

Demonstration-Conditioned Reinforcement Learning for Few-Shot Imitation



Théo Cachet



Julien Perez



Christopher Dance

Given a few demonstrations of a new, previously unseen task

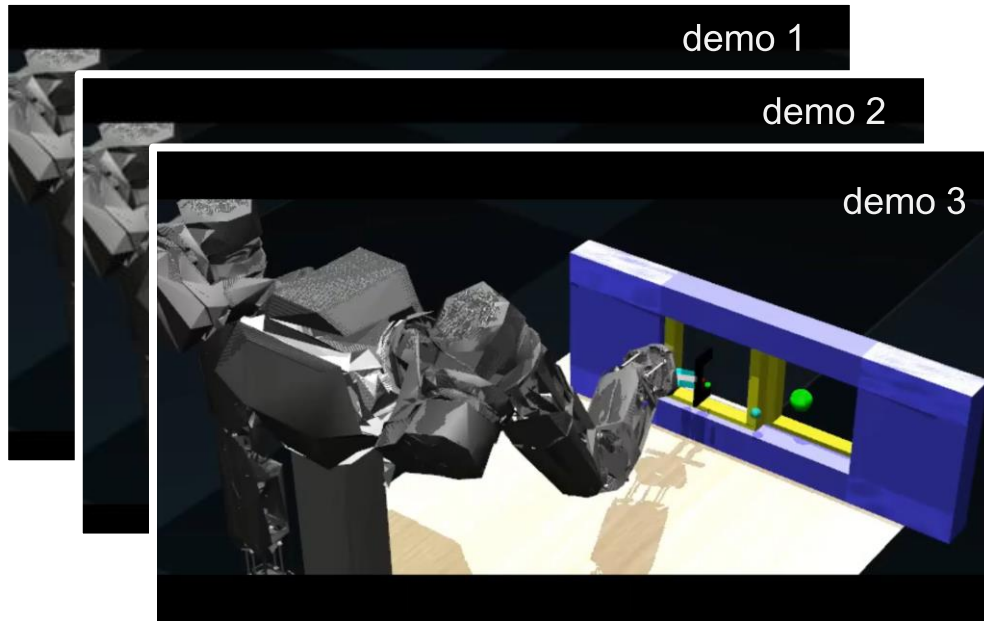


Find a policy which performs that task effectively.

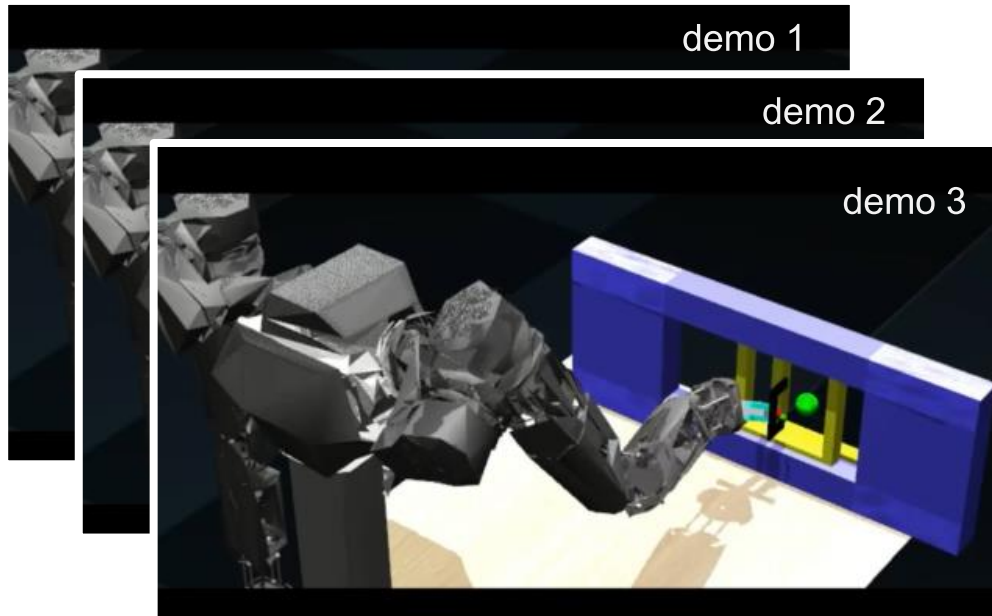
Given a few demonstrations of a new, previously unseen task



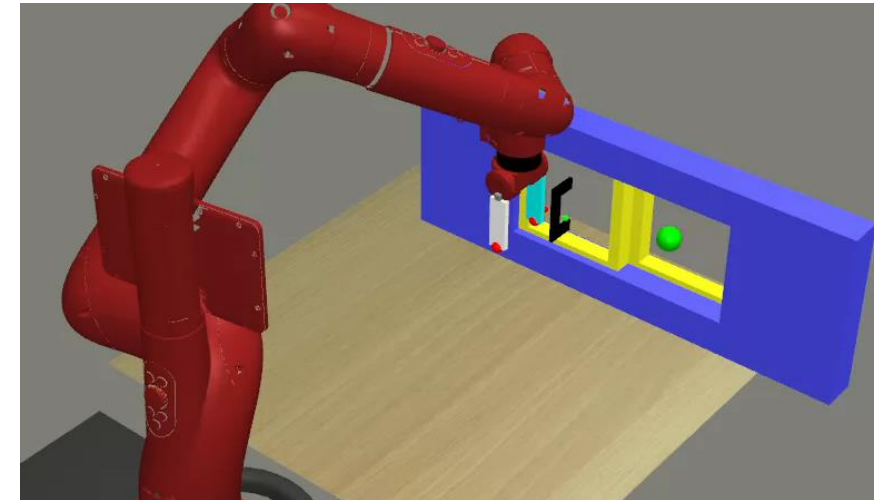
Find a policy which performs that task effectively.



Given a few demonstrations of a new, previously unseen task



Find a policy which performs that task effectively.

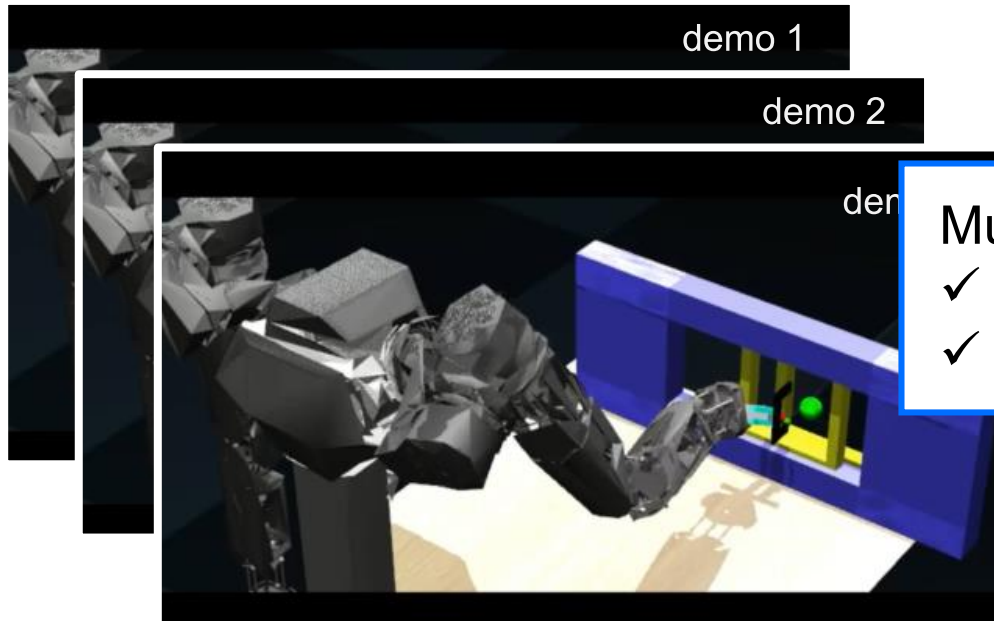


action $\sim \pi(\cdot \mid \text{history, demonstrations})$

Given a few demonstrations of a new, previously unseen task

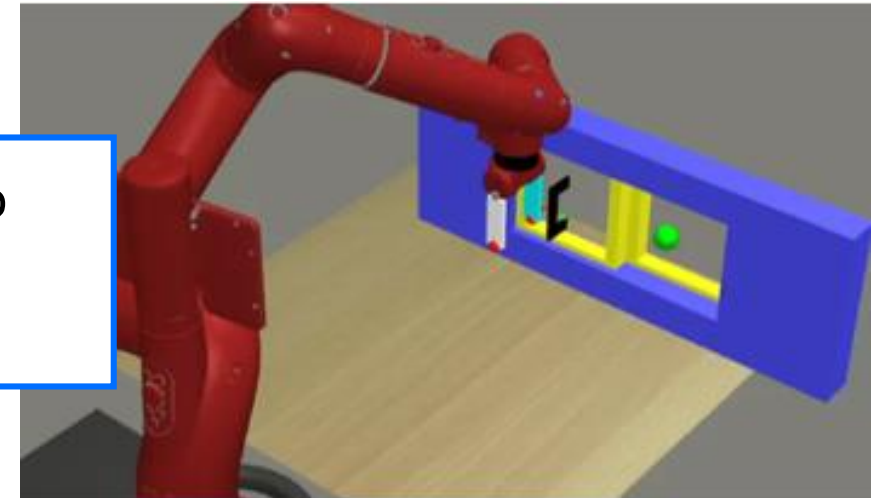


Find a policy which performs that task effectively.



Must **generalize** to

- ✓ new tasks
- ✓ new states

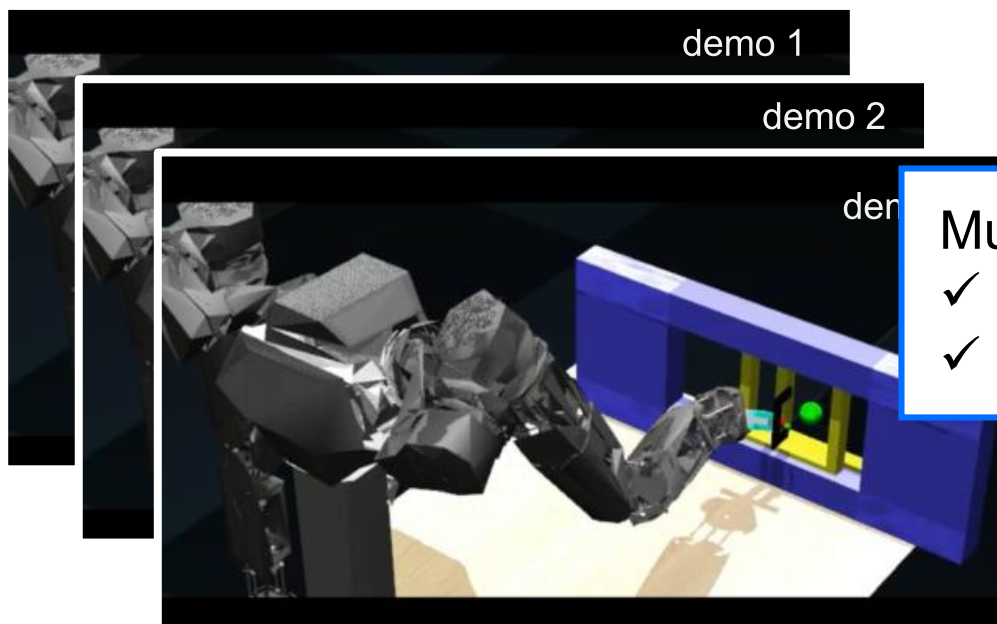


action $\sim \pi(\cdot \mid \text{history, demonstrations})$

Given a few demonstrations of a new, previously unseen task

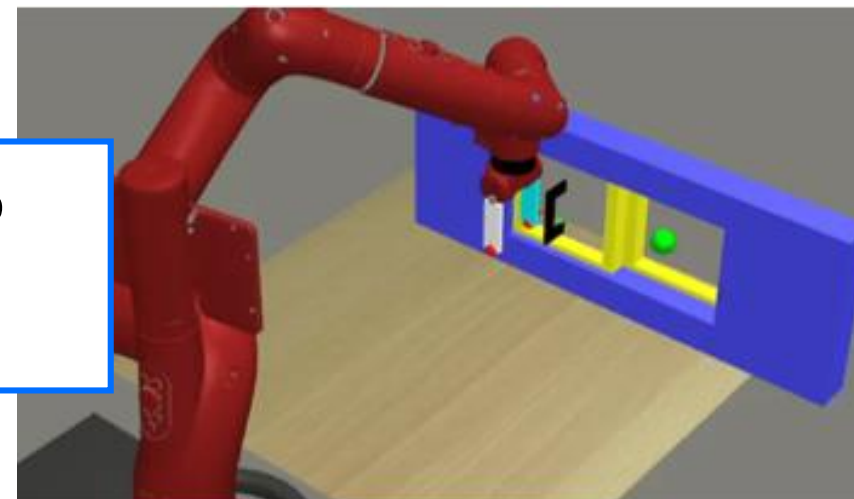


Find a policy which performs that task effectively.



Must **generalize** to

- ✓ new tasks
- ✓ new states



action $\sim \pi(\cdot \mid \text{history, demonstrations})$

Demonstration is flexibly defined:

- ✓ noisy, incomplete, sub-optimal
- ✓ no actions
- ✓ human demonstrator + robot agent

Ingredients η Distribution over tasks $\mu \sim \eta$

Ingredients

 η

Distribution over tasks $\mu \sim \eta$

 D_μ

Distribution over collections \mathbf{d} of demonstrations of task μ

Ingredients

 η

Distribution over tasks $\mu \sim \eta$

 D_μ

Distribution over collections \mathbf{d} of demonstrations of task μ

 $\pi(\cdot | h, \mathbf{d})$

Demonstration-conditioned policy given history h and demonstrations \mathbf{d}

Ingredients

 η Distribution over tasks $\mu \sim \eta$ D_μ Distribution over collections \mathbf{d} of demonstrations of task μ $\pi(\cdot | h, \mathbf{d})$ *Demonstration-conditioned policy* given history h and demonstrations \mathbf{d}

Objective

$$\max_{\pi} \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{\mathbf{d} \sim D_\mu} J_\mu(\pi(\cdot | h, \mathbf{d}))$$

↑
↑
↑
↑

demonstration-
conditioned policy
tasks
collections of
demonstrations
return

Ingredients

 η Distribution over tasks $\mu \sim \eta$ D_μ Distribution over collections \mathbf{d} of demonstrations of task μ $\pi(\cdot | h, \mathbf{d})$ *Demonstration-conditioned policy* given history h and demonstrations \mathbf{d}

Each task μ is an MDP
 $J_\mu(\pi)$ is the return for policy π

Objective

$$\max_{\pi} \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{\mathbf{d} \sim D_\mu} J_\mu(\pi(\cdot | h, \mathbf{d}))$$

↑
↑
↑
↑

demonstration-
conditioned policy
tasks
collections of
demonstrations
return

Ingredients

η

Distribution over tasks $\mu \sim \eta$

Each task μ is an MDP
 $J_\mu(\pi)$ is the return for policy π

D_μ

Distribution over collections \mathbf{d} of demonstrations of task μ

few \leftrightarrow 1 to 10

$\pi(\cdot | h, \mathbf{d})$

Demonstration-conditioned policy given history h and demonstrations \mathbf{d}

Objective

$$\max_{\pi} \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{\mathbf{d} \sim D_\mu} J_\mu(\pi(\cdot | h, \mathbf{d}))$$

↑
↑
↑
↑

demonstration-conditioned policy
tasks
collections of demonstrations
return

Ingredients

η

Distribution over tasks $\mu \sim \eta$

D_μ

Distribution over collections \mathbf{d} of demonstrations of task μ

few \leftrightarrow 1 to 10

$\pi(\cdot | h, \mathbf{d})$

Demonstration-conditioned policy given history h and demonstrations \mathbf{d}

Each task μ is an MDP
 $J_\mu(\pi)$ is the return for policy π

POMDPs

Objective

$$\max_{\pi} \mathbb{E}_{\mu \sim \eta} \mathbb{E}_{\mathbf{d} \sim D_\mu} J_\mu(\pi(\cdot | h, \mathbf{d}))$$

↑ demonstration-conditioned policy
 ↑ tasks
 ↑ collections of demonstrations
 ↑ return

DCRL maximizes the return of a demonstration-conditioned policy, averaged over a set of training tasks and corresponding demonstrations.

μ^0, \dots, μ^{N-1} may not be distinct

Train

Input Pairs $(\mathbf{d}^0, \mu^0), \dots, (\mathbf{d}^{N-1}, \mu^{N-1})$ where \mathbf{d}^i is a collection of demonstrations of task μ^i

Output A demonstration-conditioned policy π attaining $\max_{\pi} \sum_{i=0}^{N-1} J_{\mu^i}(\pi(\cdot | \cdot, \mathbf{d}^i))$

Test

Input $\left\{ \begin{array}{l} \mathbf{d} \text{ Collection of demonstrations of new, previously unseen task} \\ \pi \text{ Demonstration-conditioned policy given by DCRL} \end{array} \right.$

Repeat Observe history h_t and take action $a_t \sim \pi(\cdot | h_t, \mathbf{d})$

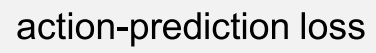
No need for reward function
No need for to explore the test env't

Behaviour cloning (BC)

Duan *et al.* '17

$$\min_{\pi} \sum_{i=0}^{N-1} \sum_{(s,a) \in \mathbf{d}^i} \ell(a, \pi(\cdot | s, \mathbf{d}^i))$$

action-prediction loss



Behaviour cloning (BC)

Duan *et al.* '17

$$\min_{\pi} \sum_{i=0}^{N-1} \sum_{(s,a) \in \mathbf{d}^i} \ell(a, \pi(\cdot | s, \mathbf{d}^i))$$

Meta-inverse RL

Yu *et al.* '19a, Goo and Niekum '19

Infer reward \hat{R}_{μ} from \mathbf{d}

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_t \gamma^t \hat{R}_{\mu}(s_t, a_t) \right]$$

Behaviour cloning (BC)Duan *et al.* '17

$$\min_{\pi} \sum_{i=0}^{N-1} \sum_{(s,a) \in \mathbf{d}^i} \ell(a, \pi(\cdot | s, \mathbf{d}^i))$$

Meta-inverse RLYu *et al.* '19a, Goo and Niekum '19Infer reward \hat{R}_{μ} from \mathbf{d}

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_t \gamma^t \hat{R}_{\mu}(s_t, a_t) \right]$$

Demonstration-conditioned RL

This paper

$$\max_{\pi} \sum_{i=0}^{N-1} J_{\mu^i}(\pi(\cdot | \cdot, \mathbf{d}^i))$$

Behaviour cloning (BC)Duan *et al.* '17

$$\min_{\pi} \sum_{i=0}^{N-1} \sum_{(s,a) \in \mathbf{d}^i} \ell(a, \pi(\cdot | s, \mathbf{d}^i))$$

Needs actions in demo's

No interaction with the test env't

Compounding errors \Rightarrow loss is $O(H^2)$
on horizon H (Rajaraman *et al.*, 2020)**Meta-inverse RL**Yu *et al.* '19a, Goo and Niekum '19Infer reward \hat{R}_{μ} from \mathbf{d}

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_t \gamma^t \hat{R}_{\mu}(s_t, a_t) \right]$$

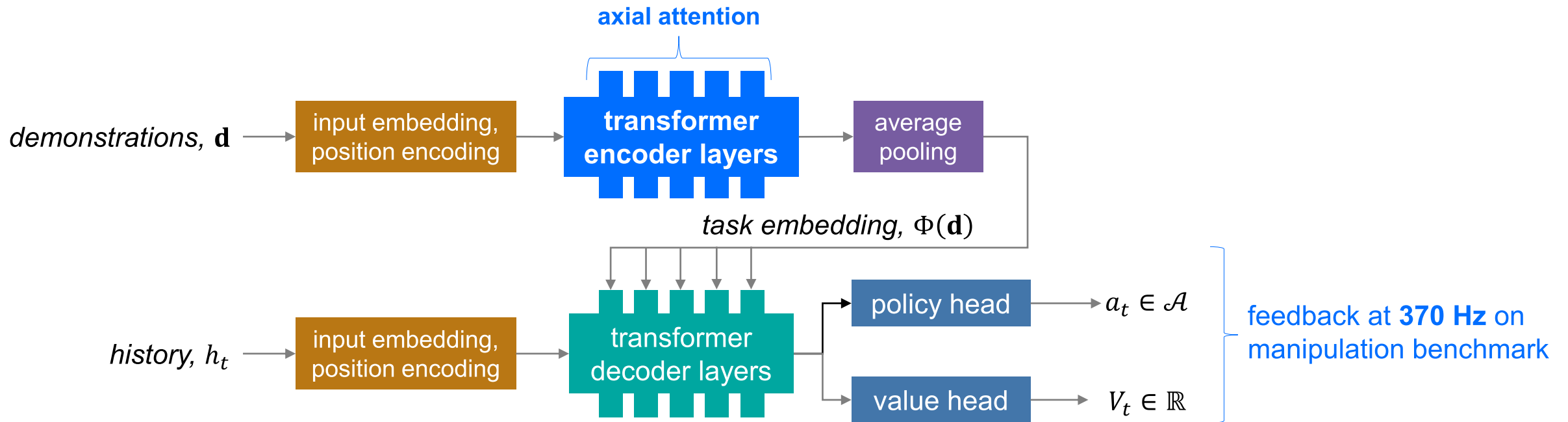
No need for actions in demo's

Must interact with the test env't**Demonstration-conditioned RL**

This paper

$$\max_{\pi} \sum_{i=0}^{N-1} J_{\mu^i}(\pi(\cdot | \cdot, \mathbf{d}^i))$$

No need for actions in demo's**No interaction with the test env't****Improves on suboptimal demo's****Copes with demonstrator domain shift****Needs reward function for training**



Literature

Attention and transformers already in use for few-shot imitation (Duan *et al.* '17, Mishra '18, James '18, Dasari '20)

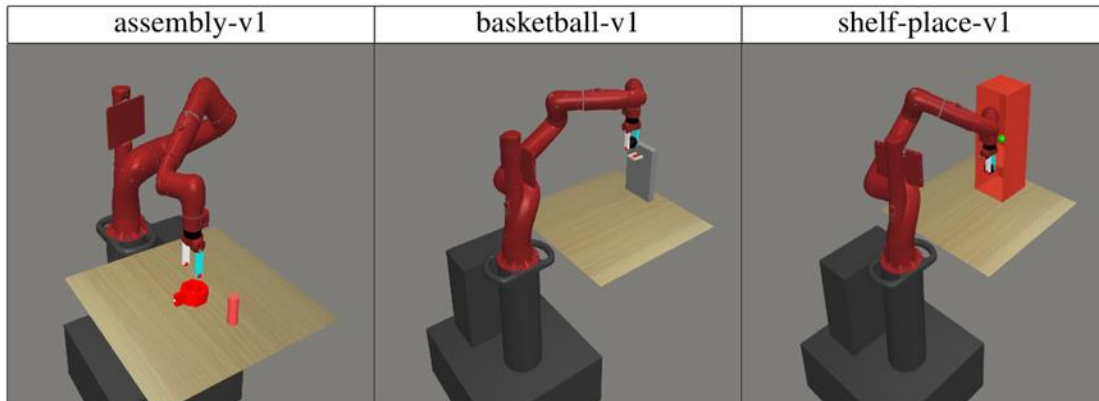
Novelty

- Cross-demonstration attention.** Process multiple demonstrations *jointly*.
- Axial attention.** Attend to one dimension of the input at a time (Ho *et al.* '18).
Reduces time and memory from $O(T^2n^2)$ to $O(Tn(T + n))$ for n time series of length T

1. Meta-World Benchmark (Yu et al., '19)

50 manipulation tasks

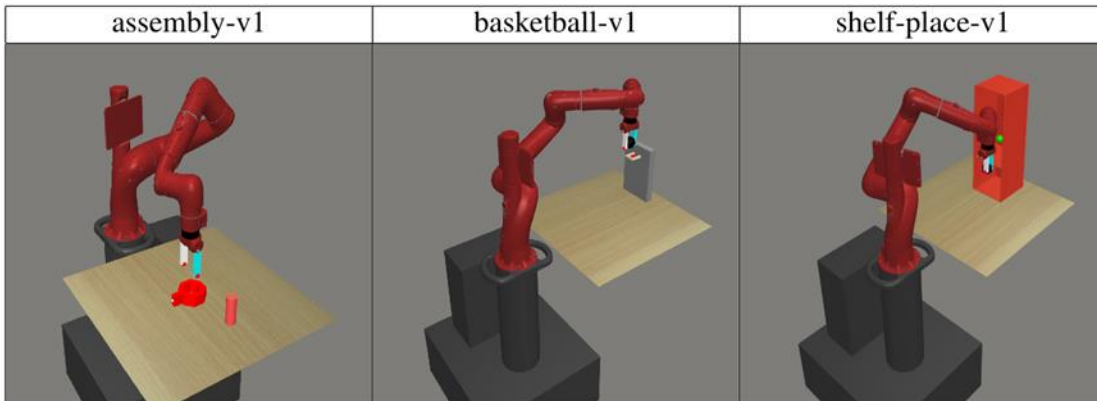
e.g. open-window, lock-door



1. Meta-World Benchmark (Yu et al., '19)

50 manipulation tasks

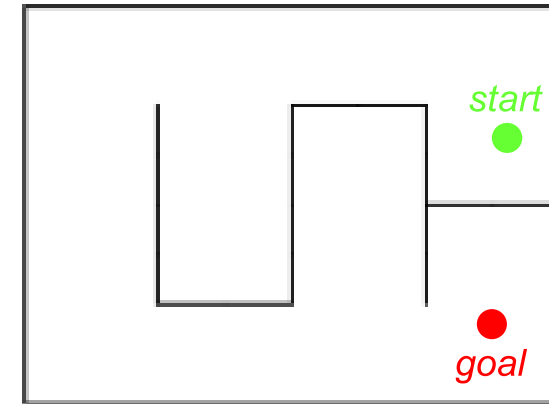
e.g. open-window, lock-door



2. Navigation Benchmark

60 mazes ↔ tasks

map of test maze



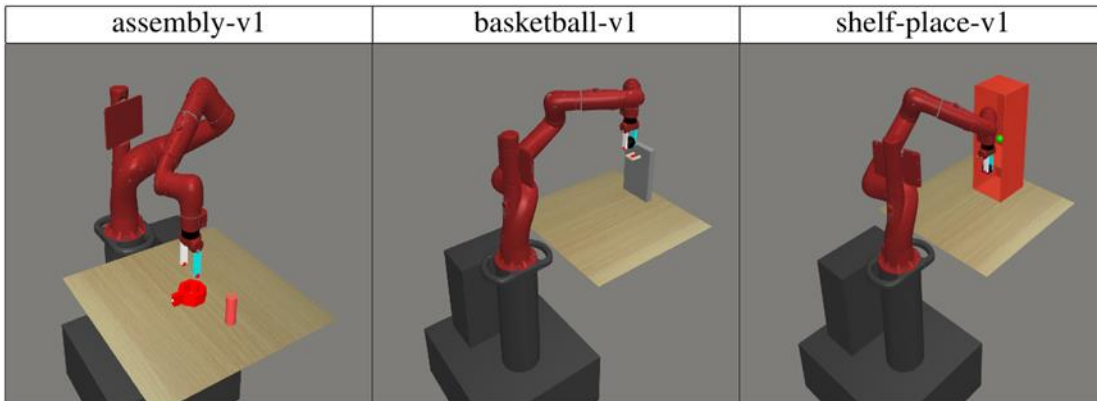
Aim

Get from start to goal state (randomized).

1. Meta-World Benchmark (Yu et al., '19)

50 manipulation tasks

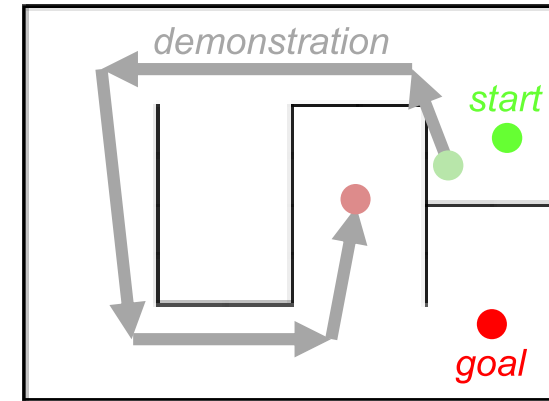
e.g. open-window, lock-door



2. Navigation Benchmark

60 mazes ↔ tasks

what the agent sees



Aim

Get from start to goal state (randomized).

Agent can't see the walls, but is penalized if it hits a wall.

So, it must infer the walls from the demonstrations.

Evaluation. All experiments we will now present had distinct training and test tasks.

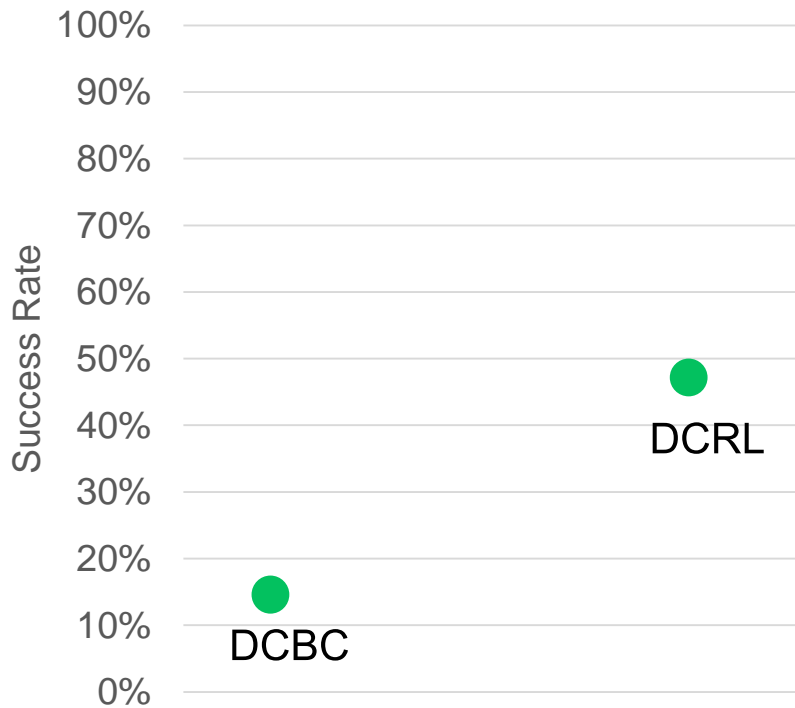
Definition (DCBC). *Demonstration-conditioned behavioural cloning (DCBC)* has the same architecture as DCRL, but is trained with a behaviour-cloning loss.

} *like Duan et al. (2016)*

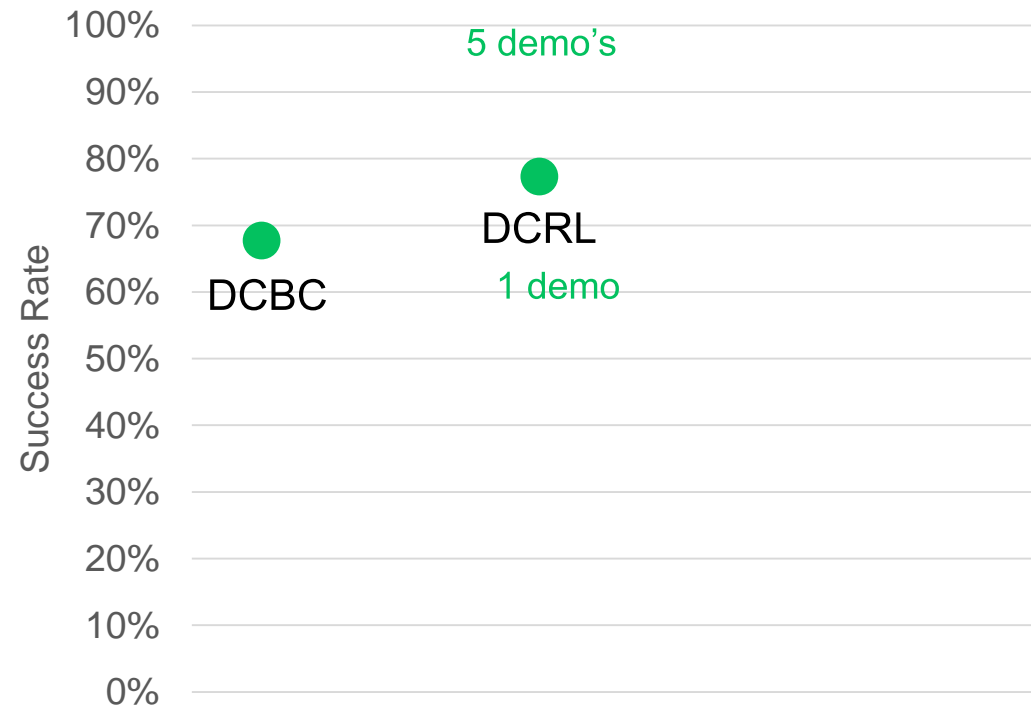
Definition (DCBC). *Demonstration-conditioned behavioural cloning (DCBC)* has the same architecture as DCRL, but is trained with a behaviour-cloning loss.

} like Duan et al. (2016)

Meta-World



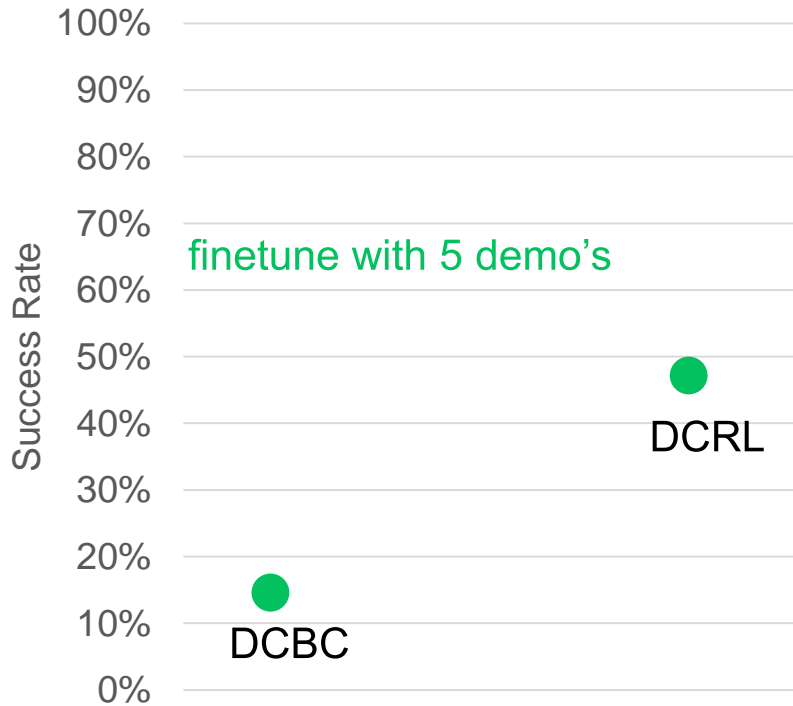
Navigation



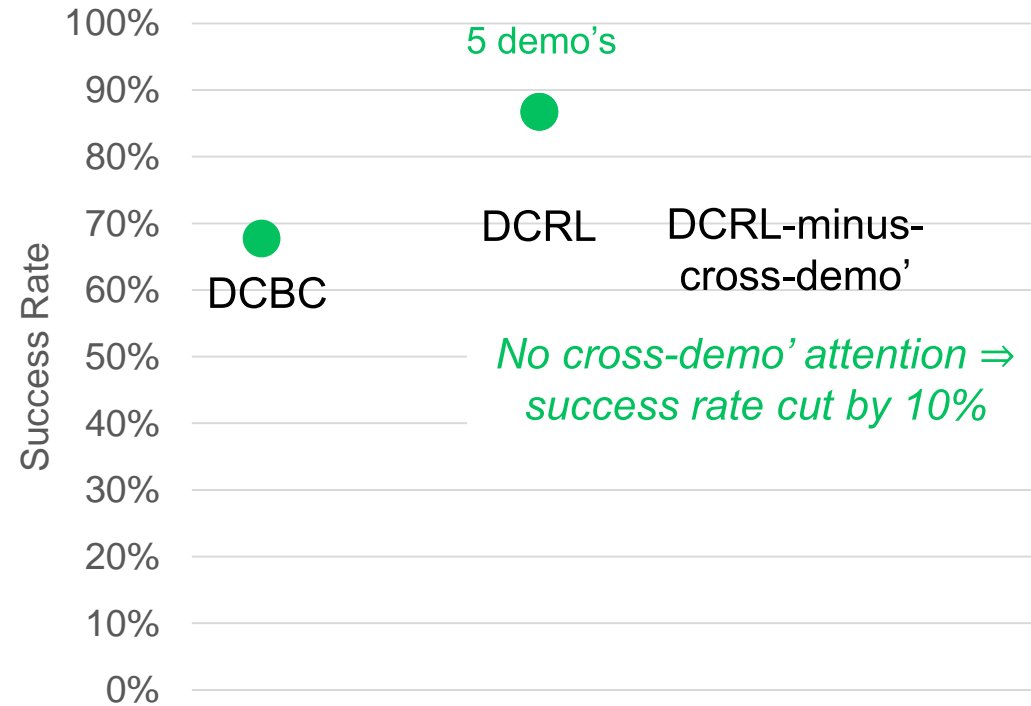
Definition (DCBC). *Demonstration-conditioned behavioural cloning (DCBC)* has the same architecture as DCRL, but is trained with a behaviour-cloning loss.

} like Duan et al. (2016)

Meta-World



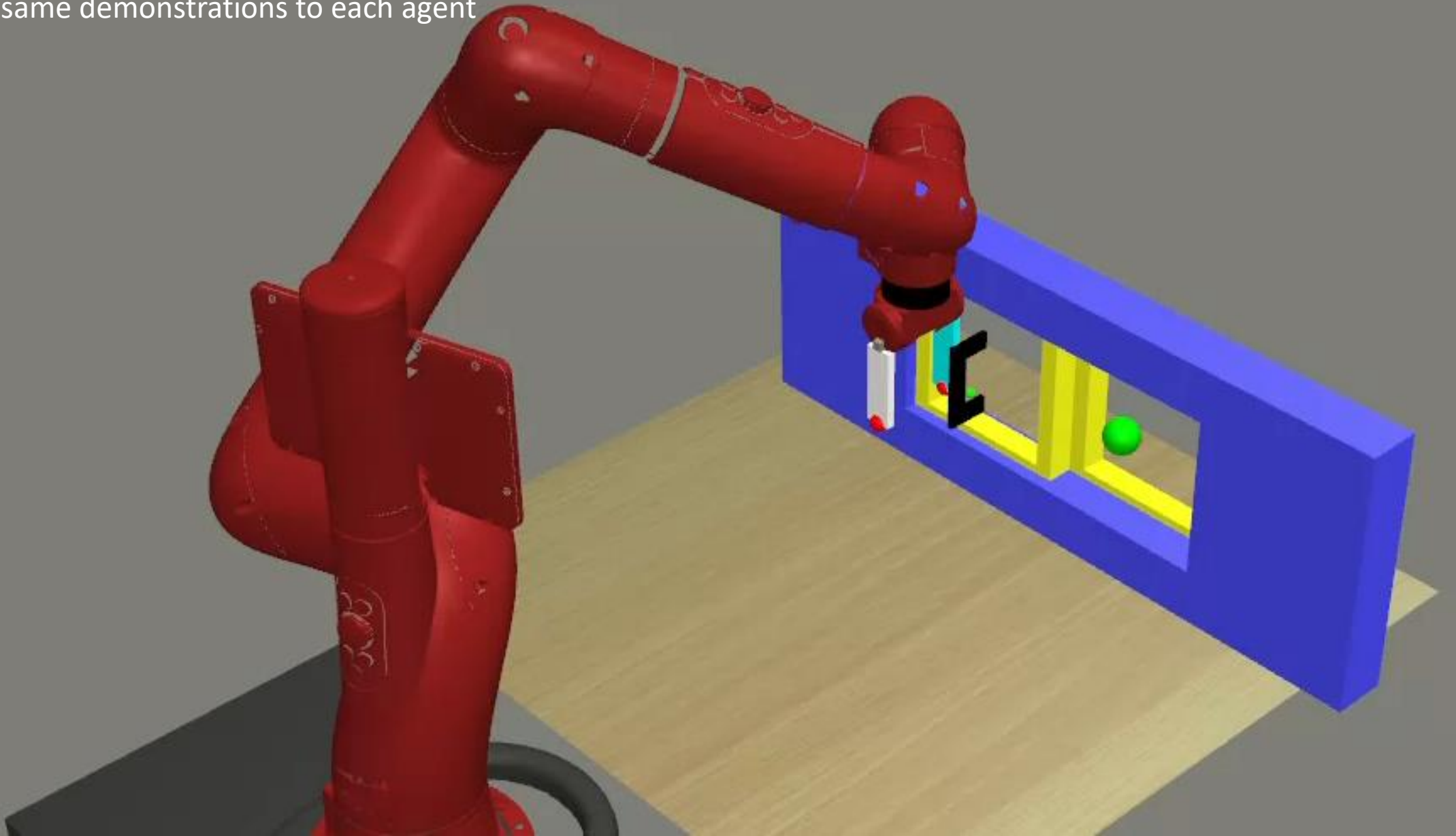
Navigation



Why does DCLR outperform DCBC?

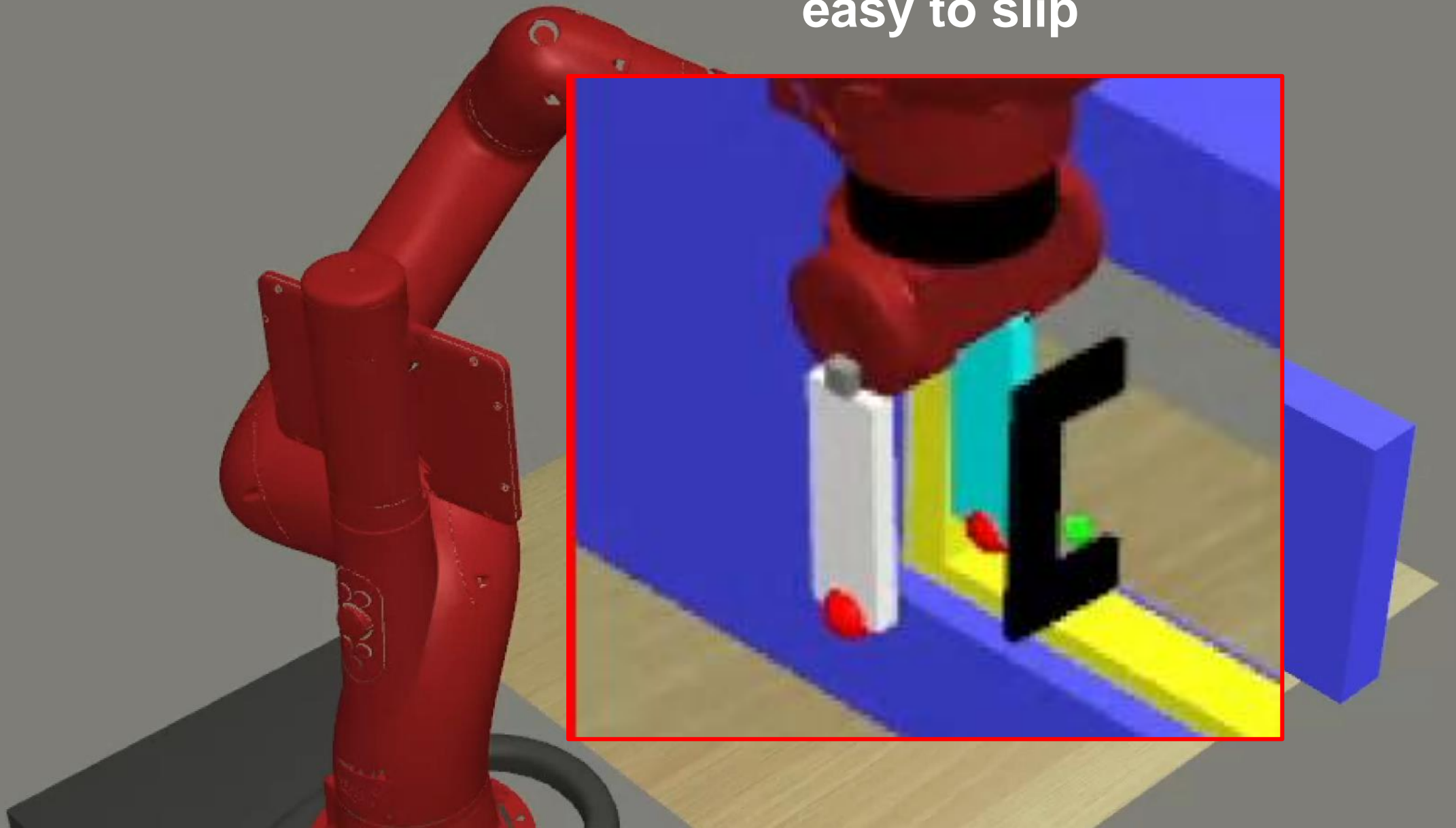
DEMONSTRATIONS

Present the same demonstrations to each agent



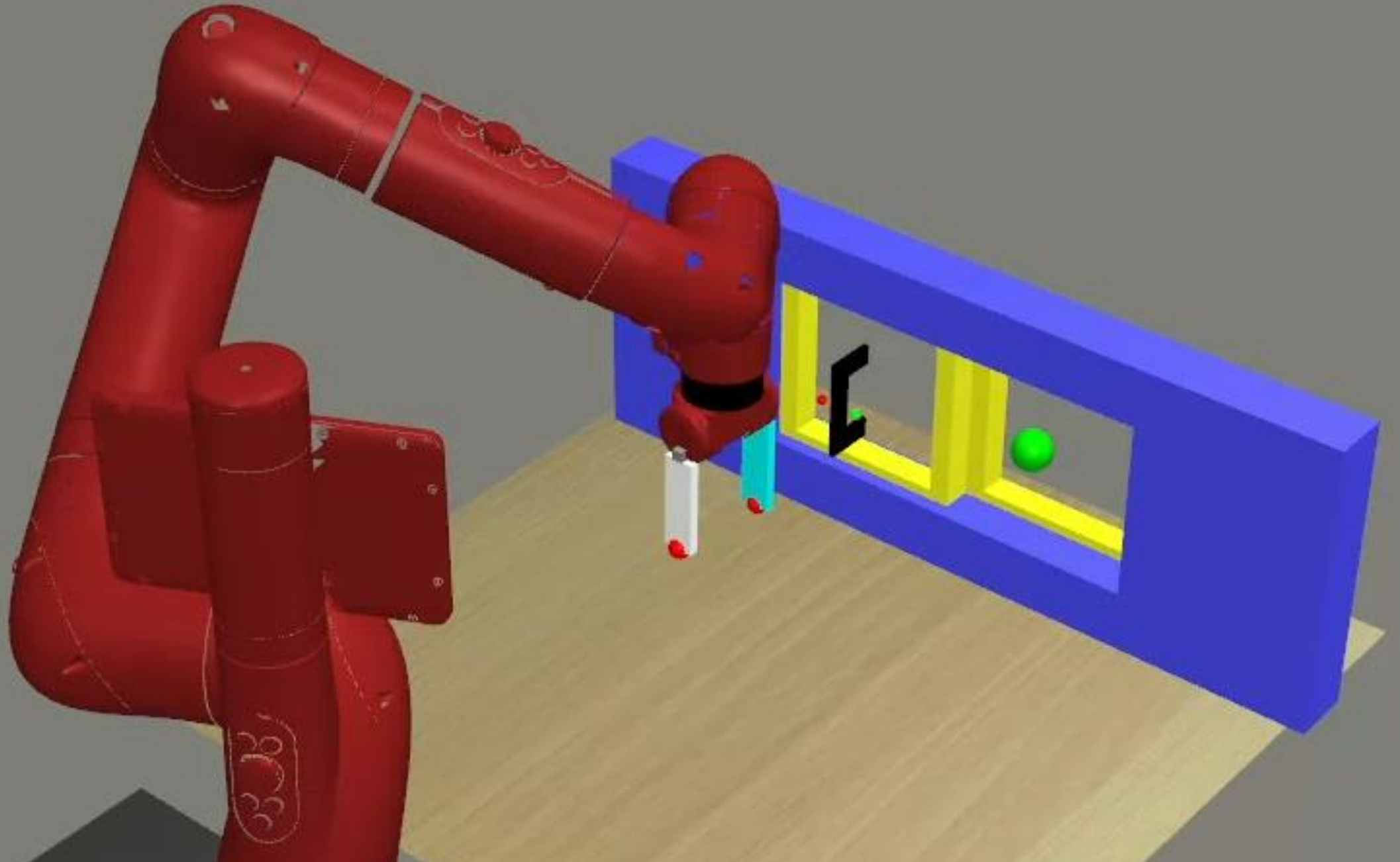
Why does DCLR outperform DCBC?

easy to slip



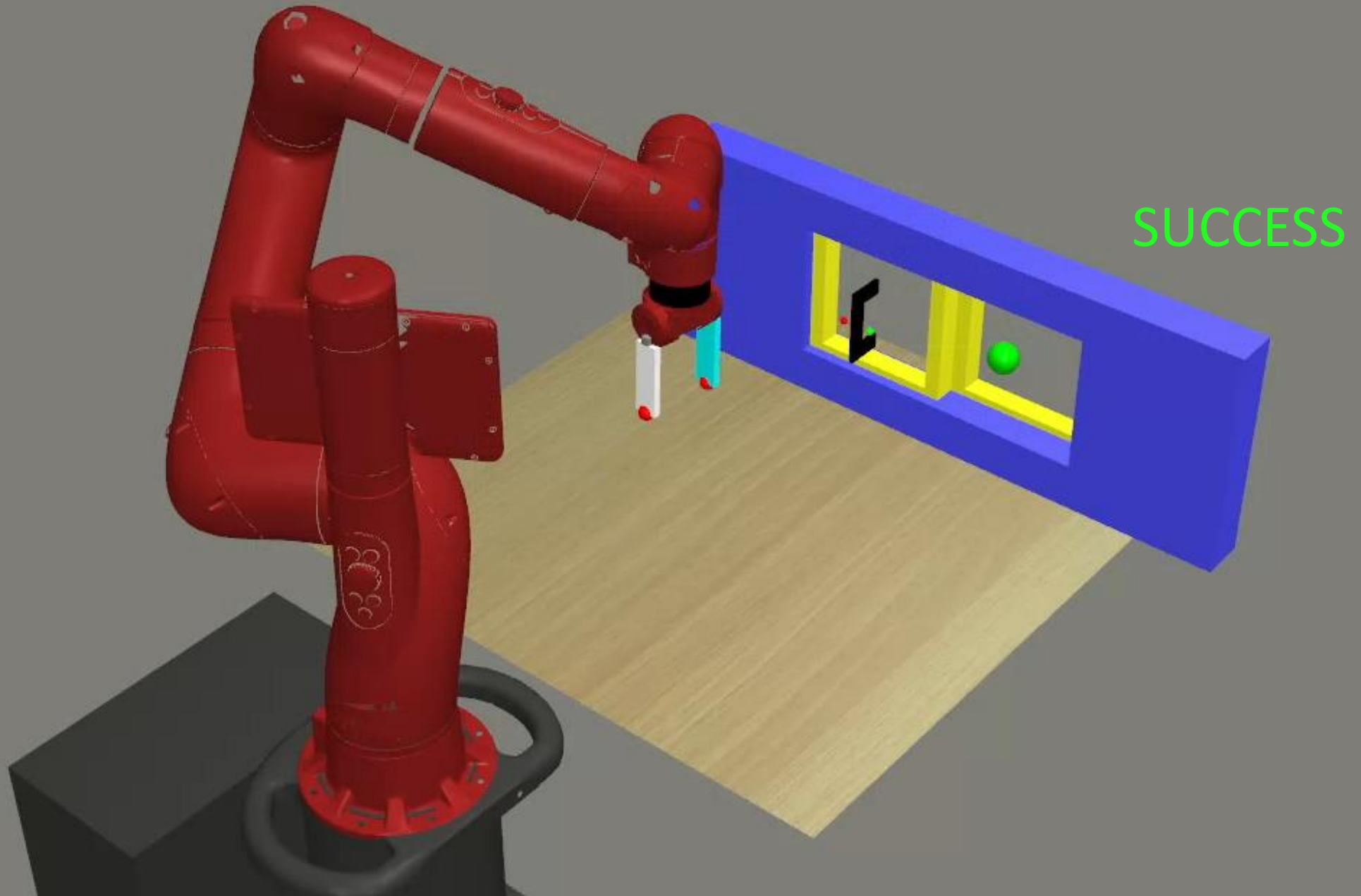
Why does DCLR outperform DCBC?

DCBC FAILURE



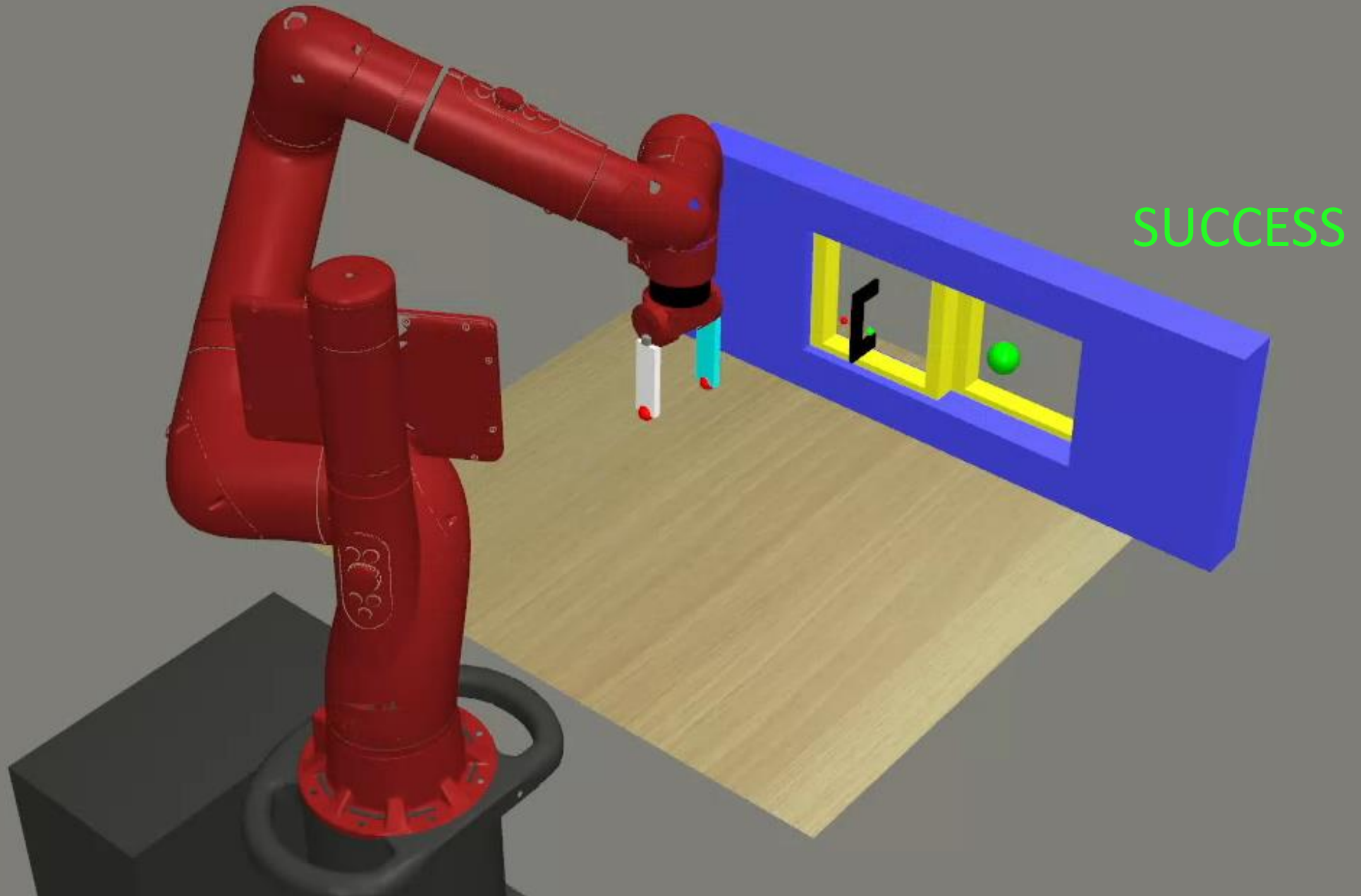
Why does DCLR outperform DCBC?

DCRL RECOVERY



Why does DCLR outperform DCBC?

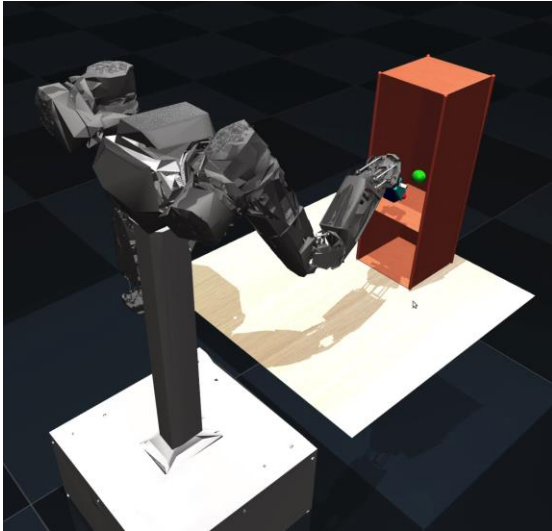
DCRL RECOVERY



Motivation. Control a robot given human demonstrations.

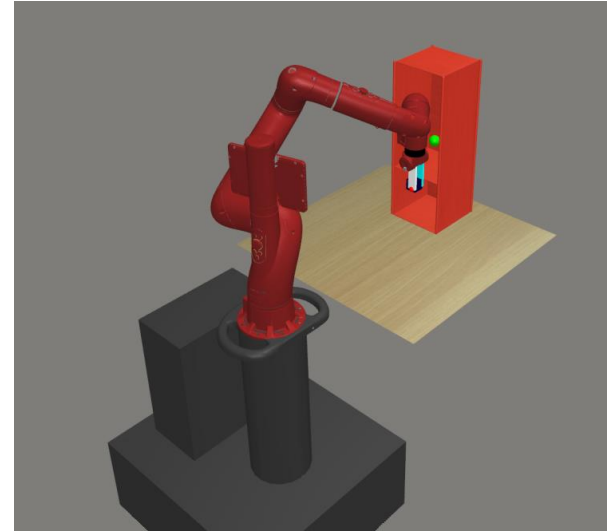
Experiment. Control Sawyer robot given demo's from an AMBIDEX robot.

AMBIDEX demo's



→ DCRL agent →

Sawyer robot doing shelf-place task



Results.

Demo's
from



# demo's	1	5	1	5
Sawyer	51%	51%	316	323
AMBIDEX	45%	48%	308	329
	success rates		average returns	

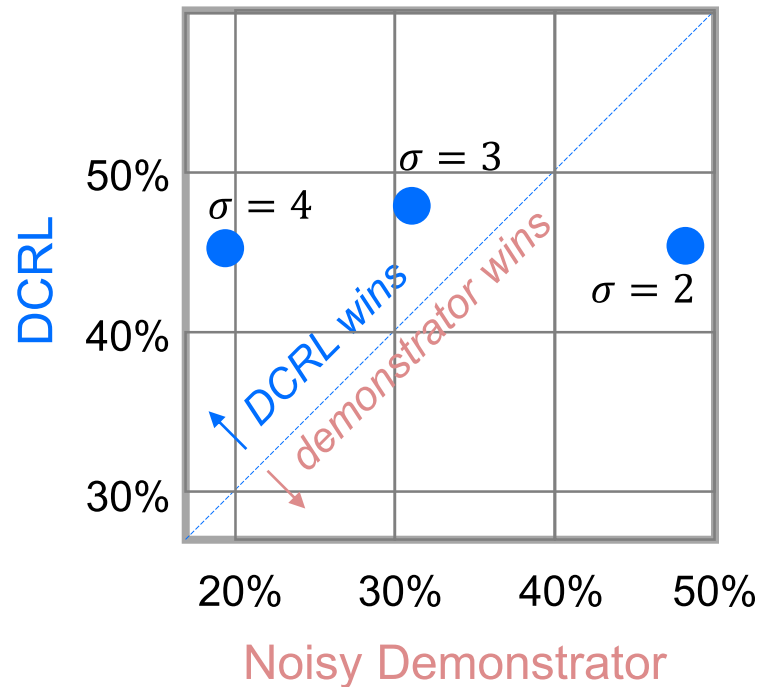
Motivation. Demonstrations are often suboptimal

Experiment. Add noise $\sim \mathcal{N}(0, \sigma^2 I_{4 \times 4})$ to demonstrator actions (only at test)

Motivation. Demonstrations are often suboptimal

Experiment. Add noise $\sim \mathcal{N}(0, \sigma^2 I_{4 \times 4})$ to demonstrator actions (only at test)

Success Rates on Meta-World



For $\sigma > 2$ DCRL outperforms the noisy demonstrator

DCRL is a new, third family of approaches to few-shot imitation
 $\{ \text{IRL, BC} \} \cup \{ \text{DCRL} \}$

Pros

- + suboptimal demo's
- + demonstrator domain shift
- + no test-time exploration

Cons

- needs training rewards

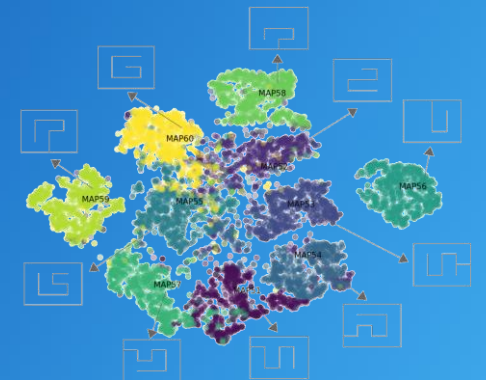
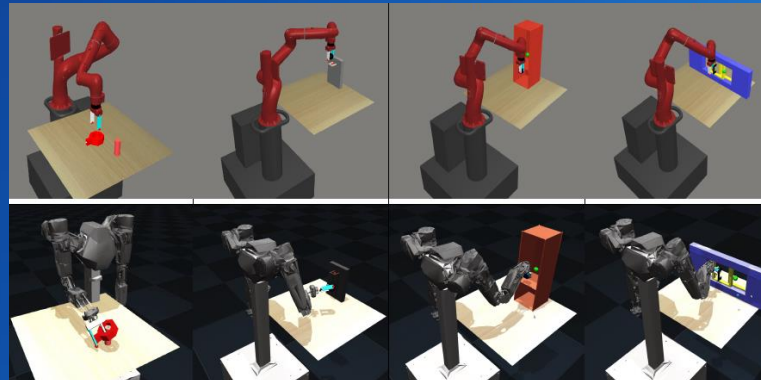
Future

- ☀ more training tasks
- ☀ better actor-critic training
- ☀ video of humans + real robot

Thank you

See the full paper for more experiments, t -SNE plots of the task embeddings, ...

See europe.naverlabs.com for videos of task execution



NAVER LABS