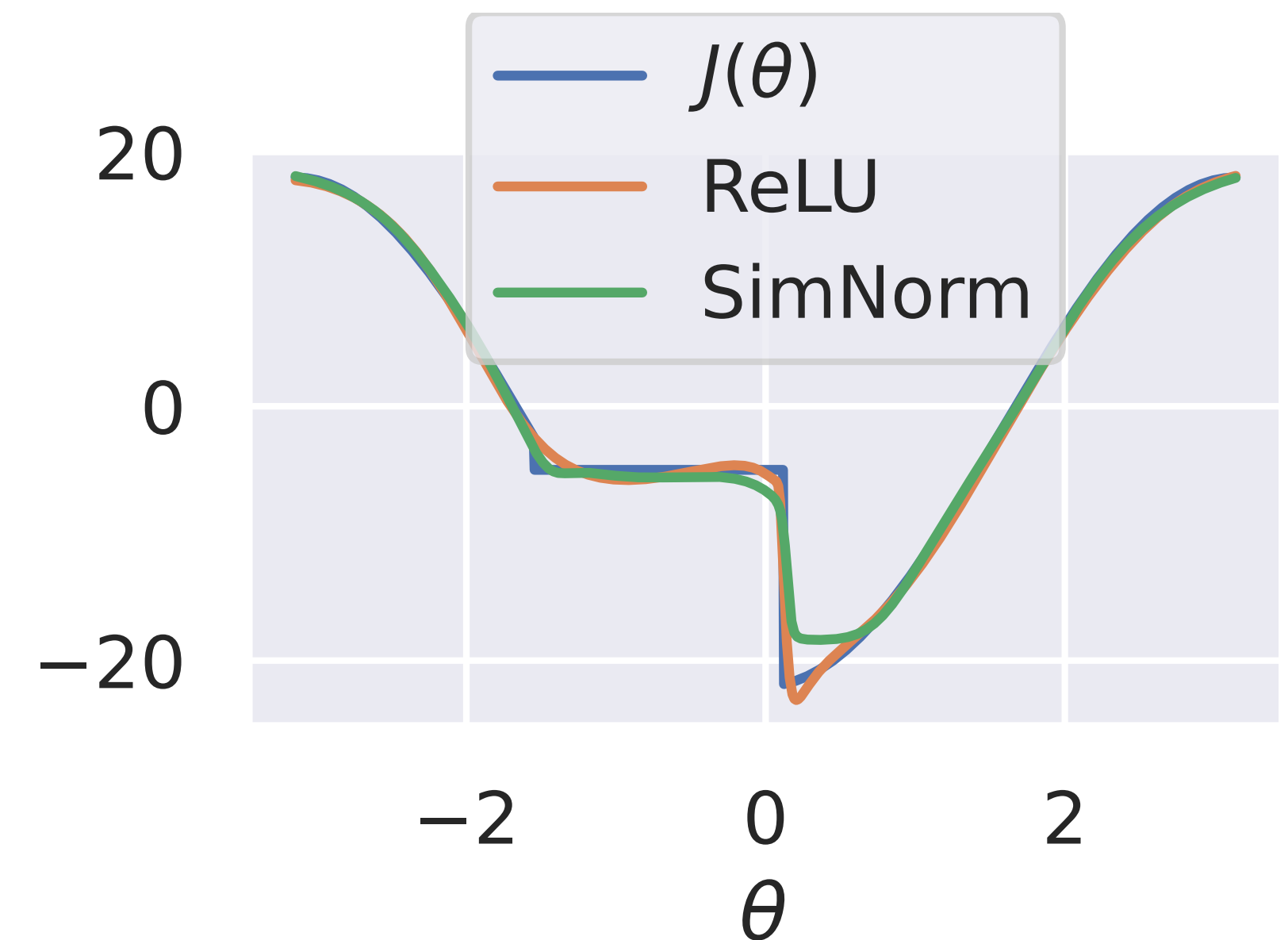
- No sampling required!

- Maps \mathbf{z} into V L -dimensional simplices

$$\text{SimNorm}(\mathbf{z}) := [\mathbf{g}_1, \dots, \mathbf{g}_L], \quad \mathbf{g}_i = \text{Softmax}(\mathbf{z}_{i:i+V})$$

- The key is not to make models accurate

- But to make them smooth
 - And have a low optimality gap



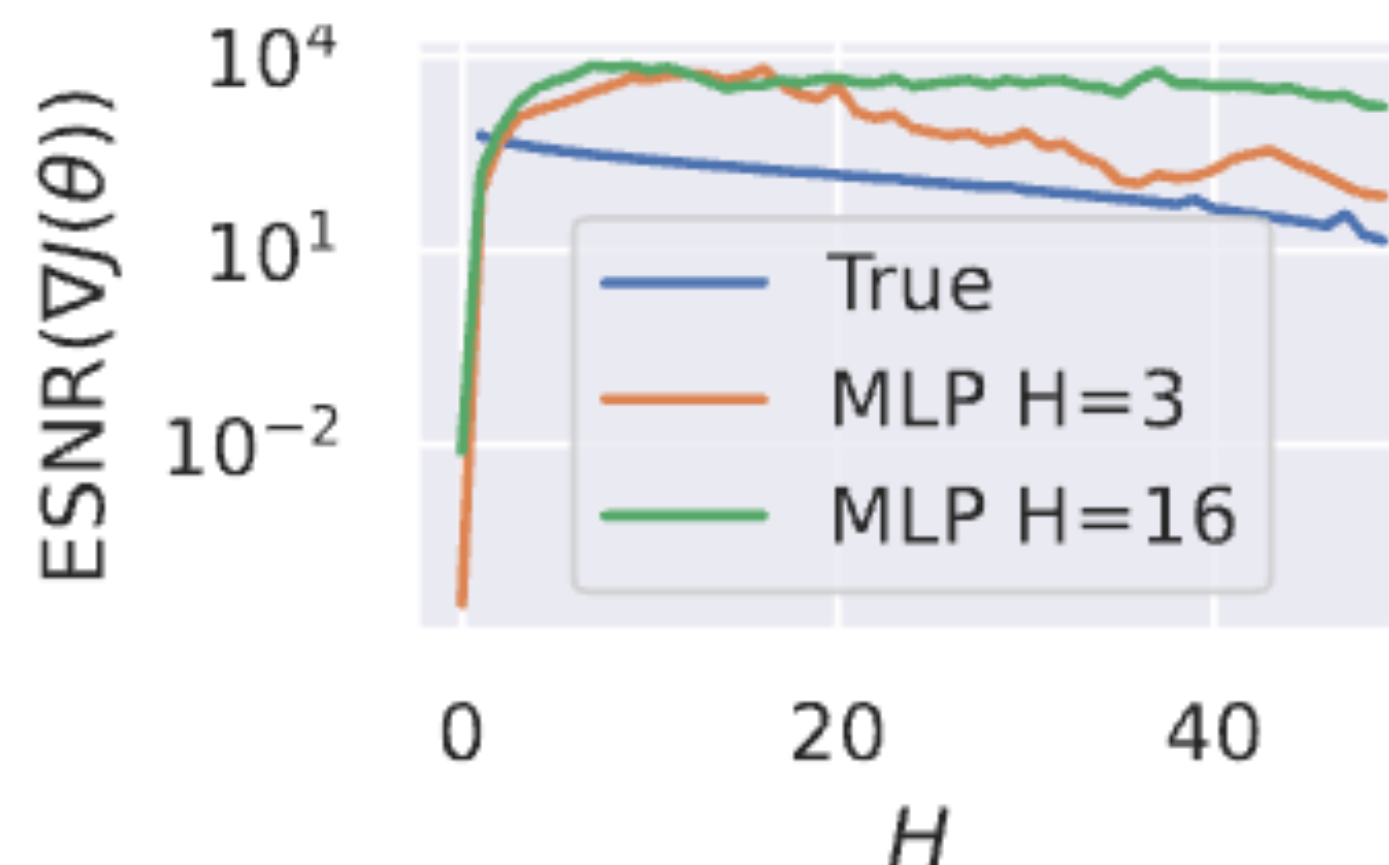
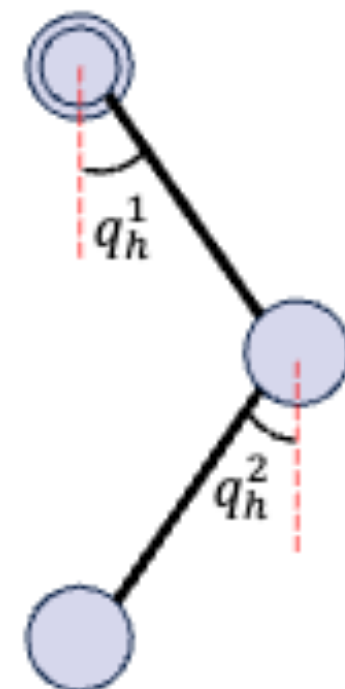
Model	Model error	Opt. gap
True	0.0	16.850
ReLU	0.707	16.046
SimNorm	1.131	3.473

(c) Model error and optimality gap.

World models provide stability

- When dealing with chaotic systems (e.g. double pendulum) gradients become less useful as we increase horizon H
- Policy gradient usefulness can be measured by variance or Expected Signal-to-Noise Ratio (ESNR)

$$\text{ESNR}(\nabla J(\boldsymbol{\theta})) = \mathbb{E} \left[\frac{\sum \mathbb{E} [\nabla^{[1]} J(\boldsymbol{\theta})]^2}{\sum \text{Var} [\nabla^{[1]} J(\boldsymbol{\theta})]} \right]$$



World models are good surrogates

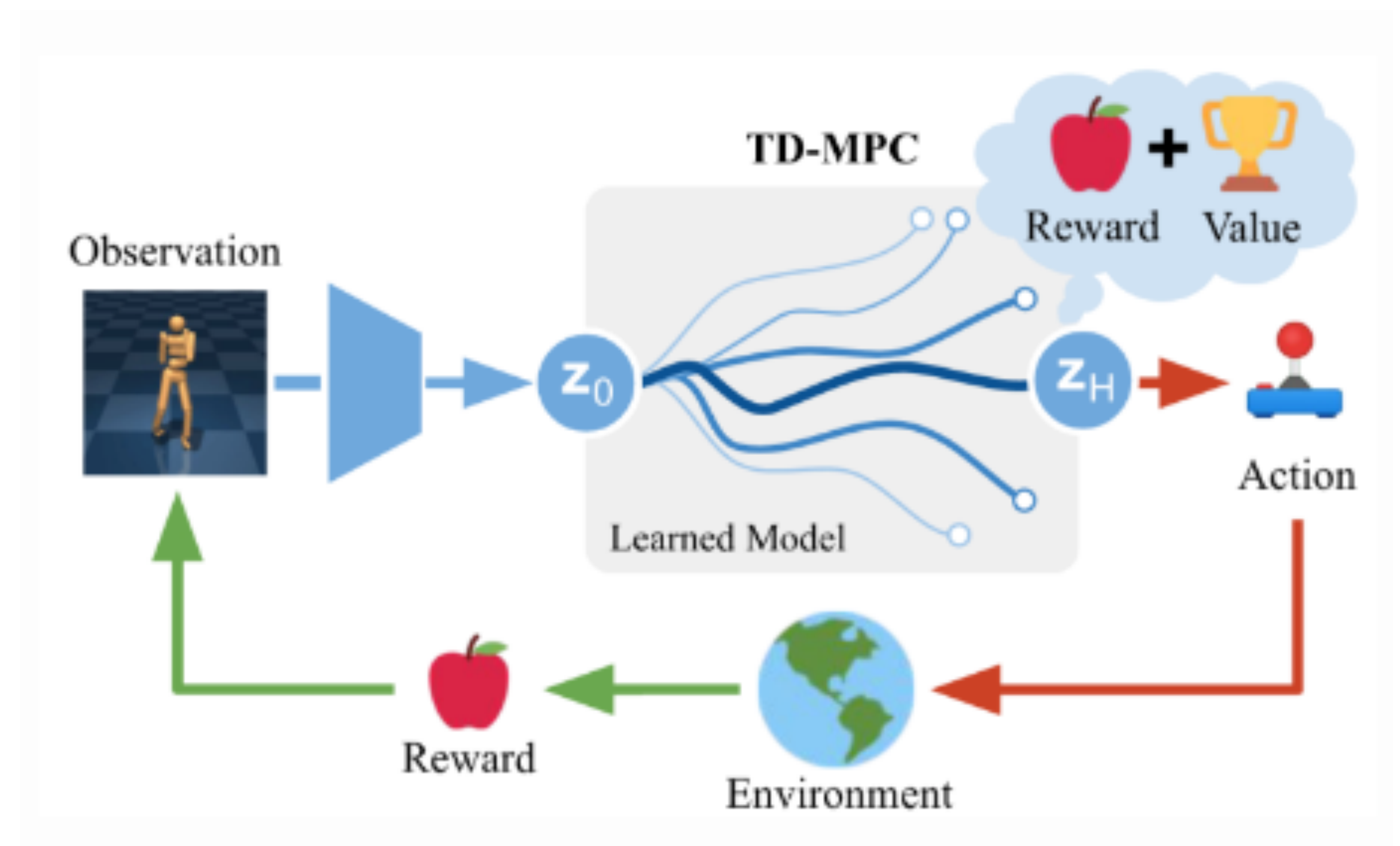
- Provide smooth problem landscapes
- Provide stable gradients through long trajectories

Takeaway: it is better to optimize over world models rather than true dynamics

TDMPC2

A scalable multi-task world model approach

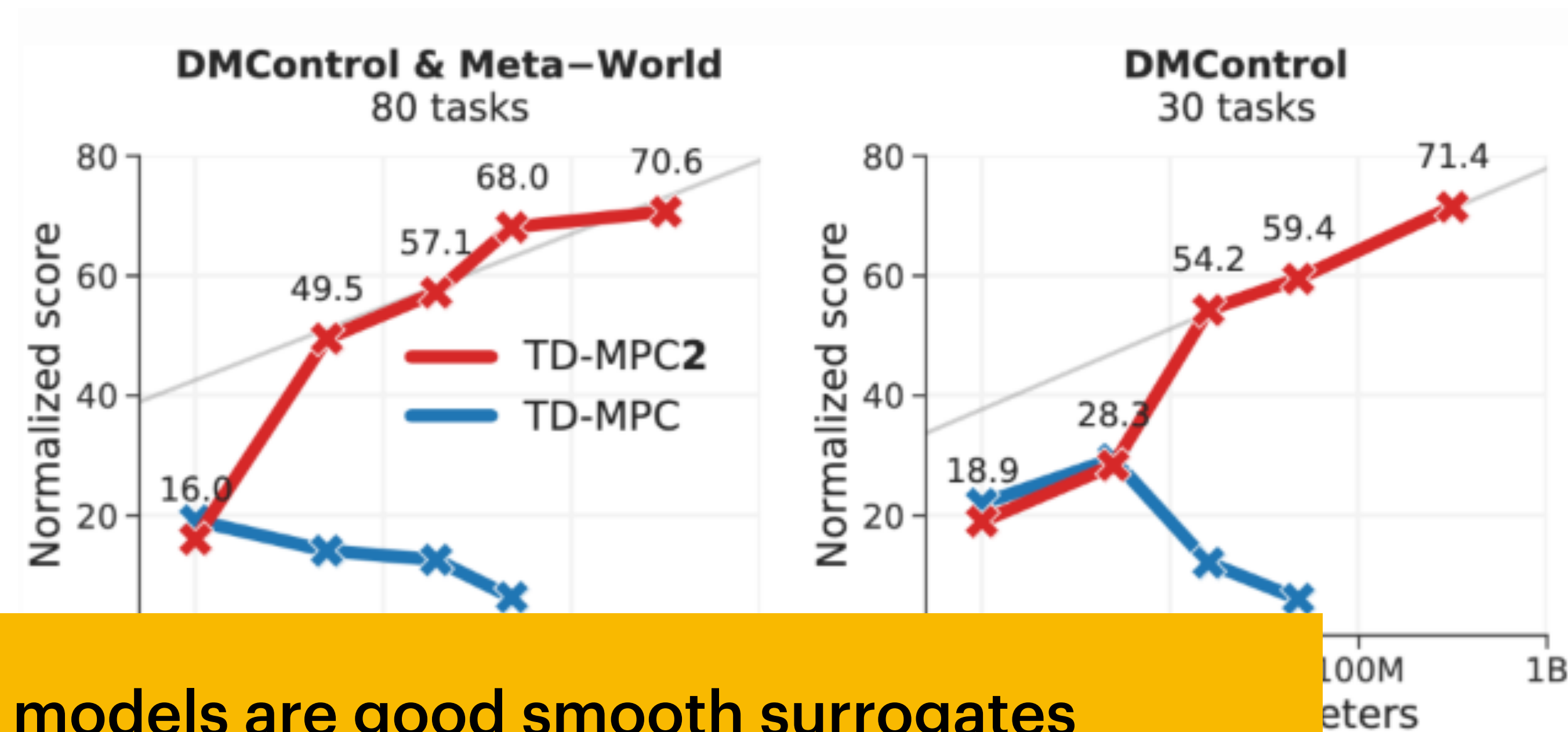
- Model-based RL approach
- Task-orientated latent dynamics model
- Learns by relevant reward, not by input reconstruction
- SAC actor-critic policy combined with online planning



TDMPC2

A scalable multi-task world model approach

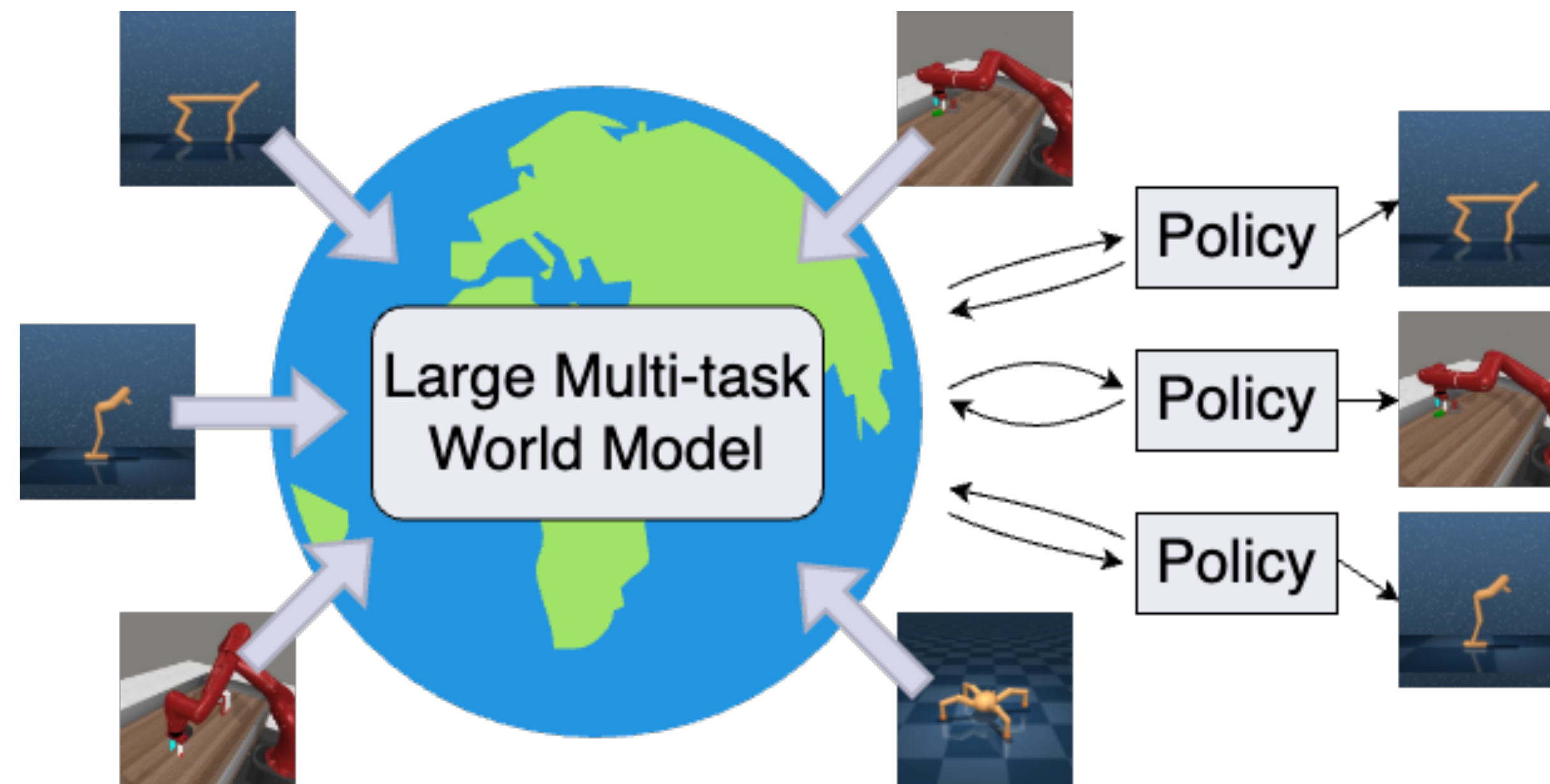
- First multi-task RL policy that scales to 80 different tasks
- Across MetaWorld and DMC
- Relies mostly on online planning



Turns out that TDMPC2 world models are good smooth surrogates

But TDMPC2 chooses to use ZoG, what if we use FoG?

PWM: Policy Learning with Large World Models



1. Regularized large models enable efficient policy learning
2. Use First-order optimization to train policies in <10m per task

PWM: Policy Learning with Large World Models

- TDMPC2 world models $(E_\phi(s, e), F_\phi(s, a, e), R_\phi(s, a, e))$
 - H=16 and $\gamma=0.99$
- Model-free critic trained with TD(λ)
- Actor trained with FoG

$$\mathcal{L}_\pi(\theta) := \mathbb{E}_{\substack{s_1 \sim \rho(\cdot) \\ a_h \sim \pi_\theta(\cdot | z_h)}} \left[\sum_{h=1}^{H-1} \gamma^h R_\phi(z_h, a_h) + \gamma^H V_\psi(z_H) \right]$$

Algorithm 1: PWM: Policy optimization through World Model

Given: Multi-task dataset \mathcal{B}

Given: γ : discount rate

Given: $\alpha_\theta, \alpha_\psi, \alpha_\phi$: learning rates

Initialize learnable parameters θ, ψ, ϕ

▷ Pre-train world model once

for N epochs **do**

$s_{1:H}, a_{1:H}, r_{1:H}, e \sim \mathcal{B}$

$\phi \leftarrow \phi + \alpha_\phi \mathcal{L}_{wm}(\phi)$

▷ Eq. 10

end

▷ Train policy on task embedding e

for M epochs **do**

$s_1 \sim \mathcal{B}$

$z_1 = E_\phi(s_1, e)$

for $h=[1, \dots, H]$ **do**

▷ Rollout

$a_h \sim \pi_\theta(\cdot | z_h)$

$r_h = R_\phi(z_h, a_h, e)$

$z_{h+1} = F_\phi(z_h, a_h, e)$

end

$\theta \leftarrow \theta + \alpha_\theta \mathcal{L}_\pi(\theta)$

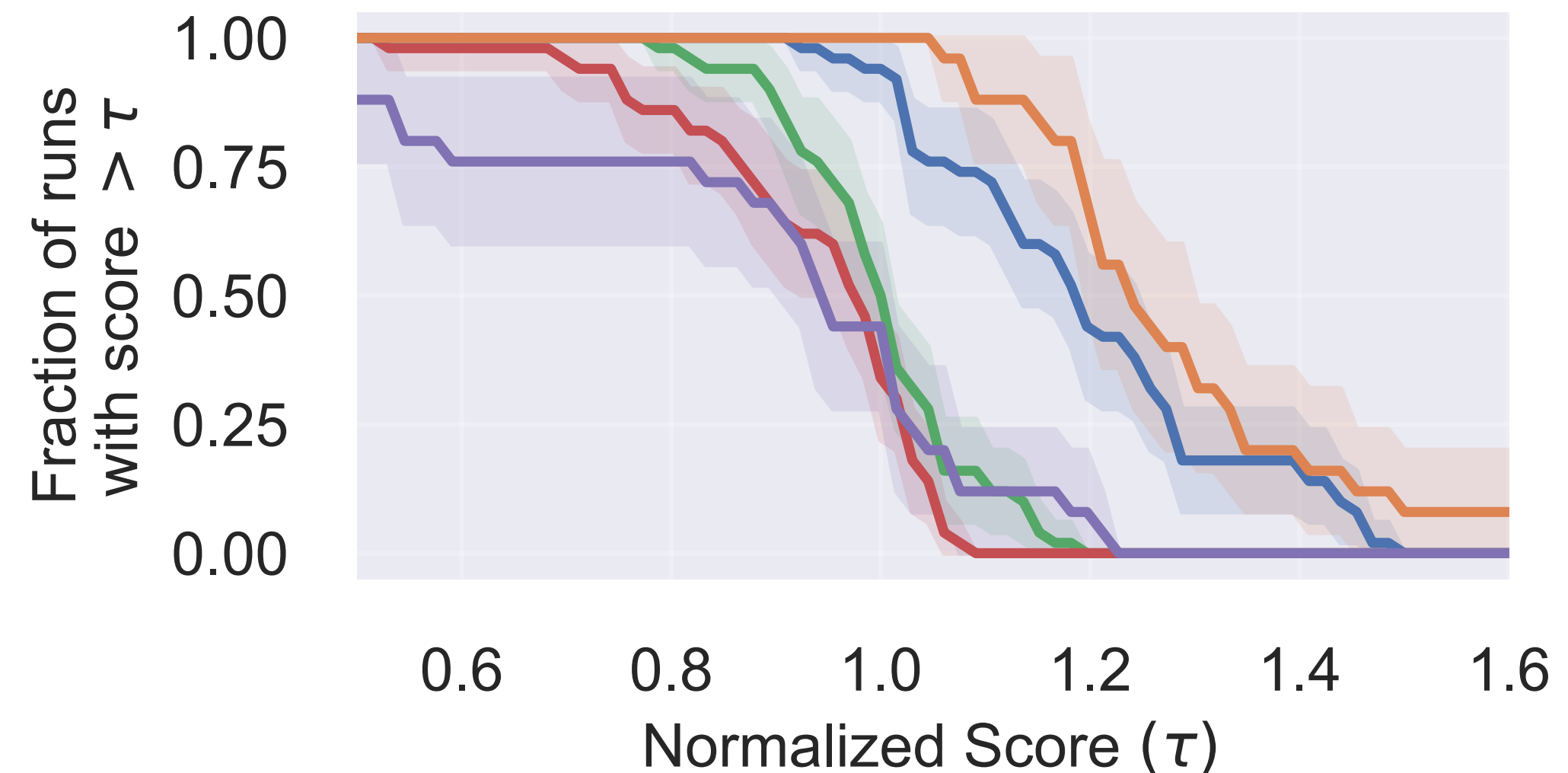
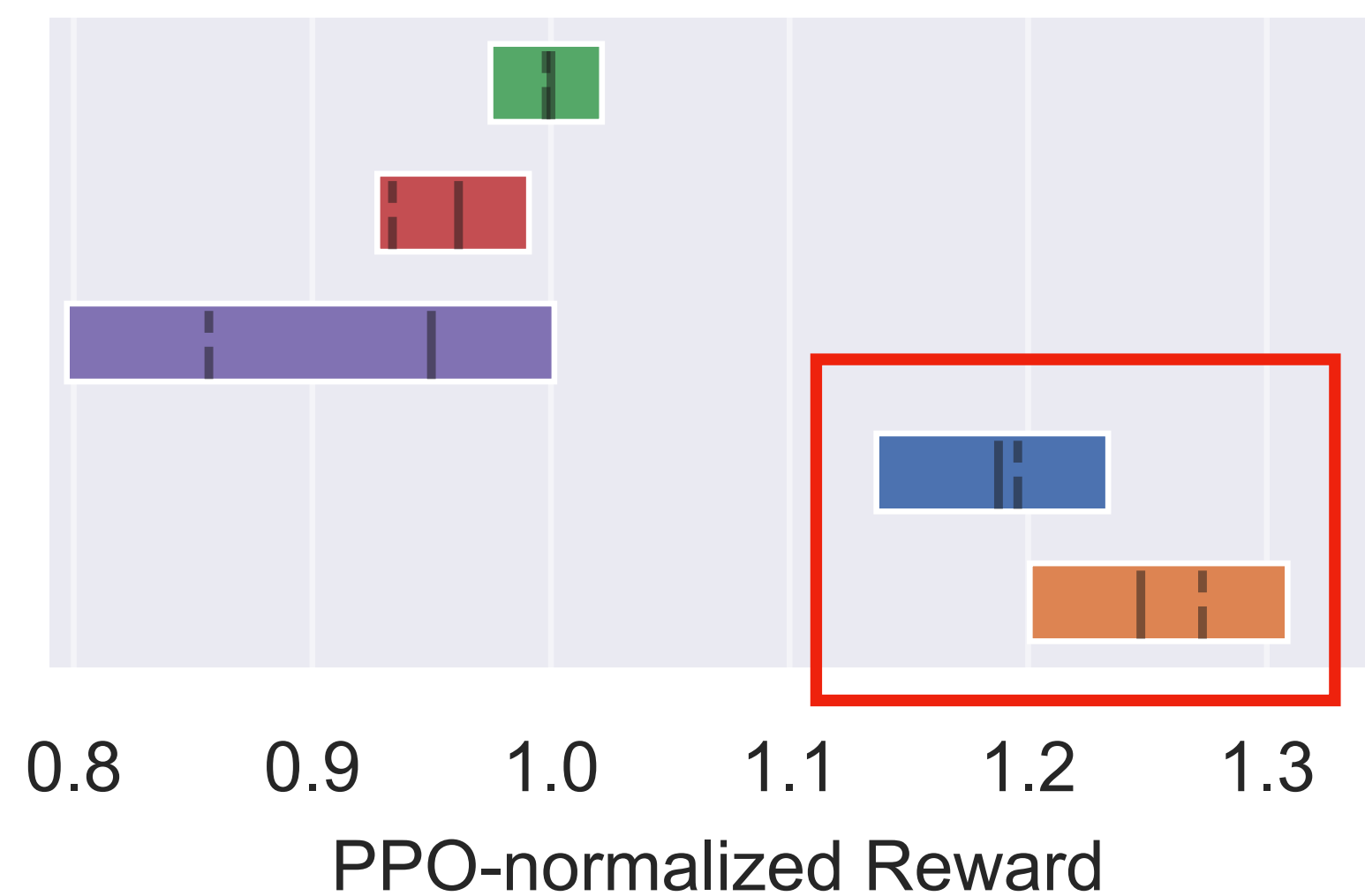
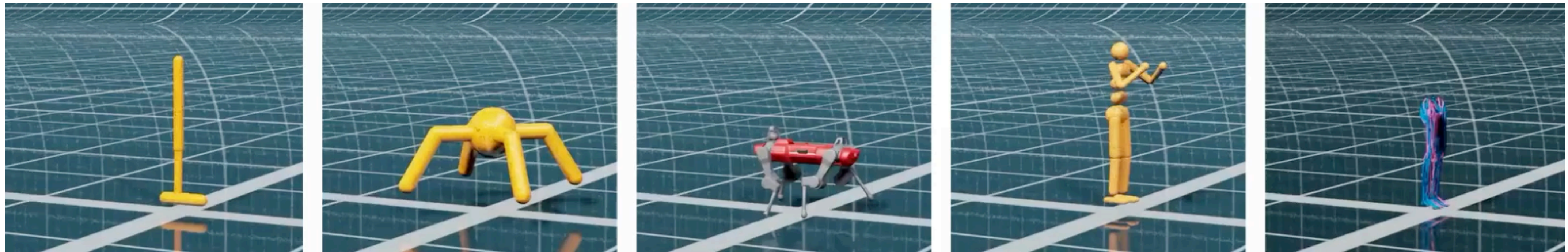
▷ Eq. 7-9

$\psi \leftarrow \psi + \alpha_\psi \mathcal{L}_V(\psi)$

▷ Eq. 6

end

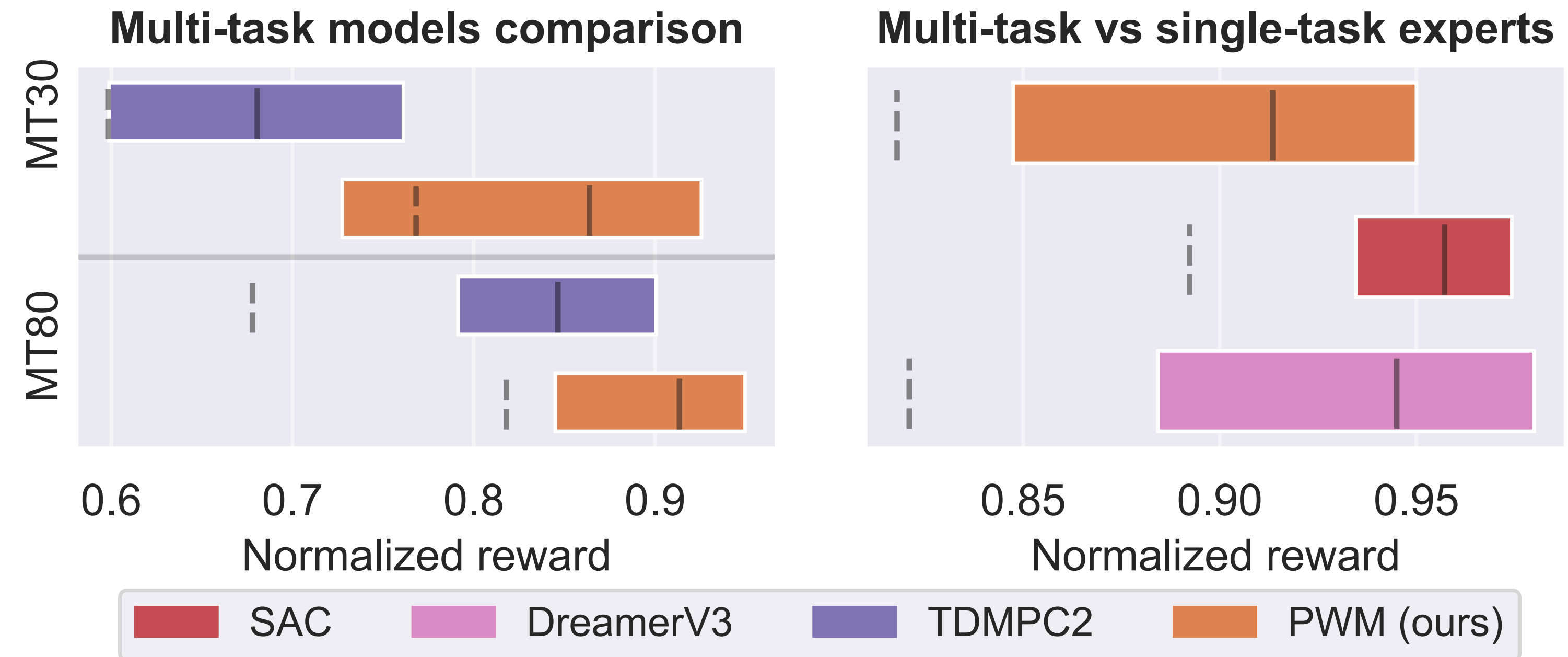
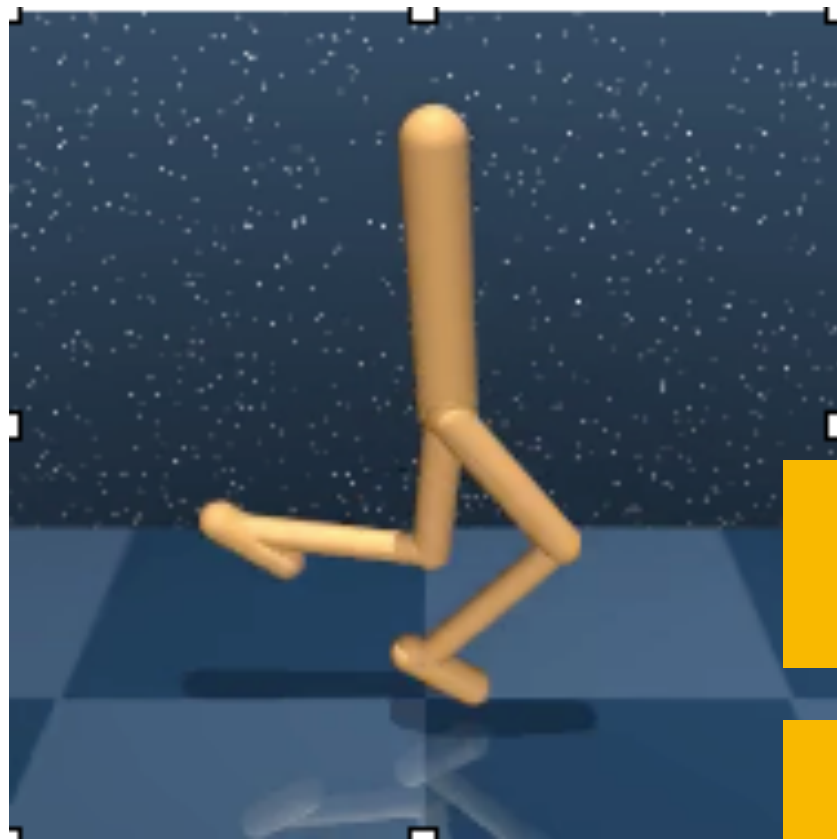
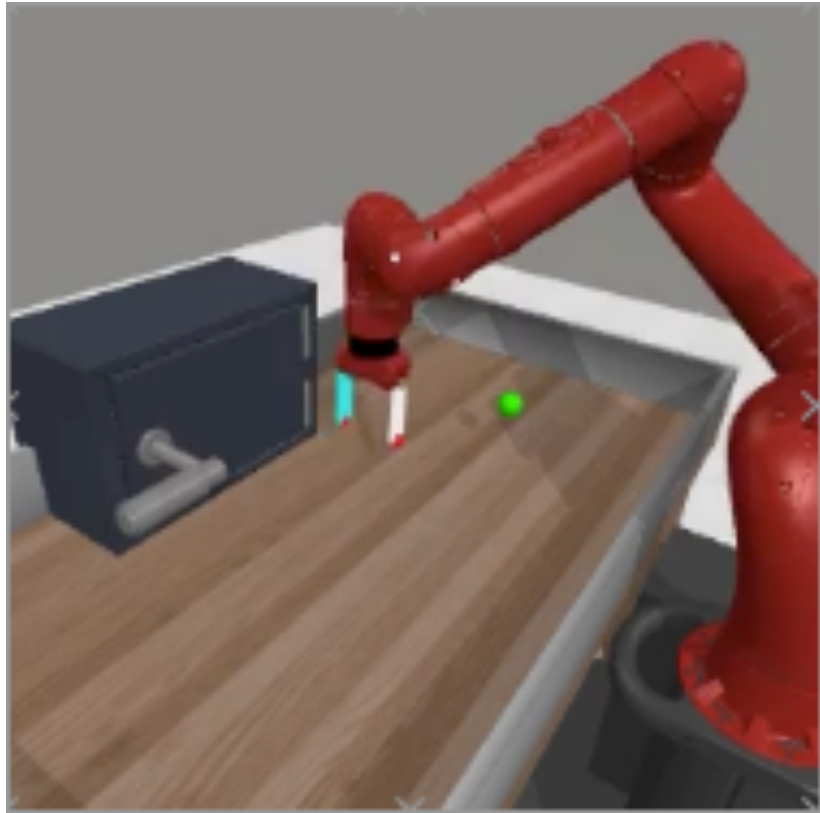
High-dimensional single-task



PPO SAC SHAC TDMPC2 PWM (ours)

Takeaway: optimizing over surrogate models obtains better policies than ground truth!

Multi-task experiments

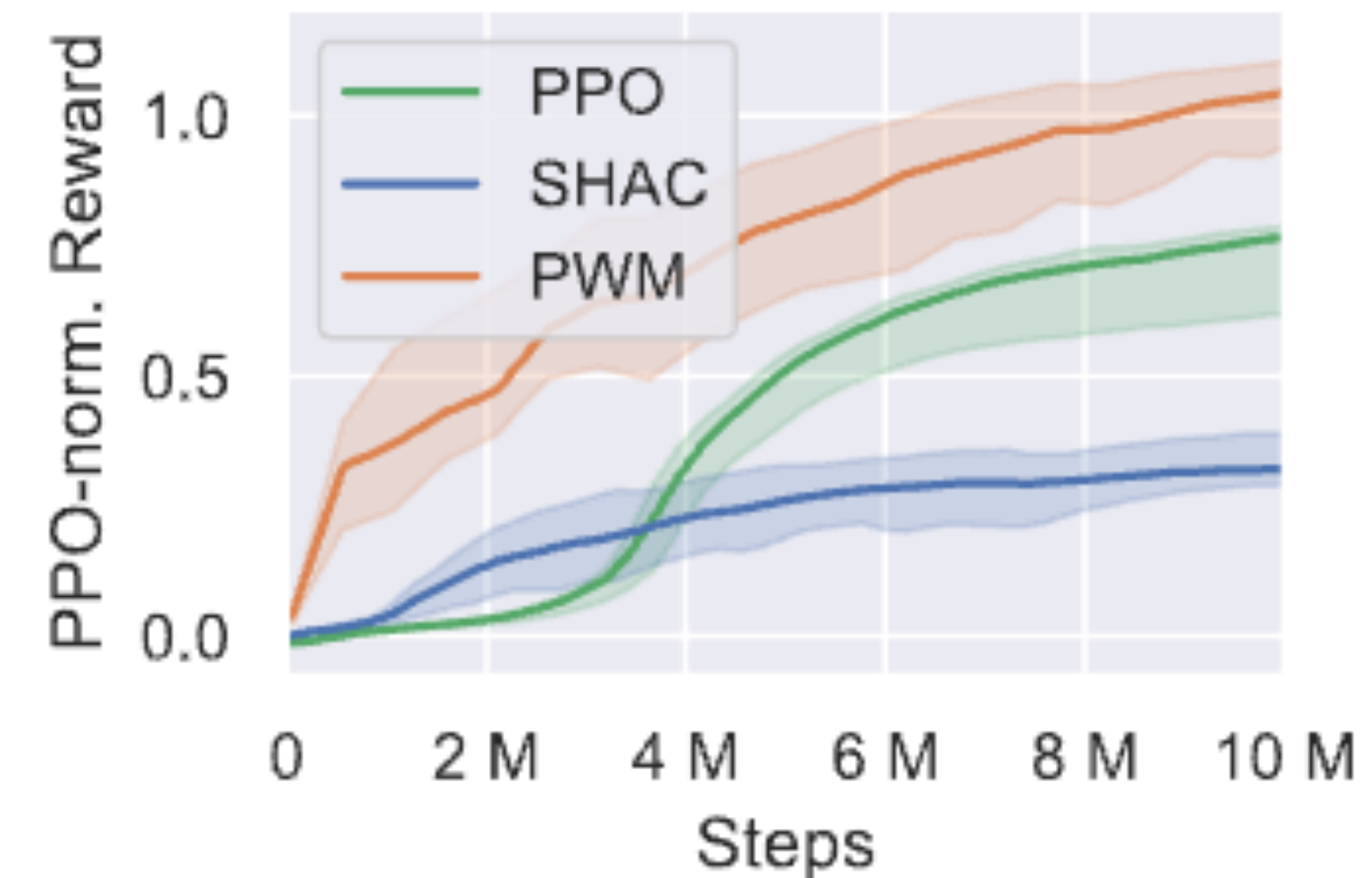


Beats TDMPC2 without the need for online planning -> more scalable

Matches single-task experts without any online interaction

Stiffness ablation

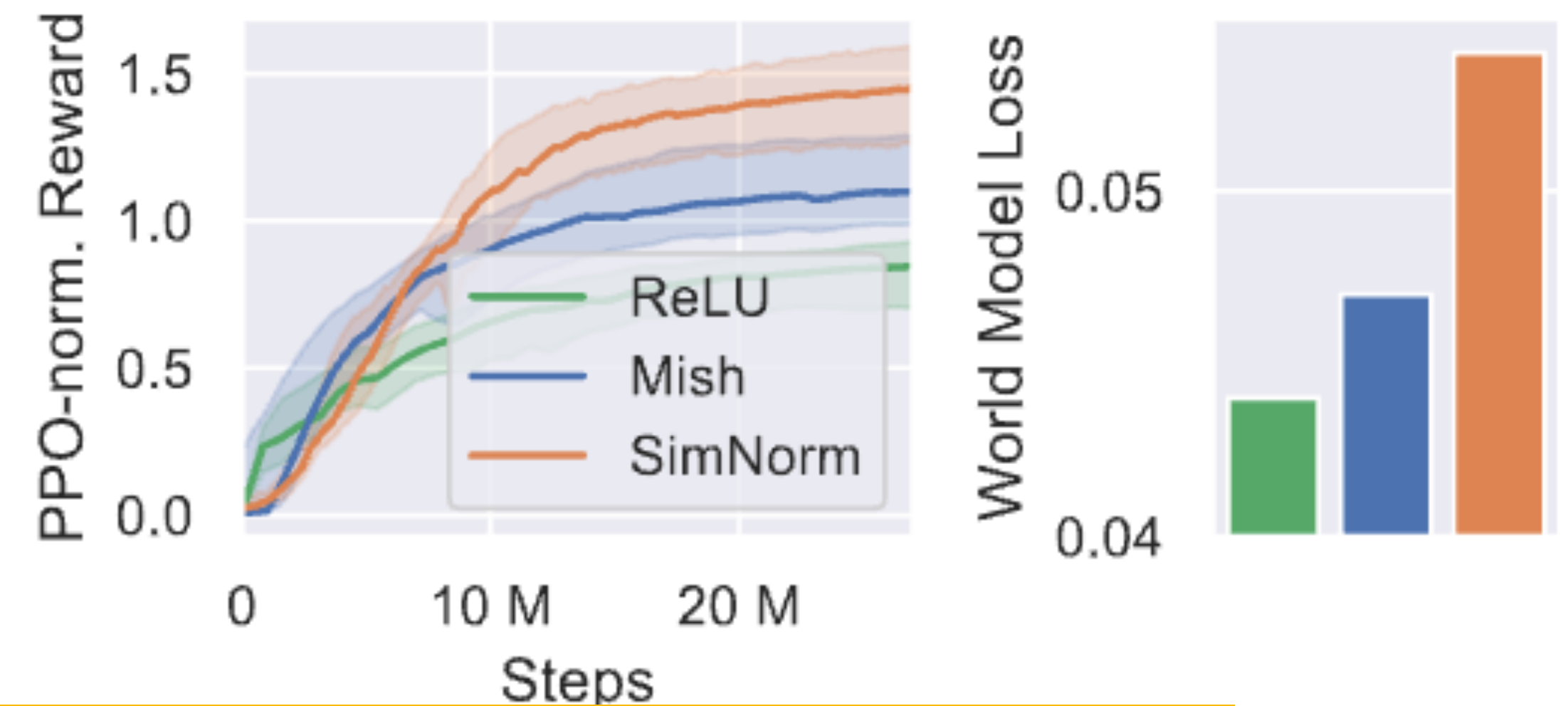
- If we increase contact stiffness, SHAC/FoG method performance decreases
- PWM sustains the same performance
- Hopper task



PWM is resilient to stiff dynamics

World Model Ablation

- More accurate models do not translate to better policy
- Actually the opposite



We should build world models for policy learning, not accuracy

Sample efficiency

1. Train world models for X time steps
2. Then train policy for 50k gradient steps



PWM is a more sample efficient policy learning technique but requires better trained world models.

(b) World model vs policy sample efficiency.

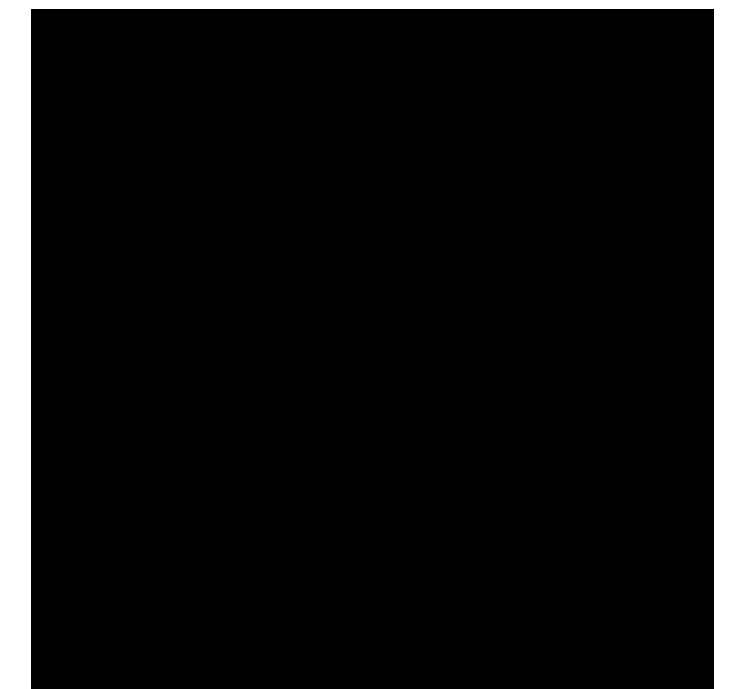
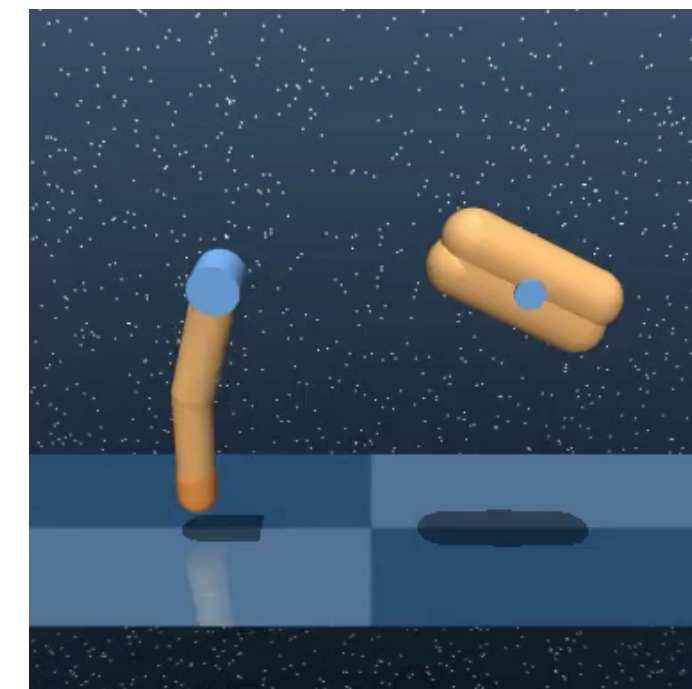
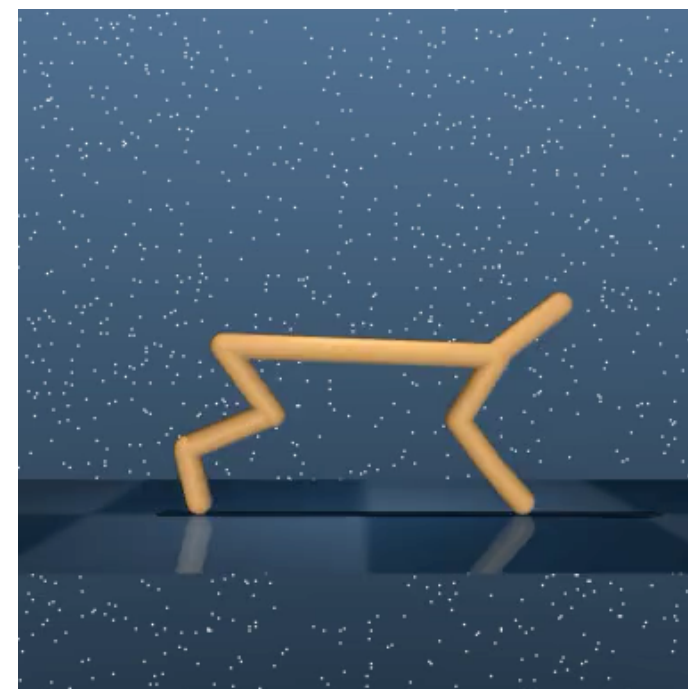
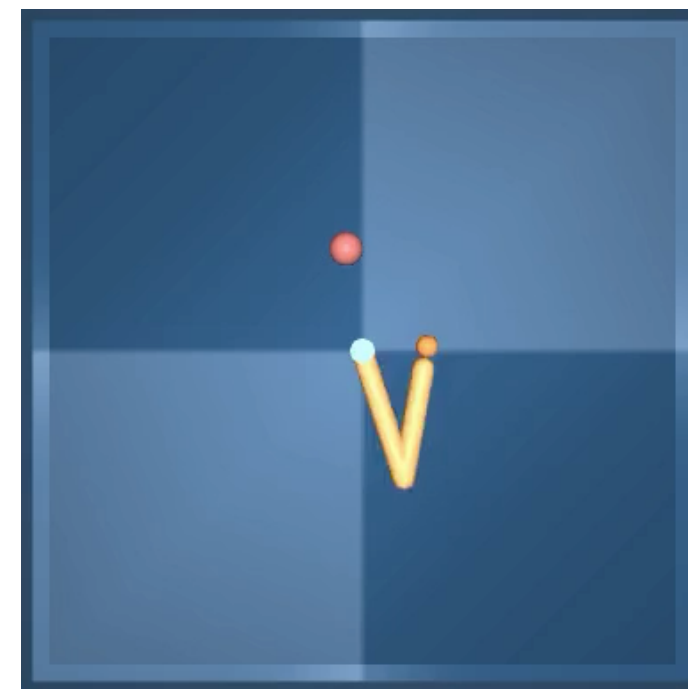
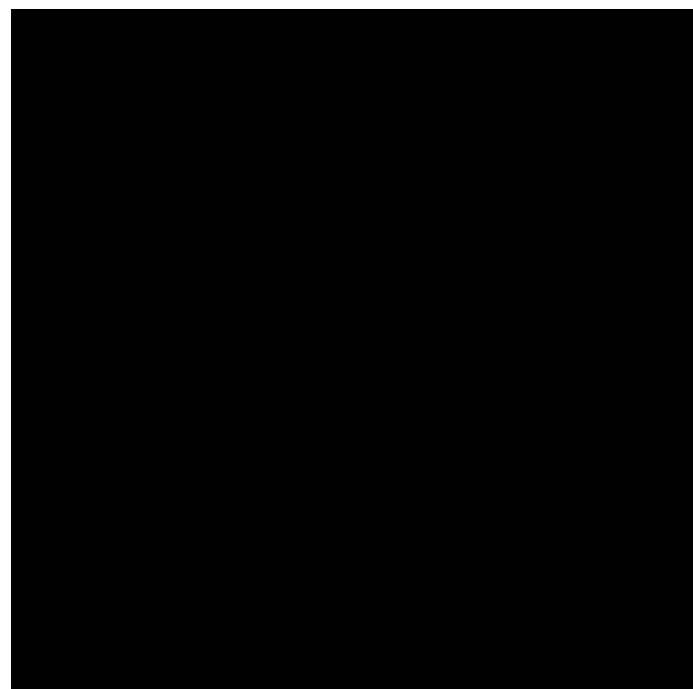
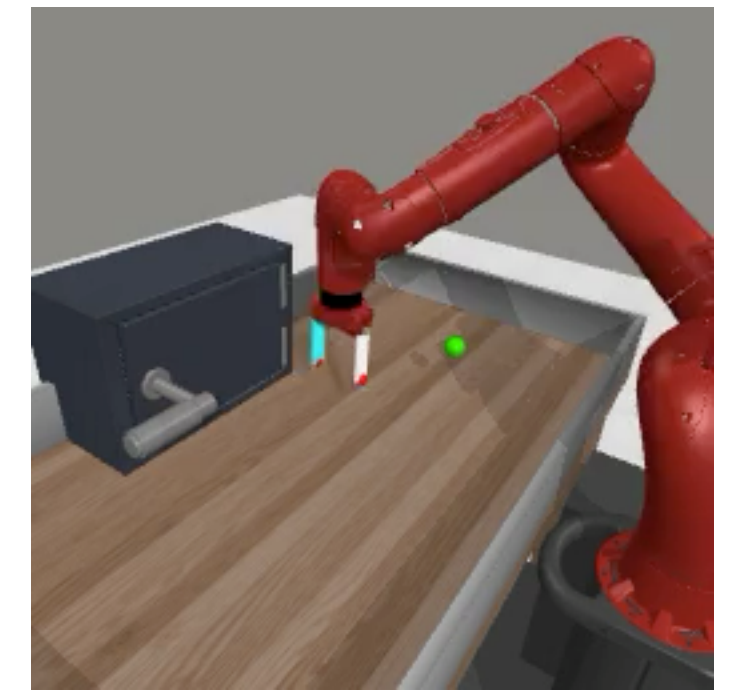
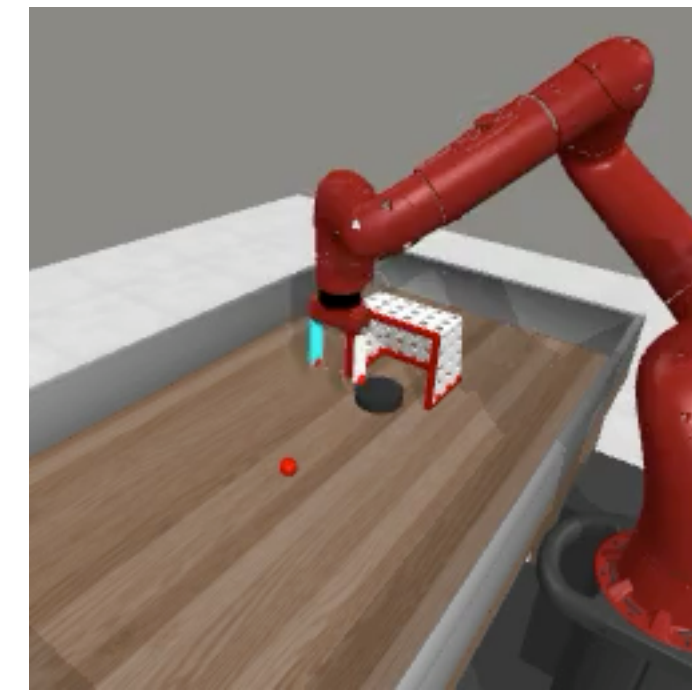
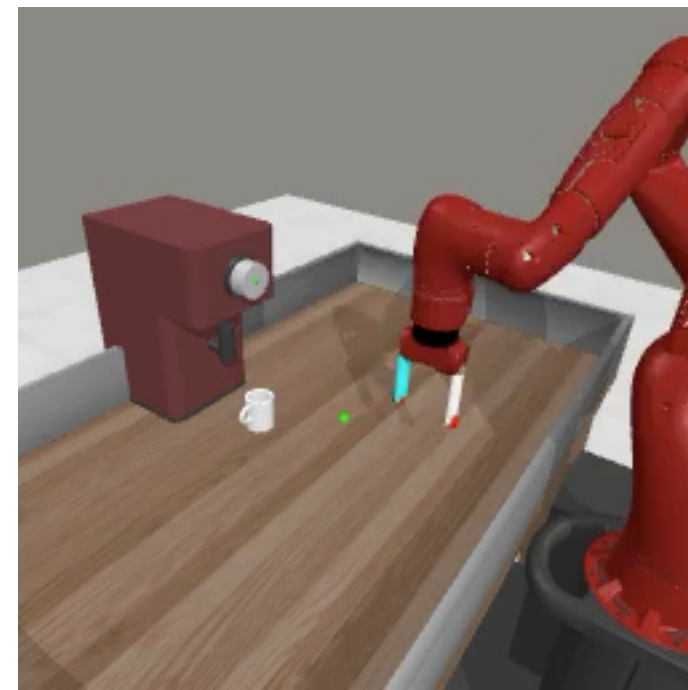
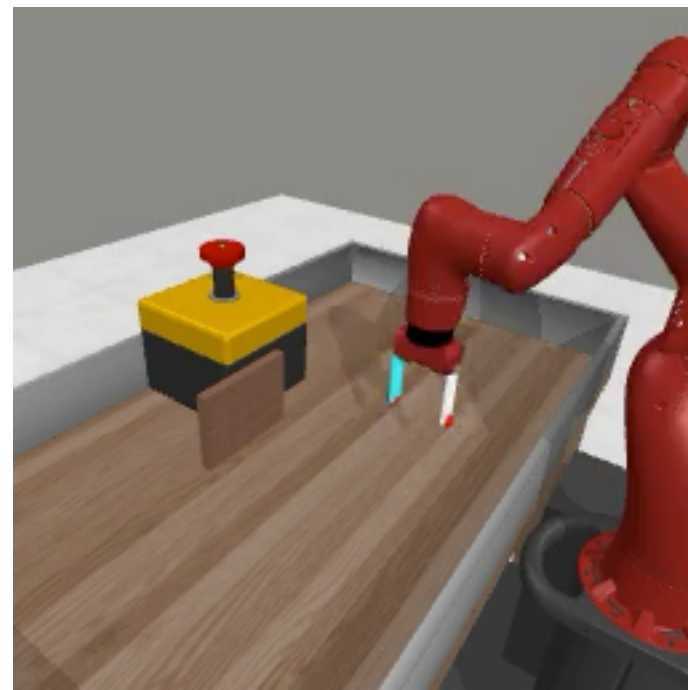
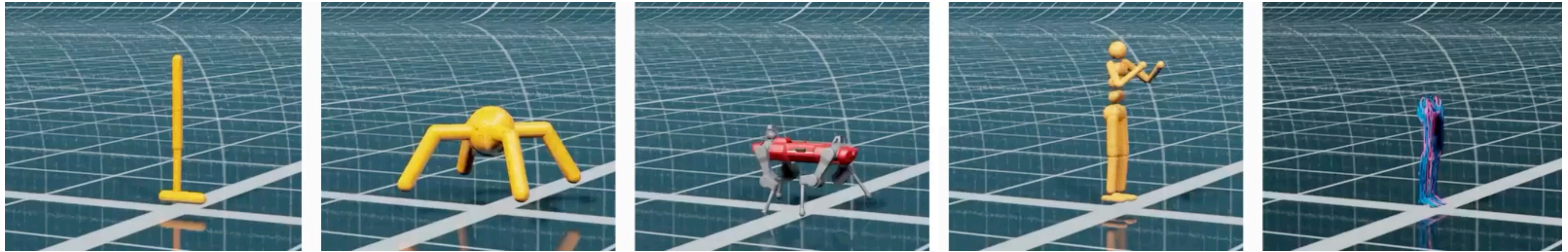
Conclusion

To get more efficient RL, we .. have to stop doing classic RL

PWM is a more scalable multi-task world model approach using FoG

When implemented correctly, world models can act as smooth surrogates of the true objective, resulting in better policies

Instead of training world models concurrently with policies, we should treat them as learned “offline simulations”



More at: <https://policy-world-model.github.io/>

Thank you

- Papers, code, data and more at: imgeorgiev.com/pwm/



Ignat
Georgiev



Varun
Giridhar



Nicklas
Hansen



Animesh
Garg