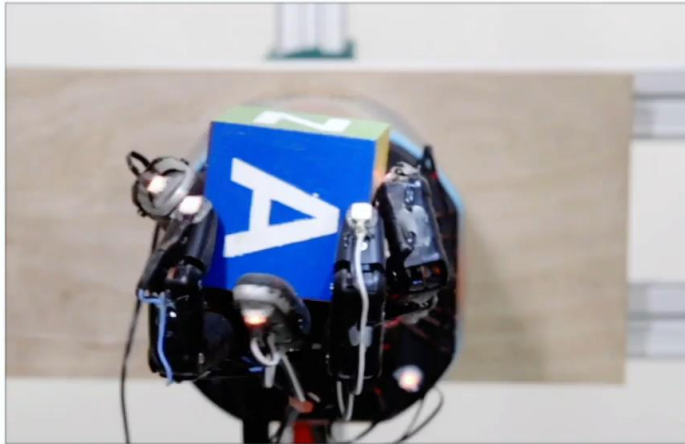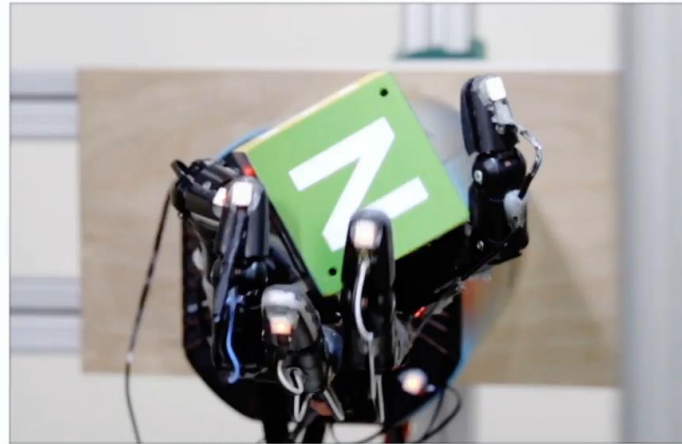# PODS: Policy Optimization via Differentiable Simulation

Miguel Zamora*, Momchil Peychev, Sehoon Ha,
Martin Vechev, Stelian Coros
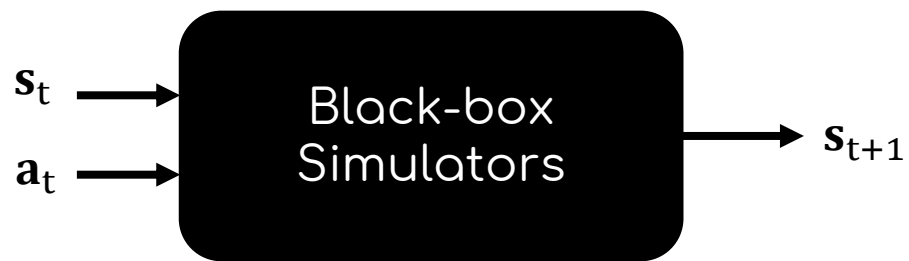
CRL

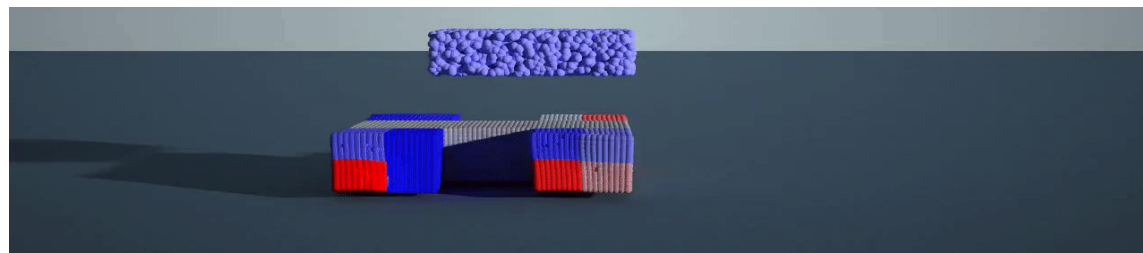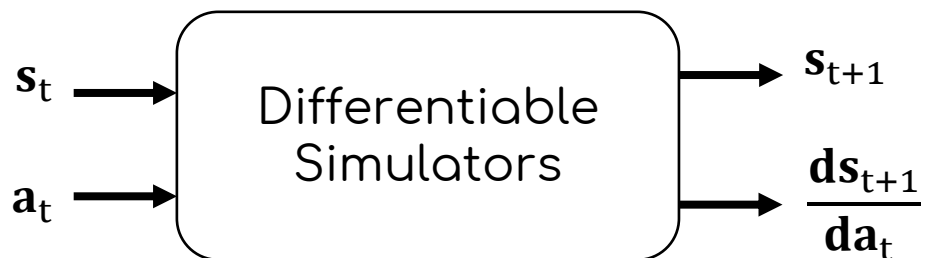**FINGER PIVOTING**           **SLIDING**           **FINGER GAITING**

Black-box Simulators

$s_t$

$a_t$

$s_{t+1}$

Goal:

How can we best use differentiable simulators to learn policies?

Differentiable Simulators

$s_t$

$a_t$

$s_{t+1}$

$\dfrac{ds_{t+1}}{da_t}$

Differentiable Cloth Simulation for Inverse Problems

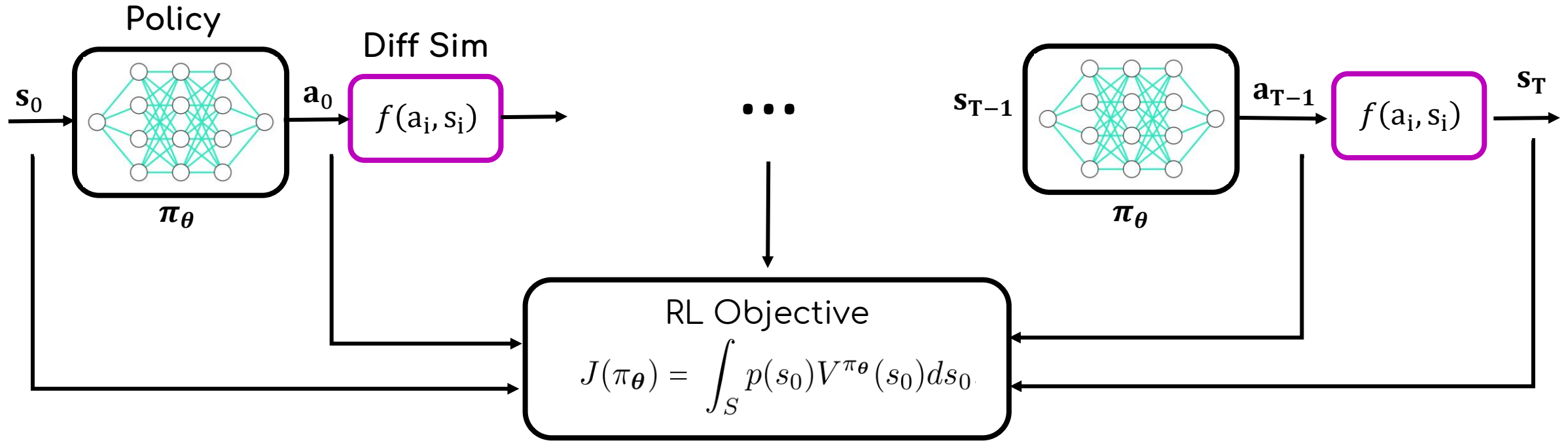DiSECt: A Differentiable Simulation Engine for Autonomous Robotic Cutting

DiffTaichi: Differentiable Programming for Physical Simulation

Target distance 5.5cm

ADD: Analytically Differentiable Dynamics for Multi-Body Systems with Frictional Contact

# Differentiable simulation as a layer?



Policy Gradient: $\nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\theta}}) = \int_S p(s_0) \nabla_{\boldsymbol{\theta}} V^{\pi_{\boldsymbol{\theta}}}(s_0) ds_0.$

$$\approx \frac{1}{k} \sum_i^k \nabla_{\boldsymbol{\theta}} V^{\pi_{\boldsymbol{\theta}}}(s_{0,i}).$$
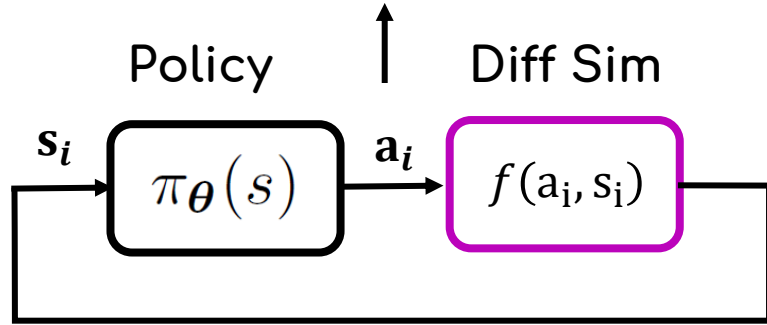
Effectively BPTT:

- Only first order
- Exploding / Vanishing gradients

# PODS

**Policy Improvement:**

$$V^{\pi_\theta}(s_0) = V^{\bar{a}}(s_0)$$

$$\bar{a} = \begin{bmatrix} a_0, a_1, \dots, a_{N-1} \end{bmatrix}$$

Policy     ↑     Diff Sim

$s_i$ → $\pi_\theta(s)$ → $a_i$ → $f(a_i, s_i)$

Improve $\bar{a}$ to get $V^{\pi_\theta}(s_0) < V^{\bar{a}}(s_0)$

$$\bar{a} = \pi_\theta + \alpha_a \frac{\mathrm{d}V^{\bar{a}}(s_0)}{\mathrm{d}\bar{a}}$$

$$\frac{\mathrm{d}V^{\bar{a}}(s_0)}{\mathrm{d}\bar{a}} = \frac{\partial V^{\bar{a}}}{\partial \bar{a}} + \frac{\partial V^{\bar{a}}}{\partial s}\frac{\mathrm{d}s}{\mathrm{d}\bar{a}}$$

**Info provided by diff sim**

**Second order improvement**

**Policy update**
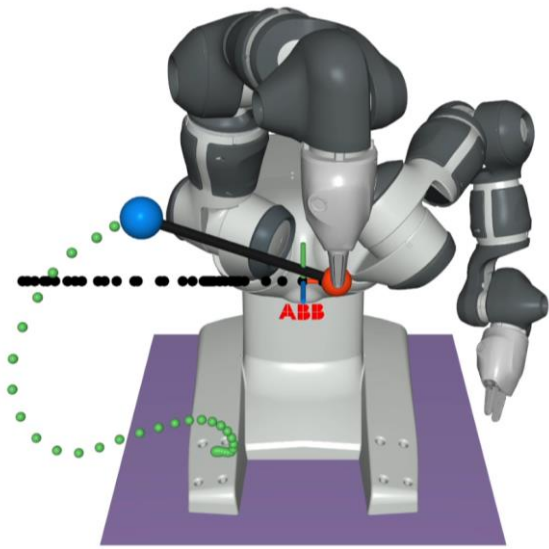
$$L_\theta = \frac{1}{k}\sum_i^k \sum_t^N \frac{1}{2}\|\pi_\theta(s_{t,i}) - a_{t,i}\|^2$$

$$\bar{a} = \pi_\theta + \alpha_a \hat{\mathbf{H}}^{-1}\frac{\mathrm{d}V^{\bar{a}}(s_0)}{\mathrm{d}\bar{a}}$$

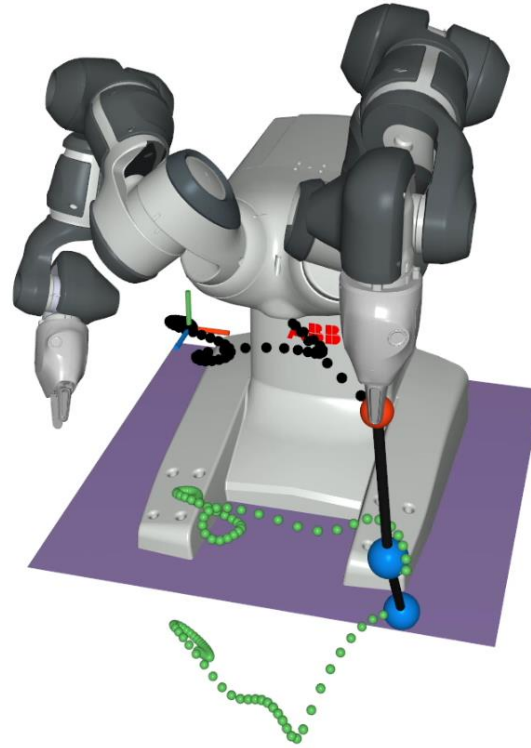**Check the paper for more details!**
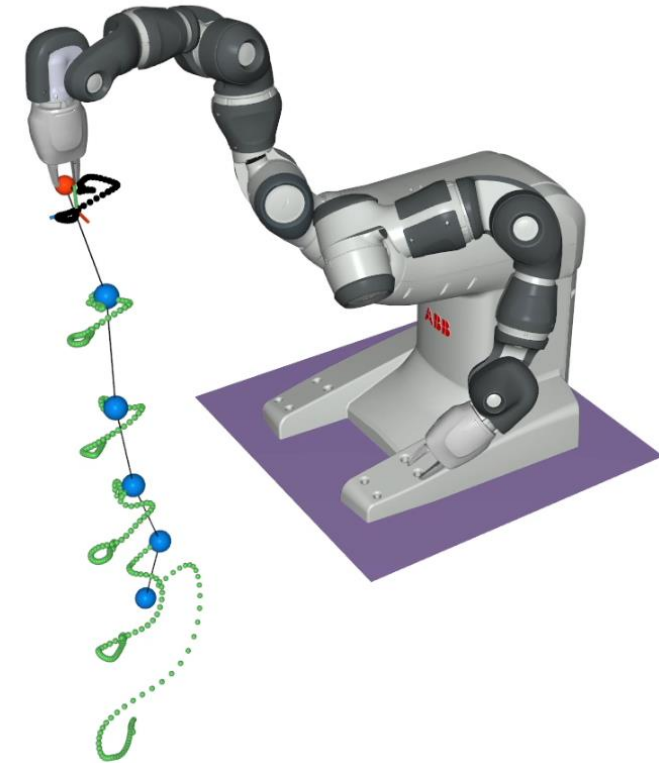
# Fine manipulation tasks



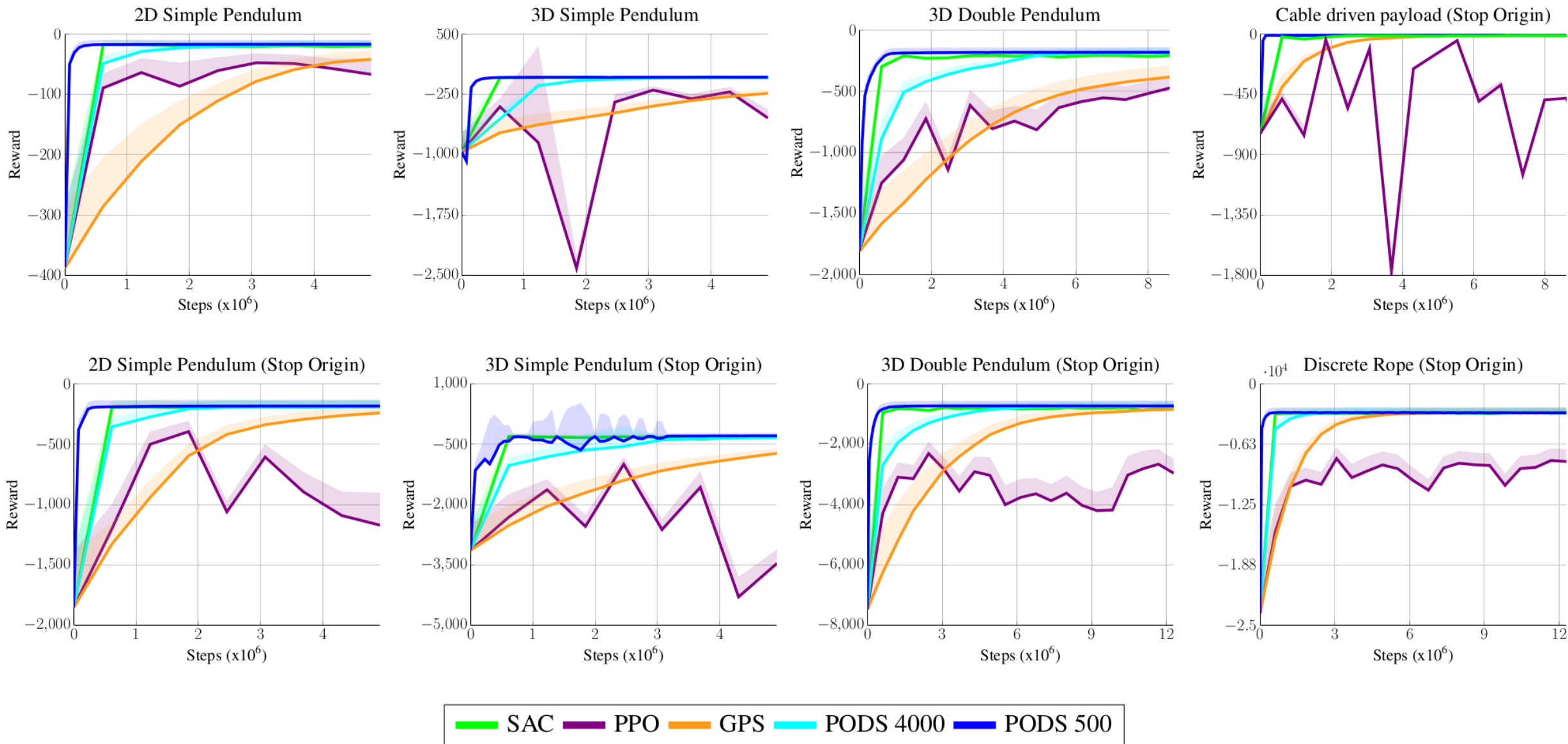2D Pendulum
Stop as fast as
possible

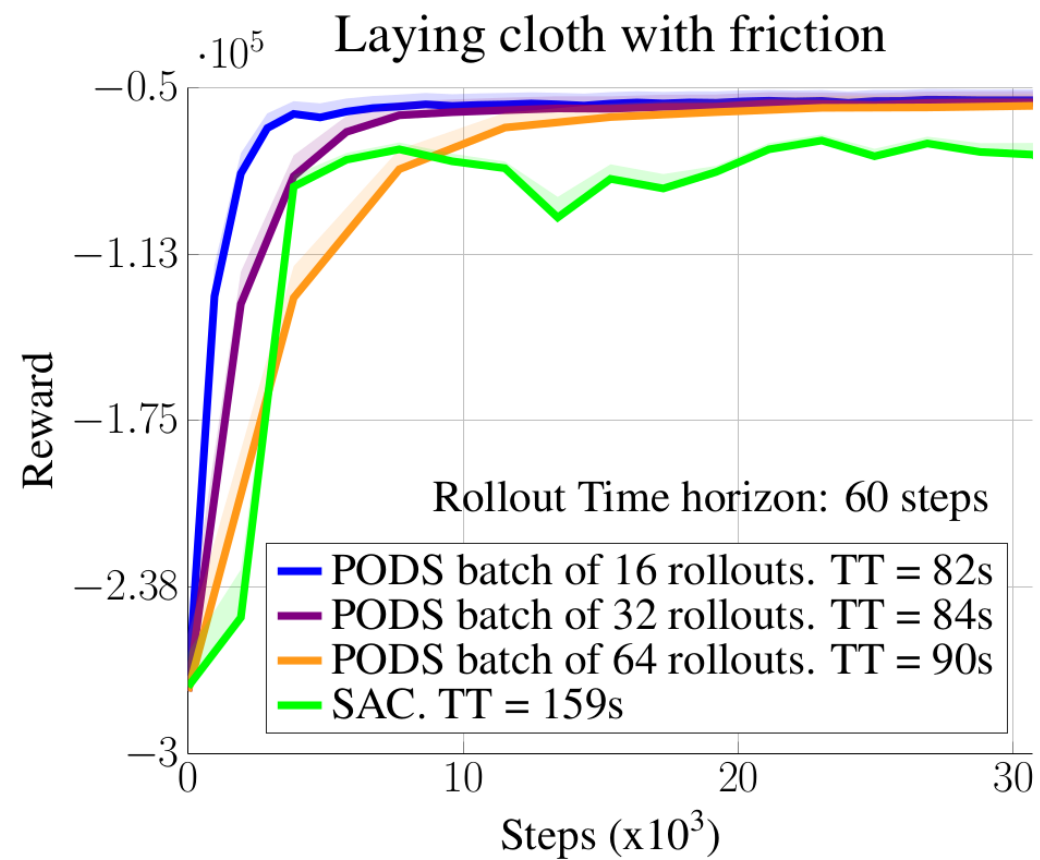3D Pendulum
Stop at origin

3D Double Pendulum
Stop at origin

Discretized Rope
Stop at origin

# PODS – Performance

Legend: SAC, PPO, GPS, PODS 4000, PODS 500

# PODS – Complexity



Sampled initial state

Intermediate state

Target location

Laying cloth with friction

Rollout Time horizon: 60 steps

- PODS batch of 16 rollouts. TT = 82s
- PODS batch of 32 rollouts. TT = 84s
- PODS batch of 64 rollouts. TT = 90s
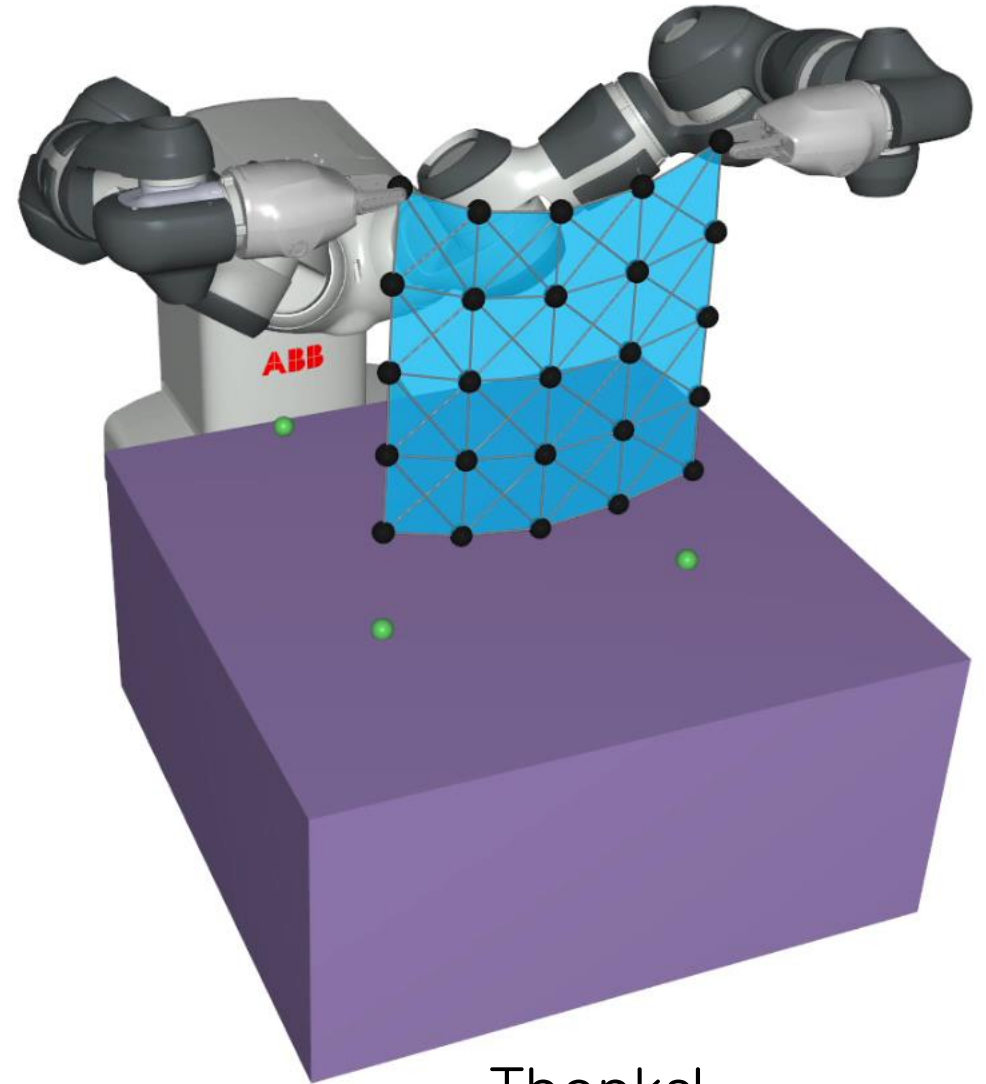- SAC. TT = 159s

Reward

Steps (x$10^3$)

# Conclusion

- Simple, fast and principled method
- Better exploit differentiable sims
- Outperformed baselines w.r.t. sample efficiency and compute time.

# Future work

- Interleave with existing RL methods (exploration).
- Leverage Inverse-RL to obtain surrogate reward function for non-smooth rewards.
- Find ceiling of complexity!



Thanks!