
Receding-Horizon Execution for Action Chunking in Offline-to-Online Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Action chunking improves offline-to-online reinforcement learning (RL) by using
2 temporally extended actions for exploration and multi-step temporal-difference
3 (TD) backups. However, a chunk length that accelerates value propagation can
4 be too long for open-loop execution, since executing the full predicted chunk
5 before replanning reduces reactivity and degrades performance. We study receding-
6 horizon execution for Q-chunking policies with a fixed prediction horizon. We train
7 a horizon-conditioned execution critic that evaluates prefixes of a predicted chunk
8 under fixed and value-based execution rules. On a long-horizon manipulation
9 task, receding-horizon execution substantially improves policies that fail under full
10 open-loop execution but does not match a fixed short-horizon-trained actor or the
11 best fixed execution horizon. Value-based horizon selection remains biased toward
12 short horizons, and state-dependent selection remains an open problem.

13 1 Introduction

14 Offline-to-online reinforcement learning (RL) uses previously collected data to improve sample
15 efficiency, but remains challenging in long-horizon, sparse-reward manipulation tasks where explo-
16 ration is difficult [Li et al., 2023, Rengarajan et al., 2022]. In imitation learning, action-chunking
17 policies predict multi-step action sequences rather than single actions, which helps capture temporally
18 coherent and potentially non-Markovian behavior in demonstration data [Zhao et al., 2023, Black
19 et al., 2026]. Q-chunking (QC-FQL) [Li et al., 2025] extends action chunking to RL by optimizing
20 the policy and critic over the action sequences. This enables unbiased multi-step temporal-difference
21 backups and temporally coherent exploration, leading to strong performance on long-horizon tasks.

22 However, executing a larger portion of a predicted chunk before replanning reduces the frequency
23 with which new observations influence control [Zhao et al., 2023, Liu et al., 2025, So et al., 2026].
24 In Q-chunking, the chunk length also determines the backup horizon used by the critic. Thus, a
25 horizon that improves value propagation can exceed the horizon that remains reliable under open-loop
26 execution [Li et al., 2025, 2026, Song et al., 2026].

27 Recent methods address this trade-off in different ways. SEAR [Nagy et al., 2026] combines receding-
28 horizon execution with multi-horizon critic targets and random replanning, but its evaluation policy
29 still uses a fixed replanning interval. DQC [Li et al., 2026] instead decouples the critic chunk length
30 from the policy chunk length, so that the critic can use long-horizon backups while the policy predicts
31 shorter chunks. However, whether the execution horizon can be selected adaptively at run time rather
32 than fixed remains an open question.

33 We study this question in the offline-to-online QC-FQL setting [Li et al., 2025]. We use multi-horizon
34 targets and random replanning from SEAR [Nagy et al., 2026], with a selector-consistent bootstrap
35 that uses the same selector at training and evaluation. We use this critic to study state-dependent
36 execution-horizon selection rather than a fixed replanning interval. Empirically, receding-horizon

37 execution substantially improves a long prediction-horizon policy that fails under full open-loop
 38 execution, but does not reach the best short-horizon-trained actor. Value-based horizon selection
 39 remains biased toward short executions, and a simple marginal selector that compares neighboring
 40 horizons does not consistently outperform the best fixed horizon. Our results suggest that multi-
 41 horizon critics partially recover long-horizon Q-chunking actors, while reliable state-dependent
 42 horizon selection remains open.

43 2 Preliminaries

44 We consider the offline-to-online QC-FQL setting [Li et al., 2025]. The actor μ_ψ with parameters ψ
 45 is a flow-matching policy that maps a state s_t and a noise variable z_t to an action chunk,

$$A_t = \mu_\psi(s_t, z_t) \in \mathbb{R}^{H \times d_a}, \quad z_t \sim \mathcal{N}(0, I), \quad (1)$$

46 where H is the prediction horizon and d_a is the action dimension. Following QC-FQL, the actor is
 47 trained by distilling a flow-matching behavior prior learned from the offline dataset, together with a
 48 Q-maximization term using the chunked critic.

49 Standard Q-chunking executes all H actions of A_t before querying the actor again, so the prediction
 50 horizon and the execution horizon coincide. We instead distinguish the prediction horizon H from the
 51 execution horizon $k \leq H$. Given a predicted action chunk A_t , the agent executes the first k actions,
 52 observes the state s_{t+k} , and replans by sampling a new chunk from the actor. In our experiments,
 53 $H = 50$ and we consider a sparse set of candidate execution horizons,

$$\mathcal{K} = \{5, 10, 15, 20, 25, 50\}. \quad (2)$$

54 For the $H = 50$ actor, full open-loop execution corresponds to $k = 50$, while receding-horizon
 55 execution uses $k < 50$.

56 3 Method

57 3.1 Multi-horizon execution critic

58 Inspired by SEAR [Nagy et al., 2026], we train a horizon-conditioned critic $Q_\phi(s_t, A_t, k)$ that
 59 estimates the return of executing the first k actions of A_t and then replanning. Unlike SEAR, our
 60 bootstrap is selector-consistent and our critic is an MLP masking actions after step k rather than a
 61 causal transformer. Q_ϕ is the ensemble-min over $M = 2$ critics (similarly $Q_{\bar{\phi}}$ for the target).

62 For the bootstrap value at the next replanning state, we sample $A' \sim \mu_\psi(\cdot | s_{t+k})$ from the current
 63 actor and apply the Q-min greedy selector (Section 3.2) at training time. The bootstrap value is

$$V(s_{t+k}) = Q_{\bar{\phi}}(s_{t+k}, A', k'), \quad (3)$$

64 where k' is the selected horizon and $\bar{\phi}$ denotes the target critic. We do not maintain a target actor.
 65 The multi-horizon target for horizon k is

$$G_t^{(k)} = \sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}), \quad (4)$$

66 and the critic is trained by averaging the squared temporal-difference error over candidate horizons,

$$\mathcal{L}_Q(\phi) = \mathbb{E} \left[\frac{1}{|\mathcal{K}|} \sum_{k \in \mathcal{K}} (Q_\phi(s_t, A_t, k) - G_t^{(k)})^2 \right]. \quad (5)$$

67 The actor follows the QC-FQL objective, with the Q-maximization term using $Q_\phi(s, \mu_\psi(s, z), k_a)$
 68 for k_a sampled uniformly from \mathcal{K} at each gradient step. Algorithm 1 summarizes the procedure.

69 3.2 Selectors

70 Fixed-horizon evaluation chooses a constant k throughout the rollout and provides a strong reference
 71 within the same actor family. The Q-min greedy selector chooses

$$k(s, A) = \arg \max_{k' \in \mathcal{K}} Q_\phi(s, A, k'). \quad (6)$$

Algorithm 1 Adaptive-horizon QC-FQL training

```
1: Input: offline dataset  $\mathcal{D}_{\text{off}}$ , candidate horizons  $\mathcal{K}$ 
2: Initialize actor  $\mu_\psi$ , critic  $Q_\phi$ , target critic  $Q_{\bar{\phi}}$ 
3: for phase  $\in \{\text{offline, online}\}$  do
4:   if phase is online then
5:      $\mathcal{D} \leftarrow \mathcal{D}_{\text{off}} \cup \mathcal{D}_{\text{buf}}$ 
6:   else  $\mathcal{D} \leftarrow \mathcal{D}_{\text{off}}$ 
7:   end if
8:   for each gradient step do
9:     Sample  $(s_t, A_t, r_{t:t+H-1}, \{s_{t+k}\}_{k \in \mathcal{K}})$  from  $\mathcal{D}$ 
10:    Construct  $G_t^{(k)}$  for each  $k \in \mathcal{K}$  and update critic on  $\mathcal{L}_Q(\phi)$ 
11:    Update actor with QC-FQL objective using  $k_a \sim \text{Unif}(\mathcal{K})$ 
12:   end for
13:   if phase is online then
14:     Collect a rollout with  $k \sim \text{Unif}(\mathcal{K})$  at each replanning step, appending to  $\mathcal{D}_{\text{buf}}$ 
15:   end if
16: end for
```

72 We also consider a marginal selector that compares neighboring horizons. Let k_i and k_{i-1} be
73 consecutive elements of \mathcal{K} . The marginal improvement of extending from k_{i-1} to k_i is

$$\Delta_i(s, A) = Q_\phi(s, A, k_i) - Q_\phi(s, A, k_{i-1}). \quad (7)$$

74 The forward variant starts from the smallest candidate horizon $k = 5$ and extends to the next horizon
75 as long as $\Delta_i \geq 0$, returning the previous horizon at the first negative gap. The reverse variant starts
76 from $k = 50$ and shrinks to the previous horizon as long as $\Delta_i \leq 0$. We use the forward and reverse
77 variants to test how the scan order affects the selected horizon.

78 4 Experiments

79 We evaluate on cube-triple task3 from OGBench [Park et al., 2025]. We first compare actors trained
80 with different prediction horizons, and then freeze the final $H = 50$ multi-horizon checkpoint to
81 compare execution rules with the same actor and critic. This separates prediction-horizon training
82 from evaluation-time execution choices.

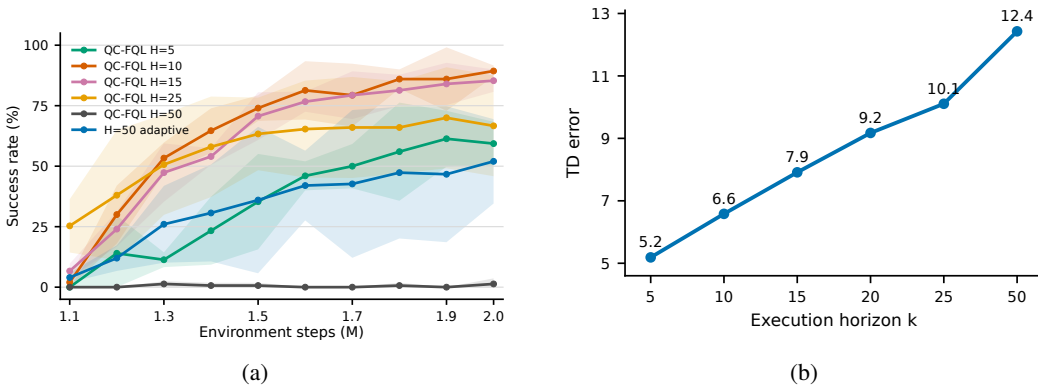


Figure 1: (a) Online evaluation curves between 1.1M and 2M environment steps for fixed- H QC-FQL actors and the $H = 50$ multi-horizon actor evaluated with adaptive execution (over 3 seeds). (b) TD error $|Q_\phi(s_t, A_t, k) - G_t^{(k)}|$, averaged over rollout states.

83 **Online learning across prediction horizons.** The $H = 50$ actor under full open-loop execution
84 stays near zero throughout online training (Figure 1a). The $H = 10$ fixed actor is the strongest
85 baseline at about 89%, while our $H = 50$ adaptive policy reaches about 52%. Receding-horizon

86 execution improves the long-horizon actor relative to full open-loop execution, but does not match
 87 the best short-horizon-trained actor.

88 **Horizon-dependent critic error.** To test whether short-horizon selection reflects only execution
 89 reliability, we compute a post-hoc bootstrap TD error on fresh rollouts from the final $H = 50$
 90 checkpoint, measuring the absolute difference between $Q_\phi(s_t, A_t, k)$ and the bootstrap target $G_t^{(k)}$
 91 in Eq. 5. The TD error increases monotonically with k and more than doubles from $k = 5$ to
 92 $k = 50$ (Figure 1b). This rollout-based diagnostic is not replay-buffer validation, but it suggests that
 93 long-horizon targets are harder to fit and may confound an interpretation of short-horizon selection as
 94 a reliable execution signal.

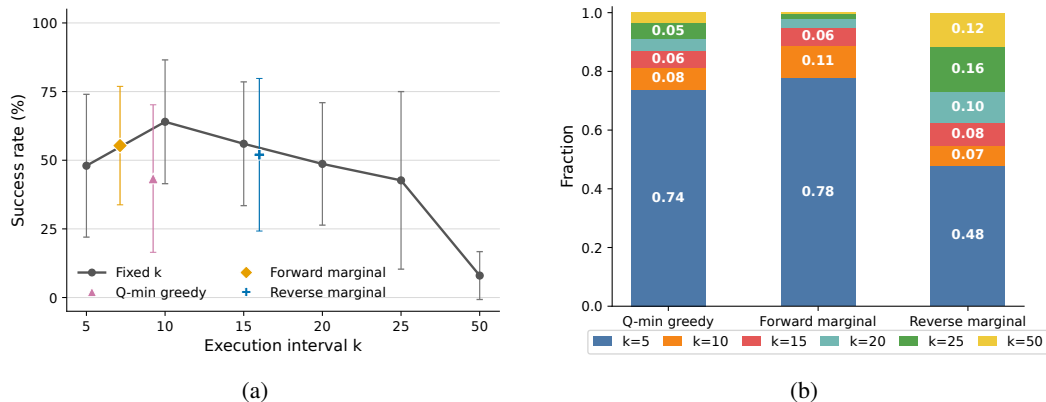


Figure 2: Selector evaluation at the final $H = 50$ multi-horizon checkpoint (mean over 3 seeds, 50 episodes per seed). (a) Success rate under fixed execution horizons (line) and adaptive selectors (markers). (b) Distribution of selected k for each adaptive selector.

95 **Selector evaluation under frozen actor and critic.** We freeze the final $H = 50$ multi-horizon
 96 checkpoint and reevaluate it under fixed- k execution and three adaptive selectors (Q-min greedy,
 97 forward marginal, reverse marginal), holding the actor and critic fixed. The best fixed horizon is
 98 $k = 10$ at about 64%, while $k = 50$ collapses to about 8% (Figure 2a), showing that execution
 99 length critically determines whether the same actor succeeds or fails. Adaptive selectors improve
 100 over $k = 50$ but consistently fall below the best fixed $k = 10$, with the marginal selectors closer to
 101 $k = 10$ than Q-min greedy.

102 **Selected horizon distribution.** Q-min greedy and forward marginal mostly select $k = 5$. Reverse
 103 marginal selects longer horizons much more often, including $k \geq 20$ (Figure 2b). However, this
 104 longer selection does not improve success and is slightly below forward marginal. The short-horizon
 105 bias of value-based selectors can be partially shifted by changing the comparison rule, but this shift
 106 alone does not close the gap to the best fixed horizon.

107 5 Discussion and Limitations

108 Our results show that the execution horizon has a large effect for action-chunking policies. Receding-
 109 horizon execution improves a long-horizon actor relative to full open-loop execution, but value-based
 110 selectors remain below both the best fixed execution horizon for the same checkpoint and the best
 111 short-horizon-trained actor.

112 The execute-then-replan critic estimates the return of executing k actions and replanning, so its value
 113 reflects both the quality of the current prefix and the benefit of an earlier decision point. We also
 114 observe that the post-hoc TD error grows with k , suggesting that long-horizon targets are harder to fit
 115 and may bias selection toward shorter horizons. Better calibrated horizon-conditioned critics and
 116 losses that balance candidate horizons are natural directions for future work.

117 **References**

- 118 Kevin Black, Manuel Y Galliker, and Sergey Levine. Real-time execution of action chunking flow
119 policies. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2026.
120 URL <https://openreview.net/forum?id=UkR2z05uww>.
- 121 Qiyang Li, Jason Zhang, Dibya Ghosh, Amy Zhang, and Sergey Levine. Accelerating exploration with
122 unlabeled prior data. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine,
123 editors, *Advances in Neural Information Processing Systems*, volume 36, pages 67434–67458. Cur-
124 ran Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper_files/paper/
125 2023/file/d53d51e88d92d3723755f6d425bc513b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/d53d51e88d92d3723755f6d425bc513b-Paper-Conference.pdf).
- 126 Qiyang Li, Zhiyuan Zhou, and Sergey Levine. Reinforcement learning with action chunking. In
127 *Advances in Neural Information Processing Systems*, 2025. URL [https://openreview.net/
128 forum?id=XUks1Y96NR](https://openreview.net/forum?id=XUks1Y96NR).
- 129 Qiyang Li, Seohong Park, and Sergey Levine. Decoupled q-chunking. In *The Fourteenth International
130 Conference on Learning Representations*, 2026. URL [https://openreview.net/forum?id=
131 aqGndZQL91](https://openreview.net/forum?id=aqGndZQL91).
- 132 Yuejiang Liu, Jubayer Ibn Hamid, Annie Xie, Yoonho Lee, Maximilian Du, and Chelsea Finn.
133 Bidirectional decoding: Improving action chunking via guided test-time sampling. In *International
134 Conference on Learning Representations*, 2025.
- 135 C. F. Maximilian Nagy, Onur Celik, Emiliyan Gospodinov, Florian Seligmann, Weiran Liao, Aryan
136 Kaushik, and Gerhard Neumann. Sear: Sample efficient action chunking reinforcement learning,
137 2026.
- 138 Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking
139 offline goal-conditioned rl. In *International Conference on Learning Representations (ICLR)*,
140 2025.
- 141 Desik Rengarajan, Gargi Vaidya, Akshay Sarvesh, Dileep Kalathil, and Srinivas Shakkottai. Rein-
142 forcement learning with sparse rewards using guidance from offline demonstration. In *International
143 Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?id=
144 YJ1WzgMVsmT](https://openreview.net/forum?id=YJ1WzgMVsmT).
- 145 Junhyuk So, Chiwoong Lee, Shinyoung Lee, Jungseul Ok, and Eunhyeok Park. Improving generative
146 behavior cloning via self-guidance and adaptive chunking. In *The Thirty-ninth Annual Conference
147 on Neural Information Processing Systems*, 2026. URL [https://openreview.net/forum?id=
148 GctsZXLCp1](https://openreview.net/forum?id=GctsZXLCp1).
- 149 Gwanwoo Song, Kwanyoung Park, and Youngwoon Lee. Chunk-guided q-learning, 2026. URL
150 <https://arxiv.org/abs/2603.13971>.
- 151 Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual
152 manipulation with low-cost hardware, 2023. URL <https://arxiv.org/abs/2304.13705>.