# On Pathologies in KL-Regularized Reinforcement Learning from Expert Demonstrations
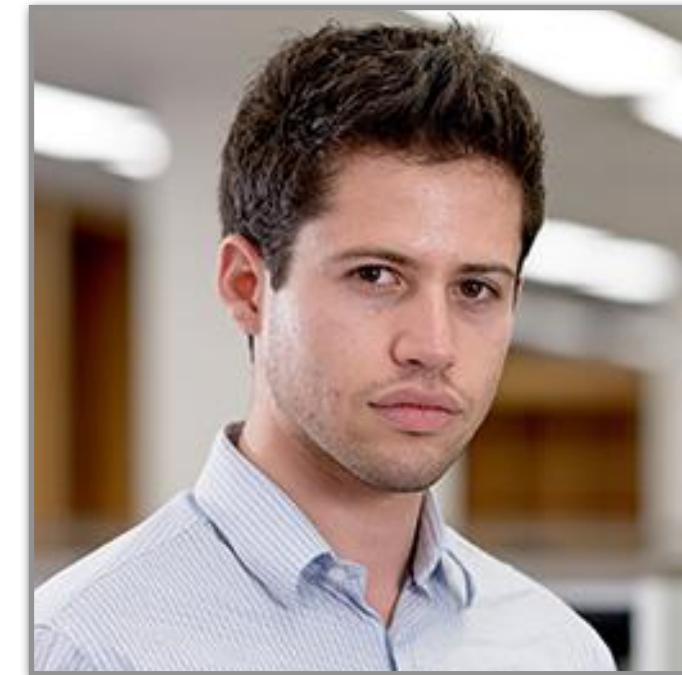
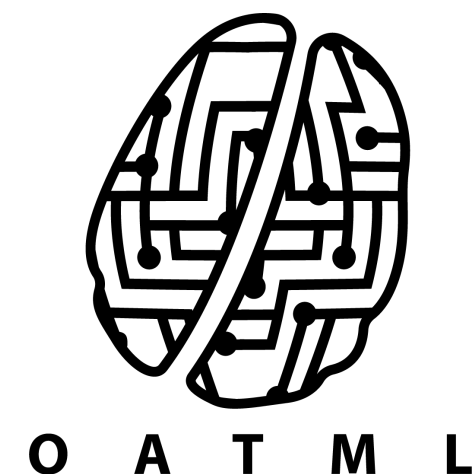**Tim G. J. Rudner***     **Cong Lu***     Michael A. Osborne     Yarin Gal     Yee Whye Teh

## Neural Information Processing Systems 2021

O A T M L

**Correspondence to**

tim.rudner@cs.ox.ac.uk

UNIVERSITY OF OXFORD

How can we use **expert demonstrations** to effectively **accelerate online training** in RL?

# KL-regularization balances fitting online data and matching a behavioral expert policy.

**KL-regularization** balances fitting online data

and matching a behavioral expert policy.

**KL-regularization** balances fitting online data and matching a **behavioral expert policy.**

# **Problem**:

**Problem**:

KL-regularized RL can suffer from pathological behavior during training.

**Problem**:

KL-regularized RL can suffer from pathological behavior during training.

**Problem**:

KM-regularized RL can suffer from
pathological behavior during training.



Expert Demonstration



Learned Behavior

**Problem**:

# KL-regularized RL can suffer from pathological behavior during training.



Expert Demonstration



Learned Behavior

How can we avoid pathological behavior that may result in poor policies?

# We show:

**We show:**

why such pathologies may occur in theory;

**We show:**

why such pathologies may occur in theory;

why they occur in practice;

**We show:**

why such pathologies may occur in theory;

why they occur in practice;

how to prevent them.

**Goal**

## Goal

▸ Learn a good policy in **as few environment interactions as possible**

## Goal

- Learn a good policy in **as few environment interactions as possible**

## How?

- Use expert demonstrations to **give agents a head start**

## Goal

‣ Learn a good policy in **as few environment interactions as possible**

## How?

‣ Use expert demonstrations to **give agents a head start**

‣ Common approach: <span style="color:red">**Behavioral cloning**</span>

   ‣ offline: $\quad \mathcal{D}_0 = \{(\mathbf{s}_n, \mathbf{a}_n)\}_{n=1}^{N} = \{\overline{\mathbf{S}}, \overline{\mathbf{A}}\} \longrightarrow \pi_0(\cdot|\mathbf{s})$

## Goal

‣ Learn a good policy in **as few environment interactions as possible**

## How?

‣ Use expert demonstrations to **give agents a head start**

‣ Common approach: Behavioral cloning + **KL regularization**

  ‣ offline: $\mathcal{D}_0 = \{(\mathbf{s}_n, \mathbf{a}_n)\}_{n=1}^N = \{\overline{\mathbf{S}}, \overline{\mathbf{A}}\} \longrightarrow \pi_0(\cdot|\mathbf{s})$

  ‣ online: $\tilde{R}(\boldsymbol{\tau}_t) = \sum_{k=t}^{\infty} \gamma^k \left[ r(\mathbf{s}_k, \mathbf{a}_k) - \alpha \mathbb{D}_{\mathrm{KL}}(\pi(\cdot \mid \mathbf{s}_k) \| \pi_0(\cdot \mid \mathbf{s}_k)) \right]$
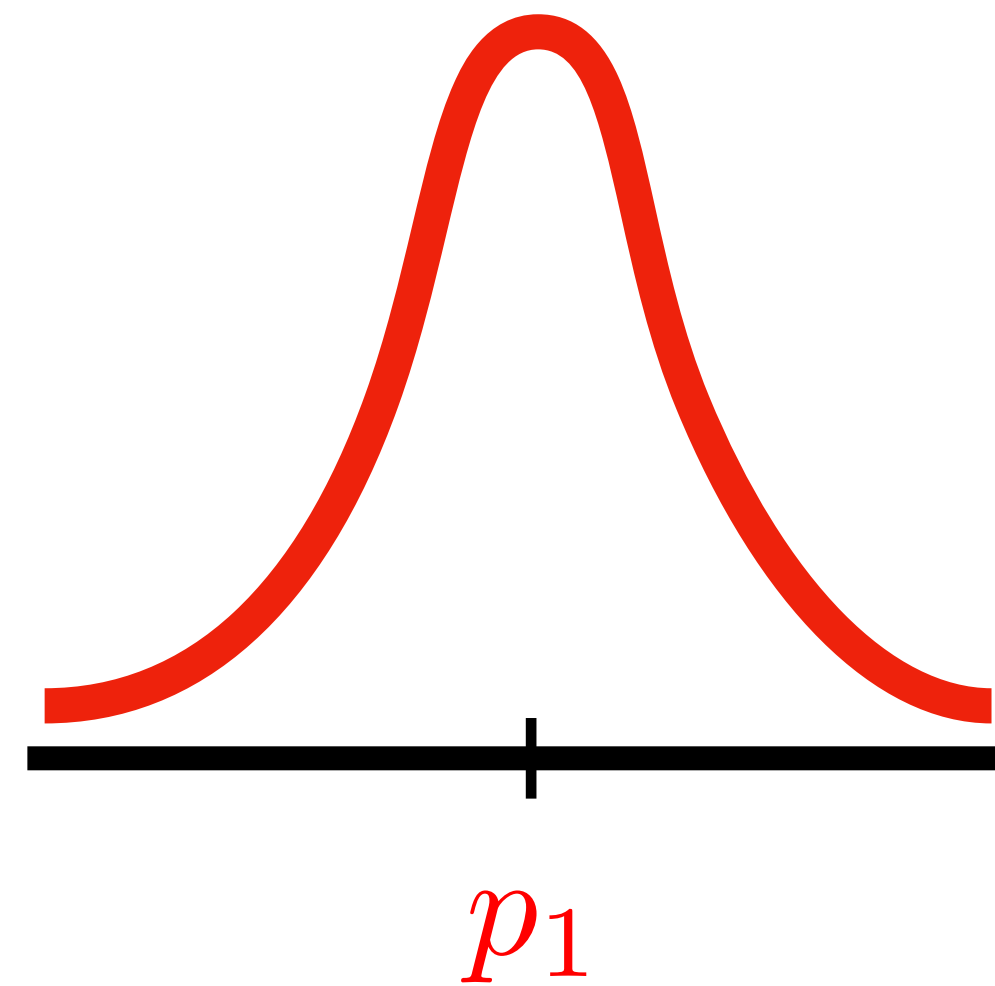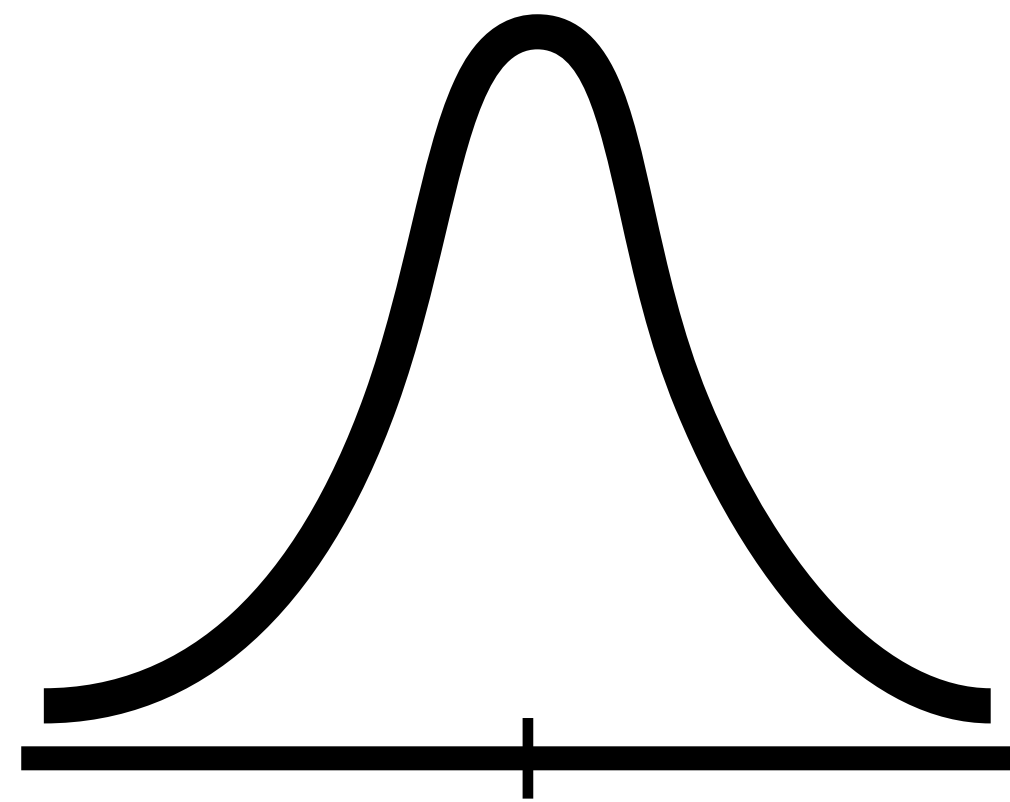
## Kullback-Leibler divergence

$$\tilde{R}\left(\boldsymbol{\tau}_t\right) = \sum_{k=t}^{\infty} \gamma^k \left[ r\left(\mathbf{s}_k, \mathbf{a}_k\right) - \alpha \textcolor{red}{\mathbb{D}_{\mathrm{KL}}\left(\pi\left(\cdot \mid \mathbf{s}_k\right) \| \pi_0\left(\cdot \mid \mathbf{s}_k\right)\right)} \right]$$

## Note!

**Kullback-Leibler divergence**

$$\tilde{R}\left(\boldsymbol{\tau}_t\right) = \sum_{k=t}^{\infty} \gamma^k \left[ r\left(\mathbf{s}_k, \mathbf{a}_k\right) - \alpha \mathbb{D}_{\mathrm{KL}}\left(\pi\left(\cdot \mid \mathbf{s}_k\right) \| \pi_0\left(\cdot \mid \mathbf{s}_k\right)\right) \right]$$

**Note!**

‣ KL divergence is well-defined (i.e., finite) **if and only if** learned policy is **absolutely continuous** w.r.t. behavioral policy

**Kullback-Leibler divergence**

$$\tilde{R}\left(\boldsymbol{\tau}_t\right) = \sum_{k=t}^{\infty} \gamma^k \left[r\left(\mathbf{s}_k, \mathbf{a}_k\right) - \alpha \mathbb{D}_{\mathrm{KL}}\left(\pi\left(\cdot \mid \mathbf{s}_k\right) \middle\| \pi_0\left(\cdot \mid \mathbf{s}_k\right)\right)\right]$$

**Note!**

‣ KL divergence is well-defined (i.e., finite) **if and only if** learned policy is **absolutely continuous** w.r.t. behavioral policy

‣ Potential failure mode: **degenerate behavioral policies**

# Could this be an issue in practice?

$$\mathbb{D}_{\mathrm{KL}}\left(q\|p_1\right) = 0$$



$p_1$

$$\mathbb{D}_{\mathrm{KL}}\left(q\|p_1\right) = 0 < \mathbb{D}_{\mathrm{KL}}\left(q\|p_2\right)$$



$p_2$

$$\mathbb{D}_{\mathrm{KL}}\left(q\|p_1\right) = 0 < \mathbb{D}_{\mathrm{KL}}\left(q\|p_2\right) < \mathbb{D}_{\mathrm{KL}}\left(q\|p_3\right)$$



$p_3$

$$\mathbb{D}_{\mathrm{KL}}\left(q\|p_1\right) = 0 < \mathbb{D}_{\mathrm{KL}}\left(q\|p_2\right) < \mathbb{D}_{\mathrm{KL}}\left(q\|p_3\right) < \mathbb{D}_{\mathrm{KL}}\left(q\|p\right) = \infty$$



$p$

## Potential Failure Mode

‣ If variance of behavioral policy tends to zero, KL blows up

# WHEN IS KL-REGULARIZED MEANINGFUL?

## Potential Failure Mode

‣ If variance of behavioral policy tends to zero, KL blows up

## Is this a problem in practice?

‣ How fast does the KL divergence blow up?

‣ Do commonly used behavioral policy have vanishingly small variance?

## Parametric policy predictive variance

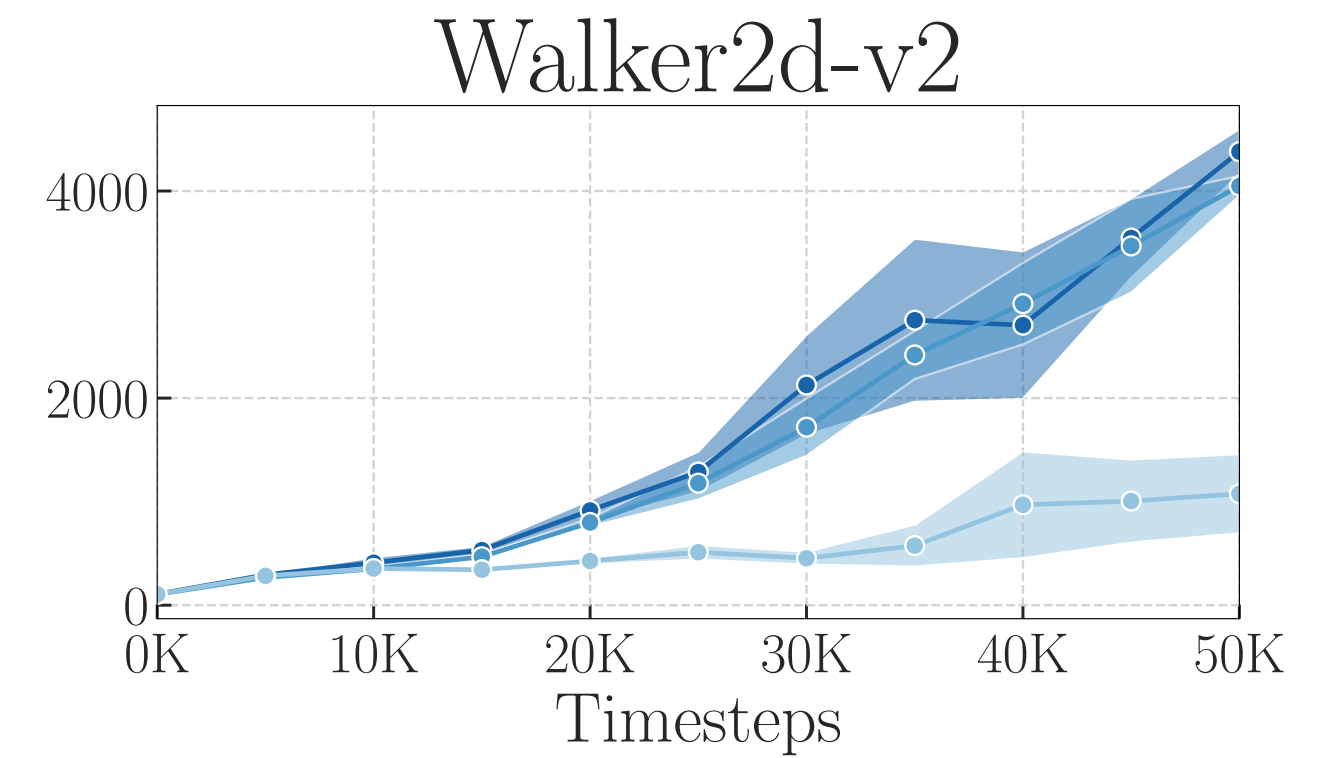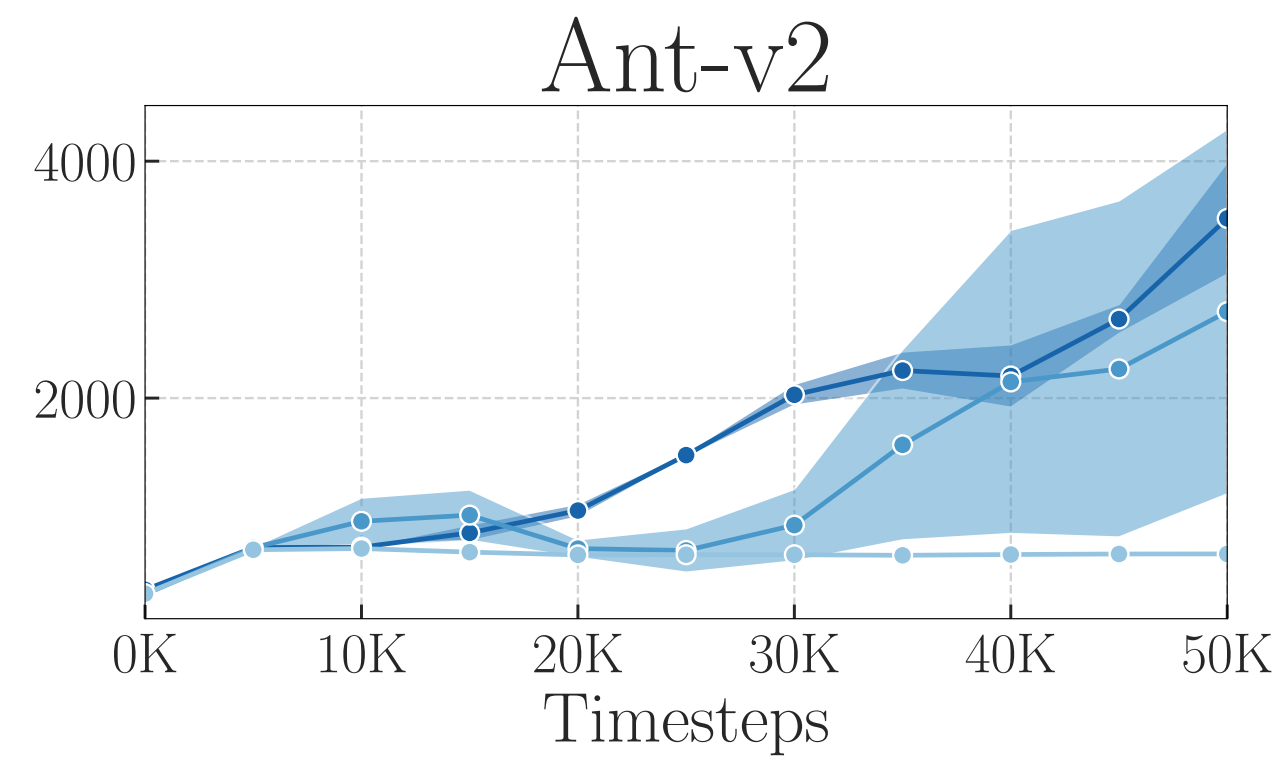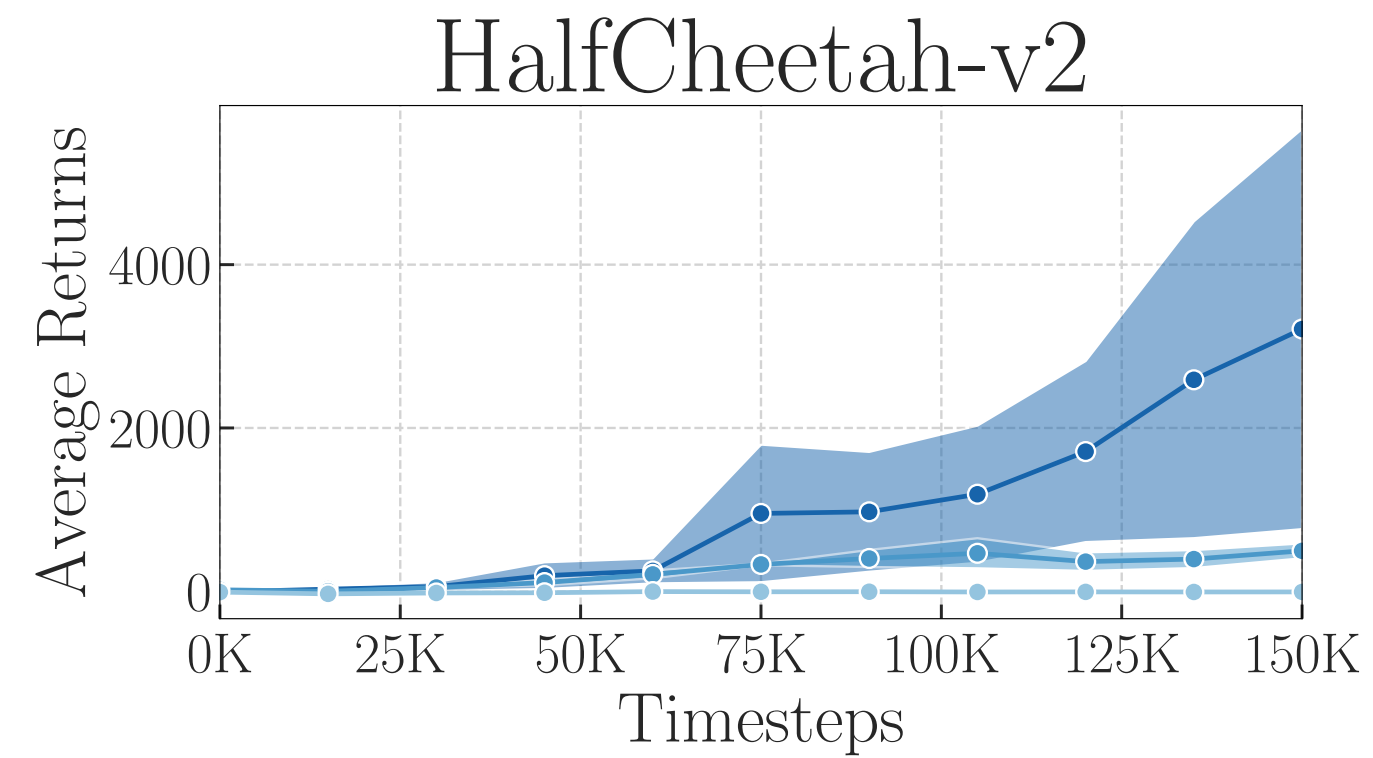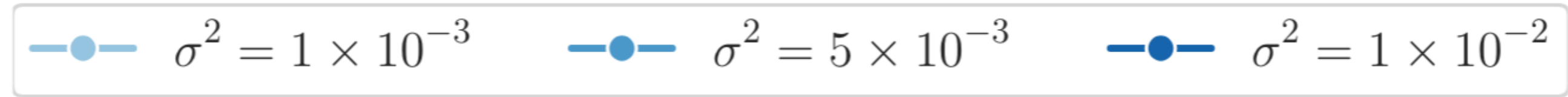‣ Collapse in predictive variance away from expert trajectories

## Parametric policy predictive variance

‣ Collapse in predictive variance away from expert trajectories

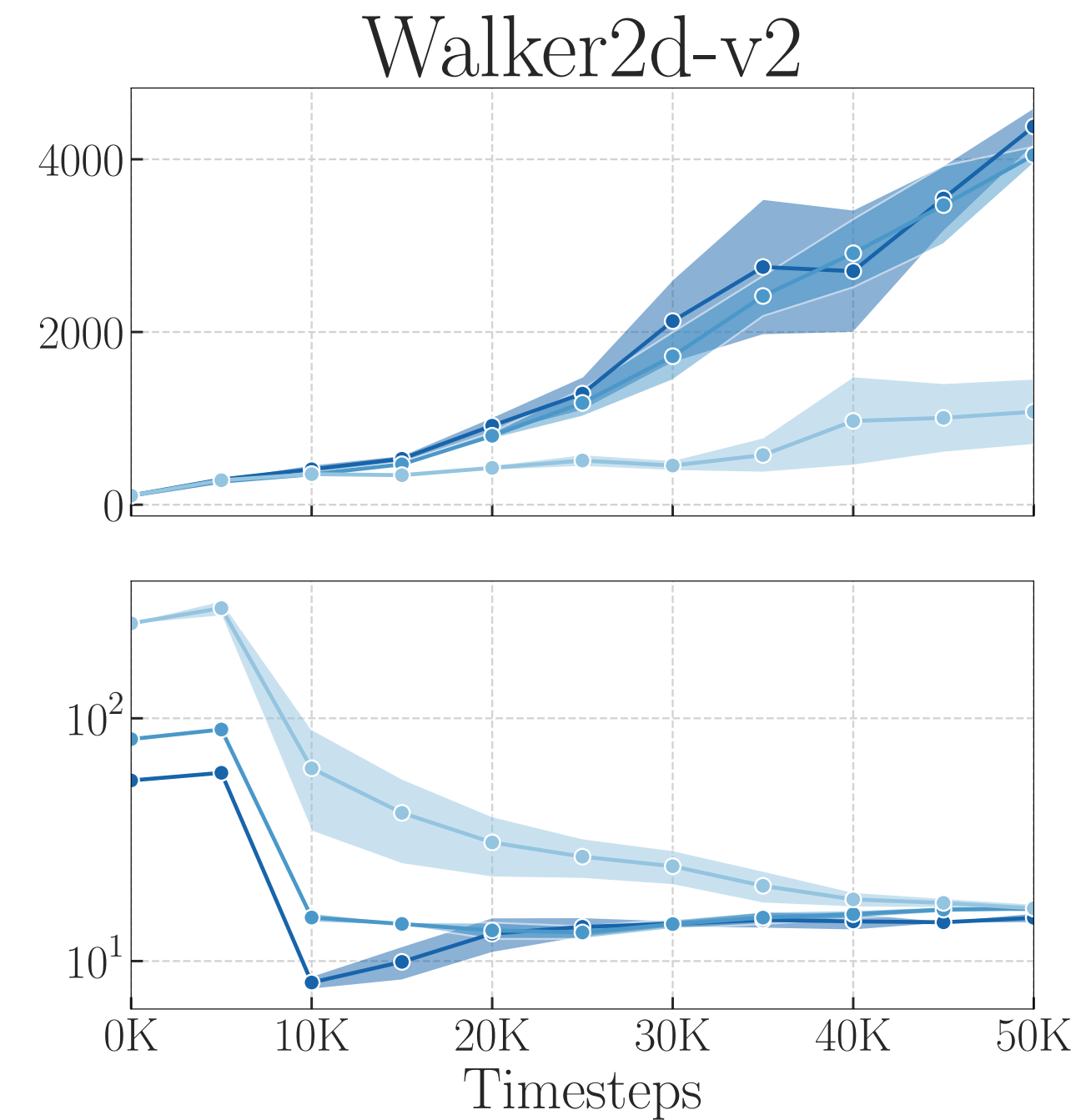# Parametric policy predictive variance

‣ Collapse in predictive variance away from expert trajectories

# Parametric policy predictive variance

‣ Collapse in predictive variance away from expert trajectories

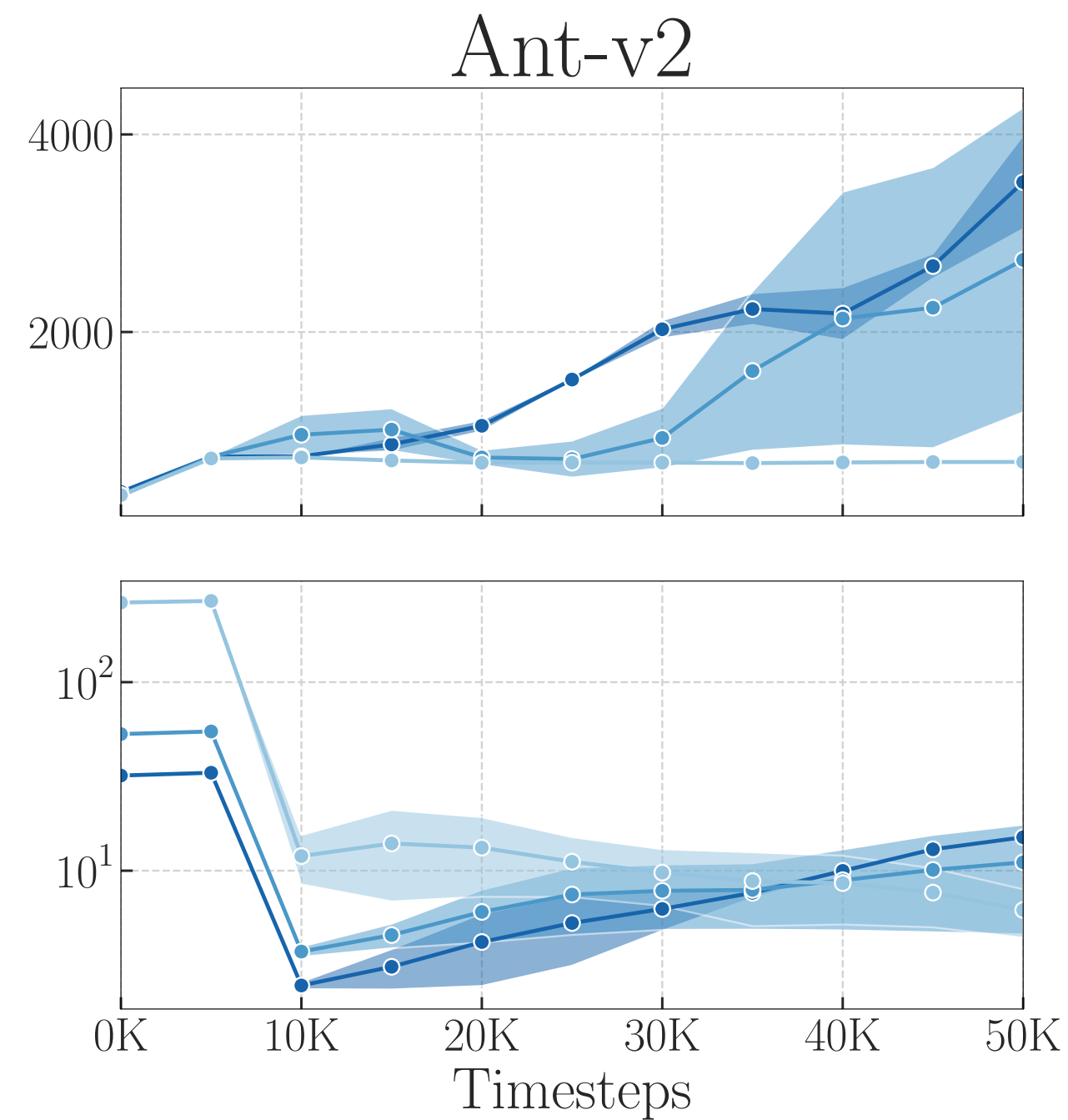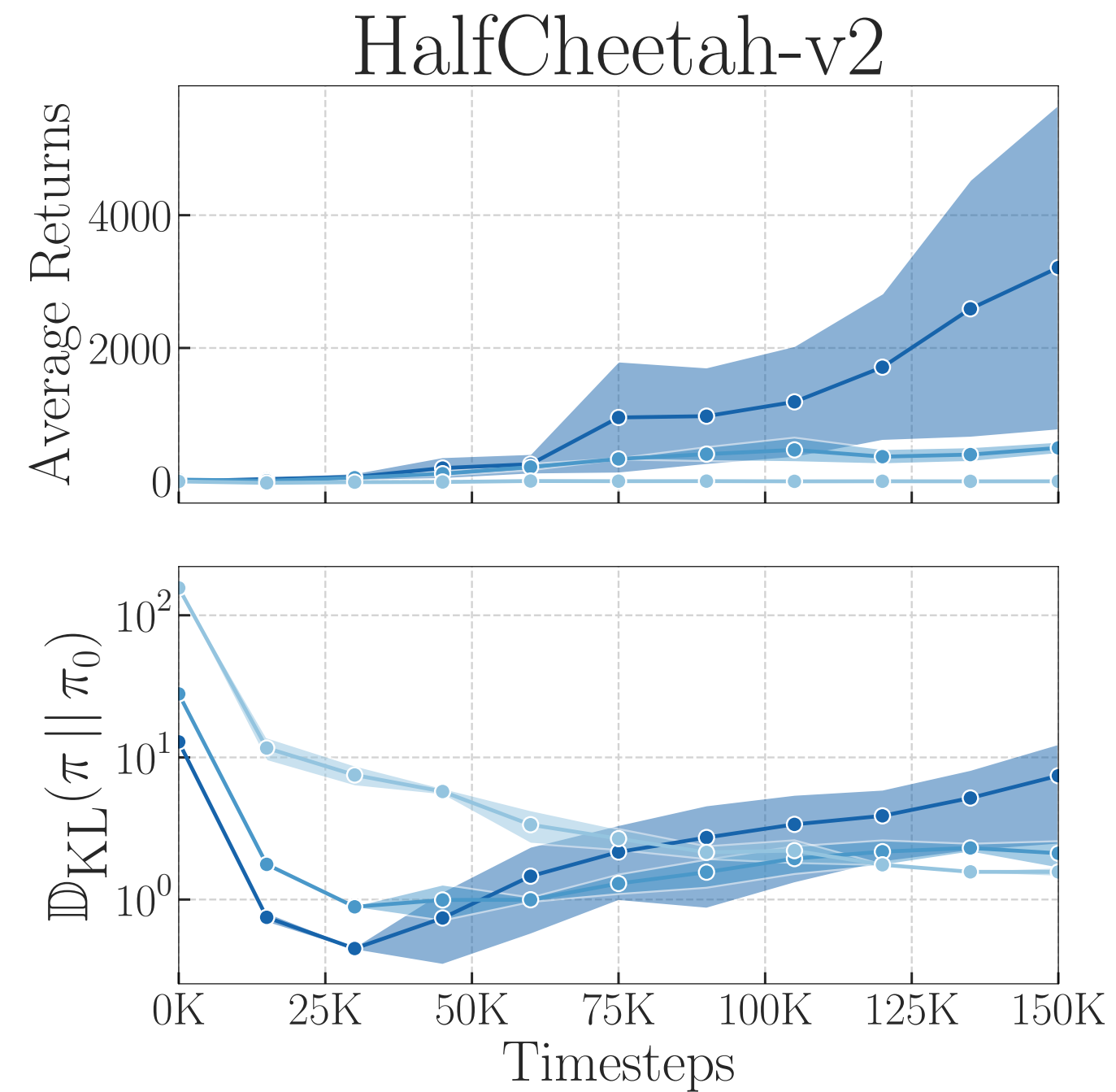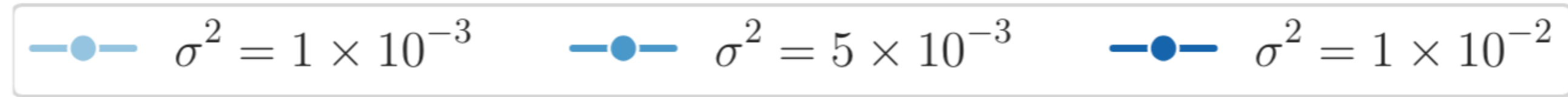$$\sigma^2 = 1 \times 10^{-3} \qquad \sigma^2 = 5 \times 10^{-3} \qquad \sigma^2 = 1 \times 10^{-2}$$

# EFFECT OF DECREASING PRIOR VARIANCE ON PERFORMANCE

$\sigma^2 = 1 \times 10^{-3}$     $\sigma^2 = 5 \times 10^{-3}$     $\sigma^2 = 1 \times 10^{-2}$

HalfCheetah-v2

Ant-v2

Walker2d-v2

# EFFECT OF DECREASING PRIOR VARIANCE ON PERFORMANCE

Legend: $\sigma^2 = 1 \times 10^{-3}$   $\sigma^2 = 5 \times 10^{-3}$   $\sigma^2 = 1 \times 10^{-2}$

HalfCheetah-v2          Ant-v2          Walker2d-v2

Average Returns

$\mathbb{D}_{\mathrm{KL}}(\pi \| \pi_0)$

Timesteps          Timesteps          Timesteps

Effect of Decreasing Prior Variance on Performance

**Proposition 1** (informal).

## Proposition 1 (informal).

‣ Let the objective function be given by

$$\tilde{R}\left(\boldsymbol{\tau}_t\right) = \sum_{k=t}^{\infty} \gamma^k \left[ r\left(\mathbf{s}_k, \mathbf{a}_k\right) - \alpha \mathbb{D}_{\mathrm{KL}}\left(\pi\left(\cdot \mid \mathbf{s}_k\right) \| \pi_0\left(\cdot \mid \mathbf{s}_k\right)\right) \right]$$

**Proposition 1** (informal).

‣ Let the objective function be given by

$$\tilde{R}\left(\boldsymbol{\tau}_t\right) = \sum_{k=t}^{\infty} \gamma^k \left[ r\left(\mathbf{s}_k, \mathbf{a}_k\right) - \alpha \mathbb{D}_{\mathrm{KL}}\left(\pi\left(\cdot \mid \mathbf{s}_k\right) \| \pi_0\left(\cdot \mid \mathbf{s}_k\right)\right) \right]$$

‣ Let the online and behavioral policies be Gaussian distributions

**Proposition 1** (informal).

‣ Let the objective function be given by

$$\tilde{R}\left(\boldsymbol{\tau}_t\right) = \sum_{k=t}^{\infty} \gamma^k \left[r\left(\mathbf{s}_k, \mathbf{a}_k\right) - \alpha \mathbb{D}_{\mathrm{KL}}\left(\pi\left(\cdot \mid \mathbf{s}_k\right) \| \pi_0\left(\cdot \mid \mathbf{s}_k\right)\right)\right]$$

‣ Let the online and behavioral policies be Gaussian distributions

‣ Let the online policy be parametrized by $\mathbf{a}_t = f_\phi\left(\epsilon_t; \mathbf{s}_t\right)$

**Proposition 1** (informal).

- Let the objective function be given by

$$\tilde{R}\left(\boldsymbol{\tau}_t\right) = \sum_{k=t}^{\infty} \gamma^k \left[ r\left(\mathbf{s}_k, \mathbf{a}_k\right) - \alpha \mathbb{D}_{\mathrm{KL}}\left(\pi\left(\cdot \mid \mathbf{s}_k\right) \| \pi_0\left(\cdot \mid \mathbf{s}_k\right)\right)\right]$$

- Let the online and behavioral policies be Gaussian distributions

- Let the online policy be parametrized by $\mathbf{a}_t = f_\phi\left(\epsilon_t; \mathbf{s}_t\right)$

- Then:

$$\left|\hat{\nabla}_\phi J_\pi(\phi)\right| \to \infty \quad \text{as} \quad \sigma_0^2 \to 0 \quad \text{with} \quad \mathcal{O}\left(\sigma_0^{-2}\left(\mathbf{s}_t\right)\right)$$

## Prevent predictive uncertainty collapse in behavioral policies

‣ Goal: increase predictive variance away from expert demonstrations

## Prevent predictive uncertainty collapse in behavioral policies

‣ Goal: increase predictive variance away from expert demonstrations

‣ Non-parametric Gaussian process behavioral policy

**Prevent predictive uncertainty collapse in behavioral policies**

- Goal: increase predictive variance away from expert demonstrations

- Non-parametric Gaussian process behavioral policy

  - Prior: $\quad A|s \sim \pi_0(\cdot|s) = \mathcal{GP}(m(s), k(s, s'))$

  - Posterior: $A|s, \mathcal{D}_0 \sim \pi_0(\cdot|s, \mathcal{D}_0) = \mathcal{GP}(\mu_0(s), \Sigma_0(s, s'))$

    - Mean: $\quad \mu_0(s) = m(s) + k(s, \bar{S})(k(\bar{S}, \bar{S}))^{-1}(\bar{A} - m(\bar{A}))$

    - Covariance: $\Sigma_0(s, s') = k(s, s') + k(s, \bar{S})k(\bar{S}, \bar{S})^{-1}k(\bar{S}, s')$

## Parametric

# Well-Calibrated Predictive Uncertainty



Parametric

Non-Parametric

# KL-Regularized RL with Non-Parametric Behavioral Policies
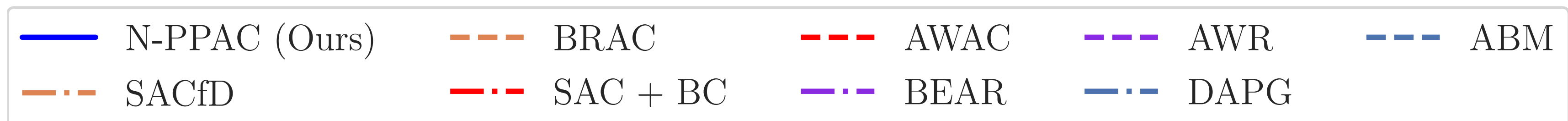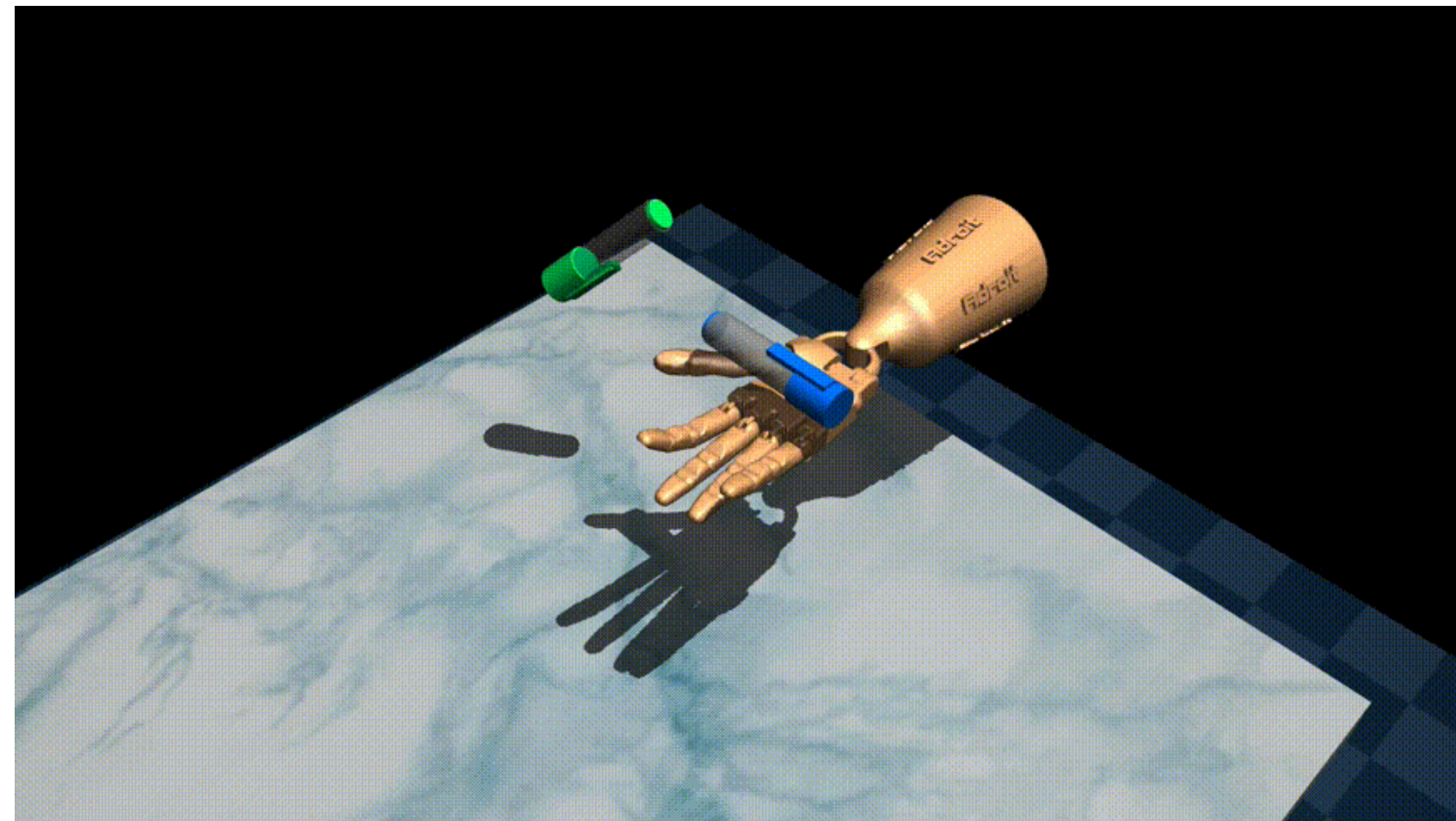
## MuJoCo Locomotion Tasks

## Dexterous Hand Manipulation Tasks



pen-binary-v0

door-binary-v0

## Dexterous Hand Manipulation Tasks



pen-binary-v0

door-binary-v0

## Dexterous Hand Manipulation: pen-binary-v0



pen-binary-v0

## Dexterous Hand Manipulation: pen-binary-v0



pen-binary-v0

Legend:
- N-PPAC (Ours)
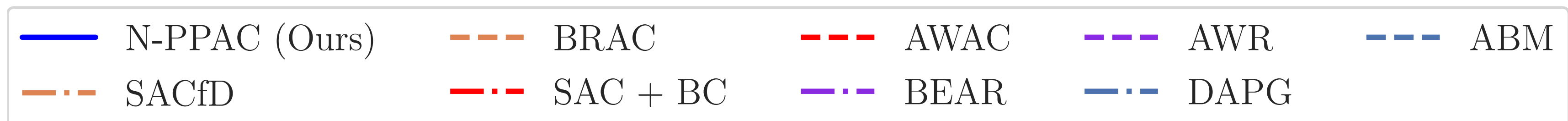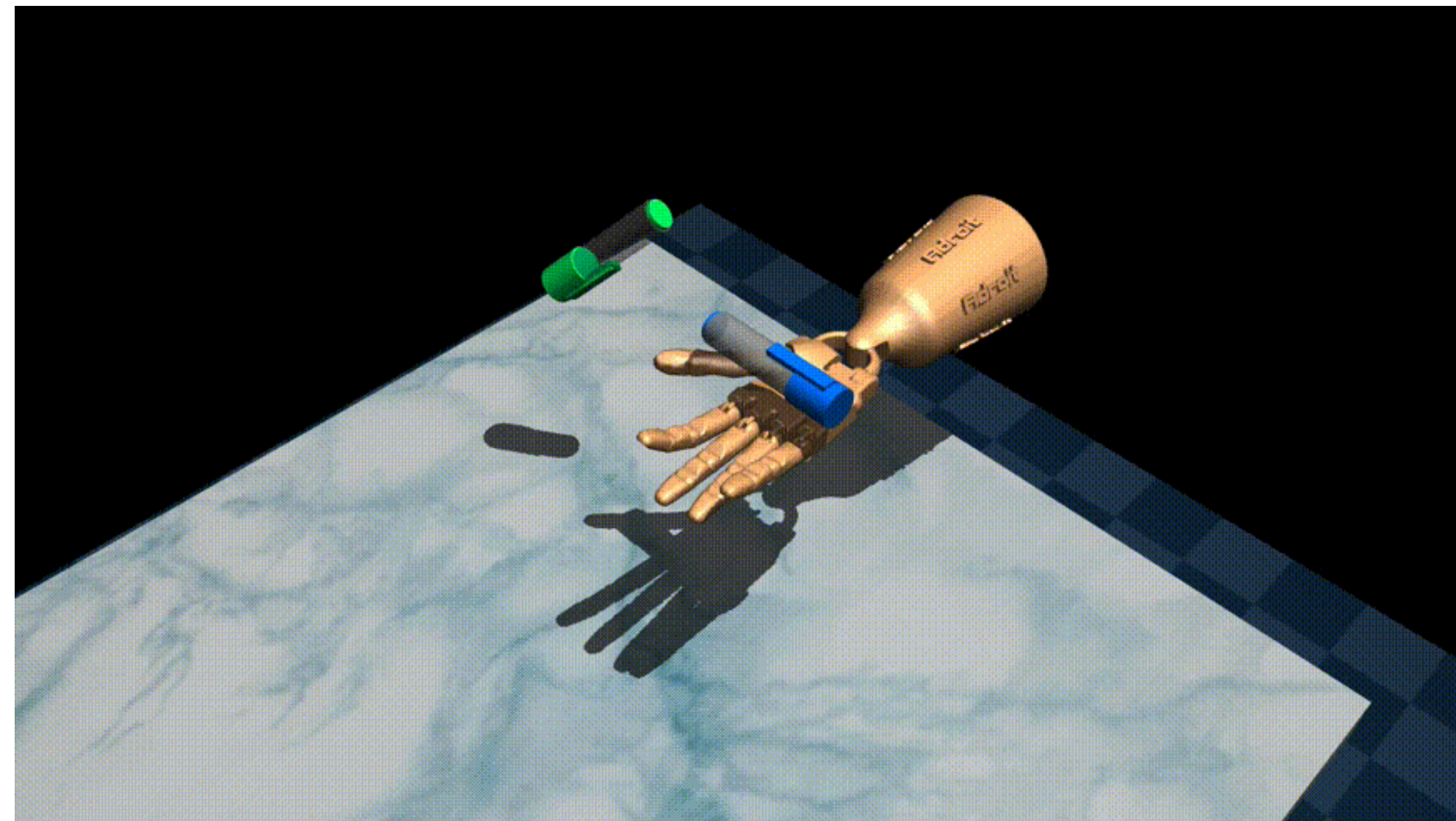- SACfD
- BRAC
- SAC + BC
- AWAC
- BEAR
- AWR
- DAPG
- ABM

# KL-Regularized RL with Non-Parametric Behavioral Policies
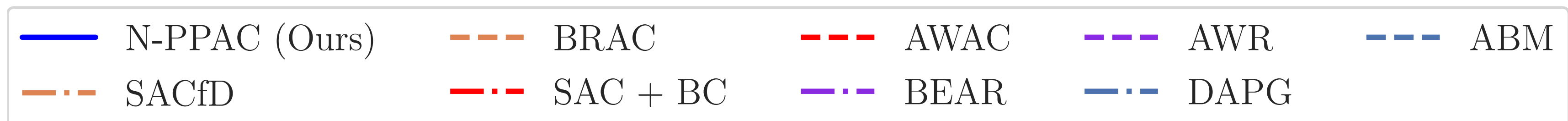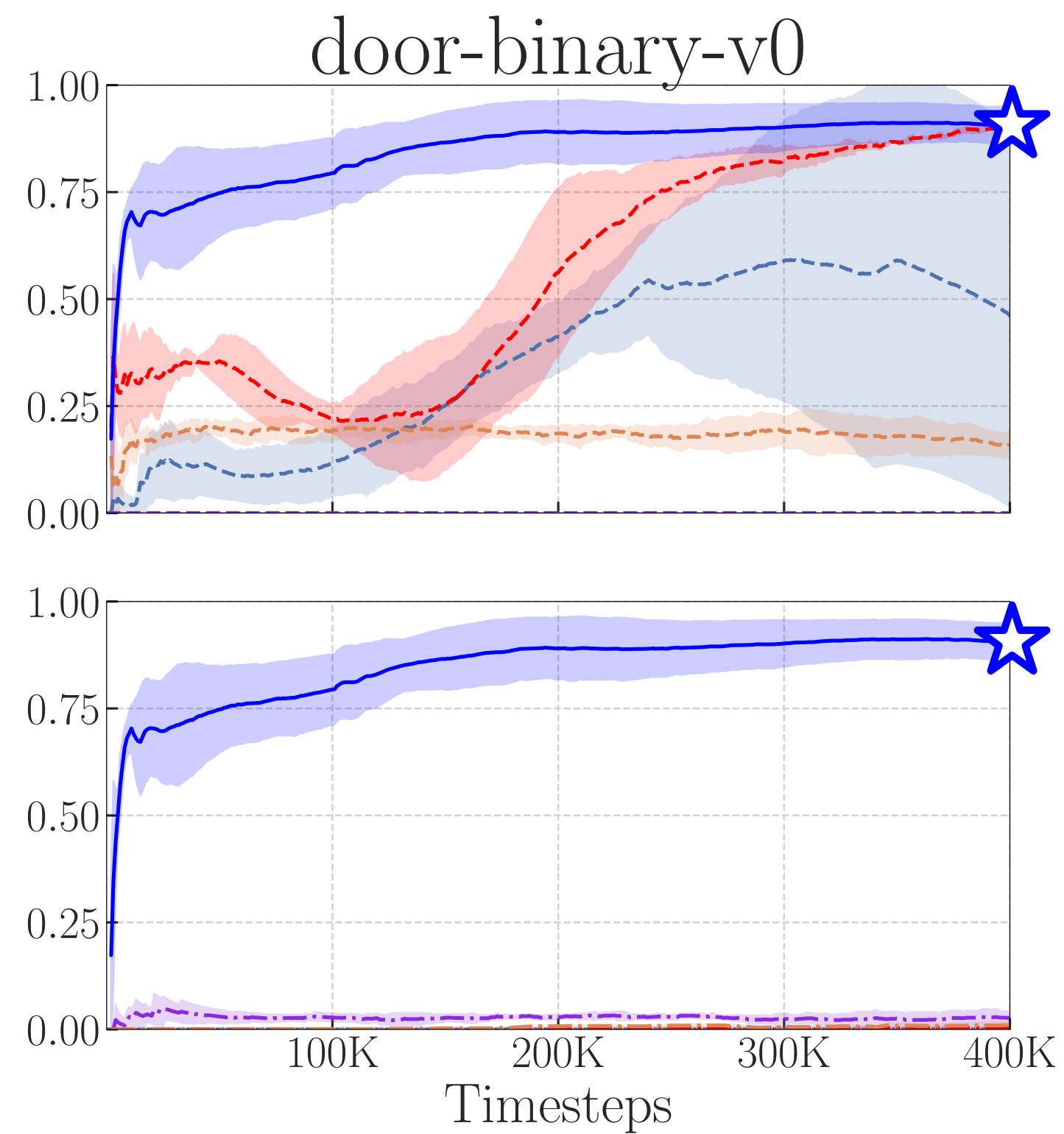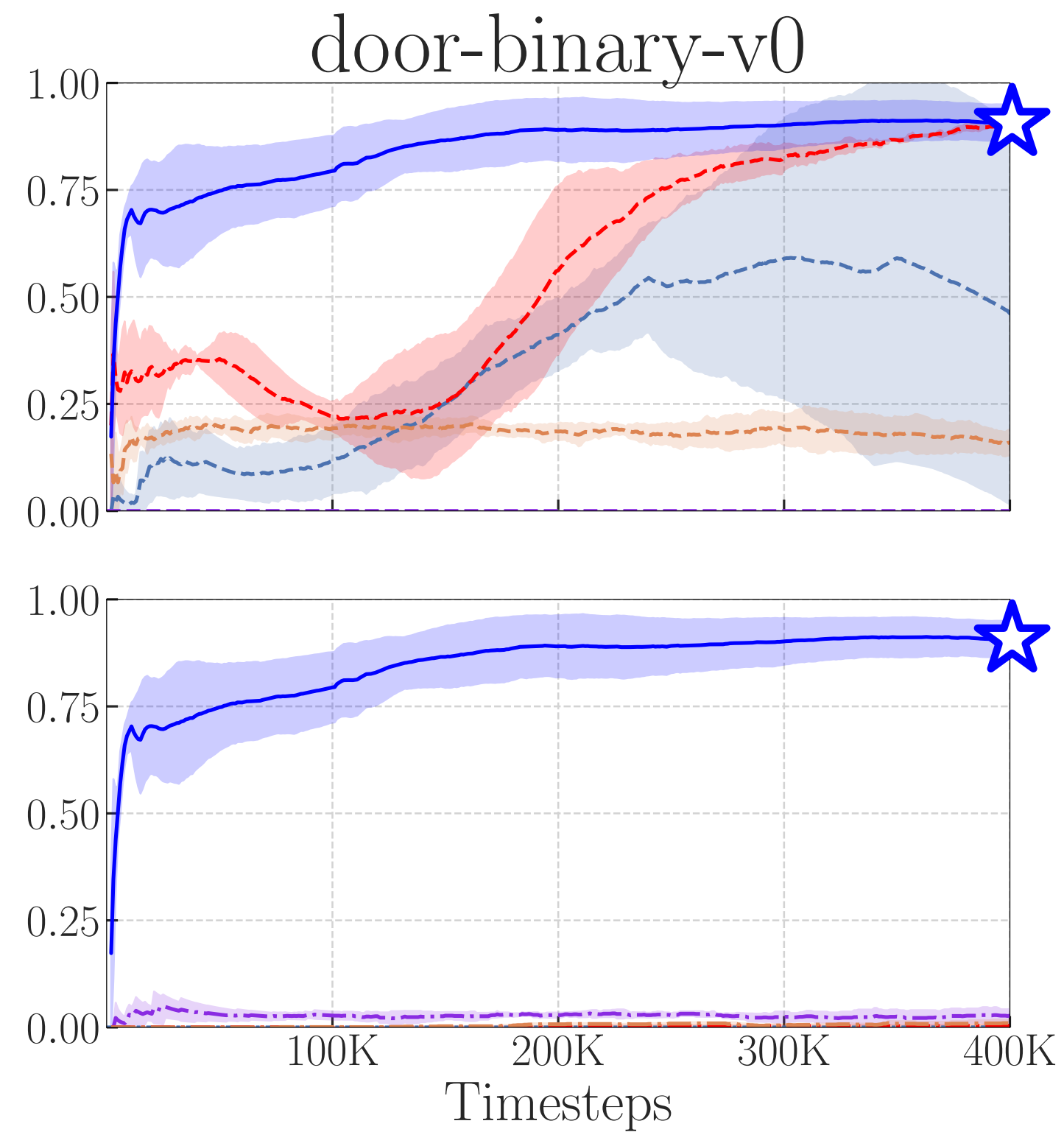
## Dexterous Hand Manipulation Example: door-binary-v0

door-binary-v0

Timesteps

N-PPAC (Ours) — BRAC — AWAC — AWR — ABM
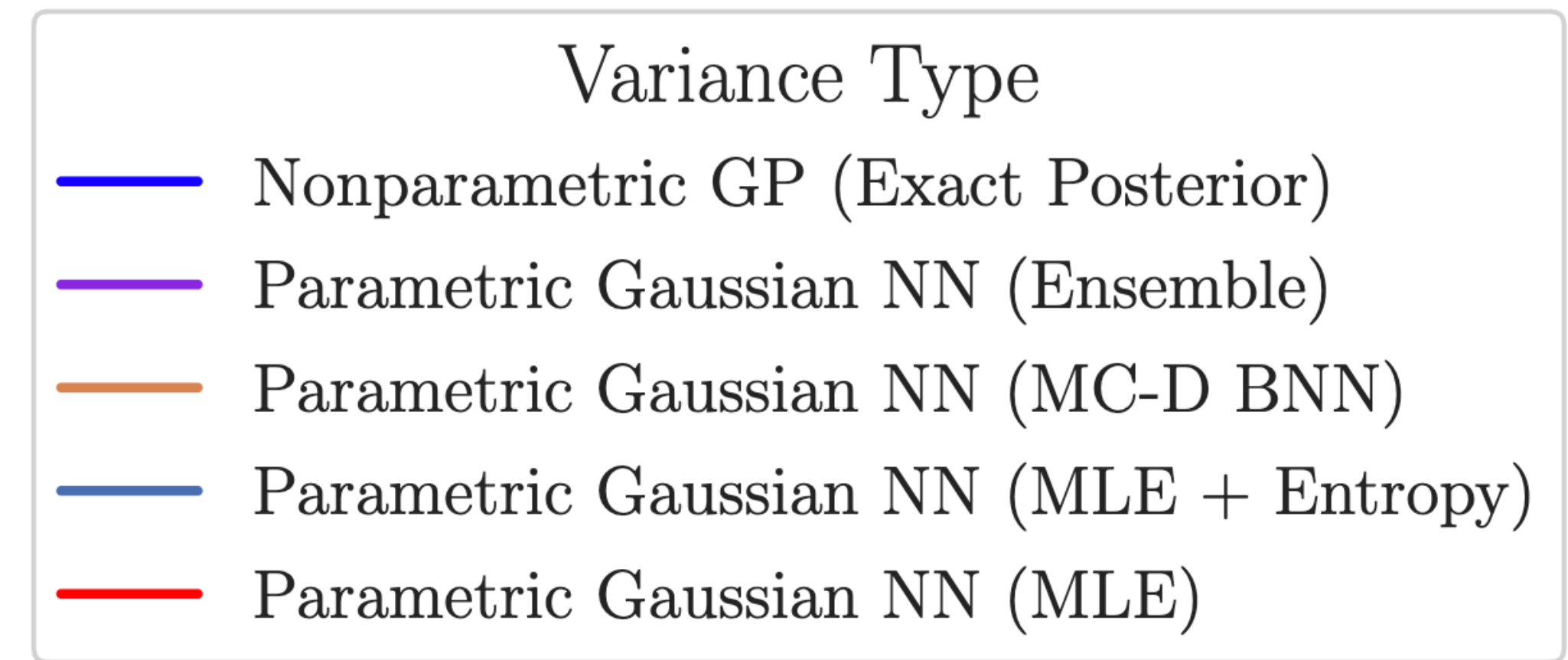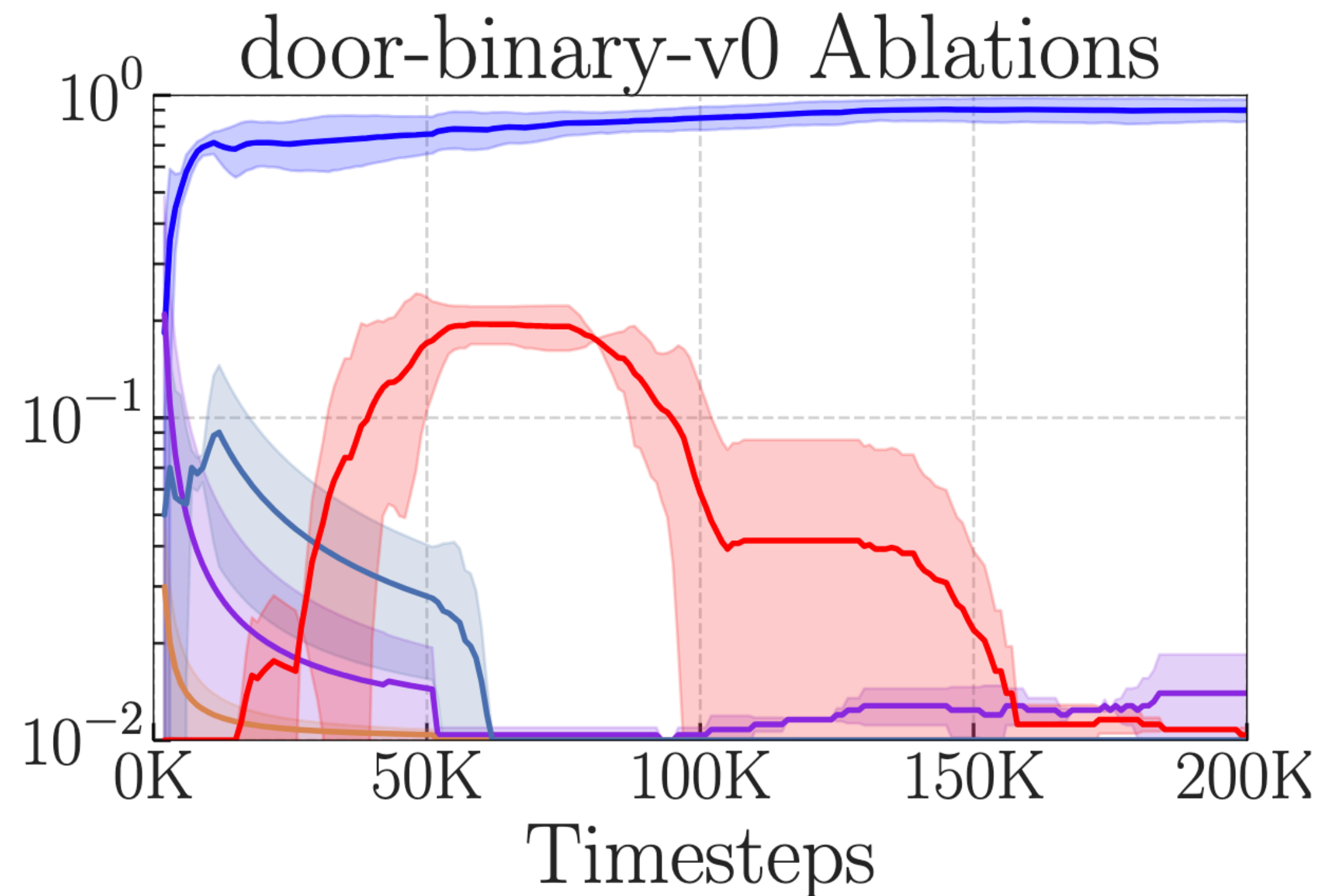SACfD — SAC + BC — BEAR — DAPG

## Dexterous Hand Manipulation Example: door-binary-v0



door-binary-v0

**Fixing the pathological training dynamics in
KL-regularized RL leads to state-of-the-art performance**

- ‣ Bayesian Neural Networks

- ‣ Deep Ensembles

- ‣ Lower-bounding Parametric Behavioral Policy Variance

**MuJoCo Locomotion Example: HalfCheetah**

**KL-regularized RL can suffer from pathological behavior during training.**

KL-regularized RL can suffer from <span style="color:red">pathological behavior</span> during training.

KL-regularized RL can suffer from <span style="color:red">pathological behavior</span> during training.

The pathology can be remedied by non-parametric behavioral policies.

KL-regularized RL can suffer from <span style="color:red">pathological behavior</span> during training.

The pathology can be <span style="color:red">remedied</span> by <span style="color:red">non-parametric</span> behavioral policies.

KL-regularized RL can suffer from pathological behavior during training.

The pathology can be remedied by non-parametric behavioral policies.

Fixing the pathology leads to state-of-the-art policies and data-efficient online training.

KL-regularized RL can suffer from pathological behavior during training.

The pathology can be remedied by non-parametric behavioral policies.

Fixing the pathology leads to state-of-the-art policies and data-efficient online training.

# Thank You!



**Tim G. J. Rudner***    **Cong Lu***    Michael A. Osborne    Yarin Gal    Yee Whye Teh

Correspondence:    tim.rudner@cs.ox.ac.uk

Project Website:    https://sites.google.com/view/nppac