

Micro-Data Reinforcement Learning for Adaptive Robots

Konstantinos Chatzilygeroudis

Adjunct Lecturer, CEID, University of Patras
& Computer Vision Team Leader, Metargus

March 26, 2021



About the speaker

- CEID Graduate: 2009-2014
- PhD in Machine Learning/Robotics (INRIA): 2015-2018
- Post-doc in Robotics (EPFL): 2018-2020
- CEID Adjunct Lecturer: Oct 2020-now
- Metargus: Jan 2021-now



CEID
COMPUTER ENGINEERING & INFORMATICS DEPARTMENT

metargus

Robotics (1)



Robotics (1)



Robotics (2)



¹Boston Dynamics, 2020

²<https://shorturl.at/oryAE>

Robotics - Limitations



¹DARPA Robotics Challenge, 2015

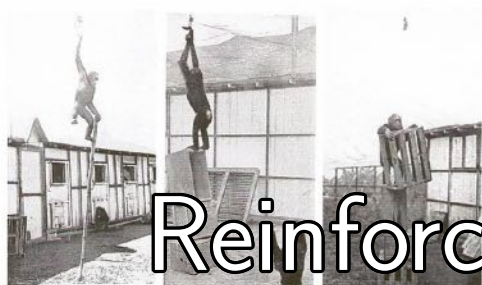
Animal Adaptation



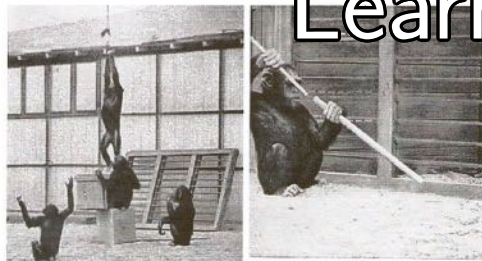
And showing how life goes on without much change!

Trial and Error Learning





Reinforcement Learning



Problem formulation

Reinforcement Learning (RL) for Robotics

We consider dynamical systems of the form:

$$\mathbf{x}_{t+1} = F(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{w} \quad (1)$$

with $\mathbf{x} \in \mathbb{R}^E$, $\mathbf{u} \in \mathbb{R}^F$, i.i.d. Gaussian system noise \mathbf{w} , and unknown transition dynamics F .

We seek to find a *policy* π , $\mathbf{u} = \pi(\mathbf{x}|\boldsymbol{\theta})$, which maximizes the *expected long-term reward* in **as little interaction time as possible** (i.e., we want a **data-efficient** algorithm):

$$J(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=1}^T r(\mathbf{x}_t) \middle| \boldsymbol{\theta} \right] \quad (2)$$

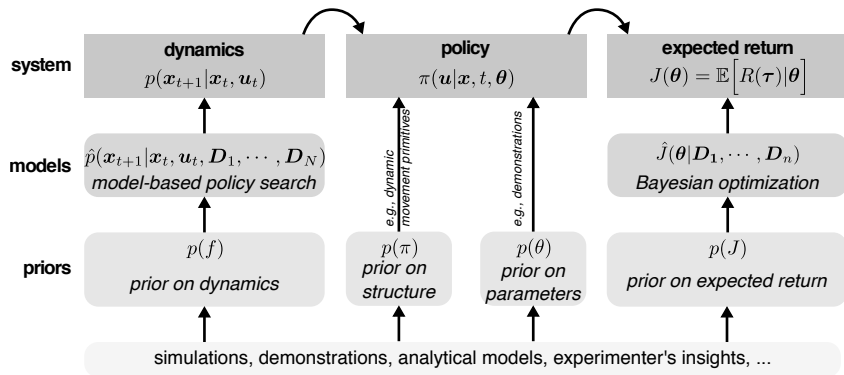
where $r(\mathbf{x}_t)$ is the immediate reward of being in state \mathbf{x} at time t .

Micro-Data Reinforcement Learning

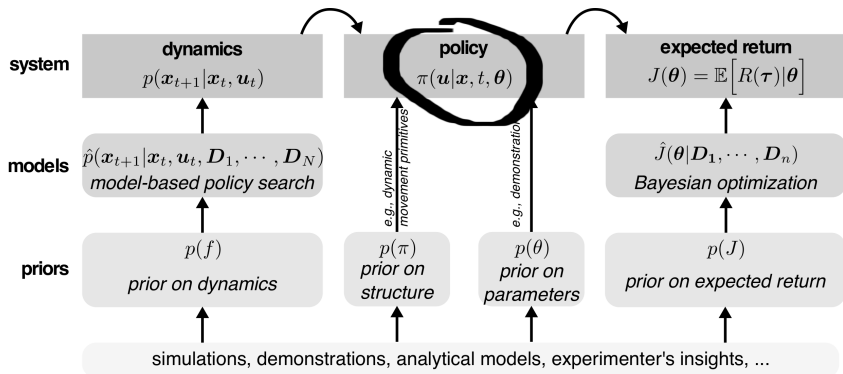


Only a few trials possible!

Strategies of Micro-Data Reinforcement Learning



Micro-Data Reinforcement Learning: No prior



State-of-the-art: Policy Gradient Algorithms

Policy search:

$$\theta^* = \underset{\theta}{\operatorname{argmax}} J(\theta) \quad (3)$$

Stochastic Policy Gradients¹:

$$\nabla_{\theta} J(\theta) = \mathbb{E} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(\mathbf{u}_t | \mathbf{x}_t, \theta) A_t \right] \quad (4)$$

where $A_t = \hat{Q}^{\pi}(\mathbf{x}_t, \mathbf{u}_t)$.

Big variance in the gradient estimation, and thus slow convergence (or even divergence)!

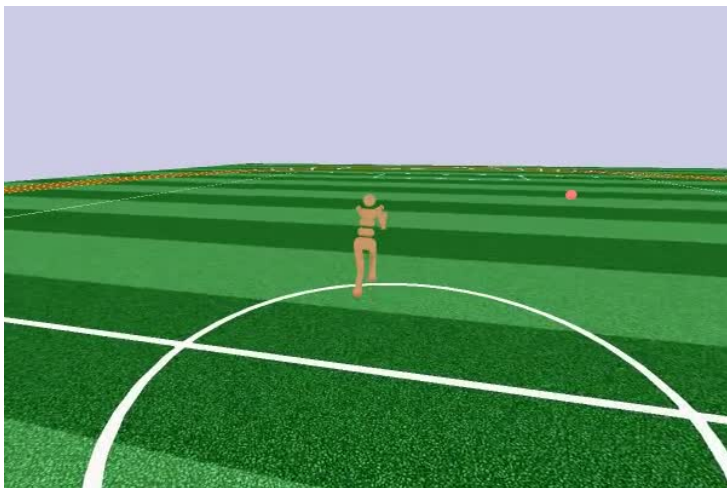
Trust Region Policy Optimization (TRPO)² and Proximal Policy Optimization (PPO)³ use an extra constraint to reduce the variance and provide monotonic improvement guarantees.

¹ Sutton, R., et al. "Policy gradient methods for reinforcement learning with function approximation", NIPS, 2000

² Schulman, J., et al. "Trust region policy optimization", ICML, 2015

³ Schulman, J., et al. "Proximal policy optimization algorithms", 2017

State-of-the-art: PPO Results

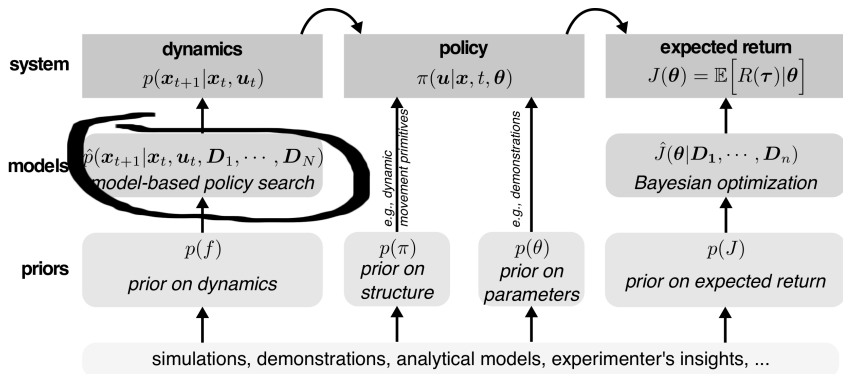


More than 80 hours of simulated training!

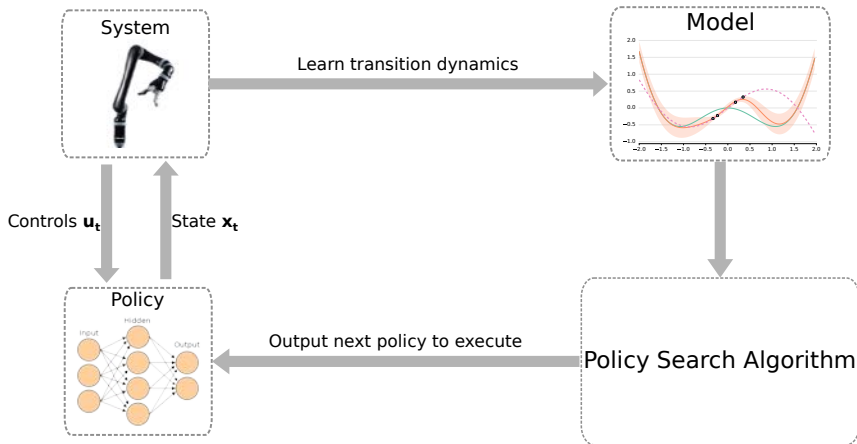
¹Schulman, J., et al. "Proximal policy optimization algorithms", 2017

²<https://blog.openai.com/openai-baselines-ppo/>

Micro-Data Reinforcement Learning: Model-based

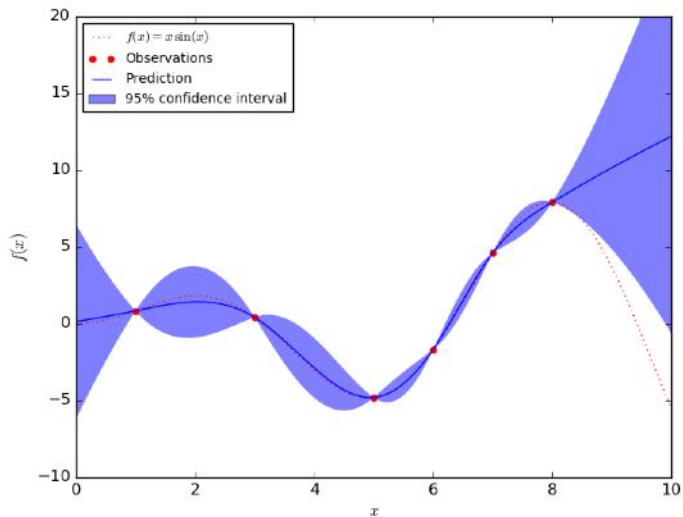


Model-based policy search



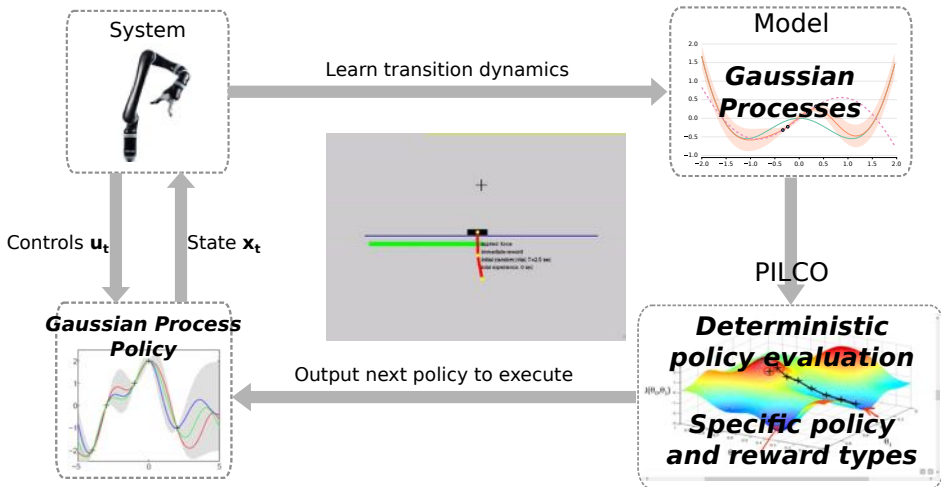
¹Deisenroth, M.P., Neumann, G., Peters, J. "A survey on policy search for robotics", Foundations and Trends in Robotics, 2013

Gaussian Processes



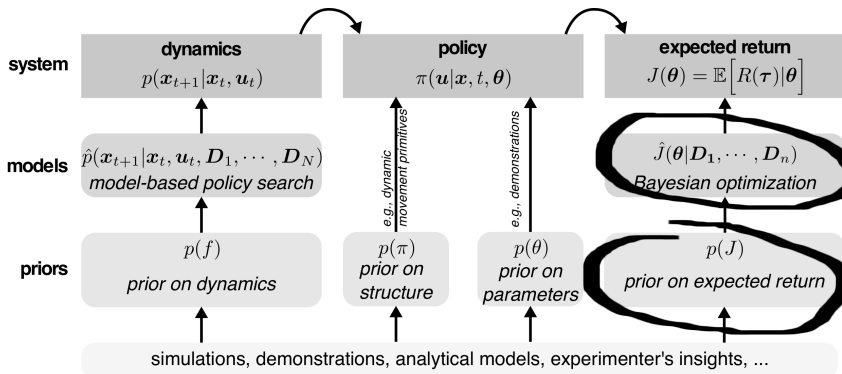
¹Rasmussen, CE. "Gaussian processes in machine learning", 2004

State-of-the-art: PILCO¹

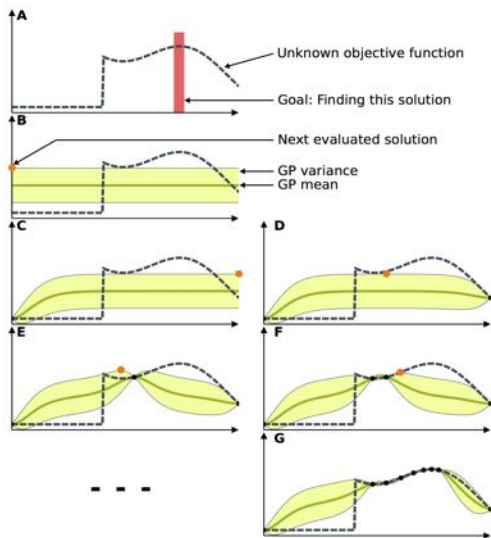


¹Deisenroth, D., et al. "Gaussian Processes for Data-Efficient Learning in Robotics and Control" IEEE Transactions on Pattern Analysis and Machine Intelligence, 2014

Micro-Data Reinforcement Learning: Surrogate Models



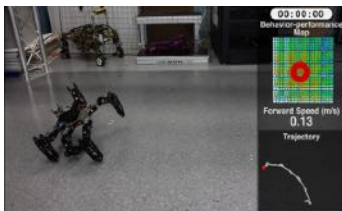
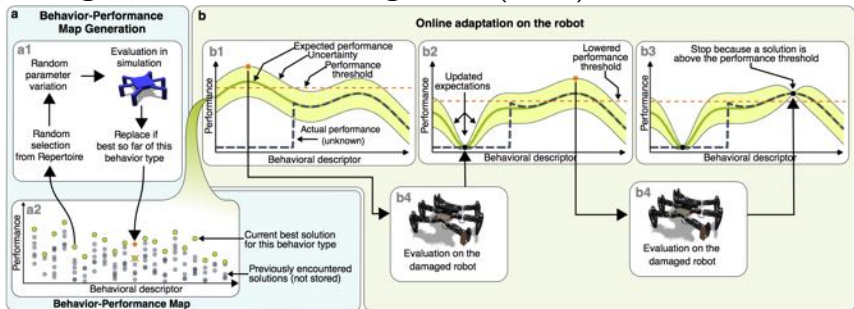
State-of-the-art: Bayesian Optimization



¹Brochu, E., Cora, V.M., De Freitas, N. "A tutorial on Bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning", 2010

State-of-the-art: Learning for Damage Recovery

Intelligent Trial and Error Algorithm (IT&E)¹



¹Cully, A. et al. "Robots that can adapt like animals", in Nature, vol. 521, no. 7553, pp. 503–507, 2015

State-of-the-art: Quadratic Programming-based Control

But we can have analytical models:

$$\begin{aligned}M(\mathbf{q})\ddot{\mathbf{q}} + \mathbf{C}_g(\mathbf{q}, \dot{\mathbf{q}}) &= \mathbf{S}\boldsymbol{\tau} + \mathbf{J}^T(\mathbf{q})\mathbf{W} \\ \mathbf{J}(\mathbf{q})\ddot{\mathbf{q}} + \dot{\mathbf{J}}(\mathbf{q}, \dot{\mathbf{q}})\dot{\mathbf{q}} &= \ddot{\mathbf{x}}\end{aligned}\quad (5)$$

And solve an optimization:

$$\begin{aligned}\min_{\boldsymbol{\chi}} \quad & -\frac{1}{2}\boldsymbol{\chi}^T \mathbf{G}\boldsymbol{\chi} + \mathbf{g}^T \boldsymbol{\chi} \\ \text{s.t.} \quad & \begin{bmatrix} \mathbf{M}(\mathbf{q}) & -\mathbf{S} & -\mathbf{J}(\mathbf{q})^T \end{bmatrix} \boldsymbol{\chi} + \mathbf{C}_g(\mathbf{q}, \dot{\mathbf{q}}) = 0\end{aligned}\quad (6)$$

where

$$\boldsymbol{\chi} = \begin{bmatrix} \ddot{\mathbf{q}} & \boldsymbol{\tau} & \mathbf{W} \end{bmatrix}^T \quad (7)$$

State-of-the-art: QP-based Control (2)



¹LARSEN Inria: https://www.youtube.com/watch?v=-An7Ju3ge0I&ab_channel=LarsenInria

State-of-the-art: Limitations

Learning Methods

Pure episodic approach:

- The robot starts in the same initial state at each episode;
- Learning process separated from operation.

Scaling issues with complex robots:

- Exponentially more data are needed as the dimensionality of state/action space increases;
- Data-efficient approaches do not take advantage of multi-core architectures.

Traditional Methods

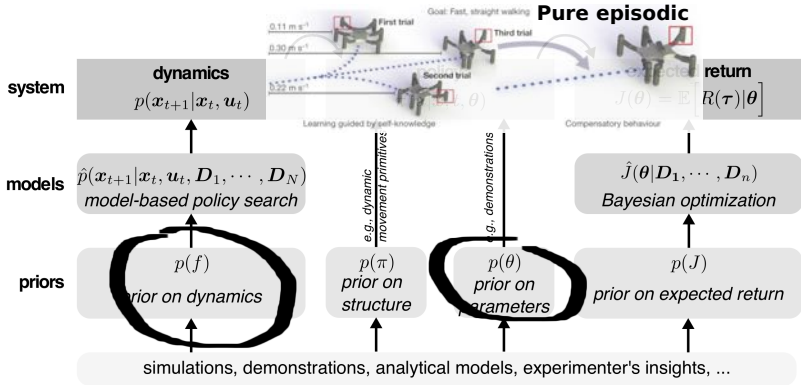
Require a lot of tuning:

- Hard to find hyper-parameters;
- Different parameters for each task.

Accurate models required:

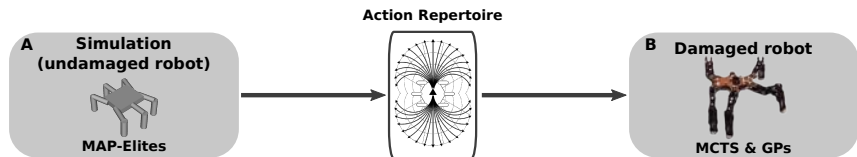
- Model-based methods (e.g., QP-control) do not work with in accurate models;
- Difficult to incorporate learning.

Reset-free Trial and Error Learning for Robot Damage Recovery



¹Chatzilygeroudis, K., Vassiliades, V. and Mouret, J. B. "Reset-free Trial-and-Error Learning for Robot Damage Recovery", RAS, 2018.

Reset-free Trial and Error Learning for Robot Damage Recovery (RTE)



Learning and correcting the repertoire

MAP-Elites¹ uses a **simulated intact** robot and **requires**:

- A parameterized policy, π_{θ}
- An action descriptor, \mathbf{a} , that describes the task space
- A performance measure, $\text{performance}(\theta)$

It provides:

- A diverse set of locally (with respect to \mathbf{a}) optimized policies, A
- A mapping between the task space, (or the set of actions A), and the policy space Θ ; i.e., $A \rightarrow \Theta$
- A mapping between actions and relative outcomes, $M : A \rightarrow O$;

We use Gaussian Processes with a non-zero mean function to correct the repertoire:

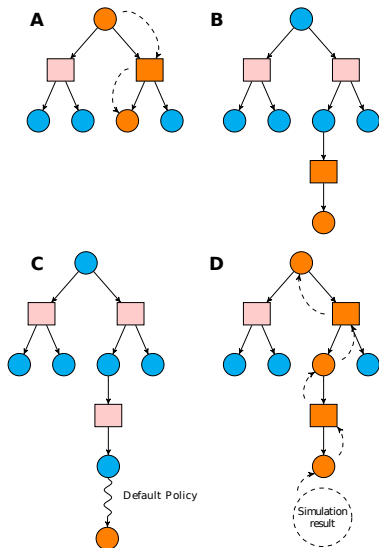
$$\begin{aligned} p(f(\mathbf{a})|D_{1:t}, \mathbf{a}) &\sim \mathcal{N}(\mu(\mathbf{a}), \sigma^2(\mathbf{a})) \\ \mu(\mathbf{a}) &= M(\mathbf{a}) + \mathbf{k}^T K^{-1} (D_{1:t} - M(\mathbf{a}_{1:t})) \\ \sigma^2(\mathbf{a}) &= k(\mathbf{a}, \mathbf{a}) - \mathbf{k}^T K^{-1} \mathbf{k} \end{aligned} \tag{8}$$

¹Mouret, J.-B., and Clune, J. "Illuminating search spaces by mapping elites.", 2015.

Planning with MCTS

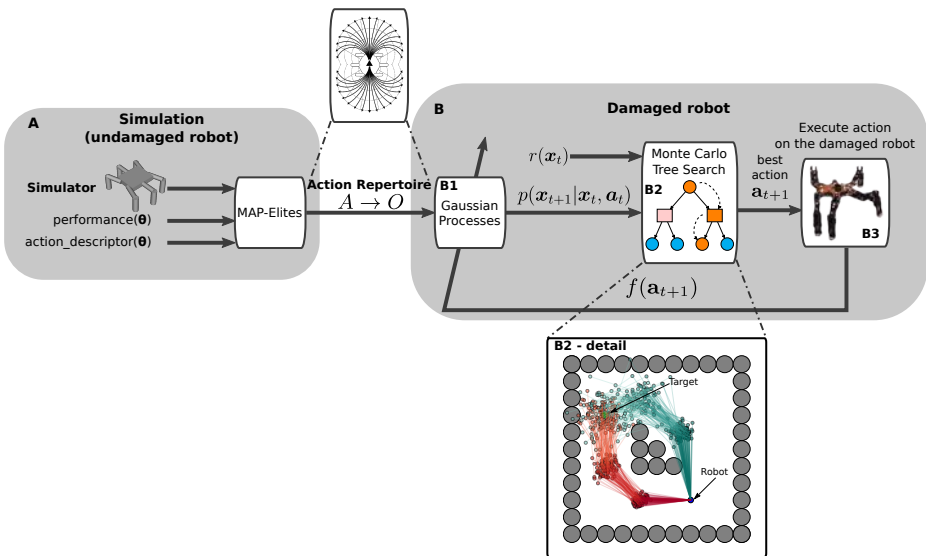
Monte Carlo Tree Search:

- is a planning algorithm that is:
 - best-first, sample-based, and
 - anytime
- finds optimal decisions
 - by taking random samples, and
 - building a search tree according to their results
- treats as a **black-box** the model of the environment (handles uncertainty)
- has successfully solved RL problems with:
 - stochastic transitions,
 - continuous state spaces, and
 - high branching factors



¹<https://github.com/resibots/mcts>

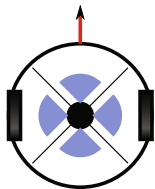
RTE Overview



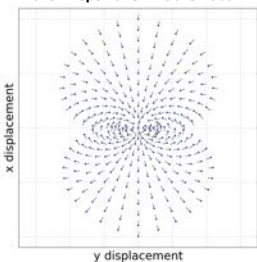
Results: Mobile robot

A velocity-controlled differential drive robot has to learn to navigate again after one of the motors only achieves half of the desired velocity.

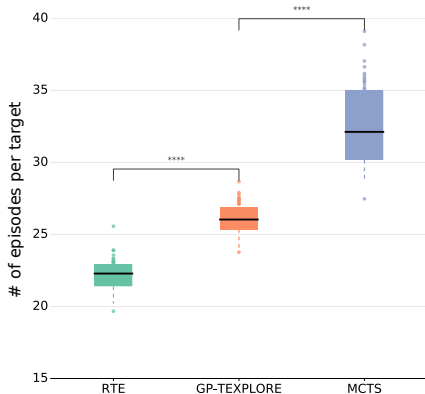
A Mobile Robot



B Action Repertoire - Mobile Robot



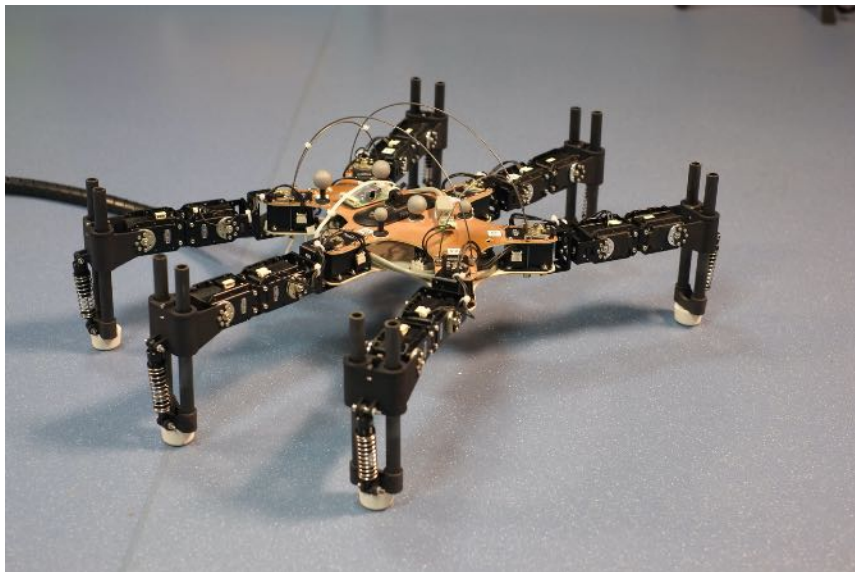
Performance - Mobile Robot



50 replicates

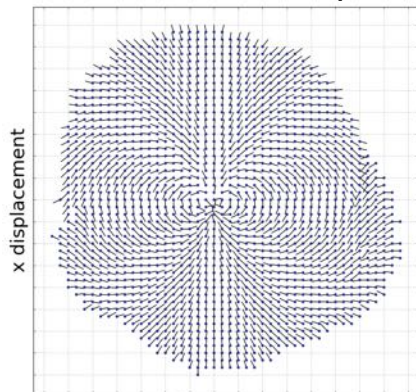
¹Hester, T., Stone, P. "TEXPLORE: real-time sample-efficient reinforcement learning for robots", Machine learning, 2013

Simulated hexapod robot

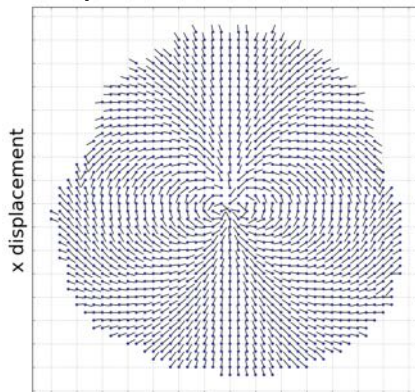


Results - Simulated hexapod robot

Action Repertoires - Hexapod Robot



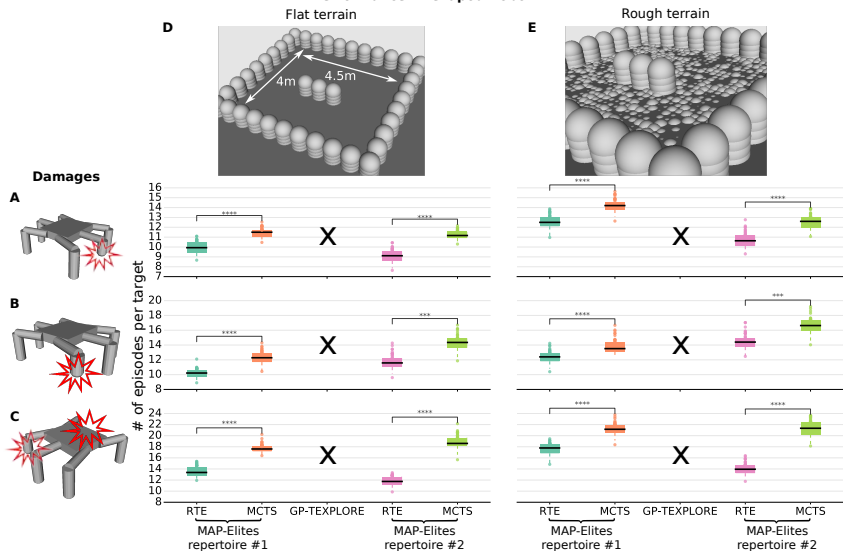
y displacement
Action Repertoire #1



y displacement
Action Repertoire #2

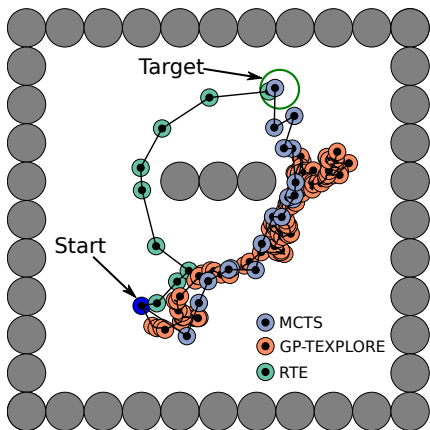
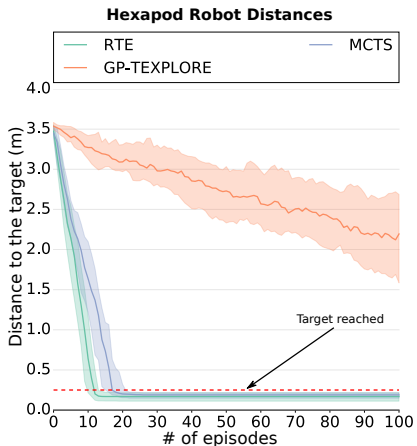
Results - Simulated hexapod robot (2)

Performance - Hexapod Robot

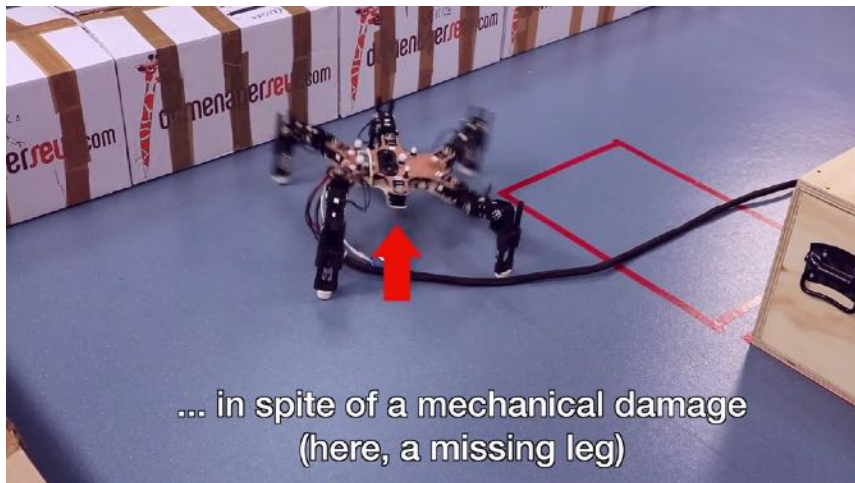


50 replicates

Results: Simulated hexapod robot (3)



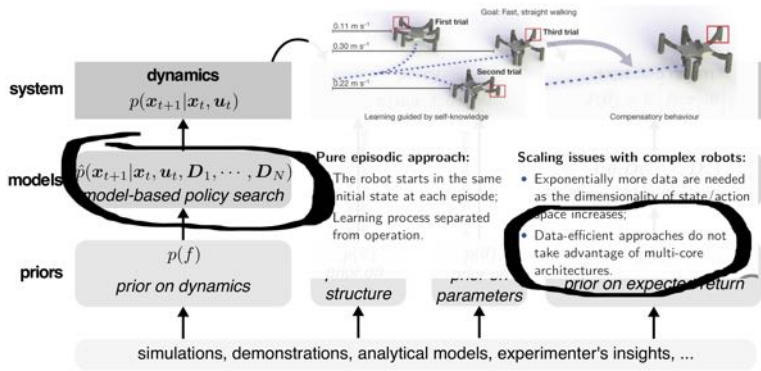
Video - Physical hexapod robot



... in spite of a mechanical damage
(here, a missing leg)

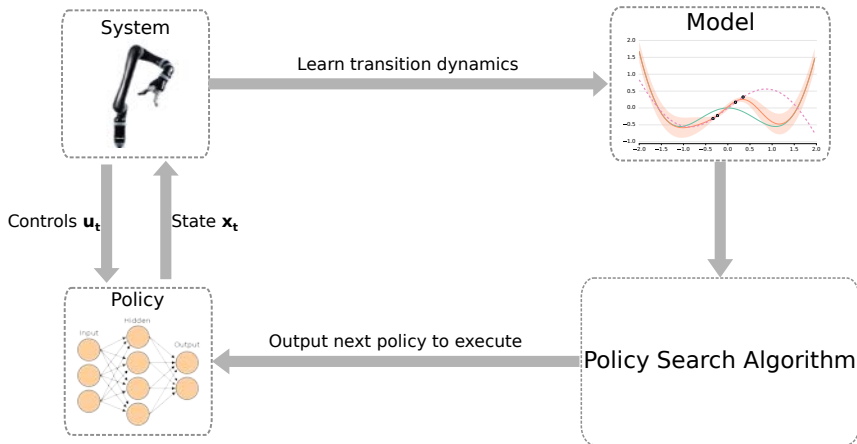
Black-Box Data-efficient Policy Search for Robotics

State-of-the-art: Limitations



¹Chatzilygeroudis, K., Rama, R., Kaushik, R., Goepf, D., Vassiliades, V., and Mouret, J. B. "Black-Box Data-efficient Policy Search for Robotics", IROS, 2017

Model-based policy search - Reminder



¹Deisenroth, M.P., Neumann, G., Peters, J. "A survey on policy search for robotics", Foundations and Trends in Robotics, 2013

Black-box Data-efficient ROBot Policy Search (Black-DROPS)

Key Idea #1 *Implicit Policy Evaluation*; treat each rollout as a noisy measurement of the *expected long-term reward*:

$$G(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + N(\boldsymbol{\theta}) \quad (9)$$

$$\begin{aligned} \mathbb{E}[G(\boldsymbol{\theta})] &= \mathbb{E}[J(\boldsymbol{\theta}) + N(\boldsymbol{\theta})] = \mathbb{E}[J(\boldsymbol{\theta})] + \mathbb{E}[N(\boldsymbol{\theta})] \\ &= J(\boldsymbol{\theta}) + \mathbb{E}[N(\boldsymbol{\theta})] \text{ (since } \mathbb{E}[\mathbb{E}[x]] = \mathbb{E}[x]) \end{aligned} \quad (10)$$

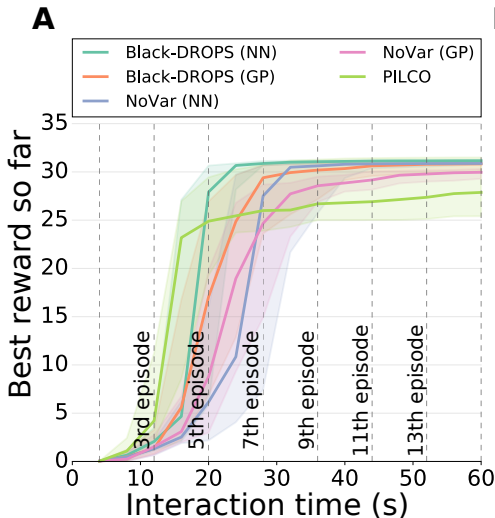
Key Idea #2 Use a black-box optimizer that:

- Is suited for noisy optimization;
- Can take advantage of multi-core computers.

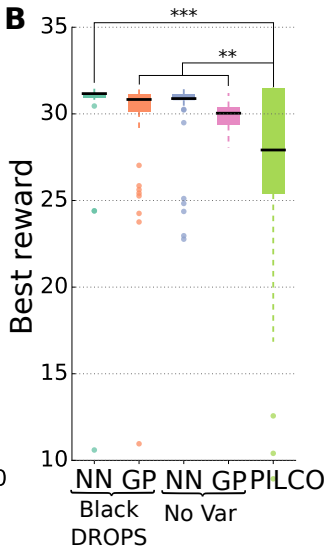
We use a modified version of IPOP-CMA-ES¹ that fulfills the above properties.

¹Hansen, N., Ostermeier, A. "Completely derandomized self-adaptation in evolution strategies", Evolutionary Computation, 2001

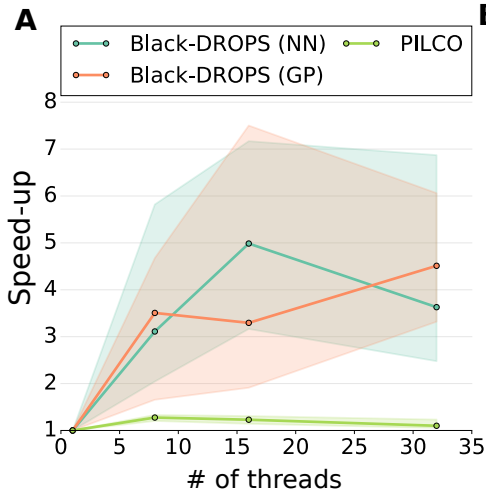
Results - Noiseless Cartpole Simulation



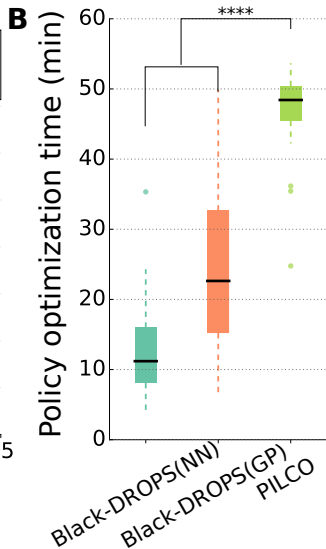
80 replicates



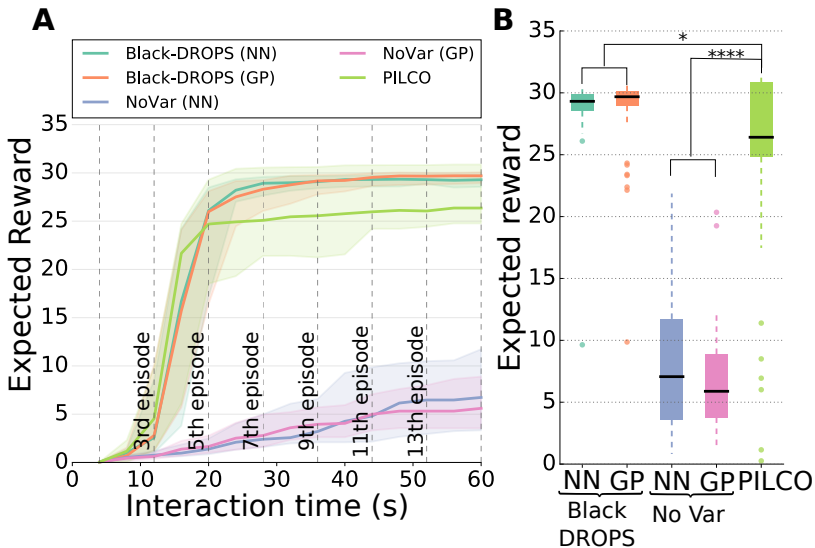
Results - Cartpole Simulation (Scaling)



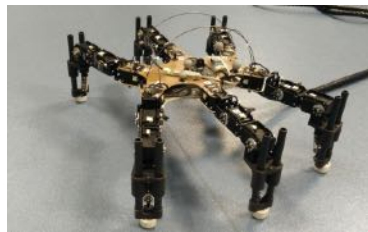
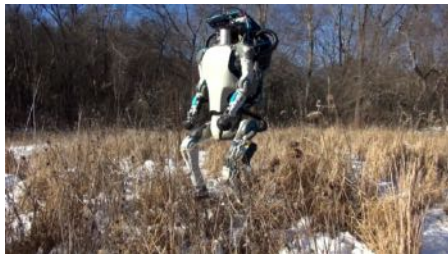
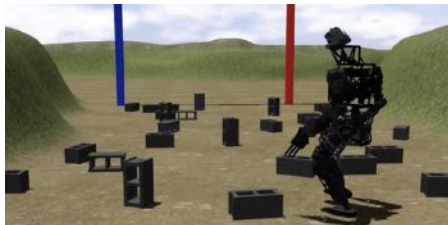
80 replicates



Results - Noisy Cartpole Simulation

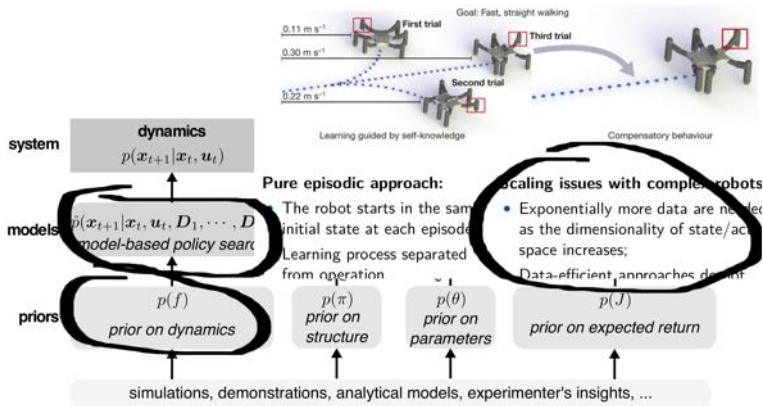


But model-based policy search does not scale!!



Using Parameterized Black-Box Priors to Scale Up Model-Based Policy Search for Robotics

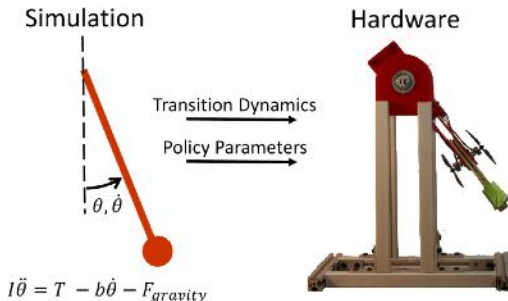
State-of-the-art: Limitations



¹Chatzilygeroudis, K., and Mouret, J.-B. "Using Parameterized Black-Box Priors to Scale Up Model-Based Policy Search for Robotics", ICRA, 2018.

Priors to the rescue!

We can use **dynamic simulators as priors for the dynamics model.**



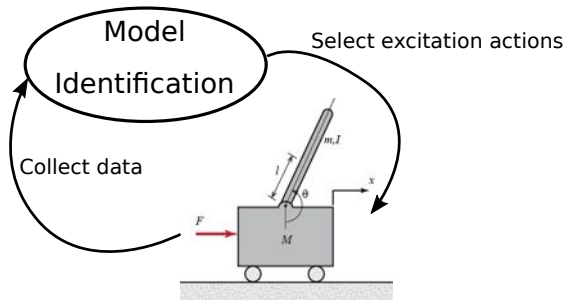
¹Cutler, M., and How, J. P. "Efficient reinforcement learning for robots using informative simulated priors." ICRA, 2015.

²Saveriano, M., Yin, Y., Falco, P., and Lee, D. "Data-Efficient Control Policy Search using Residual Dynamics Learning.", IROS, 2017.

Model Identification

When equations are available, identifying the parameters from data would be the classic approach, but:

- Assumes the actual system can be fully captured by the available equations;
- Proper excitation of the system required for non-trivial cases.



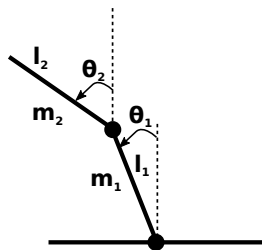
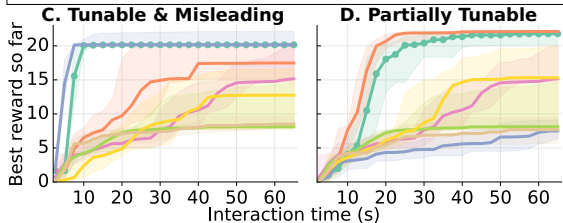
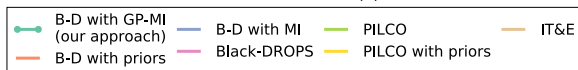
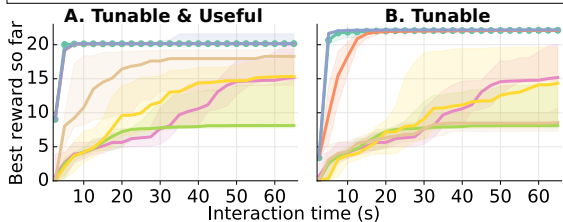
Gaussian Processes with Parameterized Black-Box Priors

- Let's re-write the equation of our dynamic system:

$$\mathbf{x}_{t+1} = \underbrace{M(\mathbf{x}_t, \mathbf{u}_t, \phi_M) + f(\mathbf{x}_t, \mathbf{u}_t, \phi_K)}_{F(\mathbf{x}_t, \mathbf{u}_t)} + \mathbf{w} \quad (11)$$

- Each ϕ_M corresponds to a different parameterization of the mean model M .
- We use GPs to model F and $M(\mathbf{x}_t, \mathbf{u}_t, \phi_M)$ is the mean function of the GPs; $f(\mathbf{x}_t, \mathbf{u}_t, \phi_K)$ captures what cannot be captured by the mean function.
- Now we **can jointly optimize** for ϕ_M and ϕ_K through Maximum Likelihood Estimation (MLE).
- The modeling procedure can balance between non-parametric modeling and model identification (we call it **GP-MI**).

Results - Pendubot Simulation



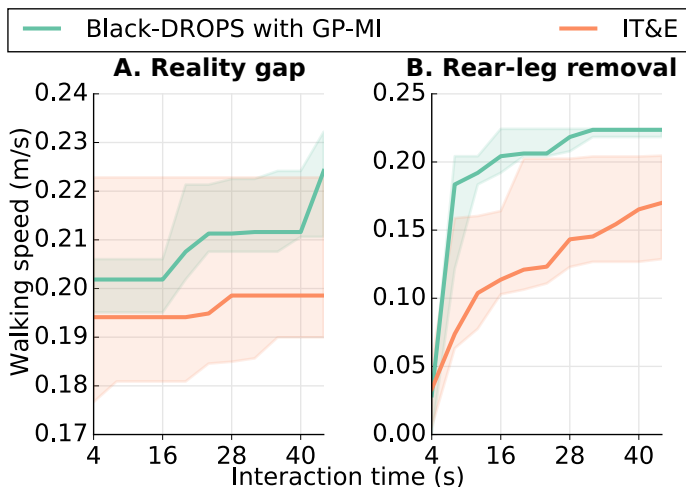
Pendubot system

30 replicates

Results - Physical Hexapod Robot Experiments



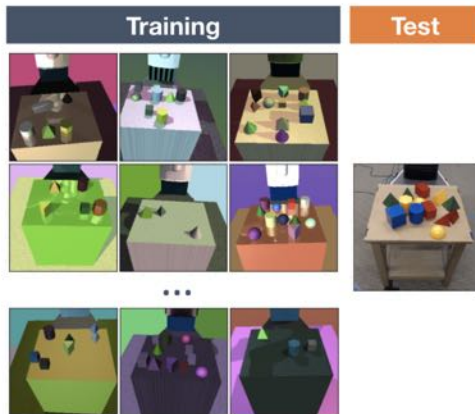
Real Hexapod Robot Experiments



5 replicates

¹Cully, A. et al. "Robots that can adapt like animals", in Nature, vol. 521, no. 7553, pp. 503–507, 2015

Sim2Real Methods



ΠΑΝΕΠΙΣΤΗΜΙΟ
ΠΑΤΡΩΝ
UNIVERSITY OF PATRAS



CEID
COMPUTER ENGINEERING & INFORMATICS DEPARTMENT

¹Tobin, J. et al. "Domain Randomization for Transferring Deep Neural Networks from Simulation to the Real World", IROS, 2017

Main intuitions:

- Simulators give us the ground truth data;
- Simulators can be fast;
- Simulators are reproducible;
- Simulators allow easy customization.

Main steps:

Domain Randomization (DR)

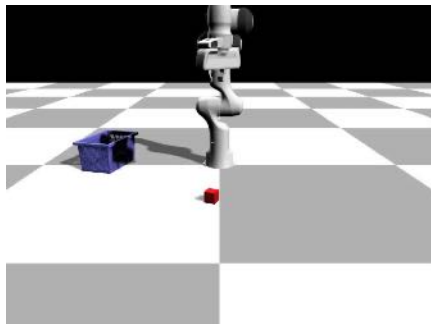
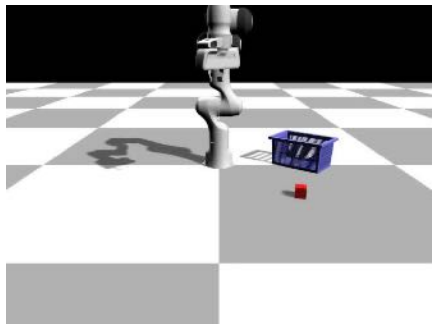
- For each episode, select a random initialization of the world dynamics;
- Optimize the policy for a few steps (e.g., via RL);
- Create a new episode.

Imitation Learning (IL)

- Solve/Learn the task in simulation;
- Collect many many variations of the solution;
- Learn the policy via supervised learning.

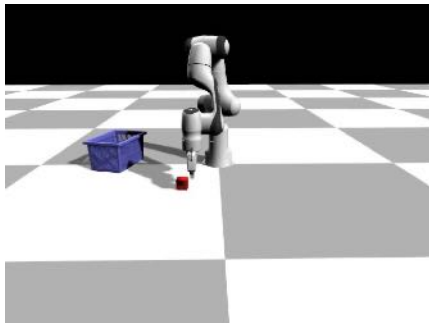
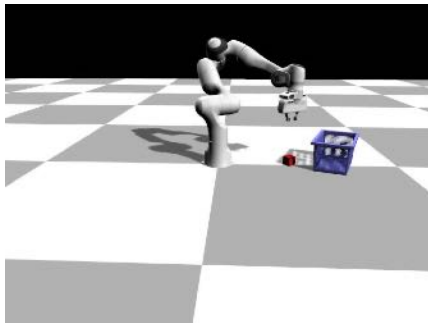
DR + IL gives best results...

Imitation Learning using a Simulator



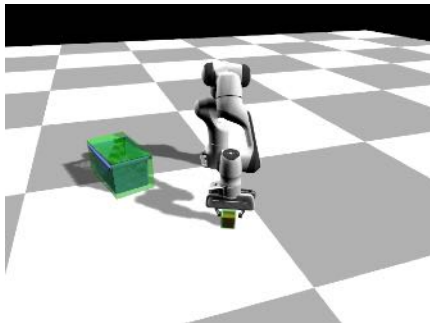
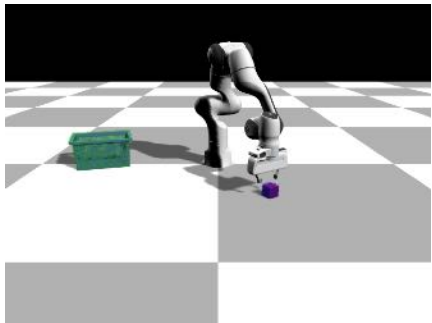
¹Konstantinos Tsinganos, CEID

Imitation Learning using a Simulator (2)



¹Konstantinos Tsinganos, CEID

Visual/Realistic Sim2Real



¹Konstantinos Tsinganos, CEID

Combining Learning with Traditional Methods & Exploration



ΠΑΝΕΠΙΣΤΗΜΙΟ
ΠΑΤΡΩΝ
UNIVERSITY OF PATRAS



EPFL

Quadratic Programming-based Control - Reminder

But we can have analytical models:

$$\begin{aligned}M(\mathbf{q})\ddot{\mathbf{q}} + \mathbf{C}_g(\mathbf{q}, \dot{\mathbf{q}}) &= \mathbf{S}\boldsymbol{\tau} + \mathbf{J}^T(\mathbf{q})\mathbf{W} \\ \mathbf{J}(\mathbf{q})\ddot{\mathbf{q}} + \dot{\mathbf{J}}(\mathbf{q}, \dot{\mathbf{q}})\dot{\mathbf{q}} &= \ddot{\mathbf{x}}\end{aligned}\tag{12}$$

And solve an optimization:

$$\begin{aligned}\min_{\boldsymbol{\chi}} & -\frac{1}{2}\boldsymbol{\chi}^T \mathbf{G}\boldsymbol{\chi} + \mathbf{g}^T \boldsymbol{\chi} \\ \text{s.t.} & \begin{bmatrix} \mathbf{M}(\mathbf{q}) & -\mathbf{S} & -\mathbf{J}(\mathbf{q})^T \end{bmatrix} \boldsymbol{\chi} + \mathbf{C}_g(\mathbf{q}, \dot{\mathbf{q}}) = 0\end{aligned}\tag{13}$$

where

$$\boldsymbol{\chi} = \begin{bmatrix} \ddot{\mathbf{q}} & \boldsymbol{\tau} & \mathbf{W} \end{bmatrix}^T\tag{14}$$

Combining QP-control with Learning

Inverse Dynamics:

$$\begin{aligned}S\tau &= M(q)\ddot{q} + C_g(q, \dot{q}) - J^T(q)W \\ \tau &= M(q)\ddot{q} + C_g(q, \dot{q}) - J^T(q)W\end{aligned}\quad (15)$$

Learning Inverse Dynamics Models:

$$\tau = f_{\text{joint}}(q, \dot{q}, \ddot{q}_{\text{desired}}) \quad (16)$$

$$f_{\text{joint}}(q, \dot{q}, \ddot{q}_{\text{desired}}) = f_{\text{analytic}} + e_{\text{joint}}(q, \dot{q}, \ddot{q}_{\text{desired}}) \quad (17)$$

But e_{joint} **needs to be linear wrt** $\mathcal{X} = [\ddot{q}_{\text{desired}} \quad \tau \quad W]^T$!

Let's linearize: We are in a state $(q_t, \dot{q}_t, \ddot{q}_t)$ and we take the first two terms of the Taylor series expansion of $e_{\text{joint}}(\cdot, \cdot, \cdot)$:

$$\begin{aligned}f_{\text{learned}}(q_t, \dot{q}_t, \ddot{q}_{t+1}) &= f_{\text{analytic}}(q_t, \dot{q}_t, \ddot{q}_{t+1}) + e_{\text{joint}}(q_t, \dot{q}_t, \ddot{q}_t) \\ &\quad + (\ddot{q}_{t+1} - \ddot{q}_t) \frac{\partial e_{\text{joint}}(q, \dot{q}, \ddot{q})}{\partial \ddot{q}} \Bigg|_{\substack{q=q_t \\ \dot{q}=\dot{q}_t \\ \ddot{q}=\ddot{q}_t}}\end{aligned}\quad (18)$$

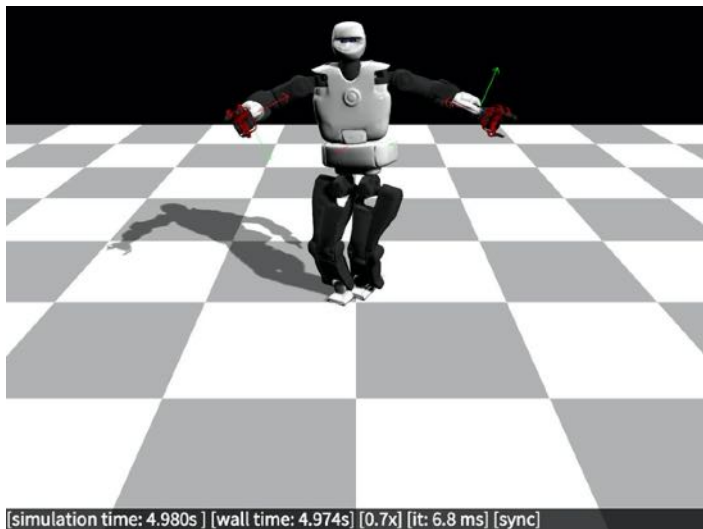
Combining QP-control with Learning (2)

Algorithm Self-correcting QP

- 1: Design the task specifications: $\mathbf{x}_d^*(t)$
 - 2: Configure the adaptive controller and the model learning procedure
 - 3: **for** $n = 1 \rightarrow N_{\text{episodes}}$ **do** ▷ For each episode
 - 4: **for** $t = 0 \rightarrow T$ **do** ▷ For each time-step
 - 5: Get $\ddot{\mathbf{x}}_r^*$ from adaptive controller
 - 6: Compute the cost function for the QP given the $\ddot{\mathbf{x}}_r^*$
 - 7: Get $\boldsymbol{\tau}$ from the QP with the updated cost function, and the learned model $f_{\text{learned}}(\mathbf{q}, \dot{\mathbf{q}}, \ddot{\mathbf{x}}_r^*)$
 - 8: Apply $\boldsymbol{\tau}$ to the robot and collect data
 - 9: Update the adaptive controller
 - 10: Inverse dynamics model learning
-

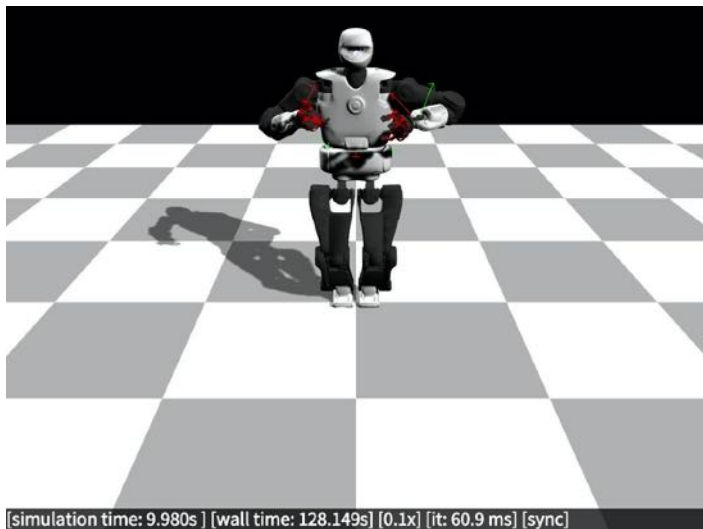
Combining QP-control with Learning (3)

Before Learning



Combining QP-control with Learning (4)

After Learning



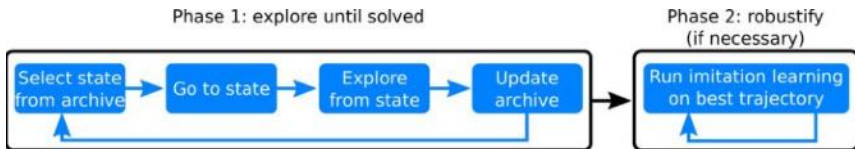
Combining QP-control with Learning (5)

Learning on a physical robot



Discovering interesting behaviors using exploration

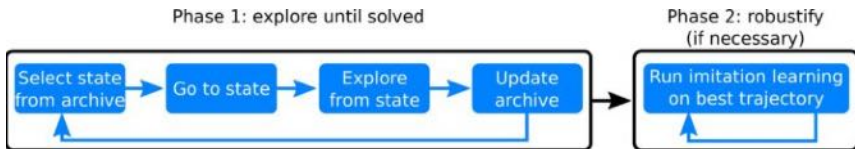
Go-Explore: a New Approach for Hard-Exploration Problems



¹Ecoffet, A. et al. "First return, then explore", Nature, 2021

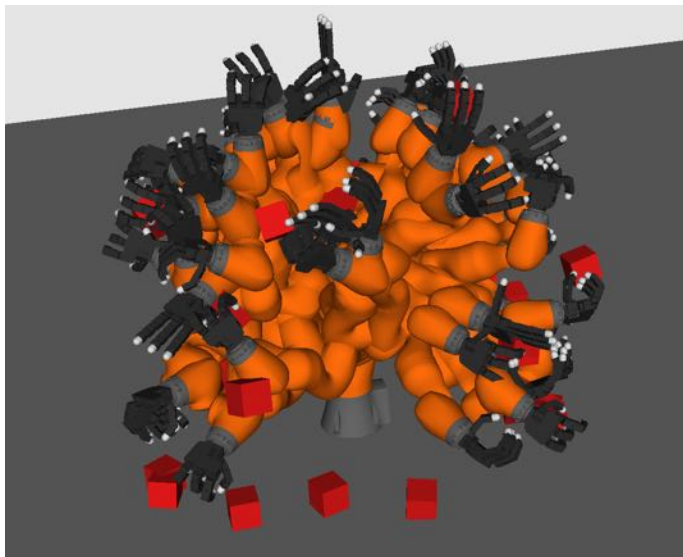
Discovering interesting behaviors using exploration

Go-Explore: a New Approach for Hard-Exploration Problems

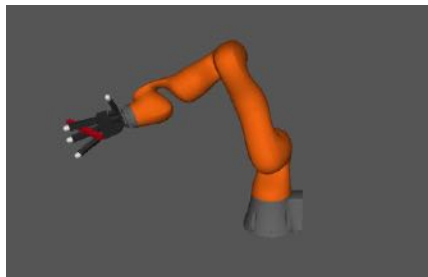
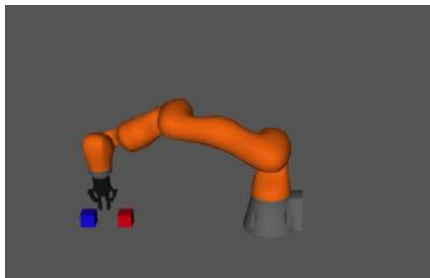
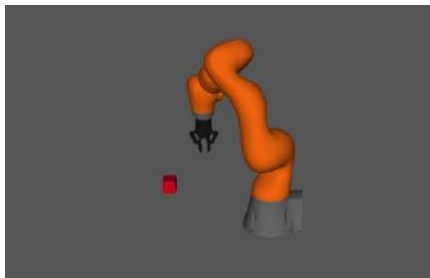


¹Ecoffet, A. et al. "First return, then explore", Nature, 2021

Discovering interesting behaviors using exploration (2)

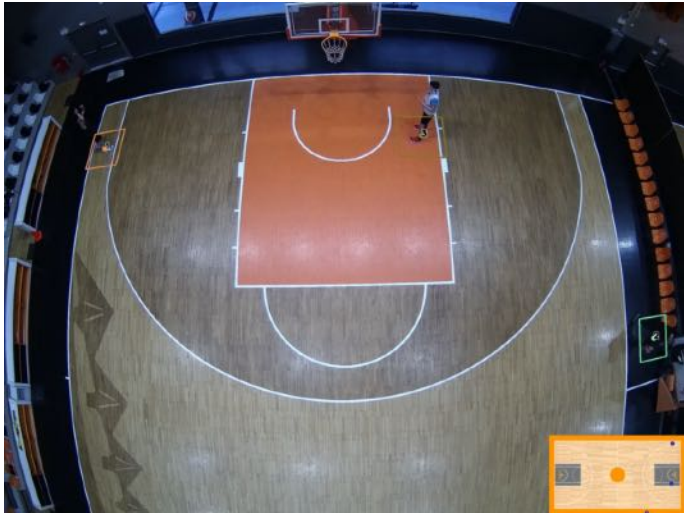


Discovering interesting behaviors using exploration (3)



¹Andrea Mussati, EPFL

Automated Sports Player Tracking and Statistics



metarqus

Thank you!
Any questions?

`costashatz@upatras.gr`

`http://costashatz.github.io`

