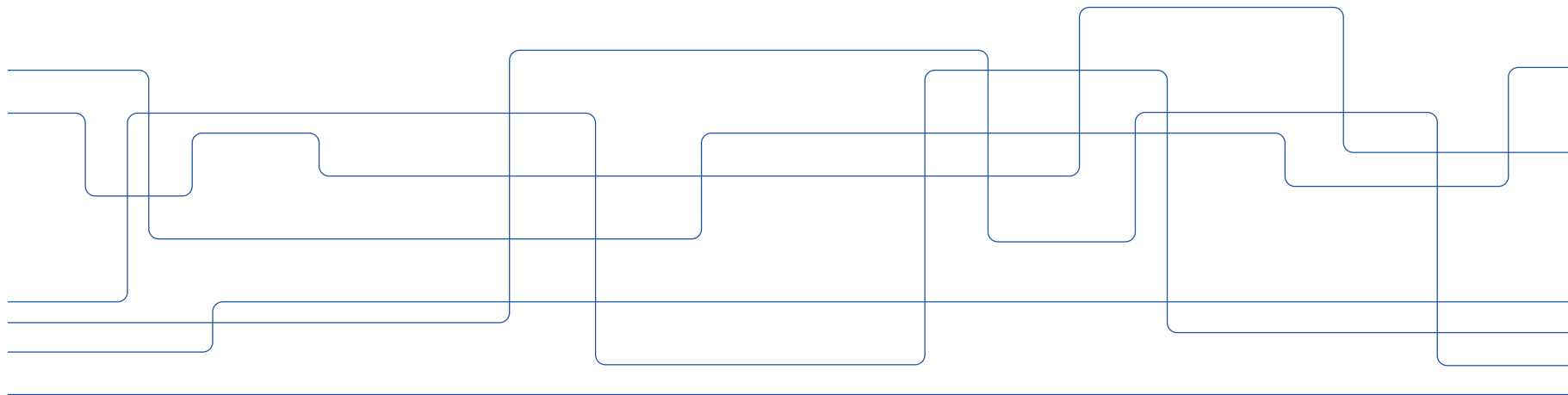




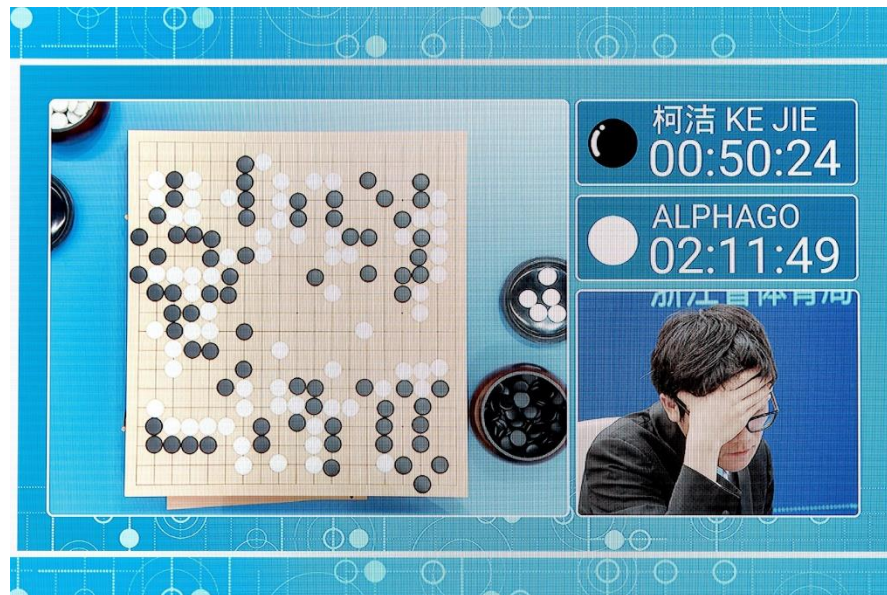
Data Efficiency for RL in Control Applications

Seminar for FDD3359 – Reinforcement Learning



Key to RL Successes

- AlphaGo
- Environment Model
 - Exact $p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$
 - Cheap to query
 - Extensive offline training
 - Online Planning
 - > *MCTS Rollouts*
 - > *Value function from offline training*



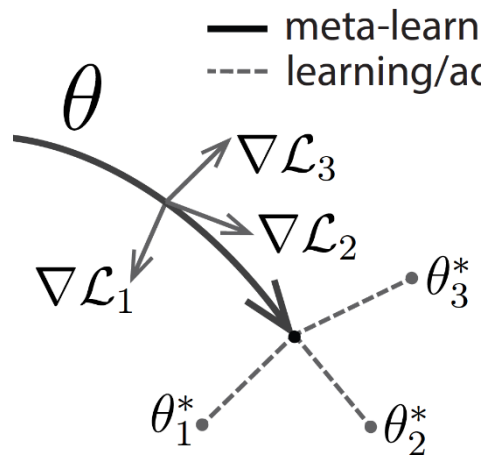
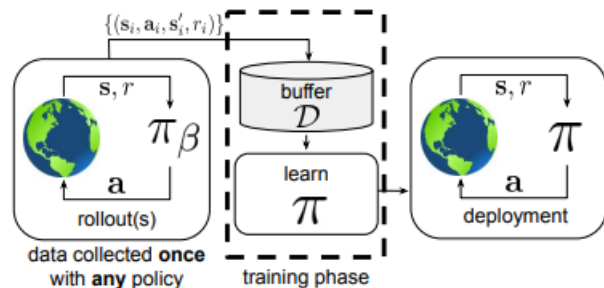
Real Systems

- Most problems
- Environment Model
 - No exact $p(s_{t+1}|s_t, a_t)$
 - Expensive/dangerous to query
 - Limited data for offline training
 - Online Planning?



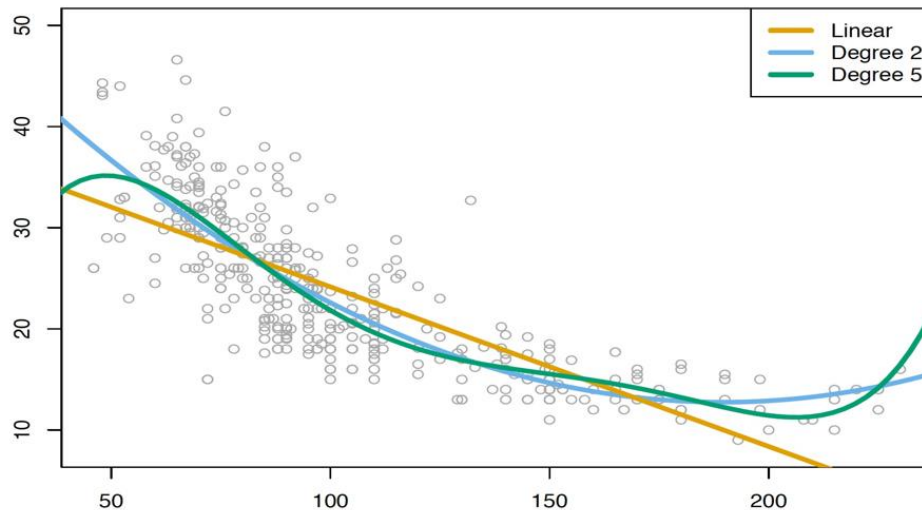
Possible Solutions

- Offline Reinforcement Learning
 - Fixed dataset
 - Avoid online environment query
- Meta Reinforcement Learning
 - Few-shot adaptation for target task
 - Limit environment query
- ...



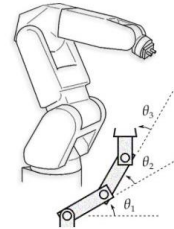
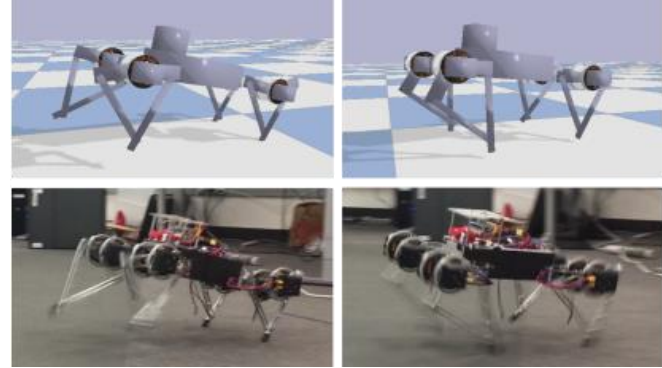
General Rule for Practising Machine Learning

- More data
 - High capacity model
 - Overfitting
- More domain knowledge
 - Choose “right” model
 - Limit the parameter space



More Data...

- Synthetic data from simulation
 - Various robot dynamics simulators
- Sim-to-real gap
 - Dynamics
 - > *Contact/actuator/friction...*
 - Perception



Learning Quadrupedal Locomotion over Challenging Terrain

JOONHO LEE^{1,*}, JEMIN HWANGBO^{1,2,†}, LORENZ WELLHAUSEN¹, VLADLEN KOLTUN³, AND MARCO HUTTER¹

¹Robotic Systems Lab, ETH Zurich, Zurich, Switzerland

²Robotics and Artificial Intelligence Lab, KAIST, Daejeon, Korea

³Intelligent Systems Lab, Intel, Santa Clara, CA, USA

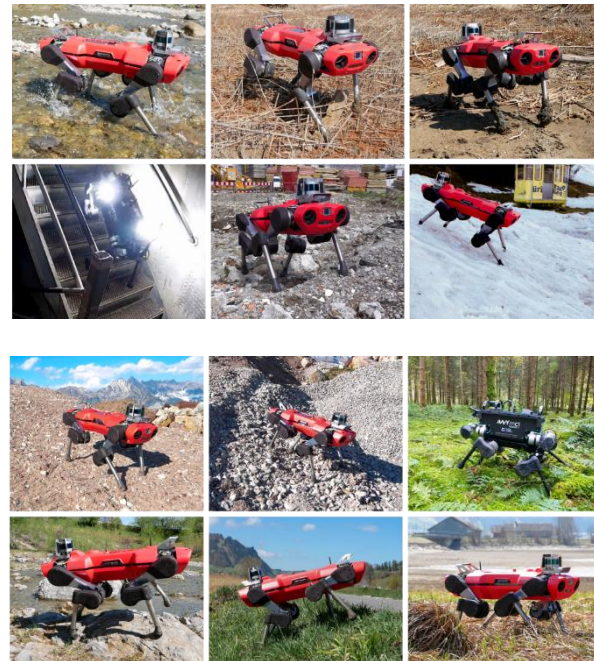
[†]Substantial part of the work was carried out during his stay at 1

^{*}Corresponding author: jolee@ethz.ch

This is the accepted version of Science Robotics Vol. 5, Issue 47, eabc5986 (2020)

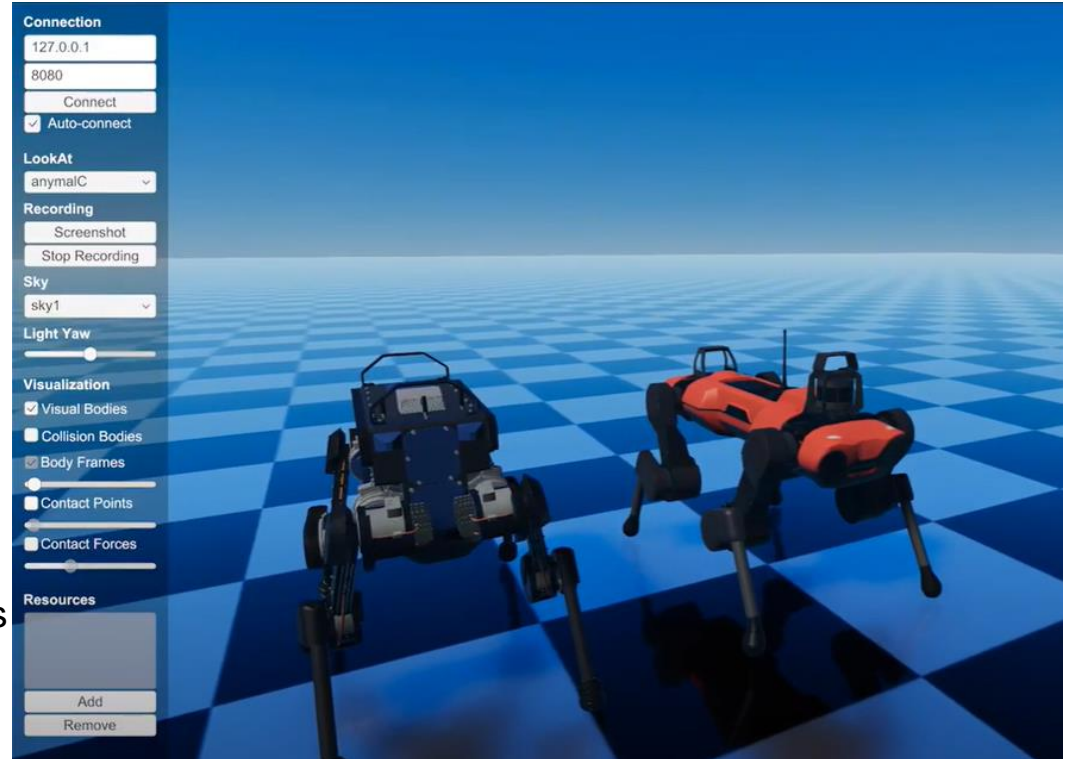
DOI: 10.1126/scirobotics.abc5986

Some of the most challenging environments on our planet are accessible to quadrupedal animals but remain out of reach for autonomous machines. Legged locomotion can dramatically expand the operational domains of robotics. However, conventional controllers for legged locomotion are based on elaborate state machines that explicitly trigger the execution of motion primitives and reflexes. These designs have escalated in complexity while falling short of the generality and robustness of animal locomotion. Here we present a radically robust controller for legged locomotion in challenging natural environments. We present a novel solution to incorporating proprioceptive feedback in locomotion control and demonstrate remarkably zero-shot generalization from simulation to natural environments. The controller is trained by reinforcement learning in simulation. It is based on a neural network that acts on a stream of proprioceptive signals. The trained controller has taken two generations of quadrupedal ANYmal robots to a variety of natural environments that are beyond the reach of prior published work in legged locomotion. The controller retains its robustness under conditions that have never been encountered during training: deformable terrain such as mud and snow, dynamic footholds such as rubble, and overground impediments such as thick vegetation and gushing water. The presented work opens new frontiers for robotics and indicates that radical robustness in natural environments can be achieved by training in much simpler domains.



More Data... from More Knowledge

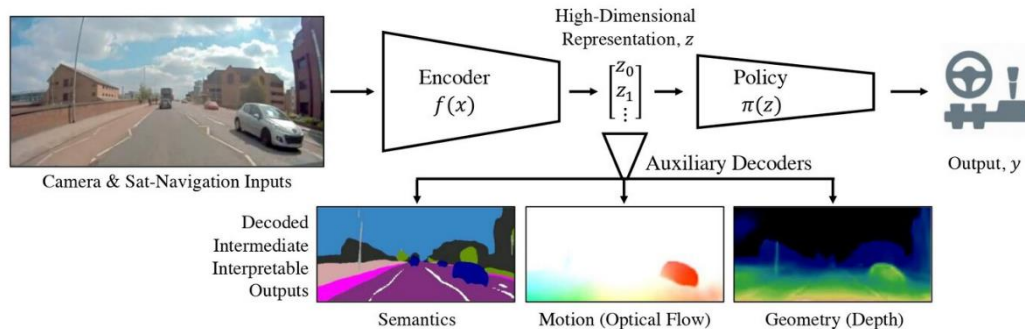
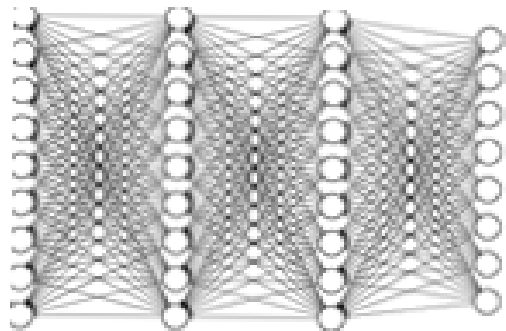
- Dedicated simulation RAI
 - Contact model
 - Actuator dynamics
 - Disturbances
 - Parameter randomization
- Terrain curriculum
- Teacher/student networks
 - Privileged simulation params



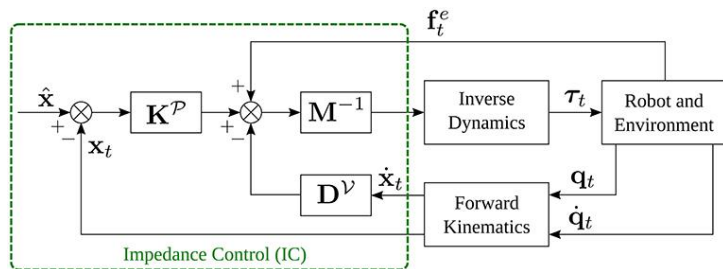
“Right” Policy

$$\pi(\mathbf{a}|\mathbf{s})$$

0.2	0.3	0.4	0.3
0.3	0.3	0.5	0.3
0.4	0.4	0.6	0.4
0.5	0.5	0.7	0.5



“Right” Policy for Controlling Robots



$$\underset{F}{\text{minimize}} \quad \sum_{k=1}^M L(\mathbf{q}[k], \mathbf{v}[k], \ddot{\mathbf{r}}[k], \boldsymbol{\lambda}[k], h[k])$$

$$\text{subject to} \quad m\ddot{\mathbf{r}}[k] = m\mathbf{g} + \sum_j \boldsymbol{\lambda}_j[k] \quad (\text{linear momentum})$$

$$\mathbf{k}[k] = \mathbf{A}_G^k(\mathbf{q}[k])\mathbf{v}[k] \quad (\text{angular momentum})$$

$$\dot{\mathbf{k}}[k] = \sum_j (\mathbf{c}_j[k] - \mathbf{r}[k]) \times \boldsymbol{\lambda}_j \quad (\text{angular momentum rate})$$

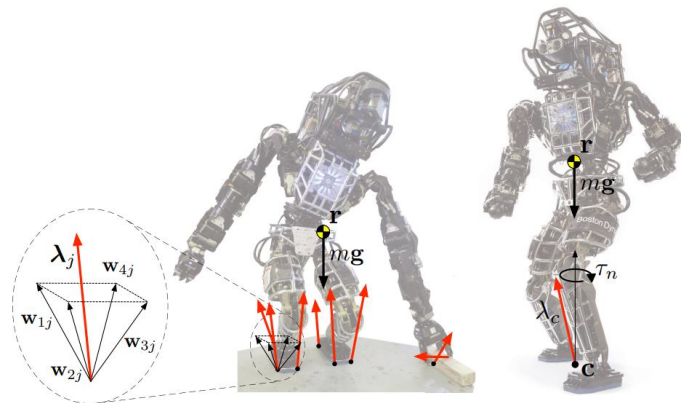
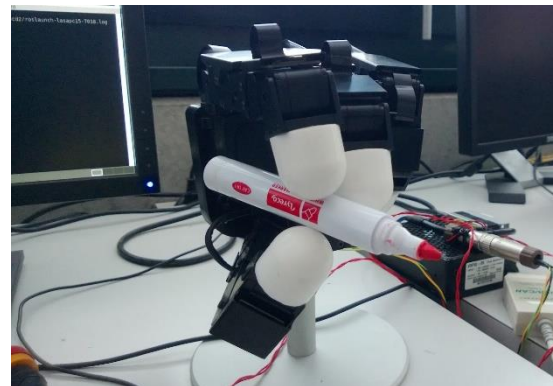
$$\forall_j \quad \boldsymbol{\lambda}_j[k] = \sum_{i=1}^{N_a} \beta_{ij}[k] \mathbf{w}_{ij} \quad (\text{friction})$$

$$\forall_{i,j} \quad \beta_{ij}[k] \geq 0$$

$$\mathbf{r}[k] = \text{COM}(\mathbf{q}[k]) \quad (\text{COM location})$$

Kinematic constraints

Time integration constraints



Perspective from Optimal Control

Reinforcement Learning

State/Action **\mathbf{s}** **\mathbf{a}**

Reward $r(\mathbf{s}, \mathbf{a})$

State transition $p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$

$$\arg \max_{\pi} \mathbb{E}_{\tau \sim p_{\pi}(\tau)} [r(\mathbf{s}_t, \mathbf{a}_t)]$$

$$\tau \sim p_{\pi} = p(\mathbf{s}_0) \prod_t p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi(\mathbf{a}_t | \mathbf{s}_t)$$

Optimal Control

State/Control **\mathbf{x}** **\mathbf{u}**

Cost $l(\mathbf{x}, \mathbf{u})$


Dynamics $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$

$$\arg \min_{\mathbf{u}} \int l(\mathbf{x}(t), \mathbf{u}(t)) dt$$

$$s.t. \quad \dot{\mathbf{x}}(t) = f(\mathbf{x}(t), \mathbf{u}(t))$$

Revisit to Value Iteration

value iteration algorithm:

- 
- A green curved arrow pointing from the first step to the second step, indicating a sequence or flow.
1. set $Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E[V(\mathbf{s}')]]$
 2. set $V(\mathbf{s}) \leftarrow \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$

$$V^*(\mathbf{x}) = \max_a [r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} (V^*(\mathbf{s}'))]$$

for optimal policy and any state

Hamilton-Jacobi-Bellman Equation

- Similar equation in continuous case

$$V^*(\mathbf{x}, t) = \min_{\mathbf{u}} \int_t l(\mathbf{x}(t), \mathbf{u}(t)) dt$$

- We have

$$-\frac{\partial V^*(\mathbf{x}, t)}{\partial t} = \min_{\mathbf{u}} [l(\mathbf{x}, \mathbf{u}) + \frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u})]$$

- vs Value Iteration

$$V^*(\mathbf{x}) = \max_a [r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})} (V^*(\mathbf{s}'))]$$


HJB Equation - Derivation

- For a small time interval to future

$$V^*(\mathbf{x}, t) = \min_{\mathbf{u}} [V^*(\mathbf{x}(t + dt), t + dt) + \int_t^{t+dt} l(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau]$$

- Taylor expansion of $V^*(\mathbf{x}(t + dt), t + dt)$

$$V^*(\mathbf{x}(t + dt), t + dt) = V^*(\mathbf{x}, t) + \frac{\partial V^*(\mathbf{x}, t)}{\partial t} dt + \frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \cdot \dot{\mathbf{x}} dt + o(dt)$$

A red arrow points downwards from the boxed $\dot{\mathbf{x}}$ term in the equation above to the definition of $\dot{\mathbf{x}}$.
$$\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$$

HJB Equation - Derivation

- Ignore high-order infinitesimal and subtract $V^*(\mathbf{x}, t)$

$$V^*(\mathbf{x}, t) = \min_{\mathbf{u}} [V^*(\mathbf{x}(t + dt), t + dt) + \int_t^{t+dt} l(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau]$$




$$0 = \frac{\partial V^*(\mathbf{x}, t)}{\partial t} dt + \min_{\mathbf{u}} \left[\frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u}) dt + \int_t^{t+dt} l(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau \right]$$

HJB Equation - Derivation

- Divide by $dt \rightarrow 0$

$$0 = \frac{\partial V^*(\mathbf{x}, t)}{\partial t} dt + \min_{\mathbf{u}} \left[\frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u}) dt + \int_t^{t+dt} l(\mathbf{x}(\tau), \mathbf{u}(\tau)) d\tau \right]$$

A red arrow points from the $\frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u}) dt$ term in the equation above to the corresponding term in the equation below.
$$-\frac{\partial V^*(\mathbf{x}, t)}{\partial t} = \min_{\mathbf{u}} \left[l(\mathbf{x}, \mathbf{u}) + \frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u}) \right]$$

- Note $\frac{\partial V^*(\mathbf{x}, t)}{\partial t} = 0$ if cost-to-go is static $V^*(\mathbf{x})$

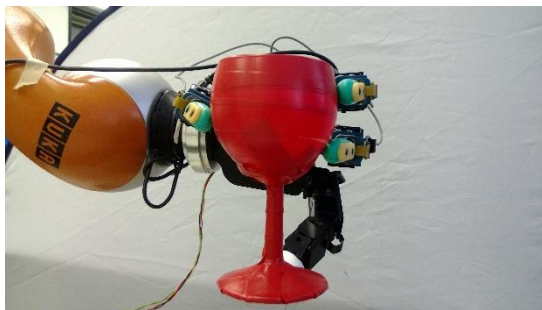
$$0 = \min_{\mathbf{u}} \left[l(\mathbf{x}, \mathbf{u}) + \frac{\partial V^*(\mathbf{x})}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u}) \right]$$

Control Affine Dynamics

- Structured $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$

$$\dot{\mathbf{x}} = f(\mathbf{x}) + g(\mathbf{x}) \cdot \mathbf{u}$$

- General motion equation of many physical systems, e.g. robots, rigid and deformable objects...



MuJoCo computes both forward and inverse dynamics in continuous time. Forward dynamics are then integrated over the specified `mjModel.opt.timestep` with the chosen [numerical integrator](#). The general equations of motion in continuous time are

$$M\dot{v} + c = \tau + J^T f \quad (1)$$

The Jacobian establishes the relationship between quantities in joint and constraint coordinates. It maps motion vectors (velocities and accelerations) from joint to constraint coordinates: the joint velocity v maps to velocity Jv in constraint coordinates. The transpose of the Jacobian maps force vectors from constraint to joint coordinates: the constraint force f maps to force $J^T f$ in joint coordinates.

The joint-space inertia matrix M is always invertible. Therefore once the constraint force f is known, we can finalize the forward and inverse dynamics computations as

$$\text{forward: } \dot{v} = M^{-1}(\tau + J^T f - c)$$

$$\text{inverse: } \tau = M\dot{v} + c - J^T f$$

Decomposition of Cost Function

- Structured $l(\mathbf{x}, \mathbf{u}) = c(\mathbf{x}) + \mathbf{u}^\top \mathbf{R} \mathbf{u}$ $\mathbf{R} \succ 0$
- Use them in HJB equation and the optimality condition of \mathbf{u} ...

$$-\frac{\partial V^*(\mathbf{x}, t)}{\partial t} = \min_{\mathbf{u}} \left[l(\mathbf{x}, \mathbf{u}) + \frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u}) \right]$$



$$\mathbf{u}^* = \arg \min_{\mathbf{u}} \left[c(\mathbf{x}) + \mathbf{u}^\top \mathbf{R} \mathbf{u} + \frac{\partial V^*}{\partial \mathbf{x}} \cdot (f(\mathbf{x}) + g(\mathbf{x}) \cdot \mathbf{u}) \right]$$

Decomposition of Cost Function

$$\mathbf{u}^* = \arg \min_{\mathbf{u}} [c(\mathbf{x}) + \mathbf{u}^\top \mathbf{R} \mathbf{u} + \frac{\partial V^*}{\partial \mathbf{x}} \cdot (f(\mathbf{x}) + g(\mathbf{x}) \cdot \mathbf{u})]$$



$$\mathbf{u}^* = -\frac{1}{2} \mathbf{R}^{-1} g(\mathbf{x})^\top \left(\frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}} \right)^\top$$

A closed-form solution...

Further Approximation – Linear Quadratic Regulator

- Lets further the control affine and quadratic approximations

$$\dot{\mathbf{x}} = f(\mathbf{x}) + g(\mathbf{x}) \cdot \mathbf{u}$$



$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$$

$$l(\mathbf{x}, \mathbf{u}) = c(\mathbf{x}) + \mathbf{u}^\top \mathbf{R}\mathbf{u}$$

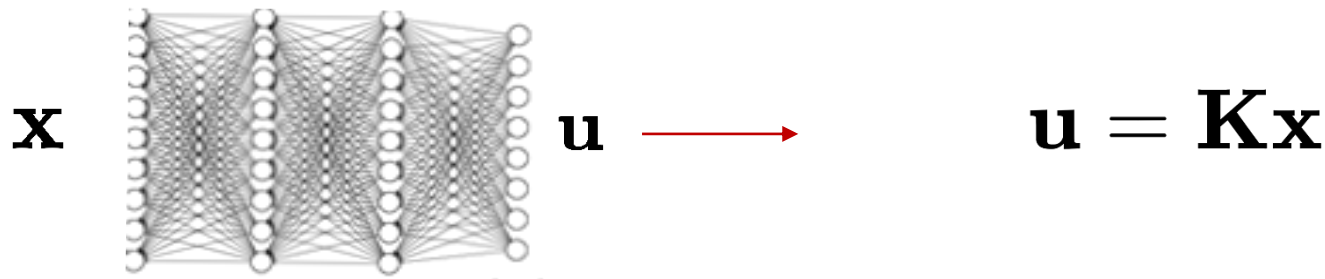
$$l(\mathbf{x}, \mathbf{u}) = \mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{u}^\top \mathbf{R}\mathbf{u}$$

- We will have quadratic $V^*(\mathbf{x}) = \mathbf{x}^\top \mathbf{P}\mathbf{x}$
- And a linear solution

$$\mathbf{u}^* = -\frac{1}{2}\mathbf{R}^{-1}g(\mathbf{x})^\top \left(\frac{\partial V^*(\mathbf{x}, t)}{\partial \mathbf{x}}\right)^\top \longrightarrow \mathbf{u}^* = -\frac{1}{2}\mathbf{R}^{-1}\mathbf{B}^\top \mathbf{P}\mathbf{x}$$

Power of Linear Policy

- So what if



- We can always assume differentiability and have local linear approximation...

$$\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u}) \quad \longrightarrow \quad \delta \dot{\mathbf{x}} = \left. \frac{\partial f}{\partial \mathbf{x}} \right|_{\mathbf{x}, \mathbf{u}} \delta \mathbf{x} + \left. \frac{\partial f}{\partial \mathbf{u}} \right|_{\mathbf{x}, \mathbf{u}} \delta \mathbf{u}$$

Power of Linear Policy

Simple random search of static linear policies is competitive for reinforcement learning

Horia Mania
hmania@berkeley.edu

Aurelia Guy
lia@berkeley.edu

Benjamin Recht
brecht@berkeley.edu

Department of Electrical Engineering and Computer Science
University of California, Berkeley

Abstract

Model-free reinforcement learning aims to offer off-the-shelf solutions for controlling dynamical systems without requiring models of the system dynamics. We introduce a model-free random search algorithm for training static, linear policies for continuous control problems. Common evaluation methodology shows that our method matches state-of-the-art sample efficiency on the benchmark MuJoCo locomotion tasks. Nonetheless, more rigorous evaluation reveals that the assessment of performance on these benchmarks is optimistic. We evaluate the performance of our method over hundreds of random seeds and many different hyperparameter configurations for each benchmark task. This extensive evaluation is possible because of the small computational footprint of our method. Our simulations highlight a high variability in performance in these benchmark tasks, indicating that commonly used estimations of sample efficiency do not adequately evaluate the performance of RL algorithms. Our results stress the need for new baselines, benchmarks and evaluation methodology for RL algorithms.

Algorithm 1 Augmented Random Search (ARS): four versions **V1**, **V1-t**, **V2** and **V2-t**

- 1: **Hyperparameters:** step-size α , number of directions sampled per iteration N , standard deviation of the exploration noise ν , number of top-performing directions to use b ($b < N$ is allowed only for **V1-t** and **V2-t**)
- 2: **Initialize:** $M_0 = \mathbf{0} \in \mathbb{R}^{p \times n}$, $\mu_0 = \mathbf{0} \in \mathbb{R}^n$, and $\Sigma_0 = \mathbf{I}_n \in \mathbb{R}^{n \times n}$, $j = 0$.
- 3: **while** ending condition not satisfied **do**
- 4: Sample $\delta_1, \delta_2, \dots, \delta_N$ in $\mathbb{R}^{p \times n}$ with i.i.d. standard normal entries.
- 5: Collect $2N$ rollouts of horizon H and their corresponding rewards using the $2N$ policies

$$\mathbf{V1:} \begin{cases} \pi_{j,k,+}(x) = (M_j + \nu \delta_k)x \\ \pi_{j,k,-}(x) = (M_j - \nu \delta_k)x \end{cases}$$

$$\mathbf{V2:} \begin{cases} \pi_{j,k,+}(x) = (M_j + \nu \delta_k) \text{diag}(\Sigma_j)^{-1/2} (x - \mu_j) \\ \pi_{j,k,-}(x) = (M_j - \nu \delta_k) \text{diag}(\Sigma_j)^{-1/2} (x - \mu_j) \end{cases}$$

for $k \in \{1, 2, \dots, N\}$.

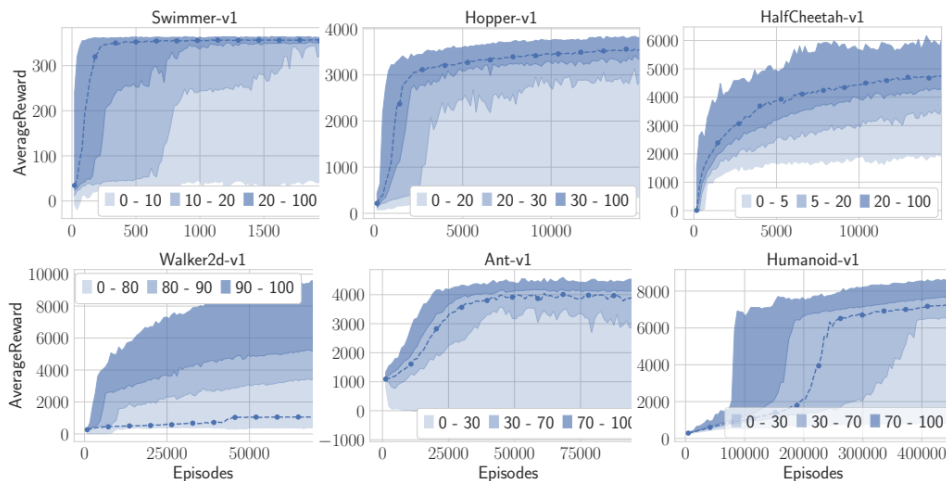
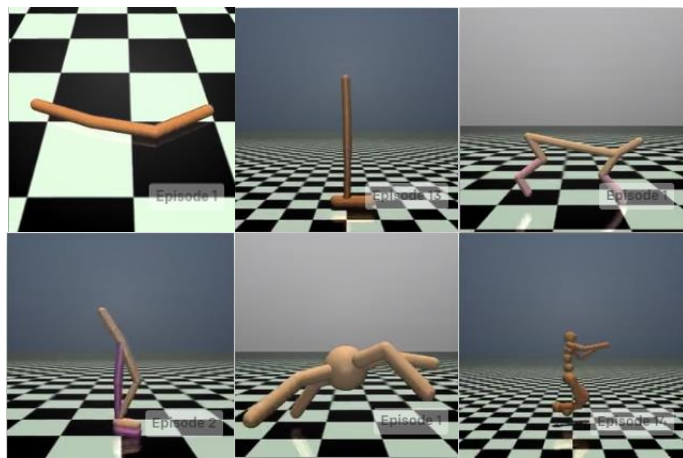
- 6: **V1-t, V2-t:** Sort the directions δ_k by $\max\{r(\pi_{j,k,+}), r(\pi_{j,k,-})\}$, denote by $\delta_{(k)}$ the k -th largest direction, and by $\pi_{j,(k),+}$ and $\pi_{j,(k),-}$ the corresponding policies.
- 7: Make the update step:

$$M_{j+1} = M_j + \frac{\alpha}{b\sigma_R} \sum_{k=1}^b [r(\pi_{j,(k),+}) - r(\pi_{j,(k),-})] \delta_{(k)},$$

where σ_R is the standard deviation of the $2b$ rewards used in the update step.

- 8: **V2:** Set μ_{j+1} , Σ_{j+1} to be the mean and covariance of the $2NH(j+1)$ states encountered from the start of training □
- 9: $j \leftarrow j + 1$
- 10: **end while**

Power of Linear Policy



Embedded Optimization as Policy

- Still have a policy with inspired structure without a closed-form solution

$$\mathbf{u}^* = \arg \min_{\mathbf{u}} \int l(\mathbf{x}(t), \mathbf{u}(t)) dt$$

$$s.t. \quad \dot{\mathbf{x}}(t) = f(\mathbf{x}(t), \mathbf{u}(t))$$

$$\mathbf{u}^* = \arg \min_{\mathbf{u}} \sum_t^T l(\mathbf{x}_t, \mathbf{u}_t)$$

$$s.t. \quad \mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$$

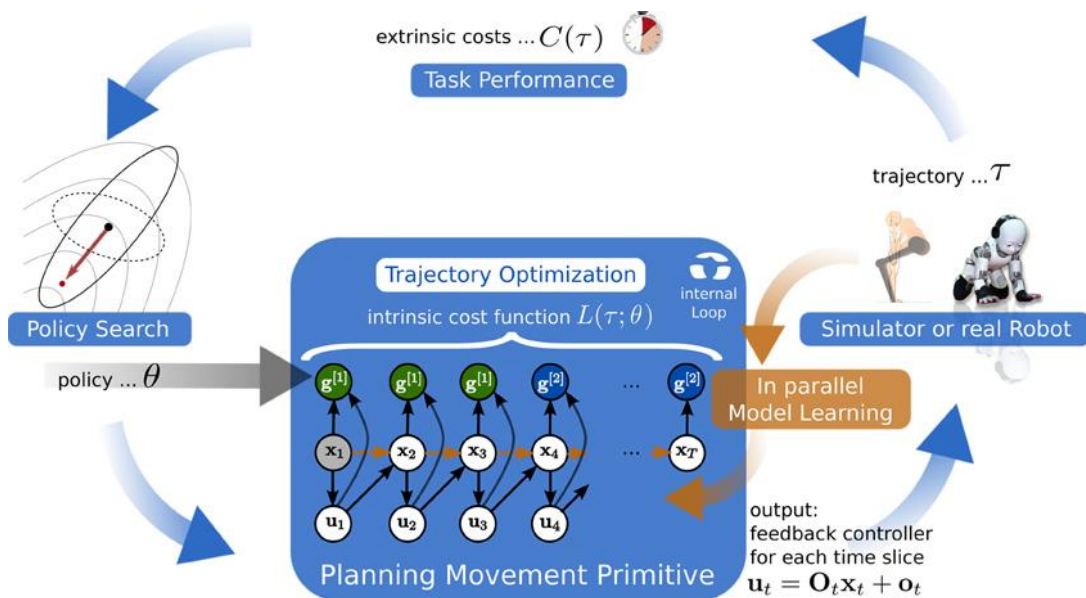


$$\nabla \mathbf{u}^* ?$$

- How about differentiating the entire optimization process for policy gradient algorithms – differentiable optimization

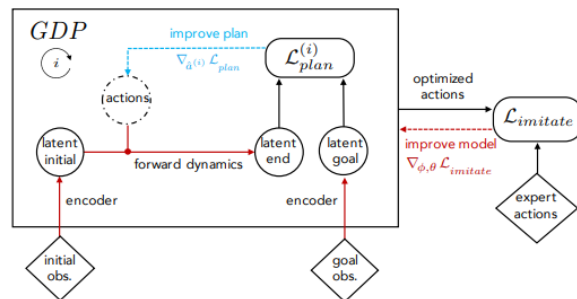
Embedded Optimization as Policy

- Using internal planning as primitives
 - Quadratic intrinsic cost
 - Inference for internal TO
 - Evolution strategy for PS

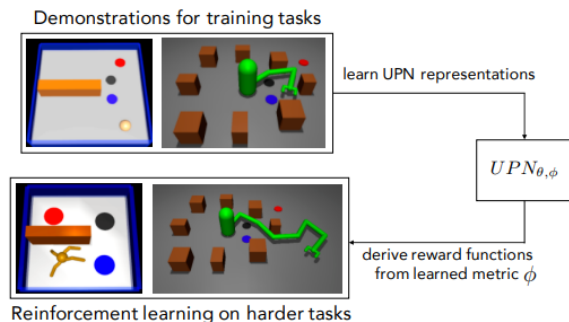


Embedded Optimization as Policy

- Unrolling optimization iteration steps
 - Leveraging automatic differentiation
 - Policy learning via imitation
 - Huber loss in the latent space
 - Using learned encoders for RL



(a) Universal Planning Network (UPN)



(b) Leveraging learned latent representations

Embedded Optimization as Policy

- Avoiding backwarding iterative steps for gradients by differentiating KKT conditions

$$\mathcal{L}(\tau, \lambda) = \sum_{t=1}^T \left(\frac{1}{2} \tau_t^\top C_t \tau_t + c_t^\top \tau_t \right) + \sum_{t=0}^{T-1} \lambda_t^\top (F_t \tau_t + f_t - x_{t+1})$$

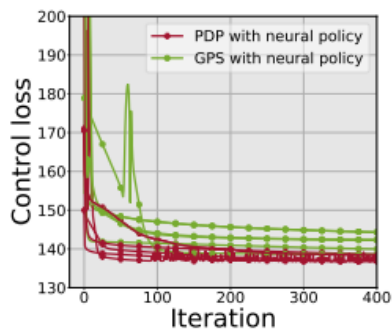
$$\nabla_{\tau_t} \mathcal{L}(\tau^*, \lambda^*) = C_t \tau_t^* + c_t + F_t^\top \lambda_t^* - \begin{bmatrix} \lambda_{t-1}^* \\ 0 \end{bmatrix} = 0$$

$$\begin{aligned} \nabla_{C_t} \ell &= \frac{1}{2} (d_{\tau_t}^* \otimes \tau_t^* + \tau_t^* \otimes d_{\tau_t}^*) & \nabla_{c_t} \ell &= d_{\tau_t}^* & \nabla_{x_{\text{init}}} \ell &= d_{\lambda_0}^* \\ \nabla_{F_t} \ell &= d_{\lambda_{t+1}}^* \otimes \tau_t^* + \lambda_{t+1}^* \otimes d_{\tau_t}^* & \nabla_{f_t} \ell &= d_{\lambda_t}^* \end{aligned}$$

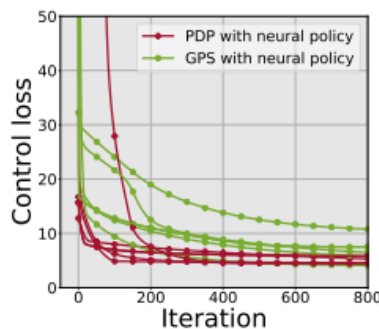
- With $\tau_t = \{\mathbf{x}_t, \mathbf{u}_t\}$ $F_t = [\mathbf{A}_t, \mathbf{B}_t]$ and linear terms c_t f_t
- Construct another LQR to solve $d_{\lambda_t}^*$ $d_{\tau_t}^*$

Embedded Optimization as Policy

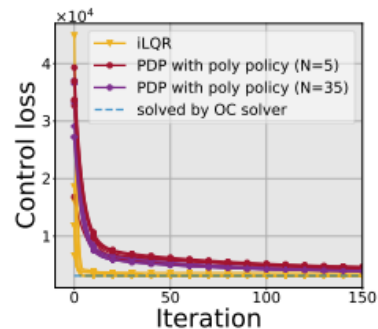
- Forwarding the policy still needs to iteratively
 - following gradient-based optimization (Jin et al, 2020)
 - or forming local LQ problems (Amos et al, 2018)
 - Solution must exist
 - Often requires parameterized dynamics model/trajectories



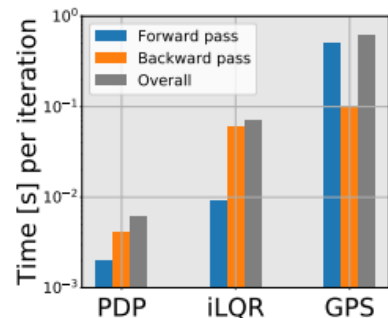
(a) Cart-pole control



(b) Robot arm control



(c) Quadrotor planning



(d) Timing results

Embedded Optimization as Policy

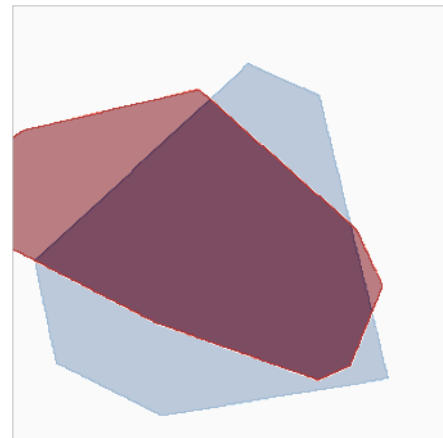
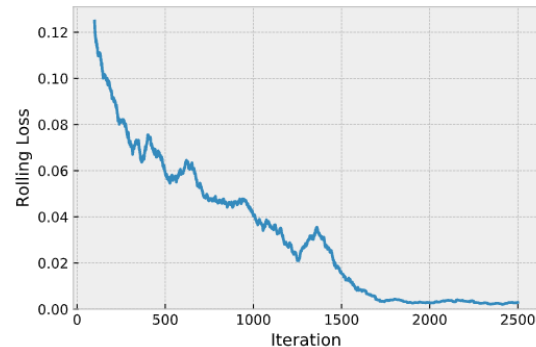
- Using convex optimization
 - Fast policy forwarding
 - No local optimality issue

- Differentiable layers

(<https://github.com/cvxgrp/cvxpylayers>)

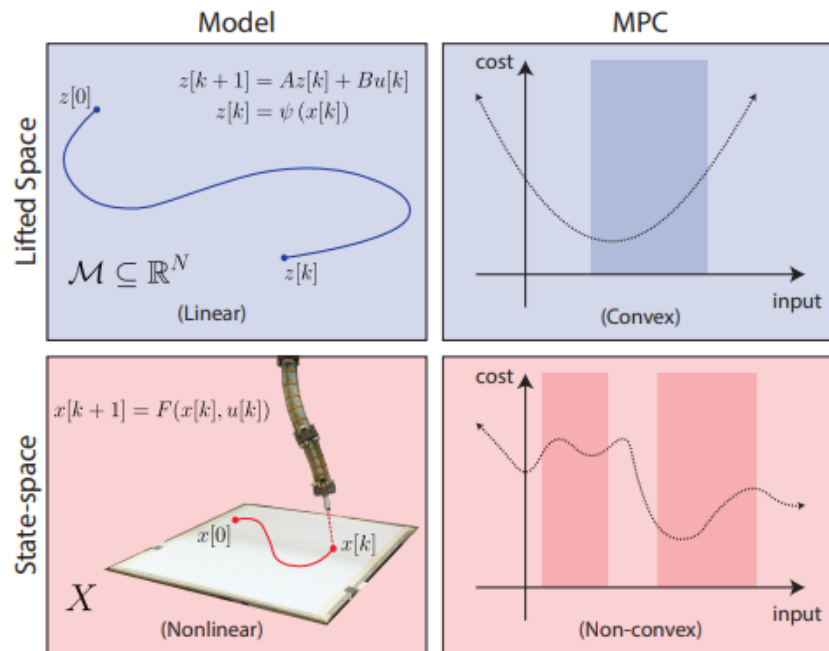
Generating data from

$$\hat{y} = \operatorname{argmin}_y \frac{1}{2} \|x - y\|_2^2$$
$$\text{s. t. } Gy \leq h.$$

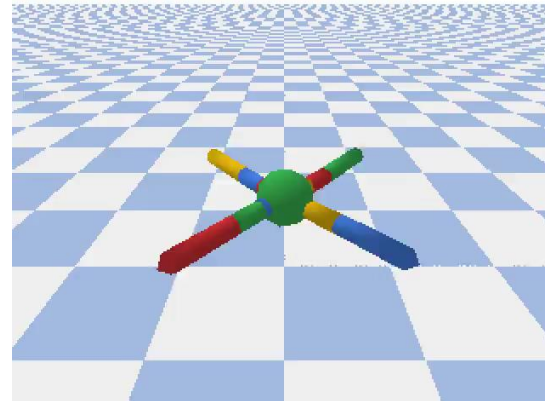
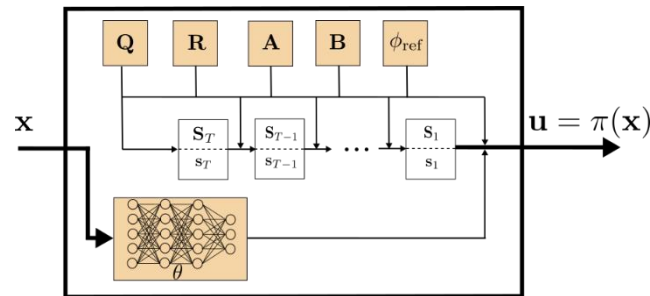


Convex Approximation – Koopman Representation

- How about non-convex dynamics constraints
- Using a high-dimensional representation such that states evolve with a linear model
- Can be LQR again using quadratic cost in the transformed representation space

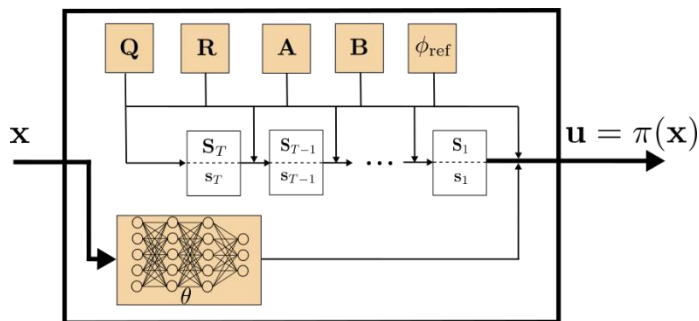


-
- AntBulletEnv-v0**
- Legend: — koopman_residual — koopman — vanilla
- Y-axis: Score (0 to 2500)
- X-axis: Epochs (0 to 800+)
- | Epochs | koopman_residual (Score) | koopman (Score) | vanilla (Score) |
|--------|--------------------------|-----------------|-----------------|
| 0 | ~500 | ~500 | ~500 |
| 200 | ~800 | ~800 | ~800 |
| 400 | ~1500 | ~1200 | ~1100 |
| 600 | ~2000 | ~1600 | ~1400 |
| 800 | ~2500 | ~2100 | ~1900 |

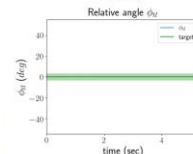
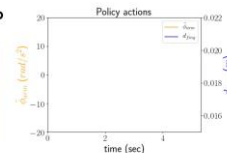
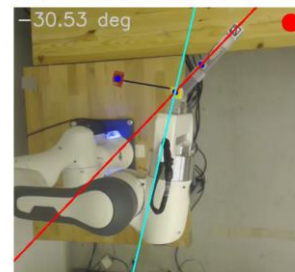


Convex Approximation – Koopman Representation

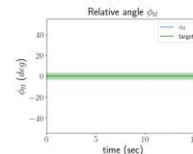
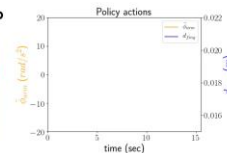
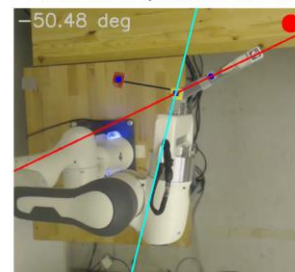
- And just neural network inference for policy evaluation



Pivot to 0° ; start at -30°



Pivot to 0° ; start at -50°



Policy with Robotic Control Structure

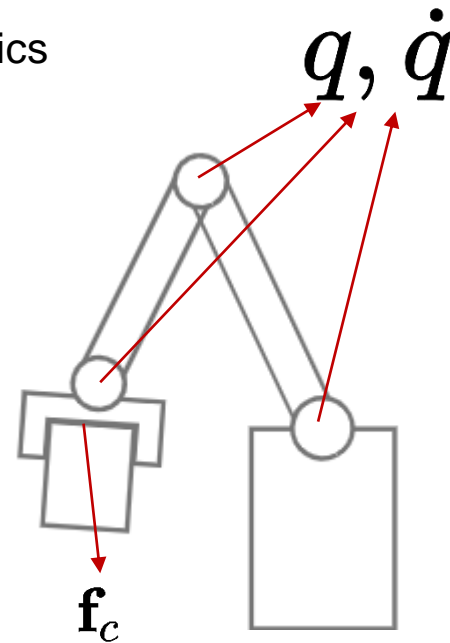
- We usually have a concrete form of manipulator dynamics

$$\dot{\mathbf{x}} = f(\mathbf{x}) + g(\mathbf{x}) \cdot \mathbf{u}$$



$$\mathbf{M}(q)\ddot{q} + \mathbf{b}(q, \dot{q}) + \mathbf{g}(q) = \mathbf{u} + \mathbf{J}_c(q)^\top \mathbf{f}_c$$

$$\mathbf{x} = \{q, \dot{q}\}$$

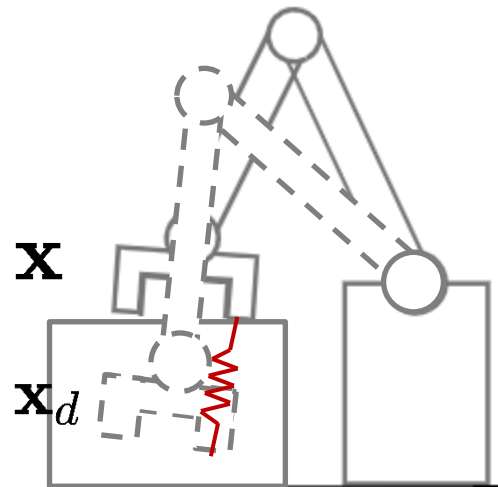


Motion Compliance

- Implicitly control of the interaction force

$$\mathbf{u} = -\mathbf{K}(\mathbf{x} - \mathbf{x}_d)$$

- Characterize robot behavior as a spring
- Policy learning is to find the right stiffness



Policy with Variable Compliance

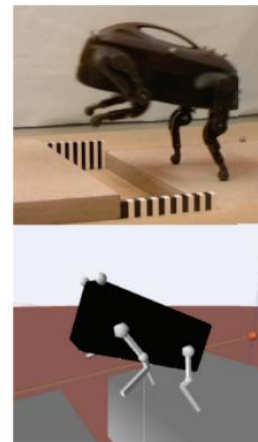
- Implicitly control of the interaction force

$$\mathbf{u} = -\mathbf{K}(\mathbf{x} - \mathbf{x}_d)$$

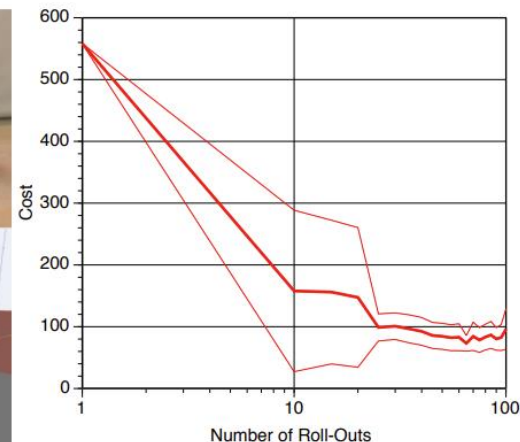


$$\mathbf{u}(t) = -\mathbf{K}(t)(\mathbf{x}(t) - \mathbf{x}_d(t))$$

- Parameterize compliance and desired trajectories for independent joint DOF
- Learning to jump with 100 trials



(a) Real & Simulated Robot Dog



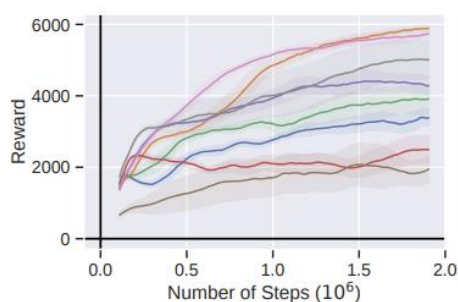
(b) Learning curve for Dog Jump with $\mathbf{PI}^2 \pm 1std$

Policy with Variable Compliance

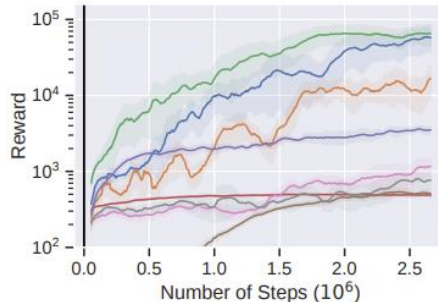
- Implicitly control of the interaction force

$$\mathbf{u} = -\mathbf{K}(\mathbf{x} - \mathbf{x}_d) \longrightarrow \mathbf{u} = -\mathbf{K}(\mathbf{x})\mathbf{x}$$

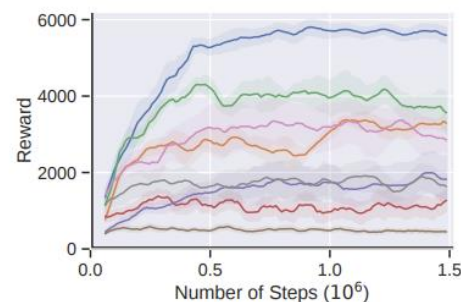
- State-dependent stiffness in the joint/end-effector space
- Neural networks to output positive \mathbf{K}



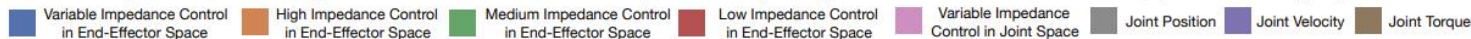
(a) Path Following



(b) Door Opening



(c) Surface Wiping

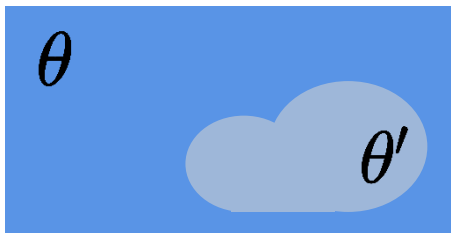


Control Guarantees for RL Policy

- Can we search policies that have some inherent guarantees?

$$\pi_{\theta} \longrightarrow \pi_{\theta'}$$

- Having some statements like $\forall \theta' \quad x(t) \in \dots$
- Stability/Convergence/Consensus...



Lyapunov Stability

For $\mathbf{u}(\mathbf{x})$ that steers trajectory of \mathbf{x} under $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$

An equilibrium \mathbf{x}_e s.t. $f(\mathbf{x}_e, \mathbf{u}(\mathbf{x}_e)) = 0$ is

- Stable if $\forall \epsilon > 0, \exists \delta > 0, \text{ s.t. } \|\mathbf{x}(0) - \mathbf{x}_e\| < \delta \rightarrow \forall t \geq 0 \|\mathbf{x}(t) - \mathbf{x}_e\| < \epsilon$

- Asymptotic stable if ...

$$\lim_{t \rightarrow \infty} \|\mathbf{x}(t) - \mathbf{x}_e\| \rightarrow 0$$

- Globally asymptotic stable (G. A. S.) if

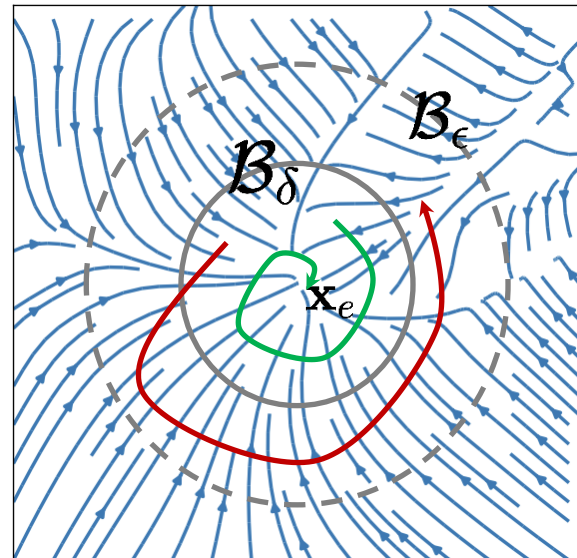
$$\forall \mathbf{x}(0) \rightarrow \lim_{t \rightarrow \infty} \|\mathbf{x}(t) - \mathbf{x}_e\| \rightarrow 0$$

Lyapunov Stability

For $\mathbf{u}(\mathbf{x})$ that steers trajectory of \mathbf{x} under $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$

An equilibrium \mathbf{x}_e s.t. $f(\mathbf{x}_e, \mathbf{u}(\mathbf{x}_e)) = 0$ is

- Stable
- Asymptotic stable
- Globally asymptotic stable (G. A. S.)



How to Certify That

If we can find a scalar function bounded on all sublevel sets and $V(\mathbf{x}) > 0 \quad \forall \mathbf{x} \neq \mathbf{x}_e$

An equilibrium \mathbf{x}_e s.t. $f(\mathbf{x}_e, \mathbf{u}(\mathbf{x}_e)) = 0$ is

- **Stable** if $\dot{V}(\mathbf{x}) = \nabla V(\mathbf{x}) \cdot f(\mathbf{x}, \mathbf{u}(\mathbf{x})) \leq 0 \quad \forall \mathbf{x} \in \mathcal{B}_e$
- **Asymptotic stable** $\dot{V}(\mathbf{x}) = \nabla V(\mathbf{x}) \cdot f(\mathbf{x}, \mathbf{u}(\mathbf{x})) < 0 \quad \forall \mathbf{x} \in \mathcal{B}_e, \mathbf{x} \neq \mathbf{x}_e, \dot{V}(\mathbf{x}_e) = 0$
- Globally asymptotic stable (G. A. S.)

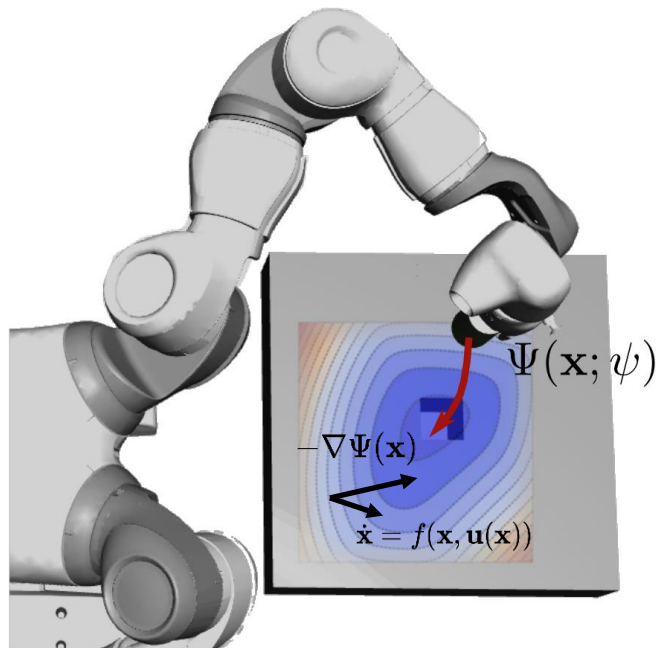
$$\dot{V}(\mathbf{x}) = \nabla V(\mathbf{x}) \cdot f(\mathbf{x}, \mathbf{u}(\mathbf{x})) < 0 \quad \forall \mathbf{x} \neq \mathbf{x}_e, \dot{V}(\mathbf{x}_e) = 0$$

How to Certify That

If we can find a scalar function bounded on all sublevel sets and $V(\mathbf{x}) > 0 \quad \forall \mathbf{x} \neq \mathbf{x}_e$

An equilibrium \mathbf{x}_e s.t. $f(\mathbf{x}_e, \mathbf{u}(\mathbf{x}_e)) = 0$ is

- Stable
- Asymptotic stable
- Globally asymptotic stable (G. A. S.)



Are General (Robotic) Controllers Stable

$$\mathbf{M}(q)\ddot{q} + \mathbf{b}(q, \dot{q}) + \mathbf{g}(q) = \mathbf{u} + \mathbf{J}_c(q)^\top \mathbf{f}_c$$

$$\mathbf{u} = -\mathbf{K}(\mathbf{x} - \mathbf{x}_d) \quad \mathbf{K} \succ 0$$



$$\mathbf{u}(t) = -\mathbf{K}(t)(\mathbf{x}(t) - \mathbf{x}_d(t))$$



$$\mathbf{u} = -\mathbf{K}(\mathbf{x})\mathbf{x}$$



$$\mathbf{u}^* = \arg \min_{\mathbf{u}} [l(\mathbf{x}, \mathbf{u}) + \frac{\partial V^*(\mathbf{x})}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u})]$$



$$0 = \min_{\mathbf{u}} [l(\mathbf{x}, \mathbf{u}) + \frac{\partial V^*(\mathbf{x})}{\partial \mathbf{x}} \cdot f(\mathbf{x}, \mathbf{u})] \quad l(\mathbf{x}, \mathbf{u}) > 0$$

$$\dot{V}(\mathbf{x})$$

$$\mathbf{u} = -\frac{\partial}{\partial q}V(q) - \mathbf{D}(\dot{\mathbf{q}})\dot{q}$$

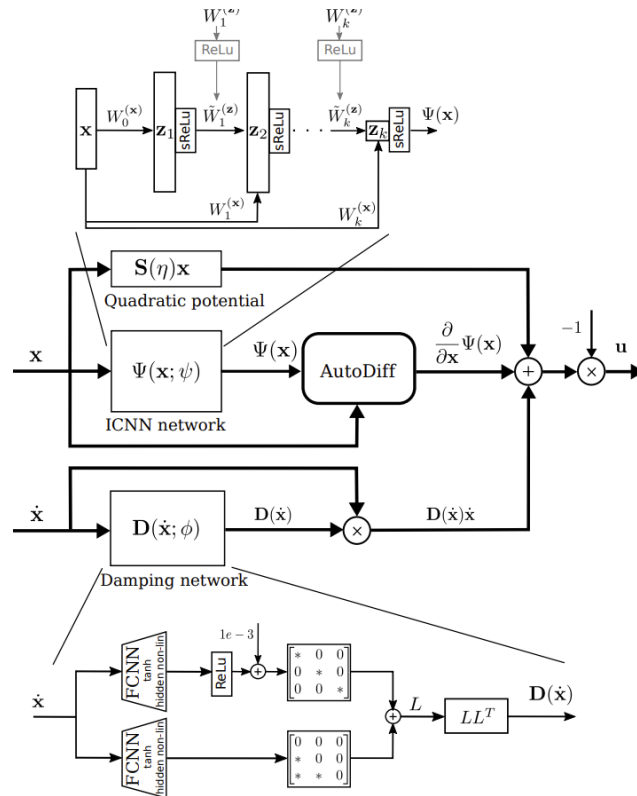
at q^* , q^* is G. A. S.

Policy biases a goal-directed behaviour

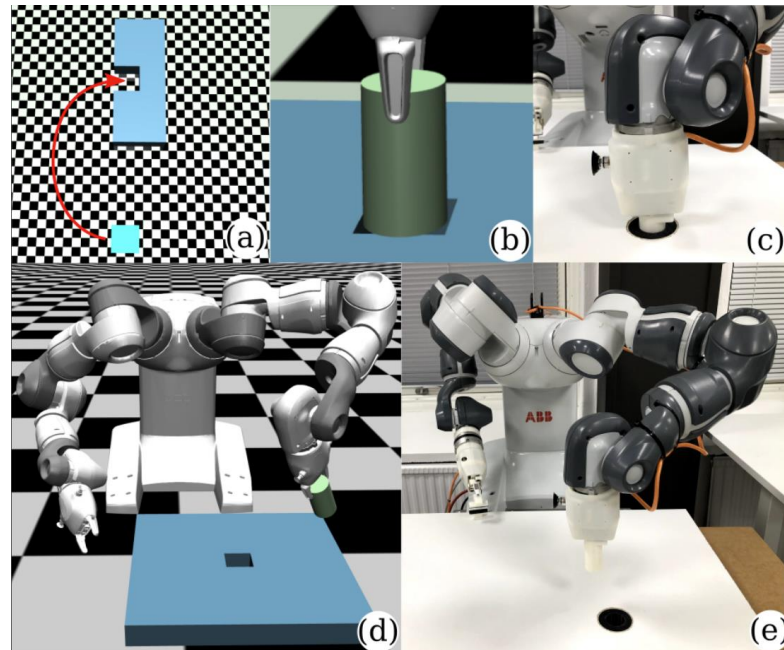
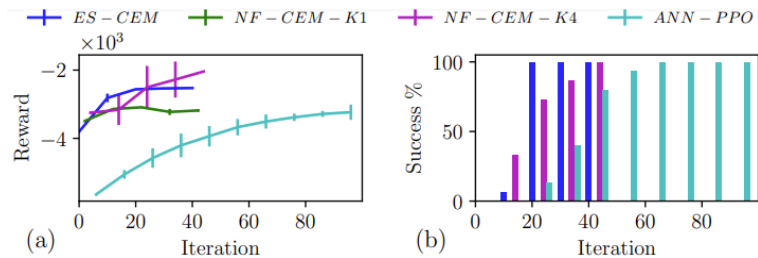
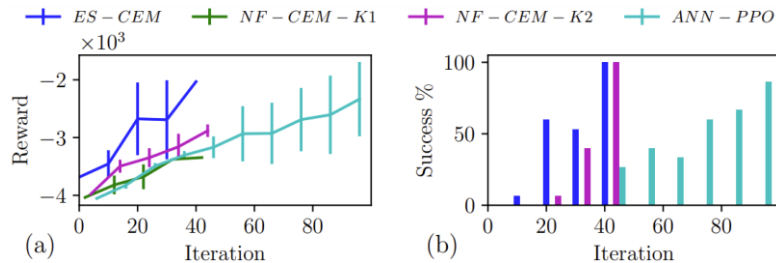


Parameterize policy with Energy-Shaping Controllers

- How can we get a convex (neural network) energy function
- Input Convex Neural Network (ICNN, Amos et al, 2017)
 - Additive quadratic for strong convexity
 - Zero bias for unique minimum at origin



How Does It Work



How Does It Work

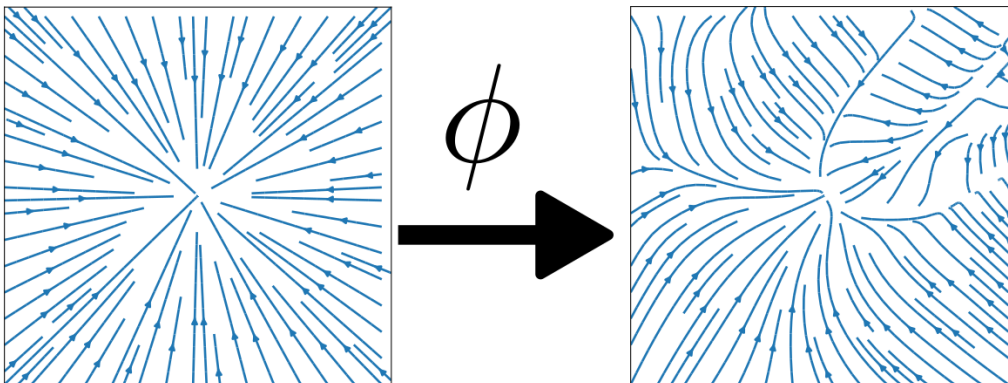
Reinforcement learning of the 2D block-insertion task

Policy: The proposed energy shaping policy

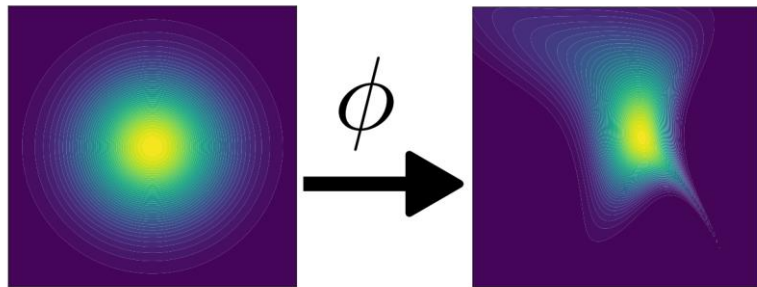
Algorithm: Cross Entropy Method

Another Way to Gain Stability Guarantees

- Lyapunov analysis is usually easy for linear control & systems
 - LQR – Quadratic function $V^*(\mathbf{x}) = \mathbf{x}^\top \mathbf{P} \mathbf{x}$
- Can we construct policy parameterization by “non-linearizing” a stable linear control/system, with guarantees reserved?



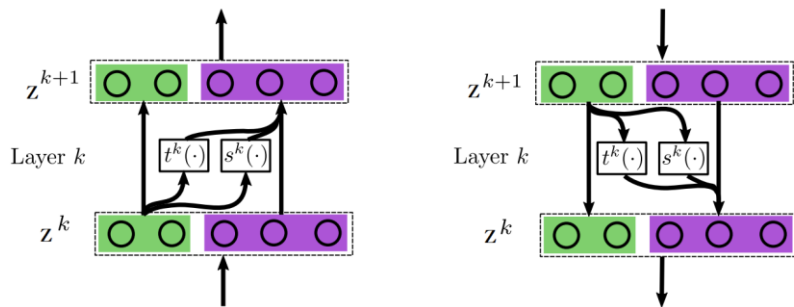
Diffeomorphism via Normalizing-Flow Models



- Construct non-Gaussian models from Gaussians

$$p(\mathbf{x}) = |\mathbf{J}_\phi(\mathbf{z})|^{-1} p(\mathbf{z}|\mu, \Sigma)$$

- Differentiable and invertible neural networks – Real NVP (Dinh et al , 2016)



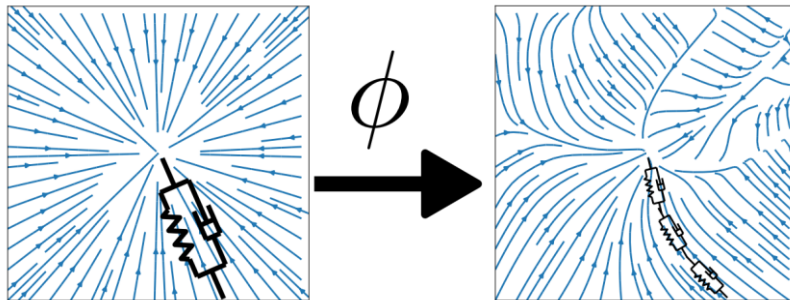
- Affine composition

$$\mathbf{z}^{k+1} = \phi^k(\mathbf{z}^k)$$

$$\mathbf{z}_{1:d^k}^{k+1} = \mathbf{z}_{1:d^k}^k$$

$$\mathbf{z}_{d^k+1:d}^{k+1} = \mathbf{z}_{d^k+1:d}^k \odot \exp(\mathbf{s}^k(\mathbf{z}_{1:d^k}^k)) + \mathbf{t}^k(\mathbf{z}_{1:d^k}^k)$$

Constructing Stable Normalizing-Flow Controller



$$\mathbf{u} = -\mathbf{S}\mathbf{x} - \mathbf{D}\dot{\mathbf{x}}$$



$$\mathbf{u} = \mathbf{g}(\mathbf{x}) - \mathbf{D}\dot{\mathbf{x}} - \mathbf{J}_\phi(\mathbf{x})^\top \mathbf{S}[\phi(\mathbf{x}) - \phi(\mathbf{x}_{\text{ref}})]$$

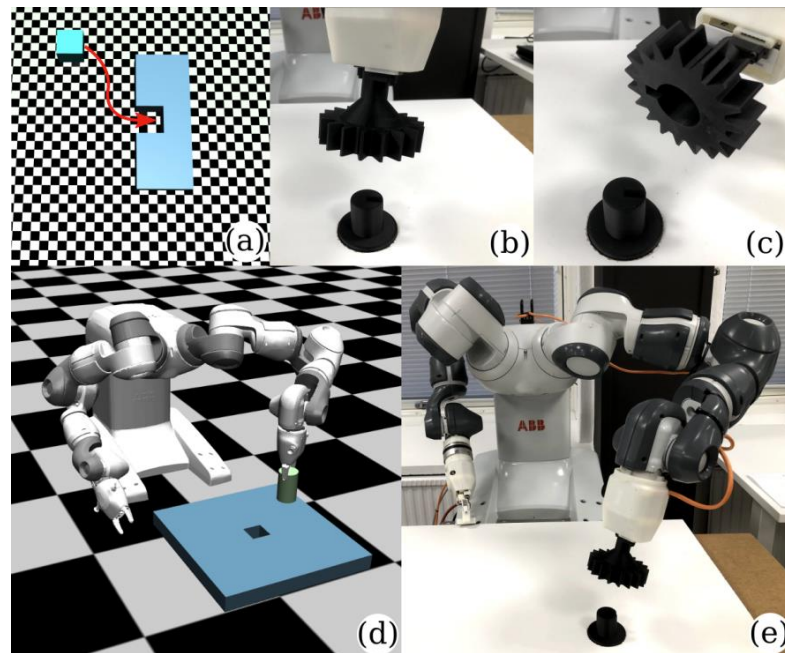
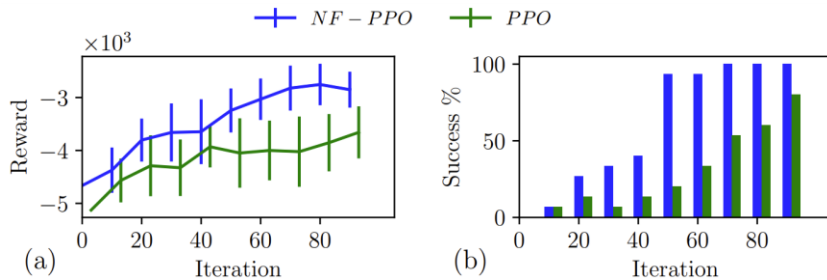
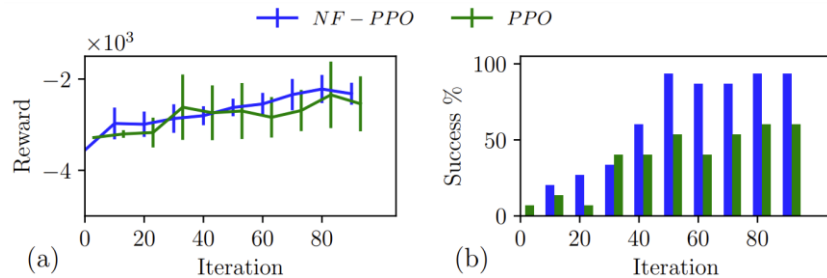
$$\mathbf{u} = -\mathbf{K}(\mathbf{x})\mathbf{x}$$



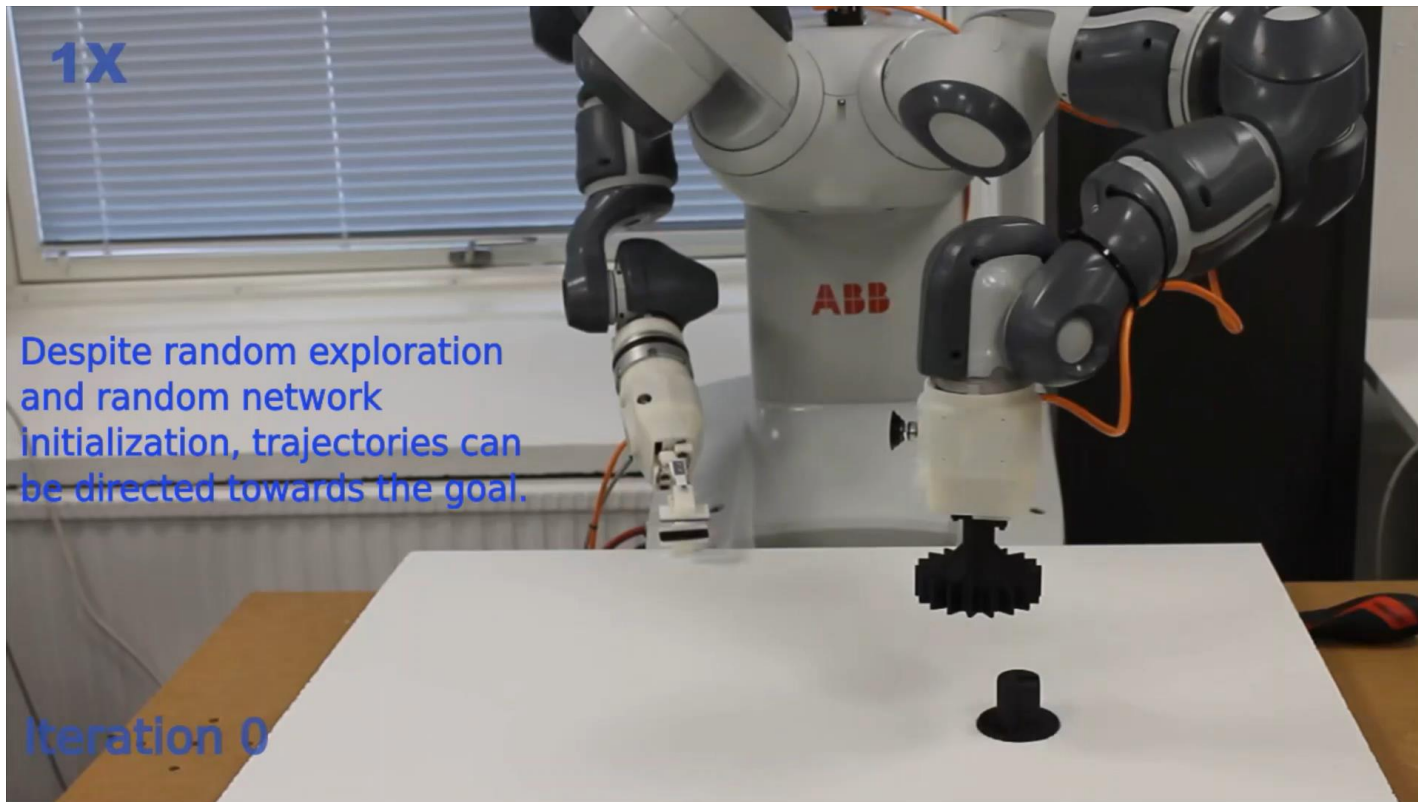
- Neural network parameterized
- G. A. S. of \mathbf{x}_{ref} under

$$\mathbf{M}(q)\ddot{q} + \mathbf{b}(q, \dot{q}) + \mathbf{g}(q) = \mathbf{u} + \mathbf{J}_c(q)^\top \mathbf{f}_c$$
- Stable variable impedance controller comparing to Martin-Martin, 2019

How Does It Work

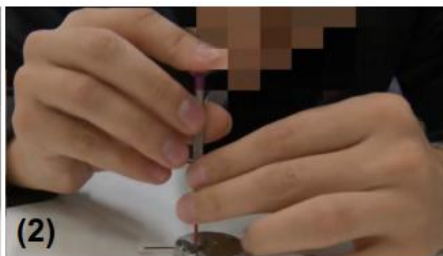
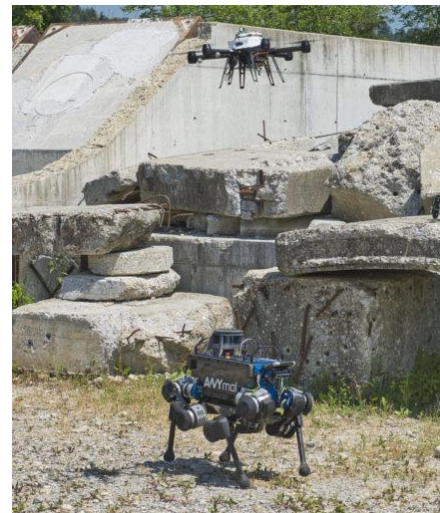


How Does It Work



From Single Robot to Multiple Robots

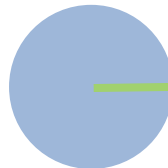
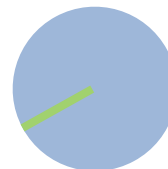
- More dexterity and capacity
 - Heterogenous robot teams
 - Anthropomorphic robots
- Larger and diverse action/policy space





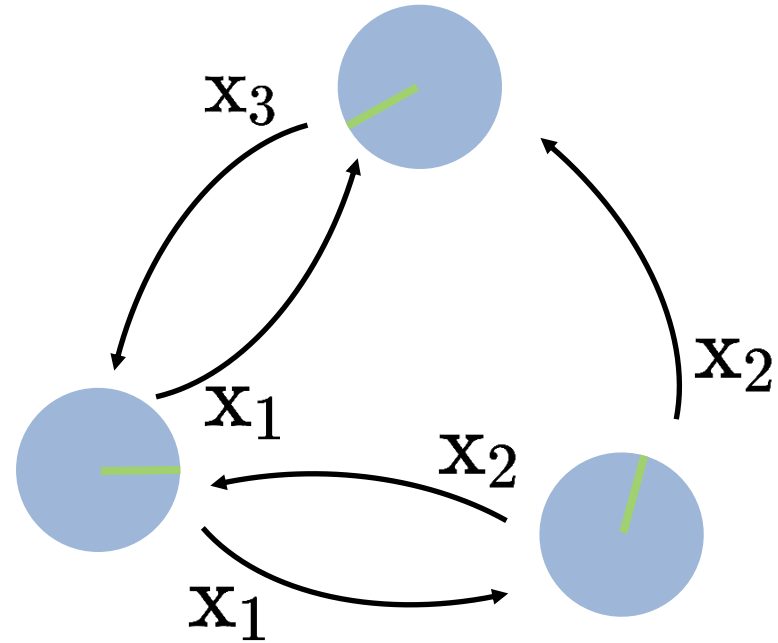
Consensus in Multi-Robot Systems

- N robots



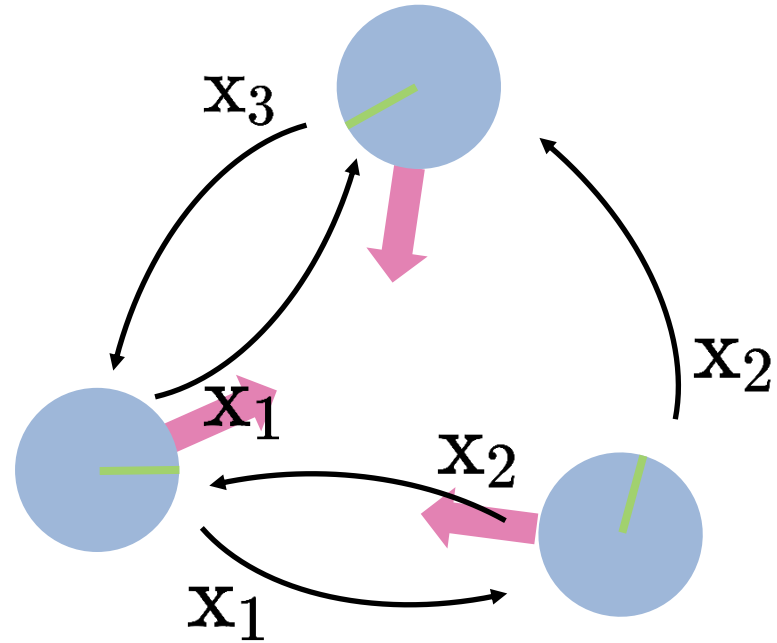
Consensus in Multi-Robot Systems

- N robots
 - Communicate with neighbouring robots



Consensus in Multi-Robot Systems

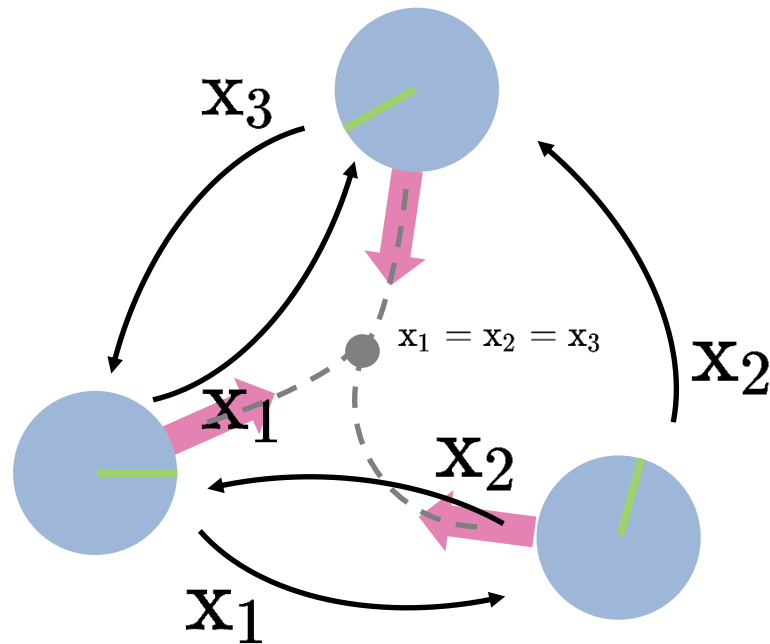
- N robots
 - Communicate with neighbouring robots
 - Apply distributed control



Consensus in Multi-Robot Systems

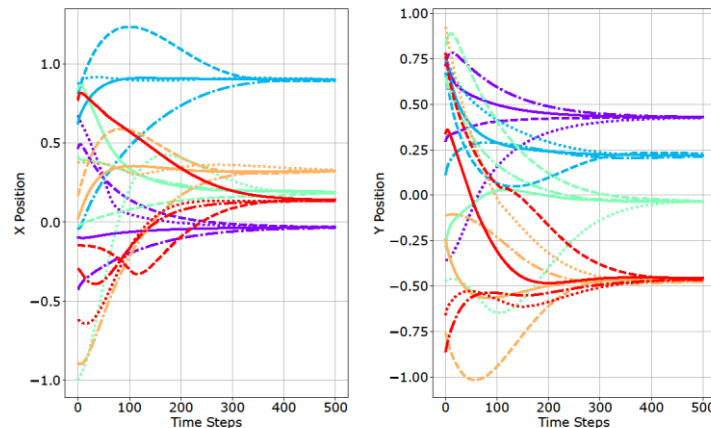
- N robots
 - Communicate with neighbouring robots
 - Apply distributed control
 - Reach agreement on communicated information

e.g. $\dot{\mathbf{x}} = \mathbf{L}\mathbf{x}$



Consensus-based Normalizing-Flow Control

- Nonlinear transient behavior
 - Allowing to learn complex rendez-vous trajectories
- Apply normalizing-flow to consensus problem: “non-linearize” basic protocol
 - N robots with 2nd-order dynamics
 - Global consensus stabilization



$$\bar{\mathbf{u}} = \mathbf{g}(\mathbf{x}) - \mathbf{D}\dot{\mathbf{x}} - \bar{\mathbf{J}}(\mathbf{x})^\top (\mathbf{L} \otimes I_n) \bar{\boldsymbol{\phi}}(\mathbf{x})$$

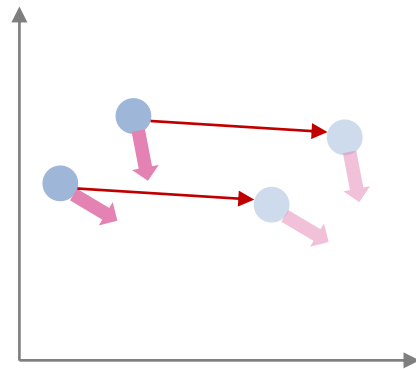
$$\mathbf{x} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top \in \mathbb{R}^{nd}$$

Symmetry in Consensus

- Translation-invariant
 - Control remains when all robots are translated by a same vector
 - Keep a consistent behavior across entire workspace

$$\dot{\mathbf{x}} = \mathbf{L}\mathbf{x}$$

$$\bar{\mathbf{u}} = \mathbf{g}(\mathbf{x}) - \mathbf{D}\dot{\mathbf{x}} - \bar{\mathbf{J}}(\mathbf{x})^\top (\mathbf{L} \otimes I_n) \bar{\phi}(\mathbf{x})$$



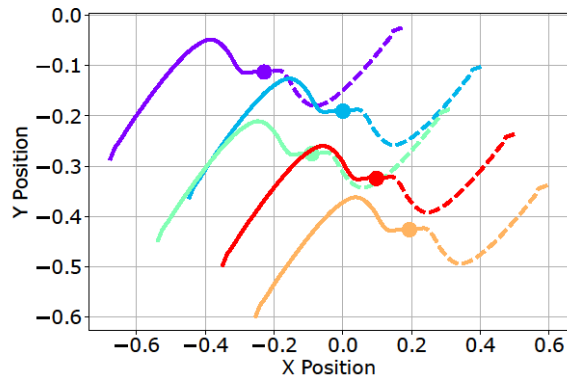
Translation-invariant Normalizing-Flow Control

- N=2 robots, e.g. dual-arm
- Global consensus stabilization

$$\tau = \mathbf{g}(\mathbf{x}) - \mathbf{D}\dot{\mathbf{x}} - \bar{\mathbf{L}}\bar{\mathbf{J}}(\bar{\mathbf{L}}\mathbf{x})^\top \bar{\mathbf{S}}\bar{\phi}(\bar{\mathbf{L}}\mathbf{x})$$

- Need antipodal-equivariant normalizing-flow

$$\phi(*) = -\phi(-*)$$

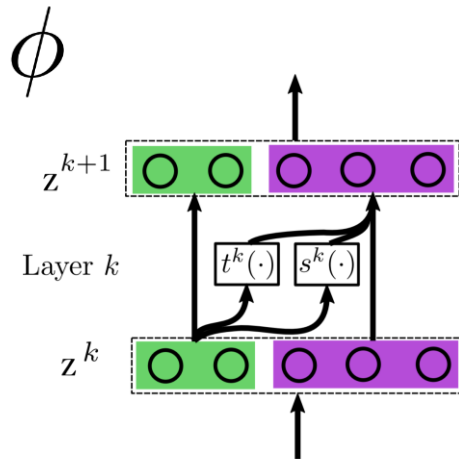


Antipodal-Equivariant Normalizing-Flows

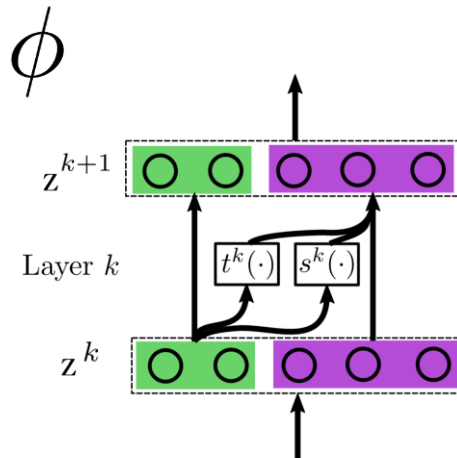
$$\mathbf{z}^{k+1} = \phi^k(\mathbf{z}^k)$$

$$\mathbf{z}_{1:d^k}^{k+1} = \mathbf{z}_{1:d^k}^k$$

$$\mathbf{z}_{d^k+1:d}^{k+1} = \mathbf{z}_{d^k+1:d}^k \odot \exp(\mathbf{s}^k(\mathbf{z}_{1:d^k}^k)) + \mathbf{t}^k(\mathbf{z}_{1:d^k}^k)$$



Need to be invariant



Antipodal-Equivariant Normalizing-Flows

$$\begin{aligned} \mathbf{z}^{k+1} &= \phi^k(\mathbf{z}^k) \\ \mathbf{z}_{1:d^k}^{k+1} &= \mathbf{z}_{1:d^k}^k \\ \mathbf{z}_{d^k+1:d}^{k+1} &= \mathbf{z}_{d^k+1:d}^k \odot \exp[\mathbf{s}^k(\mathbf{z}_{1:d^k}^k)] + \mathbf{t}^k(\mathbf{z}_{1:d^k}^k) \end{aligned}$$



Equivariant



Need to be equivariant



Need to be invariant

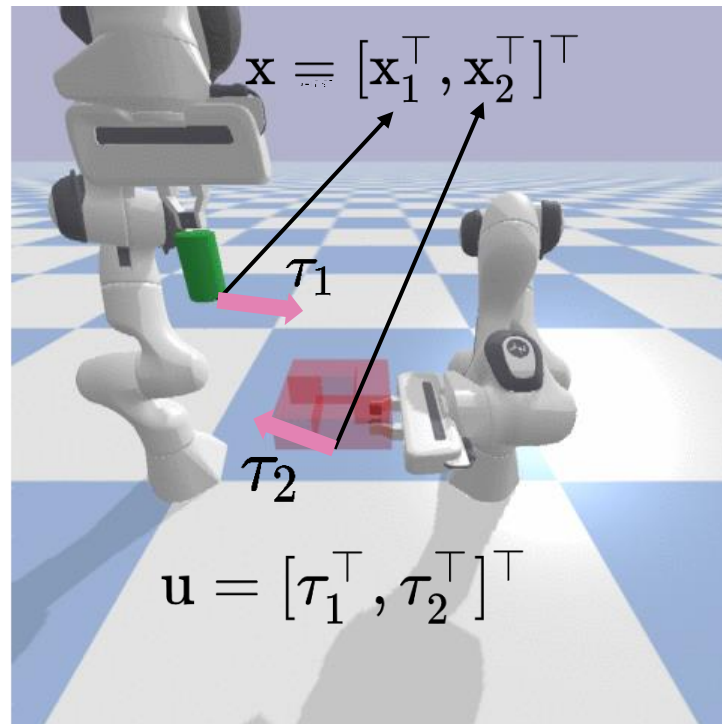
Construction with even and odd functions

$$\tilde{\mathbf{s}}_k(*) = \mathbf{s}_k(*) + \mathbf{s}_k(-*)$$

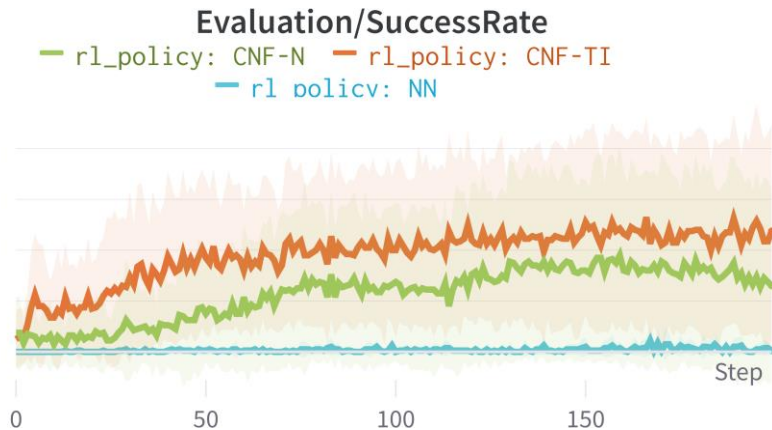
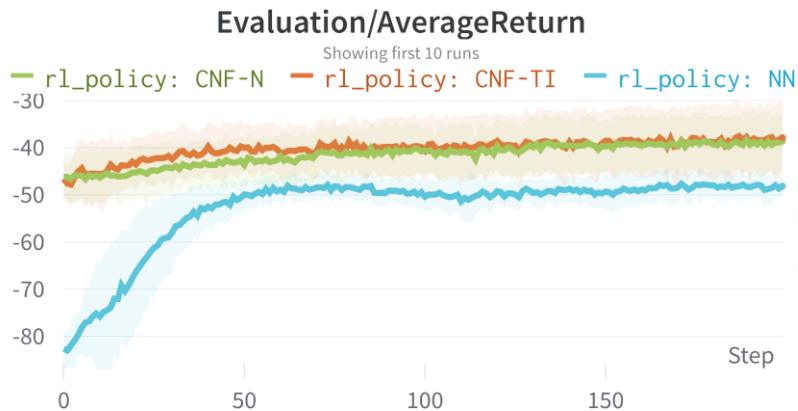
$$\tilde{\mathbf{t}}_k(*) = \mathbf{t}_k(*) - \mathbf{t}_k(-*)$$

How Does It Work

- Bimanual insertion with a clearance of 2mm
 - Observation/Action: translational position and force
 - Fixed stiffness controller for end-effector orientations
- PPO
 - Stochastic policy with controller adding Gaussian noise
 - 1e-4 learning rate
 - 10⁶ environment steps
 - Randomized initial positions
10cmx10cm area
 - 10 random seeds

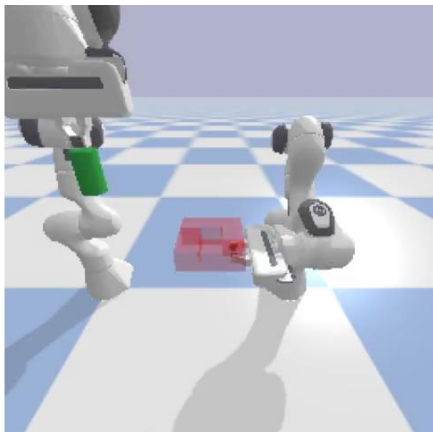


How Does It Work- Training Curves

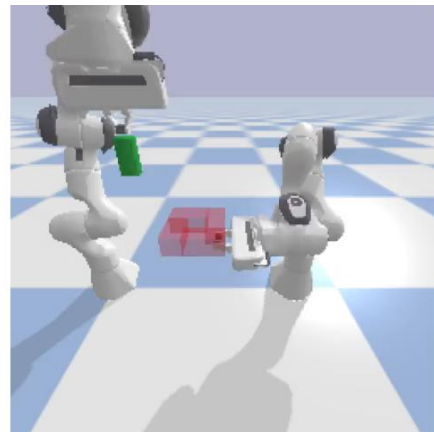
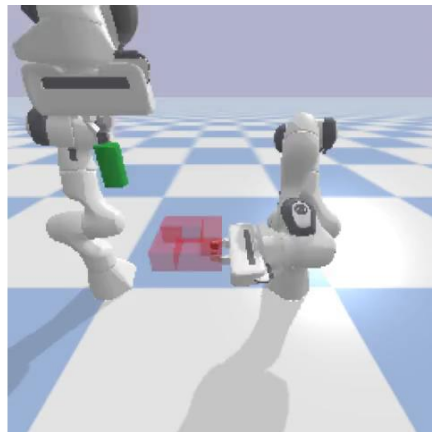


How Does It Work - Results

Basic consensus-based policy (CNF-N)

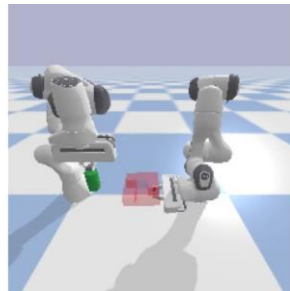
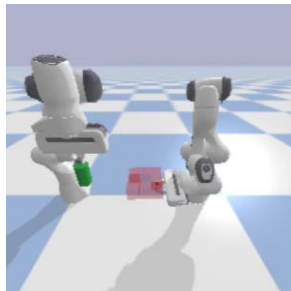
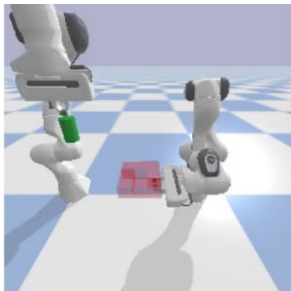
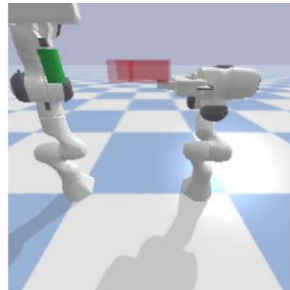
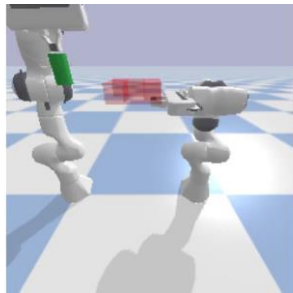
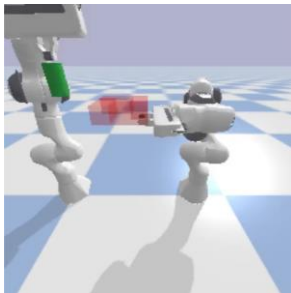


Translation-invariant policy (CNF-TI)

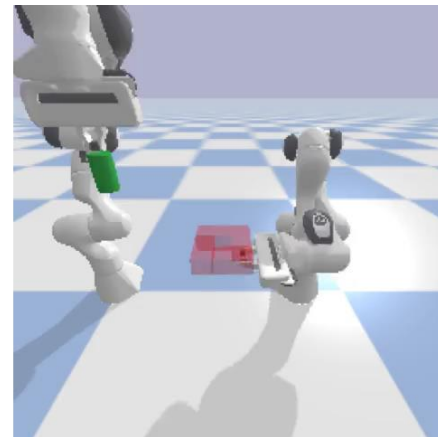


How Does It Work - Generalization and Failure Case

Generalize to novel positions



Failure case





Closing Notes

- RL (real) data efficiency gained from prior knowledge
 - In form of synthetic simulation data: zero/few-shot sim-to-real
 - In form of structured policy design
- Choosing “right” policy
 - Optimal control perspectives
 - Structure from robotics literature
 - Consideration on Guarantees
 - Train on real robots
- Design RL system with provable performance guarantees
 - Computational tractability
 - Efficiency from regulated searching space
 - Theory-grounded and certified algorithms



Closing Notes

- Connection to Chris and Alexis lectures
 - RL exploration with provable properties
 - Stability does not necessarily mean safety
 - Mostly in robot control application: continuous, steady, constraints
 - Need more powerful description of desired system behaviours



Supplemental Materials

- Benjamin Recht – An outsider's tour of reinforcement learning
 - <http://www.argmin.net/2018/06/25/outsider-rl/>
- Underactuated Robotics
 - Course notes (<https://underactuated.mit.edu/>)
- Dimitri Bertsekas
 - Reinforcement Learning and Optimal Control (<http://web.mit.edu/dimitrib/www/RLbook.html>)
 - Lessons from AlphaZero for Optimal, Model Predictive, and Adaptive Control (http://web.mit.edu/dimitrib/www/abstractdp_MIT.html)