# Towards Model-Based Reinforcement Learning on Real Robots

by
Georg Martius

AUTONOMOUS LEARNING
MAX PLANCK INSTITUTE
FOR INTELLIGENT SYSTEMS

Georg Martius <georg.martius@tue.mpg.de>
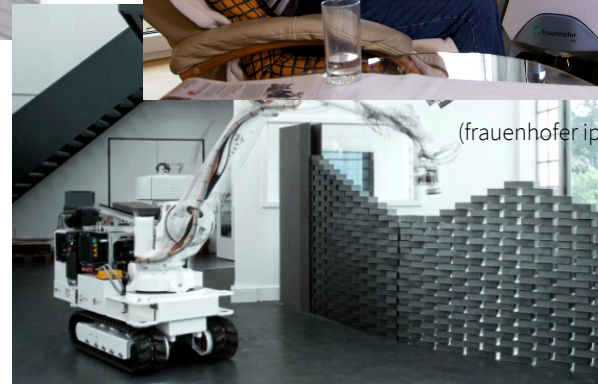
# Vision

Dexterous and versatile robots
as assistance to humans

[Kuka]

(frauenhofer ipa)

[NCCR Digital Fabrication]

[bergkvistanna karin@tuvie.com]

➢ learning

➢ adaptivity

➢ safety

Georg Martius <georg.martius@tue.mpg.de>

# Reinforcement learning achievements

## Robotics, Games: Go, Dota, Starcraft



[OpenAI 2019]



[Deepmind 2019]

## Problems

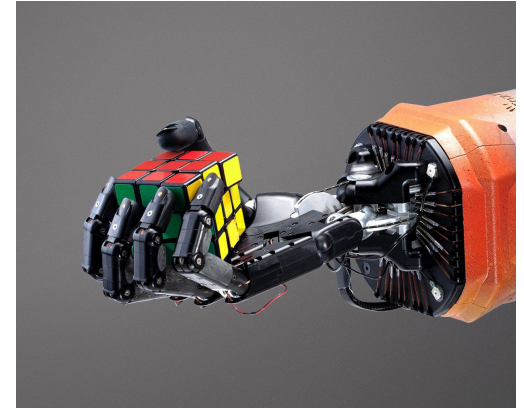need probitive amount of data
simulations
need very long time



[ionos.com]

[DARPA rescue
challenge]

Georg Martius <georg.martius@tue.mpg.de>

# Current Situation

Learned robot control in recent **research**

- using simulations of robot and environment
- trained to cope with anticipated variations

✓ works already for difficult tasks

✗ needs a high-fidelity simulation

✗ learning is inefficient (needs domain randomization)

✗ resulting controller is fixed

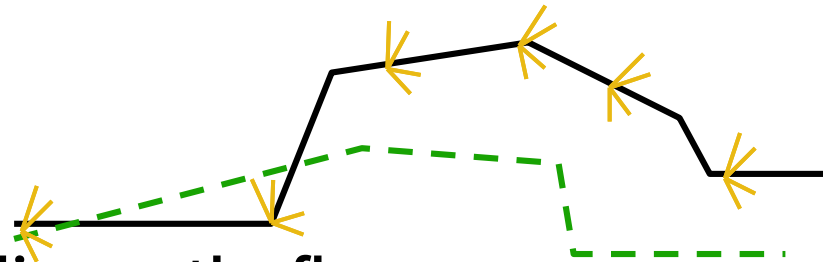✗ new task: start from scratch



[OpenAI 2019]



[Hutter lab. ETH, 2020]

## How to achieve efficient learning and online adaptation?

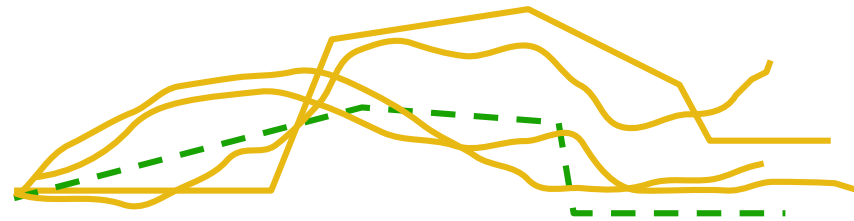# Model-based Reinforcement Learning

## Two instantiations

## Bellman operator to optimize a value function and the policy

➢ use model to collect data nearby real observation
➢ learn to solve a specific task
➢ global optimimization
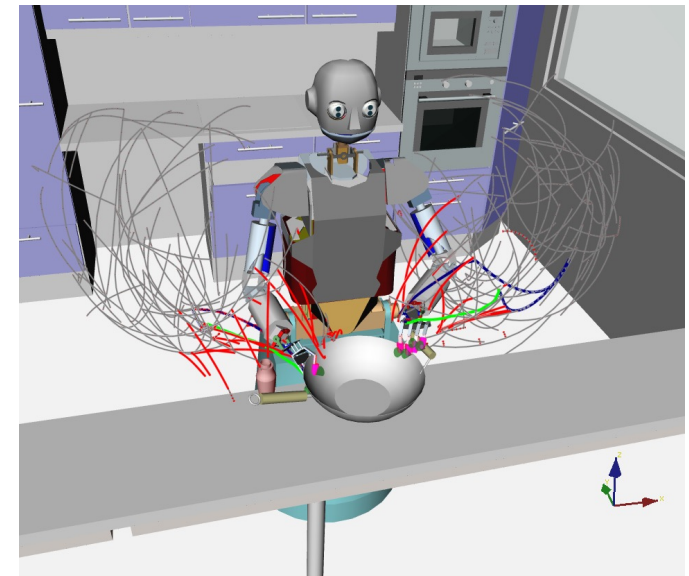
## Planning to search for a policy on the fly

➢ use model for planning
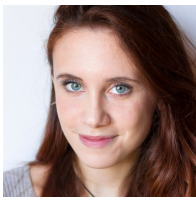➢ perform new task on the fly
➢ optimize finite horizon problem

Georg Martius <georg.martius@tue.mpg.de>

# Model-based RL with Planning

## Challenges:

- ➢ Real time planning

    - ○ previous general purpose planner are

    - ○ 1-2 orders of magnitudes too slow

- ➢ Good models + uncertainty aware

- ➢ Safety



(KIT H²T)

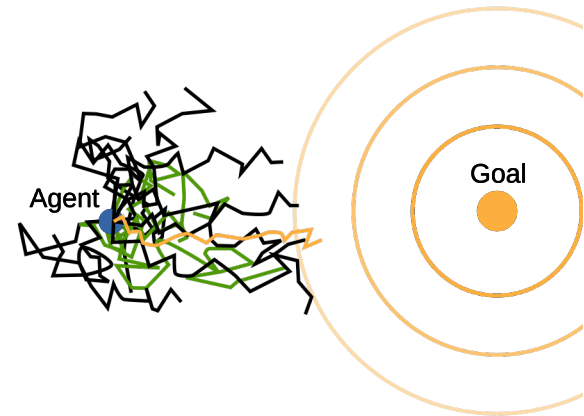**Cristina Pinneri**　**Sebastian Blaes**　**Marin Vlastelica**　**Shambhuraj Sawant**　**Georg Martius**

Georg Martius <georg.martius@tue.mpg.de>

# Model-based Planning

**Cross Entropy Method (CEM)**

➢ Sampling based optimization

$$a_1, \ldots, a_h \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

$$a_1, \ldots, a_h = \underset{a_1, \ldots, a_h}{\arg\min} J(a_1, \ldots, a_h)$$

J cost of rollout

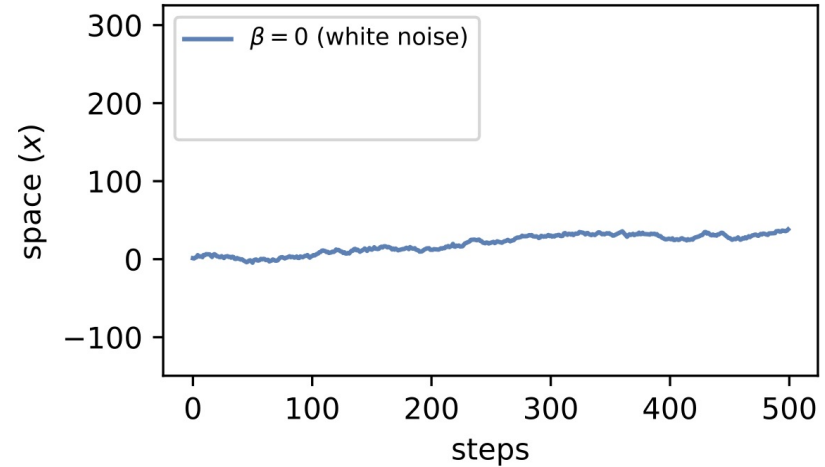Georg Martius <georg.martius@tue.mpg.de>

# Random Walk Model and Colored Noise
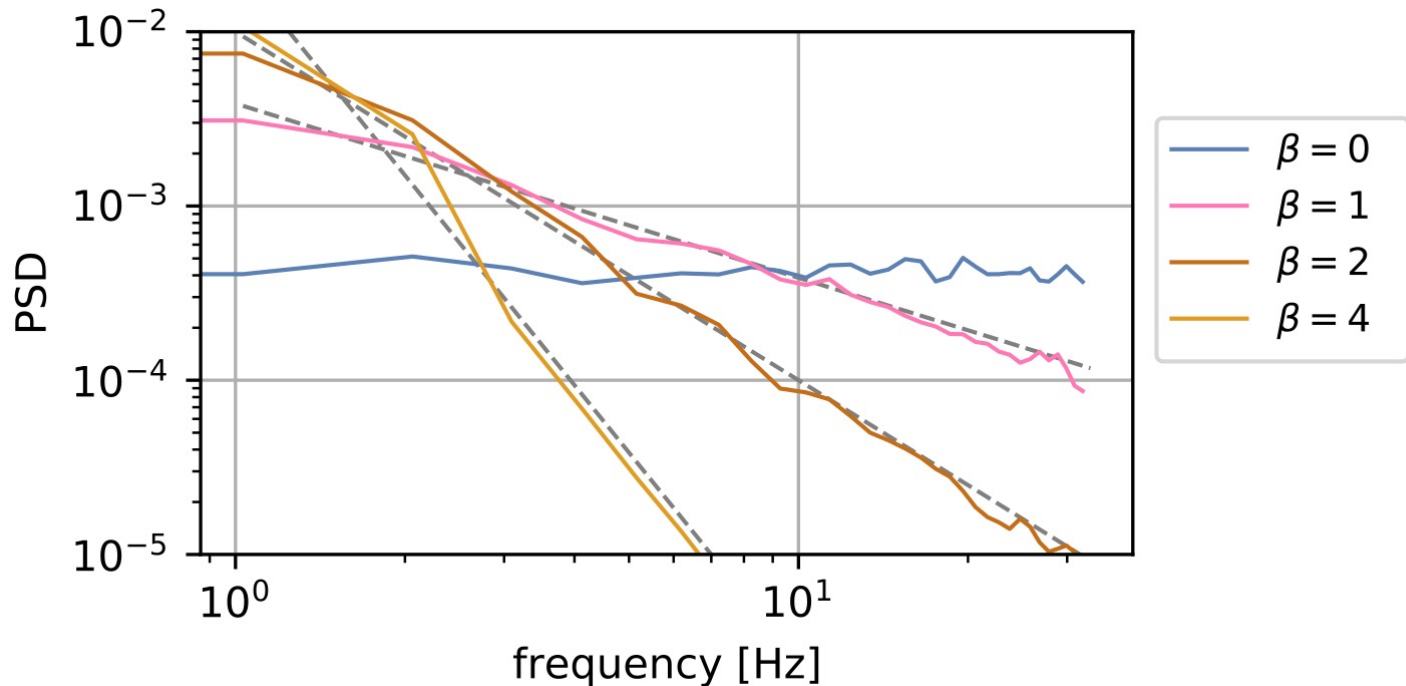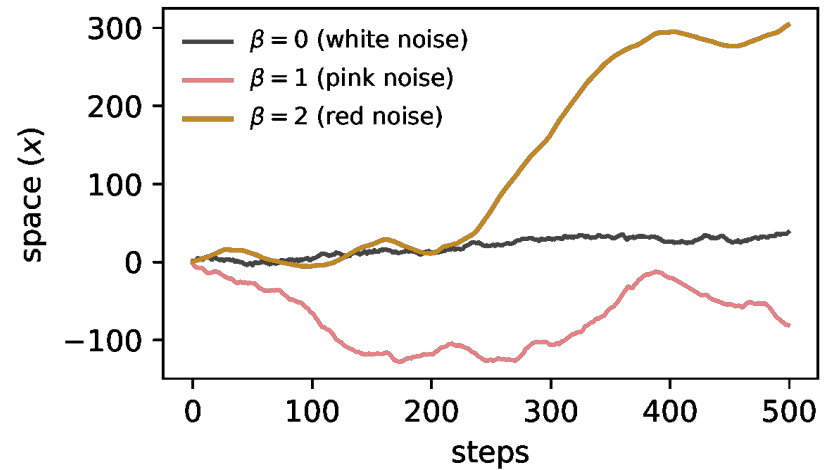
$$x_{t+1} = x_t + a_t$$

Brownian random walk

$$a_t \sim \mathcal{N}(0,1)$$

Lévy walk

$$a_t \sim \mathcal{CN}^{\beta}(0,1)$$



Power spectral density of action sequences

Georg

# Random Walk Model and Colored Noise
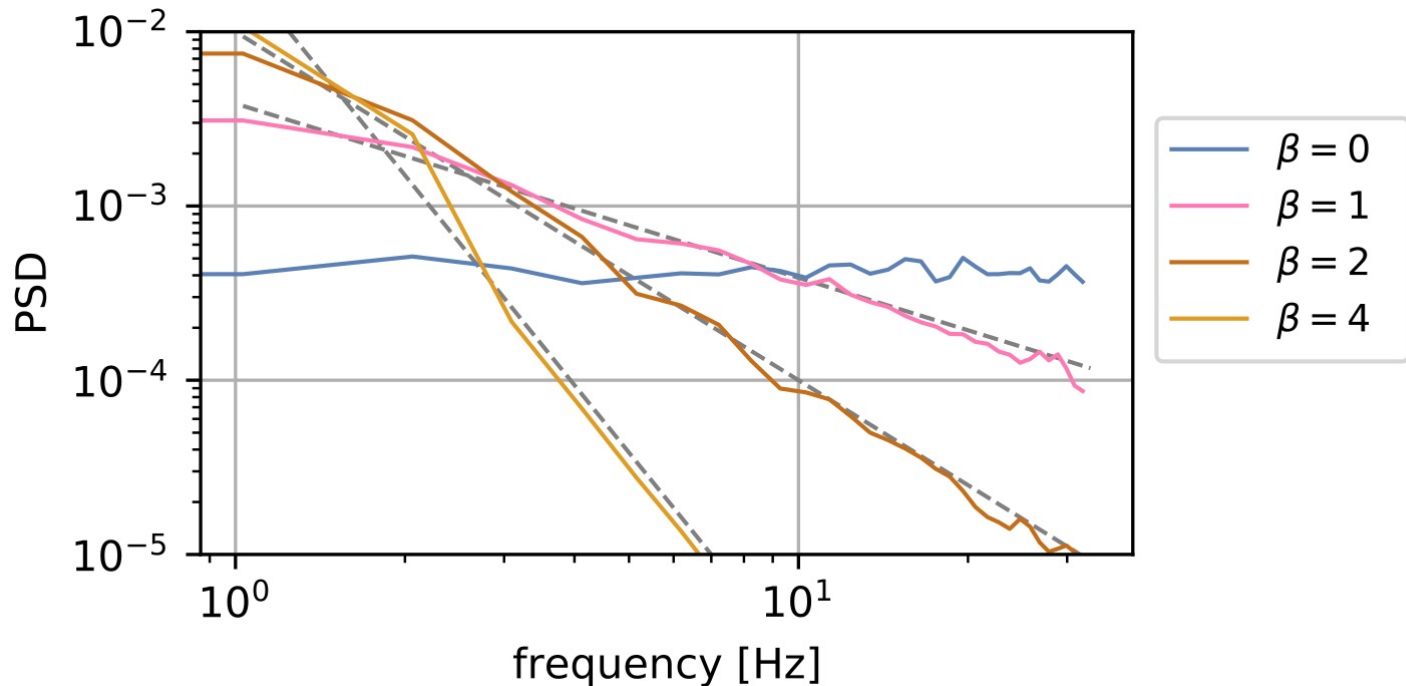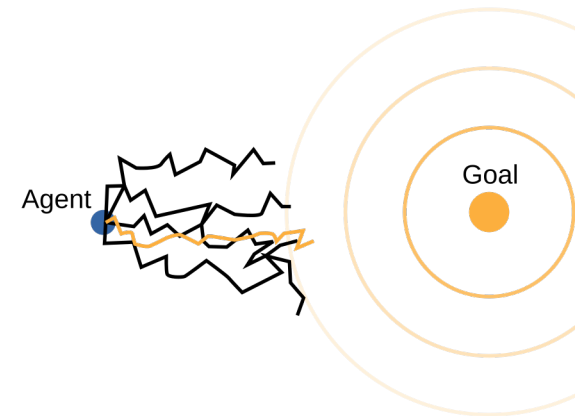
$$x_{t+1} = x_t + a_t$$

Brownian random walk

$$a_t \sim \mathcal{N}(0, 1)$$

Lévy walk

$$a_t \sim \mathcal{CN}^{\beta}(0, 1)$$



Power spectral density of action sequences

Georg

# Model-based Planning

## Cross Entropy Method (CEM)

➤ Sampling based optimization

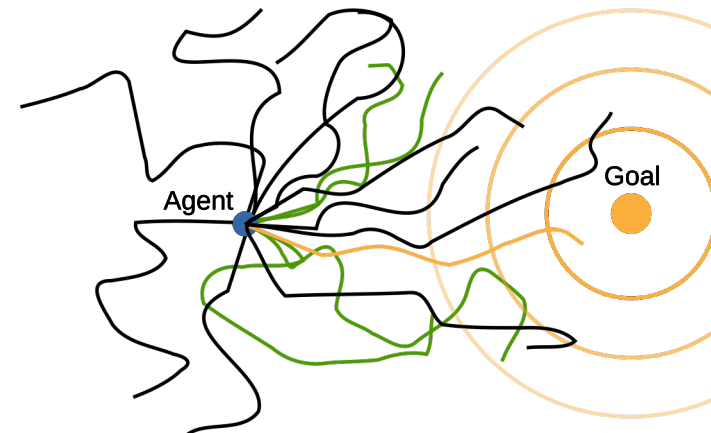$$a_{t,\ldots,t+H} \sim \mathcal{N}(\mu_i, \sigma_i^2)$$



Agent    Goal

## improved Cross Entropy Method
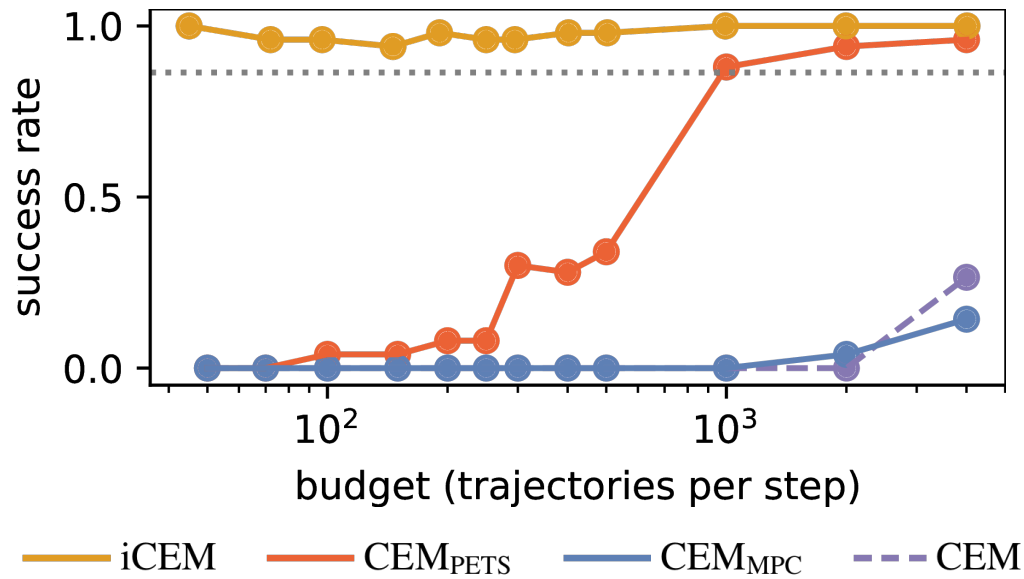
✚ Memory

✚ Colored noise: temporal correlation

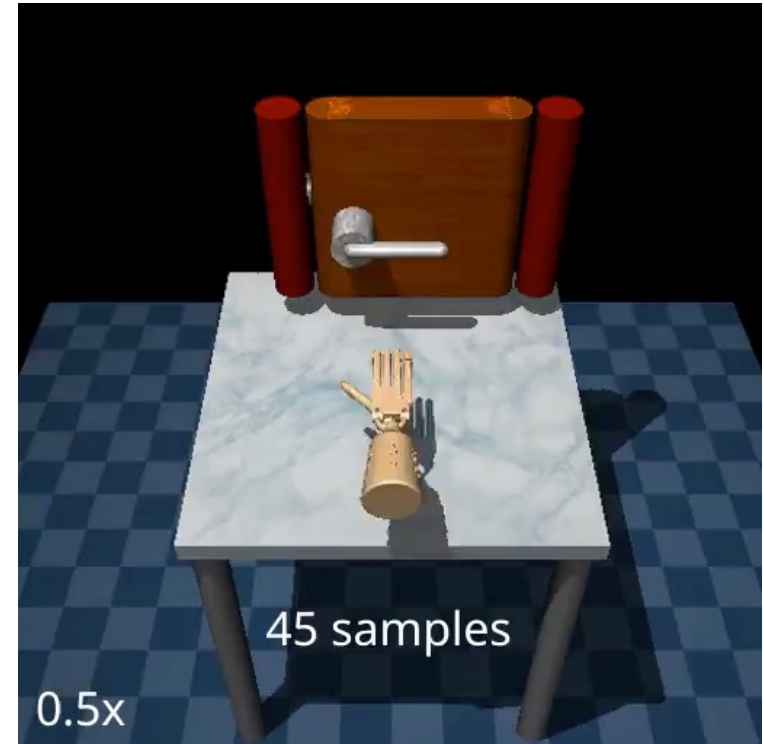   power-law spectrum: $\text{PSD} \propto \dfrac{1}{f^\beta}$

✚ small improvements



Agent    Goal

Georg Martius <georg.martius@tue.mpg.de>

# Model-based Planning

Door (sparse reward)



Ground Truth Models

iCEM · CEM$_{PETS}$ · CEM$_{MPC}$ · CEM



45 samples

0.5x

(environment from DAPG project)

Georg Martius <georg.martius@tue.mpg.de>
Pinneri, Sawant, Blaes, Achterhold, Stückler, GM. *CORL* 2020
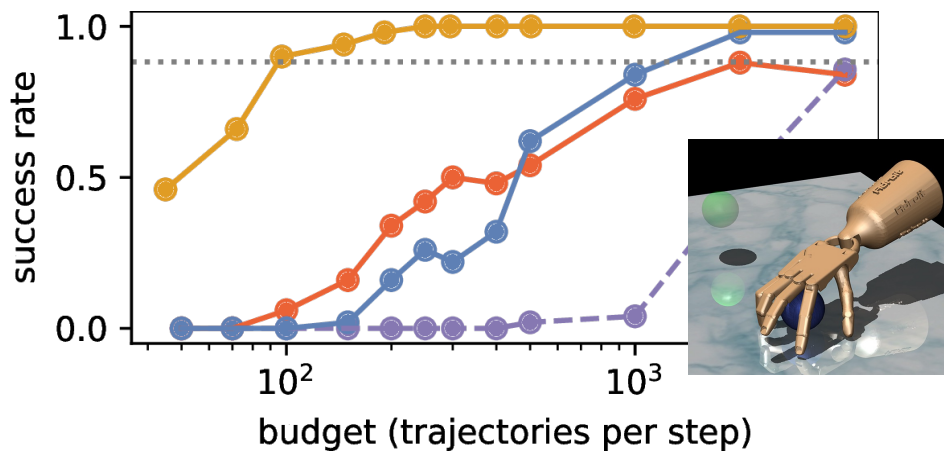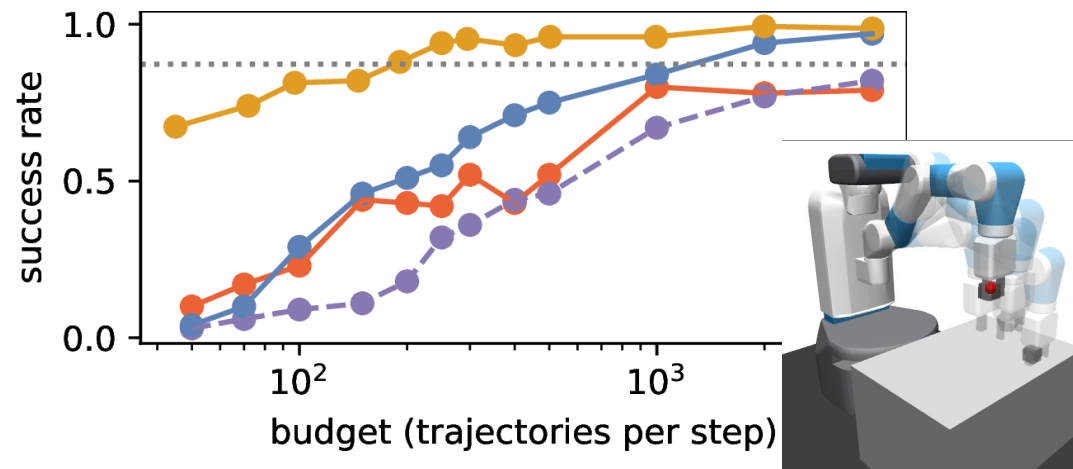
# Model-based Planning



Halfcheetah (running)

Humanoid Standup

Relocate
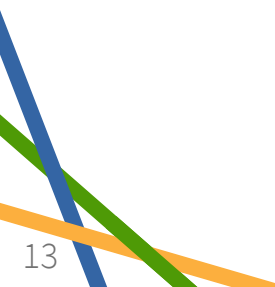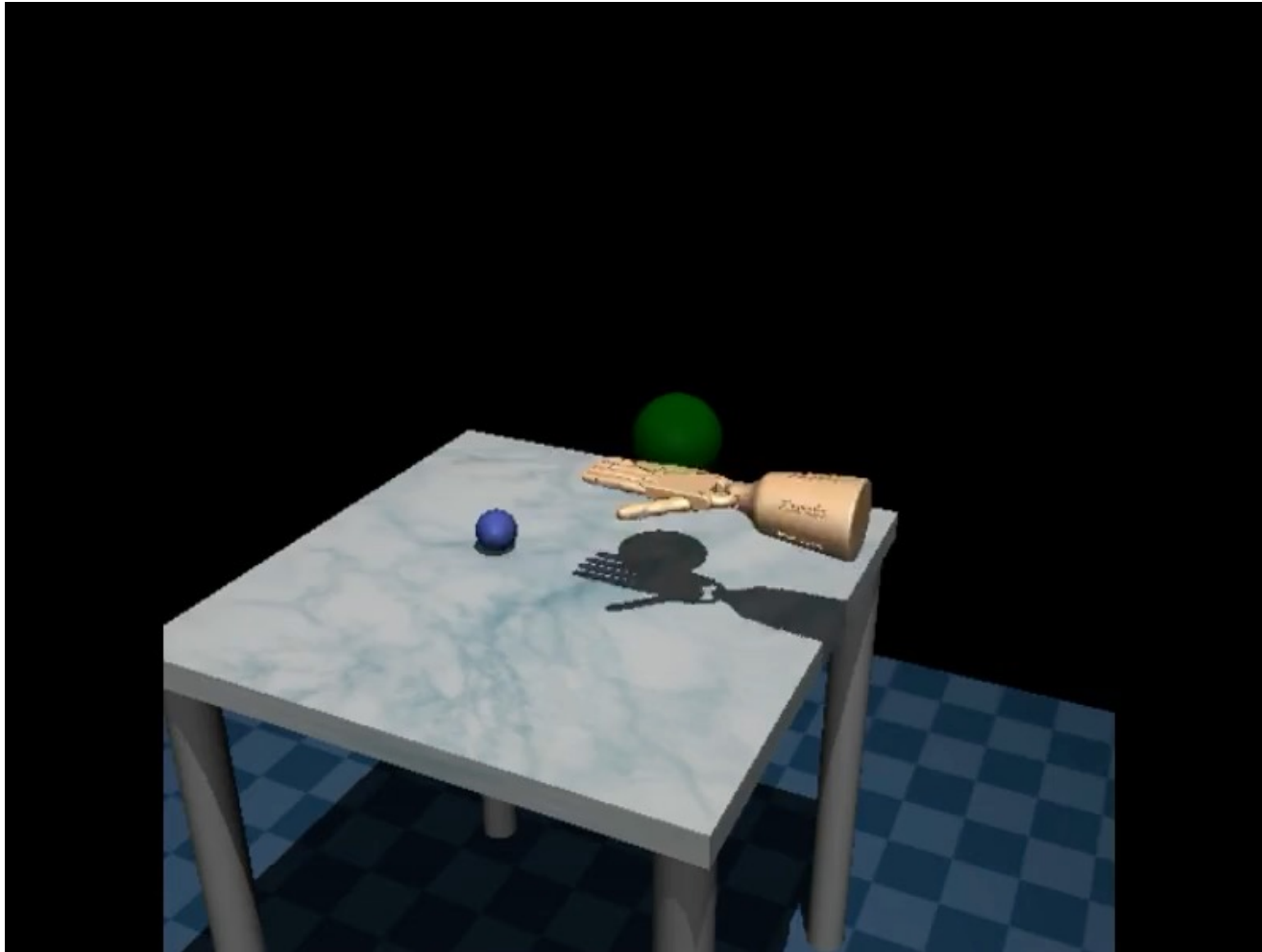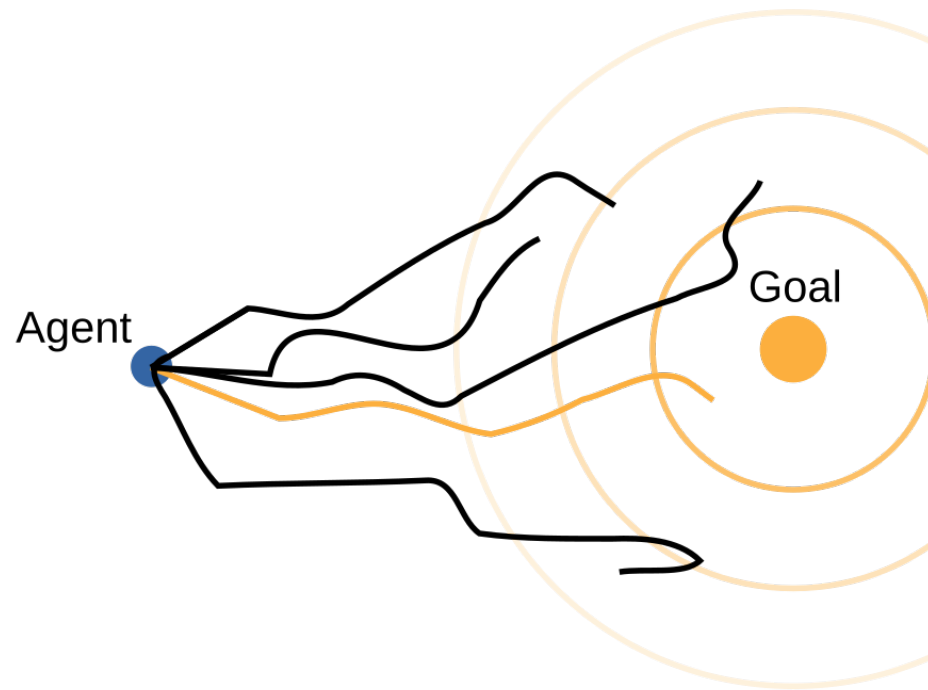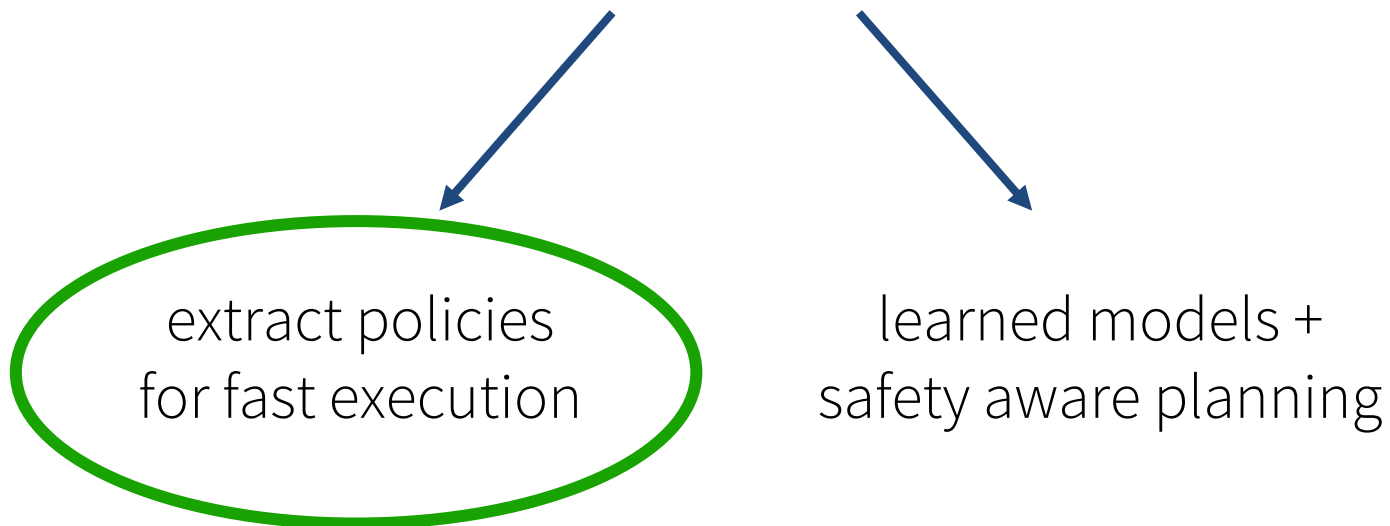
Fetch Pick & Place

iCEM    CEM$_{PETS}$    CEM$_{MPC}$    CEM    90% baseline

Georg Martius <georg.martius@tue.mpg.de>

# Action



Georg Martius <georg.martius@tue.mpg.de>

Fast planner – iCEM

extract policies
for fast execution

learned models +
safety aware planning

Georg Martius <georg.martius@tue.mpg.de>

# Learn Policy from Plans

We can create solutions for complicated control problems within seconds, but:

➢ need a lot of run-time compute

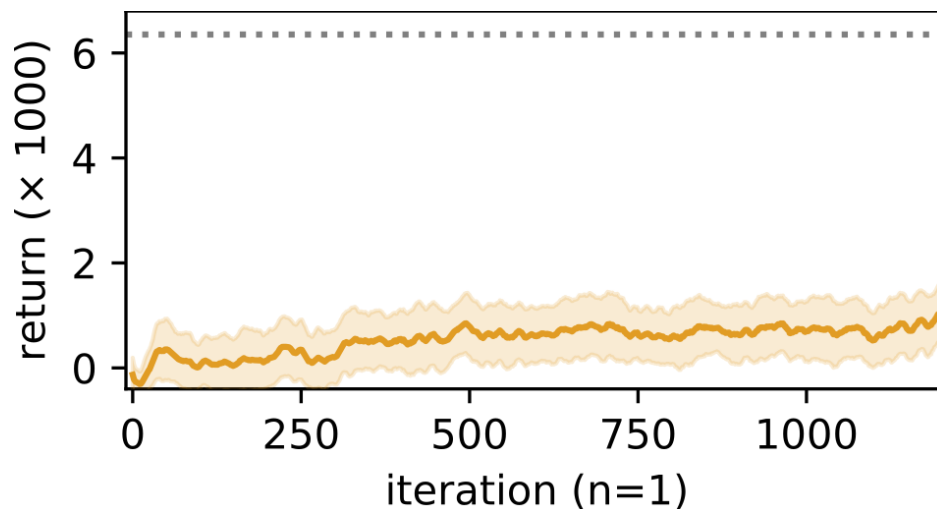➢ mostly with Ground Truth models (simulations)

**Goal of our method:**

➢ train policy from planner data

➢ make policy and planner mutually improve themselves
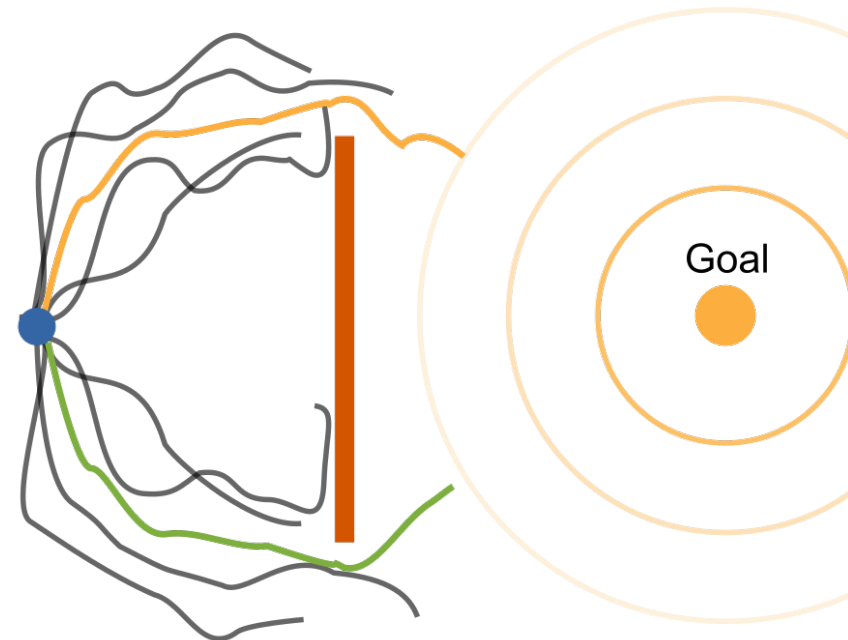
➢ solve tasks that standard RL struggles with

Georg Martius <georg.martius@tue.mpg.de>

# Challenges of learning policy from planner's data

**Let's do simple behavioral cloning:**

➢ does not work!

➢ multimodality + combounding errors
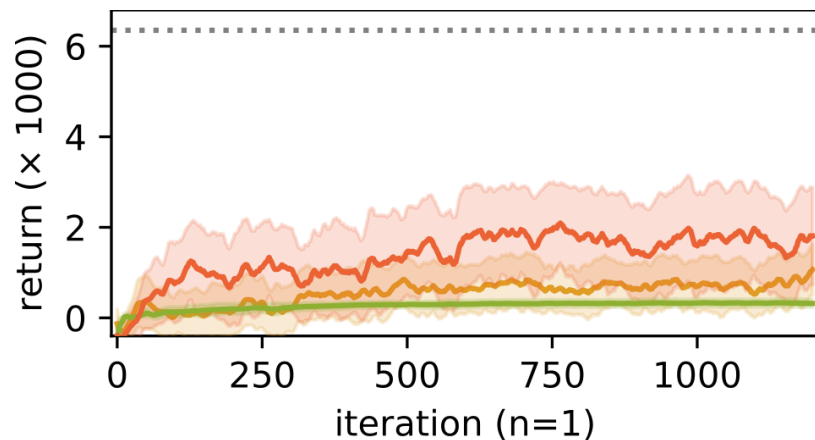


Georg Martius <georg.martius@tue.mpg.de>

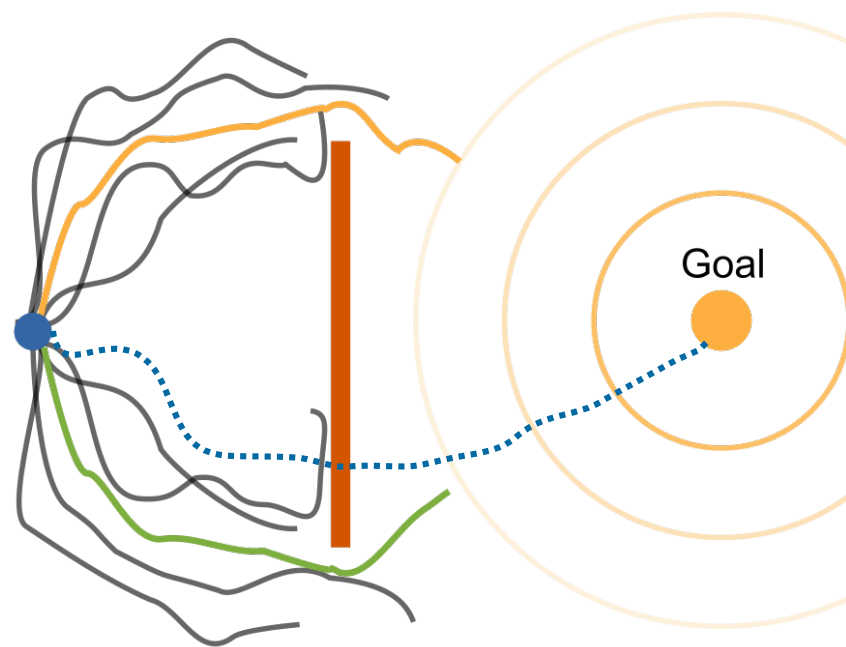# Challenges of learning policy from planner's data

**Okay, use guidance (*guided policy search*)**

$$a_1, \ldots, a_h = \arg\min_{a_1, \ldots, a_h} f(a_1, \ldots, a_h) + \lambda \sum_i \|a_i - \pi(s_i)\|$$

better, but not ideal!

Georg Martius <georg.martius@tue.mpg.de>

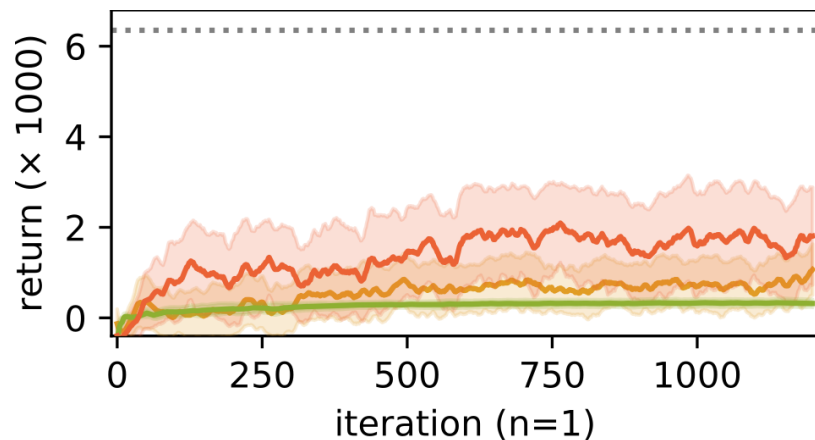Pinneri*, Sawant*, Blaes, GM. ICLR 2021

# Challenges of learning policy from planner's data

**Okay, use guidance (*guided policy search*)**

$$a_1, \ldots, a_h = \arg\min_{a_1, \ldots, a_h} f(a_1, \ldots, a_h) + \lambda \sum_i \|a_i - \pi(s_i)\|$$
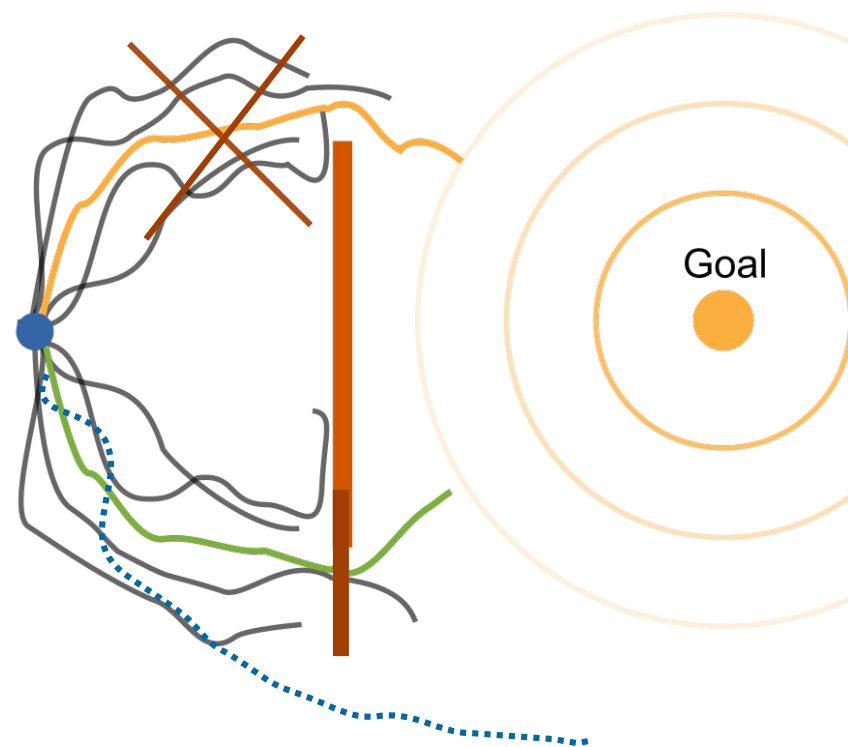
better, but not ideal!

➤ planner premature convergence

➤ combounding errors

Agent

Goal



return (× 1000)

iteration (n=1)

EX — iCEM$_\pi$-GPS — iCEM DAgger — iCEM BC ····· SAC

Georg Martius <georg.martius@tue.mpg.de>

# Challenges of learning policy from planner's data

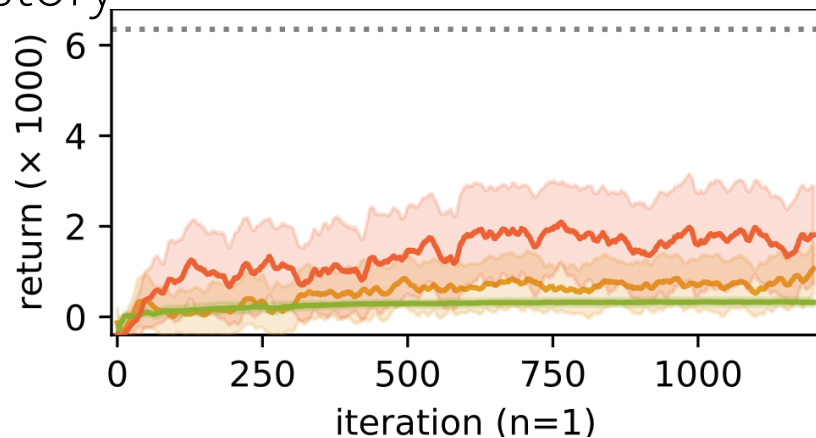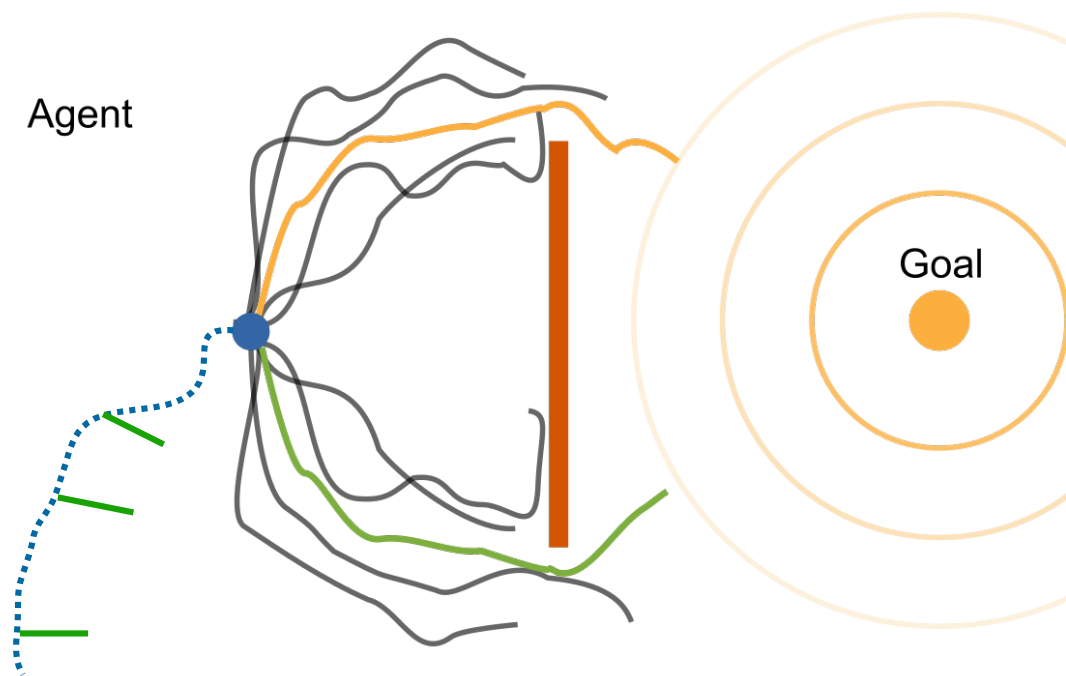**Okay, use guidance (*guided policy search*)**

$$a_1, \ldots, a_h = \underset{a_1, \ldots, a_h}{\arg \min} f(a_1, \ldots, a_h) + \lambda \sum_i \| a_i - \pi(s_i) \|$$

**+ relabeling (*DAgger*):**

➤ ask planner for better action

  ➢ Problem:
         planner is not good without
history

Agent

Goal



⎓X ── iCEM$_\pi$-GPS ── iCEM DAgger ── iCEM BC ····· SAC

# Challenges of learning policy from planner's data

**Okay, use guidance (*guided policy search*)**

$$a_1, \ldots, a_h = \arg\min_{a_1, \ldots, a_h} J(a_1, \ldots, a_h) + \lambda C^{\mathrm{aux}}(a_1, \ldots, a_h)$$

**+ relabeling (*Dagger*)**
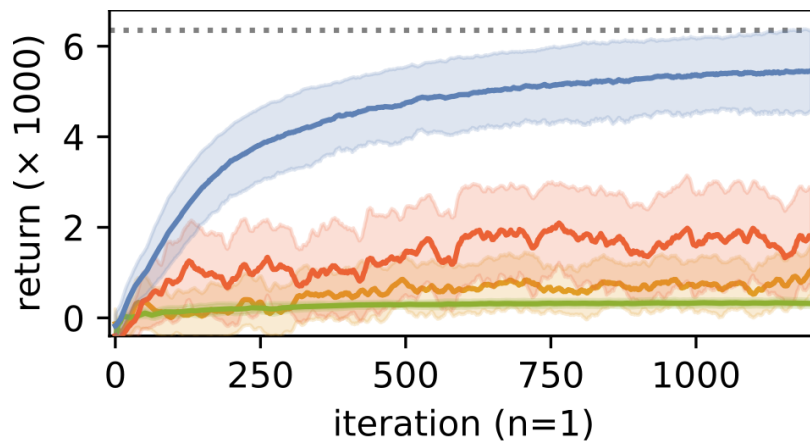
**+ adaptive λ**

don't use guidance
if no progress on actual task

$$\lambda = c \frac{\mathcal{R}(J)}{\mathcal{R}(C^{\mathrm{aux}}) + \epsilon}$$
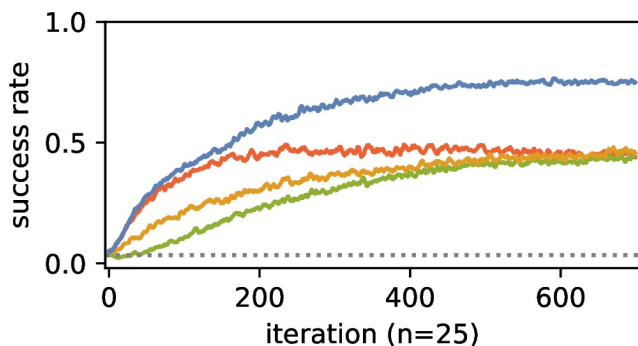
$c$ relative importance

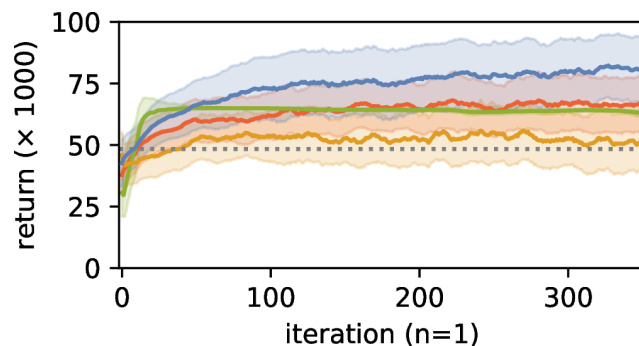$$\mathcal{R}(C) = \max_{\text{elite-set}} C - \min_{\text{elite-set}} C$$
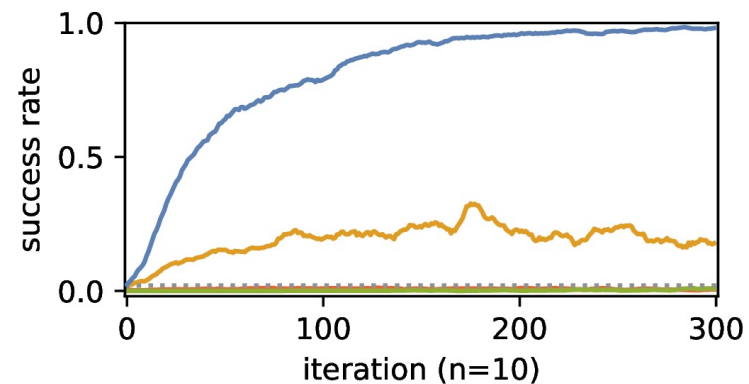
**APEX (Adaptive Policy Extraction)**



return (× 1000) vs iteration (n=1)

Legend: —— APEX    —— iCEM$_\pi$-GPS    —— iCEM DAgger    —— iCEM BC    ····· SAC

# APEX results



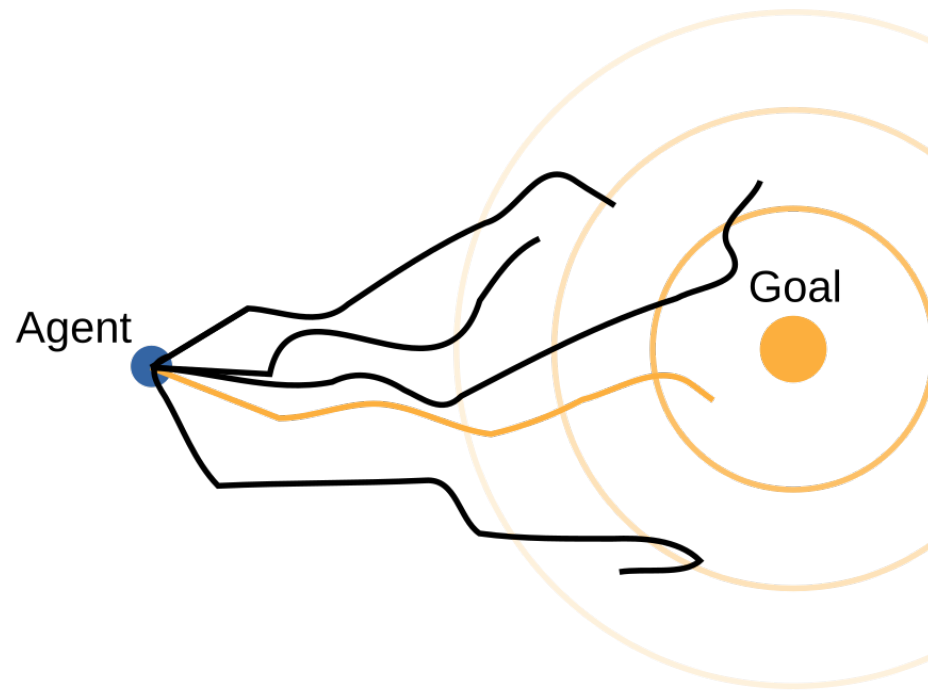Fetch Pick & Place     Humanoid Standup (half duration)     Door (sparse reward)
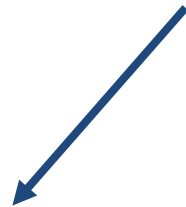
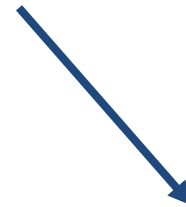APEX    iCEM$_\pi$-GPS    iCEM DAgger    iCEM BC    SAC

## Discussion:

+ strong policies for hard tasks

× behavior not perfect

× high computational costs

Georg Martius <georg.martius@tue.mpg.de>     Pinneri*, Sawant*, Blaes, GM. ICLR 2021

Fast planner – iCEM

extract policies
for fast execution

learned models +
safety aware planning

# Real-time Risk-Averse Model-based Planning

## Towards real-robots:

➢ **learned models:** adapt to real  system

➢ **safety:** do not destroy the robot / environment
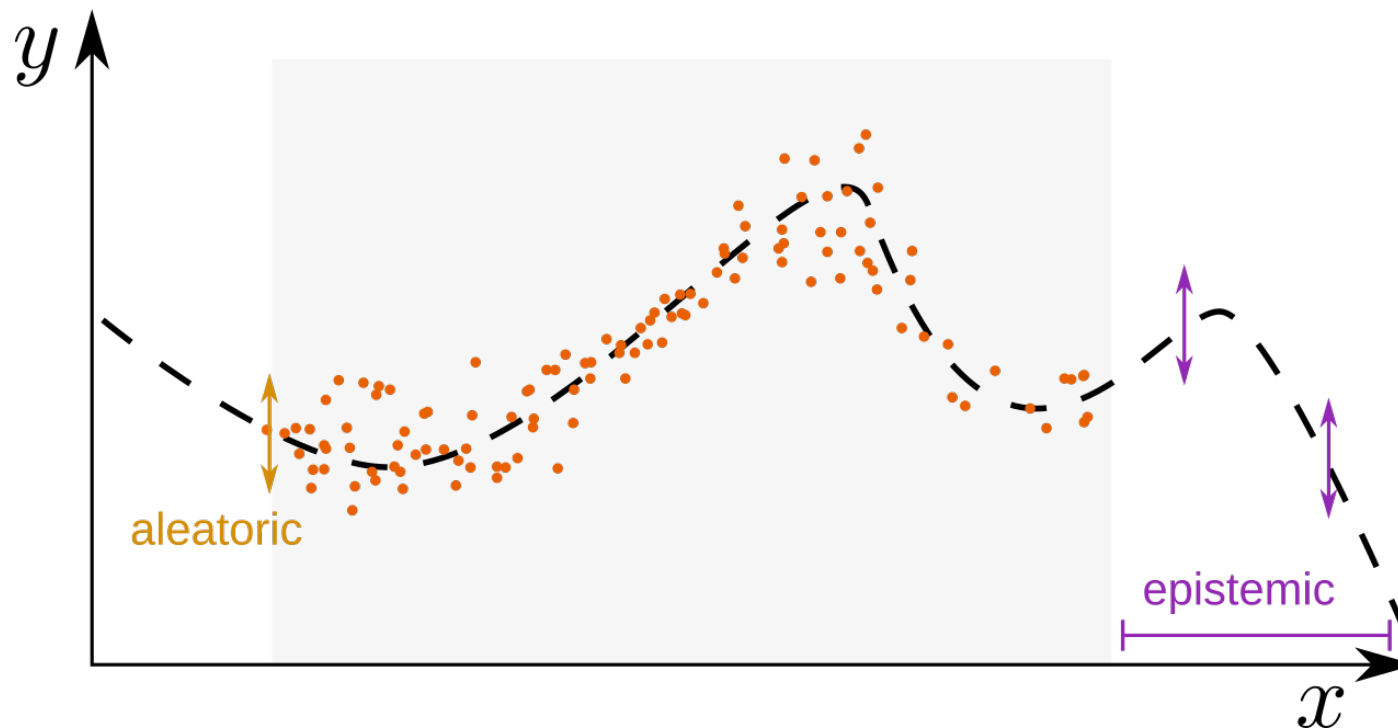
➢ **real-time:** run at > 25Hz

## Challenges:

➢ explore effective but careful

➢ need powerful dynamics models

➢ awareness of **uncertainties** required

Georg Martius <georg.martius@tue.mpg.de>

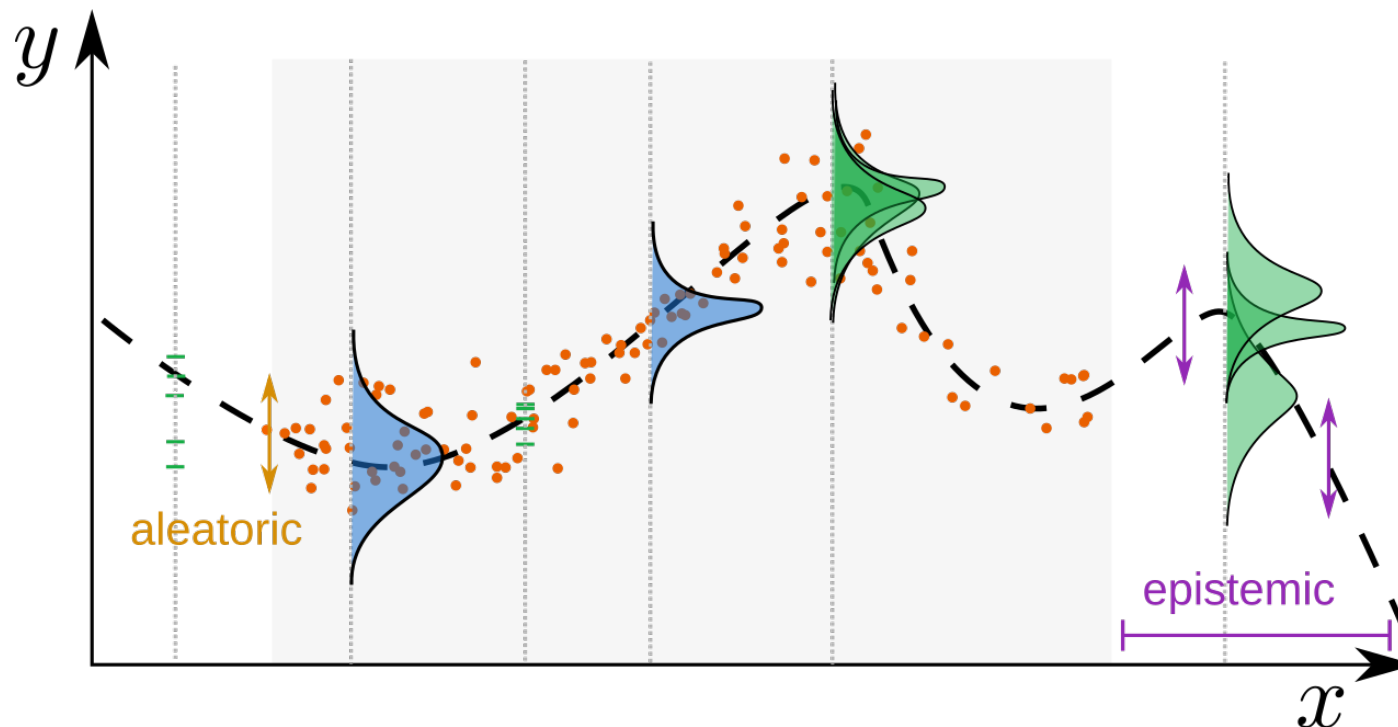Vlastelica*, Blaes*, Pinneri, GM. *CORL* 2021

# Dynamics Models with n-step Uncertainty

➤ separation of *aleatoric* and *epistemic* uncertainty

**Why?**

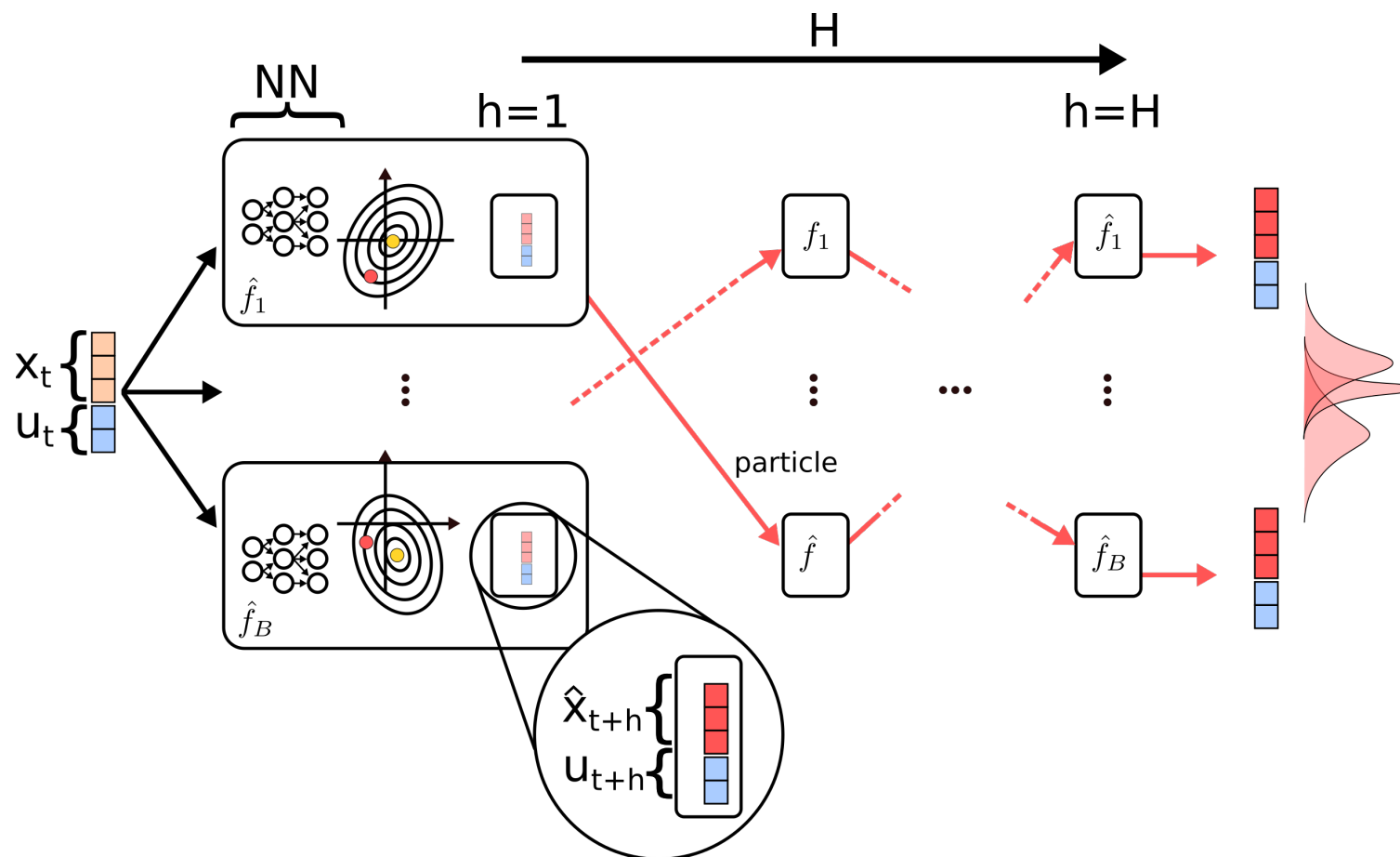➤ aleatoric: avoid

➤ epistemic: seek to reduce



Georg Martius <georg.martius@tue.mpg.de>

# Dynamics Models with n-step Uncertainty

➢ separation of *aleatoric* and *epistemic* uncertainty

**Why?**

➢ aleatoric: avoid

➢ epistemic: seek to reduce



## Ensemble of probabilistic Deep Nets

➢ good estimates of separation both types of uncertainty

Georg Martius <georg.martius@tue.mpg.de>

Vlastelica*, Blaes*, Pinneri, GM. *CORL* 2021

# Dynamics Models with n-step Uncertainty

**What about n-step preditions?**

➢ PETS: [Chua et al 2018] Probabilistic Ensemble models with Trajectory Sampling

➢ sampling at every timestep **mixes** uncertainties

Georg Martius <georg.martius@tue.mpg.de>
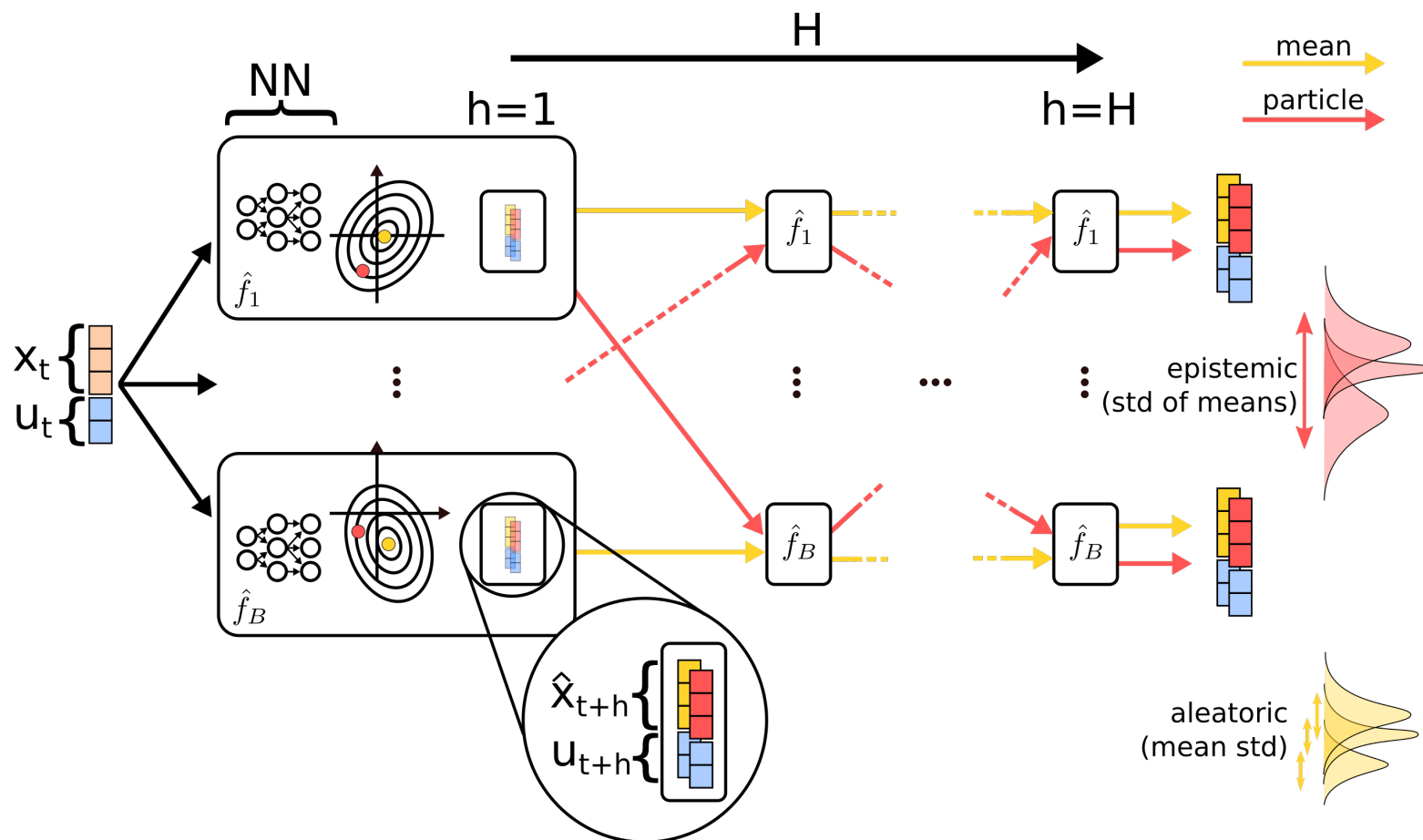
# Dynamics Models with n-step Uncertainty

**What about n-step preditions?**

➤ PETS: [Chua et al 2018] Probabilistic Ensemble models with Trajectory Sampling
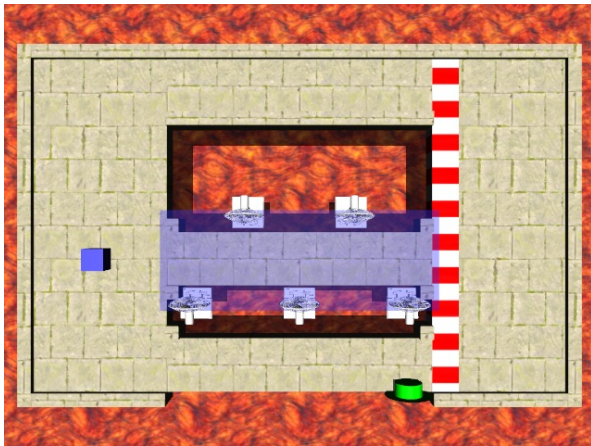
➤ PETSUS [ours] disentangle uncertainties



Georg Martius <georg.martius@tue.mpg.de>

# Dynamics Models with n-step Uncertainty

**What about n-step preditions?**

➤ PETS: [Chua et al 2018] Probabilistic Ensemble models with Trajectory Sampling

➤ PETSUS [ours] disentangles uncertainties also for n-step predictions



Georg Martius <georg.martius@tue.mpg.de>

# Efficient Exploration

**Seek to reduce epistemic uncertaint:**

$$\mathbf{a} = \arg\min_{\mathbf{a}} J(\mathbf{a}) - w^{\mathcal{E}}\mathrm{Epistemic}(\mathbf{a})$$
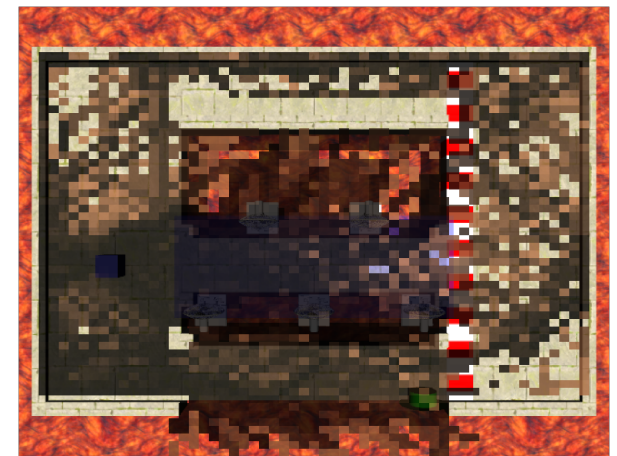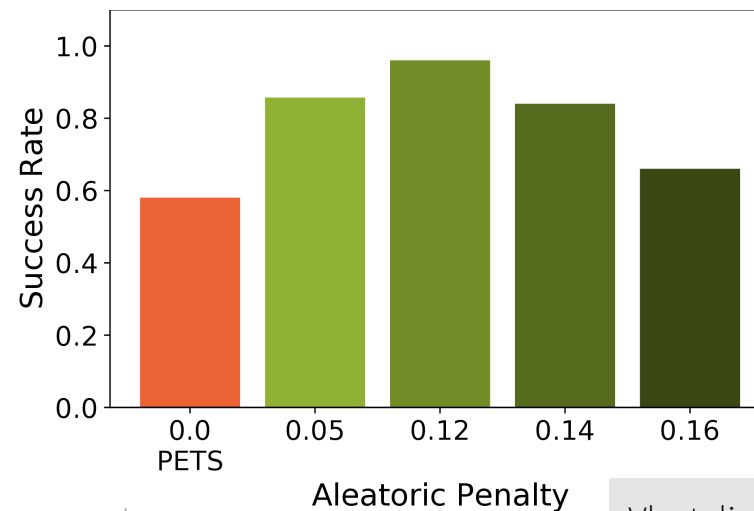
➢ Exploration bonus

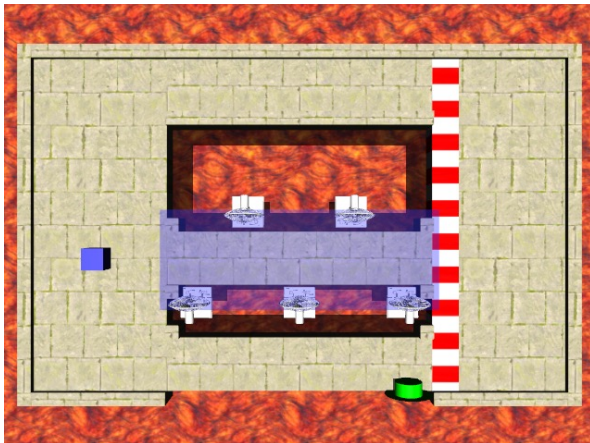toy world: Bridge Maze          no bonus $w^{\mathcal{E}} = 0$          bonus $w^{\mathcal{E}} = 0.05$

Georg Martius <georg.martius@tue.mpg.de>
Vlastelica*, Blaes*, Pinneri, GM. *CORL* 2021

# Risk-averse Behavior

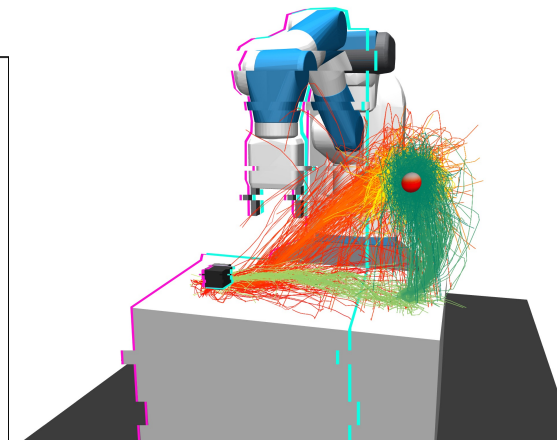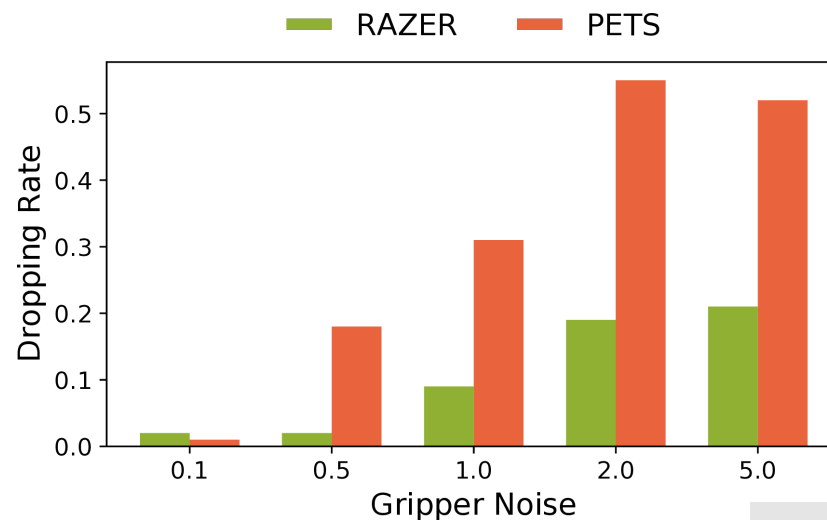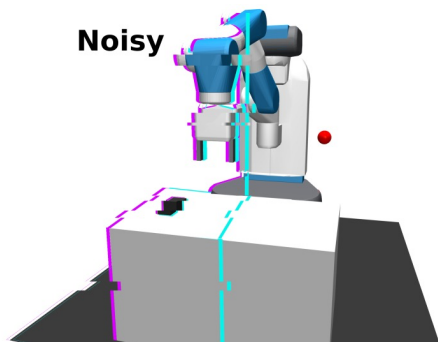Seek to reduce epistemic uncertaint:

$$\mathbf{a} = \arg\min_{\mathbf{a}} J(\mathbf{a}) - w^{\mathcal{E}}\mathrm{Epistemic}(\mathbf{a})$$

## Risk-averse planning

$$\mathbf{a} = \arg\min_{\mathbf{a}} J(\mathbf{a}) + w^{\mathcal{A}}\mathrm{Aleatoric}(\mathbf{a})$$

➢ avoid unpredictable areas: uncertainty penalty
    toy world: Bridge Maze

Georg Martius <georg.martius@tue.mpg.de>

# Risk-averse Behavior

**Seek to reduce epistemic uncertaint:**

$$\mathbf{a} = \arg\min_{\mathbf{a}} J(\mathbf{a}) + w^{\mathcal{E}}\mathrm{Epistemic}(\mathbf{a})$$

**Risk-averse planning**

$$\mathbf{a} = \arg\min_{\mathbf{a}} J(\mathbf{a}) + w^{\mathcal{A}}\mathrm{Aleatoric}(\mathbf{a})$$

➢ avoid unpredictable areas: uncertainty penalty
   Noisy Fetch Pick & Place

Georg Martius <georg.martius@tue.mpg.de>

# Safety-aware Behavior

## Probabilistic safety constraints as cost penalty

➤ compute probability of entering unsafe set $\mathbb{C}$

➤ add high penalty cost if larger $\delta$

➤ here: simple box constraints: analytic solution for probability

$$\mathbf{a} = \arg\min_{\mathbf{a}} J(\mathbf{a}) + w^{\mathcal{S}} \sum_{\Delta t=1}^{H} \llbracket p(\hat{x}_{t+\Delta t} \in \mathbb{C}) > \delta \rrbracket$$
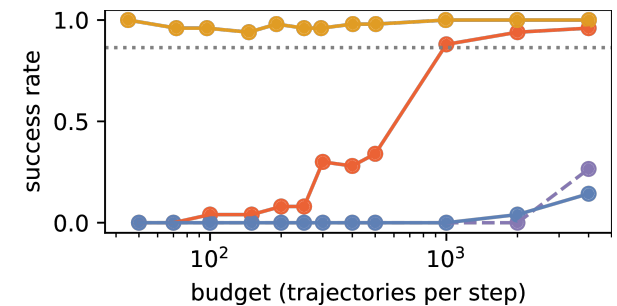
Georg Martius <georg.martius@tue.mpg.de>

# Safety-aware Behavior

## Probabilistic safety constraints as cost penalty

➢ compute probability of entering unsafe set $\mathbb{C}$

➢ add high penalty cost if larger $\delta$

➢ here: simple box constraints: analytic solution for probability

$$\mathbf{a} = \arg\min_{\mathbf{a}} J(\mathbf{a}) + w^{\mathcal{S}} \sum_{\Delta t=1}^{H} [\![ p(\hat{x}_{t+\Delta t} \in \mathbb{C}) > \delta ]\!]$$

### Risk-Averse Zero-Order Trajectory Optimization
Solo8-LeanOverObject

**RAZER**

**PETS**

(23 / 30 runs w/o violations)          (0 / 30 runs w/o violations)

# Safety-aware Behavior

## Probabilistic safety constraints as cost penalty

- ➤ compute probability of entering unsafe set $\mathbb{C}$

- ➤ add high penalty cost if larger $\delta$

- ➤ here: simple box constraints: analytic solution for probability

$$\mathbf{a} = \arg\min_{\mathbf{a}} J(\mathbf{a}) + w^{\mathcal{S}} \sum_{\Delta t=1}^{H} \llbracket p(\hat{x}_{t+\Delta t} \in \mathbb{C}) > \delta \rrbracket$$



(Solo8 robot MPI-IS)

Georg Martius <georg.martius@tue.mpg.de>

# Summary

## Fast universal planning – iCEM

- ➢ works with arbtrary models and cost functions
- ➢ 1-2 orders of magnitude faster than previous SotA



## Extracting policies is hard – APEX is one solution

- ➢ adaptive guided policy search + Dagger
- ➢ offline RL might be a better choice



## Planning with learned models – RAZER

- ➢ uncertainty estimation
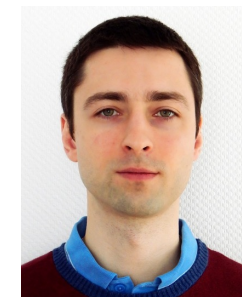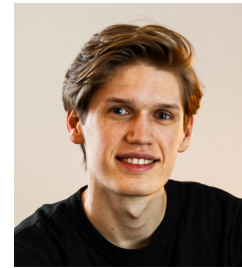- ➢ efficient exploration, risk and safety-aware behavior
- ➢ GPU implementation at 30 Hz

# Near future applications to real robots:



real-robot-challenge.com @ MPI-IS

Georg Martius <georg.martius@tue.mpg.de>

Georg Martius <georg.martius@tue.mpg.de>

**Huanbo Sun**　　**Georg Martius**
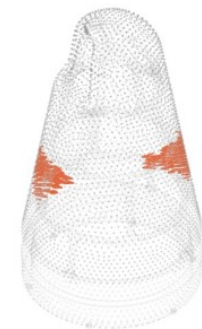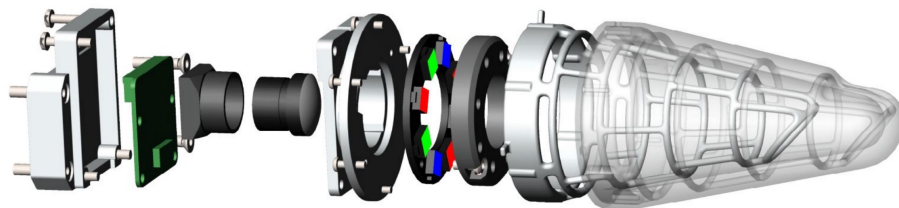
# Coming out soon: ML-driven haptic sensor

**Huanbo Sun** (Ph.D. Student)
Autonomous Learning Group

**Katherine J. Kuchenbecker**
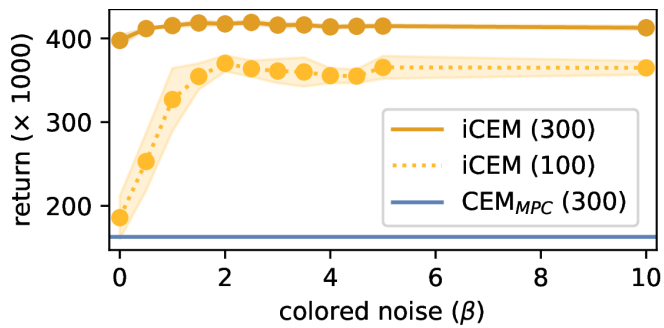Haptic Intelligence
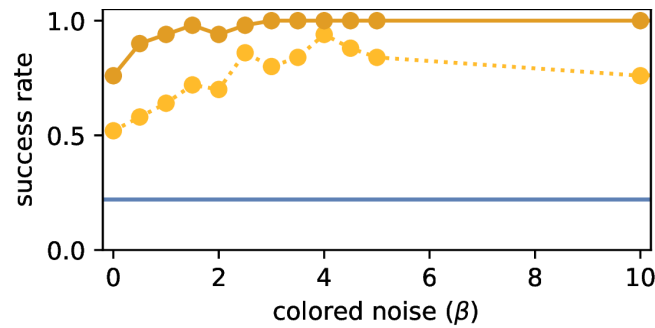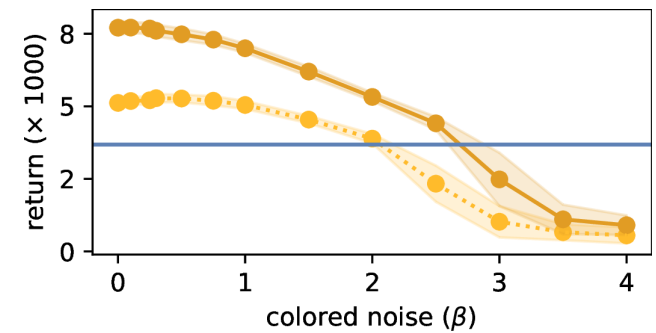Department

**Georg Martius**
Autonomous Learning Group

Georg Martius <georg.martius@tue.mpg.de>

# iCEM Sensitivity



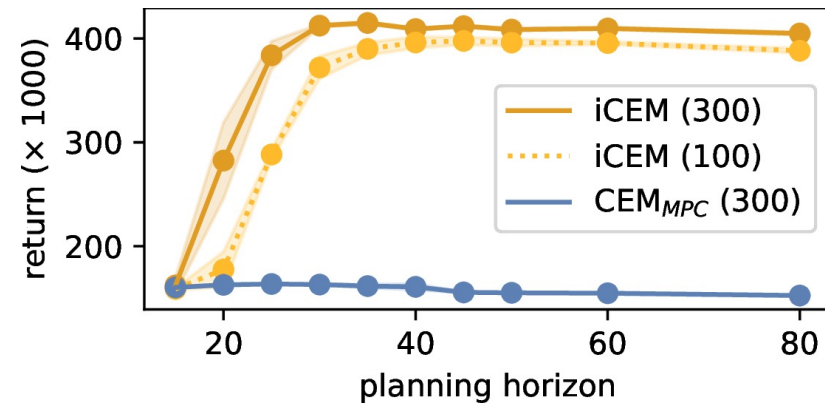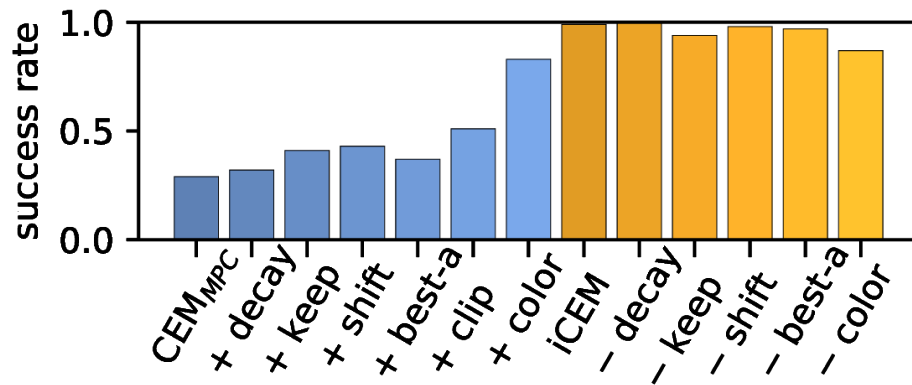Georg Martius <georg.martius@tue.mpg.de>

# iCEM Ablation



Relocate 300

Fetch Pick & Place 100