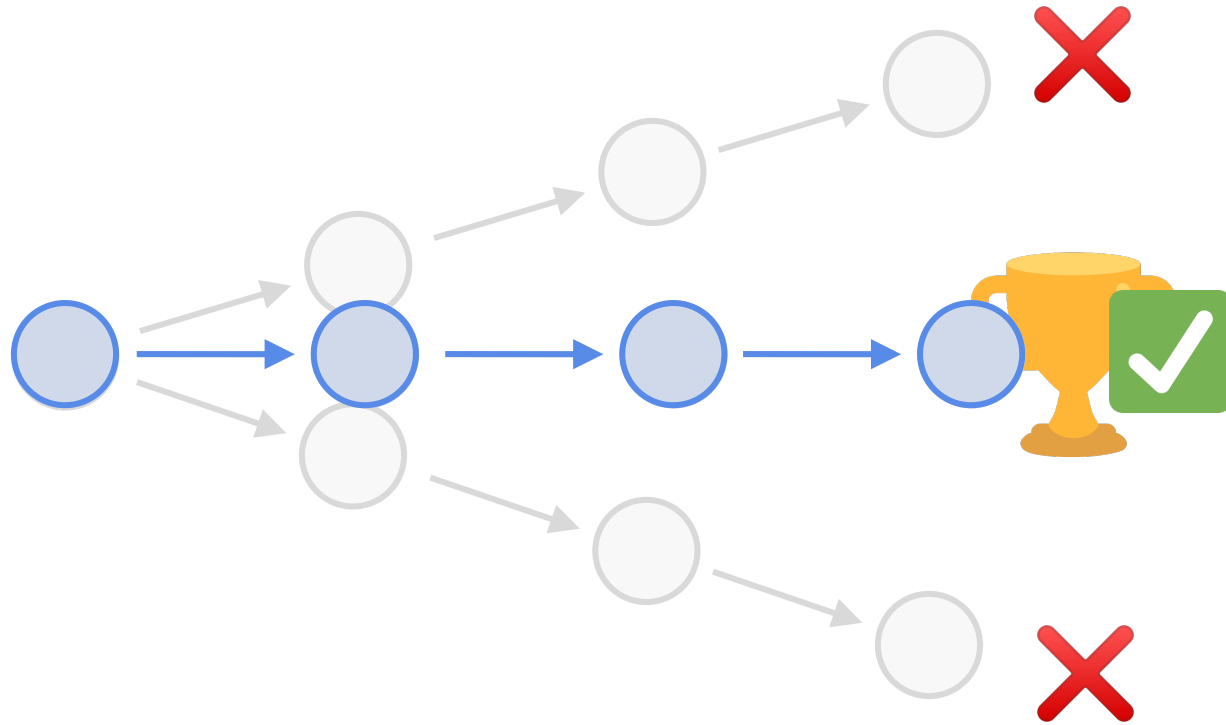


# *Scalable Offline Model-Based RL with Action Chunks*

Kwanyoung Park, Seohong Park, Youngwoon Lee, Sergey Levine

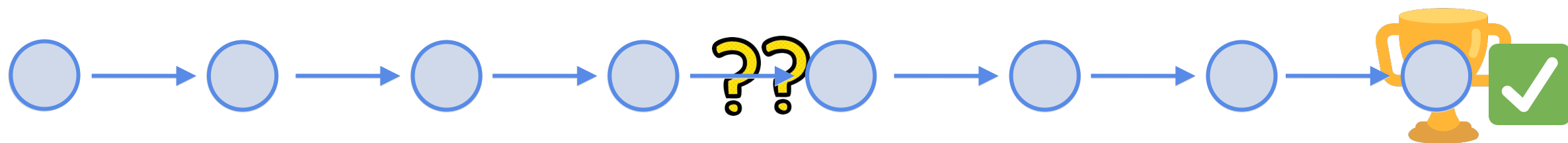
In model-based RL, the agent imagines future trajectories of the policy, efficiently finding a best policy.



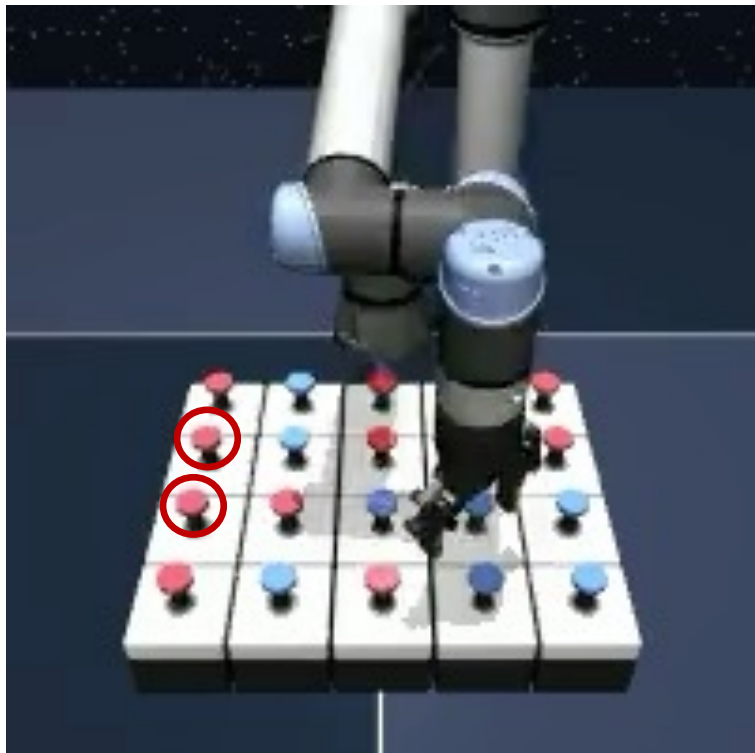
For long-horizon tasks, we need long enough model rollouts to propagate the learning signals.



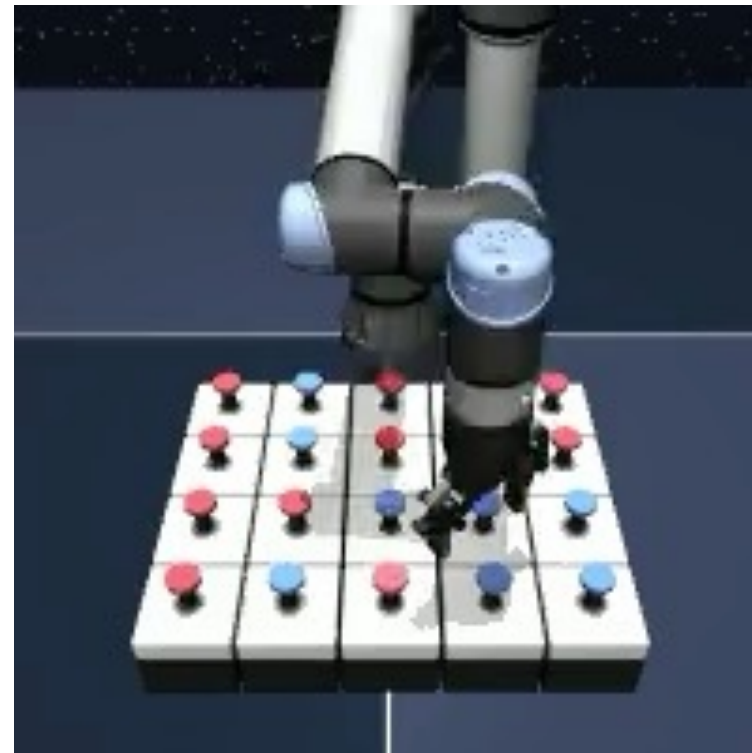
For long-horizon tasks, we need long enough model rollouts to propagate the learning signals.



However, model prediction becomes unstable when we generate long model rollouts.



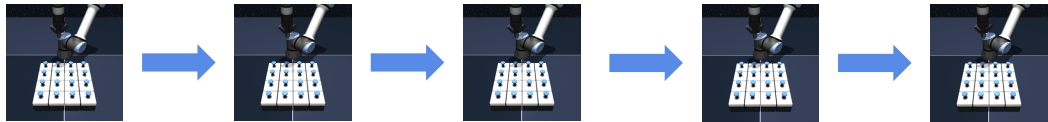
Model prediction



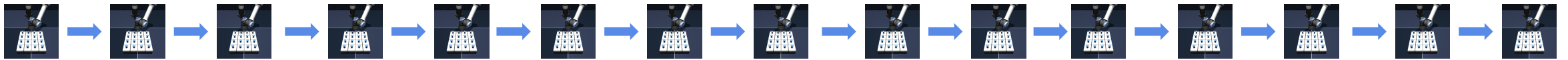
Real

Thus, prior works have been using rollout steps of 5 to 15, which is insufficient for solving long-horizon tasks.

Standard MBRL (MOPO, MOBILE)

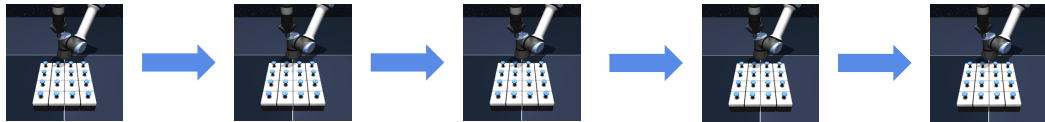


Long-horizon MBRL (LEQ, Dreamer)

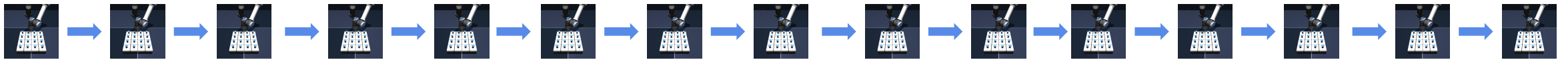


Thus, prior works have been using rollout steps of 5 to 15, which is insufficient for solving long-horizon tasks.

Standard MBRL (MOPO, MOBILE)

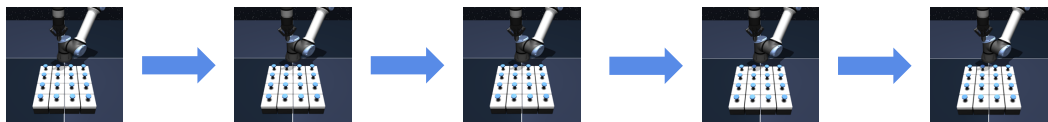


Long-horizon MBRL (LEQ, Dreamer)

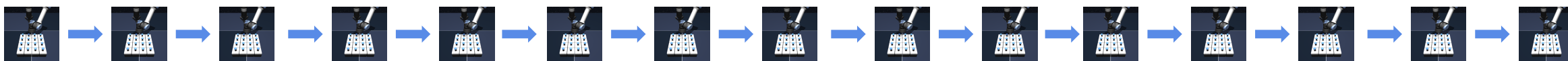


However, our method utilizes **100 steps** of rollouts, pushing the boundary of offline MBRL in long-horizon tasks.

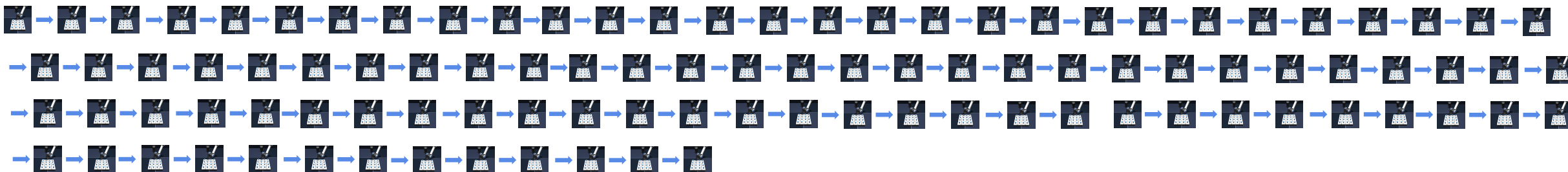
Standard MBRL (MOPO, MOBILE)



Long-horizon MBRL (LEQ, Dreamer)

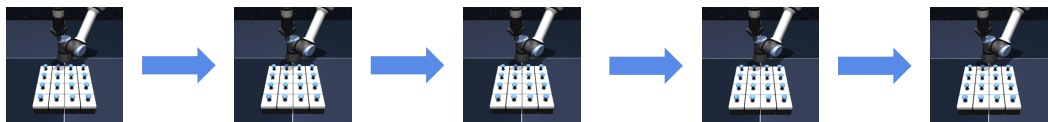


Ours (MAC)



However, our method utilizes **100 steps** of rollouts, pushing the boundary of offline MBRL in long-horizon tasks.

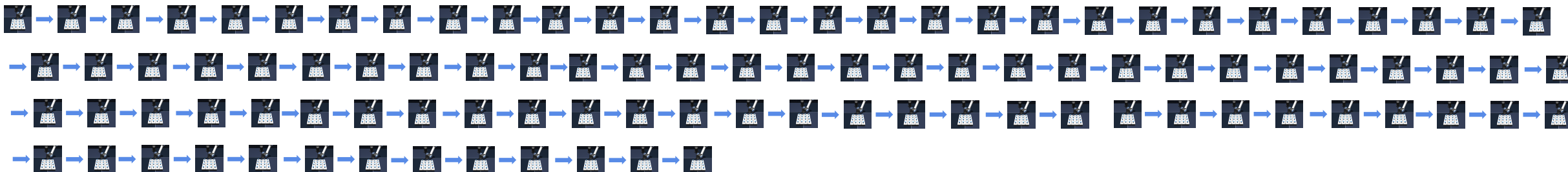
Standard MBRL (MOPO, MOBILE)



Long-horizon MBRL (LEQ, Dreamer)

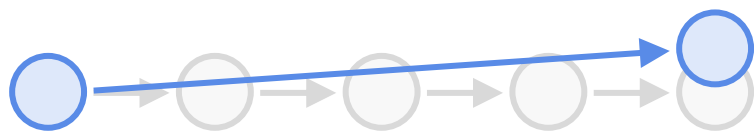


Ours (MAC)

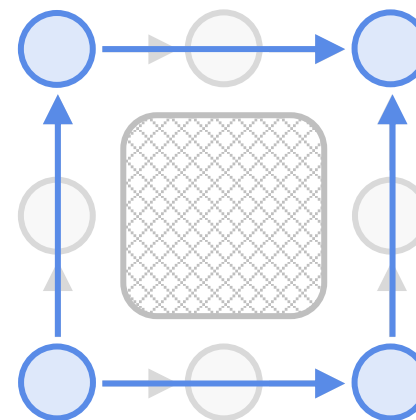


Our core recipes consists of  
1) Action chunking, 2) Flow rejection sampling.

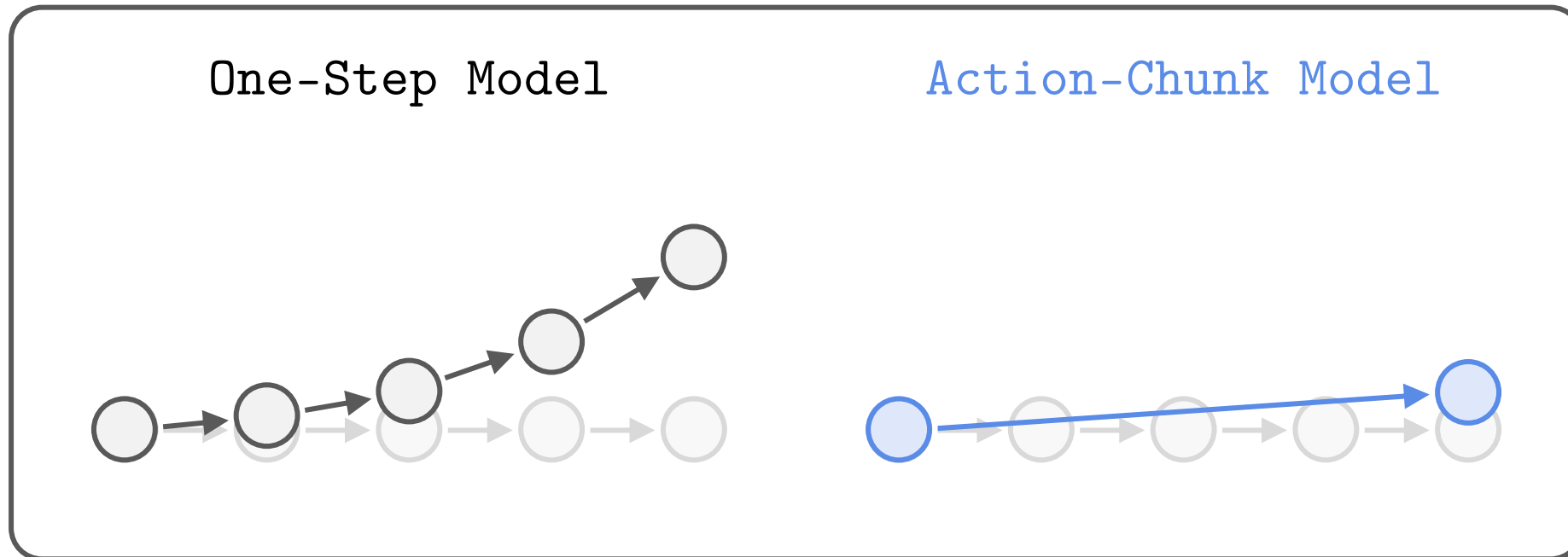
Action-Chunk Model



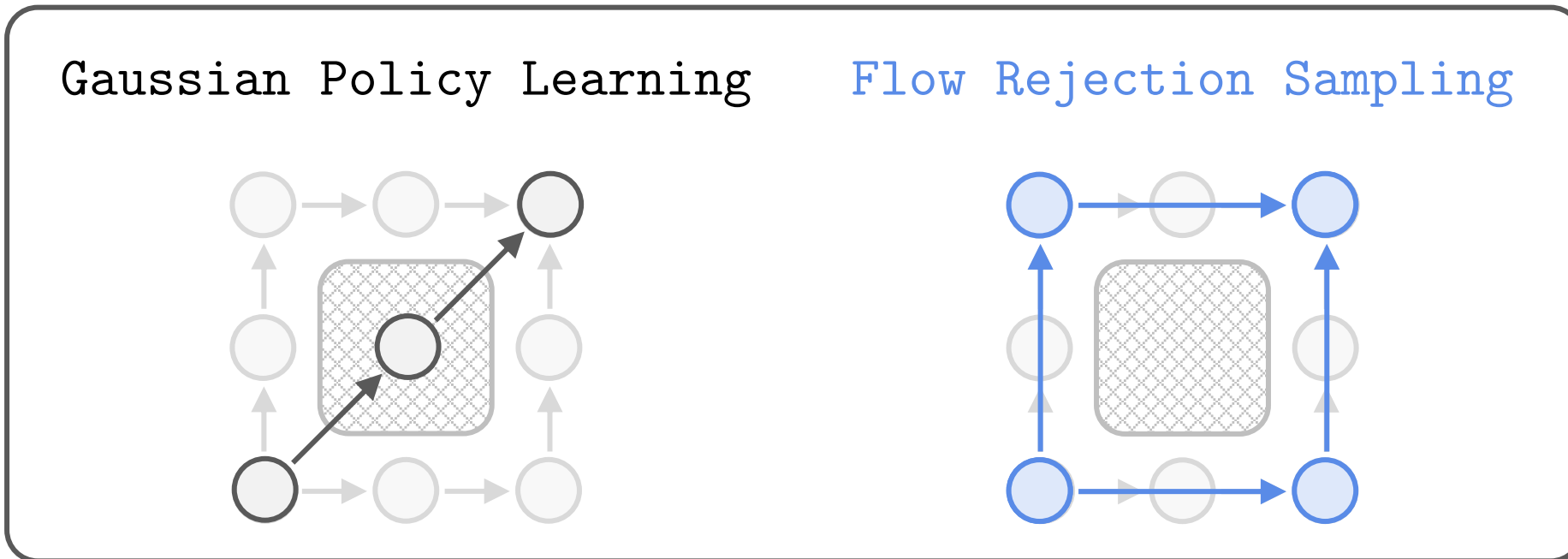
Flow Rejection Sampling



Action chunking reduces the number of model execution, mitigating the compounding error problem.



Rejection sampling from flow BC policy enables modeling multi-modal action distributions, preventing OOD actions.



Specifically, MAC consists of 3 steps:

Step 1: **Action-chunked** flow BC policy and world model

Flow BC policy:  $\pi_{\beta}(a_{t:t+h} | s_t)$

World model:  $p(\hat{r}_{t:t+h}, \hat{s}_{t+h} | s_t, a_{t:t+h})$

Step 2: Policy extraction with **rejection sampling**

$$\hat{a}_{t:t+h} = \max_i Q(s_t, a_{t:t+h}^{(i)}), \text{ where } a_{t:t+h}^{(i)} \sim \pi_{\beta}(\cdot | s_t)$$

Step 3: Q-learning with **model-generated rollouts**

$$Q(s_t, a_t) \leftarrow \gamma^{hn} Q(\hat{s}_{t+nh}, \hat{a}_{t+(n-1)h:t+nh}) + \sum_{i=0}^{n-1} \gamma^{hi} \hat{r}_{t+ih:t+(i+1)h}$$

(On-policy n-step returns with action chunks, n=10, h=10)

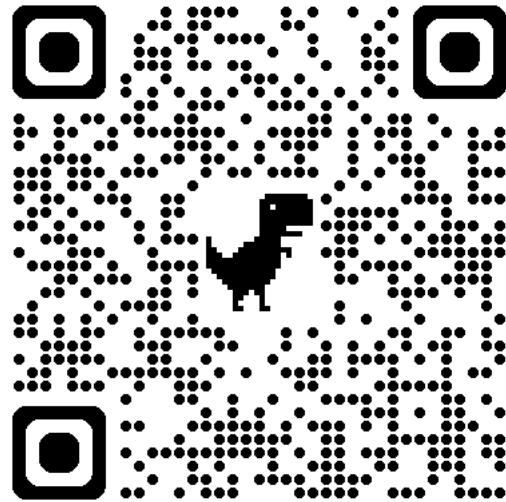
On 100M-scale datasets for long-horizon tasks, our method drastically improves the performance of offline MBRL approaches.

Environment	MOPO	MOBILE	LEQ	FMPC	MAC
humanoidmaze-medium-navigate-oracclerrep-v0	27 $\pm$ 5	23 $\pm$ 3	0 $\pm$ 0	18 $\pm$ 5	36 $\pm$ 2
humanoidmaze-giant-navigate-oracclerrep-v0	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0
cube-double-play-oracclerrep-v0	25 $\pm$ 12	15 $\pm$ 3	0 $\pm$ 0	37 $\pm$ 13	100 $\pm$ 1
cube-octuple-play-oracclerrep-v0	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	30 $\pm$ 6
puzzle-3x3-play-oracclerrep-v0	19 $\pm$ 2	15 $\pm$ 5	1 $\pm$ 1	12 $\pm$ 6	100 $\pm$ 0
puzzle-4x5-play-oracclerrep-v0	0 $\pm$ 0	0 $\pm$ 0	1 $\pm$ 3	0 $\pm$ 0	99 $\pm$ 3

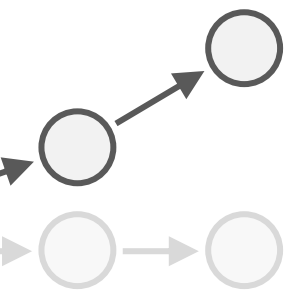
Our method also drastically improves the performance of offline MBRL in standard-scale (1M) datasets.

Environment	MOPO	MOBILE	LEQ	FMPC	MAC
cube-single-play-v0 (5 tasks)	12 $\pm$ 4	81 $\pm$ 8	0 $\pm$ 0	9 $\pm$ 5	99 $\pm$ 2
cube-double-play-v0 (5 tasks)	1 $\pm$ 1	1 $\pm$ 2	0 $\pm$ 0	3 $\pm$ 2	53 $\pm$ 4
scene-play-v0 (5 tasks)	6 $\pm$ 8	8 $\pm$ 4	0 $\pm$ 0	4 $\pm$ 4	97 $\pm$ 4
puzzle-3x3-play-v0 (5 tasks)	20 $\pm$ 0	12 $\pm$ 9	10 $\pm$ 7	1 $\pm$ 1	20 $\pm$ 0
puzzle-4x4-play-v0 (5 tasks)	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	0 $\pm$ 0	78 $\pm$ 13

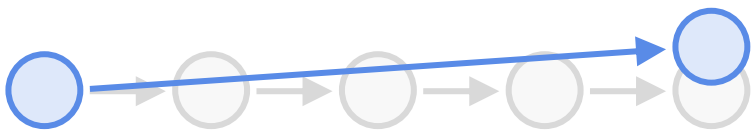
If you are interested in our project, please check the paper!



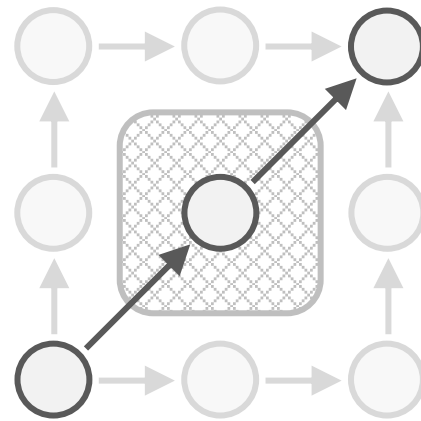
Model



Action-Chunk Model



Gaussian Policy Learning



Flow R

