

# Online adaptation *is a reality gap problem*



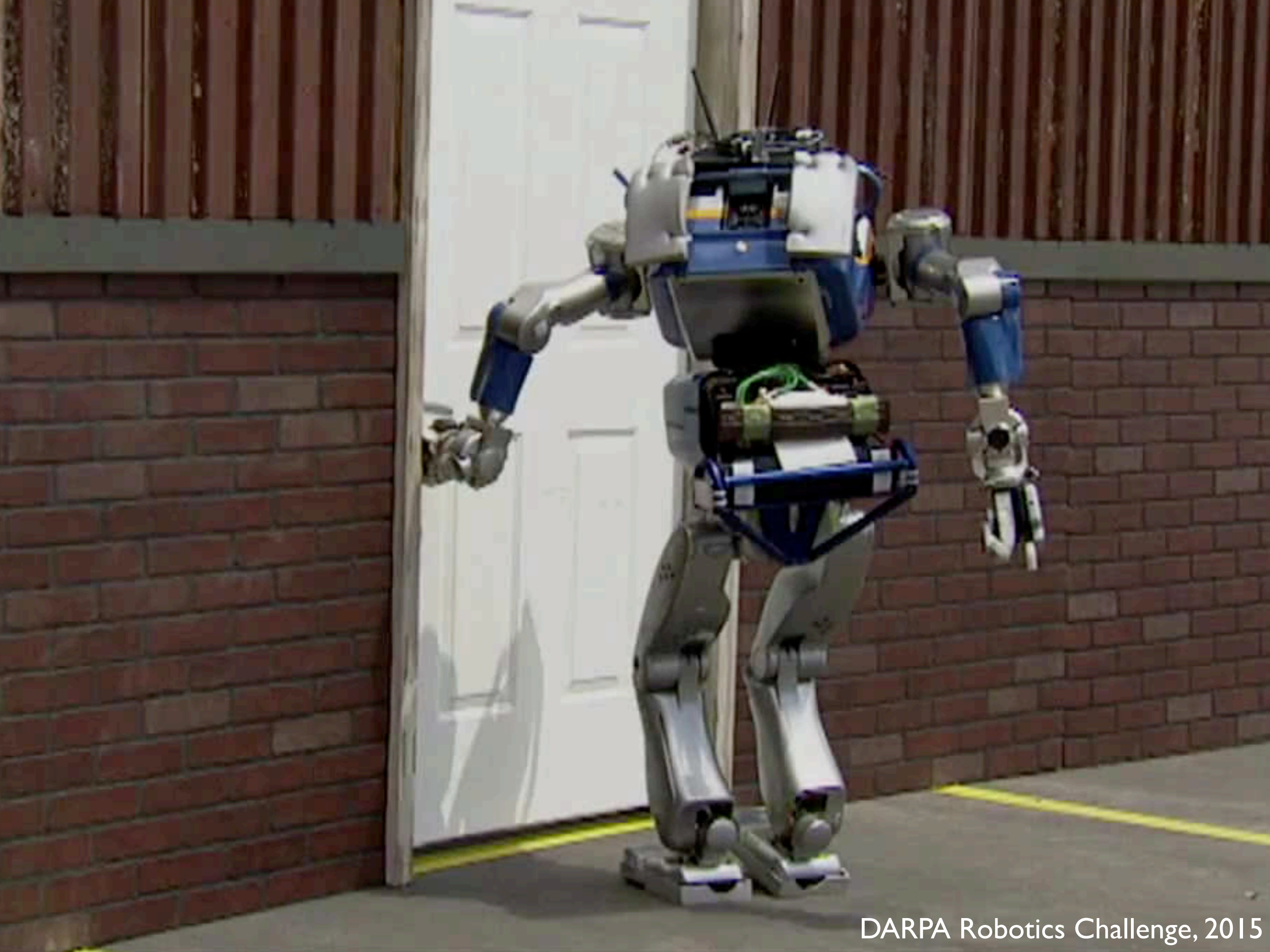
UNIVERSITÉ  
DE LORRAINE



European Research Council  
Established by the European Commission

**Jean-Baptiste Mouret (Inria)**





DARPA Robotics Challenge, 2015



**Online adaptation is a major “learning problem” of robotics  
... more important than controller design / synthesis (?)**



The issue with current robots is not that they fail...

**... it is that they do not get back on their feet and try again**

Our Atlas Robot had mean time between failures of hours or, at most, days [...]

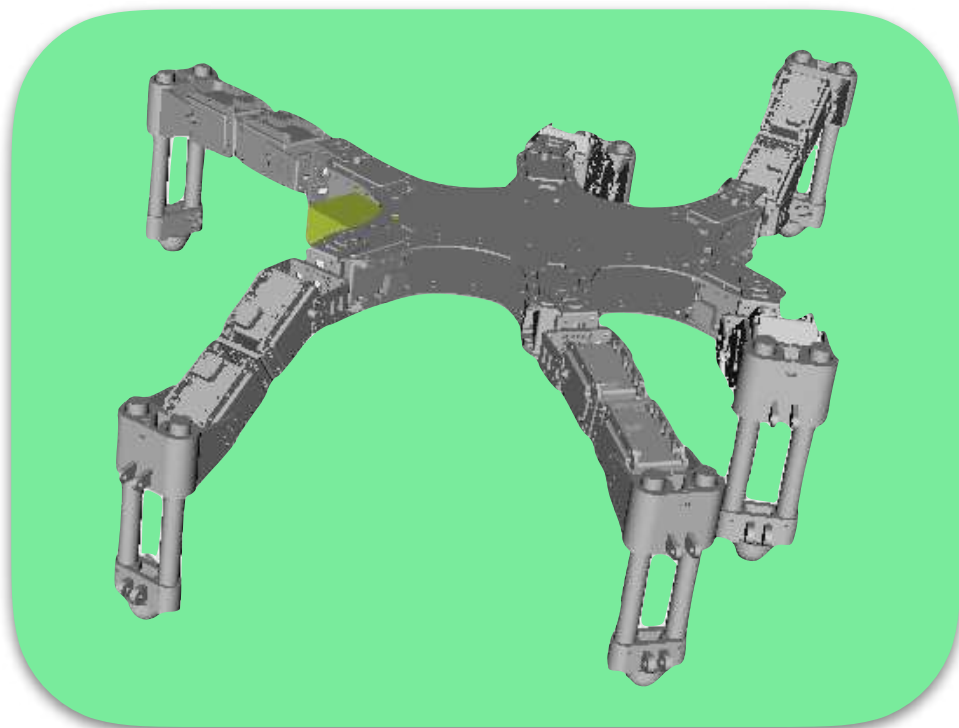
Behaviors that worked perfectly and robustly in the team’s labs did not work or were erratic when tested on an “identical” setup at a DARPA test site.

— **C. Atkeson *et al.* (2018).** What Happened at the DARPA Robotics Challenge Finals. In *The DARPA Robotics Challenge Finals: Humanoid Robots To The Rescue*. Springer.

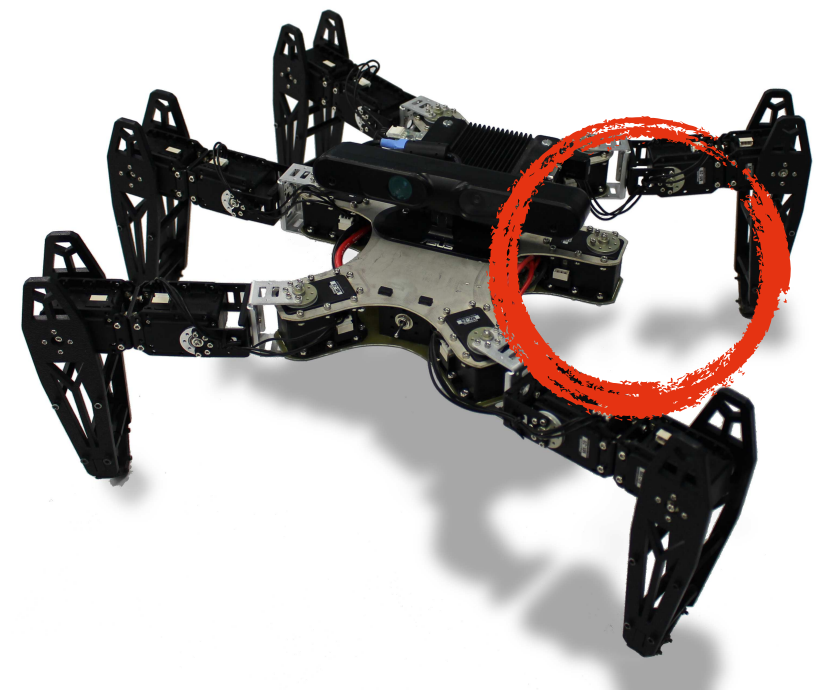


# Reality gap & online adaptation / damage recovery

- We know models of the intact robot in a nominal environment (simulator)
  - We do not know the change / damage
  - We need to adapt quickly to a change (e.g., a damage)
- ... while minimising the number of episodes on the robot
- ⇒ If we can “solve” the reality gap problem, we might solve the adaptation problem
- ⇒ we need to cross the reality gap, be data-efficient, and fast



Simulator of the intact robot



Real (damaged) robot



## Crossing the reality gap

### Learn on the real robot (data-efficiently)

Demonstrations  
(prior on policy  
parameters)

Bayesian  
optimisation

- This is never the full environment
- Only for a few parameters

### Improve *robustness* of policies

Domain  
randomization

Meta-learning

Distribution of possible  
differences between  
simulation and reality?

### Improve simulators (*models*)

Model  
identification

Model-based  
policy search

- When is it easier to learn a simulator than a policy?
- Need to wait between episodes

**Combinations are possible (and useful!)**



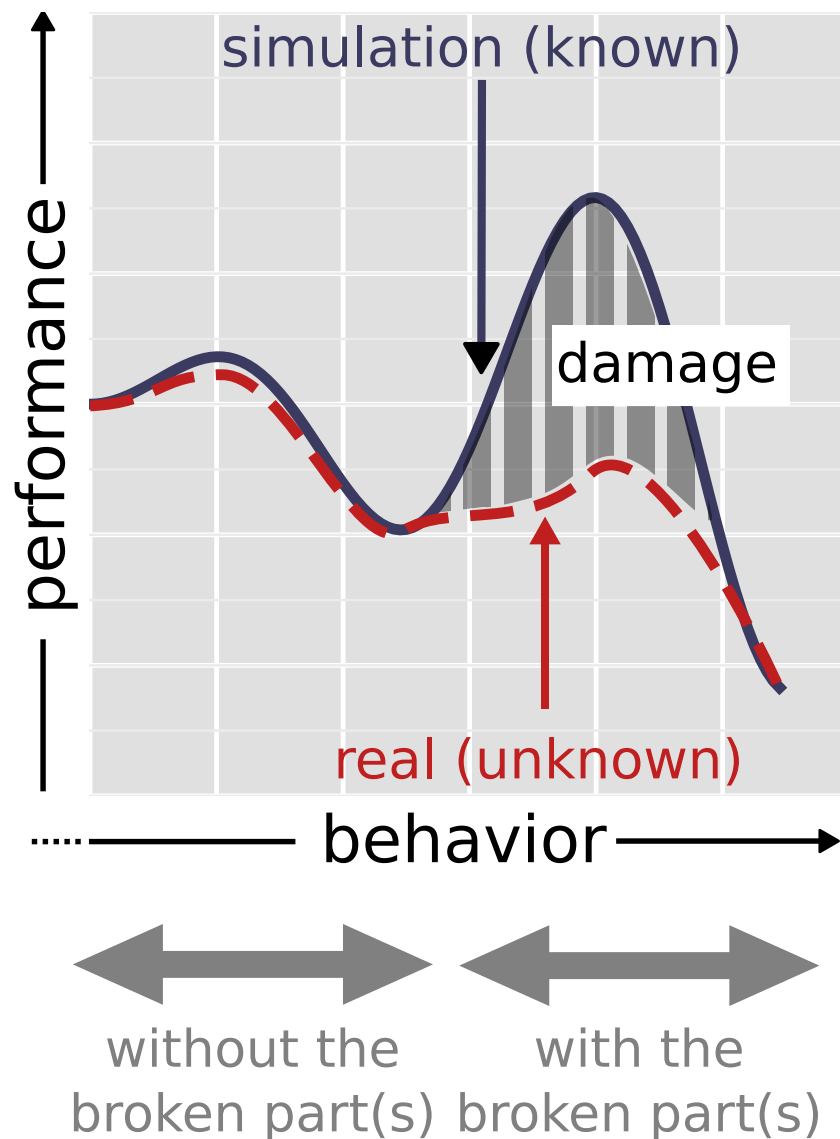
# Crossing the reality gap

Learn on the  
real robot

Improve  
*robustness*  
of policies

Improve  
simulators  
(*models*)

Learn the  
transferability  
function



**Transferability hypothesis:** *some* controllers will work similarly in simulation and in reality

## Transferability approach

- learn a model that predicts the transferability score  
~ learn the limits of the simulation
- search for a policy with a good reward and a good score

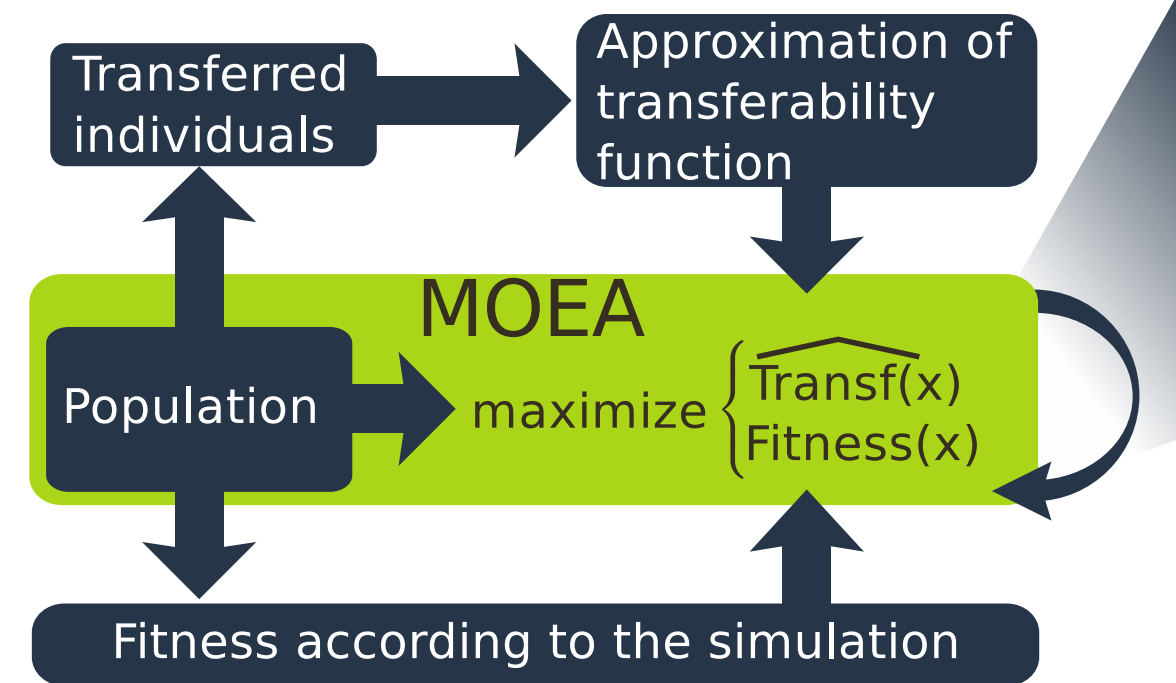
## Why?

- transferability easier to learn than the dynamics
  - ➡ know something is wrong vs knowing the correct answer
- learn a constraint (no crazy predictions)

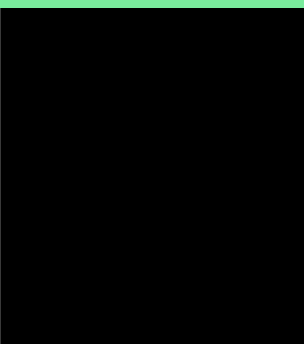
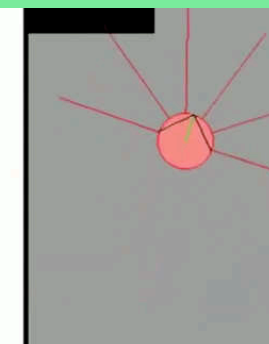
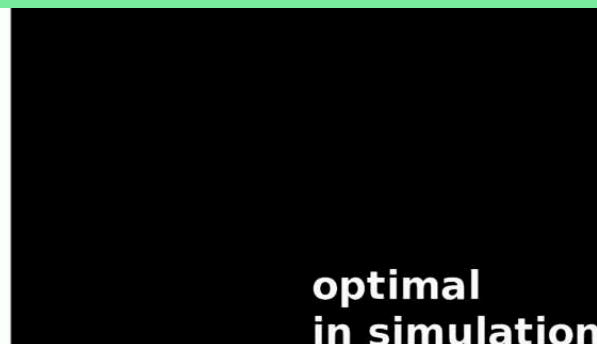
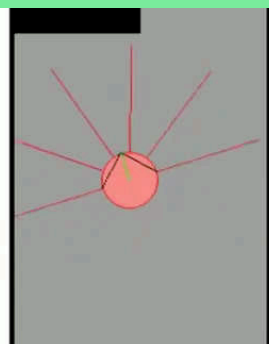


# The transferability approach

Maximize reward



## How can we adapt faster and with embedded hardware?



15 transfers (motion capture)

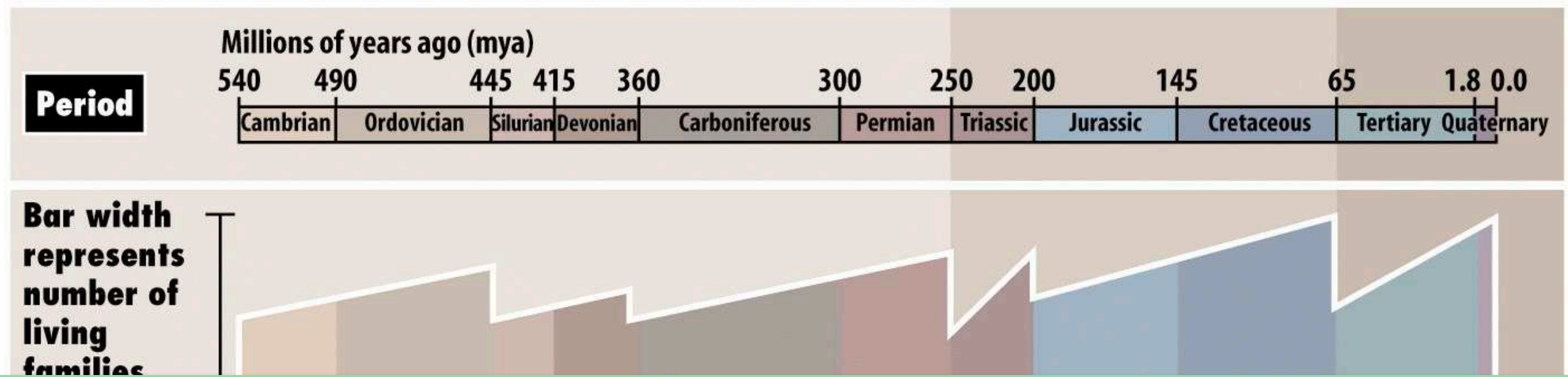
... but a lot of computation between episodes (several minutes)

Koos, S., Cully, A., & Mouret, J. B. (2013). Fast damage recovery in robotics with the T-Resilience algorithm. *The International Journal of Robotics Research*, 32(14), 1700-1723.

Koos, S., Mouret, J.-B., & Doncieux, S. (2011). The Transferability Approach : Crossing the Reality Gap in Evolutionary Robotics. *IEEE Transaction on Evolutionary Computation*.



# Adaptation to sudden changes in nature



**Diversity is preparing for the future “reality gaps”**

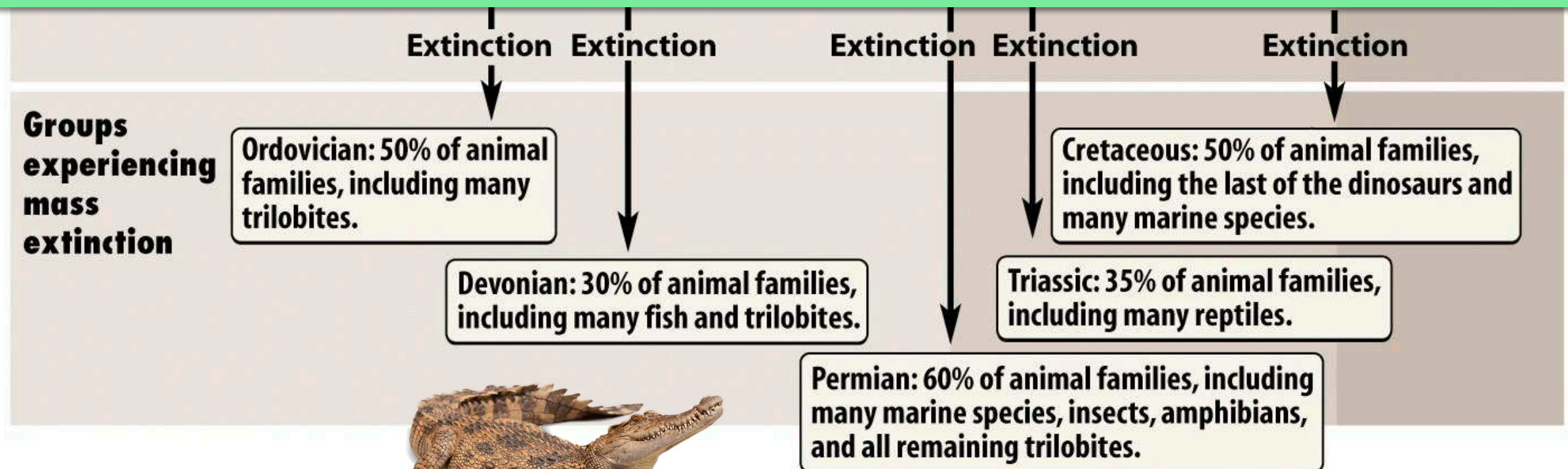


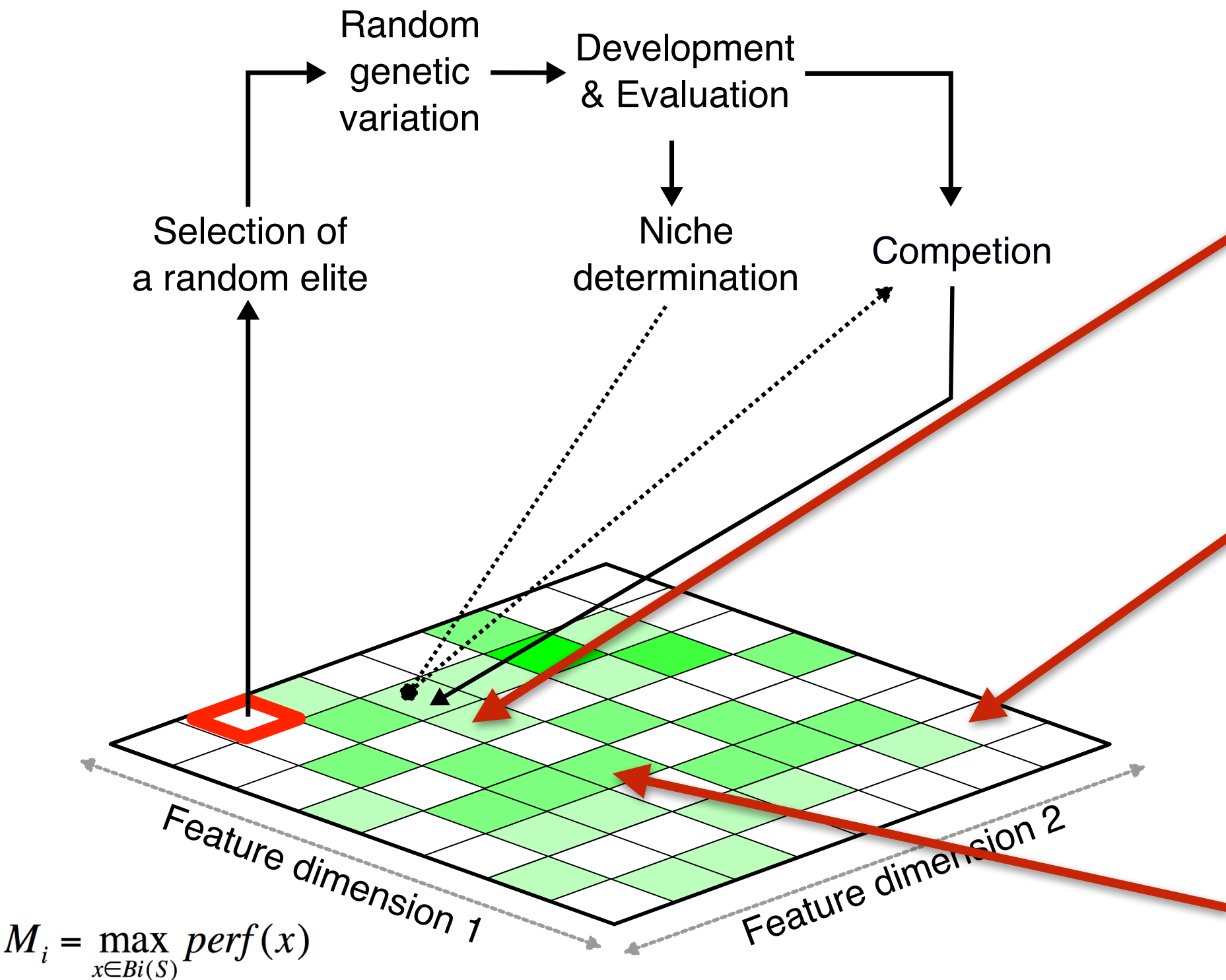
Figure 19-8 Discover Biology 3/e  
© 2006 W. W. Norton & Company, Inc.





# Generating species: The MAP-Elites algorithm

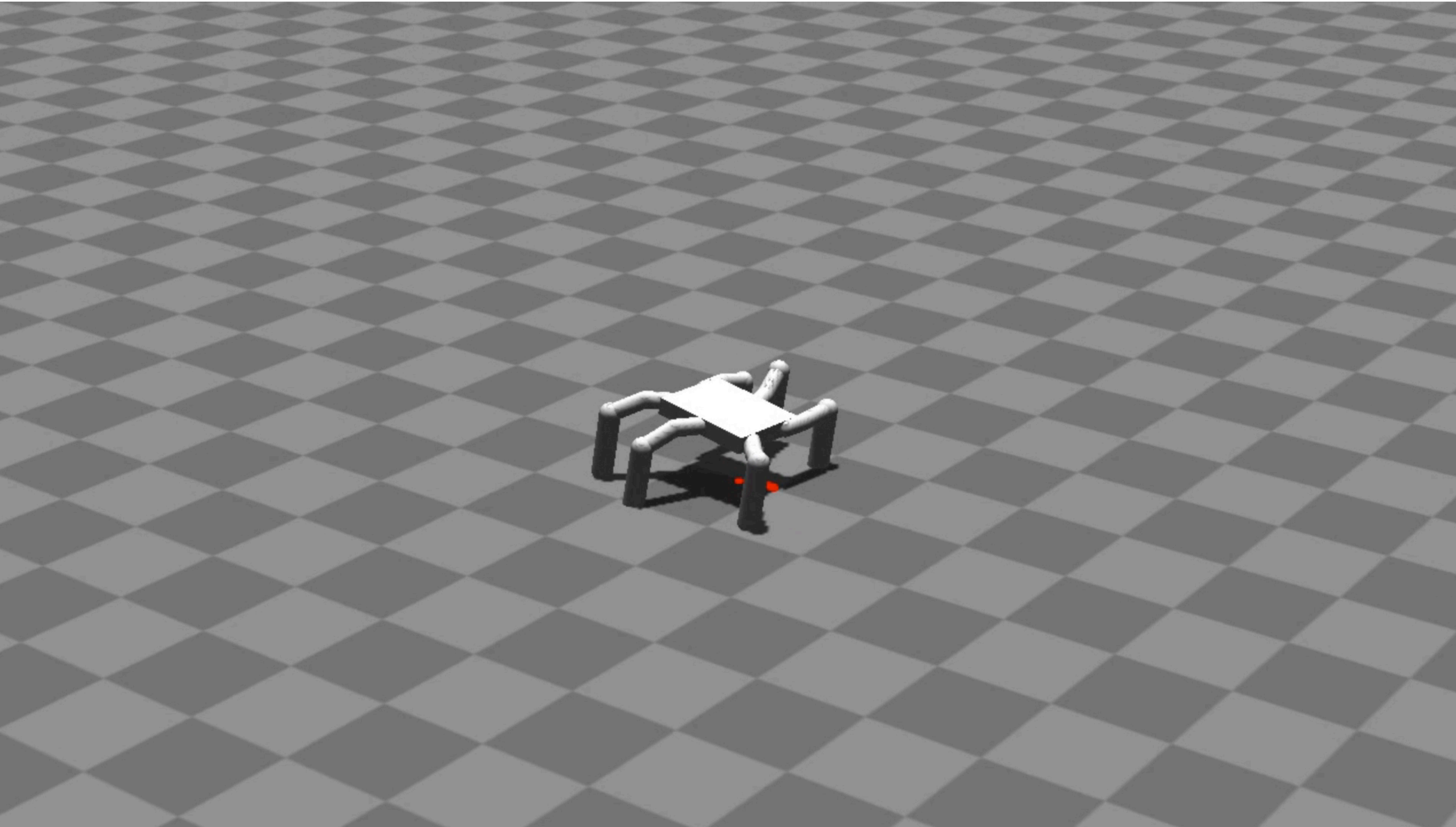
Multi-dimensional Archive of Phenotypic Elites: Quality Diversity / illumination algorithm



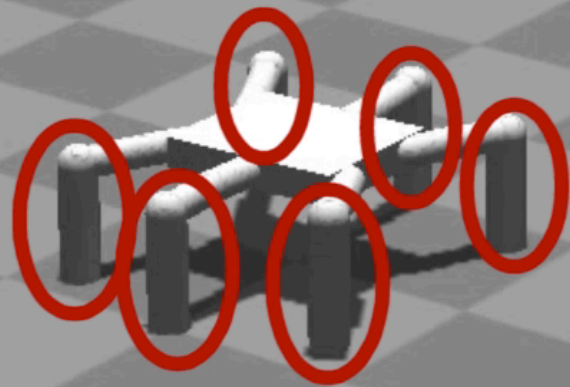
**Mouret, J.-B., and J. Clune. (2015)** "Illuminating search spaces by mapping elites." arXiv preprint arXiv:1504.04909



# MAP-Elites: 6-legged locomotion

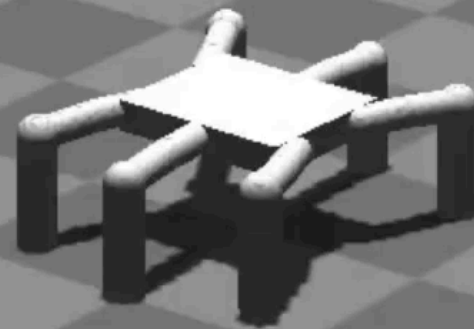


# MAP-Elites: 6-legged locomotion





# MAP-Elites: 6-legged locomotion



**Frequently uses **all** legs**

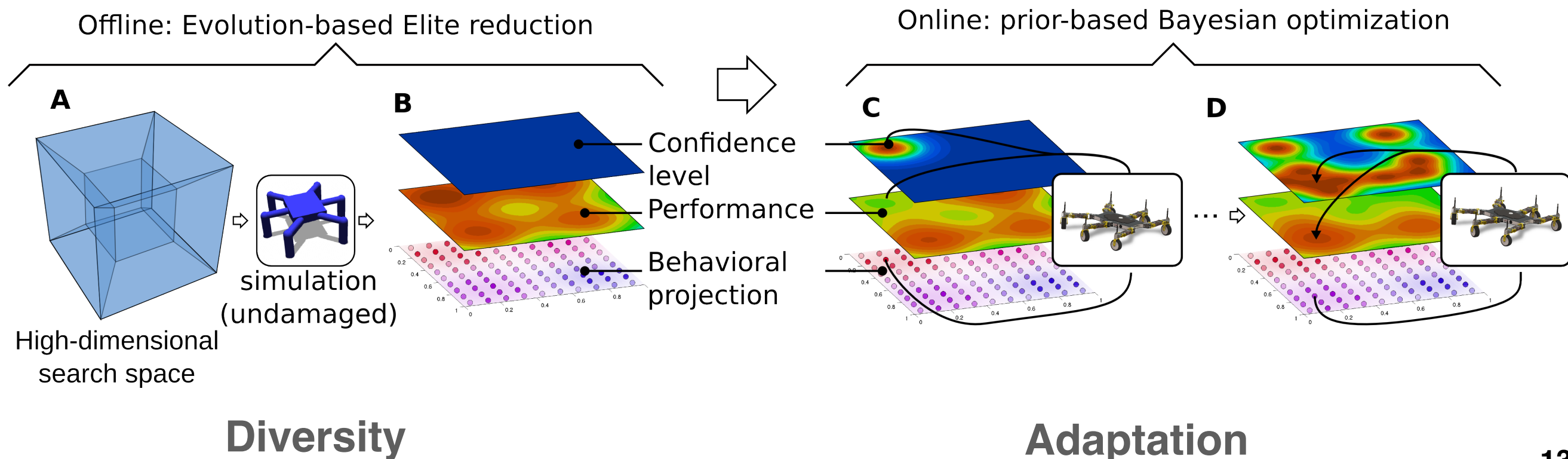
# Intelligent Trial & Error: Map-Elites + BO

The **MAP-Elites algorithm** generates the search space (prior)

- ➡ in simulation, with an intact robot
- ➡ many evaluations [simulation]
- ➡ “take the needles out of the haystack”
- ➡ provide an expected performance for the “needles”

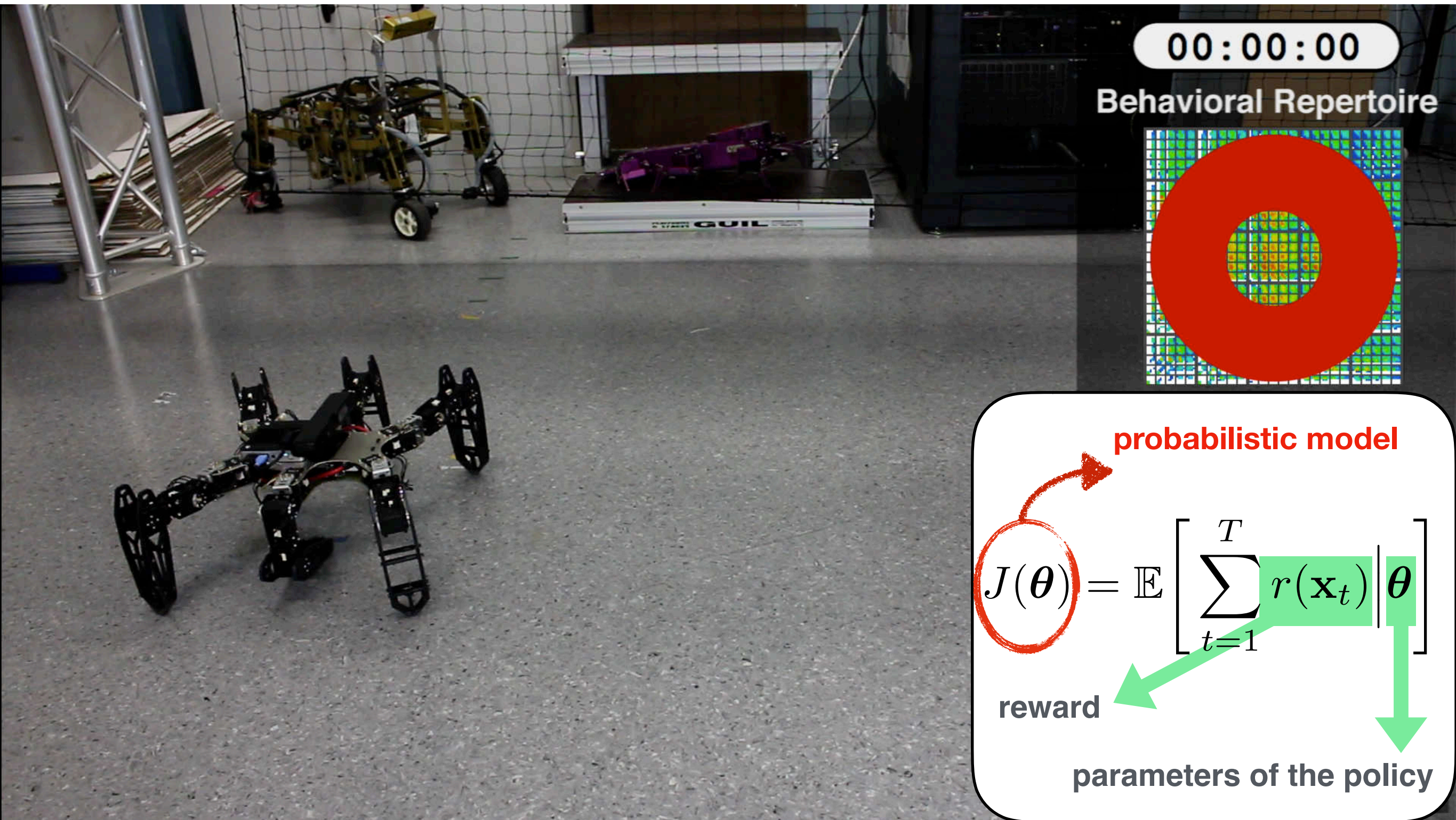
**Prior-based Bayesian optimization** does the online learning

- ➡ search only among good solutions (“needles”)
- ➡ trial-and-error
- ➡ few evaluations [real robot]





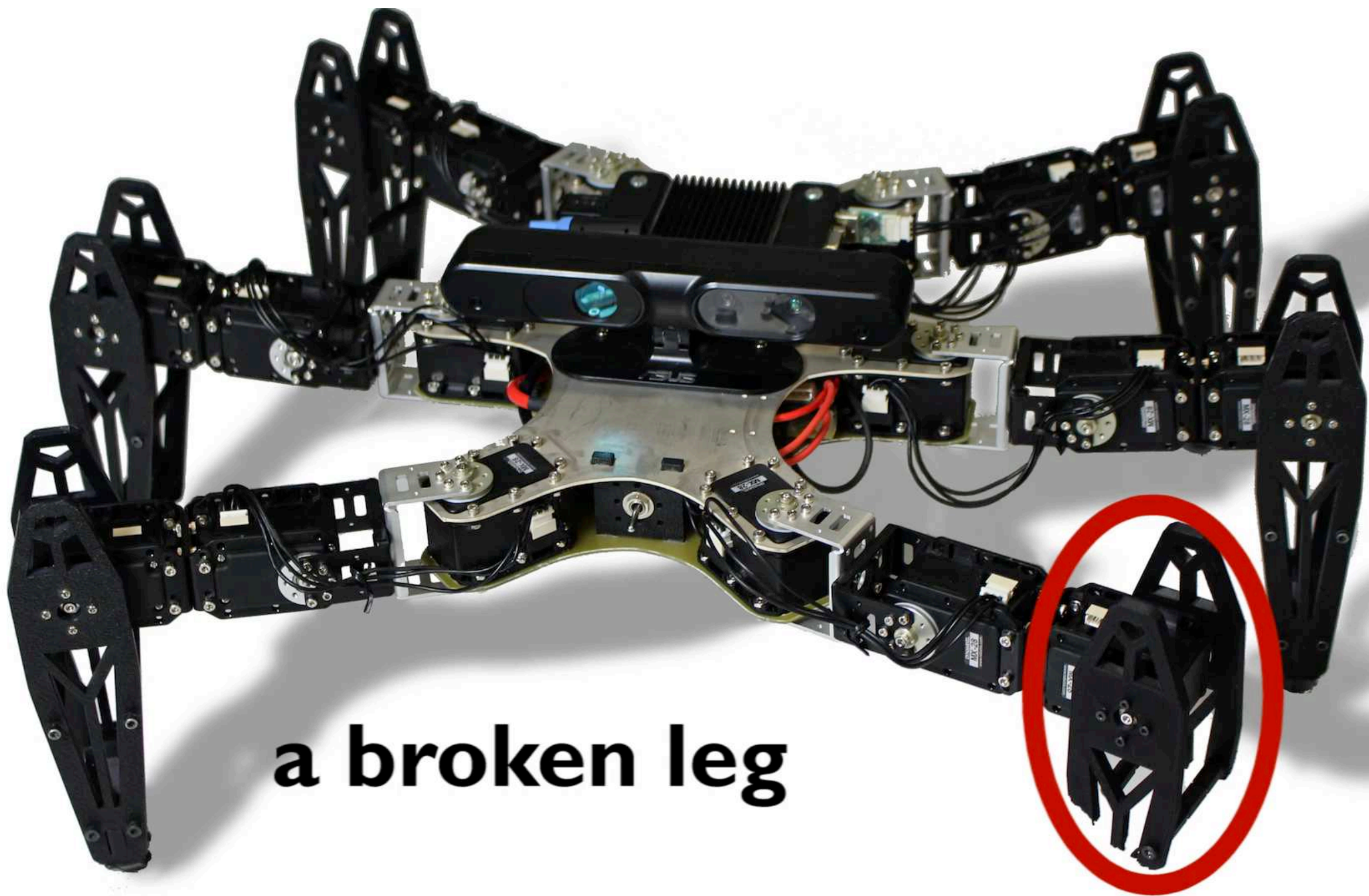
# Adaptation: Bayesian optimization



- policy : periodical signals (36 parameters)
- No information about the damage

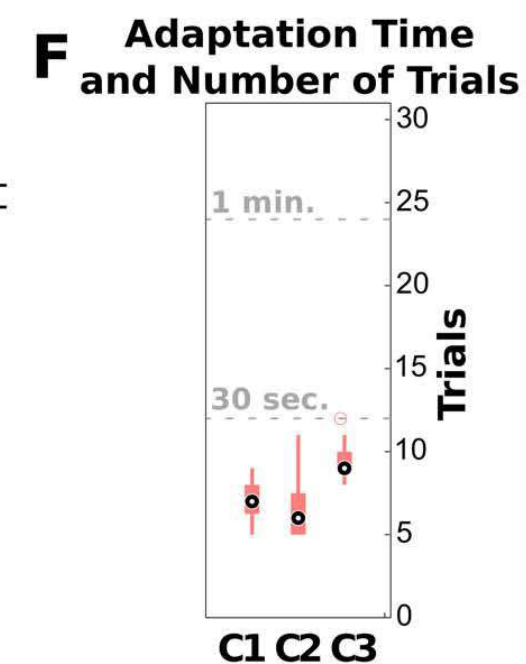
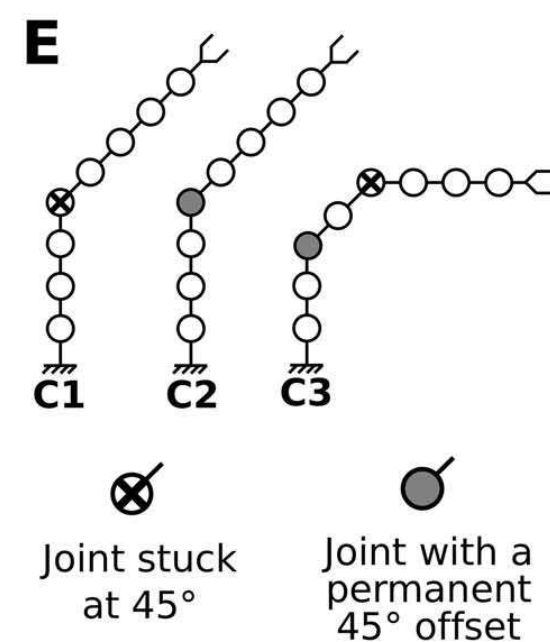
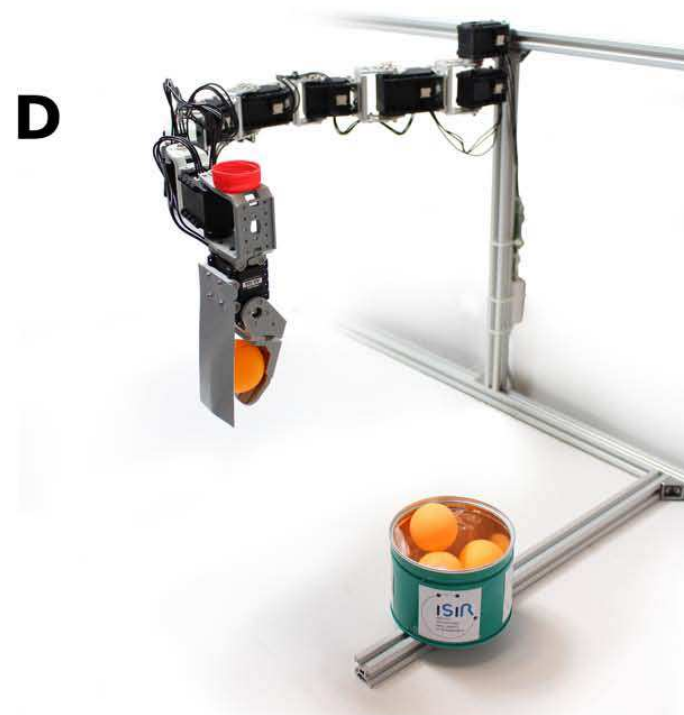
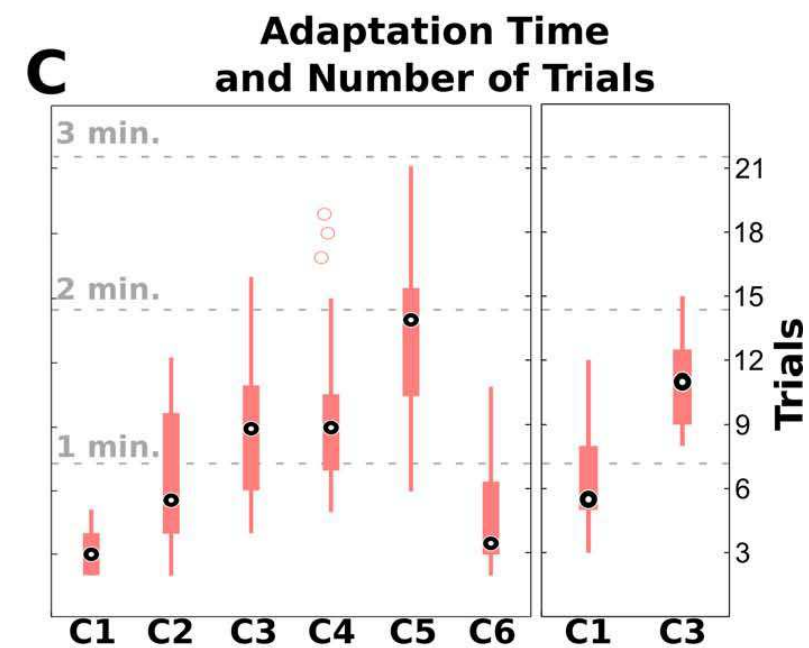
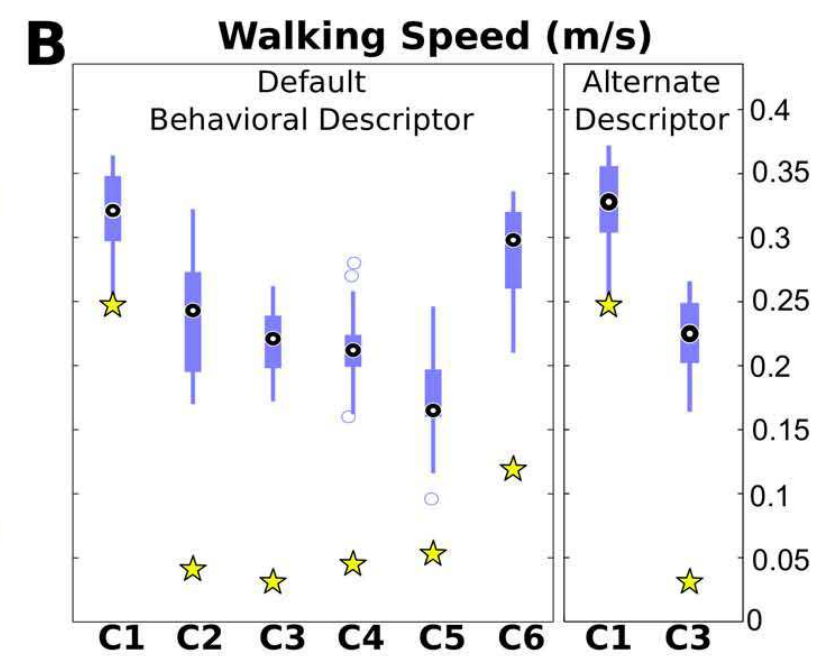
Cully, A. and Clune, J. and Tarapore, D. and Mouret, J.-B. (2015). *Robots that can adapt like animals*. Nature. Vol 521 Pages 503-507.



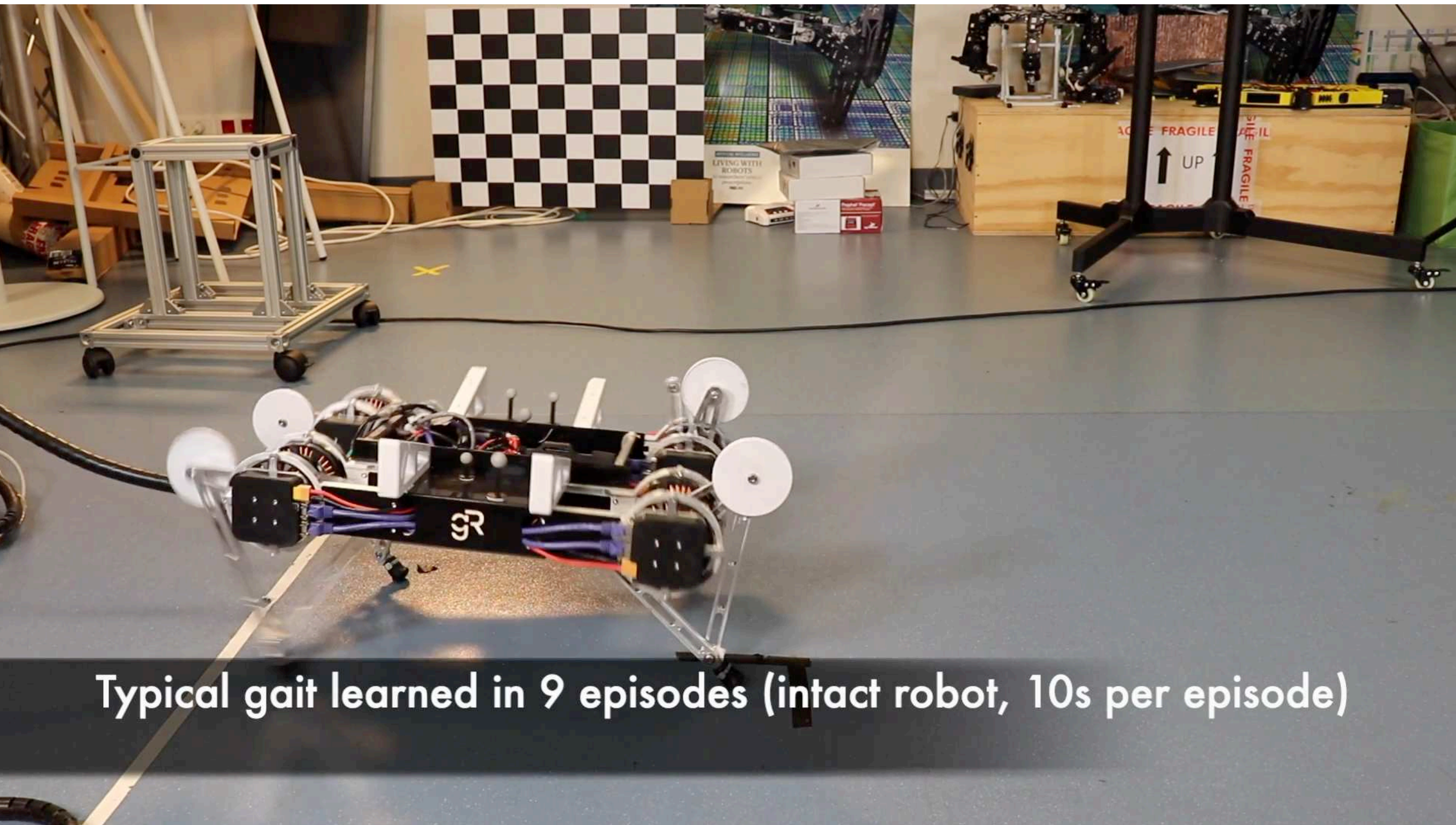


**a broken leg**









Typical gait learned in 9 episodes (intact robot, 10s per episode)

- Policy : periodical signals (24 parameters)
- No information about the damage

Dalin, P. Desreumaux, J.-B. Mouret. (2019) Learning and adapting quadruped gaits with the “Intelligent Trial & Error” algorithm. ICRA Workshop on “Learning Legged Locomotion”.





# Recent extension: Multiple priors

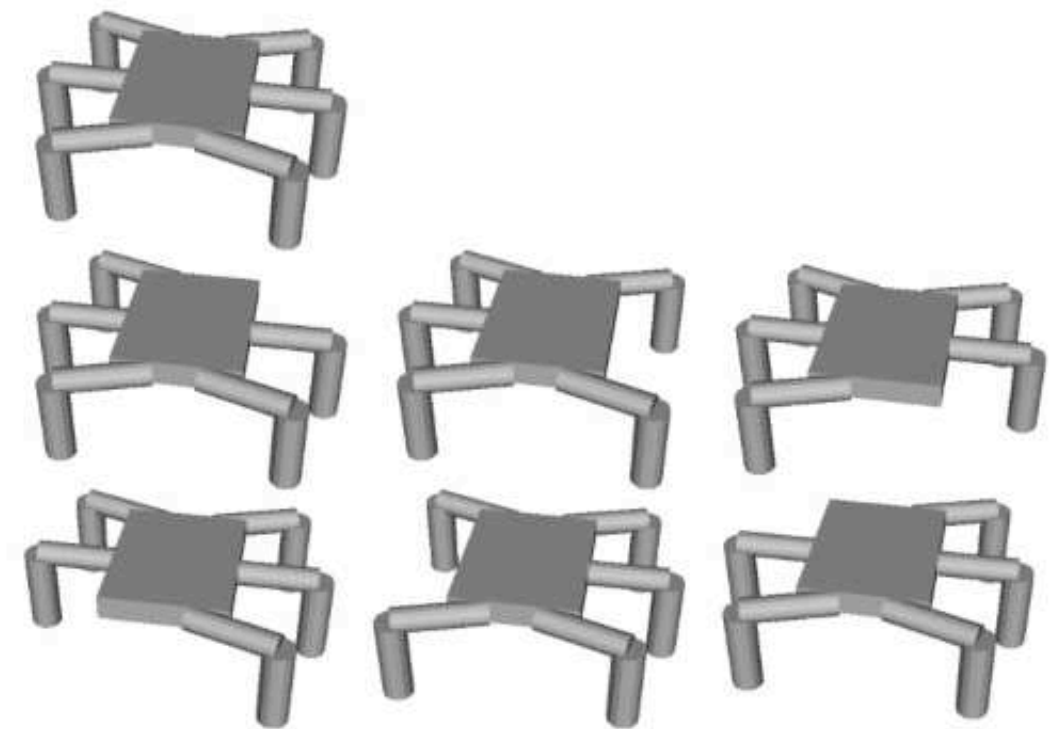
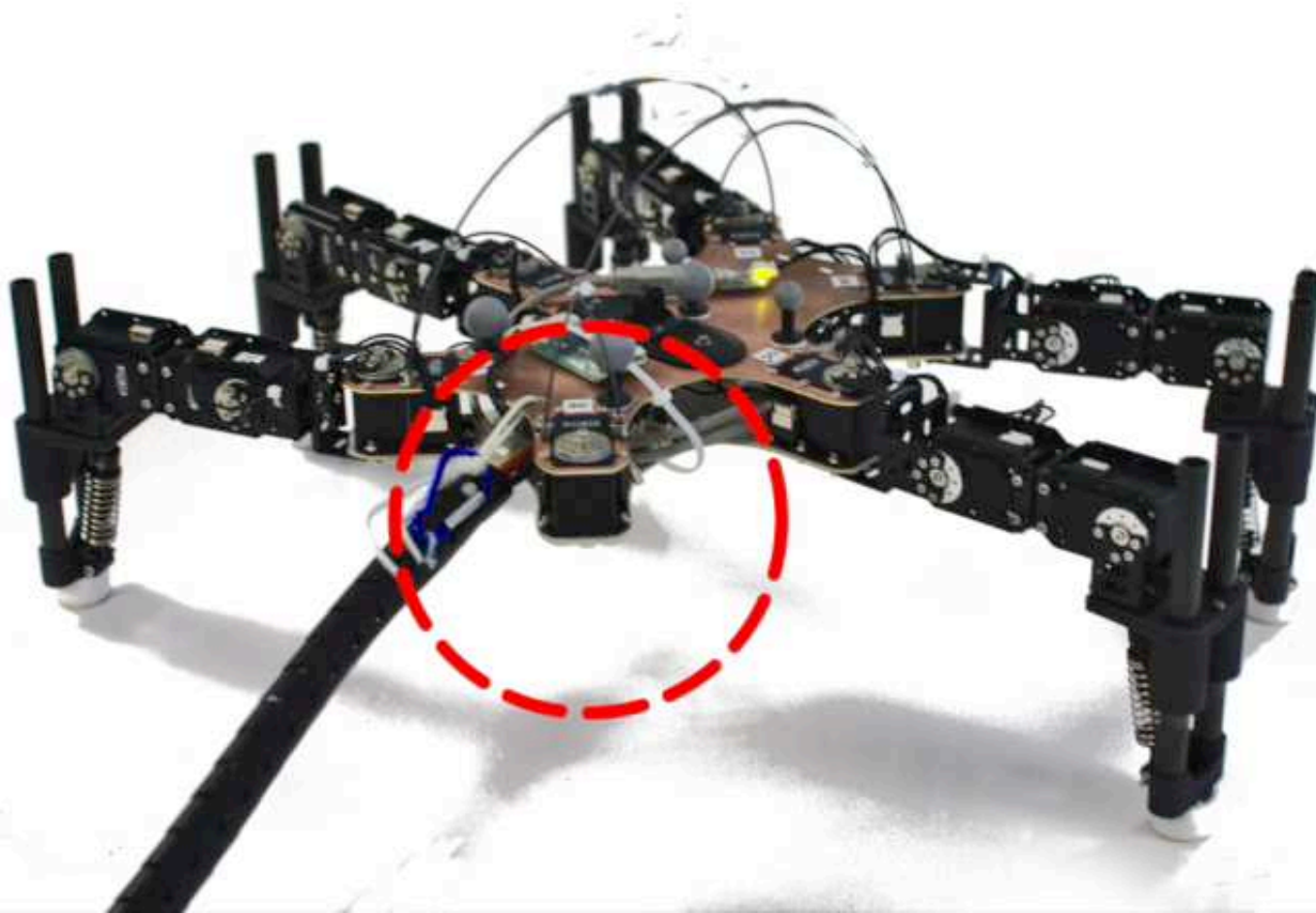
## Experiment 1 damage recovery

Unknown damage condition

105 priors (15 for each condition)

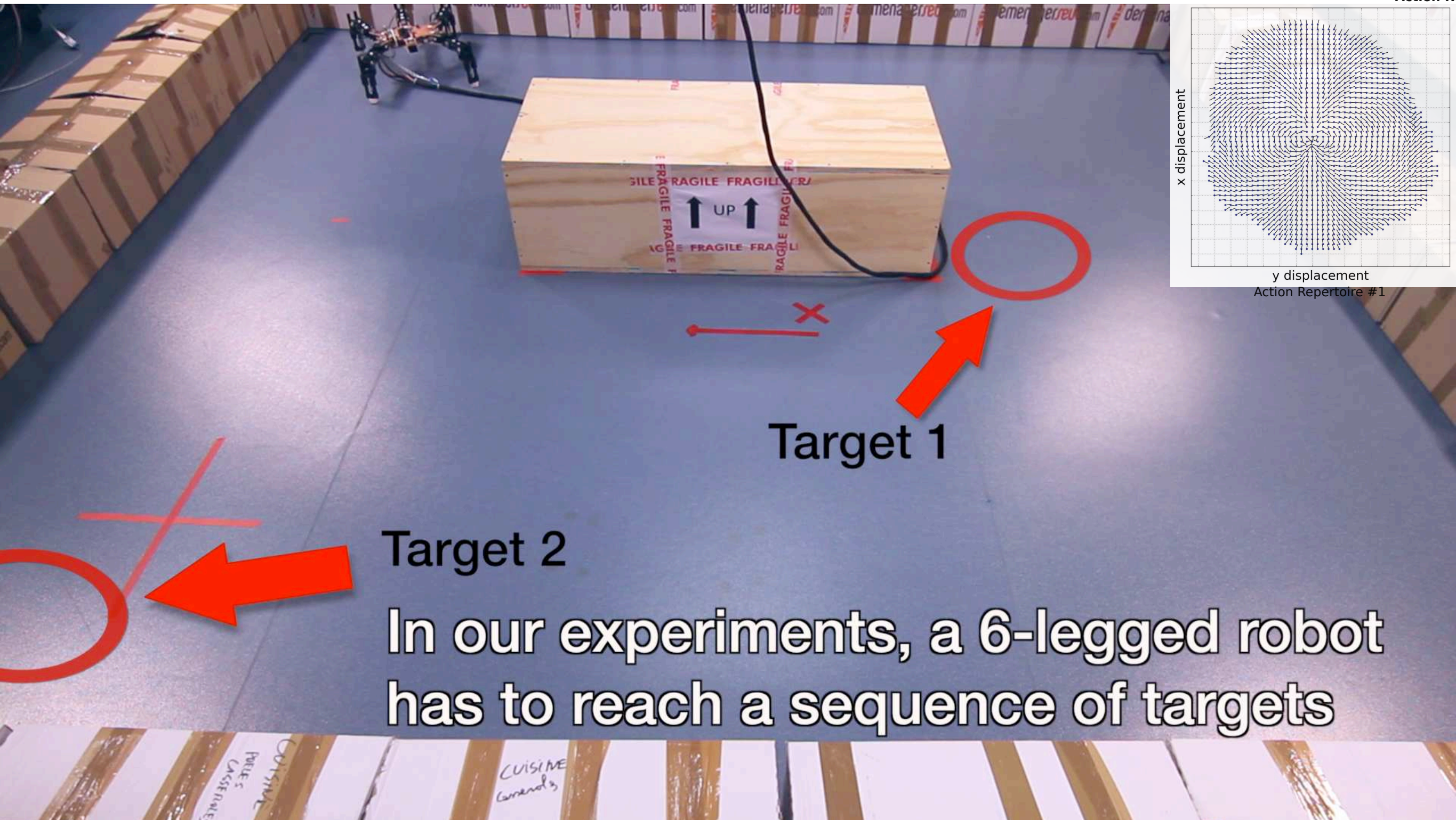
Reward: walking distance

BO with MLEI acquisition function





# Planning (MCTS) + repertoire learning (priors)

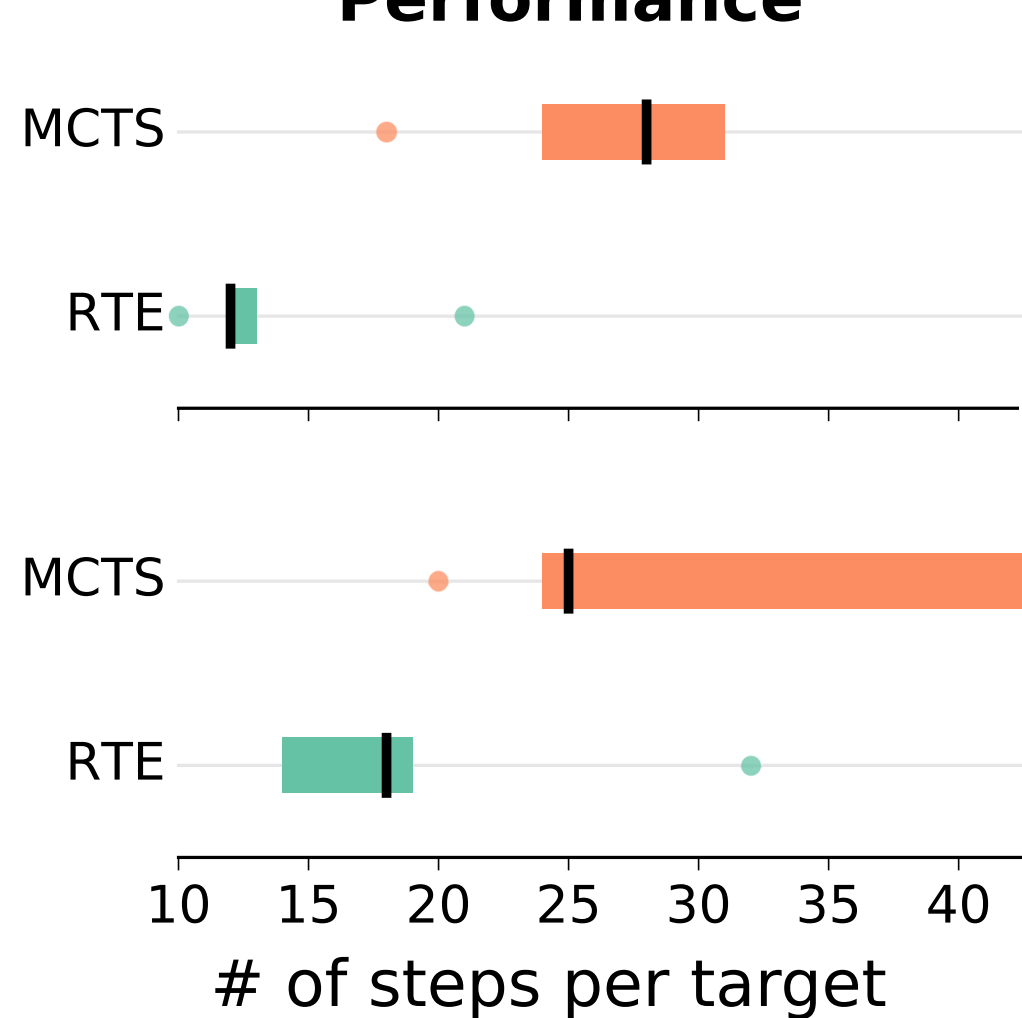


K. Chatzilygeroudis, V. Vassiliades, and J.-B. Mouret (2018). *Reset-free Trial-and-Error Learning for Data-Efficient Robot Damage Recovery*. Robotics and Autonomous Systems.



# Real robot / 5 replicates

## Performance



1 step = 3 seconds

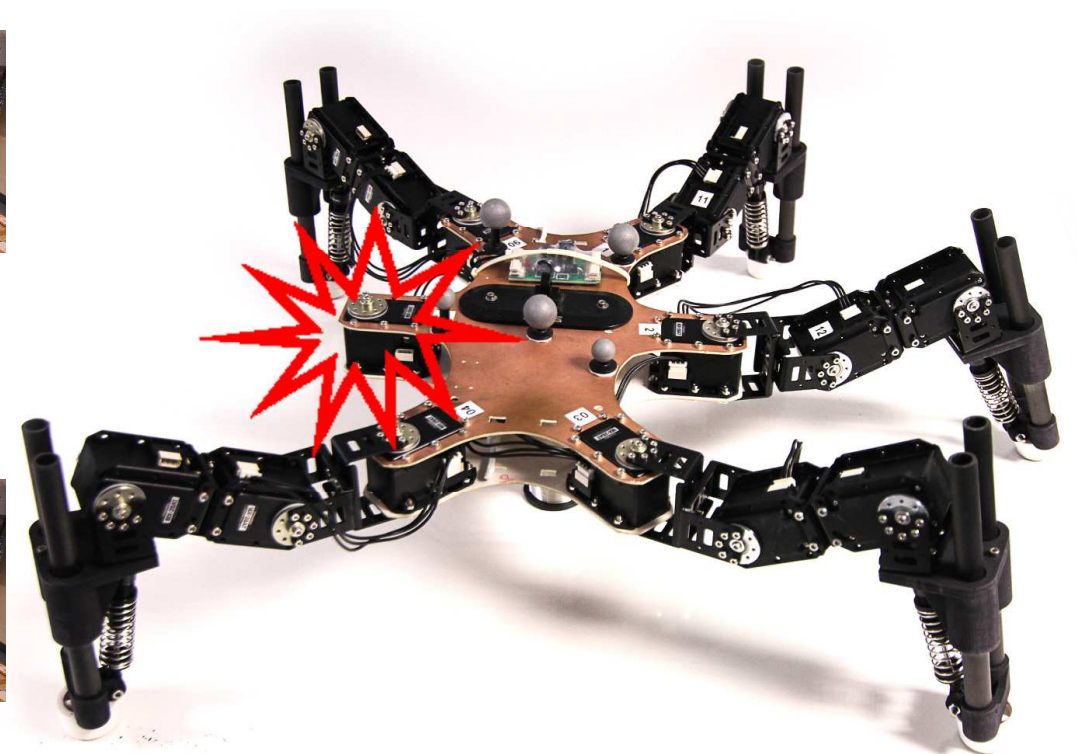
Environment #1



Environment #2



## Damaged Robot



# Bridging the gap with model-based Policy Search

- Transferability function = learning the limits of an existing simulator
- ... not far from a probabilistic model + prior

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

Gaussian processes

Simulator / model      parameters

Learning = maximize the likelihood of M+f

- ⇒ **effects that can be captured** by the simulator will be included by tuning the simulator (model identification)
- ⇒ **effects that cannot be captured** by changing the parameters are modelled by the Gaussian processes



# The Black-DROPS algorithm

Model-based  
policy search

## Black-box Data-Efficient Robot Policy Search

1. Perform a few random trials

⇒ new data

2. Learn a probabilistic model of the robot with Gaussian processes:

simulator ← Gaussian processes

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

3. Optimize with CMA-ES a policy that maximizes the long-term reward according to the model

⇒ one function evaluation = one rollout (propagate by sampling)

⇒ treat each rollout as a measurement of a noisy function

4. Evaluate the policy on the robot

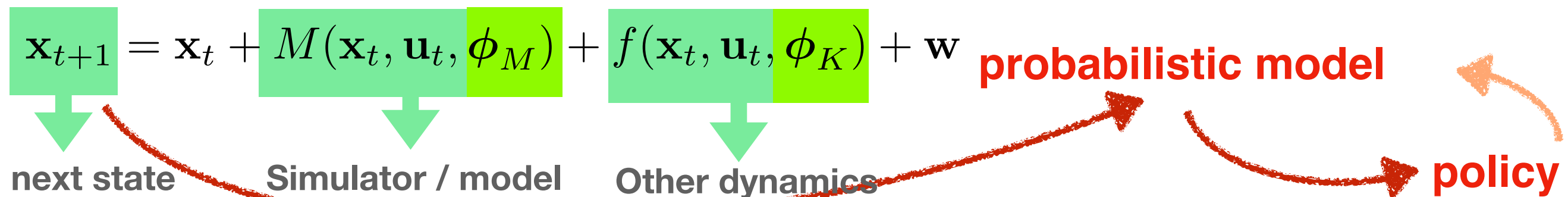
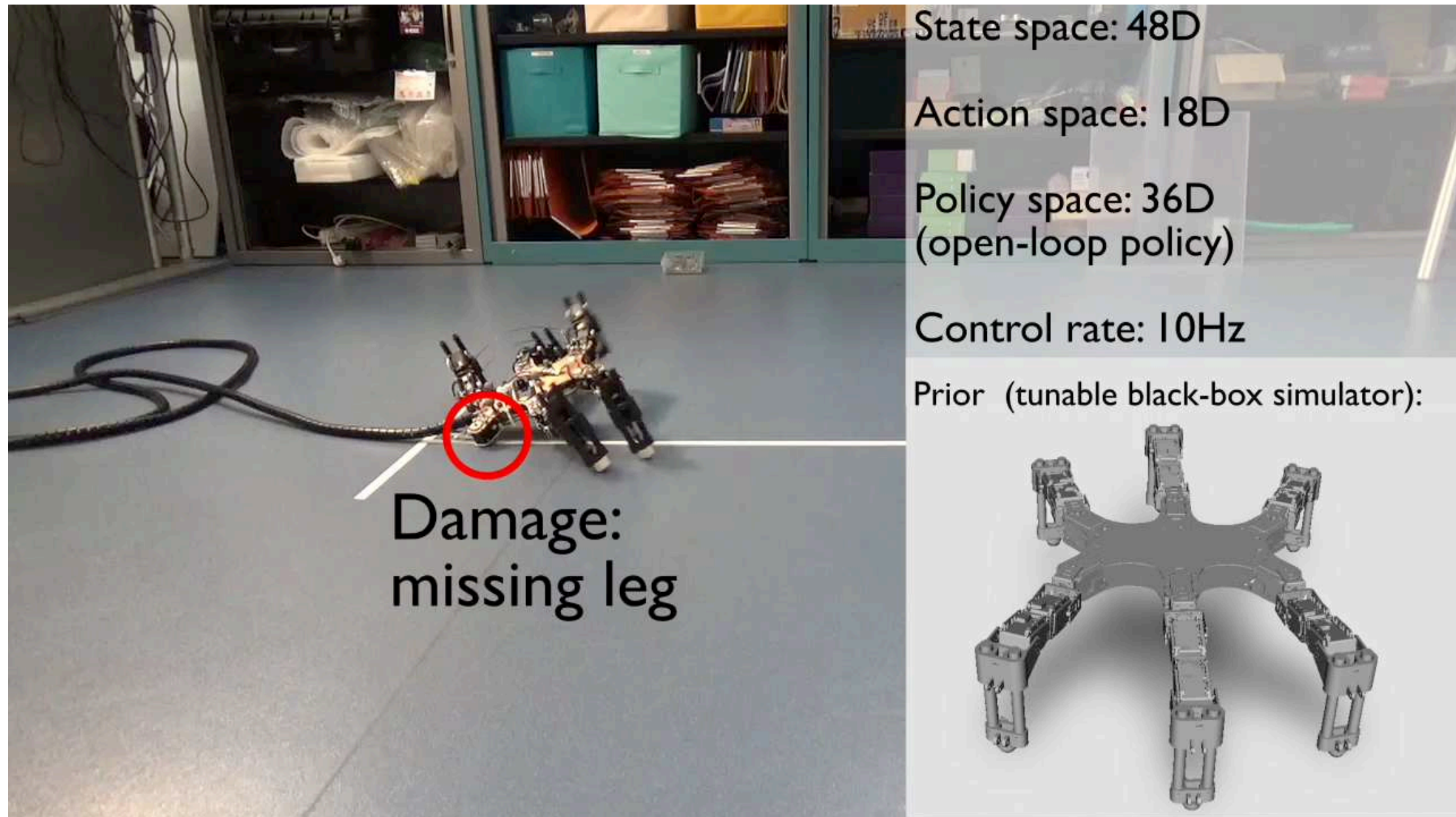
⇒ new data

Most expensive step:

- benefit a lot from parallelization
- handle uncertainty as noise ⇒ domain randomization

# Black-DROPS + priors + identification

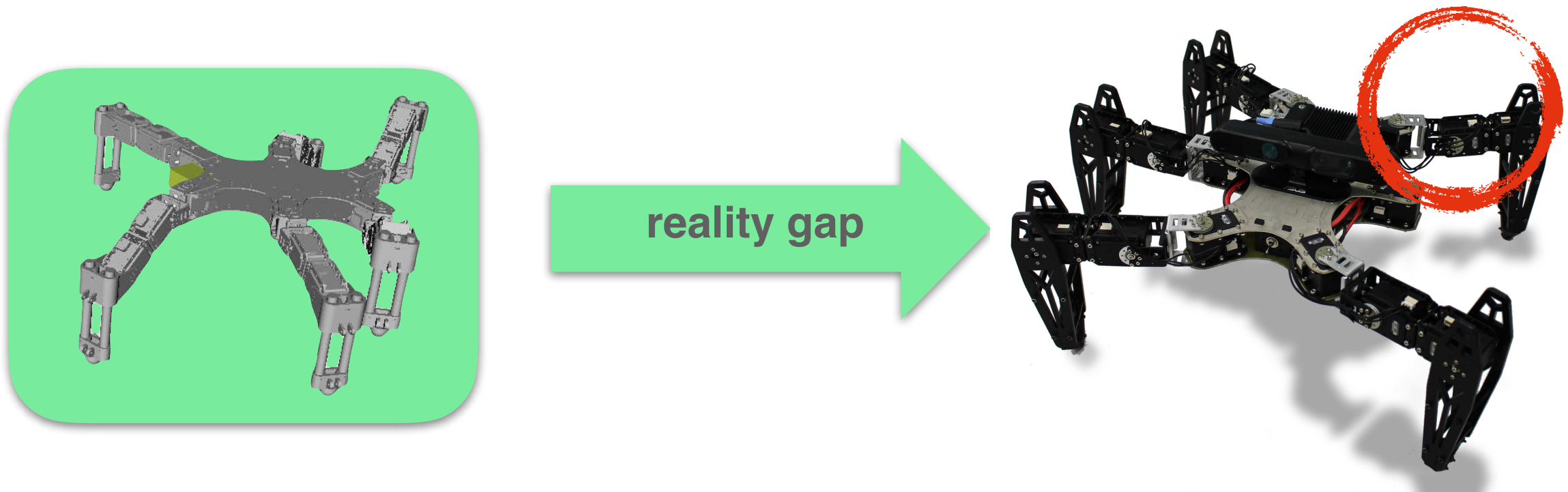
Model-based  
policy search





# Conclusions

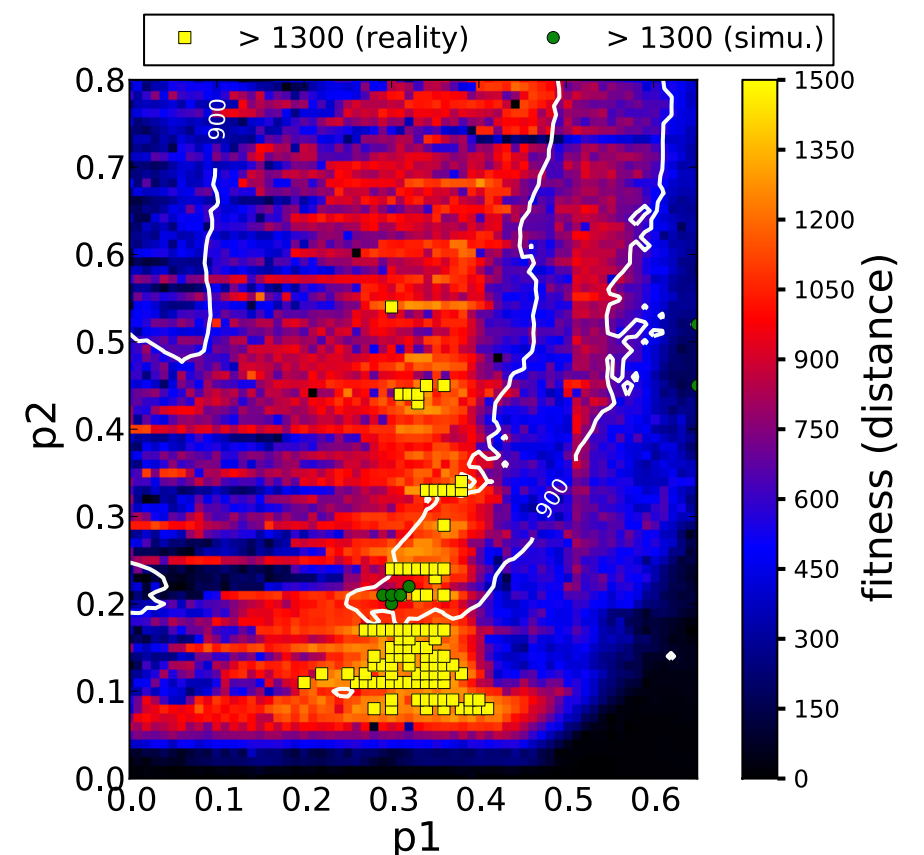
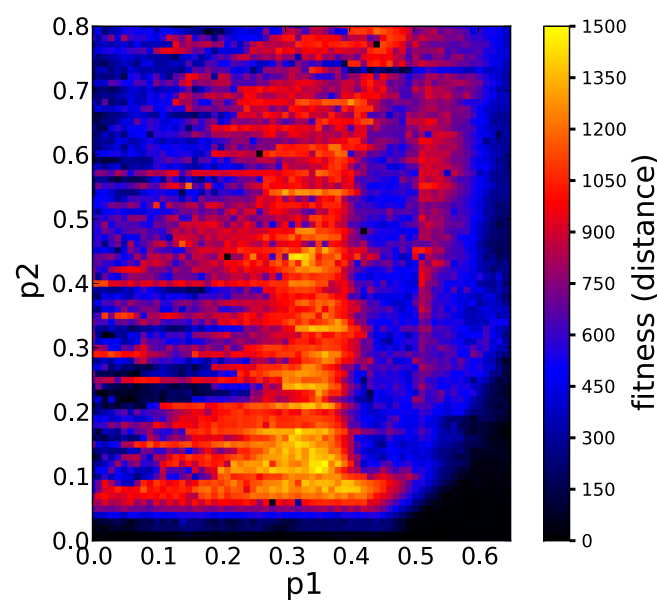
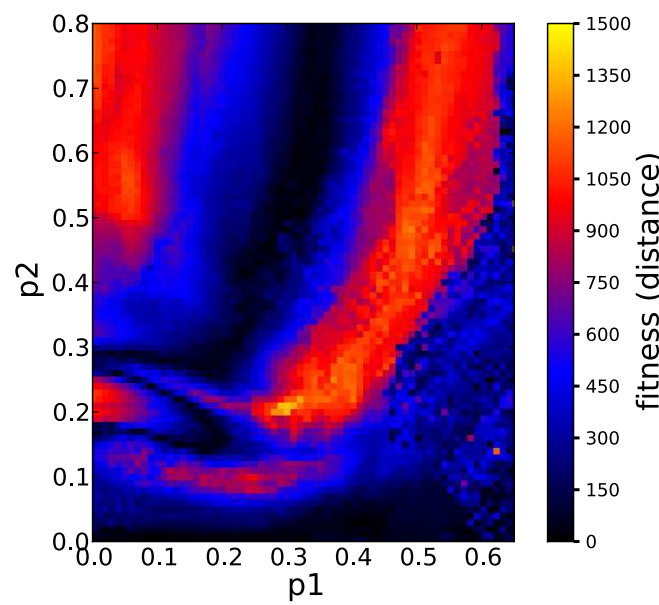
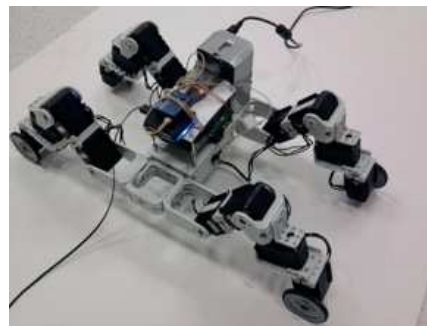
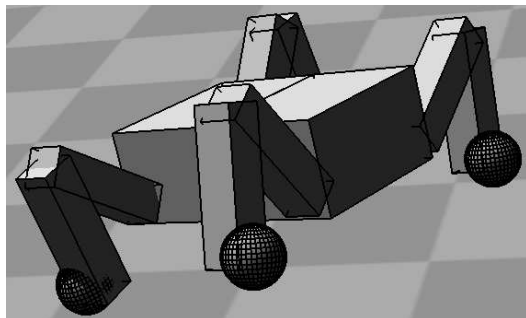
- **Adapting to damage is a reality gap problem**  
... and adaptation is critical for robotics  
... we know how to design controllers, but not how to adapt



- Diversity + Adaptation (IT&E): fast algorithm on the robot, but limited by the simulator
- Model-based Policy Search with prior and uncertainty (Black-DROPS): slow algorithm but could learn “anything”

# Conclusions (2)

- **Simulators are often right (for rigid bodies)!**
- **Simulators are good \*priors\***
- **Every simulator (model) prediction should come with a measure of transferability or uncertainty**
  - this can be learned from data
  - **crowd-source a model for bullet/dart/ode?**
- **We should map the reality gap**



Mouret, J. B., Koos, S., & Doncieux, S. (2013). Crossing the reality gap: a short introduction to the transferability approach. arXiv preprint arXiv: 1307.1870.



## Team

---

Vassilis Vassiliades  
K. Chatzilygeroudis  
Dorian Goepp  
Adam Gaier  
Brice Clément  
Rituraj Kaushik  
Jonathan Spitz  
Eloïse Dalin  
Pierre Desreumaux  
Vladislav Tempez  
Glenn Maguire  
Remi Pautrat

## Collaborators

---

Serena Ivaldi

[jean-baptiste.mouret@inria.fr](mailto:jean-baptiste.mouret@inria.fr)

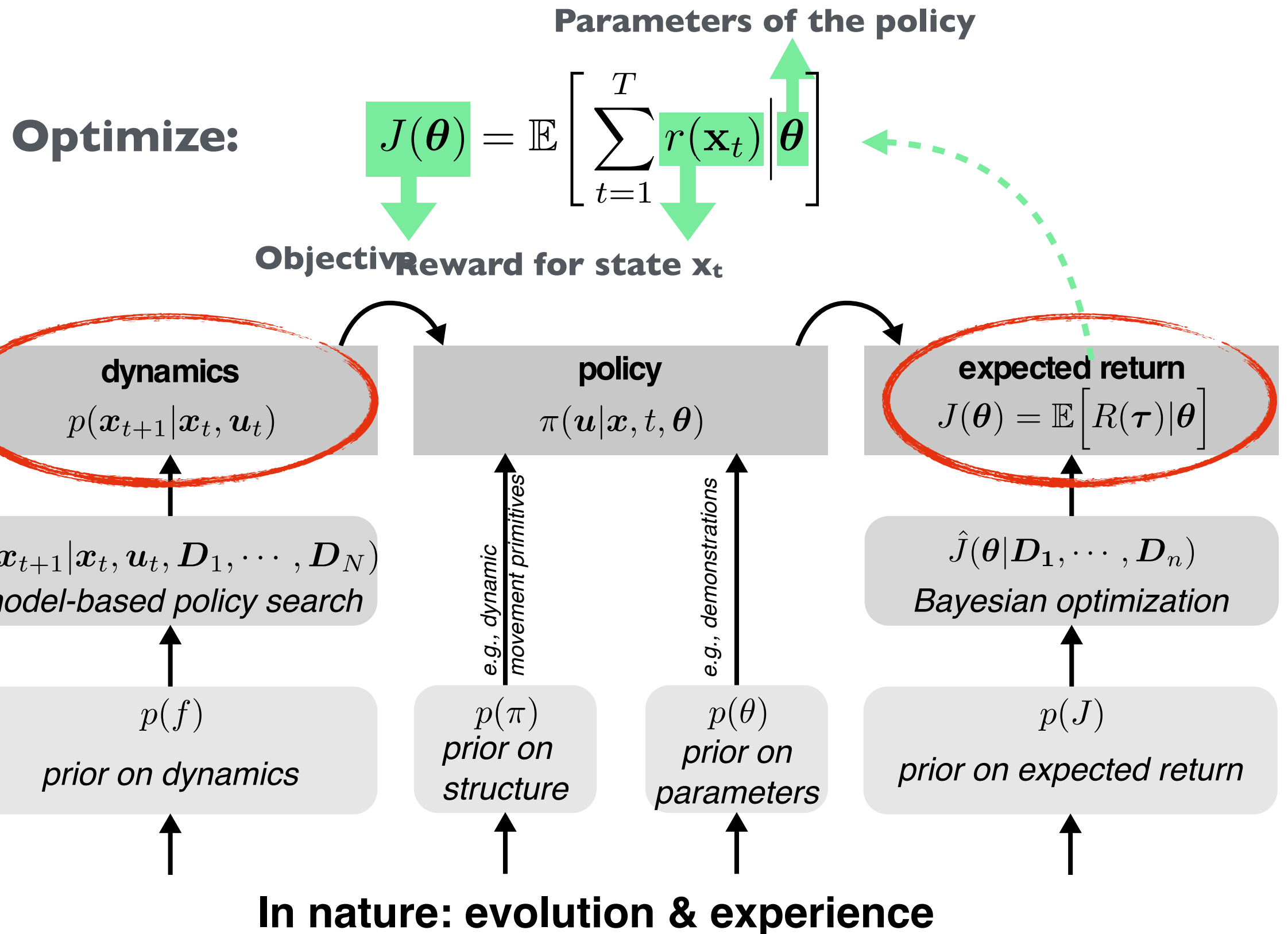
<https://members.loria.fr/jbmouret>



European Research Council

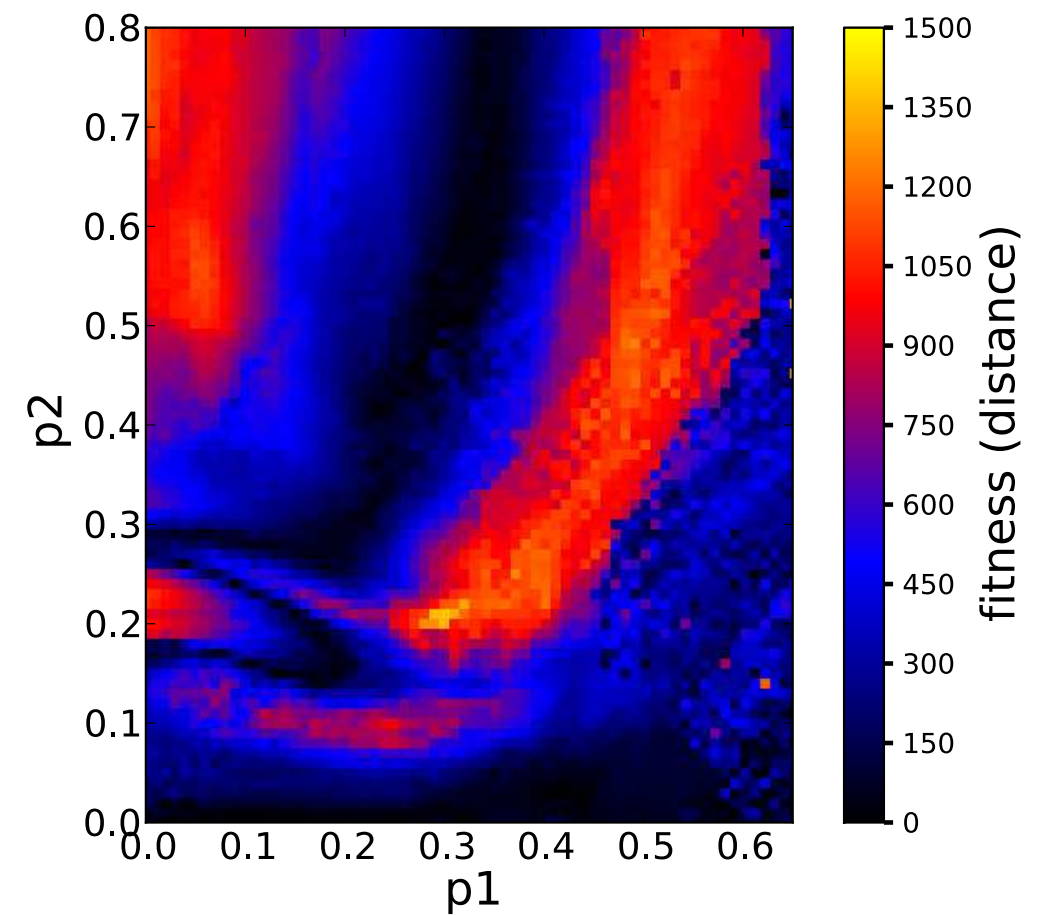
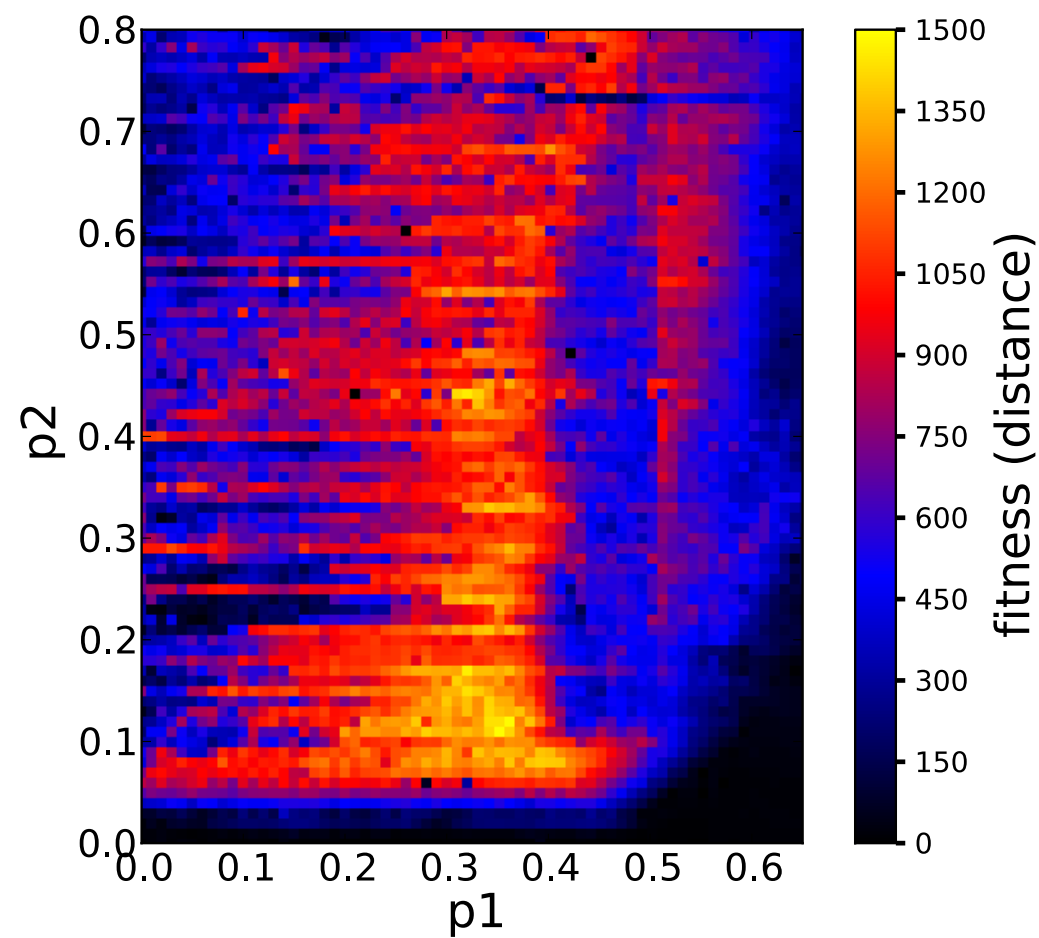
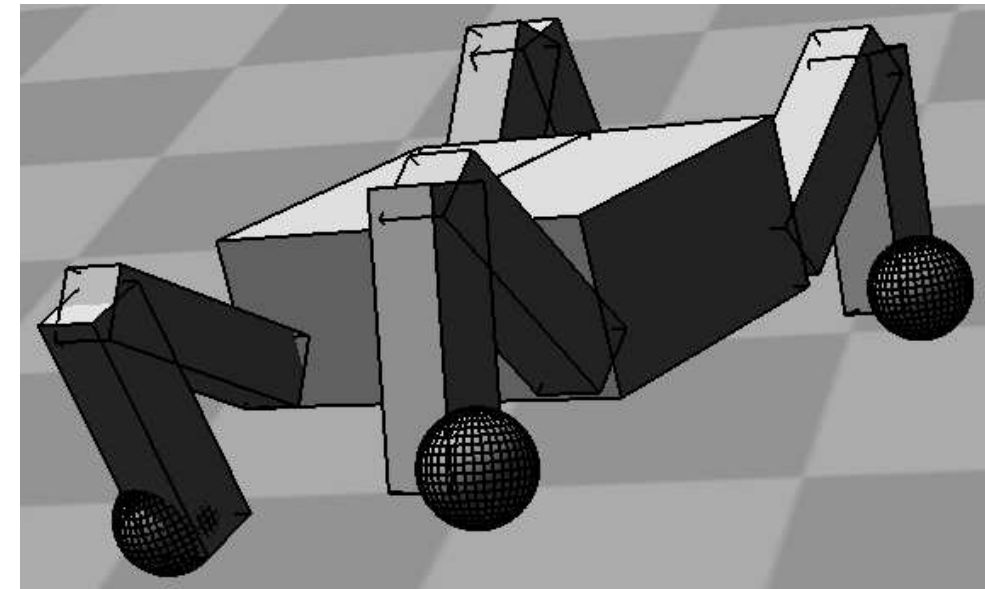
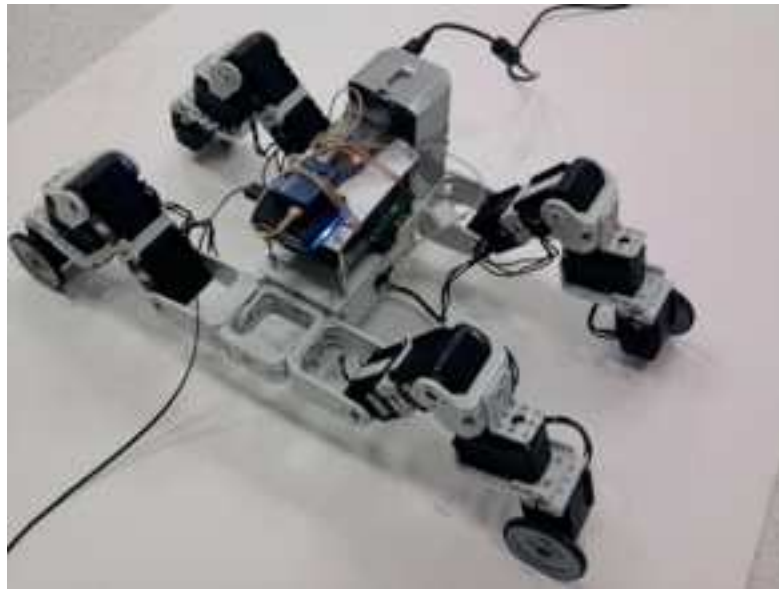
Established by the European Commission

# Micro-data policy search



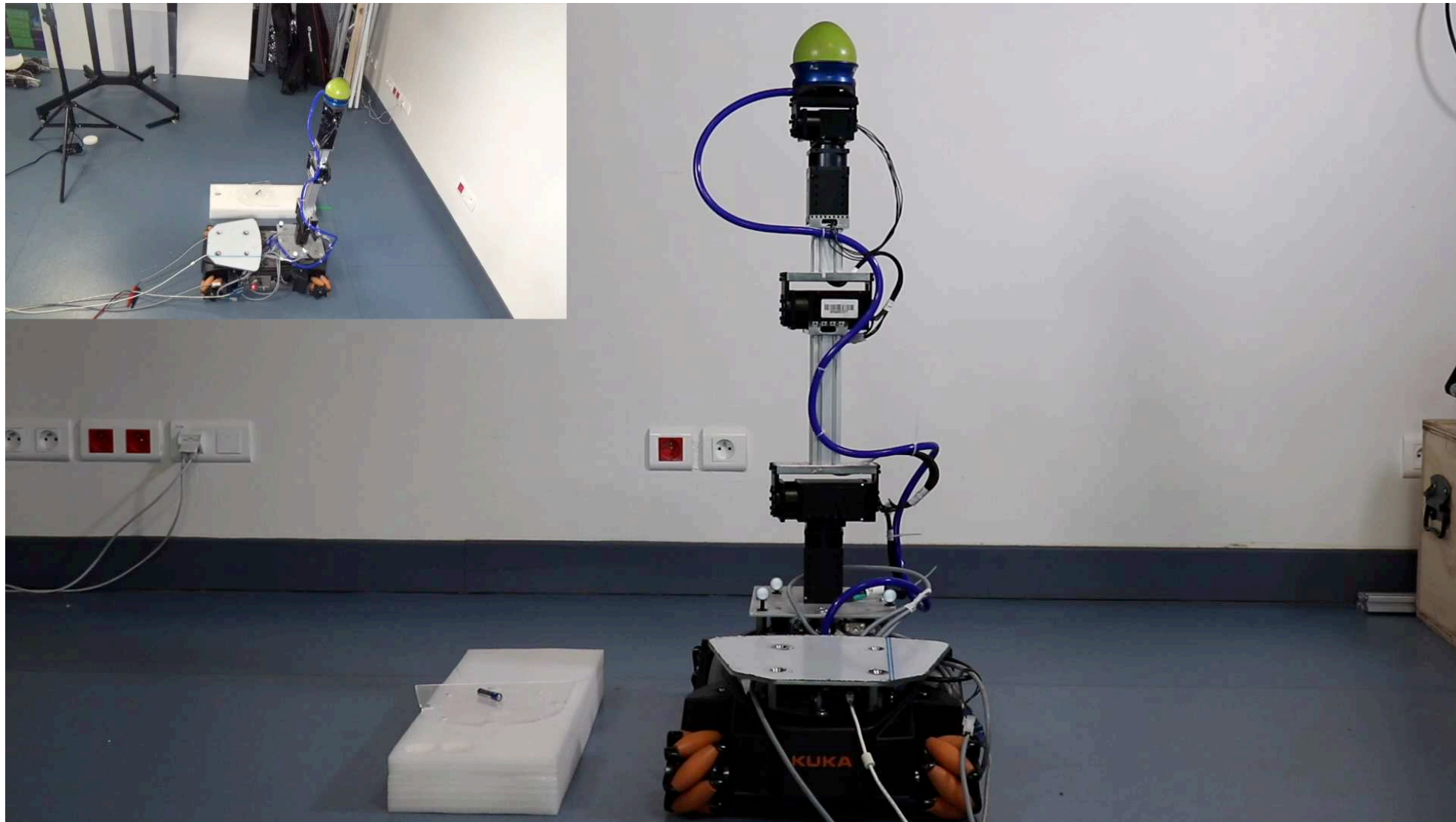
K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, J.-B. Mouret. (2018). A survey on policy search algorithms for learning robot controllers in a handful of trials. arXiv:1807.02303





Mouret, J. B., Koos, S., & Doncieux, S. (2013). Crossing the reality gap: a short introduction to the transferability approach. arXiv preprint arXiv:1307.1870.

# Strategy 1: Learning the dynamical model



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

Diagram illustrating the state transition equation:

- $\mathbf{x}_{t+1}$  (next)
- $\mathbf{x}_t$  (stat)
- $f(\mathbf{x}_t, \mathbf{u}_t)$  (dynamic)
- $\mathbf{w}$  (nois)

Probabilistic model

Optimal policy (in the model)

Chatzilygeroudis K, Rama R, Kaushik R, Goepp D, Vassiliades V, Mouret JB. (2017) Black-Box Data-efficient Policy Search for Robotics. Proc. of IEEE IROS.

Deisenroth, M. P., Fox, D., & Rasmussen, C. E. (2015). Gaussian processes for data-efficient learning in robotics and control. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2), 408-423.



# Scaling up model-based RL with priors

We cannot learn a model of a 6-legged robot

→ the state space is too large

We can add the simulator as a prior

→ we learn a marginal model (residual)

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \underbrace{M(\mathbf{x}_t, \mathbf{u}_t, \phi_M)}_{\text{Simulator / model}} + \underbrace{f(\mathbf{x}_t, \mathbf{u}_t, \phi_K)}_{\text{parameters}} + \mathbf{w}$$

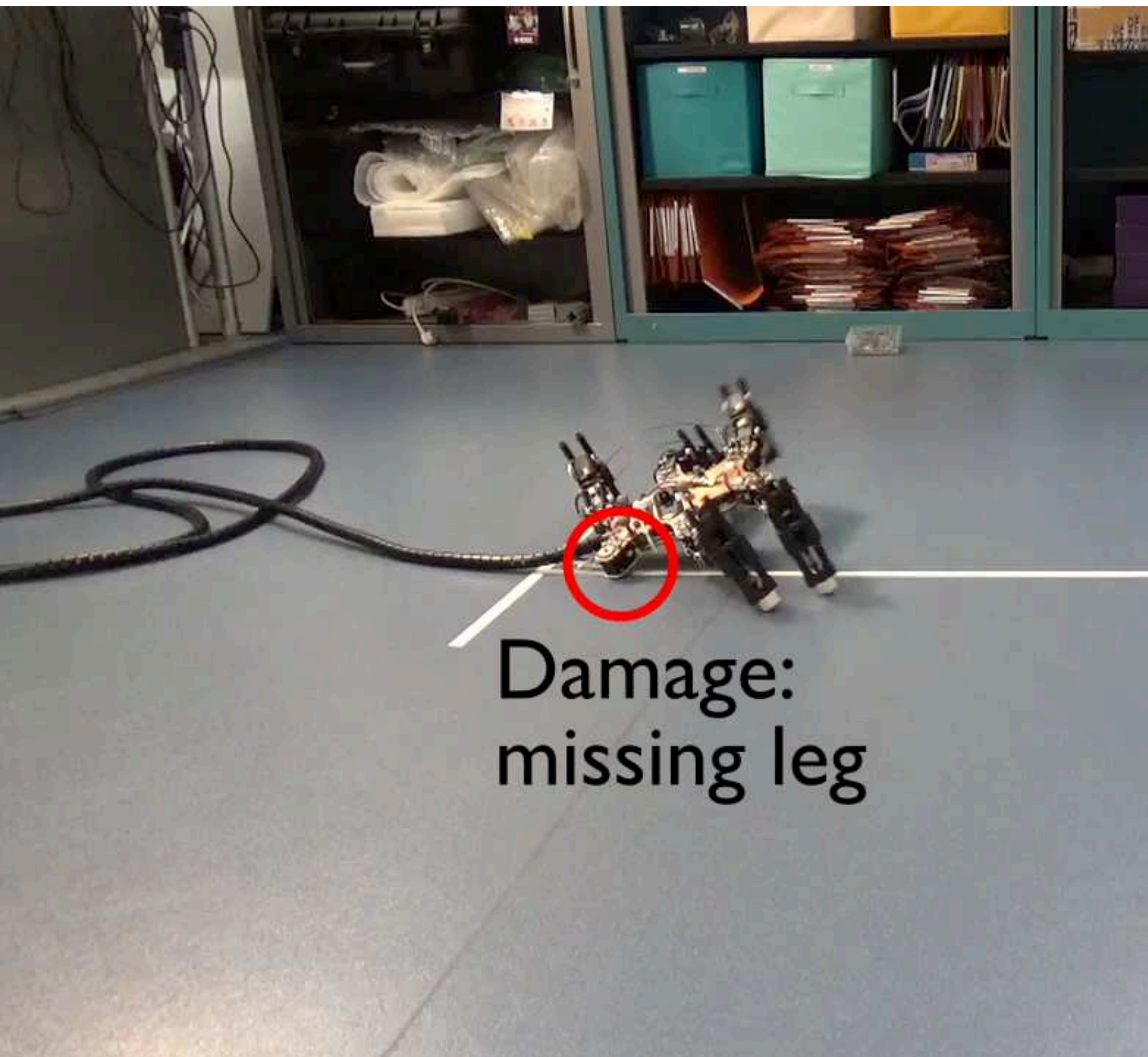
Gaussian processes

Learning = maximize the likelihood of M+f

We can combine model learning and model identification

- **effects that can be captured** by the simulator will be included by tuning the simulator (model identification)
- **effects that cannot be captured** by changing the parameters are modelled by the Gaussian processes

# Black-DROPS + priors + identification



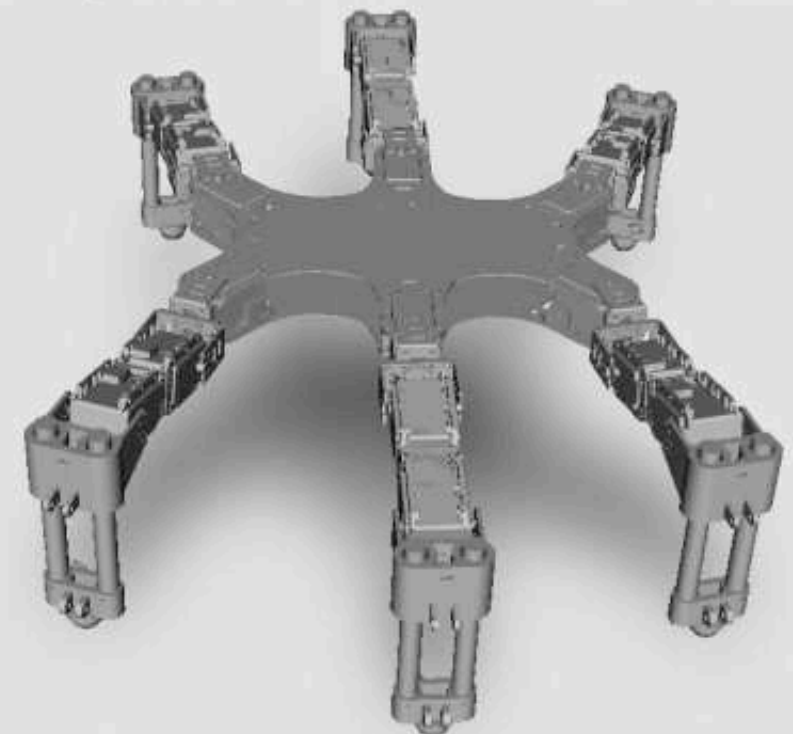
State space: 48D

Action space: 18D

Policy space: 36D  
(open-loop policy)

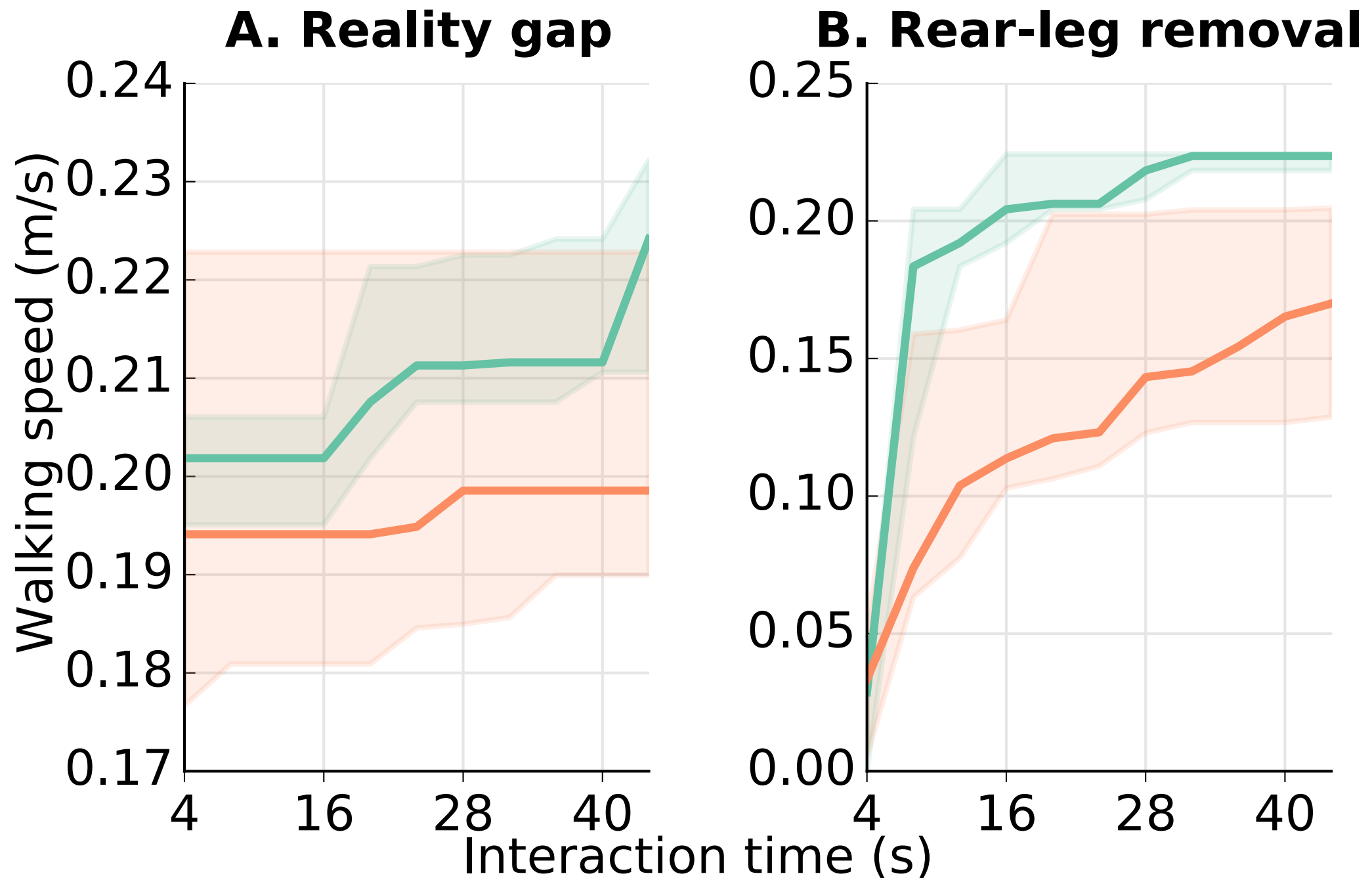
Control rate: 10Hz

Prior (tunable black-box simulator):





# Physical 6-legged robot



# Conclusions — Micro-Data learning

- We all want robots that can learn. But why? Is it *really* needed?

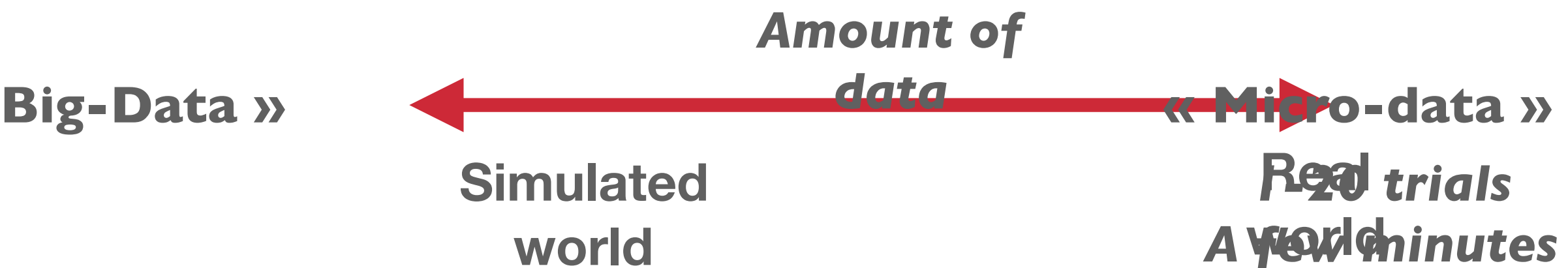
⇒ Learning in robotics is most useful for online adaptation

⇒ Damage recovery: the “killer app” for robot learning?

**... but robots need to learn in a few minutes (micro-data)**

⇒ priors: generic for robotics? automatic generation?

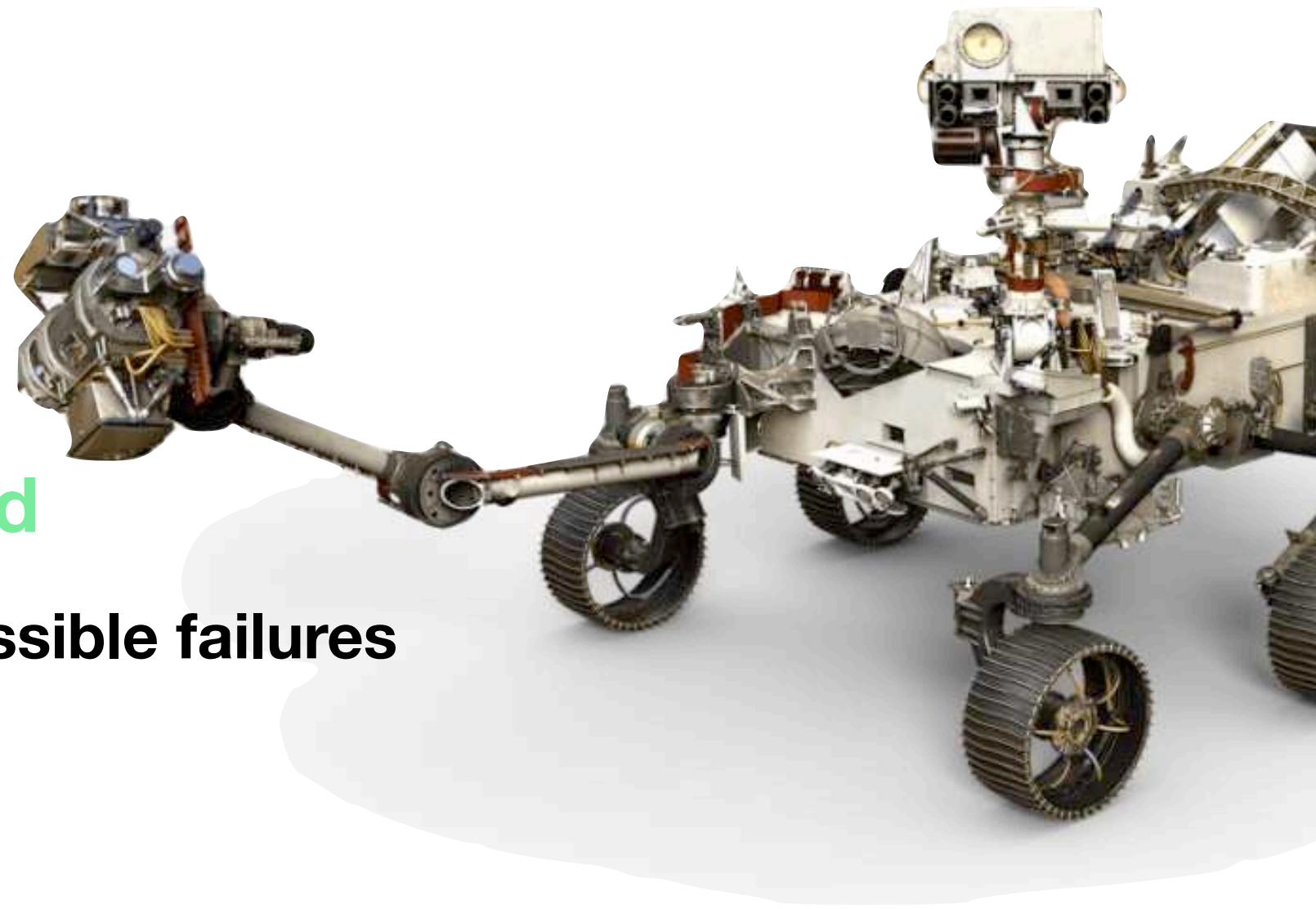
⇒ *probabilistic* models: policies robust to inaccuracies





... diagnosis is hard  
... robustness is hard

➔ need to anticipate possible failures



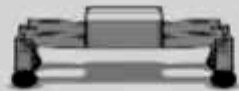
Could trial-and-error learning offer an alternative?

- ➔ no need to understand the damage to find a compensatory behaviour
- ➔ “takes a shortcut”

Where could “AI” could make a difference in robotics?

## subtitle

Control approach - 1 objective:  
covered distance



**in simulation: 1200 mm in 10 seconds**

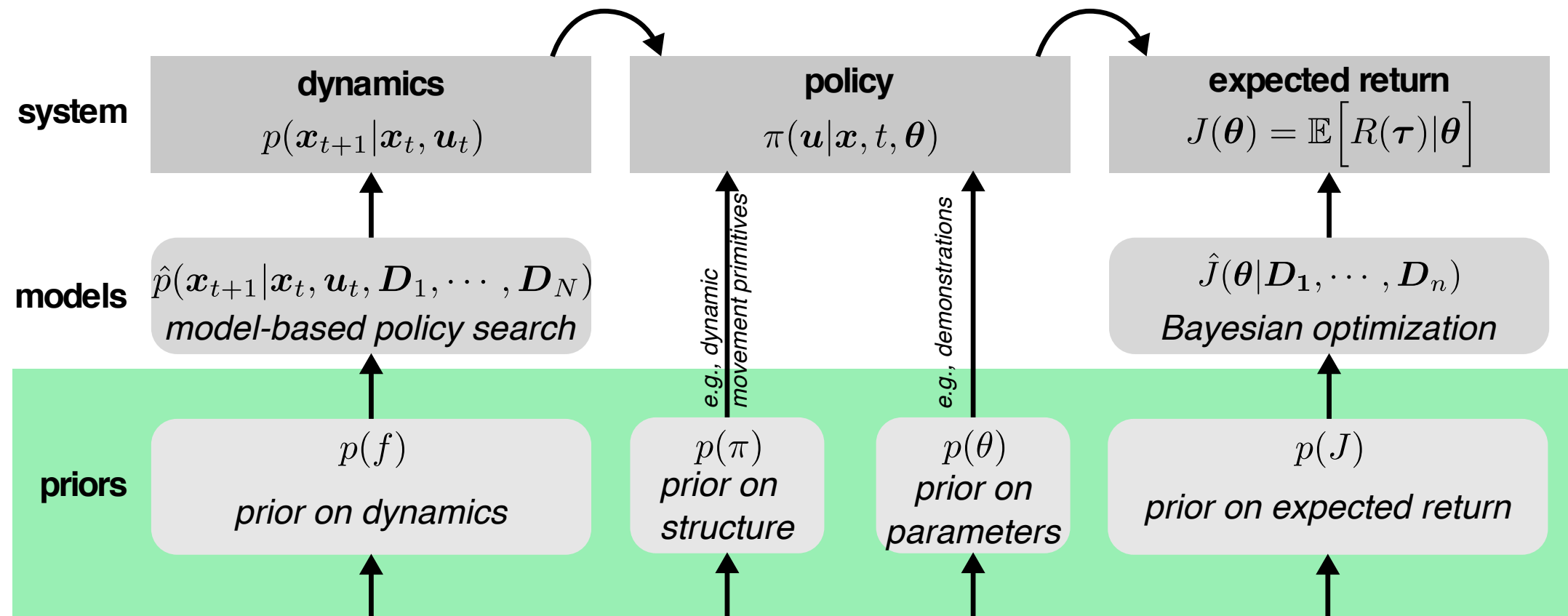
Control approach - 1 objective:  
covered distance



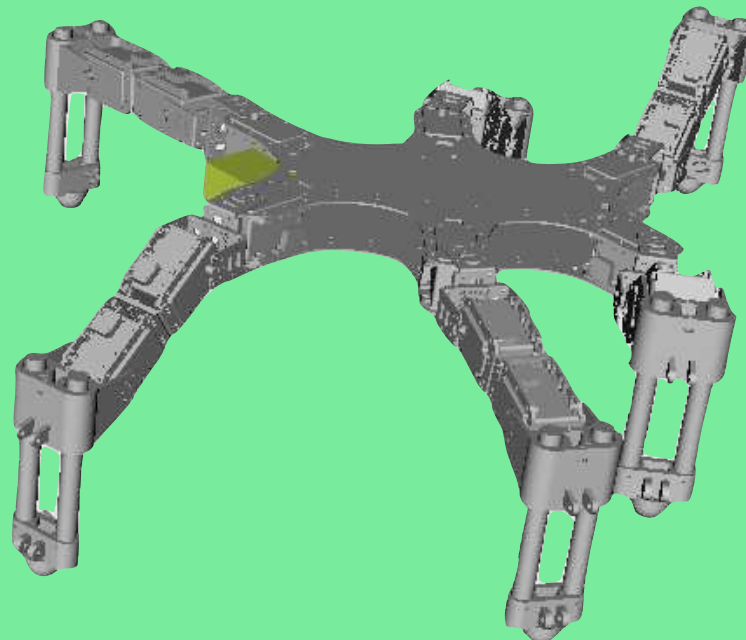
**in reality : 277 mm in 10 seconds**



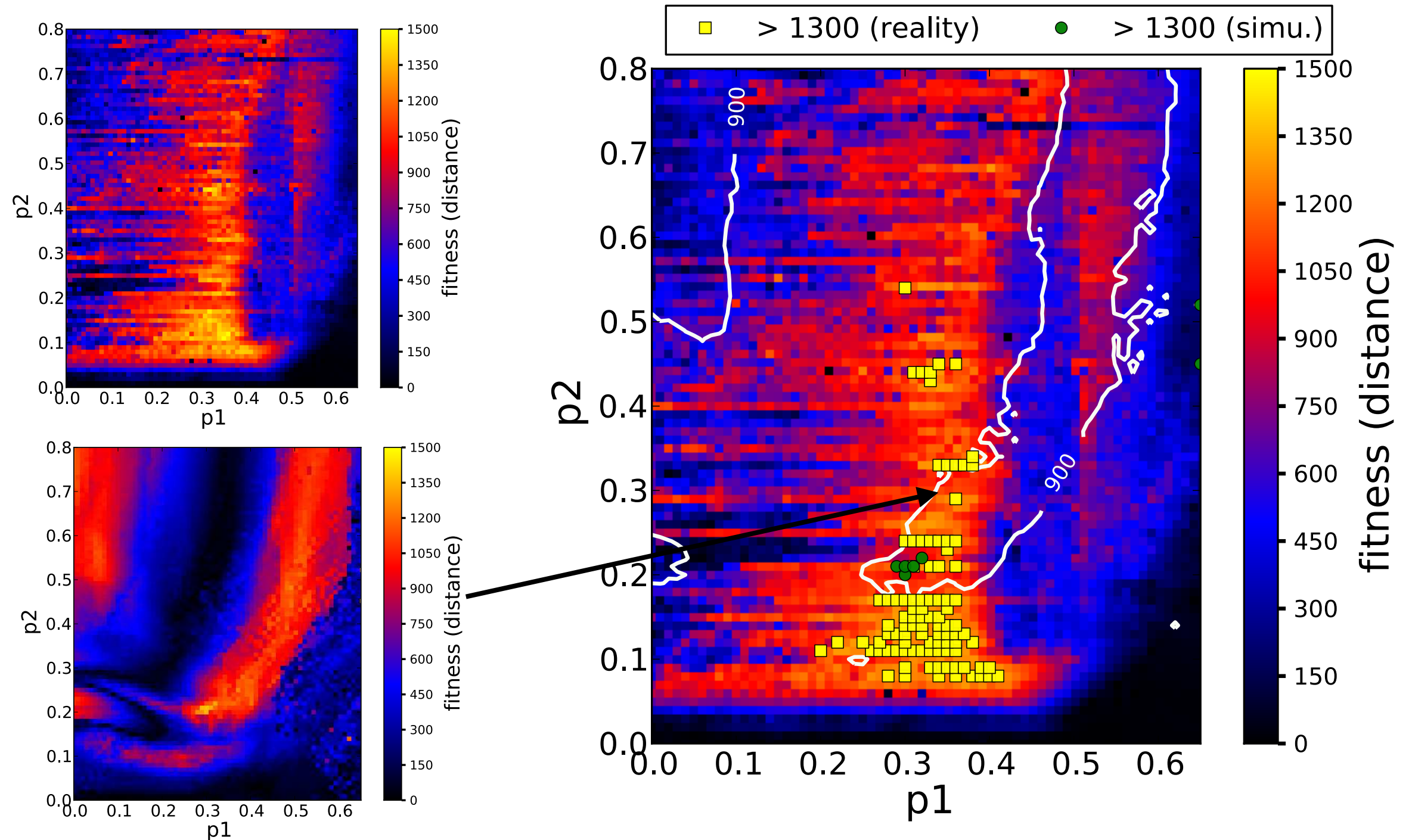
# Policy Search for damage recovery & adaptation



Model of the intact robot



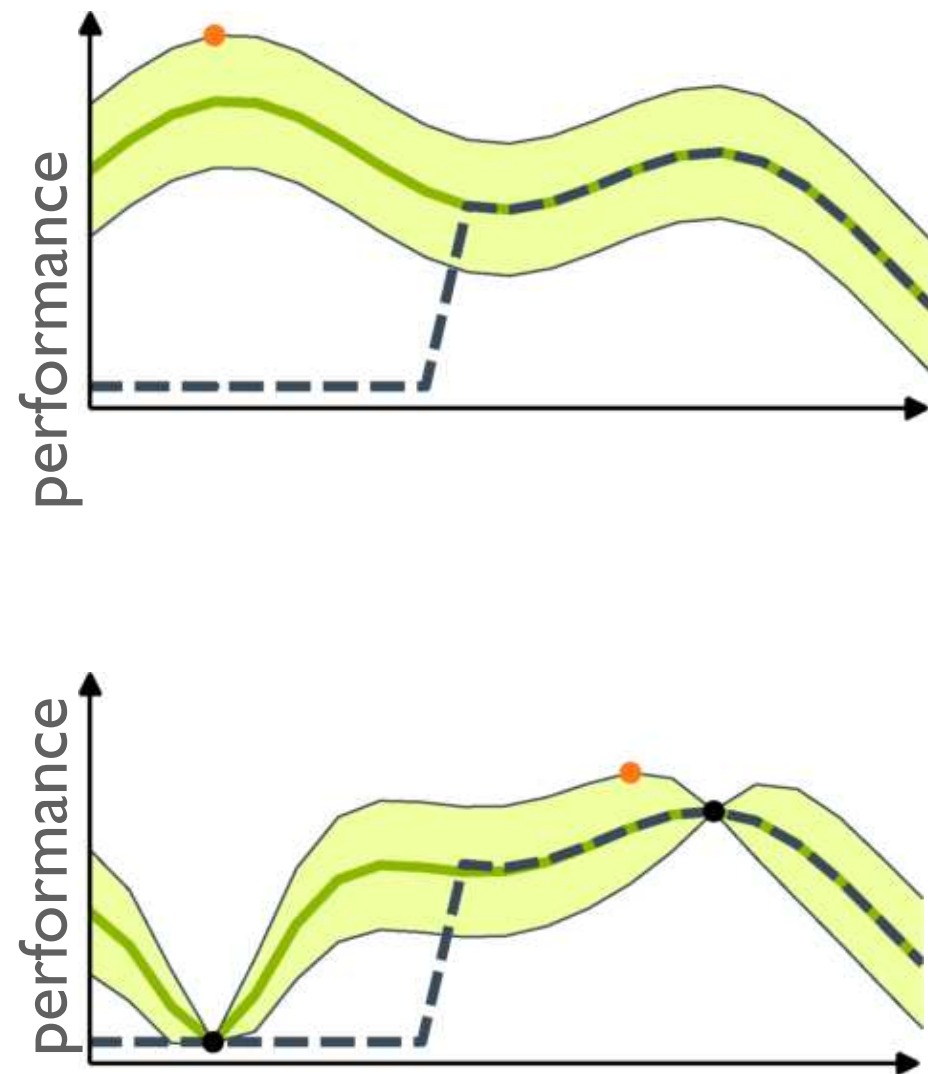
# But they can agree (sometimes)!





# BO + MAP-Elites

## “Intelligent Trial and Error”



$$P(f(\mathbf{x})|\mathbf{P}_{1:t+1}, \mathbf{x}) = \mathcal{N}(\mu_{t+1}(\mathbf{x}), \sigma_{t+1}^2(\mathbf{x}))$$

where

$$\mu_{t+1}(\mathbf{x}) = \mathcal{A}(\mathbf{x}) + \mathbf{k}^t \mathbf{K}^{-1} (\mathbf{P}_{1:t+1} - \mathcal{A}(\mathbf{y}_{1:t+1}))$$

$$\sigma_{t+1}^2(\mathbf{x}) = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}^t \mathbf{K}^{-1} \mathbf{k}$$

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{y}_1, \mathbf{y}_1) + \sigma_{noise}^2 & \cdots & k(\mathbf{y}_1, \mathbf{y}_t) \\ \vdots & \ddots & \vdots \\ k(\mathbf{y}_t, \mathbf{y}_1) & \cdots & k(\mathbf{y}_t, \mathbf{y}_t) + \sigma_{noise}^2 \end{bmatrix}$$

$$\mathbf{k} = \begin{bmatrix} k(\mathbf{x}, \mathbf{y}_1) & k(\mathbf{x}, \mathbf{y}_2) & \cdots & k(\mathbf{x}, \mathbf{y}_t) \end{bmatrix}$$

# Automatic selection of priors

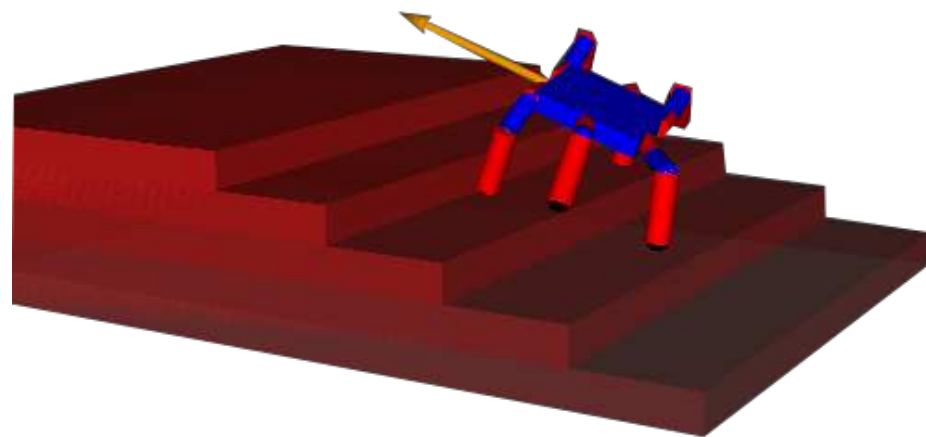
- **Goal:** relax the need to have “the right prior”
- **Concept:**
  - Generate / create many priors (e.g. 100)
  - The likelihood of a prior can be computed (prior + obs.)
  - Choose after each episode
    - prior that is likely given the optimization
    - prior that can help the optimization
- **New acquisition function:** Most likely Expected Improvement

$$\text{EIP}(\mathbf{x}, \mathcal{P}) = \text{EI}(\mathbf{x}) \times P(\mathbf{f}(\mathbf{x}_{1..t}) \mid \mathbf{x}_{1..t}, \mathcal{P}(\mathbf{x}_{1..t}))$$

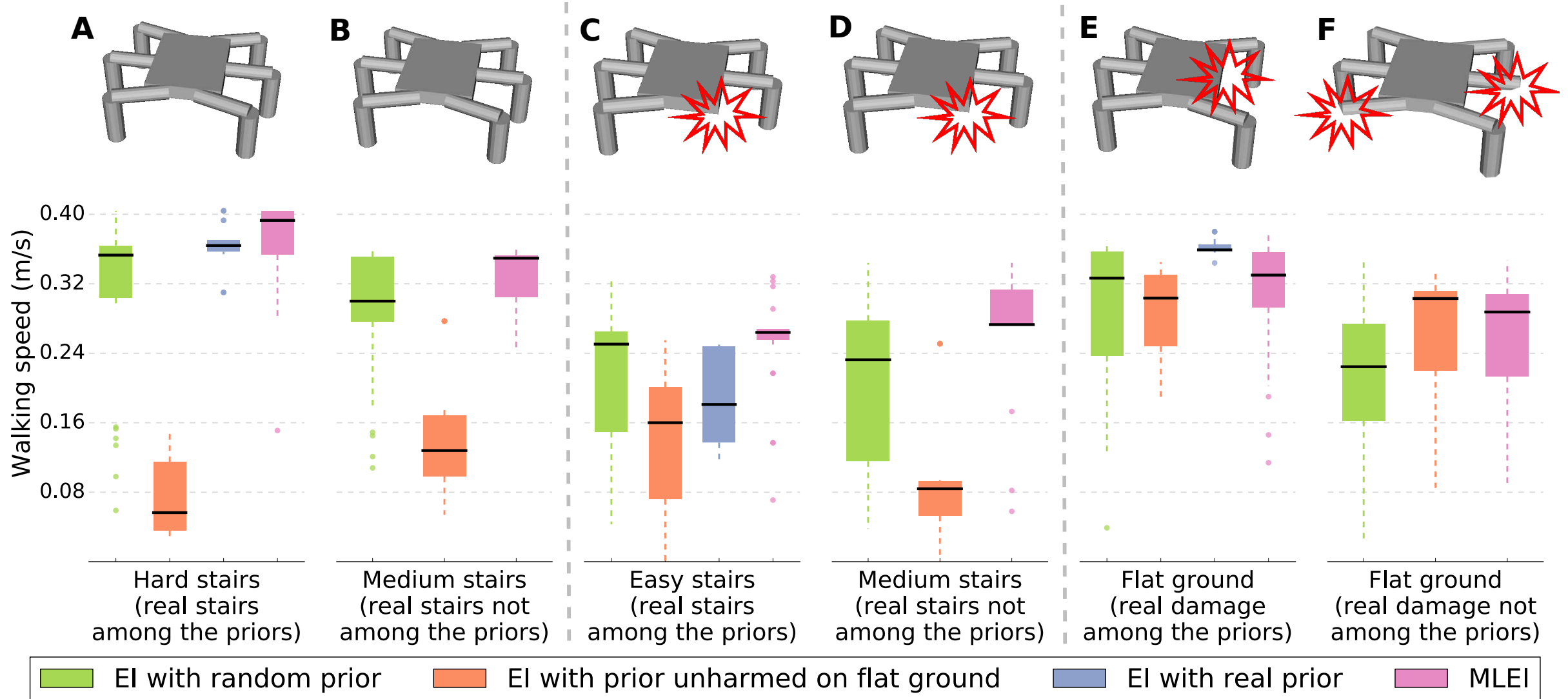
$$\text{MLEI}(\mathbf{x}, \mathcal{P}_1, \dots, \mathcal{P}_m) = \max_{p \in \mathcal{P}_1, \dots, \mathcal{P}_m} \text{EIP}(\mathbf{x}, p)$$



# Damage recovery on stairs

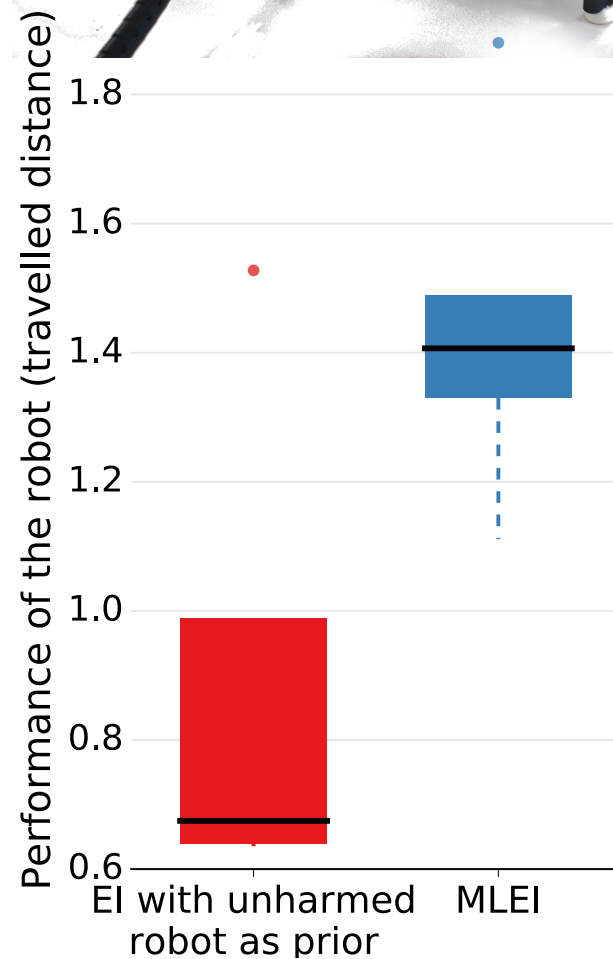
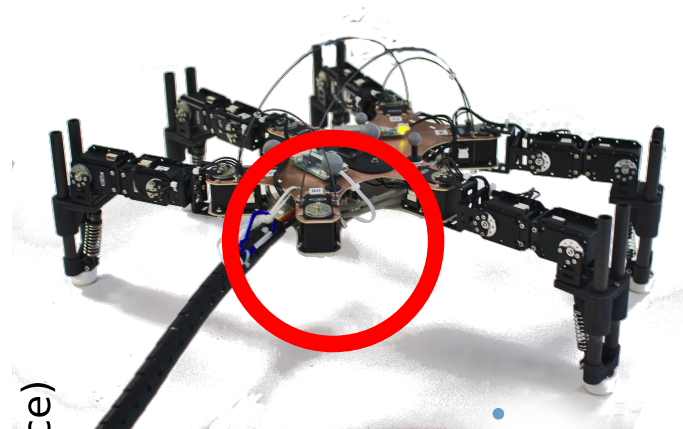


15 priors for 4 environments (stairs) = 60  
 15 priors for 6 damage conditions + intact = 105  
 5 iterations  
 30 replicates

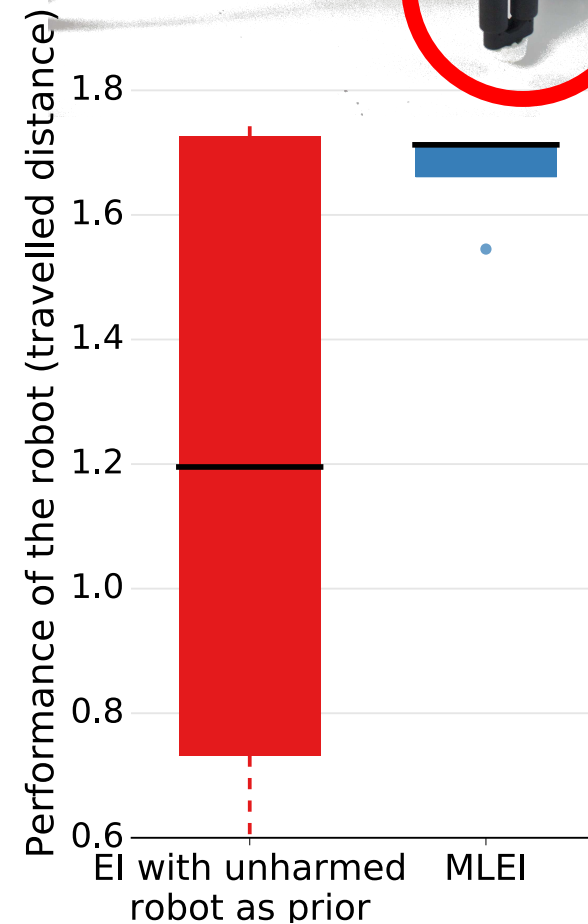
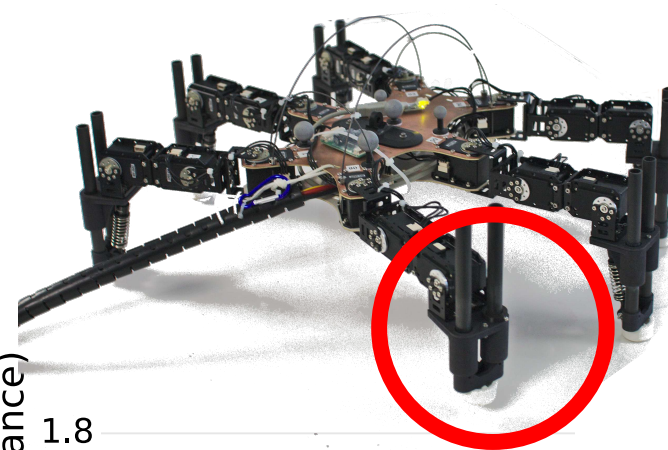


# Damage recovery

10 episodes / 5 replicates / 15 priors for 6 damage conditions + intact = 105



*damage among  
the priors*



*damage NOT among  
the priors*



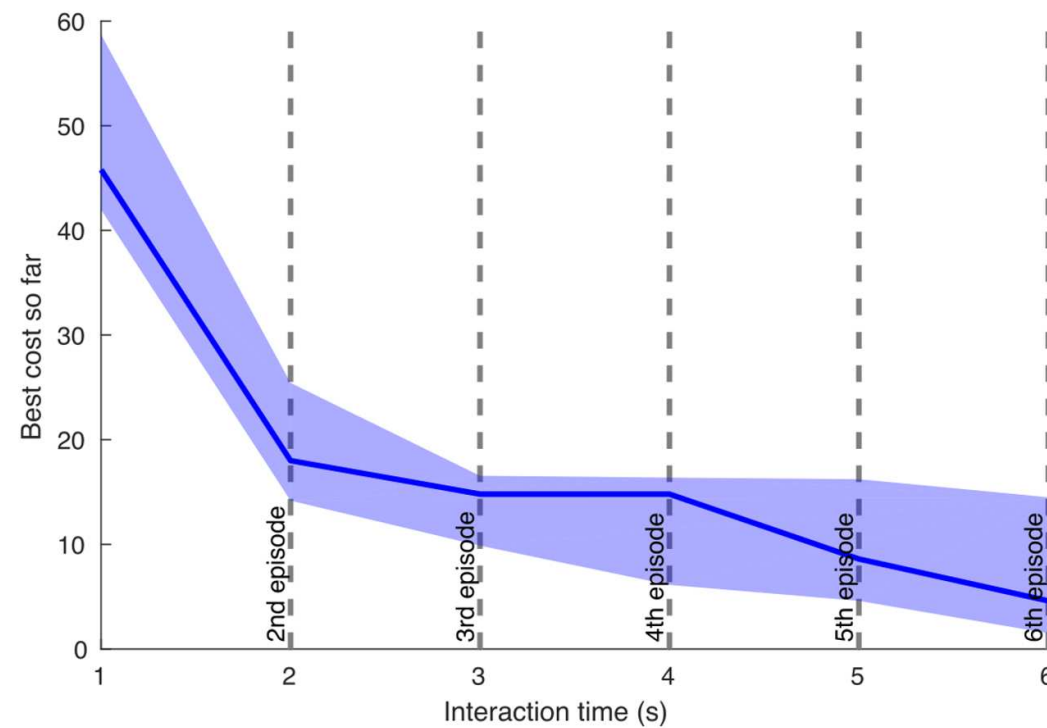
## "real" physical world



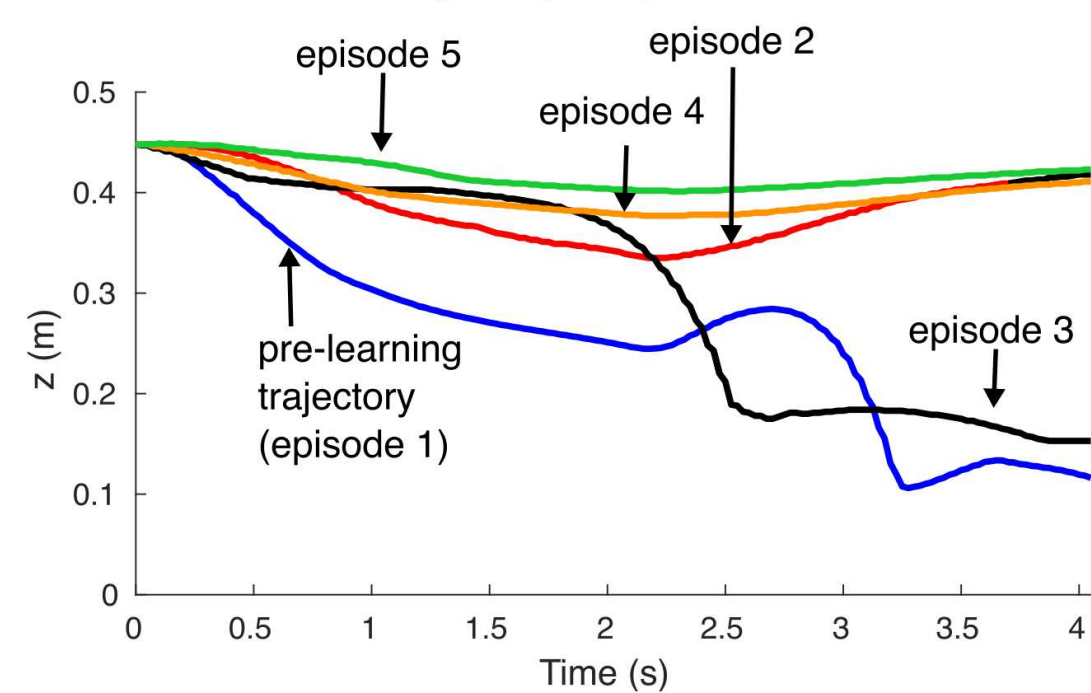
## model for the controller



**Cost of CoM Task vs Interaction time**



**Real CoM height trajectory for each episode**



Model for the controller ("QP world")



big feet

Real world

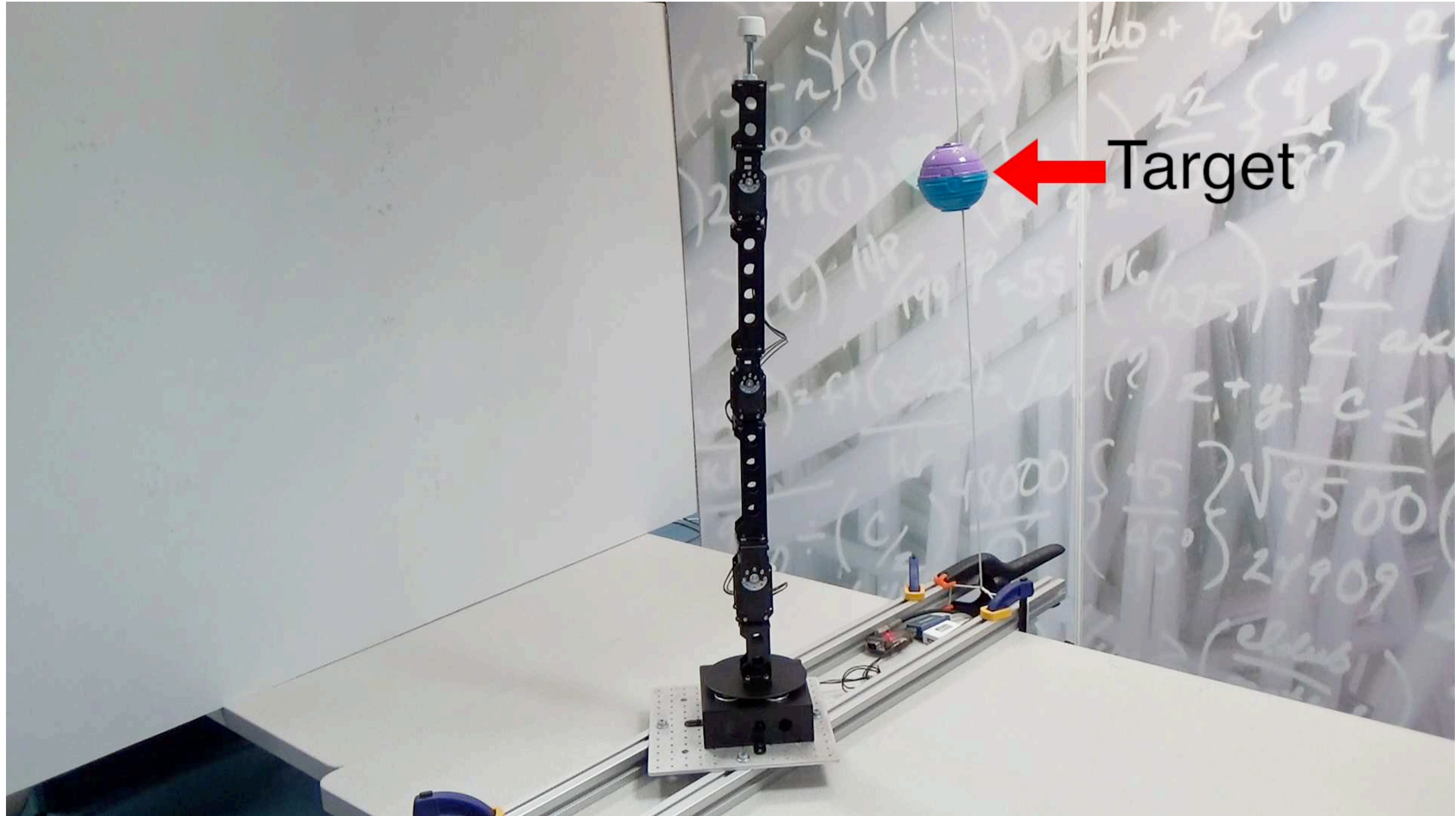


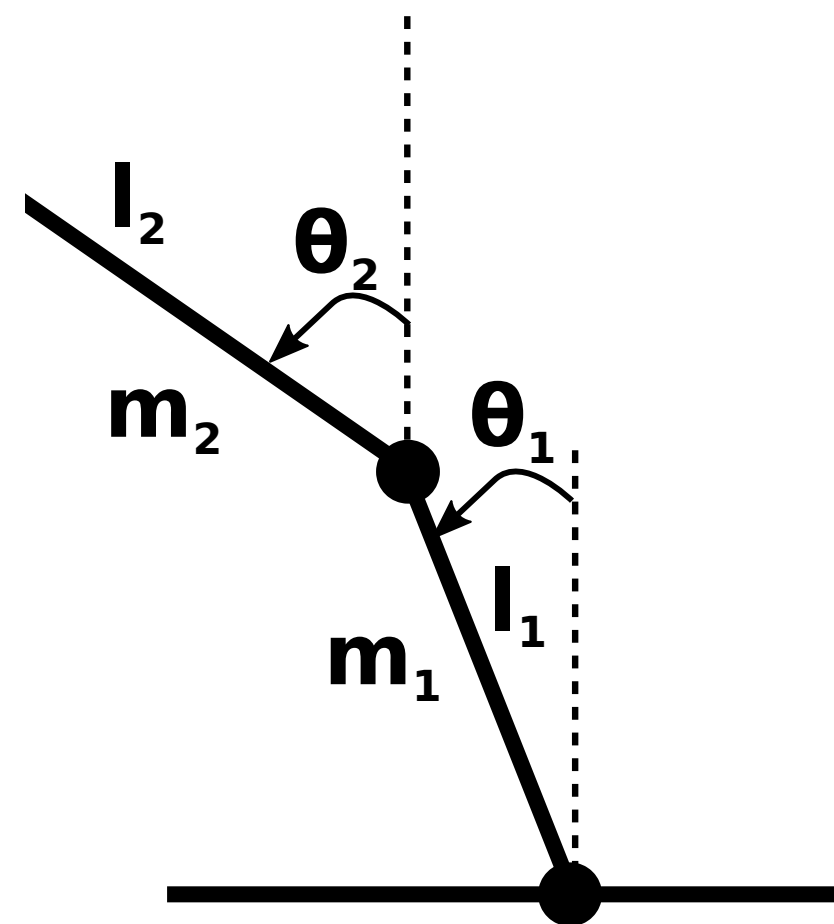
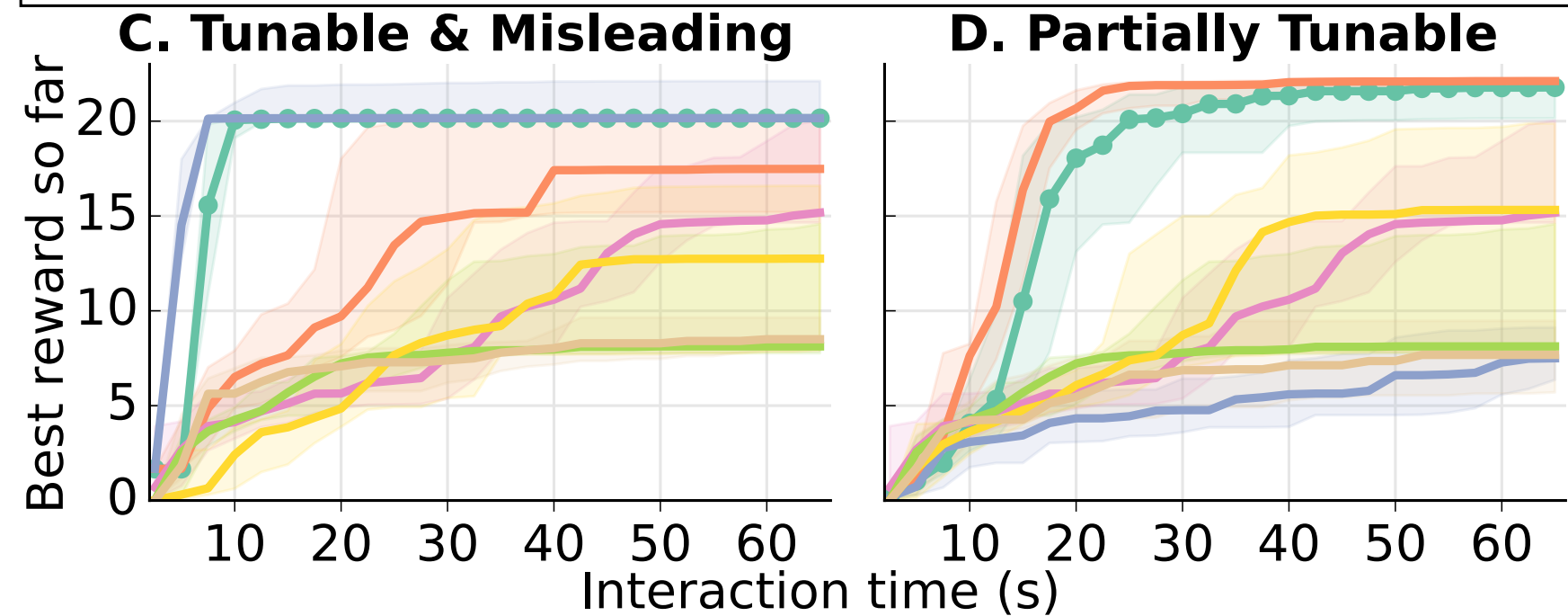
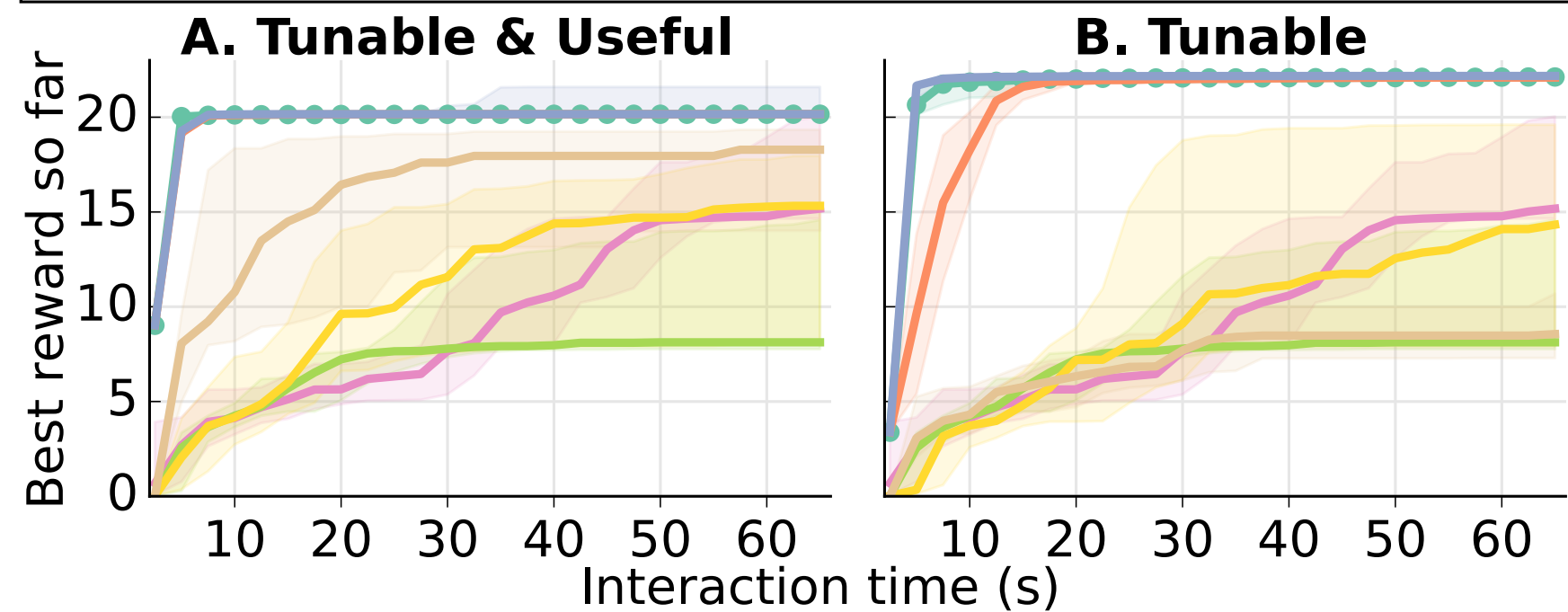
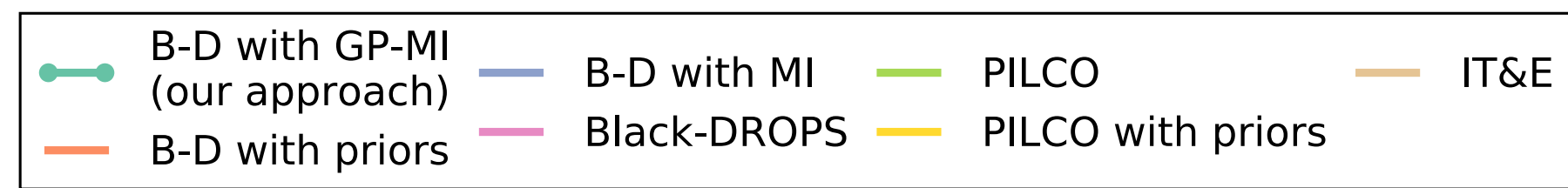
small  
feet



# Model-based policy search

The Black-DROPS algorithm / 160 parameters (neural net.)





Pendubot system