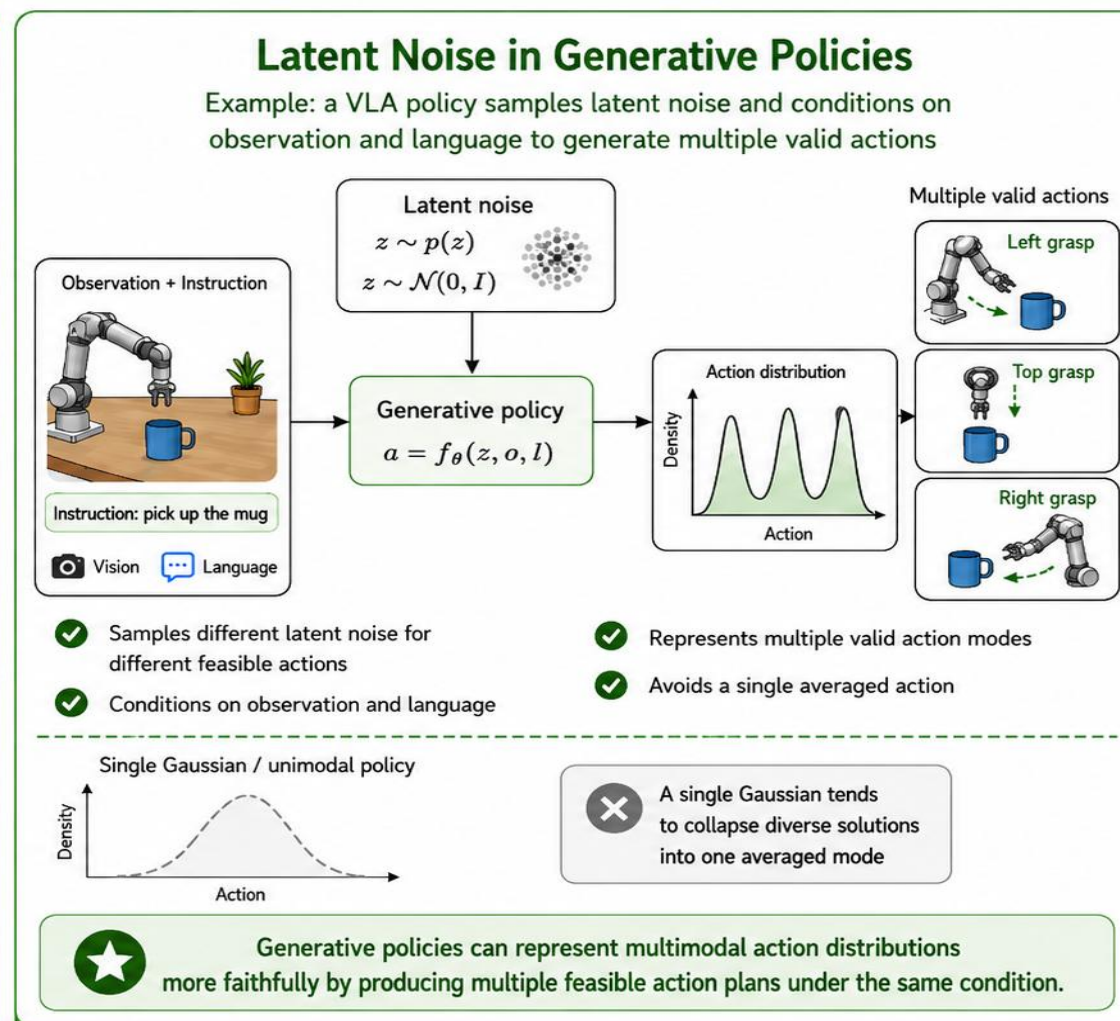
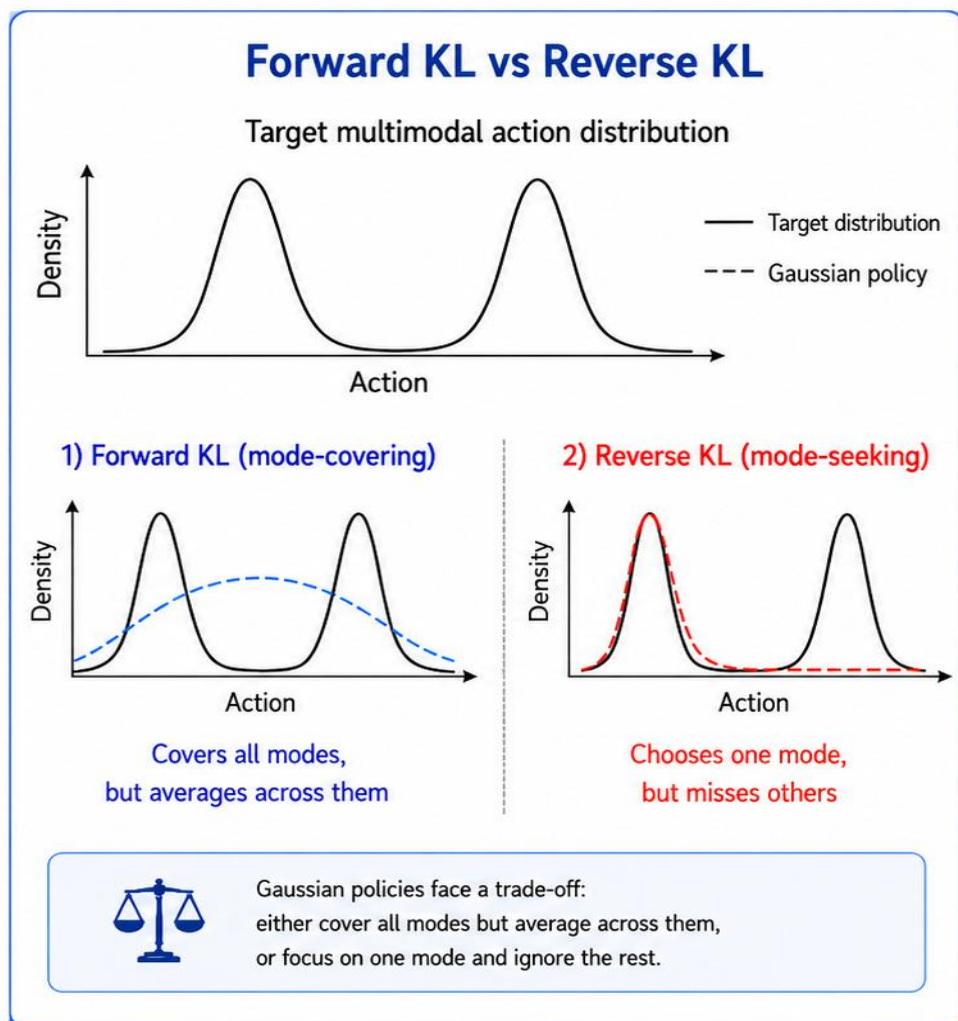


Generative Models For RL



Generative policies better capture *multimodal action distributions*.



Steering Your Diffusion Policy with Latent Space Reinforcement Learning

Andrew Wagenmaker*
UC Berkeley

Mitsuhiko Nakamoto*
UC Berkeley

Yunchu Zhang*
University of Washington

Seohong Park
UC Berkeley

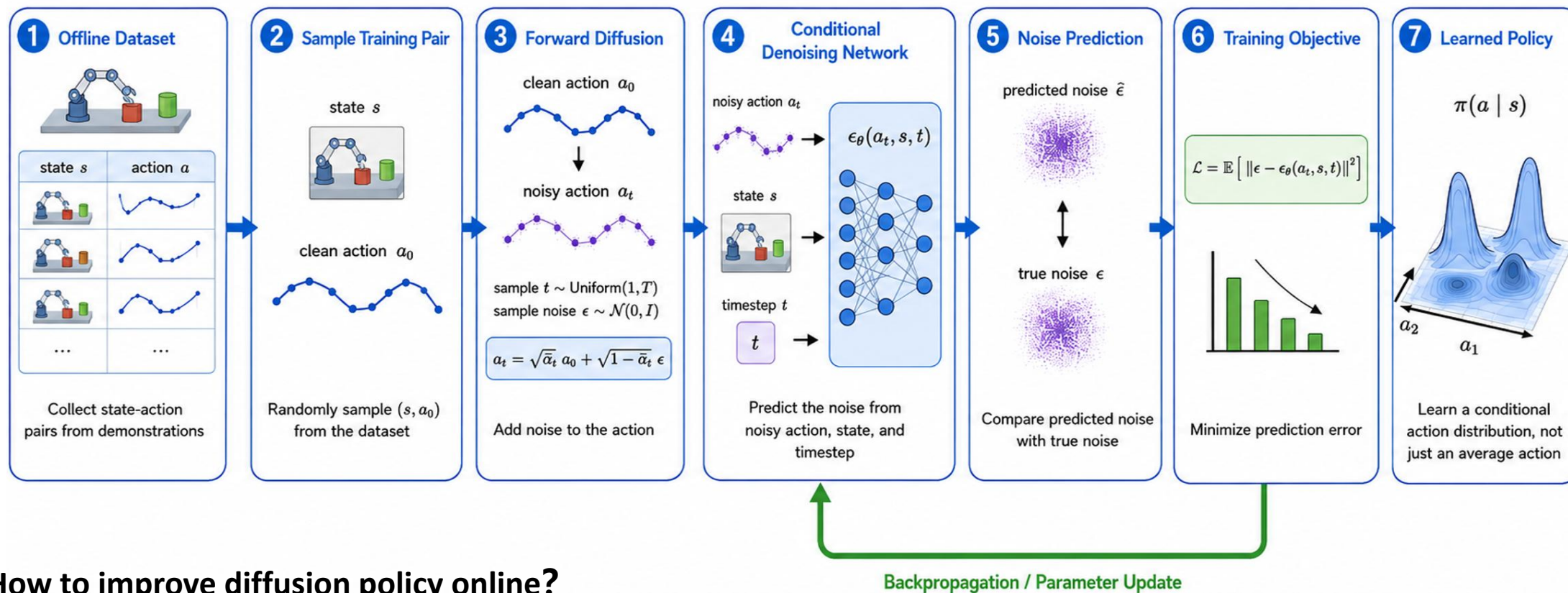
Waleed Yagoub
University of Washington

Anusha Nagabandi
Amazon

Abhishek Gupta*
University of Washington

Sergey Levine*
UC Berkeley

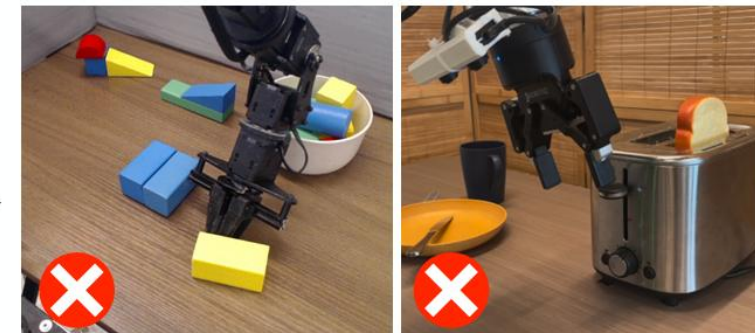
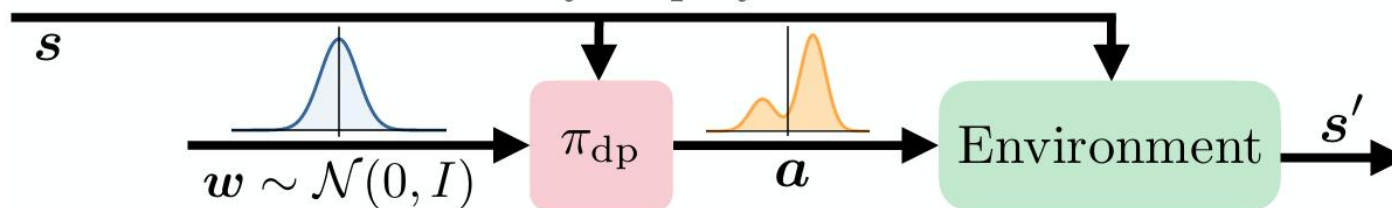
Training process of diffusion policy



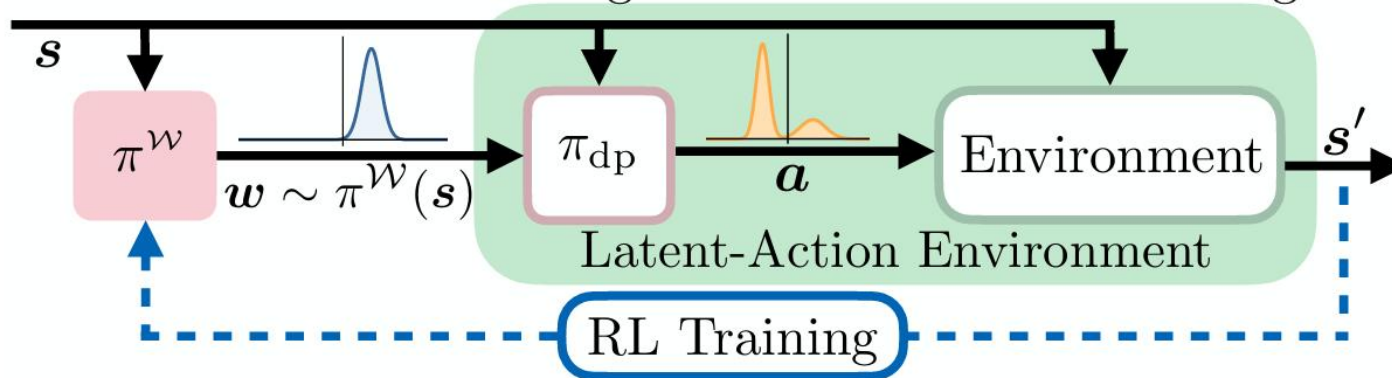
How to improve diffusion policy online?

- Collect more expert demonstrations \longrightarrow costly and time-consuming
- Online RL \longrightarrow computationally heavy, unstable, and may destroy the pretrained behavior prior

Standard Diffusion Policy Deployment



DSRL: Diffusion Steering via Reinforcement Learning



$$a \sim \pi_{dp}^w(s, w) \quad w \sim \mathcal{N}(0, I) \quad w_t \sim \pi^w(s_t)$$

Optimize the pre-trained policy by *noise manipulation*.

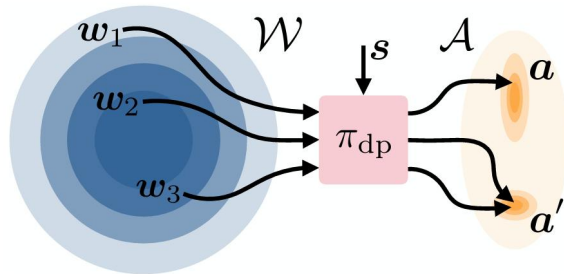
DSRL-SAC

$$w \sim \pi_{\phi}^{\mathcal{W}}(w|s)$$

$$a = \pi_{\text{dp}}^{\omega}(s, w)$$

$$(s, w, r, s')$$

$$Q^{\mathcal{W}}(s, w)$$



DSRL-NA

$$w' \neq w, \quad \pi_{\text{dp}}^{\omega}(s, w') \approx \pi_{\text{dp}}^{\omega}(s, w)$$

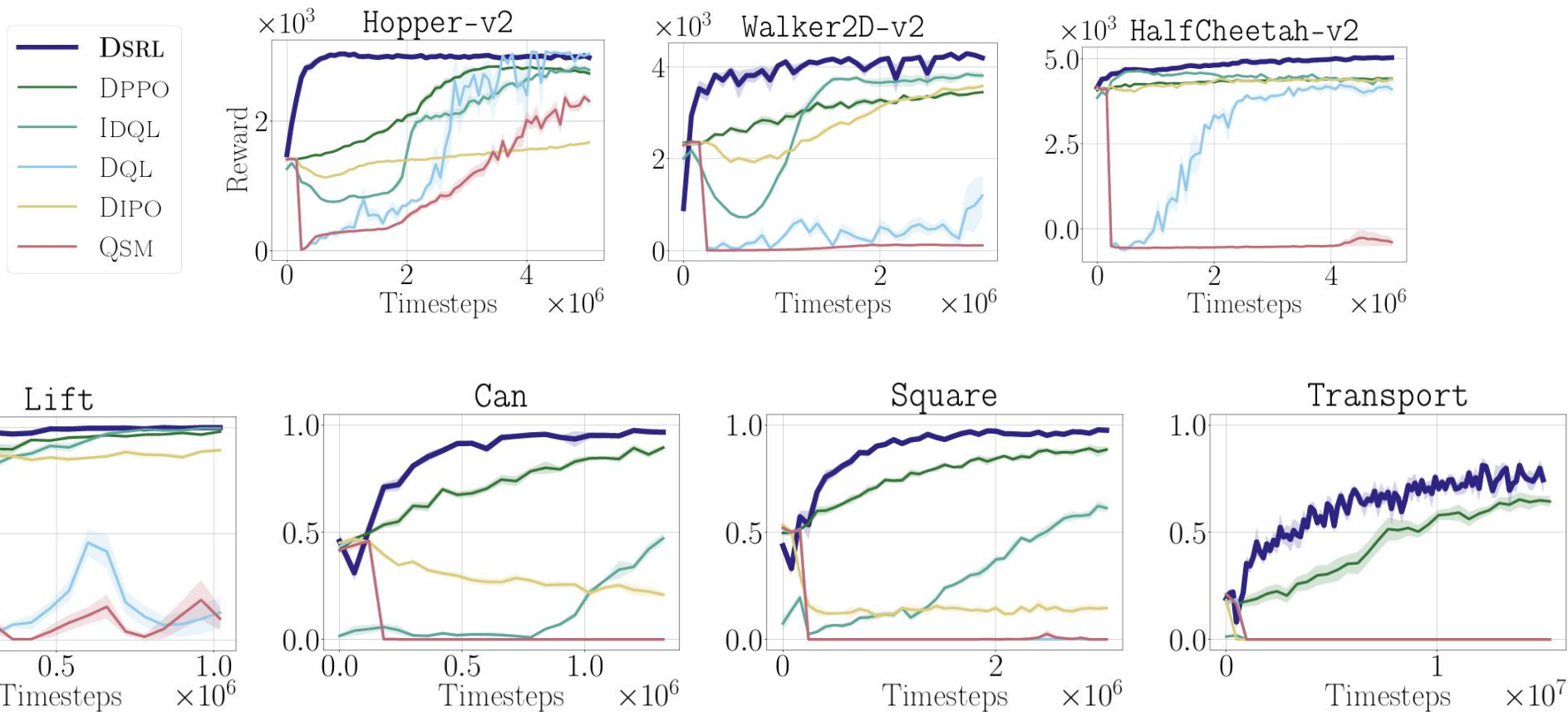
$$Q^{\mathcal{A}}(s, a) \longrightarrow Q^{\mathcal{W}}(s, w)$$

Critic distillation

Algorithm 1 Noise-Aliased Diffusion Steering via Reinforcement Learning (DSRL-NA)

- 1: **input:** pretrained diffusion policy $\pi_{\text{dp}}^{\mathcal{W}}$, offline data \mathcal{D}_{off} and/or online environment \mathcal{M}
 - 2: Initialize replay buffer $\mathfrak{B} \leftarrow \mathcal{D}_{\text{off}}$, \mathcal{A} -critic $Q^{\mathcal{A}}$, latent-noise critic $Q^{\mathcal{W}}$, latent-noise actor $\pi^{\mathcal{W}}$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Update $Q^{\mathcal{A}}$: $\min_{Q^{\mathcal{A}}} \mathbb{E}_{(s, a, r, s') \sim \mathfrak{B}, a' \sim \pi_{\text{dp}}^{\mathcal{W}}(s', \pi^{\mathcal{W}}(s'))} [(Q^{\mathcal{A}}(s, a) - r - \gamma \bar{Q}^{\mathcal{A}}(s', a'))^2]$
 - 5: Update $Q^{\mathcal{W}}$: $\min_{Q^{\mathcal{W}}} \mathbb{E}_{s \sim \mathfrak{B}, w \sim \mathcal{N}(0, I)} [(Q^{\mathcal{W}}(s, w) - Q^{\mathcal{A}}(s, \pi_{\text{dp}}^{\mathcal{W}}(s, w)))]^2$
 - 6: Update $\pi^{\mathcal{W}}$: $\max_{\pi^{\mathcal{W}}} \mathbb{E}_{s \sim \mathfrak{B}} [Q^{\mathcal{W}}(s, \pi^{\mathcal{W}}(s))]$
 - 7: **if** access to online environment \mathcal{M} **then**
 - 8: Sample latent-noise action $w_t \sim \pi^{\mathcal{W}}(s_t)$ and compute $a_t \leftarrow \pi_{\text{dp}}^{\mathcal{W}}(s_t, w_t)$
 - 9: Play a_t in \mathcal{M} , observe r_t and next state s_{t+1} , and add (s_t, a_t, r_t, s_{t+1}) to \mathfrak{B}
-

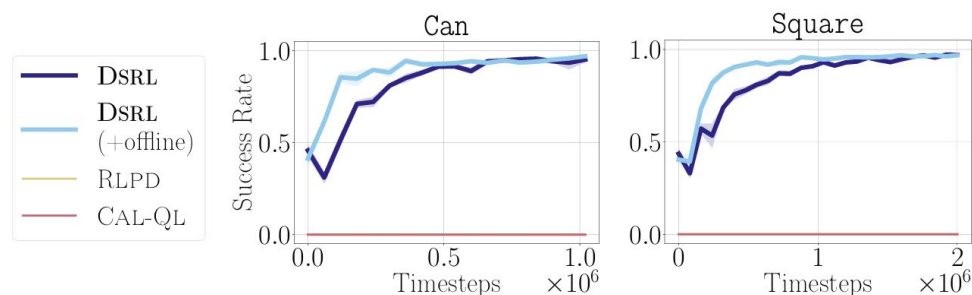
Online adaptation



Offline adaptation

Task	BC (\mathcal{N})	BC (π_{dp})	IQL	REBRAC	IDQL	SRPO	CAC	FAWAC	FBRAC	IFQL	FQL	DSRL
antmaze-large-navigate-singletask	0 \pm 0	0 \pm 0	48 \pm 9	91 \pm 10	0 \pm 0	0 \pm 0	42 \pm 7	1 \pm 1	70 \pm 20	24 \pm 17	80 \pm 8	40 \pm 29
antmaze-giant-navigate-singletask	0 \pm 0	0 \pm 0	0 \pm 0	27 \pm 22	0 \pm 0	0 \pm 0	0 \pm 0	0 \pm 0	0 \pm 1	0 \pm 0	4 \pm 5	0 \pm 0
humanoidmaze-medium-navigate-singletask	1 \pm 0	0 \pm 0	32 \pm 7	16 \pm 9	1 \pm 1	0 \pm 0	38 \pm 19	6 \pm 2	25 \pm 8	69 \pm 19	19 \pm 12	34 \pm 20
humanoidmaze-large-navigate-singletask	0 \pm 0	0 \pm 0	3 \pm 1	2 \pm 1	0 \pm 0	0 \pm 0	1 \pm 1	0 \pm 0	0 \pm 1	6 \pm 2	7 \pm 6	10 \pm 12
antsoccer-arena-navigate-singletask	1 \pm 0	0 \pm 0	3 \pm 2	0 \pm 0	0 \pm 1	0 \pm 0	0 \pm 0	12 \pm 3	24 \pm 4	16 \pm 9	39 \pm 6	28 \pm 9
cube-single-play-singletask	3 \pm 1	3 \pm 3	85 \pm 8	92 \pm 4	96 \pm 2	82 \pm 16	80 \pm 30	81 \pm 9	83 \pm 13	73 \pm 3	97 \pm 2	93 \pm 14
cube-double-play-singletask	0 \pm 0	0 \pm 0	1 \pm 1	7 \pm 3	16 \pm 10	0 \pm 0	2 \pm 2	2 \pm 1	22 \pm 12	9 \pm 5	36 \pm 6	53 \pm 14
scene-play-singletask	1 \pm 1	0 \pm 0	12 \pm 3	50 \pm 13	33 \pm 14	2 \pm 2	50 \pm 40	18 \pm 8	46 \pm 10	0 \pm 0	76 \pm 9	88 \pm 9
puzzle-3x3-play-singletask	1 \pm 1	0 \pm 0	2 \pm 1	2 \pm 1	0 \pm 0	0 \pm 0	0 \pm 0	1 \pm 1	2 \pm 2	0 \pm 0	16 \pm 5	0 \pm 0
puzzle-4x4-play-singletask	0 \pm 0	0 \pm 0	5 \pm 2	10 \pm 3	26 \pm 6	7 \pm 4	1 \pm 1	0 \pm 0	5 \pm 1	21 \pm 11	11 \pm 3	37 \pm 13

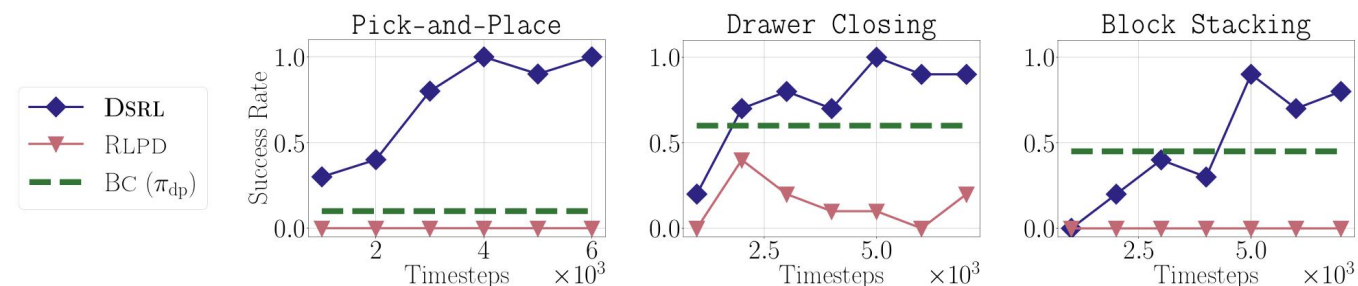
Offline-to-online adaptation



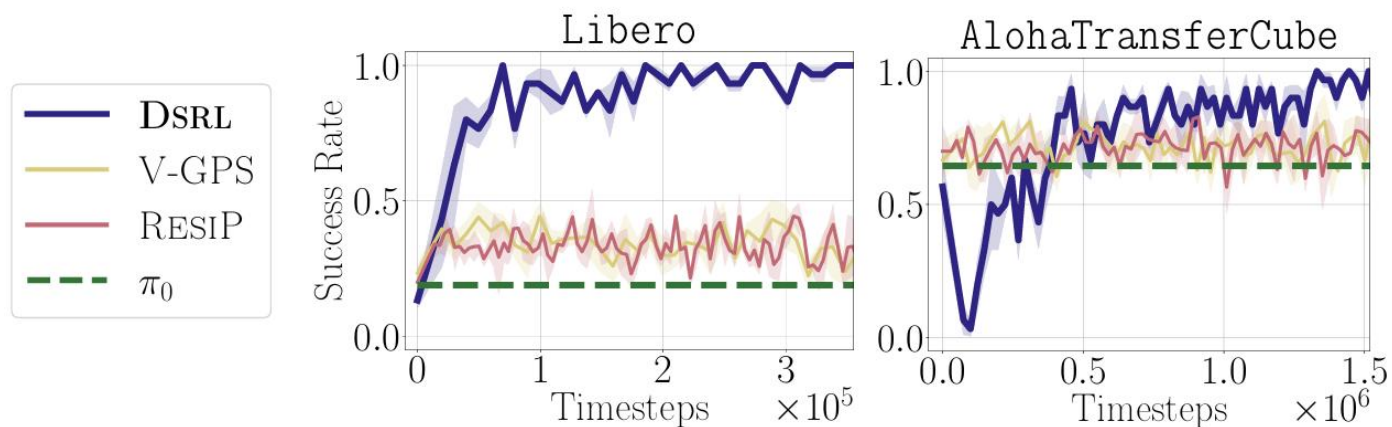
Single-task adaptation in real-world

π_{dp}	RLPD	RLPD + interventions	DSRL
2/10	0/10	0/10	9/10

Multi-task adaptation in real-world

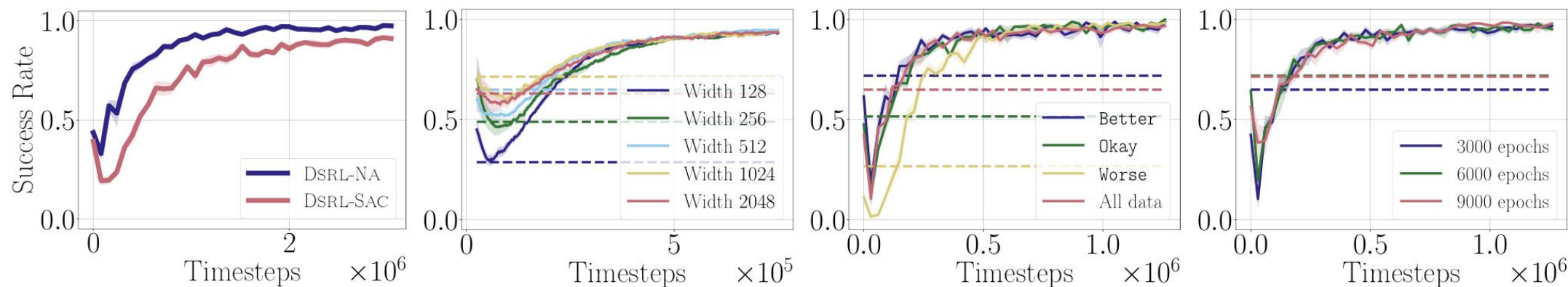


VLA adaptation



Task	π_0	DSRL
Turn on toaster	5/20	18/20
Put spoon on plate	15/20	19/20

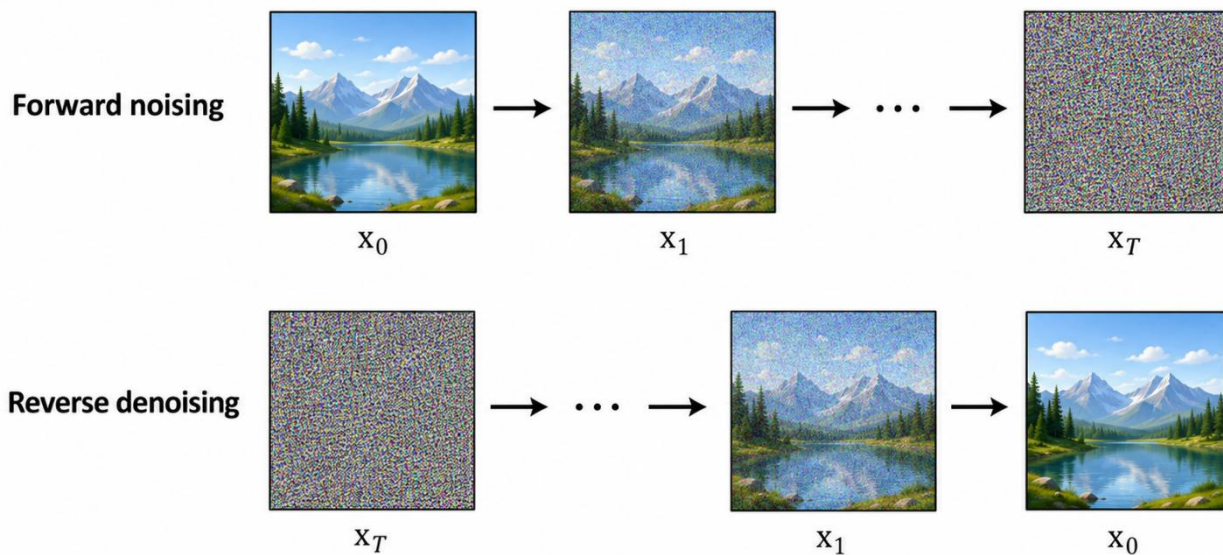
Ablation study



Flow Q-Learning

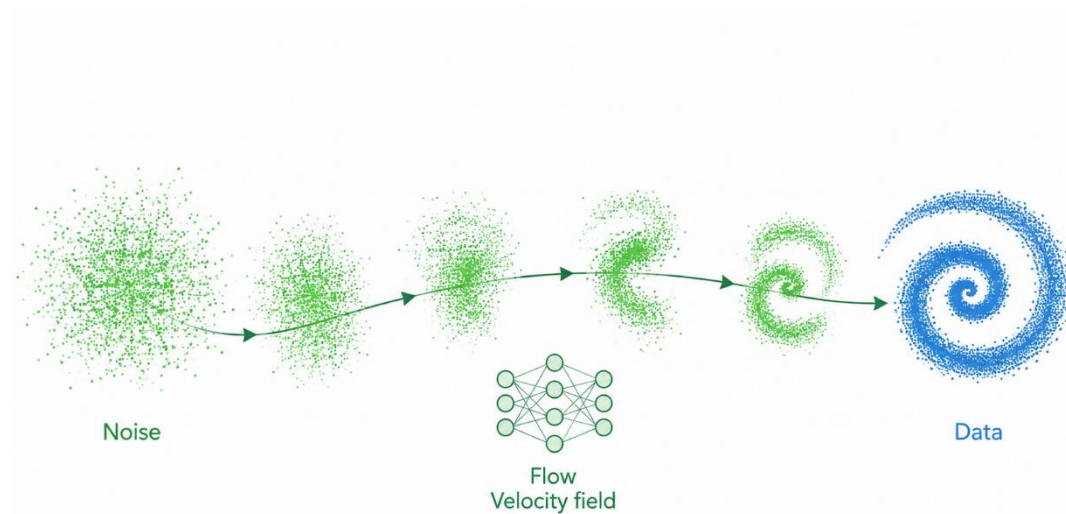
Seohong Park¹ Qiyang Li¹ Sergey Levine¹

Diffusion Model



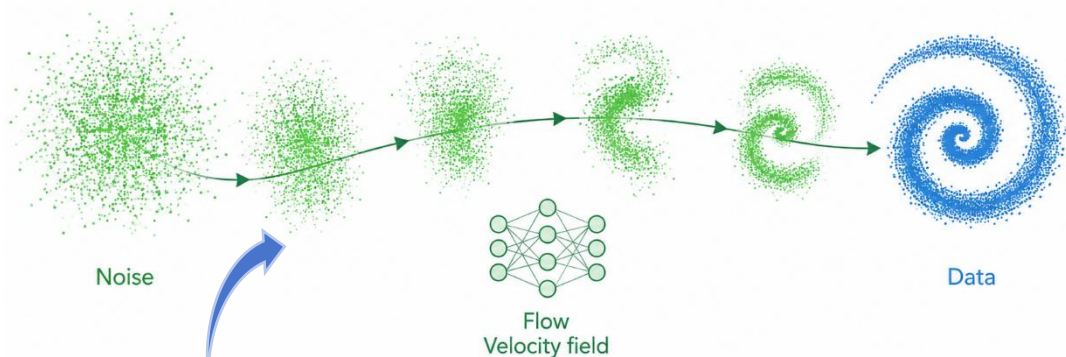
Learning objective: **reverse denoising process**

Flow Matching



Learning objective:
continuous flow from noise to data (velocity field)

Flow Matching



distribution at time t

$$\frac{dx(t)}{dt} = v_{\theta}(x(t), t)$$

velocity field

$$x_0 \sim p_0$$

$$x_1 \sim p_{\text{data}}$$

$$x_1 = x_0 + \int_0^1 v_{\theta}(t, x_t) dt$$

Intermediate distribution $\{p_t(x)\}_{t \in [0,1]}$

Velocity field $\frac{dx_t}{dt} = u_t(x_t)$

Continuity condition $\frac{\partial p_t(x)}{\partial t} + \nabla \cdot (p_t(x)u_t(x)) = 0$

Approximation $v_{\theta}(t, x) \approx u_t(x)$

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t \sim U[0,1], x \sim p_t} [\|v_{\theta}(t, x) - u_t(x)\|^2]$$

Conditional flow matching

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, x_0, x_1} [\|v_{\theta}(t, x_t) - u_t(x_t | x_0, x_1)\|^2]$$

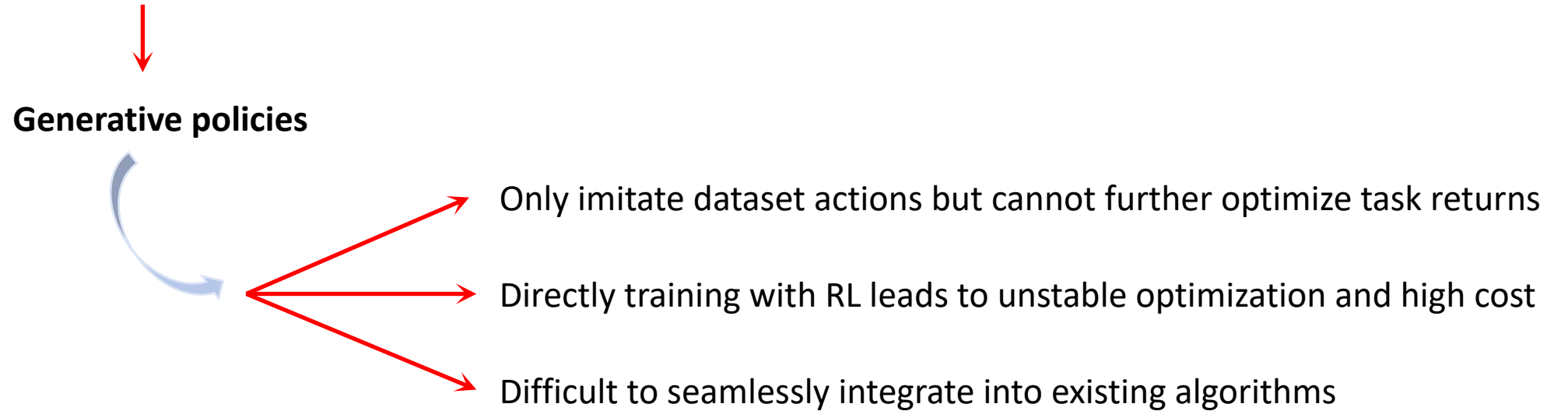
Rectified flow

$$x_t = (1 - t)x_0 + tx_1$$

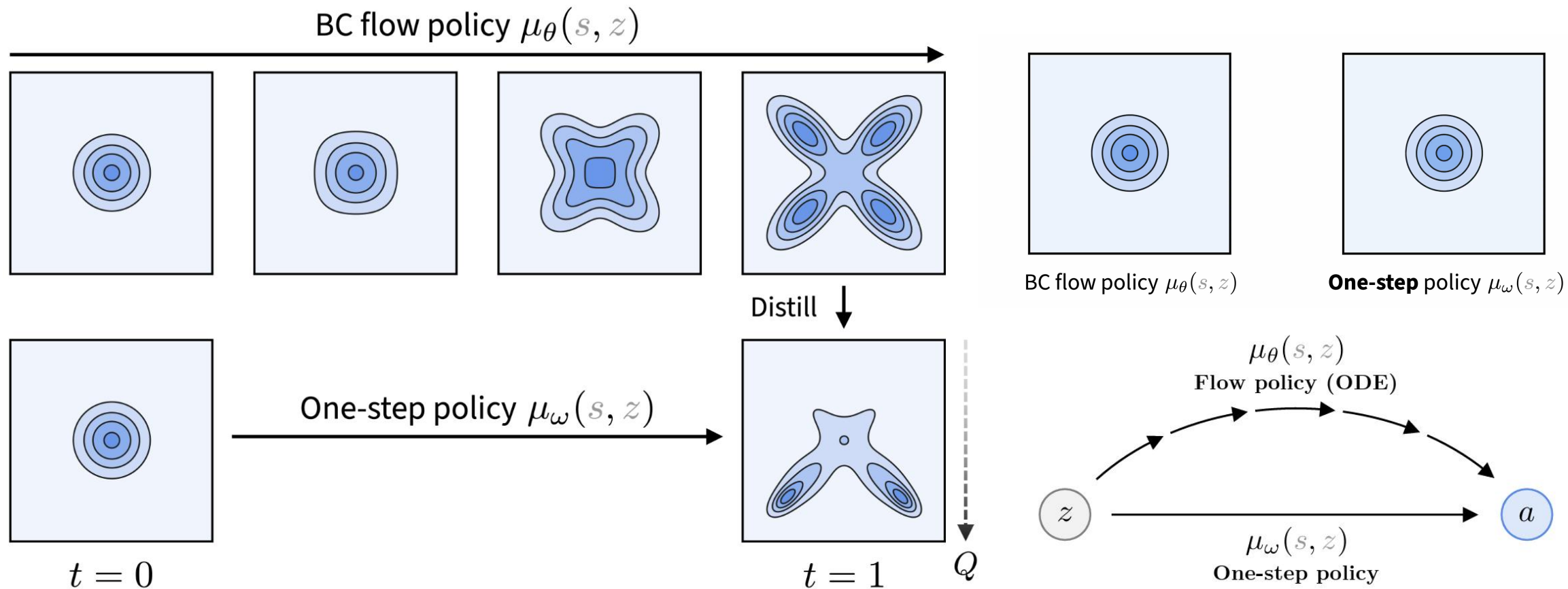
$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, x_0, x_1} [\|v_{\theta}(t, (1 - t)x_0 + tx_1) - (x_1 - x_0)\|^2]$$

Problems in offline RL

- Limited expressiveness of Gaussian policies



How to train an expressive policy that can be integrated with existing RL algorithms?



Use flow matching to model **complex action distributions**, and use Q-learning to guide a one-step policy toward **higher returns**.

Algorithm 1 Flow Q-Learning (FQL)

```

function  $\mu_\theta(s, z)$  ▷ BC flow policy
  for  $t = 0, 1, \dots, M - 1$  do
     $z \leftarrow z + v_\theta(t/M, s, z)/M$  ▷ Euler method
  return  $z$ 

```

while not converged **do**

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

▷ Train critic Q_ϕ

$z \sim \mathcal{N}(0, I_d)$

$a' \leftarrow \mu_\omega(s', z)$

Update ϕ to minimize $\mathbb{E}[(Q_\phi(s, a) - r - \gamma Q_\phi(s', a'))^2]$

▷ Train vector field v_θ in BC flow policy π_θ

$x^0 \sim \mathcal{N}(0, I_d)$

$x^1 \leftarrow a$

$t \sim \text{Unif}([0, 1])$

$x^t \leftarrow (1 - t)x^0 + tx^1$

Update θ to minimize $\mathbb{E}[\|v_\theta(t, s, x^t) - (x^1 - x^0)\|_2^2]$

▷ Train one-step policy π_ω

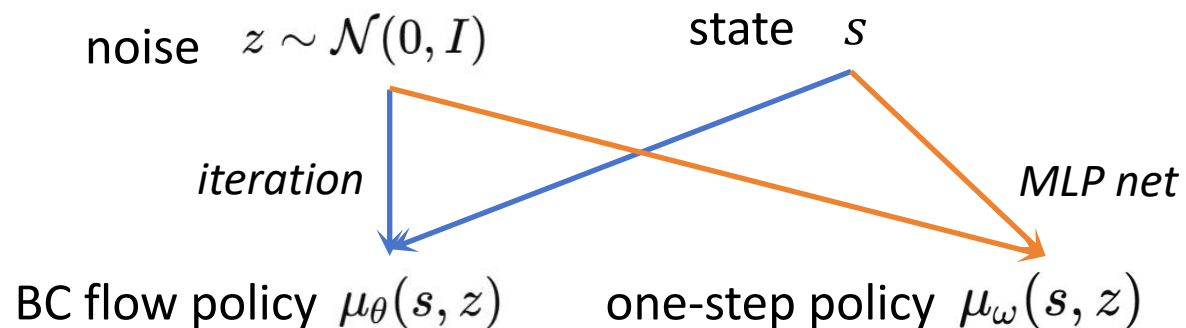
$z \sim \mathcal{N}(0, I_d)$

$a^\pi \leftarrow \mu_\omega(s, z)$

Update ω to minimize $\mathbb{E}[-Q_\phi(s, a^\pi) + \alpha \|a^\pi - \mu_\theta(s, z)\|_2^2]$

return One-step policy π_ω

$$\mathcal{L}_{\text{Distill}}(\omega) = \mathbb{E}_{\substack{s \sim \mathcal{D}, \\ z \sim \mathcal{N}(0, I_d)}} [\|\mu_\omega(s, z) - \mu_\theta(s, z)\|_2^2]$$

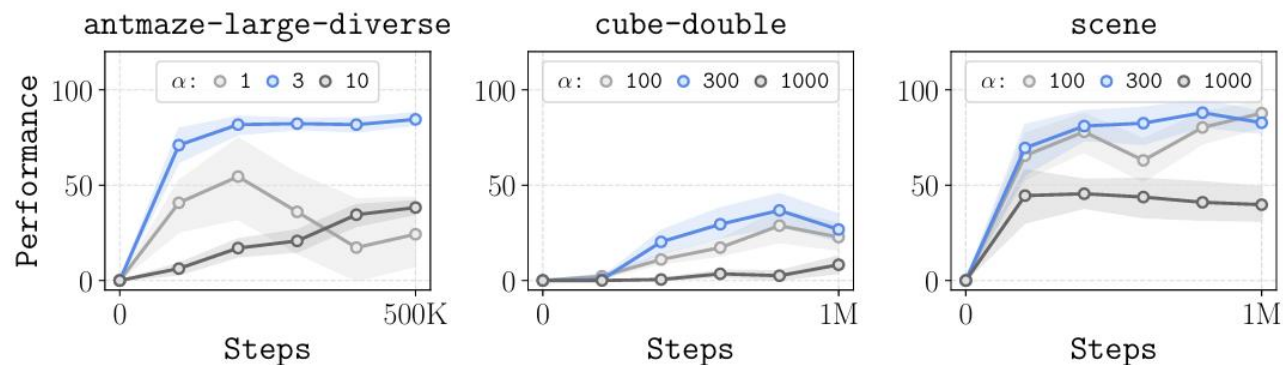


$$\begin{aligned} \mathcal{L}_{\text{Distill}}(\omega) &= \mathbb{E}_{\substack{s \sim \mathcal{D}, \\ z \sim \mathcal{N}(0, I_d)}} [\|\mu_\omega(s, z) - \mu_\theta(s, z)\|_2^2] \\ &\geq \mathbb{E}_{s \sim \mathcal{D}} \left[\inf_{\lambda \in \Lambda(\pi_\omega, \pi_\theta)} \mathbb{E}_{x, y \sim \lambda} [\|x - y\|_2^2] \right] \\ &= \mathbb{E}_{s \sim \mathcal{D}} [W_2(\pi_\omega, \pi_\theta)^2] \end{aligned}$$

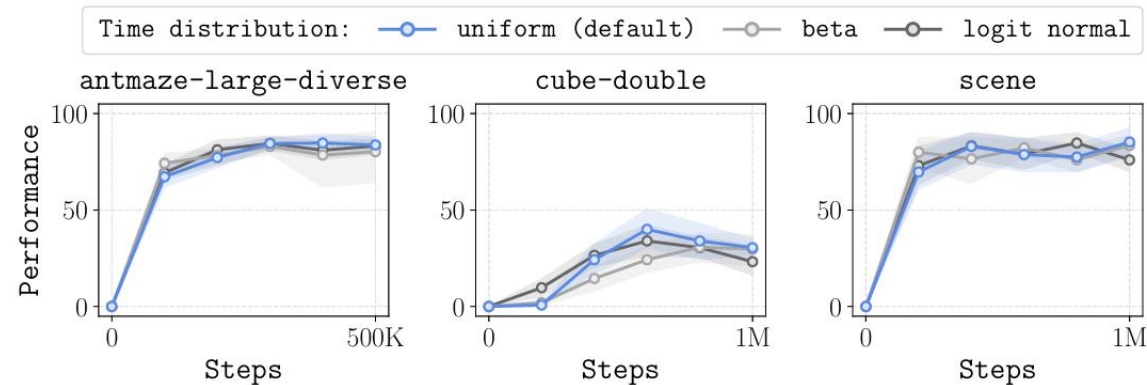
Offline performance

Task Category	Gaussian Policies			Diffusion Policies			Flow Policies			
	BC	IQL	ReBRAC	IDQL	SRPO	CAC	FAWAC	FBRAC	IFQL	FQL
OGBench antmaze-large-singletask (5 tasks)	11 \pm 1	53 \pm 3	81 \pm 5	21 \pm 5	11 \pm 4	33 \pm 4	6 \pm 1	60 \pm 6	28 \pm 5	79 \pm 3
OGBench antmaze-giant-singletask (5 tasks)	0 \pm 0	4 \pm 1	26 \pm 8	0 \pm 0	0 \pm 0	0 \pm 0	0 \pm 0	4 \pm 4	3 \pm 2	9 \pm 6
OGBench humanoidmaze-medium-singletask (5 tasks)	2 \pm 1	33 \pm 2	22 \pm 8	1 \pm 0	1 \pm 1	53 \pm 8	19 \pm 1	38 \pm 5	60 \pm 14	58 \pm 5
OGBench humanoidmaze-large-singletask (5 tasks)	1 \pm 0	2 \pm 1	2 \pm 1	1 \pm 0	0 \pm 0	0 \pm 0	0 \pm 0	2 \pm 0	11 \pm 2	4 \pm 2
OGBench antsoccer-arena-singletask (5 tasks)	1 \pm 0	8 \pm 2	0 \pm 0	12 \pm 4	1 \pm 0	2 \pm 4	12 \pm 0	16 \pm 1	33 \pm 6	60 \pm 2
OGBench cube-single-singletask (5 tasks)	5 \pm 1	83 \pm 3	91 \pm 2	95 \pm 2	80 \pm 5	85 \pm 9	81 \pm 4	79 \pm 7	79 \pm 2	96 \pm 1
OGBench cube-double-singletask (5 tasks)	2 \pm 1	7 \pm 1	12 \pm 1	15 \pm 6	2 \pm 1	6 \pm 2	5 \pm 2	15 \pm 3	14 \pm 3	29 \pm 2
OGBench scene-singletask (5 tasks)	5 \pm 1	28 \pm 1	41 \pm 3	46 \pm 3	20 \pm 1	40 \pm 7	30 \pm 3	45 \pm 5	30 \pm 3	56 \pm 2
OGBench puzzle-3x3-singletask (5 tasks)	2 \pm 0	9 \pm 1	21 \pm 1	10 \pm 2	18 \pm 1	19 \pm 0	6 \pm 2	14 \pm 4	19 \pm 1	30 \pm 1
OGBench puzzle-4x4-singletask (5 tasks)	0 \pm 0	7 \pm 1	14 \pm 1	29 \pm 3	10 \pm 3	15 \pm 3	1 \pm 0	13 \pm 1	25 \pm 5	17 \pm 2
D4RL antmaze (6 tasks)	17	57	78	79	74	30 \pm 3	44 \pm 3	64 \pm 7	65 \pm 7	84 \pm 3
D4RL adroit (12 tasks)	48	53	59	52 \pm 1	51 \pm 1	43 \pm 2	48 \pm 1	50 \pm 2	52 \pm 1	52 \pm 1
Visual manipulation (5 tasks)	-	42 \pm 4	60 \pm 2	-	-	-	-	22 \pm 2	50 \pm 5	65 \pm 2

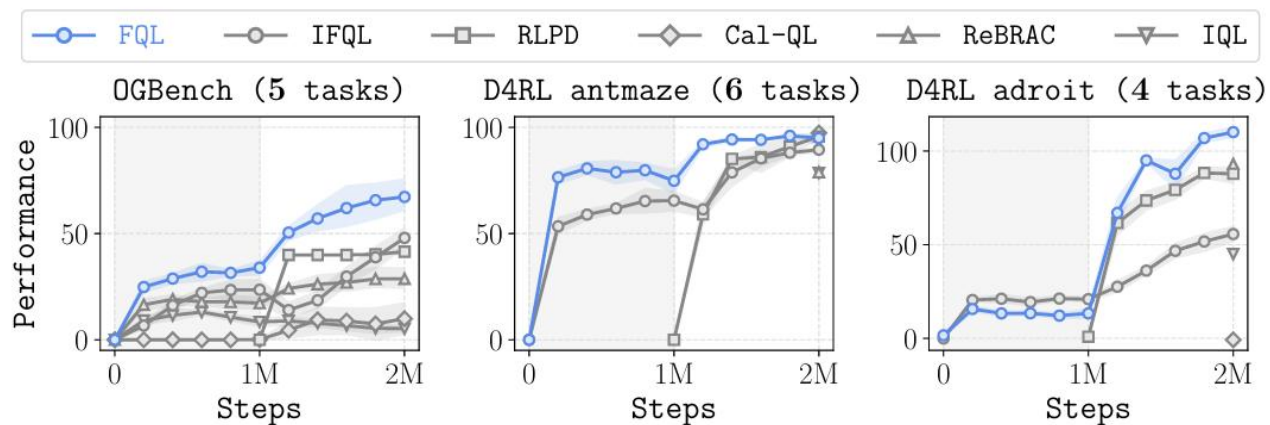
Hyperparameter sensitivity



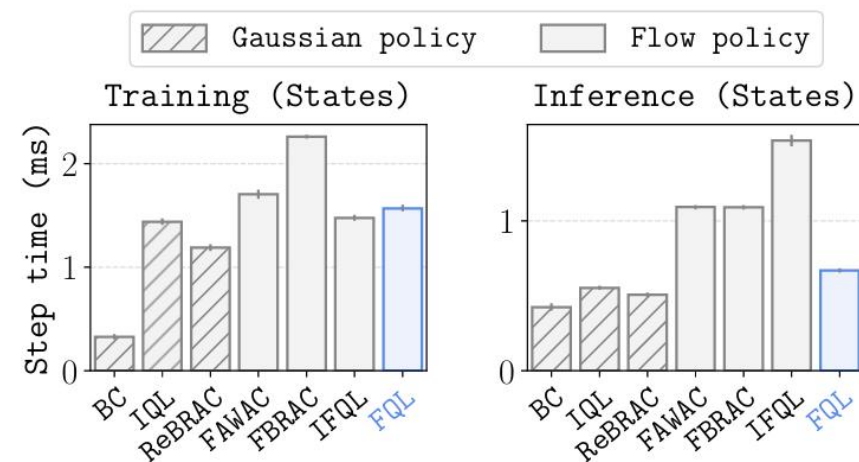
Time distribution



Offline-to-online performance



Time cost



Thanks