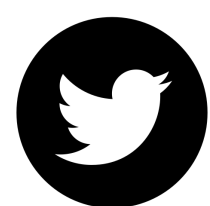# Micro-data policy search

## Learning in a handful of trials

**JB Mouret & Konstantinos Chatzilygeroudis**

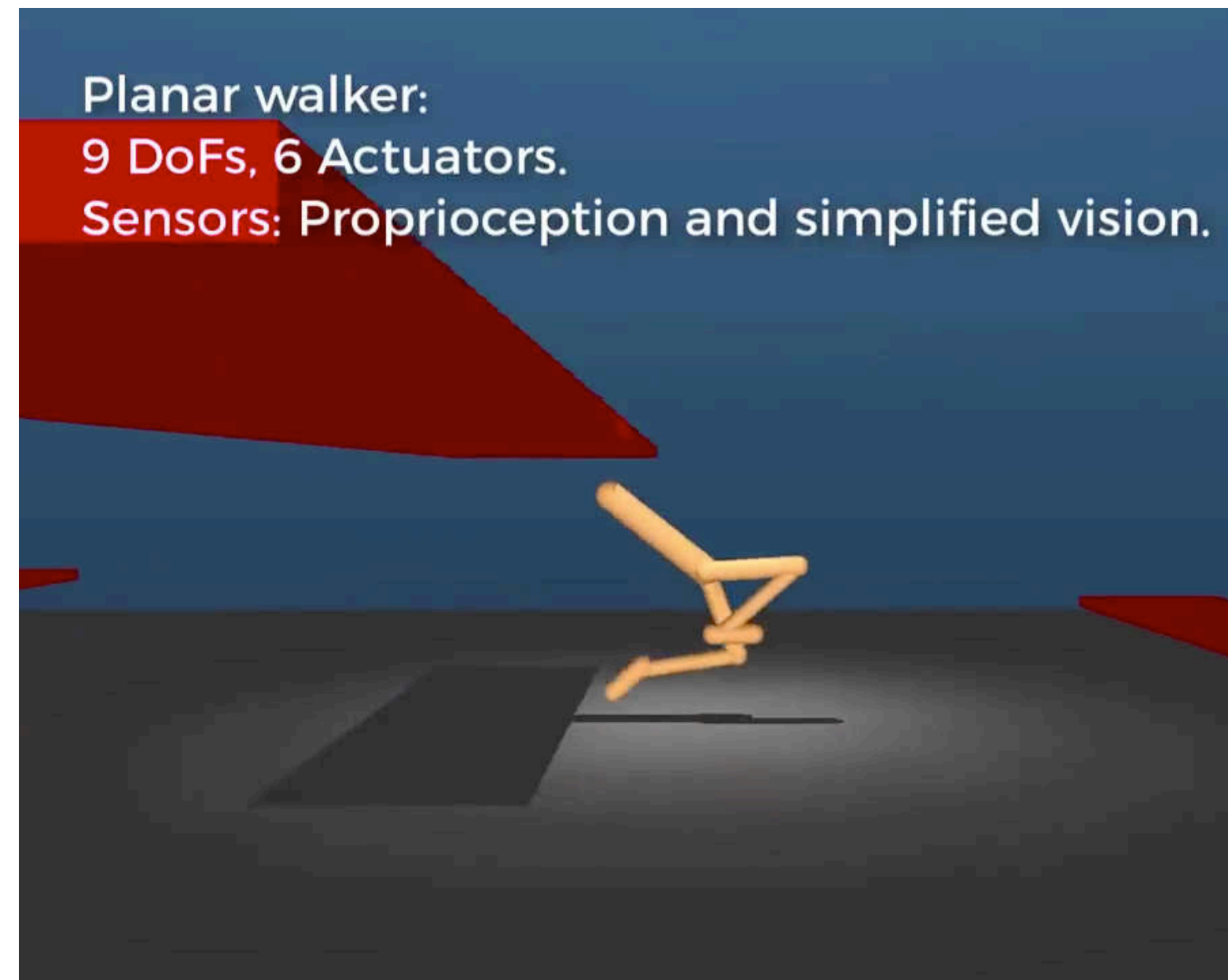**@jbmouret**　　**@jb_mouret**

jean-baptiste.mouret@inria.fr  —  https://members.loria.fr/jbmouret

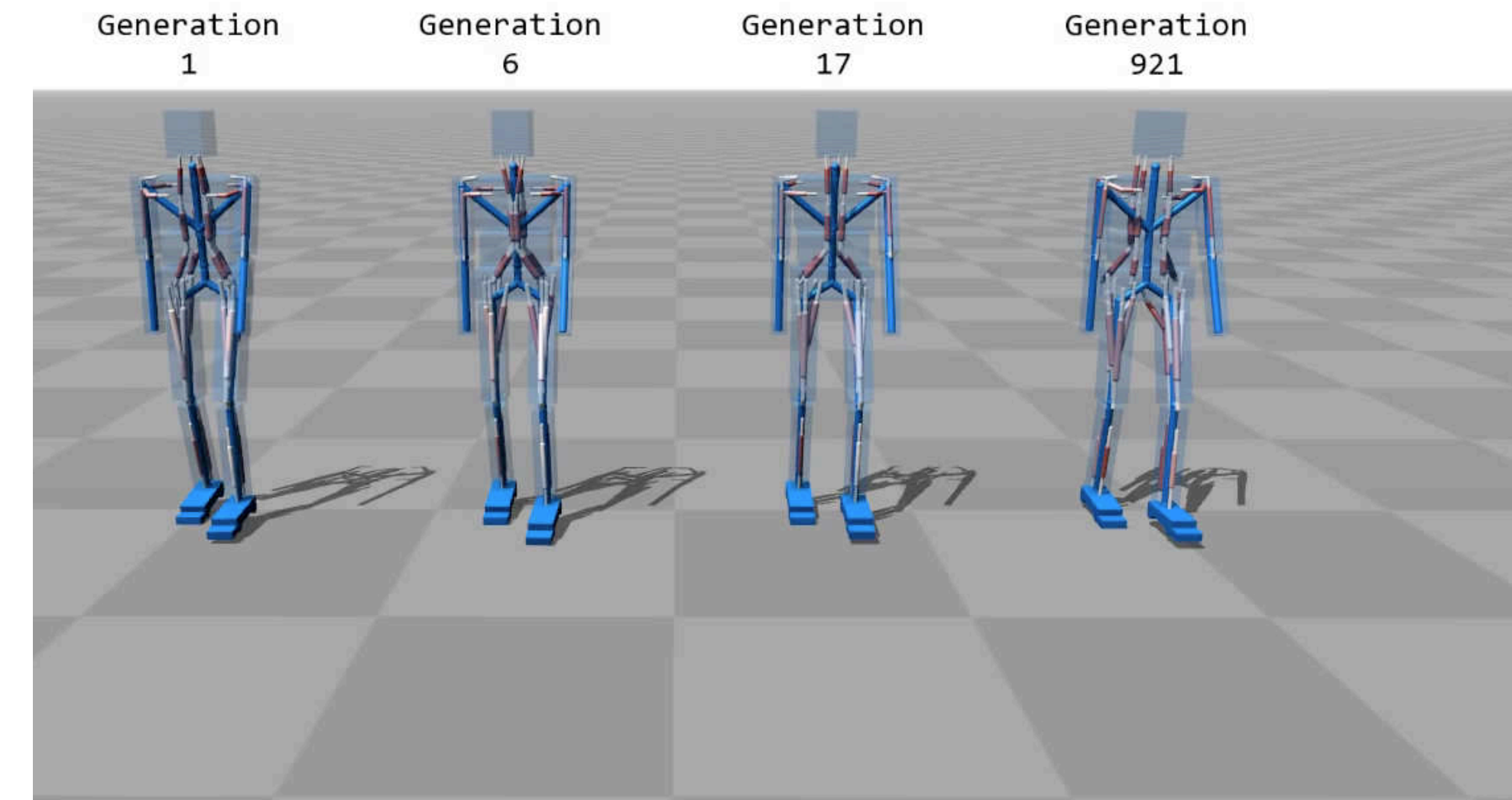# (Deep) Reinforcement Learning is nice but…

## … it needs millions of time-steps to find good solutions (data-efficiency)



**DQN for Atari games :**
38 days (original paper)



**DPPO / PPO**
~10 million training steps
~ 500k seconds at 20Hz (140 hours, or 6 days)



**Algorithm:** CMA-ES /  ~ 20000 episodes

**Mnih, V., et al (2015).** Human-level control through deep reinforcement learning. *nature*, *518*(7540), 529-533.
**Heess N, *et al.*** Emergence of locomotion behaviours in rich environments. arXiv preprint arXiv:1707.02286. 2017 Jul 7. (Google DeepMind)
**Geijtenbeek, T., van de Panne, M., & van der Stappen, A. F. (2013)**. "*Flexible muscle-based locomotion for bipedal creatures"*. ACM Transactions on Graphics (TOG), 32(6), 206.

# Can we do the same on a robot?

## … if we have a lot of time & money

**Levine, S., Pastor, P., Krizhevsky, A., Ibarz, J., & Quillen, D. (2018).** Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International Journal of Robotics Research*, *37*(4-5), 421-436.

# … this is why we use simulators!

**but:**

- Some systems are very expensive to simulate (e.g., fluid dynamics)
- There is (always) a reality gap: policies learned in simulation do not work well on the real system

- *Why do you want to learn?*
  - you have a model (simulator) but you cannot write a controller:
    - are you sure? (Model-Predictive Control, planning, …)
    … MPC/Planning is too expensive to run onboard: use learning (e.g., Guided Policy Search)
    … non-trivial sensors (e.g., vision): maybe use deep-RL (?)
    … model too complex for MPC/planning: maybe use deepRL (?)

No learning
... but model-based control!

Boston Dynamics

Unknown damage
... **no model**!

Forward Speed (m/s)
0.25

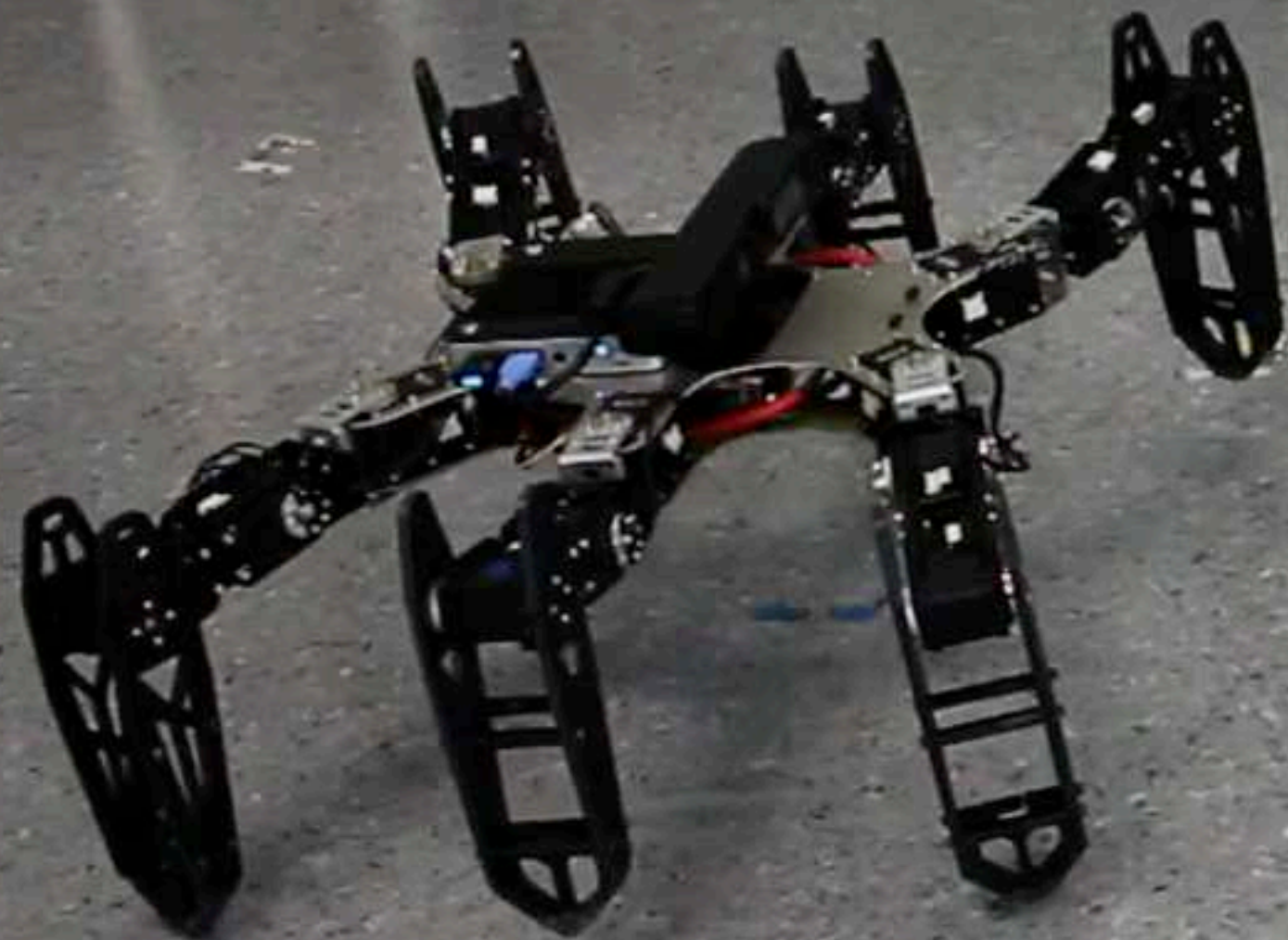Trajectory

# … this is why we use simulators!

**but:**

- Some systems are very expensive to simulate (e.g., fluid dynamics)
- There is (always) a reality gap: policies learned in simulation do not work well on the real system

- *Why do you want to learn a policy?*
  - you have a model (simulator) but you cannot write a controller:
    - are you sure? (Model-Predictive Control, planning, …)
    … MPC/Planning is too expensive to run onboard: use learning (e.g., Guided Policy Search)
    … non-trivial sensors (e.g., vision): maybe use deep-RL (?)
    … model too complex for MPC/planning: maybe use deepRL (?)

  - model/reward is (partially) unknown (e.g., damaged robot): we cannot simulate it accurately!
  - (reinforcement) learning most useful for adaptation in robotics (vs design)

- Using a simulator = having a model (possibly complex) of the system!

# learning is most useful *when the model is unknown*
## *... which means no (accurate) simulator*

AlphaGo

4.9 million games
(self-play) / 40 days

SCORE<1>  HI-SCORE SCORE<2>
  0070        0880

**Atari games :**
38 days

« Micro-data »

Amount of data

« Big-Data »

1-20 trials
A few minutes

?

Deep learning ?

- We want to minimize the **interaction time** with the system
- We consider that pure computation time is "free"

**Mouret, JB. (2016)** "Micro-Data Learning: The Other End of the Spectrum." *ERCIM News*
**K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2018).** A survey on policy search algorithms for learning robot controllers in a handful of trials. IEEE TRO 2020

8

Parameters of the policy

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

Objective        Reward for state $\mathbf{x}_t$

**system**

| **dynamics** | **policy** | **expected return** |
| $p(\boldsymbol{x}_{t+1}\vert\boldsymbol{x}_t, \boldsymbol{u}_t)$ | $\pi(\boldsymbol{u}\vert\boldsymbol{x}, t, \boldsymbol{\theta})$ | $J(\boldsymbol{\theta}) = \mathbb{E}\left[R(\boldsymbol{\tau})\vert\boldsymbol{\theta}\right]$ |

**unknown dynamics**
… but can be queried!
(simulator, robot)

**unknown parameters ($\theta$)**
… but known structure
(neural network,
motion primitive, …)

**K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2020).** A survey on policy search algorithms for learning robot controllers in a handful of trials. IEEE Transactions on Robotics,

Parameters of the policy

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

Objective          Reward for state $\mathbf{x}_t$

**system**

| **dynamics** | **policy** | **expected return** |
|---|---|---|
| $p(\boldsymbol{x}_{t+1}\|\boldsymbol{x}_t, \boldsymbol{u}_t)$ | $\pi(\boldsymbol{u}\|\boldsymbol{x}, t, \boldsymbol{\theta})$ | $J(\boldsymbol{\theta}) = \mathbb{E}\left[R(\boldsymbol{\tau})\|\boldsymbol{\theta}\right]$ |

**models**

$\hat{p}(\boldsymbol{x}_{t+1}|\boldsymbol{x}_t, \boldsymbol{u}_t, \boldsymbol{D}_1, \cdots, \boldsymbol{D}_N)$
*model-based policy search*

$\hat{J}(\boldsymbol{\theta}|\boldsymbol{D_1}, \cdots, \boldsymbol{D}_n)$
*Bayesian optimization*

<u>**learn**</u> **to predict dynamics from** $(x, u)$
… then use the predictor like
a simulator to search for $\pi(u \mid x, \theta)$

<u>**learn**</u> **to predict rewards from** $\theta$
… then search for the best
parameters directly

**K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2020).** A survey on policy search algorithms for learning robot controllers in a handful of trials. IEEE Transactions on Robotics,

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

Parameters of the policy

Objective  Reward for state $\mathbf{x}_t$

| system | | | |
|---|---|---|---|
| **dynamics** $p(\boldsymbol{x}_{t+1}|\boldsymbol{x}_t, \boldsymbol{u}_t)$ | **policy** $\pi(\boldsymbol{u}|\boldsymbol{x}, t, \boldsymbol{\theta})$ | | **expected return** $J(\boldsymbol{\theta}) = \mathbb{E}\left[R(\boldsymbol{\tau})|\boldsymbol{\theta}\right]$ |

**models**

$\hat{p}(\boldsymbol{x}_{t+1}|\boldsymbol{x}_t, \boldsymbol{u}_t, \boldsymbol{D}_1, \cdots, \boldsymbol{D}_N)$
*model-based policy search*

*e.g., dynamic movement primitives*

*e.g., demonstrations*

$\hat{J}(\boldsymbol{\theta}|\boldsymbol{D_1}, \cdots, \boldsymbol{D}_n)$
*Bayesian optimization*

**priors**

$p(f)$
*prior on dynamics*

$p(\pi)$
*prior on structure*

$p(\theta)$
*prior on parameters*

$p(J)$
*prior on expected return*

simulations, demonstrations, analytical models, experimenter's insights, ...

*part 3*    *part 1*    *part 2*

**K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2020).** A survey on policy search algorithms for learning robot controllers in a handful of trials. IEEE Transactions on Robotics,

**11**

# Outline

**1 Hand-designed policies**
*use well-designed policies with few parameters*
*use demonstrations*

**2 Bayesian optimization (a model of the reward)**

*learn to predict the reward value from $\theta$*

**3 Model-based policy search**

*learn to predict the reward value from $\theta$*

**4 Keynote: Dongheui Lee (TU Munich)**
*Efficient Motor Skills Learning in Robotics*

# Micro-data policy search
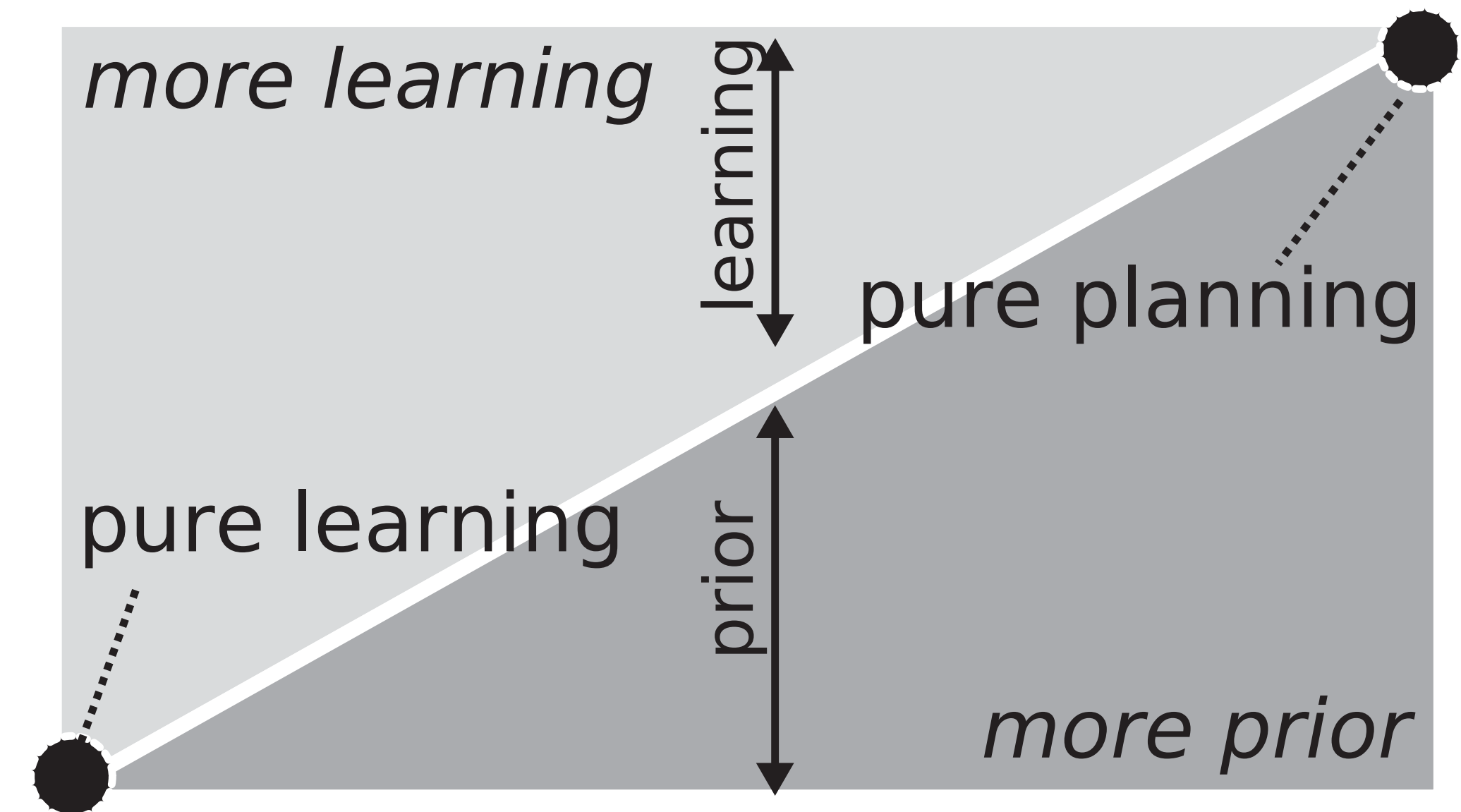
## Conclusion

**JB Mouret & Konstantinos Chatzilygeroudis**

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

Parameters of the policy

Objective

Reward for state $\mathbf{x}_t$

**system**

| dynamics | policy | expected return |
|---|---|---|
| $p(\boldsymbol{x}_{t+1}\|\boldsymbol{x}_t, \boldsymbol{u}_t)$ | $\pi(\boldsymbol{u}\|\boldsymbol{x}, t, \boldsymbol{\theta})$ | $J(\boldsymbol{\theta}) = \mathbb{E}\left[R(\boldsymbol{\tau})\|\boldsymbol{\theta}\right]$ |

**models**

$\hat{p}(\boldsymbol{x}_{t+1}|\boldsymbol{x}_t, \boldsymbol{u}_t, \boldsymbol{D}_1, \cdots, \boldsymbol{D}_N)$
*model-based policy search*

*e.g., dynamic movement primitives*

*e.g., demonstrations*

$\hat{J}(\boldsymbol{\theta}|\boldsymbol{D_1}, \cdots, \boldsymbol{D}_n)$
*Bayesian optimization*

**priors**

$p(f)$
*prior on dynamics*

$p(\pi)$
*prior on structure*

$p(\theta)$
*prior on parameters*

$p(J)$
*prior on expected return*

simulations, demonstrations, analytical models, experimenter's insights, ...

**K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2020).** A survey on policy search algorithms for learning robot controllers in a handful of trials. IEEE Transactions on Robotics,

**14**

# … what to use?

- **Low-DOF robots (<10), high-D policy (NN)**
  → small state-action space
  → model-based policy search

- **High-DOF robot, low-D policy (hand-designed)**
  → Bayesian optimization

- **Complex robots, complex policies**
  → Model-based policy search + prior from model (simulation, meta-learning)
  → BO + MAP-Elites

- **High-D raw inputs (e.g., images)**
  → Sim2real? Use a VAE for unsupervised low-dimensional input?

- **Computation time:**
  → very high with model-based approach
  → low with Bayesian optimisation (+ MAP-Elites)

**K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2020).** A survey on policy search algorithms for learning robot controllers in a handful of trials. IEEE Transactions on Robotics,

# … everything is about priors

- **There are always priors:**
  - neural network topology
  - using a neural networks vs something else
  - hyper-parameters
  - simulator of a real robot
  - design of a reward function, …

- **The best priors are generic priors**
  - e.g. convolutional neural networks are a great prior!
  - priors from simulation
  - meta-learning

- **The best priors can be overcome**

- **We should think more about our prior and work on improving them**
  (instead of hiding them)



*more learning* — learning

pure planning

pure learning — prior

*more prior*

**K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2020).** A survey on policy search algorithms for learning robot controllers in a handful of trials. IEEE Transactions on Robotics,

# Further readings

## A Survey on Policy Search Algorithms for Learning Robot Controllers in a Handful of Trials

Konstantinos Chatzilygeroudis, Vassilis Vassiliades, Freek Stulp, Sylvain Calinon, and Jean-Baptiste Mouret

*Abstract*—Most policy search (PS) algorithms require thousands of training episodes to find an effective policy, which is often infeasible with a physical robot. This survey article focuses on the extreme other end of the spectrum: how can a robot adapt with only a handful of trials (a dozen) and a few minutes? By analogy with the word "big-data," we refer to this challenge as "micro-data reinforcement learning." In this article, we show that a first strategy is to leverage prior knowledge on the policy structure (e.g., dynamic movement primitives), on the policy parameters (e.g., demonstrations), or on the dynamics (e.g., simulators). A second strategy is to create data-driven surrogate models of the expected reward (e.g., Bayesian optimization) or the dynamical model (e.g., model-based PS), so that the policy optimizer queries the model instead of the real system. Overall, all successful micro-data algorithms combine these two strategies by varying the kind of model and prior knowledge. The current scientific challenges essentially revolve around scaling up to complex robots, designing generic priors, and optimizing the computing time.

*Index Terms*—Autonomous agents, learning and adaptive systems, micro-data policy search (MDPS), robot learning.

## I. Introduction

REINFORCEMENT learning (RL) [1] is a generic framework that allows robots to learn and adapt by trial and error. There is currently a renewed interest in RL owing to recent advances in deep learning [2]. For example, RL-based agents can now learn to play many of the Atari 2600 games directly from pixels [3], [4], that is, without explicit feature engineering, and beat the world's best players at Go and chess with minimal human knowledge [5]. Unfortunately, these impressive successes are difficult to transfer to robotics because the algorithms behind them are highly data-intensive: 4.8 million games were required to learn to play Go from scratch [5], 38 days of play (real time) for Atari 2600 games [3], and, for example, about 100 h of simulation time (much more for real time) for a nine-degrees of freedom (9-DOF) mannequin that learns to walk [6].

By contrast, robots have to face the real world, which cannot be accelerated by GPUs nor parallelized on large clusters. The real world will not become faster in a few years, contrary to computers so far (Moore's law). In concrete terms, this means that most of the experiments that are successful in simulation cannot be replicated in the real world because they would take too much time to be technically feasible. As an example, Levine *et al.* [7] recently proposed a large-scale algorithm for learning hand-eye coordination for robotic grasping using deep learning. The algorithm required approximately 800 000 grasps, which were collected within a period of two months using 6–14 robotic manipulators running in parallel. Although the results are promising, they were only possible because they could afford having that many manipulators and because manipulators are easy to automate: it is hard to imagine doing the same with a farm of humanoids.

What is more, online adaptation is much more useful when it is fast than when it requires hours—or worse, days—of trial and error. For instance, if a robot is stranded in a nuclear plant and has to discover a new way to use its arm to open a door; or if a walking robot encounters a new kind of terrain for which it is required to alter its gait; or if a humanoid robot falls, damages its knee, and needs to learn how to limp: in most cases, adaptation has to occur in a few minutes or within a dozen trials to be of any use.

By analogy with the word "big-data," we refer to the challenge of learning by trial and error in a handful of trials as "micro-data RL" [8]. This concept is close to "data-efficient RL" [9], but we think it captures a slightly different meaning. The main difference is that efficiency is a ratio between a cost and benefit,

---

**K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2020).** A survey on policy search algorithms for learning robot controllers in a handful of trials. IEEE Transactions on Robotics,