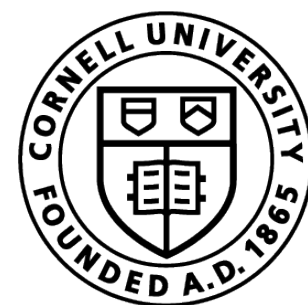# Scaling Offline RL via
# Efficient and Expressive Shortcut Models

**Nicolas Espinosa Dice**

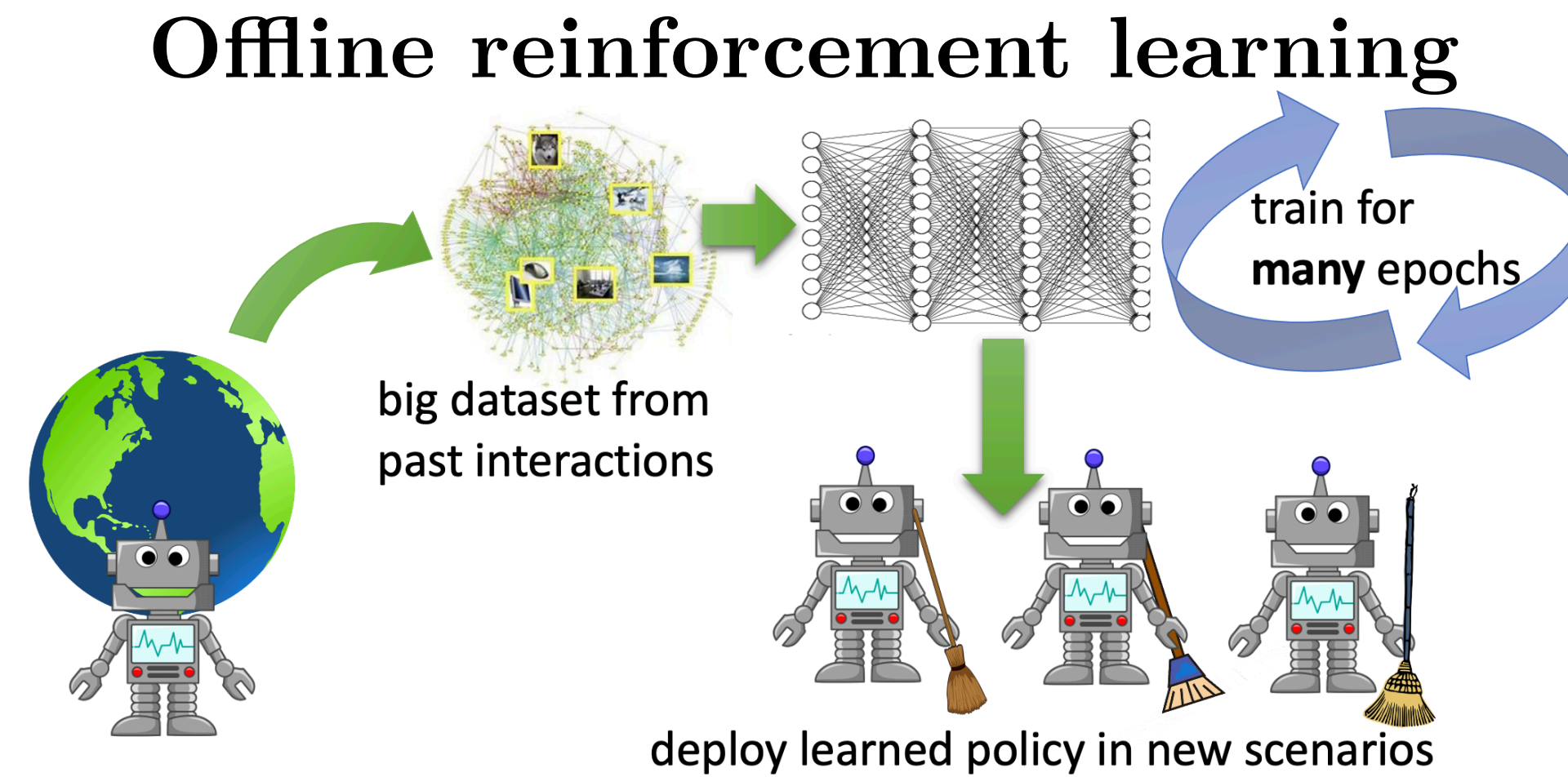*Joint work with*

Yiyi Zhang, Yiding Chen, Bradley Guo, Owen Oertell, Gokul Swamy, Kianté Brantley, Wen Sun

# How do we scale offline reinforcement learning?

**Axes of scale**

**Offline reinforcement learning**



big dataset from
past interactions

train for
**many** epochs

deploy learned policy in new scenarios

# How do we scale offline reinforcement learning?

**Axes of scale**

1. **Data**

**Offline reinforcement learning**



big dataset from past interactions

train for **many** epochs

deploy learned policy in new scenarios

# How do we scale offline reinforcement learning?

**Axes of scale**

1. Data
2. **Models**

**Offline reinforcement learning**



big dataset from past interactions

train for **many** epochs

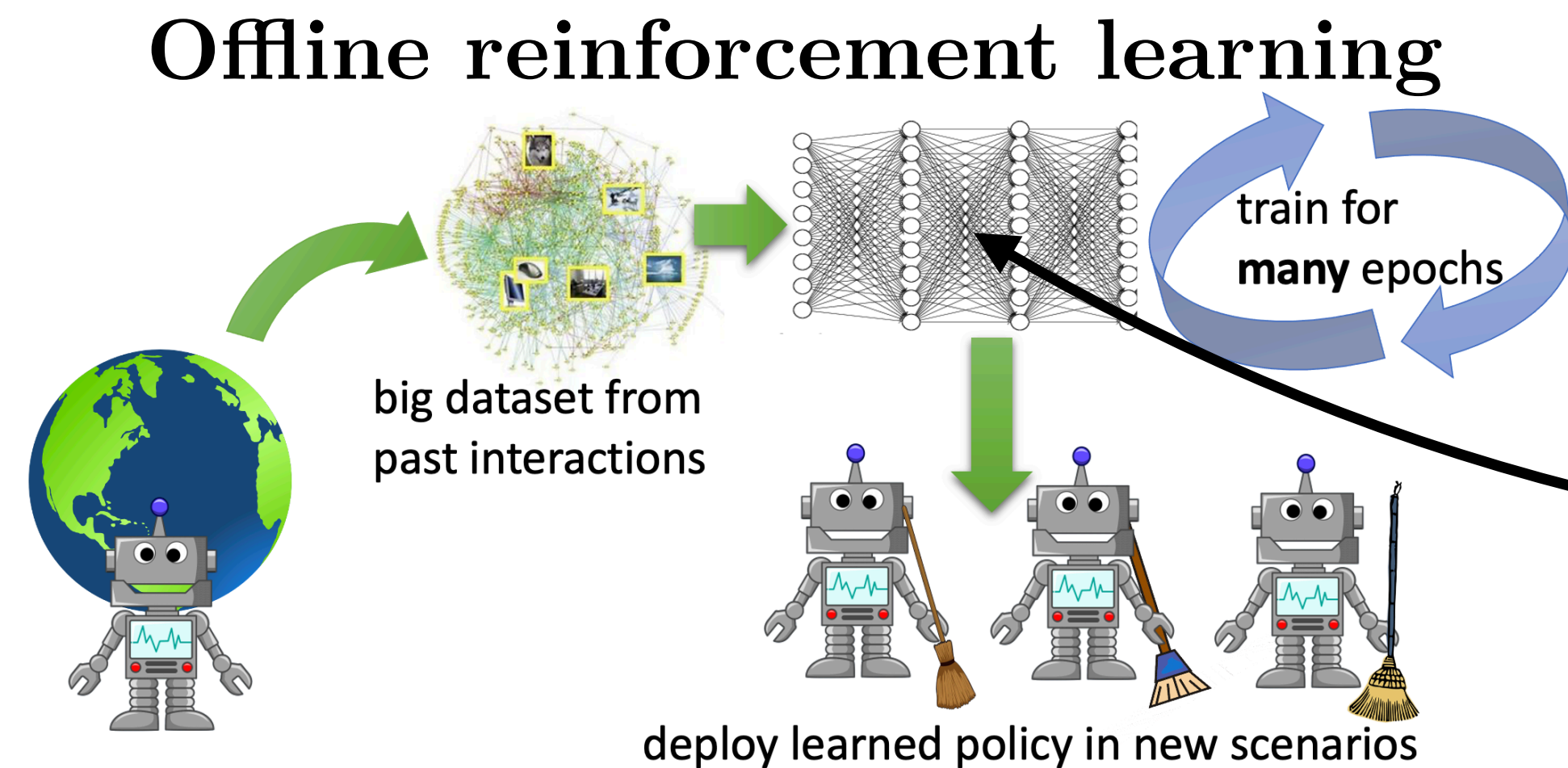deploy learned policy in new scenarios

Represent policy via modern generative models (e.g. diffusion & flow models)

# How do we scale offline reinforcement learning?

**Axes of scale**

1. Data
2. **Models**

**Offline reinforcement learning**



big dataset from past interactions

train for **many** epochs

deploy learned policy in new scenarios

Represent policy via modern generative models (e.g. diffusion & flow models)



**Advantages**: model multi-modal data
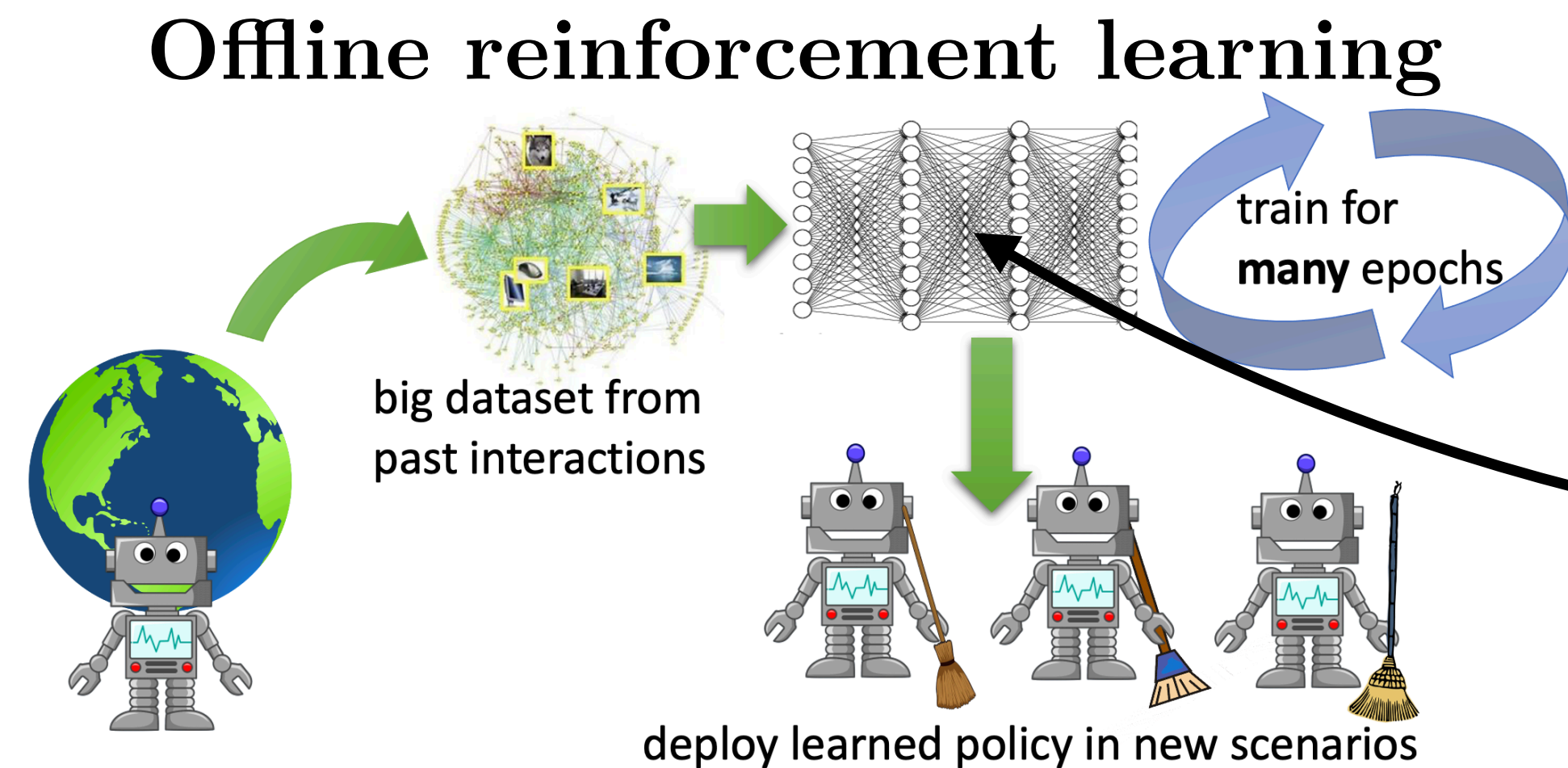
# How do we scale offline reinforcement learning?

**Axes of scale**

1. Data
2. Models
3. **Compute**

**Offline reinforcement learning**



big dataset from past interactions

train for **many** epochs

deploy learned policy in new scenarios

Represent policy via modern generative models (e.g. diffusion & flow models)



**Advantages**: model multi-modal data

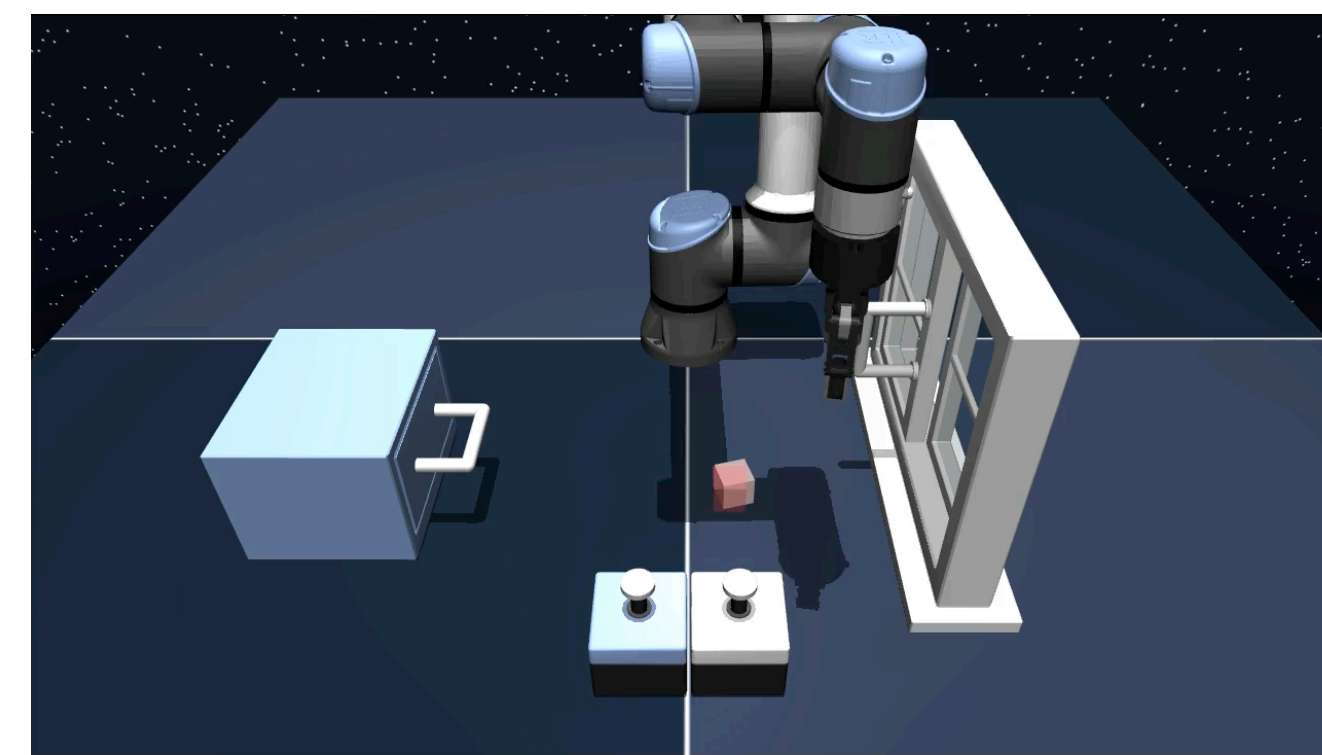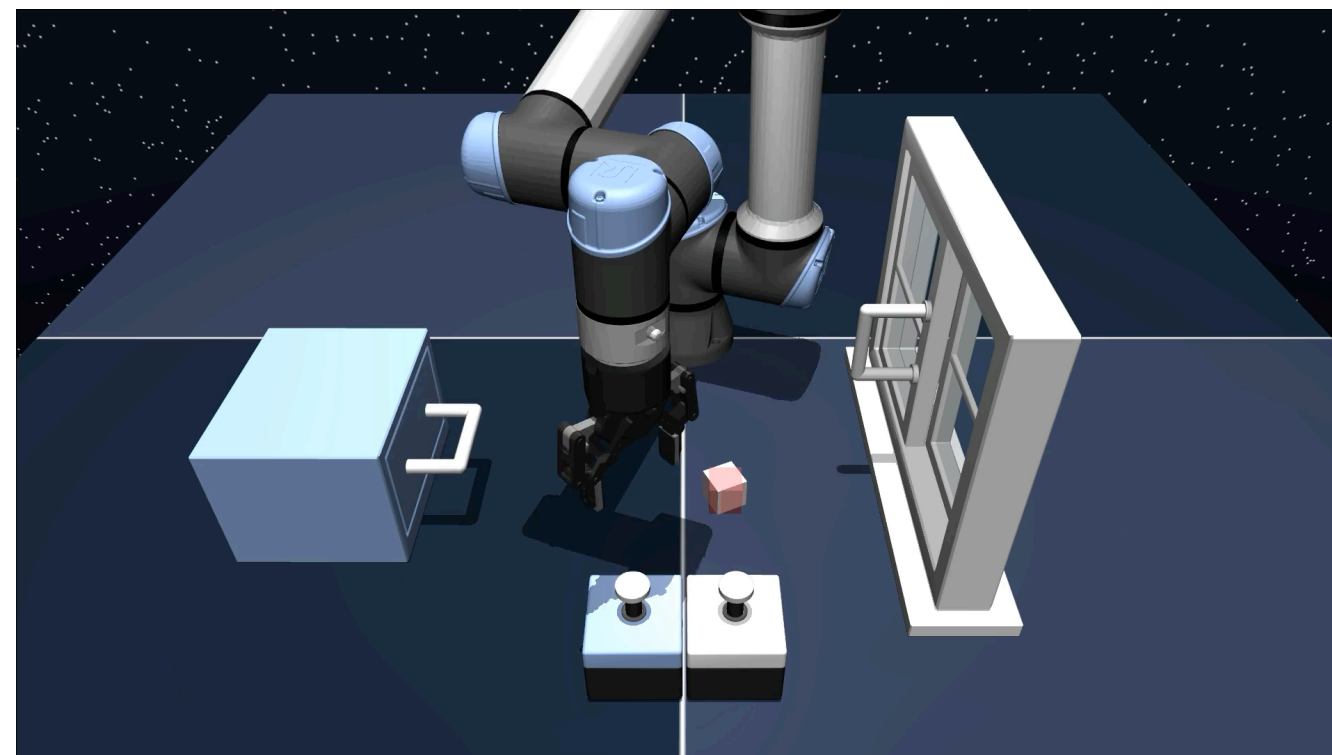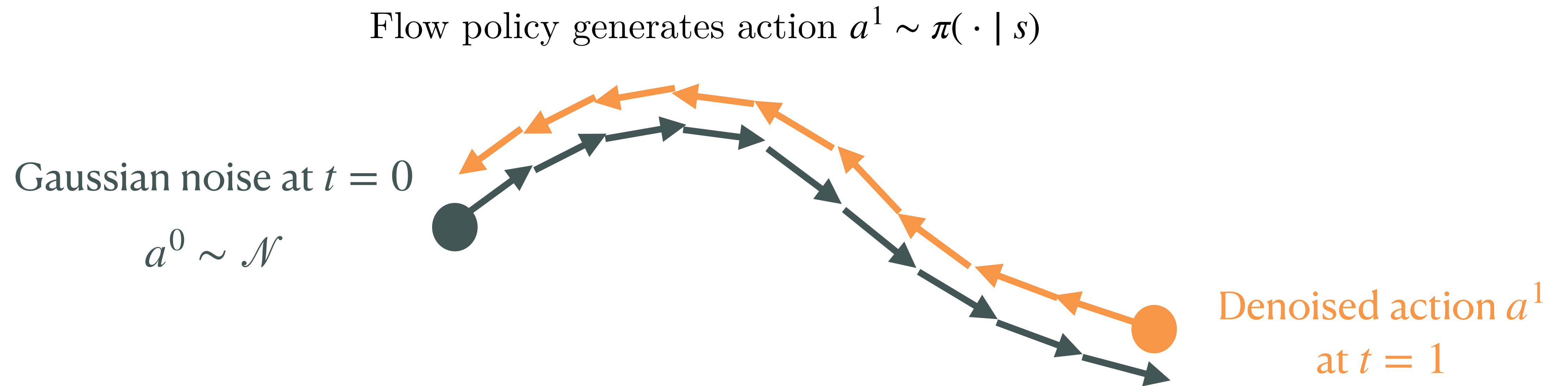# RL with diffusion/flow policies

Flow policy generates action $a^1 \sim \pi(\,\cdot\,|\,s)$

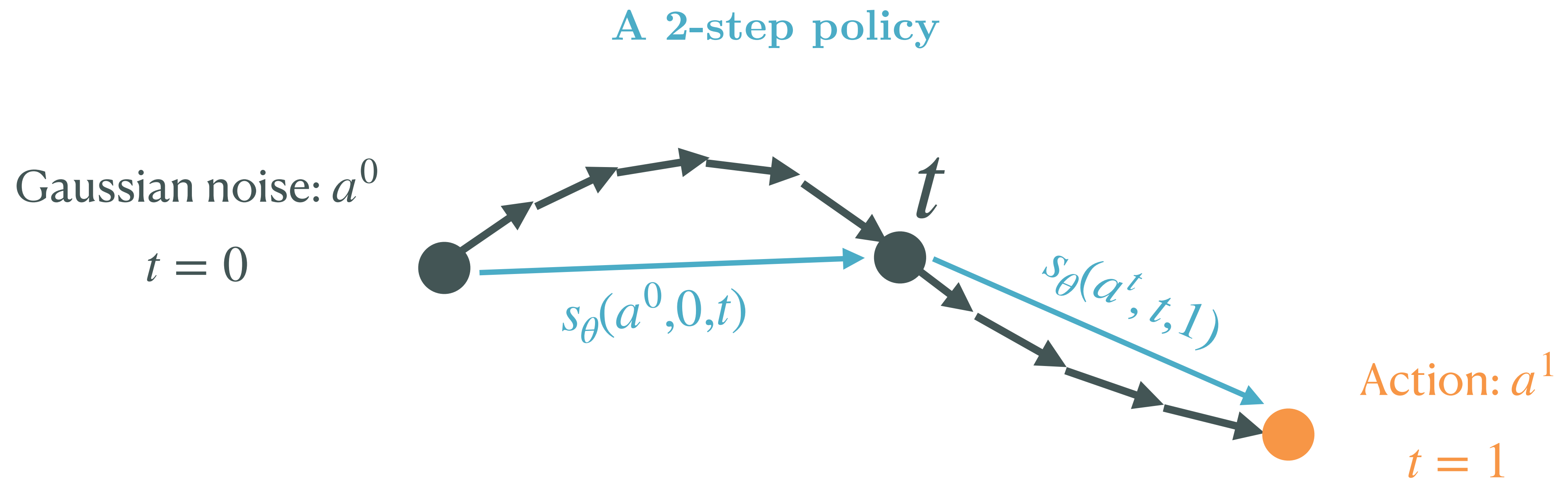Gaussian noise at $t = 0$

$a^0 \sim \mathcal{N}$

Denoised action $a^1$ at $t = 1$

Given $Q(s, a^1)$, optimize $\theta$ via backpropagation through time (BPTT)
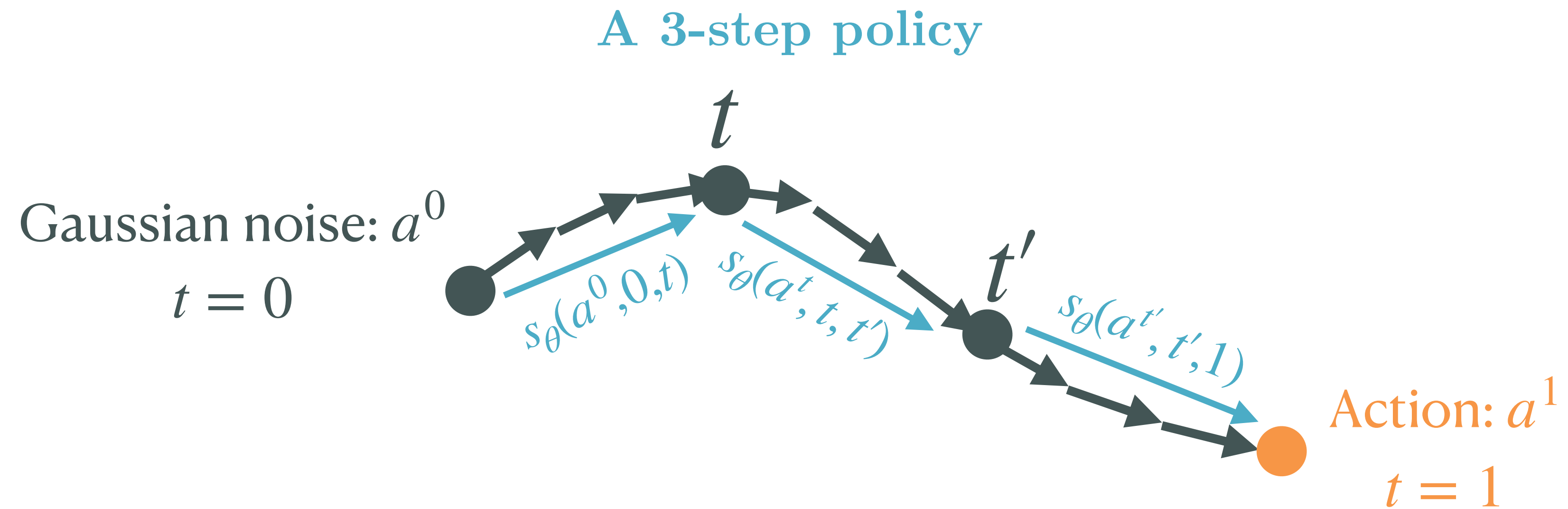
$$\frac{\partial Q}{\partial \theta} = \frac{\partial Q}{\partial a^1} \cdot \frac{\partial a^1}{\partial \theta}$$

# *Shortcut models*: a flow model with "shortcuts"

**A 2-step policy**



Gaussian noise: $a^0$

$t = 0$

$s_\theta(a^0, 0, t)$

$t$

$s_\theta(a^t, t, 1)$

Action: $a^1$

$t = 1$

$s_\theta$: a differentiable function modeling a ***jump*** for some time interval $t$

# *Shortcut models*: a flow model with "shortcuts"



**A 3-step policy**

Gaussian noise: $a^0$

$t = 0$

$s_\theta(a^0, 0, t)$

$s_\theta(a^t, t, t')$

$t$

$t'$

$s_\theta(a^{t'}, t', 1)$

Action: $a^1$

$t = 1$

**Key Advantage:**

Unlike flow models, shortcut models
are expressive with just a few steps

# Optimizing shortcut policies



**A 3-step policy**

Gaussian noise: $a^0$

$t = 0$

$s_\theta(a^0, 0, t)$

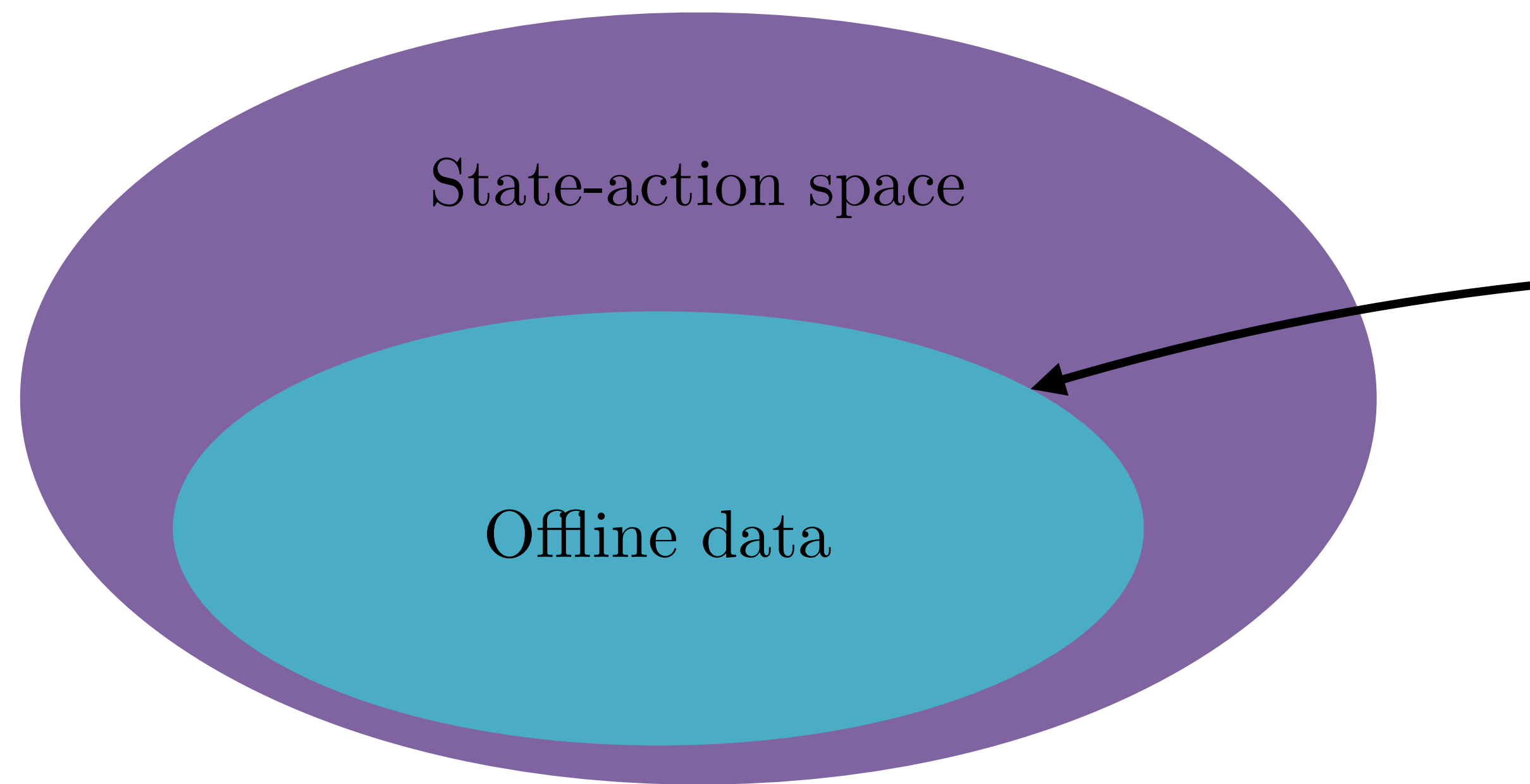$s_\theta(a^t, t, t')$

$s_\theta(a^{t'}, t', 1)$

$t$

$t'$

Action: $a^1$

$t = 1$

Given $Q(s, a^1)$, optimize $\theta$ via backpropagation k-steps

$$\frac{\partial Q}{\partial \theta} = \frac{\partial Q}{\partial a^1} \cdot \frac{\partial a^1}{\partial \theta}$$

# Training a shortcut policy in offline RL setting



**Key Challenge**
Training with offline data alone easily runs into covariate shift problems

**Goal**
Design algorithms that restrict search space inside the offline data's coverage
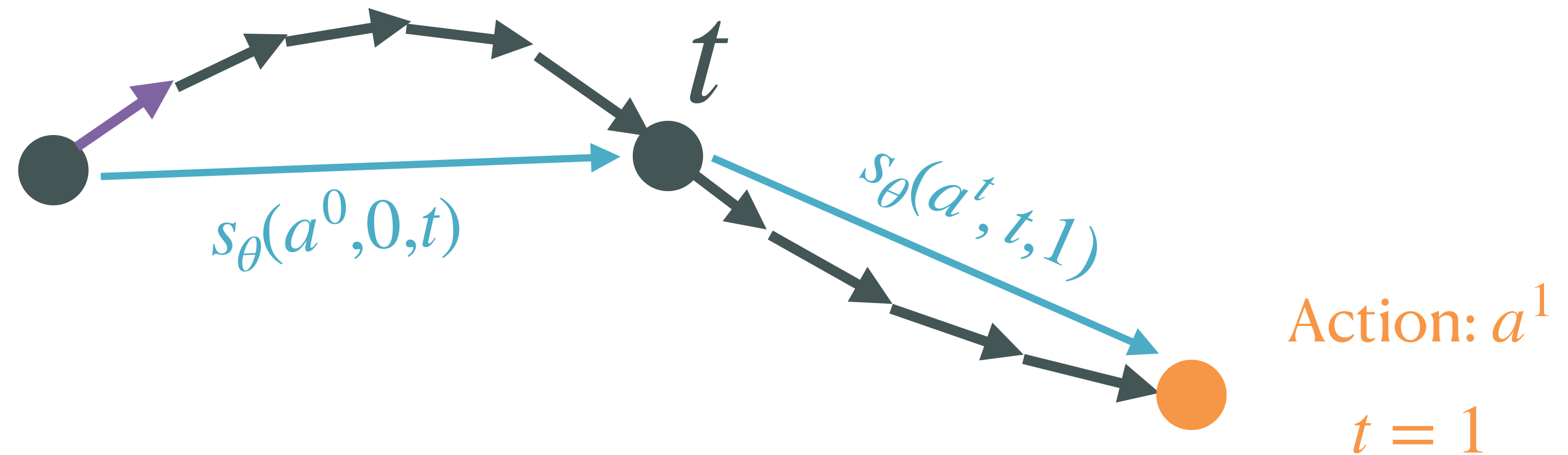
**Solution**
Regularize to the offline data

# How do we regularize to the offline data?

Learned policy

**Key idea**

Behavior policy

Match distributions of $\pi_\theta$ and $\pi_B$ through *flow matching*

$\rightarrow$ match velocity of $\pi_B$'s flow

Gaussian noise: $a^0$

$t = 0$

$t$

$s_\theta(a^0, 0, t)$
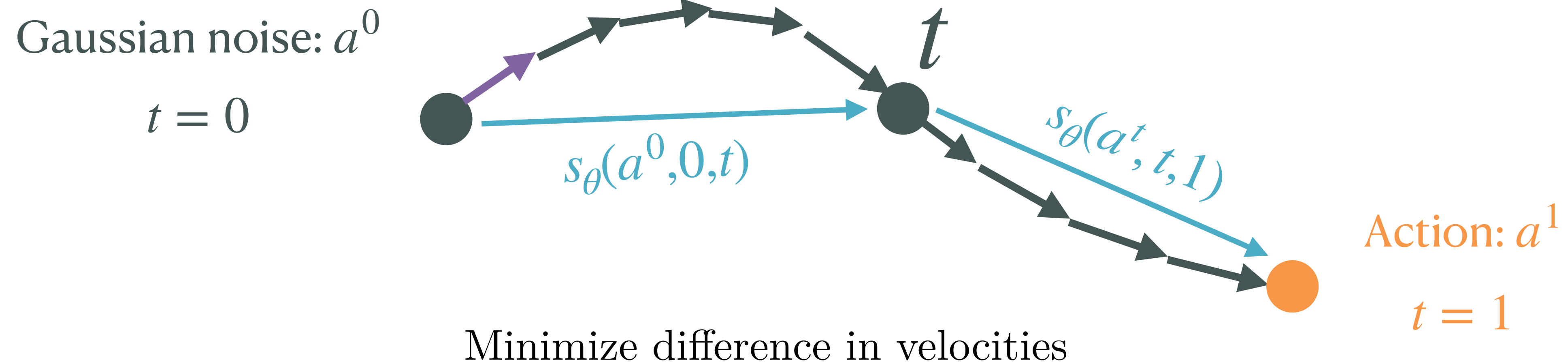
$s_\theta(a^t, t, 1)$

Action: $a^1$

$t = 1$

When $t \rightarrow 0$, $s_\theta$ models
instantaneous velocity of generator's flow

# How do we regularize to the offline data?

Learned policy

Behavior policy

**Key idea**

Match distributions of $\pi_\theta$ and $\pi_B$ through *flow matching*

$\rightarrow$ match velocity of $\pi_B$'s flow

Gaussian noise: $a^0$

$t = 0$

$s_\theta(a^0, 0, t)$

$t$

$s_\theta(a^t, t, 1)$

Action: $a^1$

$t = 1$

Minimize difference in velocities

$$\min_\theta \| \underbrace{s_\theta(a_t, t, t + \Delta)}_{\text{flow at noised action } a_t} - \underbrace{\text{velocity}(\pi_B(\,\cdot\,|\,s), t)\|_2}_{\text{approximated with } \pi_B \text{ actions}}$$

# Scalable Offline Reinforcement Learning (SORL)
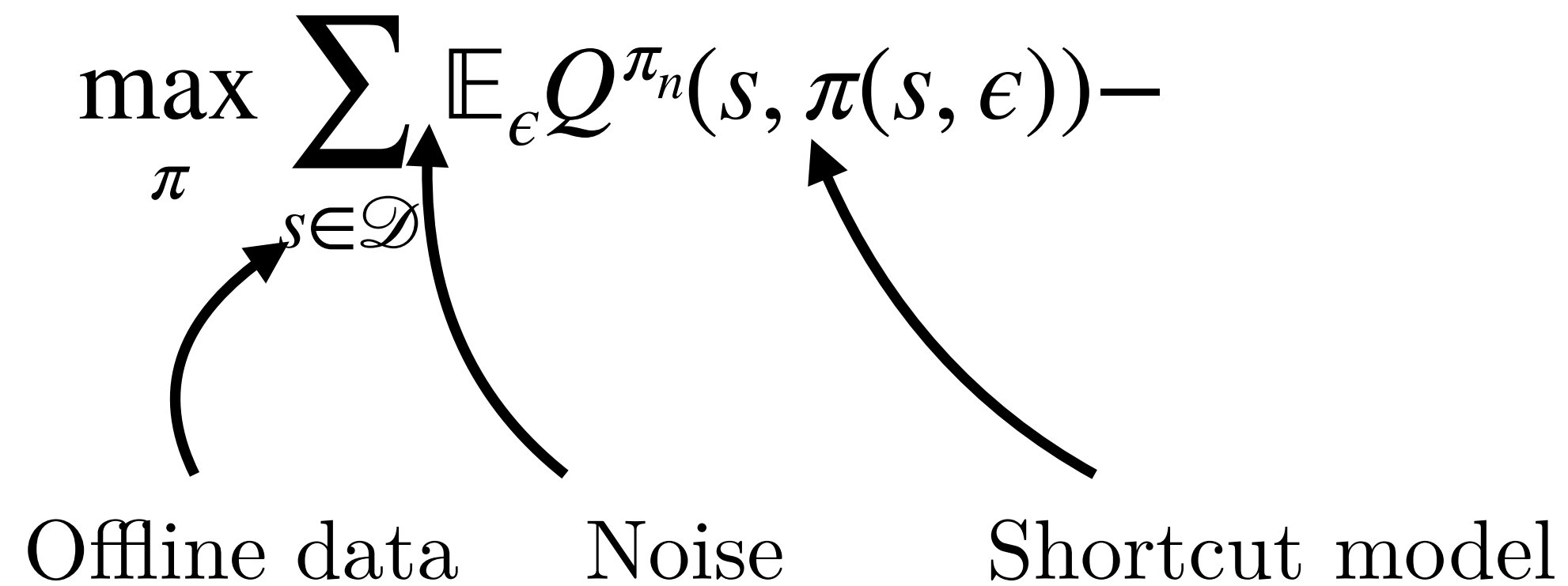
1. *Value Learning*: Train $Q^{\pi_n}$ with TD learning

# Scalable Offline Reinforcement Learning (SORL)

1. *Value Learning*: Train $Q^{\pi_n}$ with TD learning

2. *Policy Learning*: Train $\pi_{n+1}$ via the following

# Scalable Offline Reinforcement Learning (SORL)

1. *Value Learning*: Train $Q^{\pi_n}$ with TD learning

2. *Policy Learning*: Train $\pi_{n+1}$ via the following

$$\max_\pi \sum_{s \in \mathscr{D}} \mathbb{E}_\epsilon Q^{\pi_n}(s, \pi(s, \epsilon)) -$$

Offline data     Noise     Shortcut model

# Scalable Offline Reinforcement Learning (SORL)

1. *Value Learning*: Train $Q^{\pi_n}$ with TD learning

2. *Policy Learning*: Train $\pi_{n+1}$ via the following

$$\max_{\pi} \sum_{s \in \mathscr{D}} \mathbb{E}_{\epsilon} Q^{\pi_n}(s, \pi(s, \epsilon)) - \alpha_1 \, \text{FlowMatching}(\pi, \pi_B)$$

Offline data     Noise     Shortcut model     Regularization to offline data
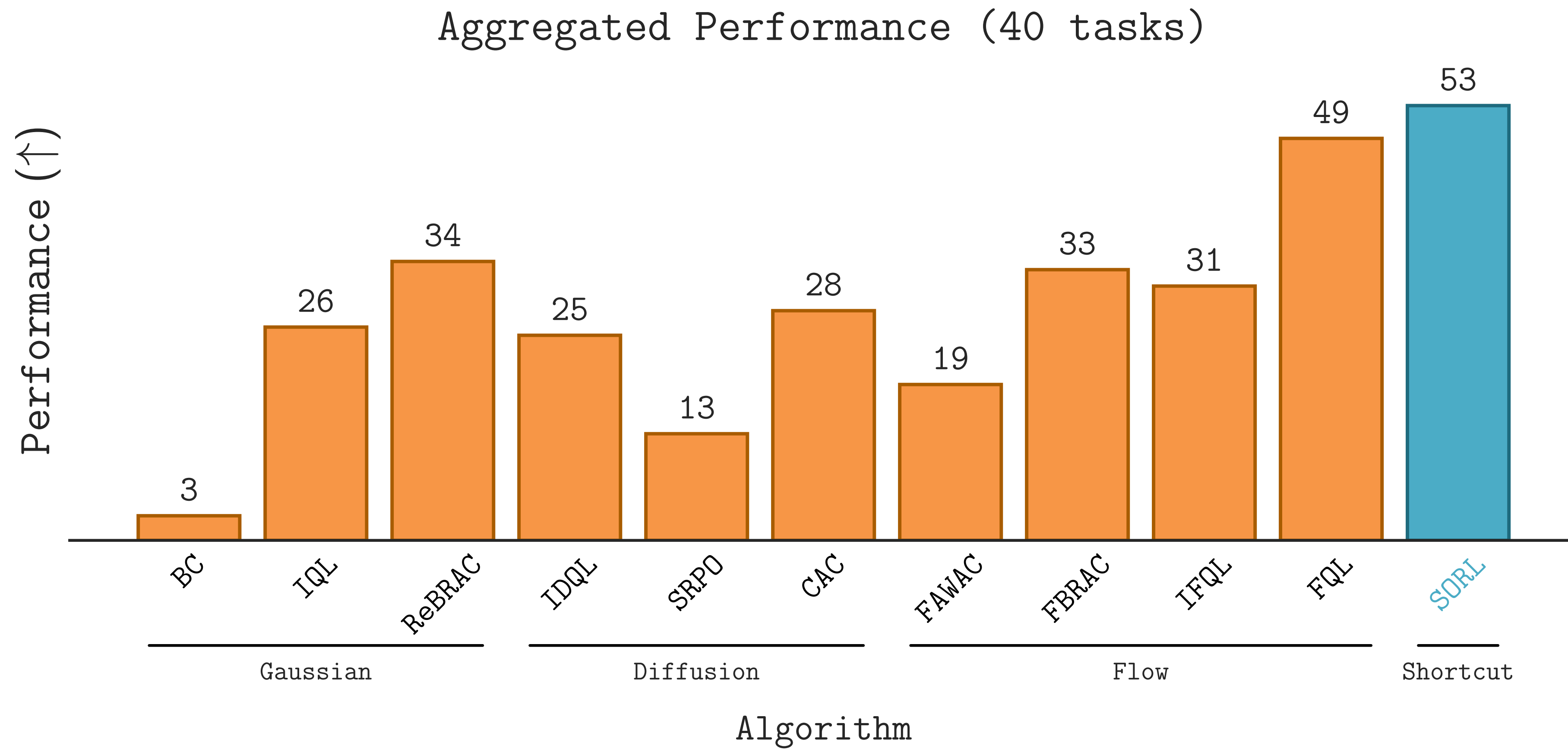
# Scalable Offline Reinforcement Learning (SORL)

1. *Value Learning*: Train $Q^{\pi_n}$ with TD learning

2. *Policy Learning*: Train $\pi_{n+1}$ via the following

$$\max_{\pi} \sum_{s \in \mathscr{D}} \mathbb{E}_{\epsilon} Q^{\pi_n}(s, \pi(s, \epsilon)) - \alpha_1 \text{FlowMatching}(\pi, \pi_B) - \alpha_2 \text{SelfConsistency}(\pi)$$

Offline data    Noise    Shortcut model    Regularization to offline data    Consistency between shortcuts and ODE

# SORL outperforms 10 baselines across 40 OGBench tasks



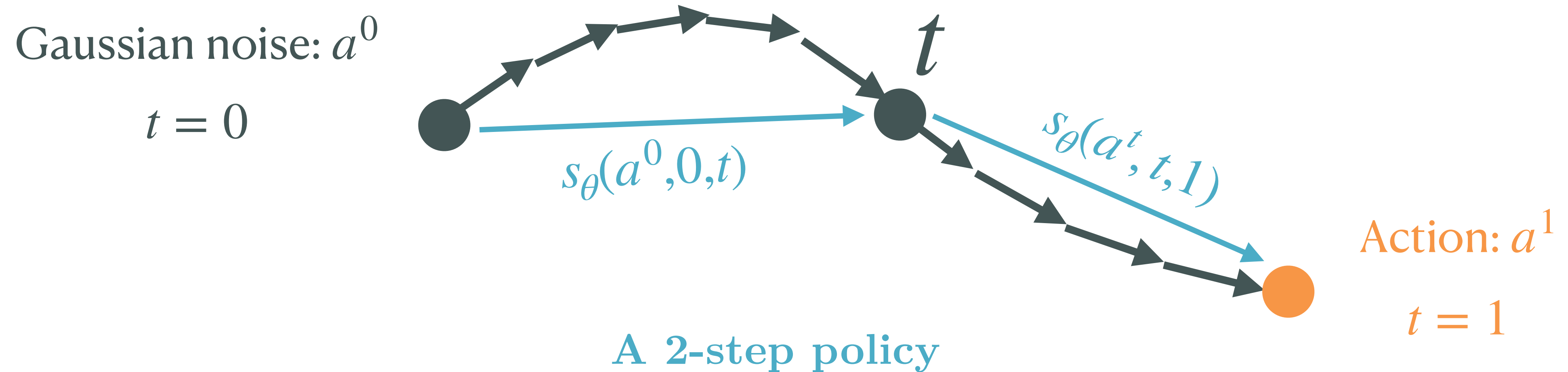Aggregated Performance (40 tasks)

# Axes of Scale in Offline RL

✅ Data

✅ Models

? Compute

# SORL enables test-time *sequential scaling*

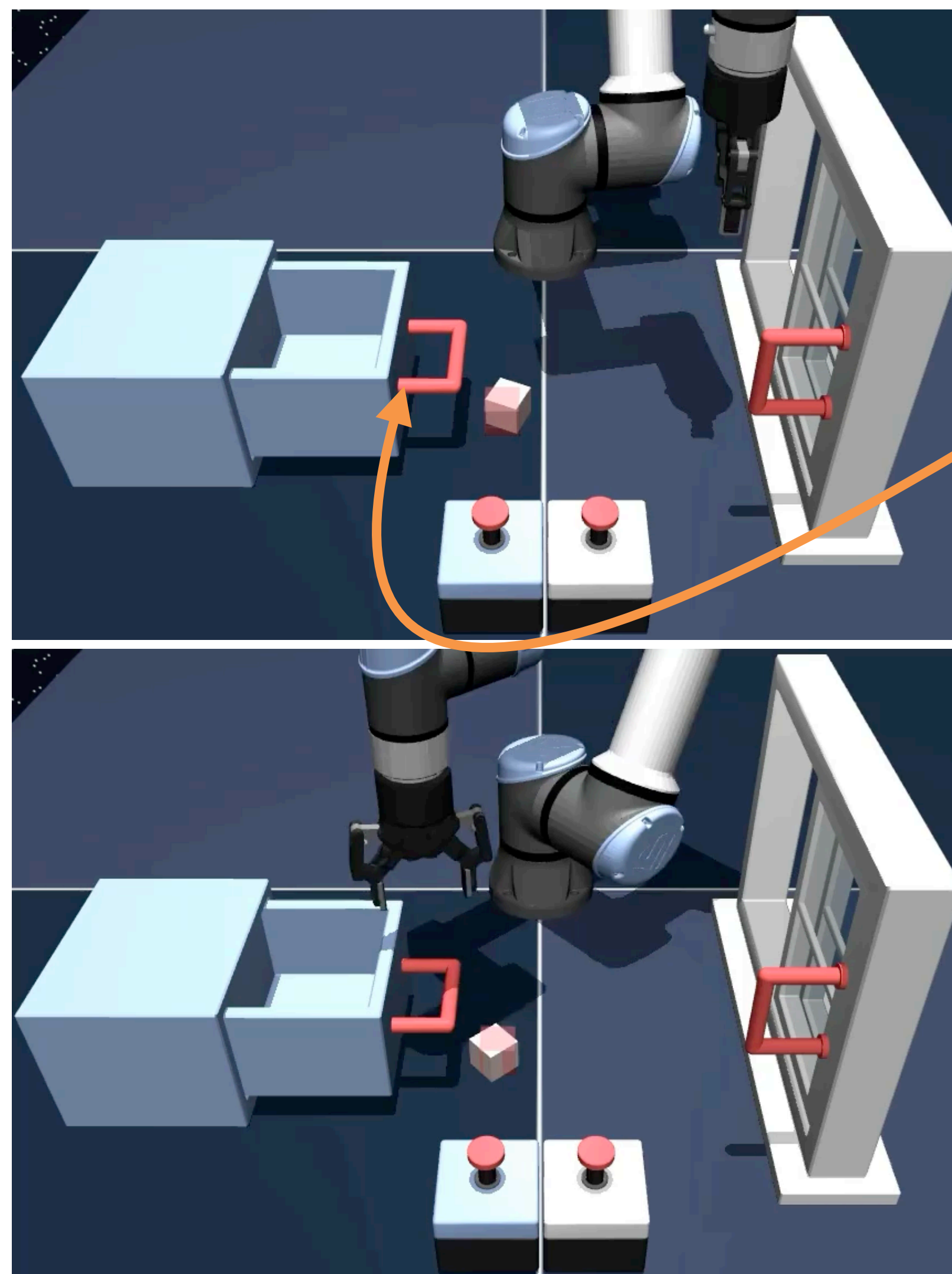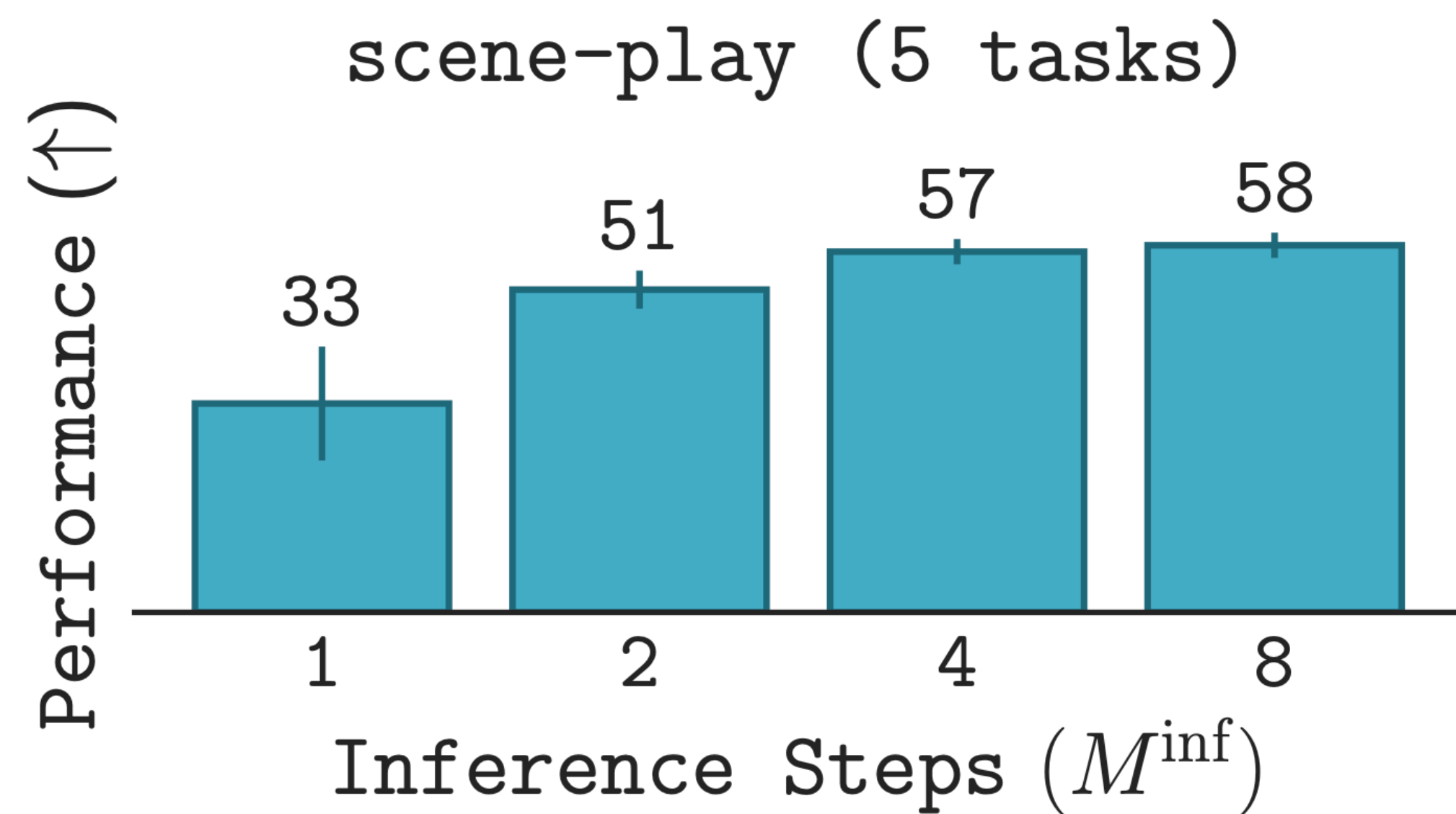At test-time, leverage more computation by using more "jumps"



Gaussian noise: $a^0$

$t = 0$

$t$

$s_\theta(a^0, 0, t)$

$s_\theta(a^t, t, 1)$

Action: $a^1$

$t = 1$

**A 2-step policy**

# SORL enables test-time *sequential scaling*

At test-time, leverage more computation by using more "jumps"



Gaussian noise: $a^0$

$t = 0$

$s_\theta(a^0, 0, t)$

$t$

$s_\theta(a^t, t, t')$

$t'$

$s_\theta(a^{t'}, t', 1)$

Action: $a^1$

$t = 1$

**A 3-step policy**

**Policies with more steps can model richer distributions, enabling better optimization of reward**

# Q: Is sequential scaling helpful?



scene-play (5 tasks)

Performance ($\uparrow$)

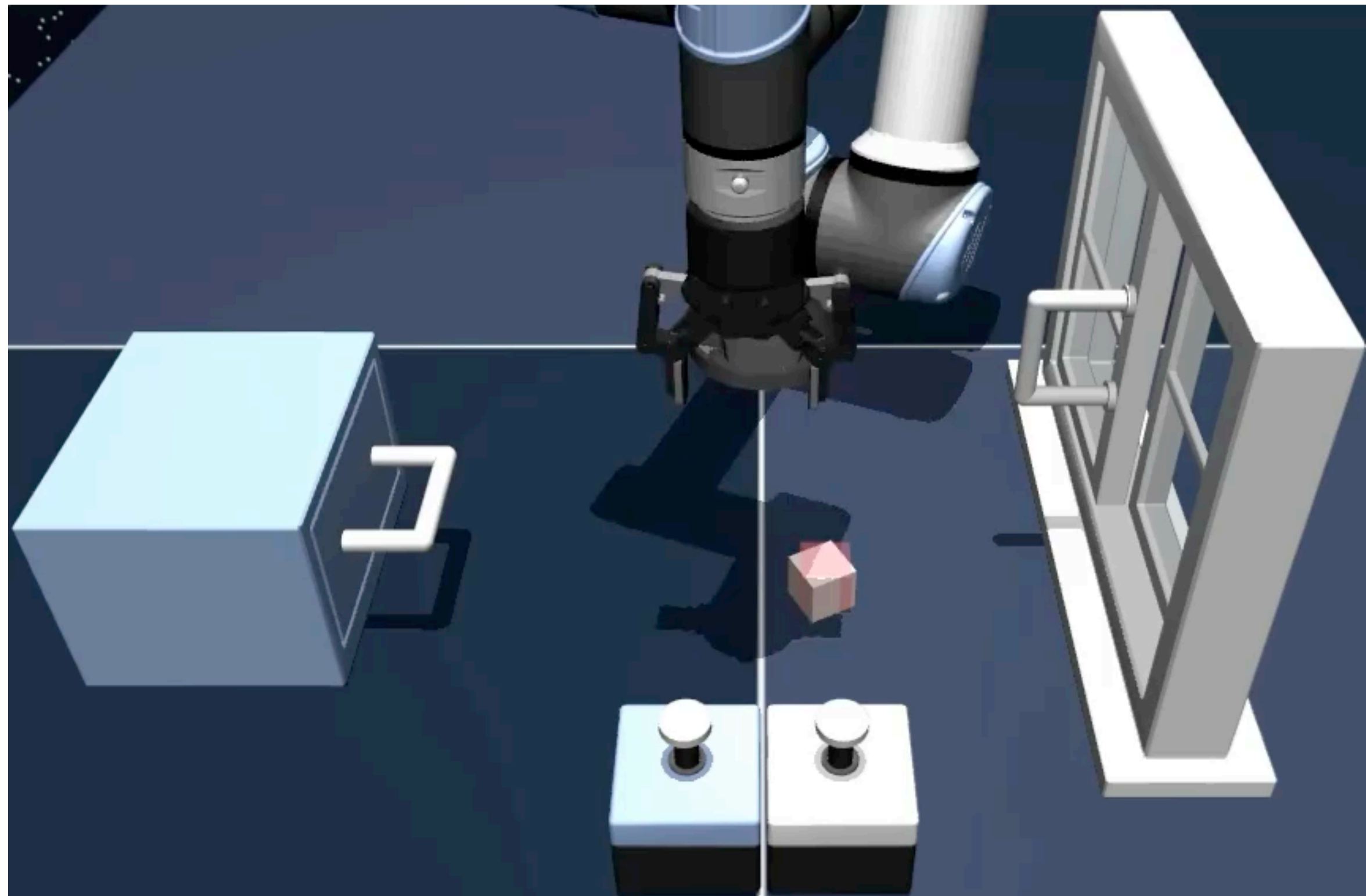33    51    57    58

Inference Steps ($M^{\mathrm{inf}}$)
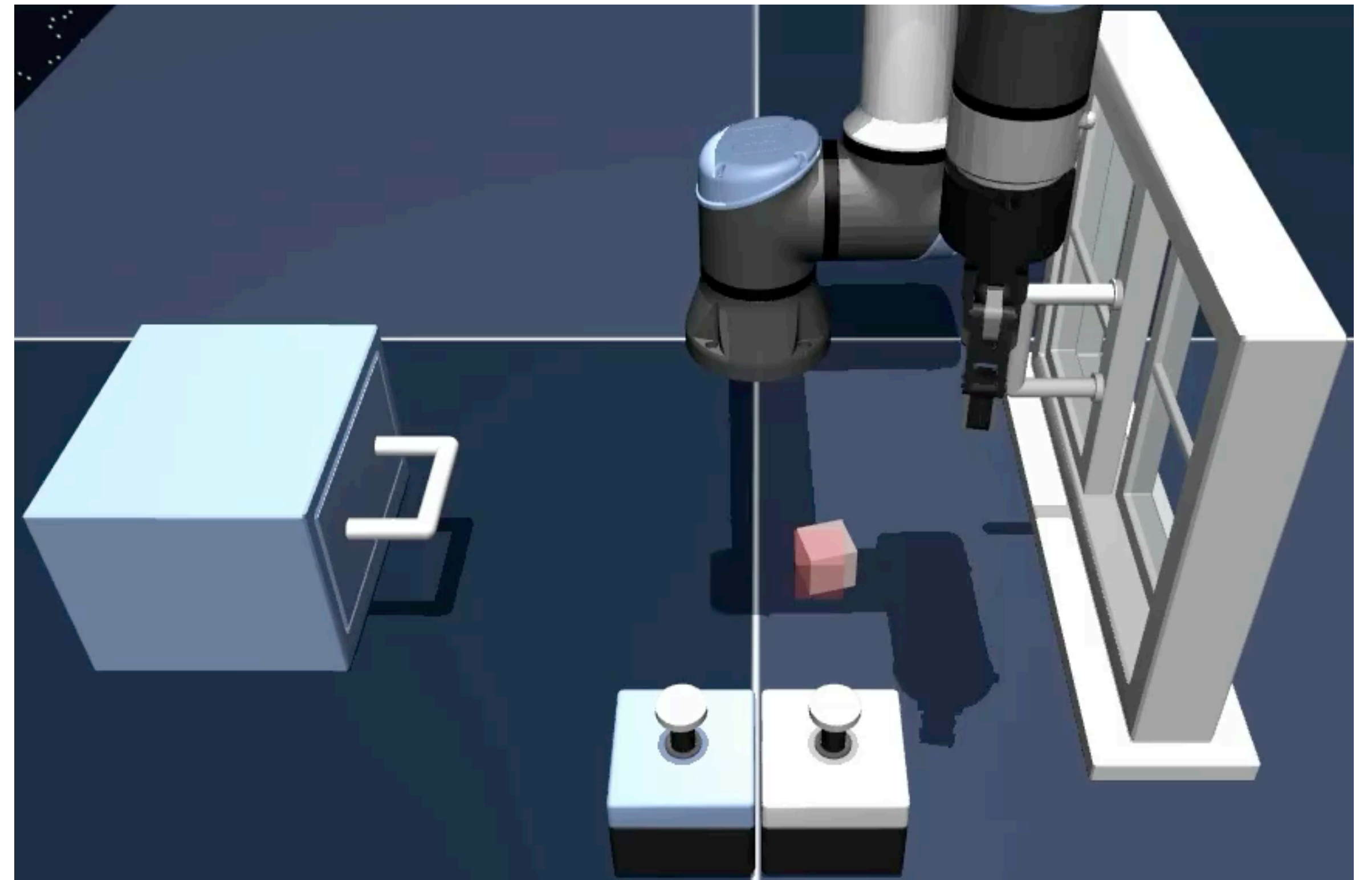
*1-step policy* struggles to close door

*8-step policy* solves task smoothly

# Q: Is sequential scaling helpful?

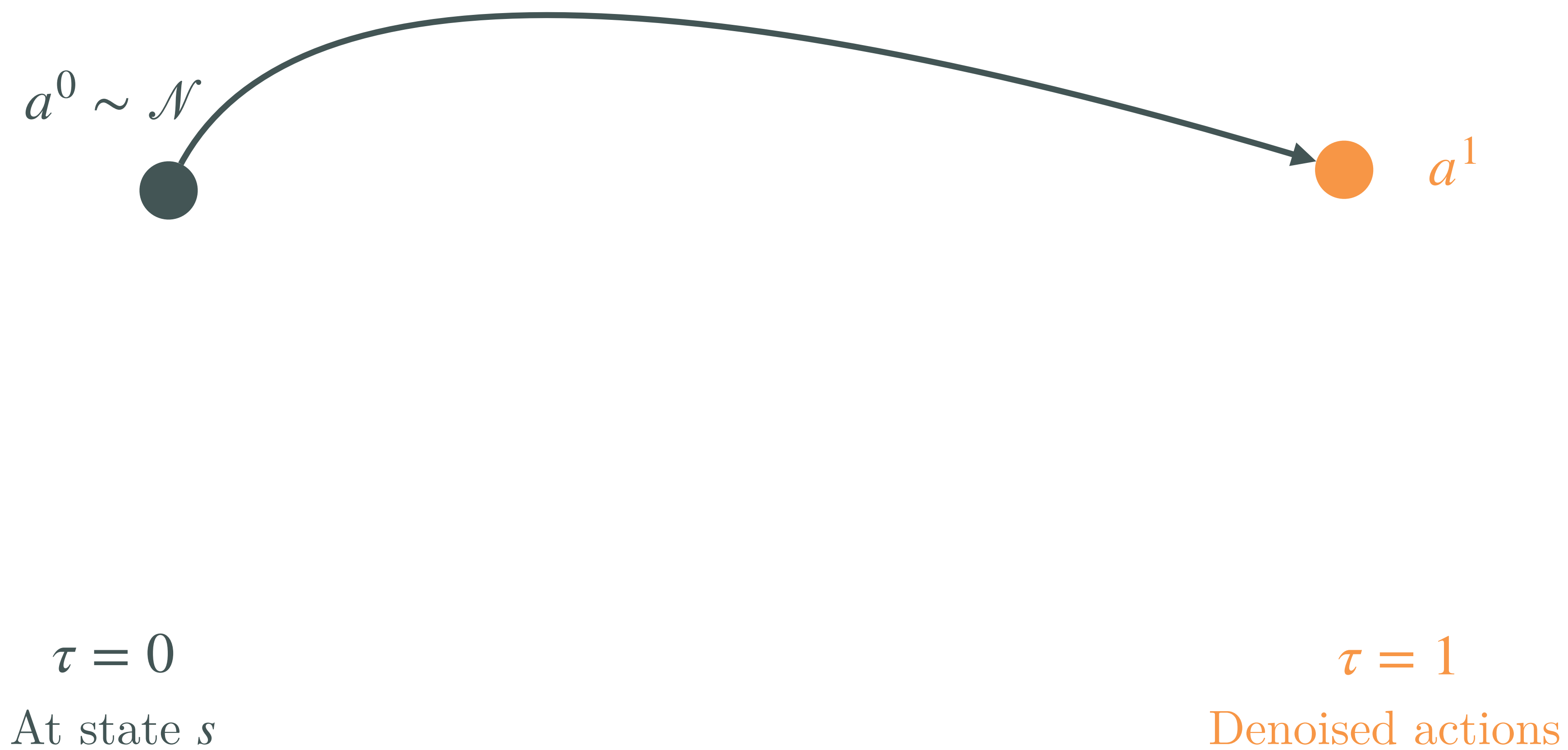Qualitatively, more inference steps leads to more precise actions



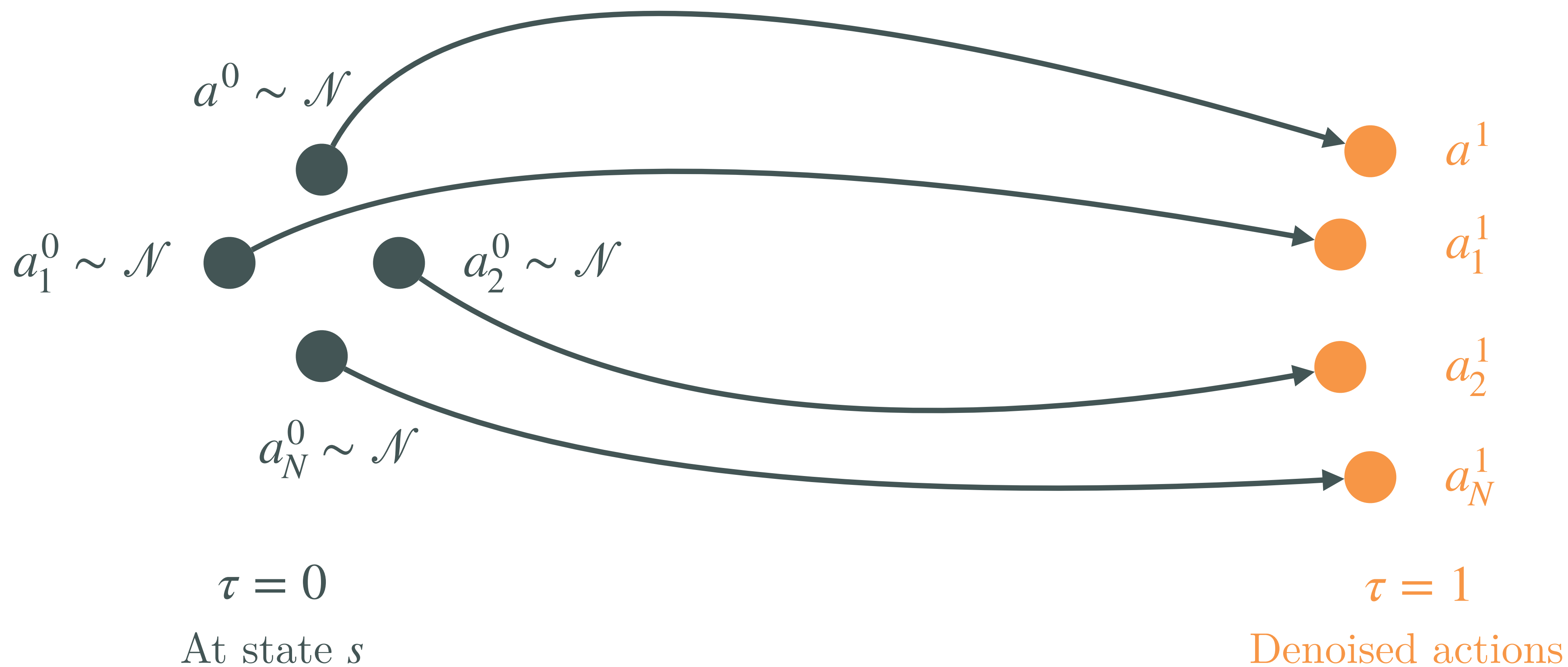1-step policy *pushes* handle

8-step policy *grips* handle

# SORL enables test-time *parallel scaling*

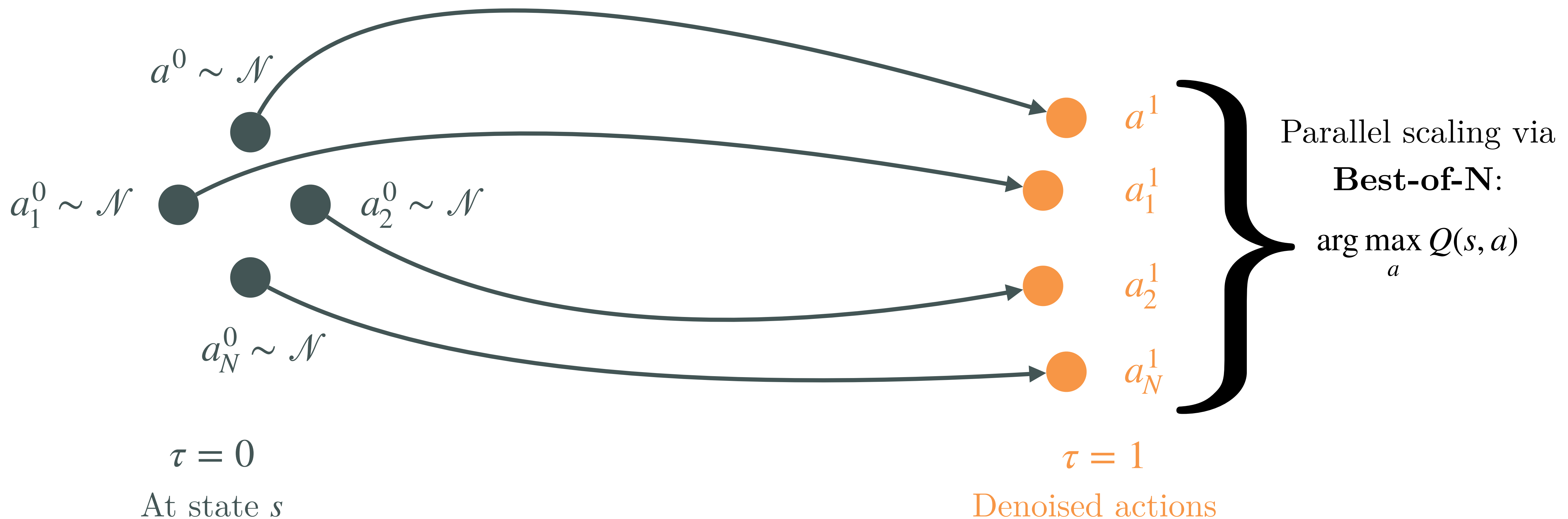# SORL enables test-time *parallel scaling*



$a^0 \sim \mathcal{N}$

$a^1$

$\tau = 0$

At state $s$

$\tau = 1$

Denoised actions

# SORL enables test-time *parallel scaling*

Search over multiple action candidates



$a^0 \sim \mathcal{N}$

$a_1^0 \sim \mathcal{N}$

$a_2^0 \sim \mathcal{N}$

$a_N^0 \sim \mathcal{N}$

$a^1$

$a_1^1$

$a_2^1$

$a_N^1$

$\tau = 0$

At state $s$

$\tau = 1$

Denoised actions

# Q: Is parallel scaling helpful?

Unlike LLMs, improvement from BoN is not obvious

scene-play (5 tasks)



Policy is trained with 2 steps,
but cannot sequentially
extrapolate to 4 steps

# Q: Is parallel scaling helpful?

Unlike LLMs, improvement from BoN is not obvious

**But BoN enables sequential extrapolation**



scene-play (5 tasks)

4-step + BoN enables sequential extrapolation

# Takeaways

Shortcut model is expressive while being suitable for RL training

Sequential and parallel scaling in general improves performance

In additional to scaling dataset, **leveraging rich generative models and scaling test-time computation** is a promising direction for RL

nico-espinosadice.github.io/projects/sorl

# Scaling Offline RL via
# Efficient and Expressive Shortcut Models

## Nicolas Espinosa Dice



*Joint work with*

Yiyi Zhang, Yiding Chen, Bradley Guo, Owen Oertell, Gokul Swamy, Kianté Brantley, Wen Sun