

# Causal Flow Q-Learning for Robust Offline Reinforcement Learning



**Mingxuan Li**



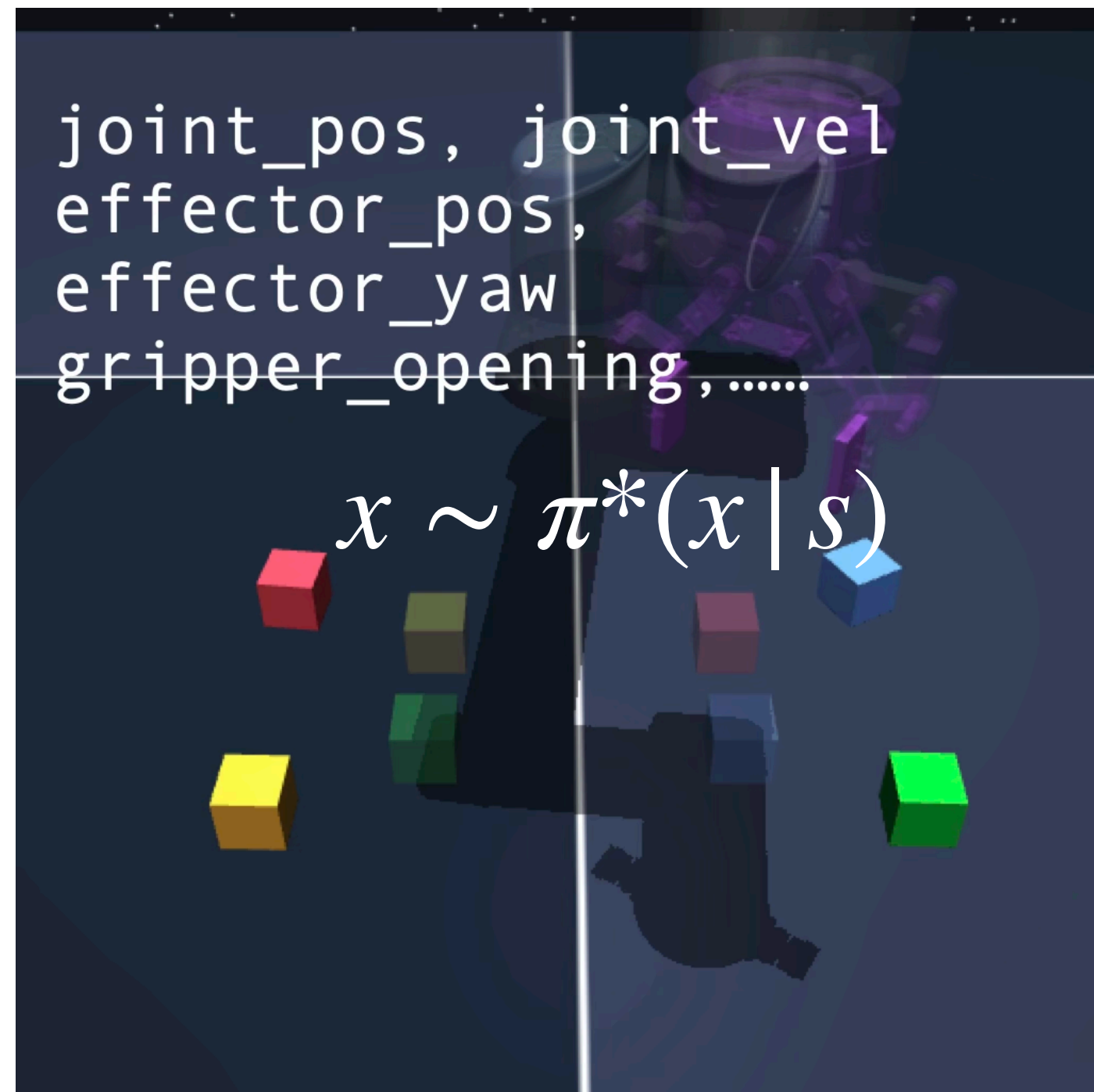
**Junzhe Zhang**



**Elias Bareinboim**

# Challenges of Confounding in Offline RL

In pixel-based offline RL tasks, the dataset is composed of expert generated **actions using structured state observations** and **blurry image observations**.



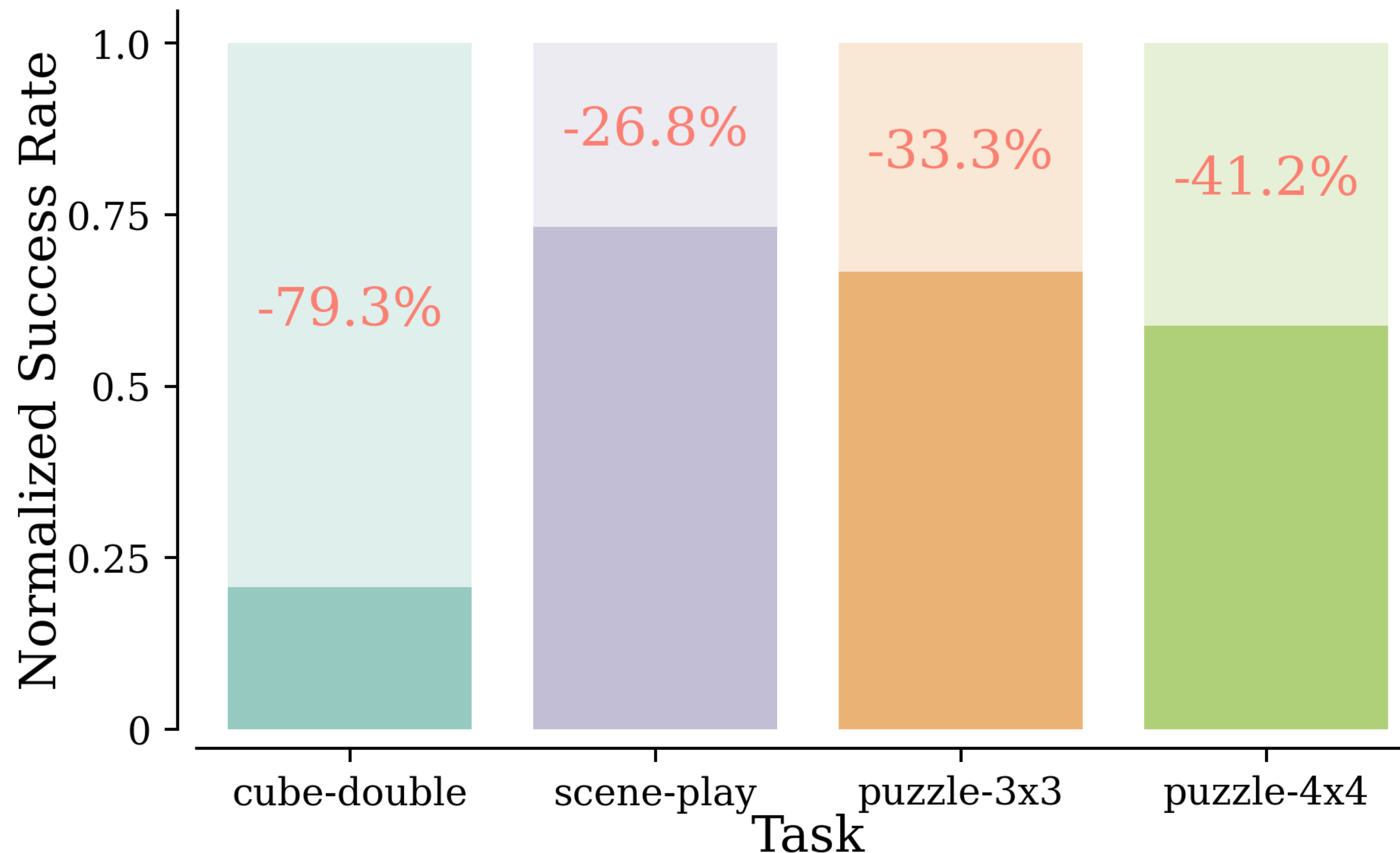
+



=  $(o, x)$ ?

# Challenges of Confounding in Offline RL

This **mismatch (confounding)** caused SOTA offline RL algorithms' performance drop sharply in pixel-based tasks, even though they achieve high success rate in the corresponding state-based version.

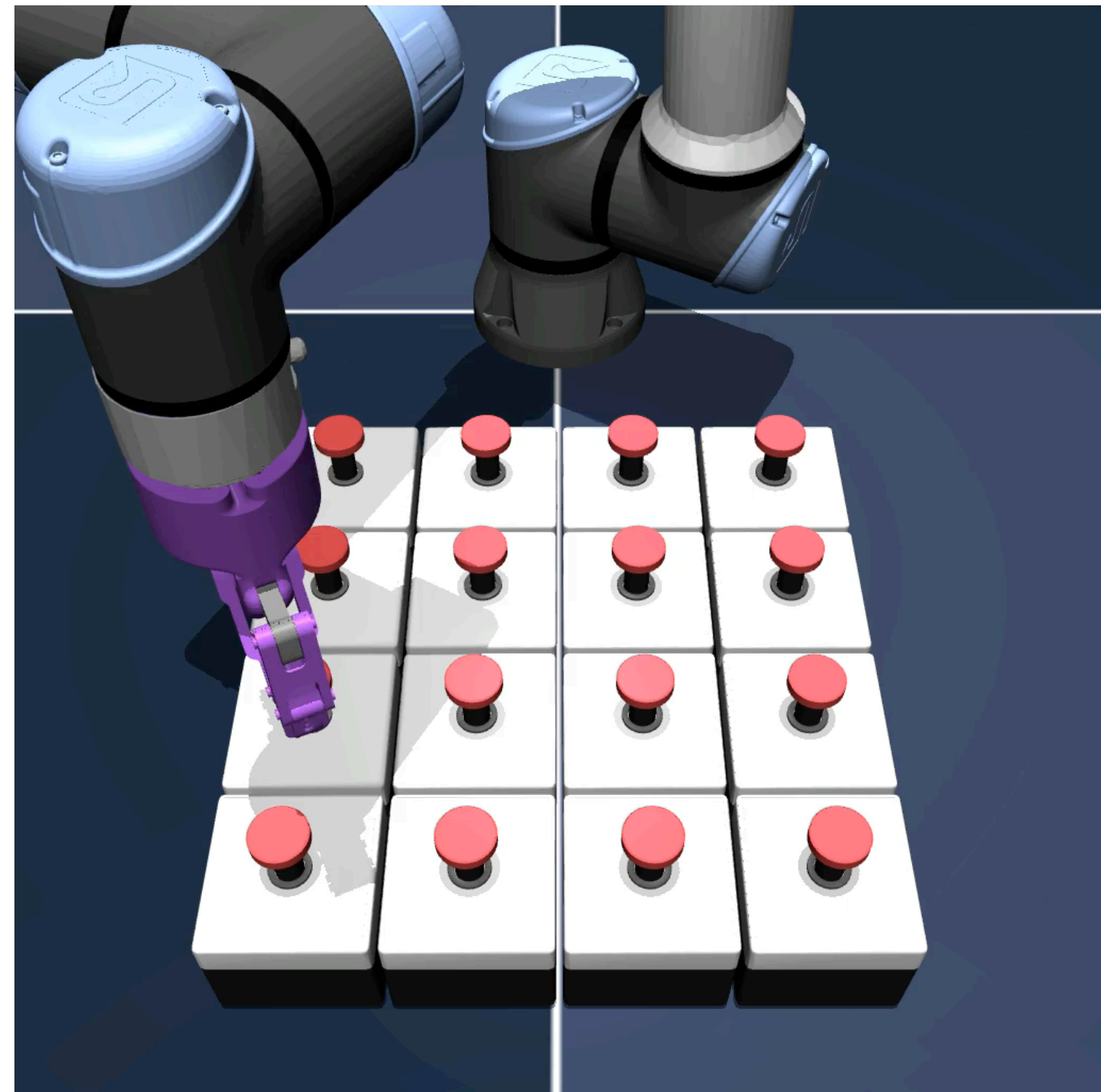


# Challenges of Confounding in Offline RL

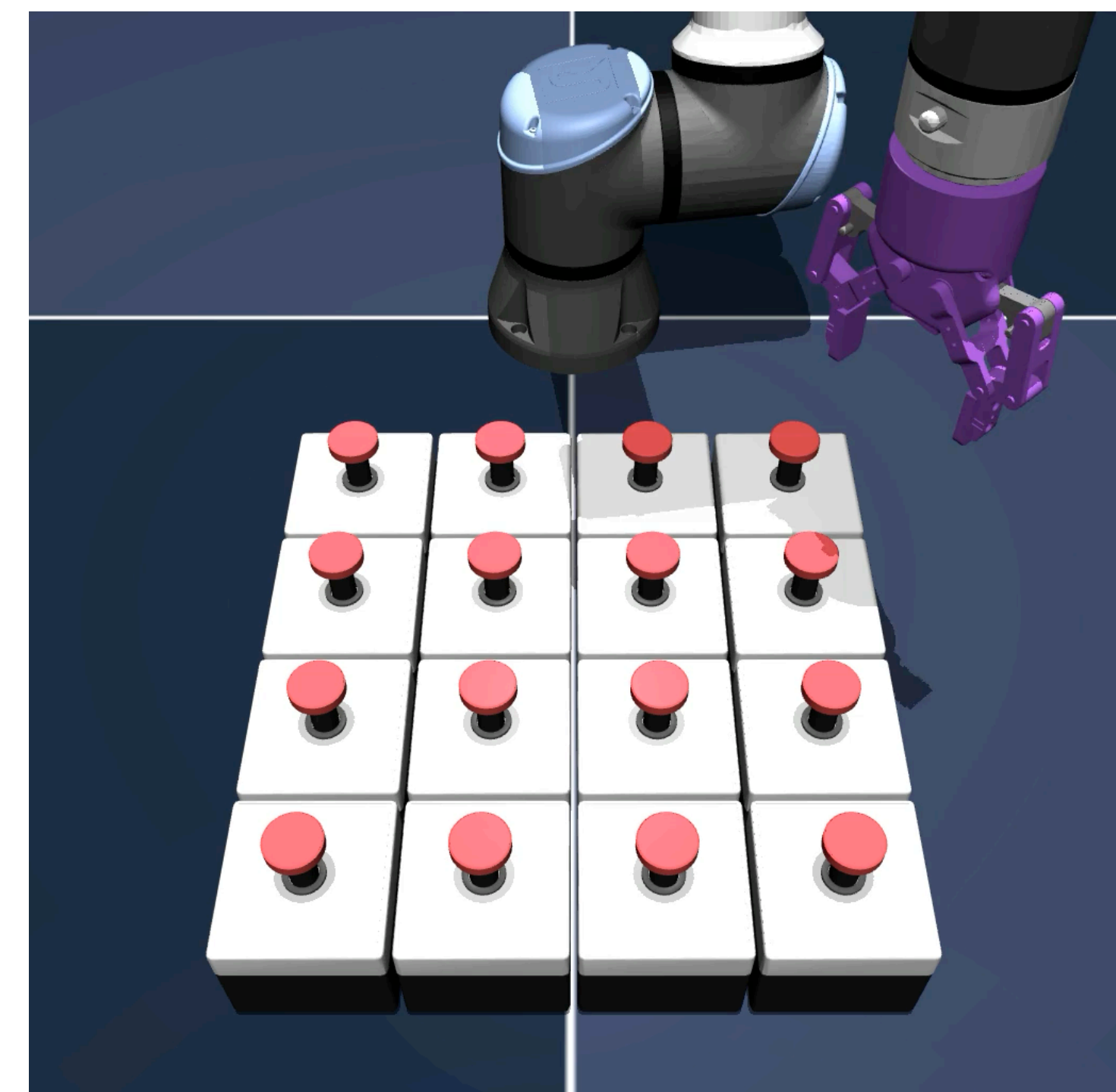
We develop a confounding robust offline RL objective that,

- is grounded in causal language, and
- improves success rate significantly in both offline & offline-to-online pixel tasks.

FQL



Causal-FQL



Avg. Success  
Rate +20%

# Modeling the Problem with Confounded MDP (CMDP)

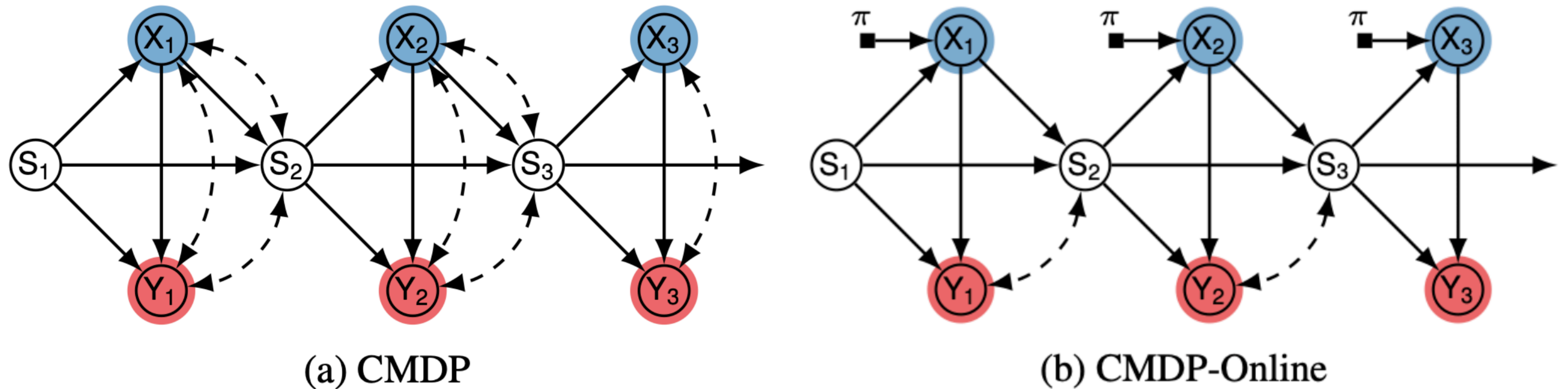
**Definition 1 (Confounded MDP).** A Confounded Markov Decision Process (CMDP)  $\mathcal{M}$  is a tuple of  $\langle \mathcal{S}, \mathcal{X}, \mathcal{Y}, \mathcal{U}, \mathbb{F}, \mathbb{P} \rangle$  where,

- $\mathcal{S}, \mathcal{X}, \mathcal{Y}$  are, respectively, the space of observed states, actions, and rewards;
- $\mathcal{U}$  is the space of unobserved exogenous noise;
- $\mathbb{F}$  is a set consisting of the transition function  $\tau_h : \mathcal{S} \times \mathcal{X} \times \mathcal{U} \mapsto \mathcal{S}$ , behavioral policy  $\beta_h : \mathcal{S} \times \mathcal{U} \mapsto \mathcal{X}$ , and reward function  $r_h : \mathcal{S} \times \mathcal{X} \times \mathcal{U} \mapsto \mathcal{Y}$  for every time step  $h = 1, \dots$ ;
- $\mathbb{P}$  is a set of distributions  $P_h$  over the unobserved domain  $\mathcal{U}$  for every time step  $h = 1, \dots$ .

# CMDP - Graphical Model

Those bi-directed dashed arrows are confounders indicating the variables that are,

1. Observable to expert; and
2. Not observable to the learner.



# Bounding Worst Case via Causal Bellman Equation

Directly maximizing the Q-values from such confounded data results in poor performance.

We propose the Causal Bellman Equation to lower bound the optimal Q-values.

**Theorem 1 (Causal Bellman Equation).** For a CMDP environment  $\mathcal{M}$  with reward  $Y_h \geq a, a \in \mathbb{R}$ , the optimal value of interventional policies,  $Q^*(s, x), \forall s \in \mathcal{S}$ , is lower bounded by  $\underline{Q}_*(s, x)$  satisfying the Causal Bellman Optimality Equation, for every step  $h = 1, \dots, H$ ,

$$\underline{Q}_*(s, x) = P(x | s) \left( \widetilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \widetilde{\mathcal{T}}(s, x, s') \max_{x'} \underline{Q}_*(s', x') \right) + P(\neg x | s) \left( a + \gamma \min_{s'} \max_{x'} \underline{Q}_*(s', x') \right)$$

where  $P(x | s) = P(X_t = x | S_t = s)$  is the propensity score calculated from the offline dataset,  $\widetilde{T}, \widetilde{R}$  are estimated transition distribution and rewards, separately.

# Causal Bellman Equation in the Continuous Space

- There are two policies involved in estimating value bounds, one **behavioral**, one **learner's policy** we are trying to train.

$$\underline{Q}_*(s, x) = P(x | s) \left( \tilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \tilde{\mathcal{T}}(s, x, s') \max_{x'} \underline{Q}_*(s', x') \right) + P(\neg x | s) \left( a + \gamma \min_{s'} \max_{x'} \underline{Q}_*(s', x') \right)$$

$x \sim \pi(\cdot | s)$

- For continuous actions,  $P(x|s)$  is zero. Thus, we instead interpret  $P(x|s)$  as an action  $x$  is sampled from the **behavioral policy**.
- And this can be thought of as a binary classification problem and can be learned by a neural network classifier,

$$P(x | s) = D_{\theta}(s, x)$$

# Approximate Causal Bellman Equation via Ensembles

- We propose to fit an ensemble of neural networks to represent different possible “realities” (SCMs) corresponding to the same offline dataset,

$$Q_{\theta_i}, i \in [1, K]$$

- Taking min over the ensembles and plug in the discriminator,

$$\underline{Q}_*(s, x) = P(x | s) \left( \widetilde{\mathcal{R}}(s, x) + \gamma \sum_{s', x'} \widetilde{\mathcal{T}}(s, x, s') \max_{x'} \underline{Q}_*(s', x') \right) + P(\neg x | s) \left( a + \gamma \min_{s'} \max_{x'} \underline{Q}_*(s', x') \right)$$



$$\underline{\hat{Q}}_*(s, x) = \frac{1}{K} D_{\theta}(x | s) \sum_{i=1}^K Q_{\theta_i}(s, x) + (1 - D_{\theta}(x | s)) \min_{i \in [1, K]} Q_{\theta_i}$$

# Causal Flow Q-Learning (Algorithm)

While not converged:

Sample a batch  $\{(s, x, r, s')\} \sim \mathcal{D}$

Train critic ensemble  $Q_{\theta_i}$

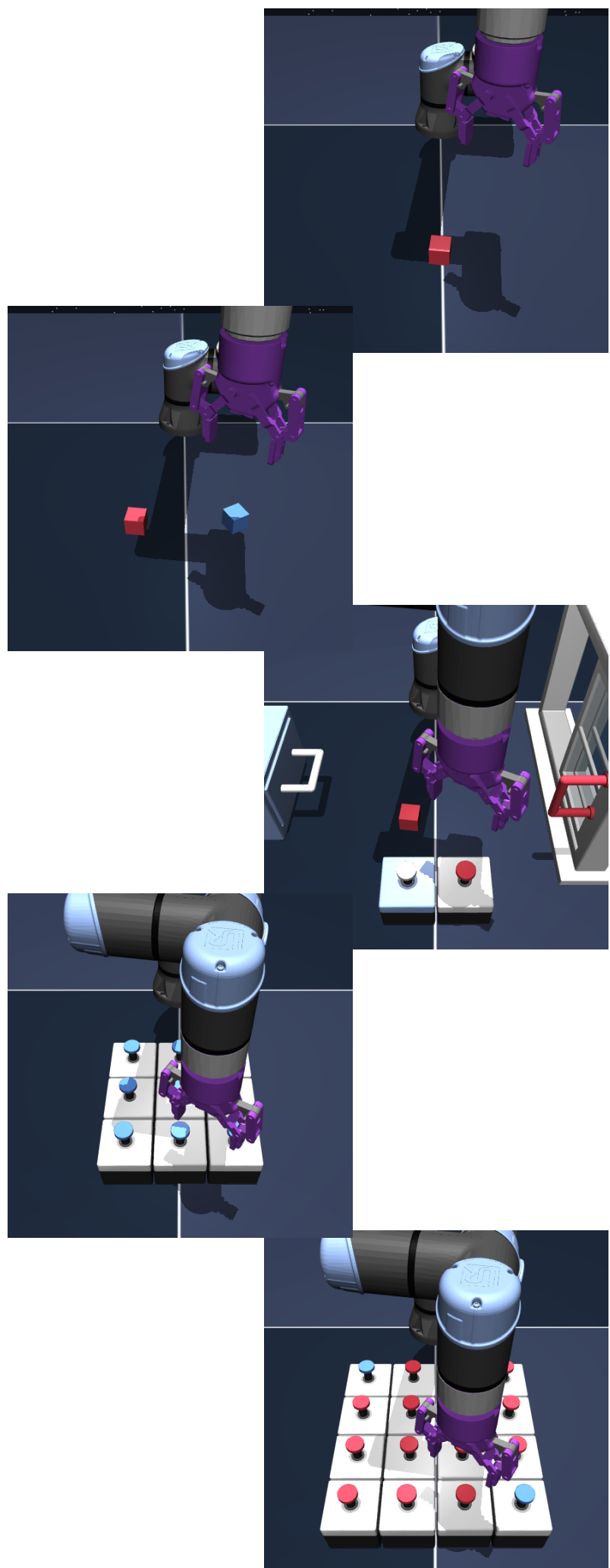
Train flow policy  $\mu_{\omega}$  and discriminator  $D_{\theta}$

Train one-step target policy:

$$\mathcal{L}_{\pi}(\theta) = \mathbb{E}_{s \sim \mathcal{D}, x^{\pi} \sim \pi_{\theta}}[-\underline{\hat{Q}}_*(s, x^{\pi})] + \alpha \mathcal{L}_{\text{Distill}}(\omega)$$

$$\underline{\hat{Q}}_*(s, x) = \frac{1}{K} D_{\theta}(x | s) \sum_{i=1}^K Q_{\theta_i}(s, x) + (1 - D_{\theta}(x | s)) \min_{i \in [1, K]} Q_{\theta_i}$$

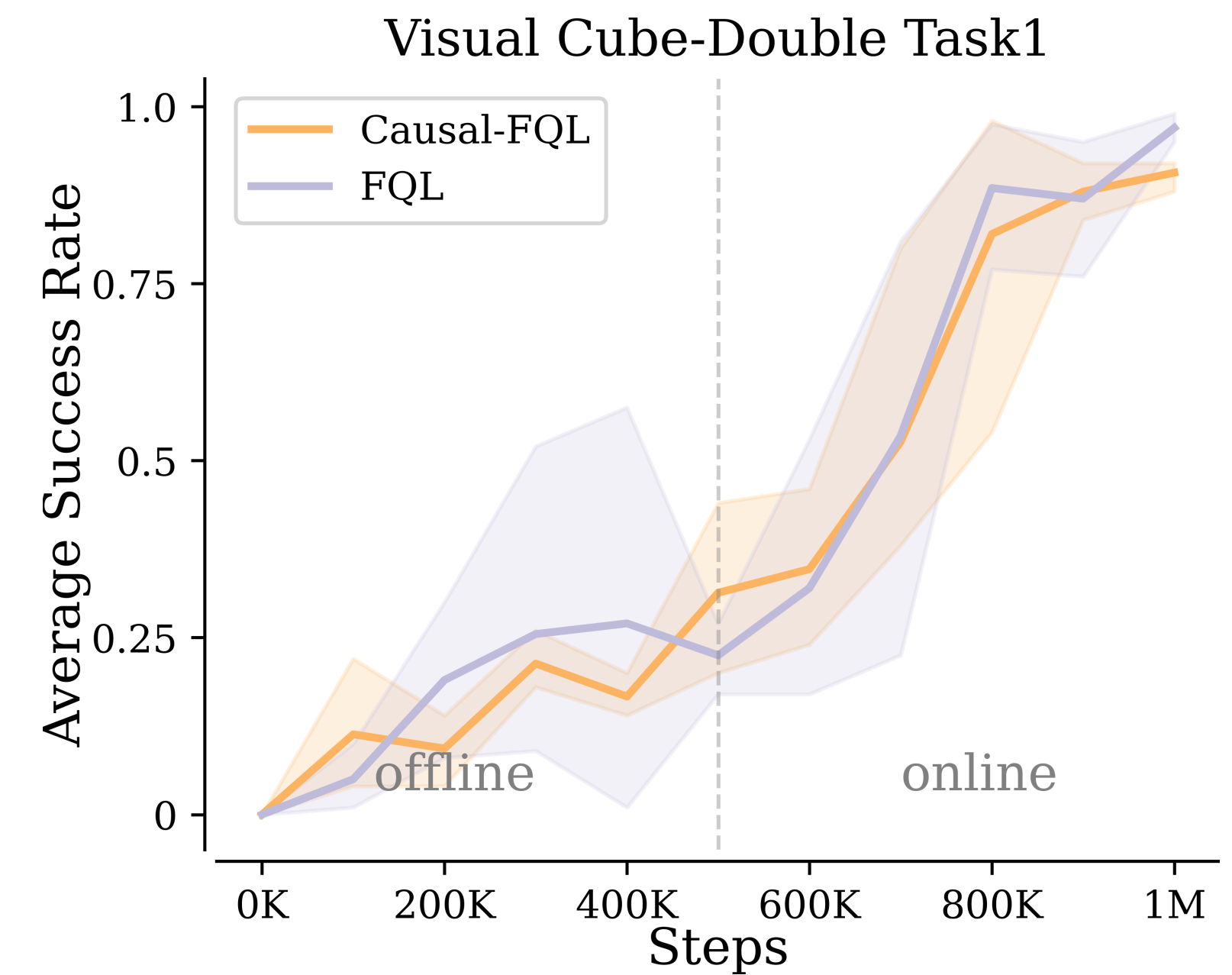
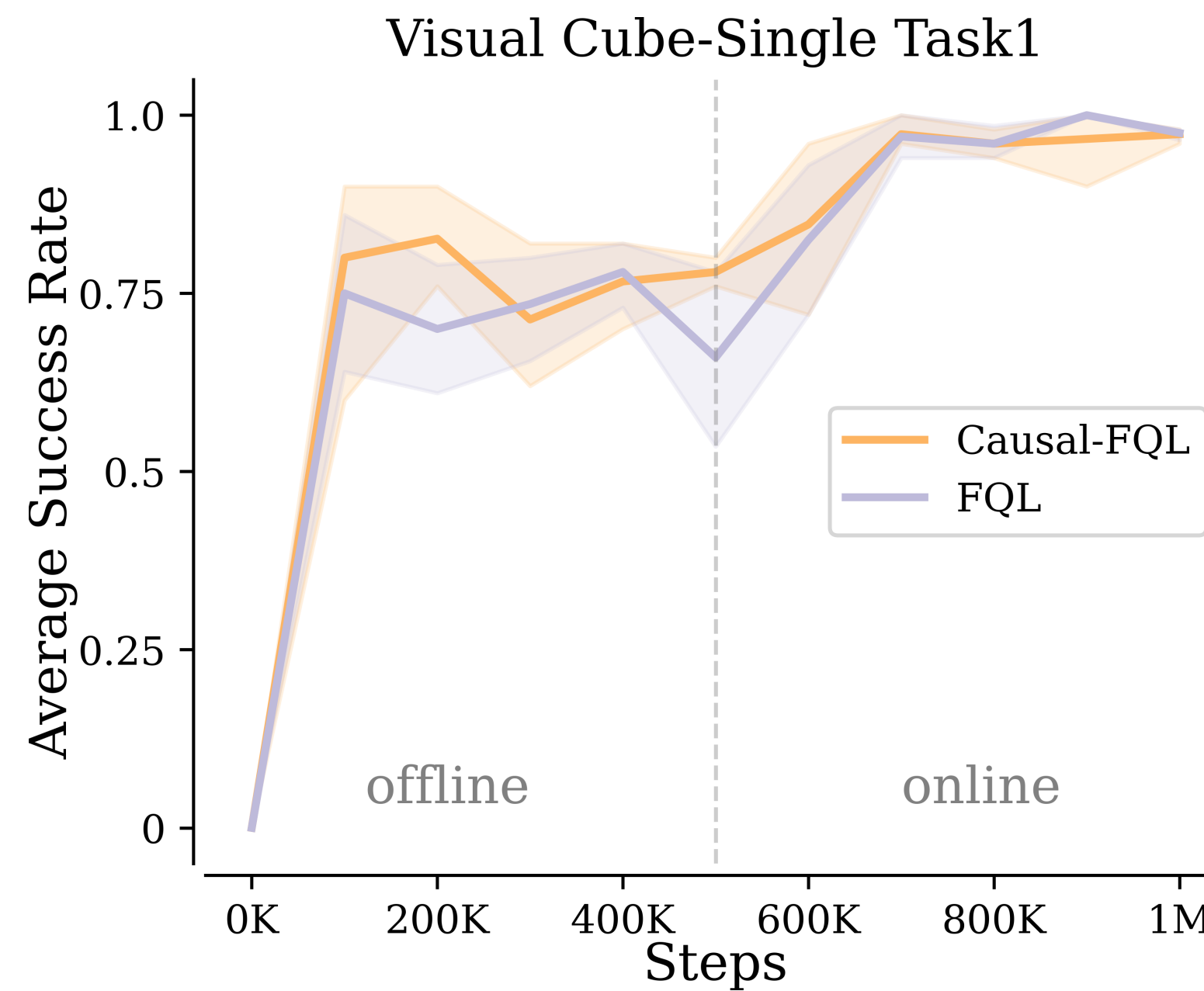
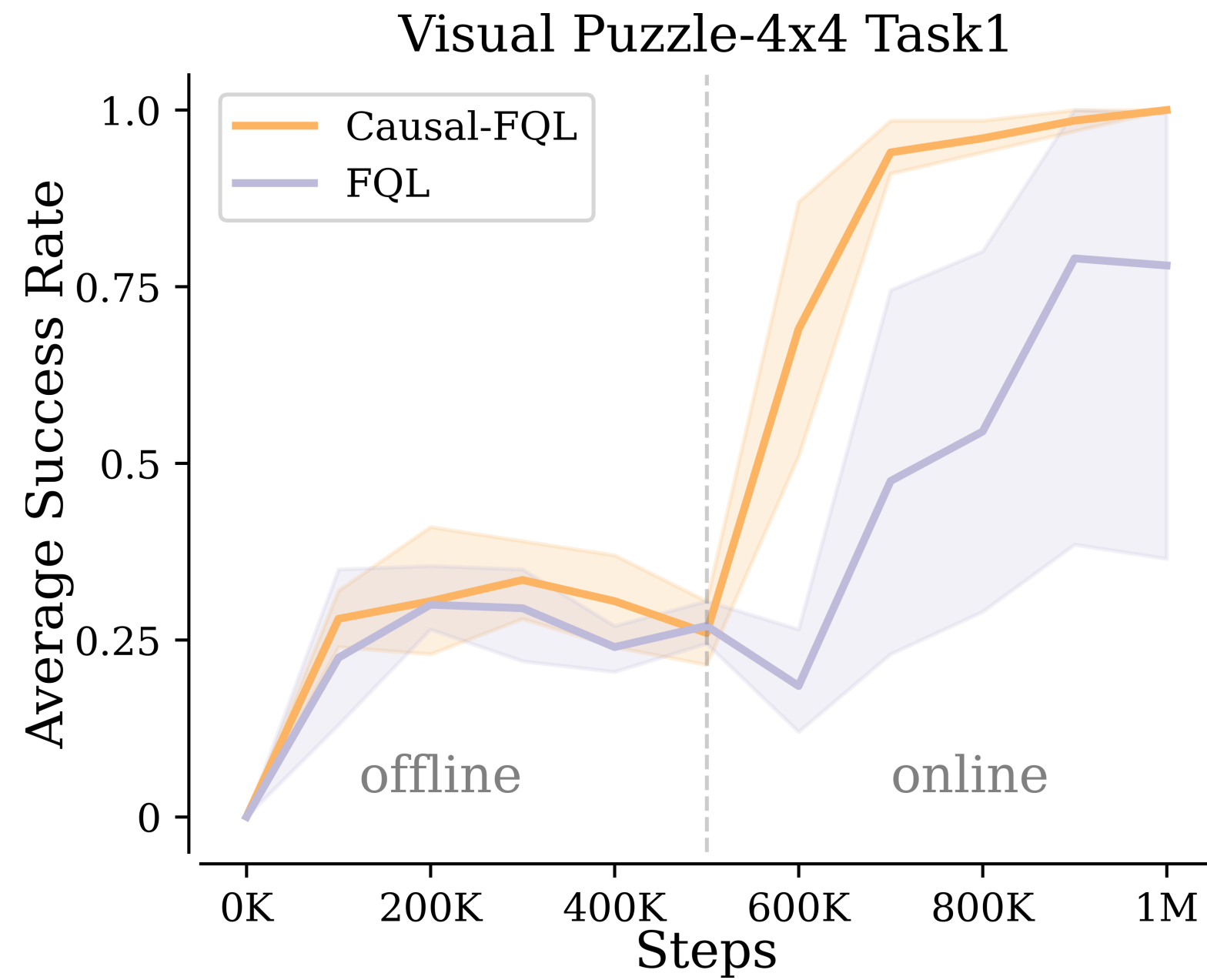
# Experimental Results - Offline



	Gaussian Policies			Flow Policies			
	IQL	ReBRAC	FBRAC	IQN	IFQL	FQL	Causal-FQL
visual-cube-single-play-singletask-task1-v0 (*)	70 ± 12	<b>83 ± 6</b>	55 ± 8	5 ± 4	49 ± 7	81 ± 12	<b>85 ± 6</b>
visual-cube-single-play-singletask-task2-v0	59 ± 13	<b>85 ± 6</b>	70 ± 11	58 ± 9	45 ± 11	58 ± 12	<b>81 ± 5</b>
visual-cube-single-play-singletask-task3-v0	64 ± 35	<b>88 ± 4</b>	81 ± 4	83 ± 11	67 ± 6	73 ± 16	<b>84 ± 5</b>
visual-cube-single-play-singletask-task4-v0	68 ± 12	<b>78 ± 9</b>	56 ± 19	<b>78 ± 3</b>	34 ± 13	63 ± 8	<b>74 ± 5</b>
visual-cube-single-play-singletask-task5-v0	59 ± 14	<b>76 ± 9</b>	59 ± 10	64 ± 9	27 ± 10	49 ± 12	<b>79 ± 7</b>
visual-cube-double-play-singletask-task1-v0 (*)	34 ± 23	4 ± 4	6 ± 2	4 ± 1	8 ± 6	23 ± 4	<b>43 ± 11</b>
visual-cube-double-play-singletask-task2-v0	<b>2 ± 1</b>	0 ± 0	2 ± 2	0 ± 0	0 ± 0	0 ± 0	<b>2 ± 1</b>
visual-cube-double-play-singletask-task3-v0	<b>7 ± 4</b>	2 ± 2	2 ± 1	0 ± 0	1 ± 1	4 ± 2	4 ± 2
visual-cube-double-play-singletask-task4-v0	<b>1 ± 1</b>	0 ± 0	0 ± 0	0 ± 0	0 ± 0	0 ± 0	<b>1 ± 1</b>
visual-cube-double-play-singletask-task5-v0	<b>11 ± 2</b>	0 ± 0	0 ± 0	1 ± 1	2 ± 1	4 ± 1	4 ± 2
visual-scene-play-singletask-task1-v0 (*)	<b>97 ± 2</b>	<b>98 ± 4</b>	46 ± 4	<b>95 ± 2</b>	86 ± 10	<b>98 ± 3</b>	<b>100 ± 0</b>
visual-scene-play-singletask-task2-v0	21 ± 16	30 ± 15	0 ± 0	79 ± 15	0 ± 0	86 ± 8	<b>92 ± 3</b>
visual-scene-play-singletask-task3-v0	12 ± 9	10 ± 7	10 ± 3	<b>31 ± 14</b>	19 ± 2	22 ± 6	22 ± 1
visual-scene-play-singletask-task4-v0	1 ± 0	0 ± 0	0 ± 0	0 ± 0	0 ± 0	1 ± 1	<b>2 ± 1</b>
visual-scene-play-singletask-task5-v0	<b>0 ± 0</b>	<b>0 ± 0</b>	<b>0 ± 0</b>	<b>0 ± 0</b>	<b>0 ± 0</b>	<b>0 ± 0</b>	<b>0 ± 0</b>
visual-puzzle-3x3-play-singletask-task1-v0 (*)	7 ± 15	88 ± 4	7 ± 2	84 ± 1	<b>100 ± 0</b>	94 ± 1	<b>98 ± 1</b>
visual-puzzle-3x3-play-singletask-task2-v0	0 ± 0	<b>12 ± 1</b>	0 ± 0	6 ± 2	0 ± 0	0 ± 0	1 ± 0
visual-puzzle-3x3-play-singletask-task3-v0	0 ± 0	1 ± 1	0 ± 0	1 ± 0	2 ± 1	0 ± 0	<b>3 ± 2</b>
visual-puzzle-3x3-play-singletask-task4-v0	1 ± 1	0 ± 1	0 ± 0	3 ± 0	1 ± 0	5 ± 4	<b>8 ± 3</b>
visual-puzzle-3x3-play-singletask-task5-v0	0 ± 0	0 ± 0	0 ± 0	1 ± 0	0 ± 0	1 ± 2	<b>19 ± 4</b>
visual-puzzle-4x4-play-singletask-task1-v0 (*)	0 ± 0	26 ± 6	0 ± 0	21 ± 3	8 ± 15	33 ± 6	<b>37 ± 4</b>
visual-puzzle-4x4-play-singletask-task2-v0	0 ± 0	0 ± 0	1 ± 1	<b>14 ± 2</b>	1 ± 1	1 ± 2	4 ± 1
visual-puzzle-4x4-play-singletask-task3-v0	0 ± 0	0 ± 0	0 ± 0	9 ± 3	9 ± 15	16 ± 6	<b>18 ± 4</b>
visual-puzzle-4x4-play-singletask-task4-v0	0 ± 0	1 ± 1	0 ± 0	<b>12 ± 4</b>	2 ± 2	2 ± 1	<b>6 ± 2</b>
visual-puzzle-4x4-play-singletask-task5-v0	0 ± 0	0 ± 0	1 ± 1	<b>11 ± 4</b>	0 ± 0	1 ± 1	1 ± 1

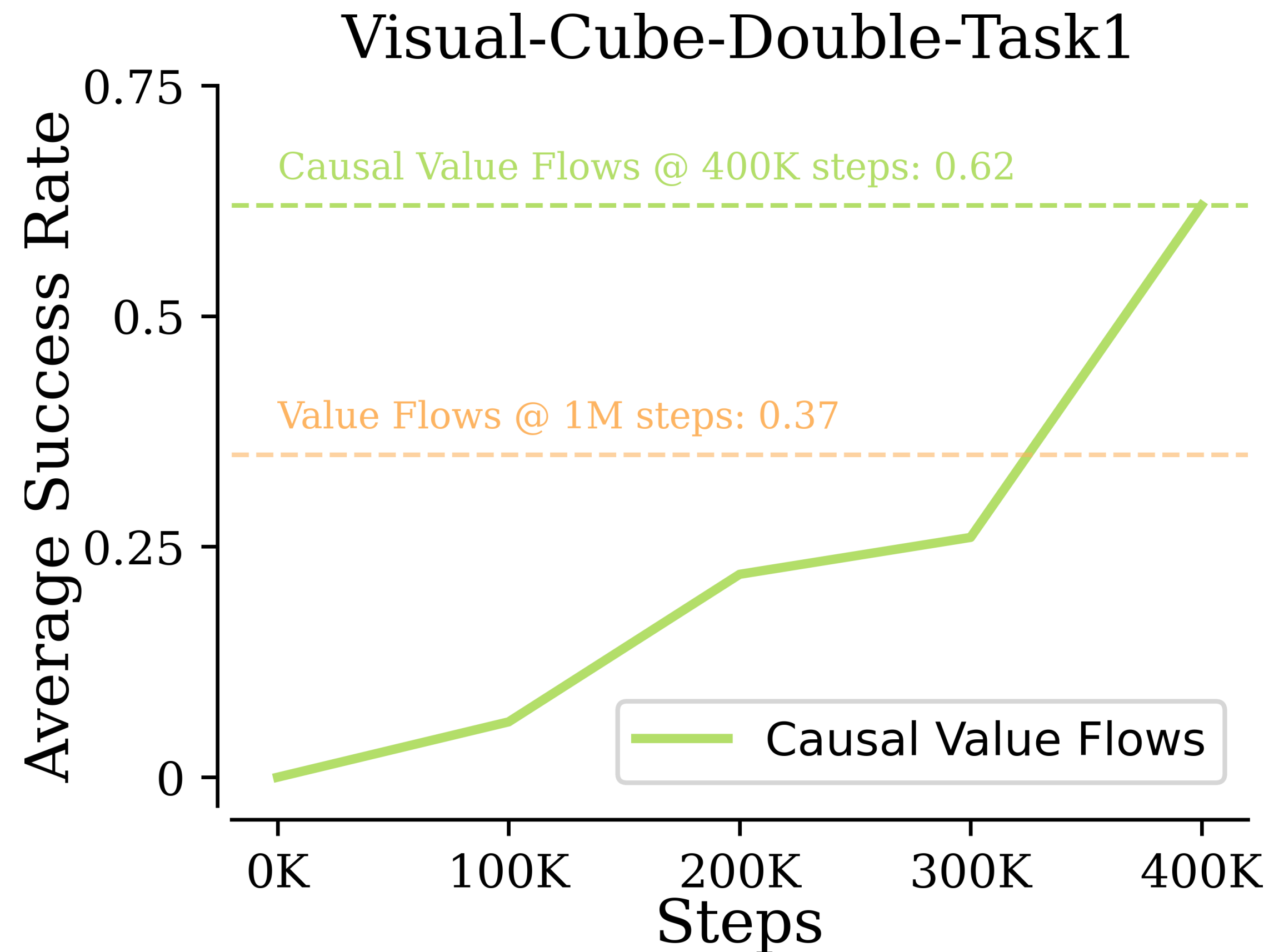
# Experimental Results - Offline to Online

On puzzle-4x4 task1, our method is the first to achieve 100% success rate.



# Experimental Results - Causal Value Flows

Our proposed causal offline RL objective can also be implemented upon other algorithms. We use value flows<sup>1</sup>, a more recent SOTA, to demonstrate this.



<sup>1</sup>Value Flows, ICLR 2026