

# MULTISTEP QUASIMETRIC LEARNING FOR SCALABLE GOAL-CONDITIONED REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

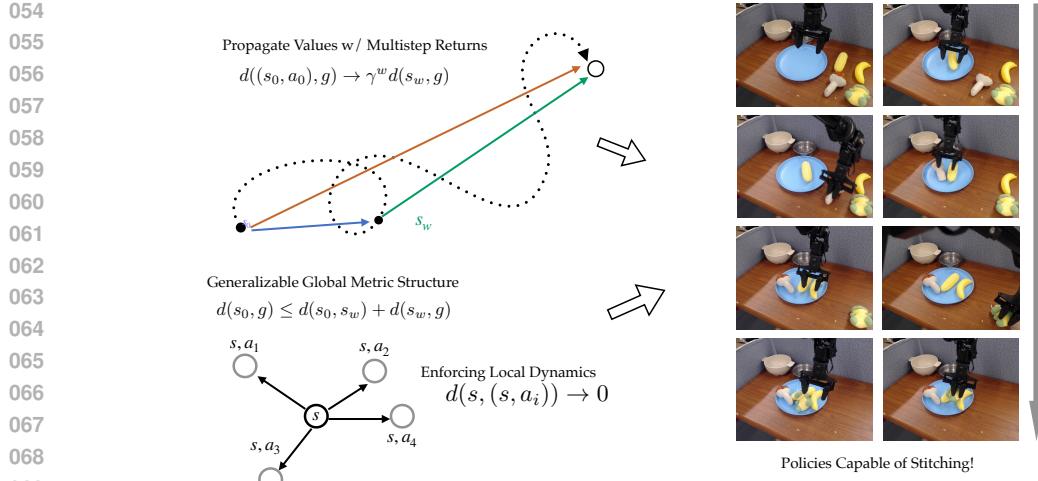
The problem of learning how to reach goals in an environment has been a long-standing challenge in for AI researchers. Effective goal-conditioned reinforcement learning (GCRL) methods promise to enable reaching distant goals without task-specific rewards by stitching together past experiences of different complexity. Mathematically, there is a duality between the notion of optimal goal-reaching value functions (the likelihood of success at reaching a goal) and temporal distances (transit times states). Recent works have exploited this property by learning quasimetric distance representations that stitch long-horizon behaviors using the inductive bias of their architecture. These methods have shown promise in simulated benchmarks, reducing value learning to a shortest-path problem. But quasimetric, and more generally, goal-conditioned RL methods still struggle in complex environments with stochasticity and high-dimensional (visual) observations. There is a fundamental tension between the local dynamic programming (TD backups, temporal distances) that enables optimal shortest-path reasoning in theory and the statistical global MC updates (multistep returns, suboptimal in theory). We show how these approaches can be integrated into a practical GCRL method that fits a quasimetric distance using a multistep Monte-Carlo return. We show our method outperforms existing GCRL methods on long-horizon simulated tasks with up to 4000 steps, even with visual observations. We also demonstrate that our method can enable stitching in the real-world robotic manipulation domain (Bridge setup). Our approach is the first end-to-end GCRL method that enables multistep stitching in this real-world manipulation domain from an unlabeled offline dataset of visual observations.<sup>1</sup>

## 1 INTRODUCTION

It is natural for humans to use inherent ideas of distances to represent task progress: a GPS will tell you how far you are from the destination, and a cookbook will tell you how long a recipe will take. Humans can generalize tasks by taking the shortest route possible (stitching) and combining several learned tasks together in sequence in a new environment (combinatorial generalization). All of these problems can be thought of as reaching goals. The AI problem of reaching goals presents rich structure (formally, an optimal substructure property) that can be exploited to decompose hard problems into easy problems, and reinforcement learning (RL) has been used to address such problems. However, many past attempts at leveraging this property in modern high-dimensional, stochastic settings are incoherent or inconsistent. Current approaches tend to separate TD learning (local updates) (Mnih et al., 2013; Kostrikov et al., 2022; Kumar et al., 2020) and MC learning (global value propagation) (Eysenbach et al., 2022; Myers et al., 2024) into separate categories, with each method basing its approach based on only one, but not the other.

However, the effectiveness of both these methods degrade from a combination of increasing horizon length for TD methods (Park et al., 2025b) or difficulties in finding the optimal temporal distance (Park et al., 2024a). In conjunction, separate challenges exist if one were to train a policy using offline RL in real world. Under controlled datasets in simulated benchmarks, these algorithms may return a good goal-reaching policy, recent works have shown that additional work are needed for it to successfully train and deploy a policy in real life (Zheng et al., 2024).

<sup>1</sup>Website and code: <https://anonymous.4open.science/w/mqe-paper-686D/>



074 Our main contribution lies in **Multistep Quasimetric Estimation** (MQE), an offline GCRL method  
075 that incorporates both multistep value learning and quasimetric architectures without needing explicit  
076 hierarchy. We make the crucial observation that under the constraint of a quasimetric architecture, the  
077 regression objective for multistep returns are compatible with objectives for one-step consistency and  
078 action invariance, which are key insights for finding the optimal Q function. To our knowledge, MQE  
079 is the first method that is capable of leveraging multistep TD returns with global value propagation  
080 via quasimetric architectures.

081 By leveraging such a unique combination of the benefits of TD and MC methods, MQE allows the  
082 learned policy to (i): display a much stronger level of horizon generalization compared to previous  
083 methods, which allow us for demonstrating the desired “stitching” behavior, (ii): provide a stable  
084 method for optimizing towards an optimal Q function, and (iii): such stability in training allows it to  
085 be applied in real-world robot learning problems without additional work in engineering. MQE can  
086 demonstrate SoTA performance on tasks that require complex control and long-horizon reasoning  
087 (up to 21 DoF and 4000 timesteps respectively), and in real world robotic settings, MQE displays  
088 compositional generalization behaviors that are not seen in previous RL algorithms.

## 2 RELATED WORKS

091 Our work build upon previous works in offline RL and temporal distance learning.

092 **Goal-conditioned Reinforcement Learning** Our findings particularly concern the discussions  
093 seen in recent work regarding the stitching and horizon generalization capabilities of GCRL. Recent  
094 findings (Park et al., 2025b) have shown that while it is easy to scale offline RL on short-horizon  
095 tasks, it is difficult to learn long-horizon tasks that require more complex reasoning within agent,  
096 which can easily deviate from the true optimal distance due to compounding TD errors.

097 **Offline RL for Robotics** In the same vein, we also focused on how we can apply these distance  
098 learning techniques to robotics. While reinforcement learning has been used to obtain highly capable  
099 specialist policies across various embodiments (Luo et al., 2025; Ball et al., 2023; Seo et al., 2025;  
100 Smith et al., 2023), behavior cloning still remains the most capable method for training generalist  
101 policies (O’Neill et al., 2024). Efforts have been made in allowing self-supervised and offline RL in  
102 robotics (Zheng et al., 2024), however, directly training a policy with offline RL remains difficult,  
103 as many researchers have instead used other ways to incorporate RL, such as rejection sampling  
104 (Nakamoto et al., 2025; Wagenmaker et al., 2025) or dataset curation (Mark et al., 2024; Xu et al.,  
105 2024). The focus of this work lies more on designing and implementing an offline RL method that  
106 can work in real world instead of only focusing on techniques for such training.

108    **Temporal Distance Learning** Our work is partially inspired by successor representations (Dayan,  
 109    1993; ?). Previous works have shown that by using contrastive learning, (Myers et al., 2024) Other  
 110    works have also shown that using contrastive learning as a way to parameterize distance learning may  
 111    also recover behavior Q function.  
 112

113    **Multistep RL** RL using multistep returns have been widely used in online RL and offline-to-online  
 114    RL settings (?Tian et al., 2025). This is desirable because in on-policy settings, performing RL with  
 115    multistep return is correct (Munos et al., 2016). However, in an offline setting, the theory behind such  
 116    correctness breaks down (watkins1989learning, 1989), as the learning becomes inherently off-policy.  
 117    However, methods that use the entire trajectory and attempt to recover the behavior Q function can  
 118    also be seen as a multistep RL algorithm (Eysenbach et al., 2022), however they also suffer the same  
 119    problem of unable to learn  $Q^*$ .  
 120

### 121    3 PRELIMINARIES

123    In this section, we concretely define temporal distances and learning objective.  
 124

125    **Notation.** We consider a controlled Markov process (CMP)  $\mathcal{M}$  with state space  $\mathcal{S}$ , action space  
 126     $\mathcal{A}$ , transition dynamics  $P(s'|s, a)$ , and discount factor as  $\gamma$ . We consider goal-reaching policies  
 127     $\pi(a|s, g) : \mathcal{S} \times \mathcal{S} \rightarrow \mathcal{A} \in \Pi$ . We denote the behavior policy that is used to populate the offline  
 128    dataset as  $\pi_\beta$ . In lieu of rewards, we optimize the maximum discounted likelihood of a policy  
 129    reaching the goal, in which we can represent the Q-function and value function as:  
 130

$$132 \quad Q_g^\pi(s, a) \triangleq \mathbb{E}_{\{\mathbf{s}_t, \mathbf{a}_t\} \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t P(\mathbf{s}_t = g \mid \mathbf{s}_0 = s, \mathbf{a}_0 = a) \right]. \quad (1)$$

$$136 \quad V_g(s) \triangleq \mathbb{E}_{\{\mathbf{s}_t\} \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t P(\mathbf{s}_t = g \mid \mathbf{s}_0 = s) \right]. \quad (2)$$

139    Equivalently, we can define the optimal Q-function as  $Q_g^*(s, a) \triangleq \max_{\pi \in \Pi} Q_g^\pi(s, a)$ . Previous work  
 140    have shown that using contrastive leaning (Eysenbach et al., 2022; van den Oord et al., 2019) can  
 141    recover the behavior distance, but not the optimal distance.  
 142

143    Prior work on MC learning (Myers et al., 2024; Eysenbach et al., 2022; 2021) has demonstrated that  
 144    instead of only using the end of trajectories as the goal state, we sample future states via a geometric  
 145    distribution as described in Eq. (3). This allows us to classify any future state in trajectory as goals,  
 146    providing a robust way of learning goal-reaching policies.  
 147

$$148 \quad \mathbf{s}_t^+ \triangleq \mathbf{s}_{t+K}, K \sim \text{Geom}(1 - \gamma). \quad (3)$$

150    **Quasimetric distance representations.** Traditionally, offline RL algorithms use neural networks to  
 151    represent the critic and value function  $Q_g(s, a)$  and  $V_g(s)$  (Kumar et al., 2020; Kostrikov et al., 2022).  
 152    Separately, other works have been using dot products  $Q(s, a, g) \triangleq \varphi(s, a)^\top \psi(g)$  (Eysenbach et al.,  
 153    2022; Zheng et al., 2024) or geometric norms  $\|\varphi(s) - \psi(g)\|_k$  for suitable values of  $k$  (Eysenbach  
 154    et al., 2024; Park et al., 2024c). We use quasimetric architectures (Wang et al., 2023; Pitis et al.,  
 155    2020) to parameterize our value-functions, which by construction (i): positivity ( $d(x, y) > 0$ ), (ii):  
 156    triangle inequality, and (iii): identity.

157    To enforce the quasimetric properties of successor distances, we use Metric Residual Networks  
 158    (MRN) (Liu et al., 2023) that directly operate on vector-based representations, which is defined by  
 159    Eq. (4). MRN splits the representation into  $N$  equally sized ensembles, and in each part, takes the  
 160    sum of the asymmetric component (maximum of ReLU) and a symmetric component ( $l_2$  norm) of  
 161    the difference between the two embeddings  $x$  and  $y$ . For simplicity, we denote  $d(x, y)$  as  $d_{\text{MRN}}(x, y)$   
 in the following sections.

162  
 163  
 164       $d_{\text{MRN}}(x, y) \triangleq \frac{1}{N} \sum_{k=1}^N \max_{m=1 \dots M} \text{ReLU}(x_{kM+m} - y_{kM+m}) + \|x_{kM+m} - y_{kM+m}\|_2$       (4)  
 165  
 166

## 4 MULTISTEP QUASIMETRIC ESTIMATION

In this section, we develop the framework of MQE based on quasimetric architectures, which can be broken down into 3 parts: (1) propagating the learned temporal distance locally using a multistep return with geometrically sampled future states as waypoints, (2) imposing additional objectives that need to be satisfied to learn  $Q^*$ , and (3) practical implementation details for extracting the goal-reaching policy. To the best of our knowledge, MQE is the **first** method to effectively use multistep returns under quasimetric distance representations, and achieve superior results to previous methods that use TD returns, contrastive learning for value estimation, or hierarchical methods.

**Definitions.** Argued by (Myers et al., 2024), define the relations between the distance metric and the underlying Q/value function of the environment:

$$Q_g(s, a) \triangleq \exp(-d(\varphi(s, a), \psi(g))), \quad V_g(s) \triangleq \exp(-d(\psi(s), \psi(g))). \quad (5)$$

We also adopt the following shorthand :

$$d((s, a), g) \triangleq d(\varphi(s, a), \psi(g)), \quad d(s, g) \triangleq d(\psi(s), \psi(g)). \quad (6)$$

### 4.1 MULTISTEP RETURNS WITH QUASIMETRIC ARCHITECTURE

The first design principle of MQE is to use multistep returns under a quasimetric architecture. To that end, we start with the fitted one-step Q iteration that operates on quasimetric distance representations, in which  $\rightarrow$  denotes regressing from the LHS to the RHS:

$$e^{-d((s, a), g)} \rightarrow \mathbb{E}_{\{(s, a, s')\} \sim \mathcal{D}} [\gamma \cdot e^{-d(s', g)}]. \quad (7)$$

This is similar to optimizing the critic in IQL (Kostrikov et al., 2022), which recovers  $V^*$  in goal-conditioned settings, the main difference being that we use a quasimetric architecture to represent the Q and value function. We may omit the reward signal, as the value function  $V$  is defined to have a value of 1 when  $s = g$ , abstracting away the need of a reward. We denote this regression objective as  $\mathcal{T}$ .

Our insight for the  $\mathcal{T}$  objective described in Eq. (7) is that, instead of applying this invariance to only the future state  $s' \sim P(s'|s, a)$ , we can extend this principle to any state between the current state and goal. This transforms a one-step optimization into a on-policy multistep optimization procedure. To do so, we first define the shorthand  $s_t^w$ , which refers to a “**waypoint**” between the current state and goal state. Empirically, we find that using a geometric distribution capped at the index of the future state work the best.

$$s_t^w \triangleq s_{t+k'}, k' \sim \min(\text{Geom}(1 - \lambda), K) \quad (8)$$

We now optimize the same objective as in Eq. (7), but across any such waypoint we sample. To account for the multistep nature of this objective, we modify Eq. (7) below to accommodate such changes.

$$e^{-d((s, a), g)} \rightarrow \mathbb{E}_{\{(s_t, a_t), s_t^w\} \sim \mathcal{D}} [\gamma^{k'} \cdot e^{-d(s_t^w, g)}]. \quad (9)$$

We denote this new objective as  $\mathcal{T}_\beta$ , as sampled future state is constrained by  $\pi_\beta$ . This is indeed similar to the n-step returns we see in RL (Munos et al., 2016; Li et al., 2025), although we do not sample future states with a fixed number of steps. In practice, we can use any loss function to make the LHS equal to the RHS (in expectation) concerning Eq. (7) and Eq. (9). We use a form of Bregman divergence, as it does not incur vanishing gradients when the two distance have become close in value, and regresses  $d$  to  $d'$  in expectation (Banerjee et al., 2004).

$$D_T(d, d') \triangleq \exp(d - d') - d' \quad (10)$$

216 By using this loss, we can concretely define both  $\mathcal{L}_{\mathcal{T}_\beta}$  and  $\mathcal{L}_{\mathcal{T}}$  that can optimize  $\mathcal{T}$  and  $\mathcal{T}_\beta$  objectives:  
 217

$$218 \quad \mathcal{L}_{\mathcal{T}_\beta}(\phi, \psi; \{s_i, a_i, s_i^w, g_i\}_{i=1}^N) = \sum_{i=1}^N \sum_{j=1}^N D_T(d((s_i, a_i), g_j), d(s_i^w, g_j)) - k' \log \gamma. \quad (11)$$

$$222 \quad \mathcal{L}_{\mathcal{T}}(\phi, \psi; \{s_i, a_i, s'_i, g_i\}_{i=1}^N) = \sum_{i=1}^N \sum_{j=1}^N D_T(d((s_i, a_i), g_j), d(s'_i, g_j)) - \log \gamma. \quad (12)$$

225 These two losses, when applied to our quasimetric network, will allow us to propagate the value in  
 226 either an one-step or multistep manner.  
 227

## 228 4.2 ENFORCING OPTIMALITY CONSTRAINTS IN DISTANCE REPRESENTATION

230 **Remark:**  $\mathcal{L}_{\mathcal{T}_\beta}$  (**weakly**) enforces  $\mathcal{L}_{\mathcal{T}}$ . While optimizing  $\mathcal{L}_{\mathcal{T}_\beta}$  allows us to learn a multistep backup  
 231 under the quasimetric architecture, it is *biased* towards the behavior policy when  $k' \neq 1$ , as the state  
 232 on the waypoint  $s_t^w$  must be conditioned on the trajectories made by the behavior policy. This is in  
 233 contrast to the one-step objective described in Eq. (7), which allows us to learn the optimal value  
 234 function in theory (Myers et al., 2024; Kostrikov et al., 2022), but is hard to use only  $\mathcal{L}_{\mathcal{T}}$  in practice  
 235 due to the extremely local nature of the one-step value propagation (Myers et al., 2025b).

236 The geometric distribution provides some remedy, as  $k' = 1$  is the most likely outcome of sampling.  
 237 As a result, we can slightly enforce  $\mathcal{L}_{\mathcal{T}}$ , as some of the sampled waypoints will be equal to  $s'_t$ . We  
 238 discuss ways to manually increase the weighting of  $\mathcal{L}_{\mathcal{T}}$  in Section 4.4.

240 **Enforcing the action invariance property.** We also note that the optimal Q and value function  
 241 observe the below property (Sutton & Barto, 2018):  
 242

$$243 \quad V_g^*(s) \propto \max_{a \in \mathcal{A}} Q_g^*(s, a). \quad (13)$$

245 This can be satisfied, under our construction, if  $\psi(s) = \varphi(s, a), \forall a \in \mathcal{A}$  (Myers et al., 2025b) (action  
 246 invariance). This can be done with the same MRN distance formulation due to the identity property  
 247 of the quasimetric network, as described in Eq. (14).  
 248

$$249 \quad d(\psi(s), \varphi(s, a)) \rightarrow 0 \quad (14)$$

251 While this objective can be optimized by taking the gradient with just the MRN distance, we modify  
 252 the action invariance such that we use the sum of squares of the MRN distance between  $s$  and  $(s, a)$ .  
 253 The insight is that now the loss will scale with the magnitude of such deviation, which removes the  
 254 need of hyperparameter tuning for an appropriate multiplier as well as stabilizing training dynamics.  
 255 We can now formally define  $\mathcal{L}_{\mathcal{I}}$  objective as:  
 256

$$257 \quad \mathcal{L}_{\mathcal{I}}(\varphi, \psi; \{s_i, a_i\}_{i=1}^N) = \sum_{i,j=1}^N (d(\psi(s_i), \varphi(s_i, a_j)))^2. \quad (15)$$

## 260 4.3 POLICY EXTRACTION

262 We extract the goal-conditioned policy  $\pi(s, g) : \mathcal{S}^2 \rightarrow \mathcal{A}$  using learned distance by using Behavior-  
 263 Regularized Deep Deterministic Policy Gradient (DDPG+BC) (Fujimoto & Gu, 2021):  
 264

$$265 \quad \mathcal{L}_\mu(\pi; \{s_i, a_i, g_i\}_{i=1}^N) = \mathbb{E} \left[ \left[ \sum_{i,j=1}^N d((s_i, \pi(s_i, g_j)), g_j) \right] - \alpha \log \pi(a_i | s_i, g_i) \right]. \quad (16)$$

268 Given that a smaller  $d$  measure correspond to a higher Q value, we can maximize the Q values by  
 269 minimizing the distance produced by our quasimetric network. We tune the BC coefficient  $\alpha$  per  
 environment. We provide more hyperparameter details in Appendix B.

270 4.4 IMPLEMENTATION DETAILS & ALGORITHM  
271

272 Although  $\mathcal{L}_{\mathcal{T}_\beta}$  becomes  $\mathcal{L}_{\mathcal{T}}$  in the case where  $k' = 1$ , we can adjust the scale of the two losses  
273 by explicitly up-weighting the probability of sampling the next state as waypoint via a Bernoulli  
274 distribution. In practice, we designate  $s_t^w = s'_t$  with probability  $p$ , and with probability  $1 - p$ , we  
275 designate  $s_t^w$  as is. In this instance, if we designate  $p = 1$ ,  $\mathcal{L}_{\mathcal{T}_\beta}$  becomes  $\mathcal{L}_{\mathcal{T}}$ . We find using  $p = 0.2$   
276 across all environments allows the network to both propagate local distances fast and maintain the  
277 ability to optimize towards  $Q^*$  using one-step consistency.

278 We concisely define our final learning objective in Algorithm 1. Unlike previous works in hierarchical  
279 RL (Nachum et al., 2018; Park et al., 2024b), we randomly sample these waypoints, and we learn  
280 a single critic  $\mathcal{Q}$  that operates on  $\mathcal{S}$  and  $\mathcal{A}$  and a single goal-reaching policy  $\pi_\mu$ . As a result, our  
281 method does not contain any hierarchical components and is simpler to implement than other horizon  
282 reduction methods.

283  
284 **Algorithm 1:** Multistep Quasimetric Reinforcement Learning

---

285 **Require:** Dataset  $\mathcal{D}$ , Batch size  $B$ , training iteration  $T$ , Probability  $p$   
286 1: Initialize Quasimetric network  $\mathcal{Q}$  with parameters  $(\varphi(s, a), \psi(s))$ , goal-reaching policy  $\pi_\mu$   
287 2: **for**  $t = 1 \dots T$  **do**  
288 3:   Sample  $\{s_i, a_i, s'_i, \tilde{s}_i^w, g_i\}_{i=1}^B \sim \mathcal{D}$  (Eq. (8))  
289 4:  
290     For each element in batch, choose  $s_i^w \sim \begin{cases} s'_i & \text{with probability } p, \\ \tilde{s}_i^w & \text{with probability } 1 - p. \end{cases}$   
291 5:     Update  $\mathcal{Q}$  with multistep backup by minimizing  $\mathcal{L}_{\mathcal{T}_\beta}(\varphi, \psi; \{s_i, a_i, s_i^w\}_{i=1}^B)$  (Eq. (9))  
292 6:     Update  $\mathcal{Q}$  with optimality constraints by minimizing  $\mathcal{L}_{\mathcal{I}}(\varphi, \psi; \{s_i, a_i\}_{i=1}^B)$  (Eq. (15))  
293 7:     Update policy  $\pi_\mu$  with DDPG+BC by minimizing  $\mathcal{L}_\mu(\pi_\mu; \{s_i, a_i, g_i\}_{i=1}^B)$  (Eq. (16))  
294 8: **return**  $\pi_\mu$

---

295 5 EXPERIMENTS  
296

300 Our goal of experiments is to understand the benefits MQE brings when it comes to enabling a policy  
301 to generalize compositionally (execute multiple tasks seen in training separately together) and in  
302 terms of horizon (generalize over a longer task when a similar but shorter task was seen in training  
303 set). To that end, we ask the following research questions:

- 304 1. How much does MQE improve the horizon generalization abilities of agents compared to  
305 prior works?  
306 2. What qualitative improvements does MQE bring in terms of compositional generalization?  
307 3. What are the important hyperparameters to ensure the success of MQE?

310 **Experiment setup** Our experiments use challenging, long-horizon problems in offline RL bench-  
311 marks as well as real-world settings. We use OGBench (Park et al., 2025a) for our experiments on  
312 simulated benchmarks and the BridgeData setup (Walke et al., 2024) for our real-world evaluation.  
313

314 5.1 EVALUATING MQE ON OFFLINE GCRL BENCHMARK  
315

316 We evaluate MQE in both locomotion and manipulation in OGBench (Park et al., 2025a). For  
317 locomotion tasks, in addition to the three standard sized mazes, we also designed a “colossal”-sized  
318 maze. This maze is 50% larger than that of the “giant”-sized mazes currently available on OGBench,  
319 and it requires as many as 4000 steps for an agent to traverse through the entire maze (see Fig 2). We  
320 employ 13 state-based environments and 5 pixel-based environments (each pixel-based environment  
321 takes in a  $64 \times 64 \times 3$  observation) with 5 tasks each, bringing a total of **90** tasks to evaluate in our  
322 OGBench setup.

323 We compare against the following baselines: GCIQL (Kostrikov et al., 2022), CRL (Eysenbach et al.,  
324 2022), QRL (Wang et al., 2023), HIQL (Park et al., 2024b), nSAC+BC (Park et al., 2025b; Haarnoja

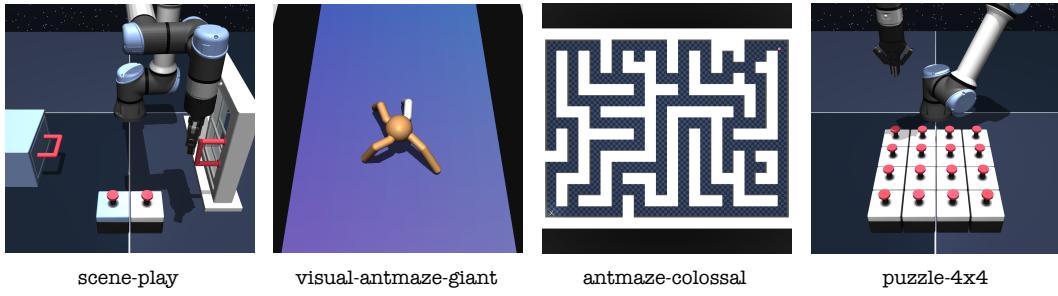


Figure 2: We provide visualizations of tasks from various state and pixel-based environments for OGBench. Note that Antmaze-colossal is 50% larger than any other mazes available on OGBench, and in `stitch` datasets, test the agent’s ability to generalize over a horizon that is up to **1000%** larger.

Table 1: OGBench Evaluation

Dataset	Methods							
	MQE	TMD	CMD	CRL	QRL	GCIQL	HIQL	n-SAC+BC
<code>pointmaze_giant_navigate</code>	<b>72.8</b> ( $\pm 2.5$ )	39.9( $\pm 5.2$ )	45.3( $\pm 3.7$ )	27.4( $\pm 3.4$ )	<b>68.5</b> ( $\pm 2.8$ )	0.0( $\pm 0.0$ )	45.9( $\pm 3.0$ )	0.0( $\pm 0.0$ )
<code>pointmaze_giant_stitch</code>	28.3( $\pm 1.1$ )	9.1( $\pm 1.0$ )	8.1( $\pm 0.6$ )	0.0( $\pm 0.0$ )	<b>49.7</b> ( $\pm 2.3$ )	0.0( $\pm 0.0$ )	0.0( $\pm 0.0$ )	0.4( $\pm 0.1$ )
<code>antmaze_large_explore</code>	<b>44.1</b> ( $\pm 4.2$ )	0.9( $\pm 0.2$ )	0.8( $\pm 0.3$ )	0.3( $\pm 0.1$ )	0.0( $\pm 0.0$ )	0.4( $\pm 0.1$ )	3.9( $\pm 1.8$ )	0.2( $\pm 0.1$ )
<code>antmaze_giant_stitch</code>	<b>35.1</b> ( $\pm 1.7$ )	2.7( $\pm 0.6$ )	2.0( $\pm 0.5$ )	0.0( $\pm 0.0$ )	0.4( $\pm 0.2$ )	0.0( $\pm 0.0$ )	1.8( $\pm 0.6$ )	9.2( $\pm 2.1$ )
<code>antmaze_colossal_navigate</code>	<b>48.6</b> ( $\pm 2.4$ )	22.3( $\pm 1.1$ )	22.5( $\pm 3.1$ )	14.6( $\pm 1.8$ )	0.0( $\pm 0.0$ )	0.0( $\pm 0.0$ )	0.0( $\pm 0.0$ )	0.3( $\pm 0.1$ )
<code>antmaze_colossal_stitch</code>	<b>27.6</b> ( $\pm 2.9$ )	0.0( $\pm 0.0$ )	0.2( $\pm 0.1$ )	0.0( $\pm 0.0$ )	0.0( $\pm 0.0$ )	0.0( $\pm 0.0$ )	0.0( $\pm 0.0$ )	0.5( $\pm 0.3$ )
<code>humanoidmaze_giant_navigate</code>	<b>46.5</b> ( $\pm 2.5$ )	9.2( $\pm 1.1$ )	5.0( $\pm 0.8$ )	0.7( $\pm 0.1$ )	0.4( $\pm 0.1$ )	0.5( $\pm 0.1$ )	12.5( $\pm 1.5$ )	3.2( $\pm 0.5$ )
<code>humanoidmaze_giant_stitch</code>	<b>26.5</b> ( $\pm 1.3$ )	6.3( $\pm 0.6$ )	0.2( $\pm 0.1$ )	1.5( $\pm 0.5$ )	0.4( $\pm 0.1$ )	1.5( $\pm 0.1$ )	3.3( $\pm 0.7$ )	1.7( $\pm 0.1$ )
<code>cube_double_play</code>	35.6( $\pm 1.9$ )	13.1( $\pm 2.3$ )	0.2( $\pm 0.1$ )	1.5( $\pm 0.5$ )	0.4( $\pm 0.1$ )	<b>40.2</b> ( $\pm 1.7$ )	6.4( $\pm 0.7$ )	19.1( $\pm 0.3$ )
<code>cube_triple_noisy</code>	<b>3.9</b> ( $\pm 0.8$ )	2.1( $\pm 0.6$ )	1.5( $\pm 0.5$ )	2.7( $\pm 0.5$ )	<b>3.4</b> ( $\pm 0.4$ )	1.8( $\pm 0.2$ )	2.6( $\pm 0.4$ )	1.4( $\pm 0.3$ )
<code>puzzle_4x4_play</code>	18.7( $\pm 2.3$ )	10.0( $\pm 1.4$ )	0.2( $\pm 0.1$ )	1.5( $\pm 0.5$ )	0.4( $\pm 0.1$ )	<b>25.7</b> ( $\pm 1.1$ )	7.4( $\pm 0.7$ )	11.4( $\pm 0.9$ )
<code>scene_play</code>	<b>69.1</b> ( $\pm 2.1$ )	21.9( $\pm 1.9$ )	1.2( $\pm 0.4$ )	18.6( $\pm 0.8$ )	5.4( $\pm 0.3$ )	51.3( $\pm 1.5$ )	38.2( $\pm 0.9$ )	17.6( $\pm 1.4$ )
<code>scene_noisy</code>	<b>30.8</b> ( $\pm 1.6$ )	19.6( $\pm 1.7$ )	4.0( $\pm 0.7$ )	1.2( $\pm 0.3$ )	9.1( $\pm 0.7$ )	25.9( $\pm 0.8$ )	25.2( $\pm 1.3$ )	19.1( $\pm 2.2$ )
<code>visual_scene_play</code>	27.2( $\pm 3.1$ )	20.7( $\pm 2.5$ )	16.1( $\pm 2.2$ )	9.6( $\pm 0.6$ )	5.4( $\pm 0.3$ )	12.2( $\pm 0.8$ )	<b>49.9</b> ( $\pm 0.6$ )	7.1( $\pm 1.2$ )
<code>visual_cube_triple_play</code>	<b>19.8</b> ( $\pm 0.9$ )	<b>17.9</b> ( $\pm 1.3$ )	<b>18.9</b> ( $\pm 1.1$ )	16.9( $\pm 1.1$ )	16.3( $\pm 0.3$ )	15.2( $\pm 0.6$ )	<b>21.0</b> ( $\pm 0.2$ )	<b>21.1</b> ( $\pm 2.4$ )
<code>visual_cube_double_noisy</code>	25.9( $\pm 1.6$ )	14.2( $\pm 1.3$ )	0.3( $\pm 0.3$ )	6.0( $\pm 1.4$ )	6.1( $\pm 1.2$ )	21.6( $\pm 0.9$ )	<b>59.4</b> ( $\pm 1.6$ )	22.7( $\pm 1.1$ )
<code>visual_puzzle_4x4_play</code>	11.3( $\pm 1.6$ )	9.8( $\pm 3.6$ )	7.2( $\pm 0.4$ )	9.6( $\pm 3.2$ )	0.0( $\pm 0.0$ )	16.2( $\pm 2.2$ )	<b>60.1</b> ( $\pm 20.4$ )	10.3( $\pm 2.6$ )
<code>visual_antmaze_giant_stitch</code>	<b>26.9</b> ( $\pm 3.1$ )	14.5( $\pm 2.5$ )	<b>22.3</b> ( $\pm 1.9$ )	0.1( $\pm 0.1$ )	0.0( $\pm 0.0$ )	0.0( $\pm 0.0$ )	0.2( $\pm 0.1$ )	7.6( $\pm 1.1$ )
Overall	<b>32.3</b> ( $\pm 0.5$ )	13.0( $\pm 0.5$ )	8.7( $\pm 0.3$ )	6.2( $\pm 0.3$ )	9.8( $\pm 0.3$ )	11.8( $\pm 0.2$ )	18.7( $\pm 1.2$ )	8.5( $\pm 0.3$ )

We **bold** the best performance. Success rate (%) is presented with the standard error across eight seeds for state-based environments and four seeds for pixel-based environments. All datasets contain 5 separate tasks each. We record the aggregate across all 5 tasks.

et al., 2018), CMD (Myers et al., 2024), and TMD (Myers et al., 2025b). These methods use either only TD learning (GCIQL, QRL, nSAC+BC), MC value estimation via contrastive learning (CRL, CMD, TMD), use horizon reduction techniques (nSAC+BC with value horizon, HIQL with policy horizon), or use a quasimetric architecture for distance learning (QRL, CMD, TMD). We detail how these methods are implemented in Appendix B. By comparing against these methods, we can gain a better outlook on *what* advantage MQE has over other works that learn distances only locally, globally, or in a hierarchical manner.

Table 1 shows the performance of MQE across state- and pixel-based environments on OGBench. Overall, MQE performed at a much better level compared to all baselines. We note the significant gains over previous methods on many extremely challenging locomotion environments, and that in manipulation environments, where GCIQL dominates the performance, MQE still closes the gap considerably, and provided the best performance over in the `scene` environment. We also demonstrate that MQE is able to extract better goal-reaching policies from both `explore` and `noisy` datasets, which for methods that only obtained  $Q_\beta$ , was difficult for solve, showing that MQE is capable of extracting optimal policies from suboptimal demonstrations.

In visual environments, MQE maintains its strong empirical performance, only being bested by HIQL, which does explicit policy horizon reduction (Park et al., 2025b) and was outperformed by other methods in state-based environments due to HIQL learning an additional subgoal representation that helps to stabilize training. Nevertheless, MQE is still able to outperform all other methods, including n-SAC+BC, which also conducts horizon reduction.

378 5.2 EVALUATING MQE IN REAL WORLD  
379

380 While the OGBench environments focus on  
 381 learning long-horizon tasks using mixed quality  
 382 data in a single environment, we can use the real-  
 383 world BridgeData tasks evaluate a more sophisti-  
 384 cated kind of compositionality: the BridgeData  
 385 tasks consist of individual object manipulation  
 386 primitives (e.g., picking up a banana and placing  
 387 it on a plate), and our evaluation tasks are signif-  
 388 icantly longer, requiring the composition of mul-  
 389 tiple tasks in the dataset (e.g., placing four dif-  
 390 ferent objects on the plate) (Walke et al., 2024).  
 391 Critically, the dataset *does not contain any ex-*  
*392 ample trajectories* that compose multiple tasks  
 393 in this way. Accomplishing this sort of temporal  
 394 composition is an important goal in offline RL,  
 395 because it allows “stitching” long behaviors out  
 396 of shorter chunks. However, most current meth-  
 397 ods struggle with this challenge, especially in  
 398 real-world domains with many objects, scenes,  
 399 and scenarios, where this sort of stitching neces-  
 400 sarily requires both temporal compositionality  
 and generalization.

401 As a result, a policy that is truly able to perform compositional generalization should be able to  
 402 compose multiple tasks at once without external guidance. We designed the tasks on BridgeData with  
 403 this specific intention. Instead of tasking the policy to complete a single pick and place (abbreviated  
 404 as PnP), we evaluate a policy’s performance with PnP of up to 4 objects, which, to our knowledge,  
**405 has never been completed without the use of hierarchical policies or high level planners.** We  
 406 also evaluate the policy on tasks requiring dependencies, with the policy being tasked with opening a  
 407 drawer and then placing the item within the drawer all conditioned by one image of an opened drawer  
 408 with an item inside. Such a task is also challenging, as **only one previous work** (Myers et al., 2025a)  
 409 has shown a non-zero success rate when using an end-to-end policy to the best of our knowledge.

410 We use a 6DoF, 5Hz, WidowX250 manipulator for our robot learning tasks and we train and deploy a  
 411 policy  $\pi(a|s, g)$  conditioned on observations and goal images. We compare against the following  
 412 methods: ding2019goal, GCIQL, and TRA (Myers et al., 2025a). GCIQL follows the same structure  
 413 as OGBench training, and both ding2019goal and TRA is a successor feature learning method that  
 414 aligns present and future representations for compositional generalization. TRA is designed for  
 415 following both goal images and language instructions, but since we only use goal images as the  
 416 modality, we denote TRA-g as the specific modality we’re testing. We provide more details on policy  
 417 training in Appendix C.1.

418 As in (Black et al., 2024), we use task progress to measure the effectiveness of these policies due to  
 419 the long-horizon nature of these tasks. We detail more on the experimental setup and how we assign  
 420 these points in Appendix C.1.

421 **What qualitative improvements does MQE provide in terms of compositional generalization?**  
 422 Fig. 4 reports the *overall task progress* on 2 single-stage tasks and 4 tasks requiring compositionality.  
 423 We provide both the binary success rate and further analysis of policy rollout in Appendix C. Here, we  
 424 observe that while MQE helps with single-stage tasks (single PnP, open drawer) against ding2019goal,  
 425 both TRA and GCIQL can still perform competitively. However, as the number of tasks needed to be  
 426 performed in sequence increases, we see that MQE is able to retain a relatively high task progress,  
 427 while both GCIQL and TRA-g’s performance regressed.

428 Taking a look at the two most difficult tasks, we have quadruple PnP and drawer open and  
 429 place. These tasks are the most challenging since quadruple PnP required the agent to reason 4  
 430 consecutive primitives together, and drawer open and place required the agent to complete the first  
 431 task (open the drawer) before completing the second task (putting the mushroom in the drawer).

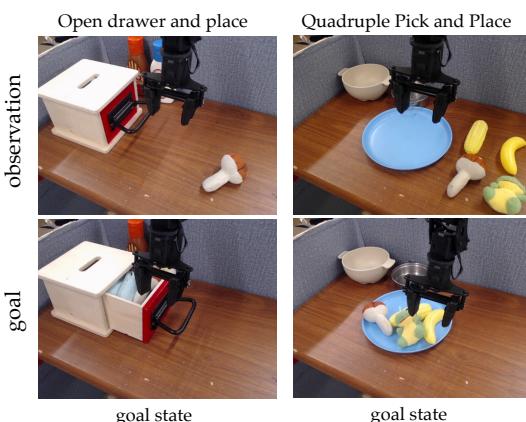


Figure 3: We evaluate MQE on multi-stage manip-  
 ulation tasks on BridgeData. Below are examples  
 of the starting observations and goal image being  
 passed in.

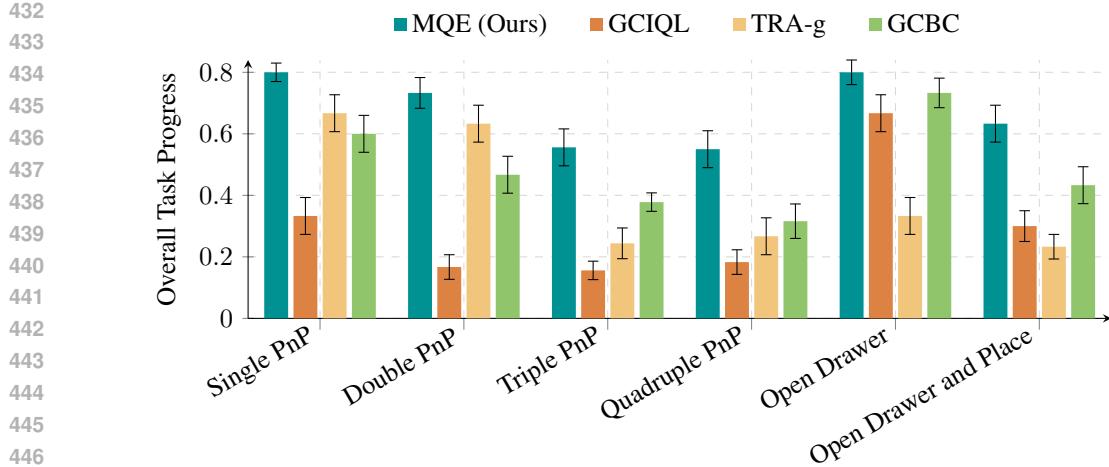


Figure 4: Task progress on BridgeData tasks; plotted with both the mean and the standard error bars.

Among all methods that we have tested, only MQE and TRA-g displayed positive success rate, as demonstrated in Table 5. We provide more details on policy rollouts in Appendix C.

### 5.3 ABLATION STUDIES

In this section, we explore design choices involved for MQE. Given the two components needed for critic loss and an additional sampling of waypoints, it is natural for one to ask whether they are conflicting objectives and *how* they contribute to the strong empirical observation. To that end, we investigate design decisions for MQE, and ask the following questions:

- What is the best empirical distribution for sampling the waypoint  $s_t^w$ ?
- What role does the invariance objective  $\mathcal{I}$  play in MQE?
- What is the relationship between  $\lambda$  and  $p$ , the two hyperparameters that enforce  $\mathcal{L}_{\mathcal{T}_\beta}$  and  $\mathcal{L}_{\mathcal{T}}$  separately?

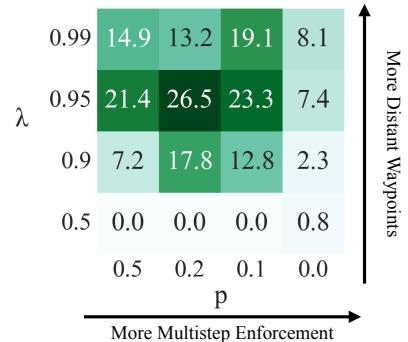
Fig. 5 provides an illustration of success rate over pairs of  $(p, \lambda)$ . The figure suggests that both hyperparameters need to be relatively high in value, which indicates that MQE needs: (1) a far enough waypoint for value to quickly propagate, and (2) a high enough  $p$  to ensure that local consistencies are being respected. This is promising, as we do not need to decide to decide a trade-off between the two hyperparameters, and these two hyperparameters are in fact not mutually interfering.

We provide additional, more minor ablations in Appendix E.

## 6 CONCLUSION

We introduced Multistep Quasimetric Estimation (MQE), a novel method that combines the benefits of fast value propagation via multistep backup and the global constraint of quasimetric distances. MQE is able to solve extremely challenging and long-horizon tasks in simulated benchmarks as well as in real life.

**Limitations and Future Work.** While MQE achieves strong performance, we sample the waypoints based on empirical analyses. This could incur more computational costs when finding the

Figure 5: Success rate of MQE on humanoidmaze\_giant\_stitch using  $\alpha = 0.01$ , averaged over 4 seeds.

optimal way of sampling the waypoint for environments that are outside of our evaluation range. Future work can investigate the theoretical connection between sampling waypoints and successor distances, investigate the effect of such policy learning on different policy classes such as autoregressive (i.e. transformer) or flow policies, and apply the same method across methods beyond offline RL in scenarios such as offline-to-online RL.

**Ethics Statement** In this work, we studied ways in which we enable a policy to better perform compositional and horizon generalization in reinforcement learning contexts. While our implementations in this paper do not constitute real-world effects due to the smaller networks we have chosen, future works using larger and more expressive policies should pay attention to the potential failure mode of our algorithm and examine the interpretability of the learned policy.

In this work, we ran the text through LLMs to fix grammar and spelling mistakes. We did not use LLMs to generate any content in the paper that were not already ideated and fully drafted by the authors.

**Reproducibility Statement** To ensure reproducibility, we have uploaded our implementation on OGBench as supplementary material, and have stated the relevant hyperparameters needed to reproduce the results in the experiment section.

## REFERENCES

- Philip J. Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient Online Reinforcement Learning With Offline Data. arXiv:2302.02948, 2023. URL <https://arxiv.org/abs/2302.02948>.
- Arindam Banerjee, Srujana Merugu, Inderjit Dhillon, and Joydeep Ghosh. Clustering With Bregman Divergences. In *SIAM International Conference on Data Mining*, pp. 234–245. Society for Industrial and Applied Mathematics, April 2004. ISBN 978-0-89871-568-2 978-1-61197-274-0. doi: 10.1137/1.9781611972740.22. URL <https://pubs.siam.org/doi/10.1137/1.9781611972740.22>.
- Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky.  $\pi_0$ : A Vision-Language-Action Flow Model for General Robot Control. arXiv:2410.24164, 2024. URL <https://arxiv.org/abs/2410.24164>.
- Peter Dayan. Improving Generalisation for Temporal Difference Learning: The Successor Representation. *Neural Computation*, (4):613–624, 1993. doi: 10.1162/neco.1993.5.4.613.
- Danny Driess, Jost Tobias Springenberg, Brian Ichter, Lili Yu, Adrian Li-Bell, Karl Pertsch, Allen Z. Ren, Homer Walke, Quan Vuong, Lucy Xiaoyang Shi, and Sergey Levine. Knowledge Insulating Vision-Language-Action Models: Train Fast, Run Fast, Generalize Better. arXiv:2505.23705, 2025. URL <https://arxiv.org/abs/2505.23705>.
- Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. C-Learning: Learning to Achieve Goals via Recursive Classification. In *International Conference on Learning Representations*. arXiv, 2021. URL <https://arxiv.org/abs/2011.08909>.
- Benjamin Eysenbach, Tianjun Zhang, Ruslan Salakhutdinov, and Sergey Levine. Contrastive Learning as Goal-Conditioned Reinforcement Learning. In *Neural Information Processing Systems*, volume 35, pp. 35603–35620. arXiv, 2022. URL <https://arxiv.org/abs/2206.07568>.
- Benjamin Eysenbach, Vivek Myers, Ruslan Salakhutdinov, and Sergey Levine. Inference via Interpolation: Contrastive Representations Provably Enable Planning and Inference. In *Neural Information Processing Systems*, December 2024. URL <https://arxiv.org/abs/2403.04082>.
- Scott Fujimoto and Shixiang Shane Gu. A Minimalist Approach to Offline Reinforcement Learning. *Neural Information Processing Systems*, 34:20132–20145, 2021. URL <https://arxiv.org/abs/2106.06860>.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning With a Stochastic Actor. arXiv:1801.01290, 2018. URL <https://arxiv.org/abs/1801.01290>.

- 540 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline Reinforcement Learning With Implicit  
 541 Q-Learning. In *International Conference on Learning Representations*. arXiv, 2022. URL  
 542 <http://arxiv.org/abs/2110.06169>.
- 543 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative Q-Learning for  
 544 Offline Reinforcement Learning. arXiv:2006.04779, 2020. URL <https://arxiv.org/abs/2006.04779>.
- 545 Qiyang Li, Zhiyuan Zhou, and Sergey Levine. Reinforcement Learning With Action Chunking.  
 546 arXiv:2507.07969, 2025. URL <https://arxiv.org/abs/2507.07969>.
- 547 Bo Liu, Yihao Feng, Qiang Liu, and Peter Stone. Metric Residual Networks for Sample Efficient  
 548 Goal-Conditioned Reinforcement Learning. arXiv:2208.08133, January 2023. URL <https://arxiv.org/abs/2208.08133>.
- 549 Jianlan Luo, Zheyuan Hu, Charles Xu, You Liang Tan, Jacob Berg, Archit Sharma, Stefan Schaal,  
 550 Chelsea Finn, Abhishek Gupta, and Sergey Levine. SERL: A Software Suite for Sample-Efficient  
 551 Robotic Reinforcement Learning. arXiv:2401.16013, 2025. URL <https://arxiv.org/abs/2401.16013>.
- 552 Max Sobol Mark, Tian Gao, Georgia Gabriela Sampaio, Mohan Kumar Srirama, Archit Sharma,  
 553 Chelsea Finn, and Aviral Kumar. Policy Agnostic RL: Offline RL and Online RL Fine-Tuning of  
 554 Any Class and Backbone. arXiv:2412.06685, December 2024. URL <https://arxiv.org/abs/2412.06685>.
- 555 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wier-  
 556 stra, and Martin Riedmiller. Playing Atari With Deep Reinforcement Learning. arXiv:1312.5602,  
 557 2013. URL <https://arxiv.org/abs/1312.5602>.
- 558 Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc G. Bellemare. Safe and Efficient Off-  
 559 Policy Reinforcement Learning. arXiv:1606.02647, 2016. URL <https://arxiv.org/abs/1606.02647>.
- 560 Vivek Myers, Chongyi Zheng, Anca Dragan, Sergey Levine, and Benjamin Eysenbach. Learning  
 561 Temporal Distances: Contrastive Successor Features Can Provide a Metric Structure for Decision-  
 562 Making. In *International Conference on Machine Learning*, 2024. URL <https://arxiv.org/pdf/2406.17098.pdf>.
- 563 Vivek Myers, Bill Chunyuan Zheng, Anca Dragan, Kuan Fang, and Sergey Levine. Temporal  
 564 Representation Alignment: Successor Features Enable Emergent Compositional in Robot  
 565 Instruction Following. arXiv:2502.05454, 2025a. URL <https://arxiv.org/abs/2502.05454>.
- 566 Vivek Myers, Bill Chunyuan Zheng, Benjamin Eysenbach, and Sergey Levine. Offline Goal Con-  
 567 ditioned Reinforcement Learning With Temporal Distance Representations. In *Neural Infor-  
 568 mation Processing Systems*, 2025b. URL <https://tmd-website.github.io/static/pdf/myers2025offline.pdf>.
- 569 Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-Efficient Hierarchical Reinforce-  
 570 ment Learning. arXiv:1805.08296, 2018. URL <https://arxiv.org/abs/1805.08296>.
- 571 Mitsuhiko Nakamoto, Oier Mees, Aviral Kumar, and Sergey Levine. Steering Your Generalists:  
 572 Improving Robotic Foundation Models via Value Guidance. arXiv:2410.13816, 2025. URL  
 573 <https://arxiv.org/abs/2410.13816>.
- 574 Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham  
 575 Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex  
 576 Herzog, Alex Irpan, Alexander Khazatsky, et al. Open X-Embodiment: Robotic Learning Datasets  
 577 and RT-X Models. In *International Conference on Robotics and Automation*. arXiv, May 2024.  
 578 URL <http://arxiv.org/abs/2310.08864>.
- 579 Seohong Park, Kevin Frans, Sergey Levine, and Aviral Kumar. Is Value Learning Really the Main  
 580 Bottleneck in Offline RL? arXiv:2406.09329, June 2024a. URL <https://arxiv.org/abs/2406.09329>.
- 581 Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. HIQL: Offline Goal-  
 582 Conditioned RL With Latent States as Actions. arXiv:2307.11949, 2024b. URL <https://arxiv.org/abs/2307.11949>.
- 583 Seohong Park, Tobias Kreiman, and Sergey Levine. Foundation Policies With Hilbert Representations.  
 584 arXiv:2402.15567, 2024c. URL <https://arxiv.org/abs/2402.15567>.

- 594 Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. OGBench: Benchmarking  
 595 Offline Goal-Conditioned RL. In *International Conference on Learning Representations*. arXiv,  
 596 2025a. URL <https://arxiv.org/abs/2410.20092>.
- 597 Seohong Park, Kevin Frans, Deepinder Mann, Benjamin Eysenbach, Aviral Kumar, and Sergey  
 598 Levine. Horizon Reduction Makes RL Scalable. arXiv:2506.04168, 2025b. URL <https://arxiv.org/abs/2506.04168>.
- 600 Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. FiLM: Visual  
 601 Reasoning With a General Conditioning Layer. arXiv:1709.07871, 2017. URL <https://arxiv.org/abs/1709.07871>.
- 604 Silviu Pitis, Harris Chan, Kiarash Jamali, and Jimmy Ba. An Inductive Bias for Distances: Neural  
 605 Nets That Respect the Triangle Inequality. arXiv:2002.05825, 2020. URL <https://arxiv.org/abs/2002.05825>.
- 607 Younggyo Seo, Carmelo Sferrazza, Haoran Geng, Michal Nauman, Zhao-Heng Yin, and Pieter  
 608 Abbeel. FastTD3: Simple, Fast, and Capable Reinforcement Learning for Humanoid Control.  
 609 arXiv:2505.22642, 2025. URL <https://arxiv.org/abs/2505.22642>.
- 610 Laura Smith, Yunhao Cao, and Sergey Levine. Grow Your Limits: Continuous Improvement With  
 611 Real-World RL for Robotic Locomotion. arXiv:2310.17634, 2023. URL <https://arxiv.org/abs/2310.17634>.
- 613 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press,  
 614 second edition, 2018. URL <http://incompleteideas.net/book/the-book-2nd.html>.
- 616 Dong Tian, Ge Li, Hongyi Zhou, Onur Celik, and Gerhard Neumann. Chunking the Critic: A  
 617 Transformer-Based Soft Actor-Critic With N-Step Returns. arXiv:2503.03660, 2025. URL  
 618 <https://arxiv.org/abs/2503.03660>.
- 619 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation Learning With Contrastive  
 620 Predictive Coding. arXiv:1807.03748, 2019. URL <https://arxiv.org/abs/1807.03748>.
- 621 Andrew Wagenmaker, Mitsuhiro Nakamoto, Yunchu Zhang, Seohong Park, Waleed Yagoub, Anusha  
 622 Nagabandi, Abhishek Gupta, and Sergey Levine. Steering Your Diffusion Policy With Latent  
 623 Space Reinforcement Learning. arXiv:2506.15799, 2025. URL <https://arxiv.org/abs/2506.15799>.
- 625 Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyi Zheng, Tony Zhao,  
 626 Philippe Hansen-Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang, Chelsea Finn, and  
 627 Sergey Levine. BridgeData V2: A Dataset for Robot Learning at Scale. arXiv:2308.12952, 2024.  
 628 URL <https://arxiv.org/abs/2308.12952>.
- 629 Tongzhou Wang, Antonio Torralba, Phillip Isola, and Amy Zhang. Optimal Goal-Reaching  
 630 Reinforcement Learning via Quasimetric Learning. arXiv:2304.01203, 2023. URL <https://arxiv.org/abs/2304.01203>.
- 632 C. J. C. H. Watkins. Learning From Delayed Rewards, 1989.
- 634 Charles Xu, Qiyang Li, Jianlan Luo, and Sergey Levine. RLDG: Robotic Generalist Policy Distillation  
 635 via Reinforcement Learning. arXiv:2412.09858, 2024. URL <https://arxiv.org/abs/2412.09858>.
- 637 Chongyi Zheng, Benjamin Eysenbach, Homer Walke, Patrick Yin, Kuan Fang, Ruslan Salakhutdinov,  
 638 and Sergey Levine. Stabilizing Contrastive RL: Techniques for Robotic Goal Reaching From  
 639 Offline Data. In *International Conference on Learning Representations*. arXiv, 2024. URL  
 640 <https://arxiv.org/abs/2306.03346>.
- 641
- 642
- 643
- 644
- A WEBSITE AND CODE
- 645
- 646 We provide the full implementation of MQE, TMD, CMD, and the new antmaze-colossal mazes  
 647 on OGBench <https://anonymous.4open.science/r/ogbench-D4D9/>. We provide website for MQE at  
<https://anonymous.4open.science/w/mqe-paper-686D/>.

648    **B OGBENCH EXPERIMENT DETAILS**  
 649

650    **B.1 BASELINE DETAILS**  
 651

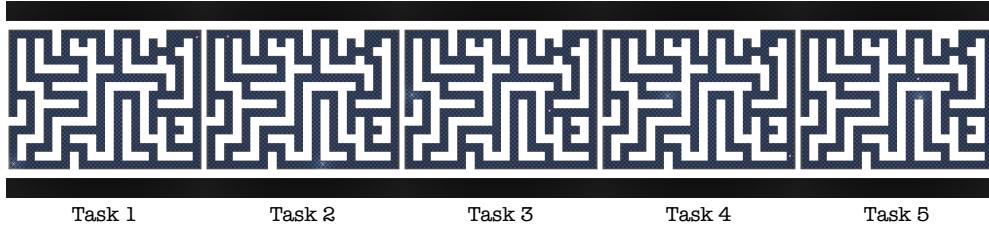
652    Here, we briefly describe the inner workings of each baseline on OGBench.

653    **GCIQL** Goal-conditioned Implicit Q-learning (GCIQL) uses expectile regression to learn a value  
 654    function. **QRL** Quasimetric RL learns a quasimetric distance using bootstrapping under a quasimetric  
 655    architecture, and constrains the distance with one-step cost in *deterministic* settings. **CRL** Contrastive  
 656    RL (CRL) uses binary cross entropy to regress the critic (defined as a dot product) towards goals  
 657    that are future states and repel those that are not. **CMD** Contrastive Metric Distillation (CMD) uses  
 658    InfoNCE loss (van den Oord et al., 2019) to recover  $Q_\beta$ . **TMD** Temporal Metric Distillation (TMD)  
 659    (Myers et al., 2025b) use contrastive learning to learn the behavior Q-function, and then tightens  
 660    the bound to optimize towards  $Q^*$ . CMD and TMD enforce the quasimetric property of successor  
 661    distance architecturally.

662    Additionally, we compare MQE against two horizon reduction methods, n-step Goal-Conditioned  
 663    Soft Actor-Critic with Behavior Cloning (n-SAC+BC) and Hierarchical Implicit Q learning (HIQL),  
 664    which explicitly reduce the policy and value horizon separately. **n-SAC+BC** n-SAC+BC is equivalent  
 665    to SAC+BC, but with  $n$ -step updates. **HIQL** trains the same value function as GCIQL, but the  
 666    agent extracts a hierarchical policy. All policies trained on OGBench are designed to return a  
 667    multimodal Gaussian distribution  $\mathcal{N}(\mu; \Sigma)$ , and during inference time, the neural network produces  
 668    the distribution, and the policy samples from that as action.

669  
 670    **B.2 TASK VISUALIZATION**

671    We provide the tasks used for `antmaze-colossal` environment on Fig. 6. The maze itself is 24  
 672    blocks in height and 18 blocks in width, which is 50% larger than the giant mazes in each dimension.  
 673



674  
 675  
 676  
 677    Figure 6: Task visualizations from `antmaze-colossal` environment. The ant occupies the starting  
 678    position, and it must reach the red dot to complete the task.  
 679  
 680

681  
 682    **B.3 IMPLEMENTATION DETAILS & HYPERPARAMETERS**  
 683

684    Table 2 details the common hyperparameters for all methods on OGBench. Table 3

685    For visual environments on OGBench, we define the ResNet encoder for state-action representation  
 686    and state representation  $f_\varphi(s, a)$ ,  $f_\psi(s)$ . Empirically, we find that when calculating the  $\mathcal{L}_\mathcal{I}$  loss, we  
 687    freeze gradients for the ResNet encoder, in a manner similar to (Driess et al., 2025). The action  
 688    invariance loss objective now becomes:

689  
 690  
 691    
$$\mathcal{L}_\mathcal{I} (\varphi, \psi; \{s_i, a_i\}_{i=1}^N) = \sum_{i,j=1}^N (d(\psi(\bar{f}_\psi(s_i)), \varphi(\bar{f}_\varphi(s_i), a_j)))^2. \quad (17)$$
  
 692  
 693

694    where  $\bar{f}$  denotes freezing the gradient. We find this to be helpful to not let the learned representations  
 695    collapse during the beginning of training, and that it maintained a good level of generalization ability  
 696    during training.

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
Table 2: Network configuration for MQE on OGBench.

Configuration	Value
batch size	256
latent dimension size	512
encoder MLP dimensions	(512, 512, 512)
policy MLP dimensions	(512, 512, 512)
layer norm in encoder MLPs	True
visual encoder (visual- envs)	impala-small
MRN components	8
Discount ( $\gamma$ )	0.995
Waypoint discount ( $\lambda$ )	0.95
Probability $p$ for directly imposing $\mathcal{L}_T$	0.2

Table 3: BC coefficient  $\alpha$  for each environment

Environment	$\alpha$
antmaze-navigate	0.1
antmaze-stitch	0.03
antmaze-explore	0.003
humanoidmaze	0.01
pointmaze	0.03
{cube, puzzle, scene (play)}	1.0
{visual - {cube, puzzle, scene (play)}}	3.0
*-noisy	0.1
visual-antmaze	0.3

We use a batch size of 256 for MQE, which, when being trained on an A6000 GPU, takes around 2.5 hours to finish both training and evaluation.

#### B.4 POLICY EVALUATION

We maintain an evaluation procedure similar to that of OGBench’s. We evaluate each policy 50 times the last 3 training epochs (800k, 900k, 1M step for state-based environments, 300k, 400k, 500k in pixel-based environments). For each of these evaluation epochs, use 8 seeds for every state-based environment and 4 seeds for every pixel-based environment.

#### B.5 FULL RESULT TABLE

We record the full success rate of all tasks across OGBench in Table 1.

### C BRIDGEDATA EXPERIMENT DETAILS

We evaluate each policy with a total of 15 trials each for every task. For each of the pick and place tasks, we assign 1 point for each successful pick and place (i.e. move a desired object to the desired location). Therefore, a policy can earn a maximum of  $i$  points for each PnP task that manipulates  $i$  objects.

For open drawer and place, we assign one point if the policy opens the drawer to the extent where the mushroom can be placed, but does not pull the drawer off the base completely, and assign another point if the policy is able to put the mushroom in the drawer.

756  
757  
758 Table 4: Network configuration for MQE on BridgeData.  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768

Configuration	Value
latent dimension size	256
encoder MLP dimensions	(256, 256, 256)
policy MLP dimensions	(256, 256, 256)
layer norm in encoder MLPs	True
MRN components	8
Discount ( $\gamma$ )	0.98
Waypoint discount ( $\lambda$ )	0.95
Probability $p$ for directly imposing $\mathcal{L}_T$	0.2

769 C.1 TRAINING CONFIGURATIONS  
770

771 We use a pretrained ResNet34 as the backbone of the policy and encoders. We do not share the  
 772 actor and critic encoder for both GCIQL and MQE. We co-train the actor and the critic for a total of  
 773 500,000 total steps with a batch size of 128, which takes around 60 hours when the model is trained  
 774 on 4 TPU pods. For TRA-g, we produce the embedding of two separate encoders  $\phi, \psi$ , and align each  
 775 other using symmetric InfoNCE loss. During inference, we use FiLM (Perez et al., 2017) to embed  
 776 the learned goal representation into the actor, staying consistent with the implementation from (Myers  
 777 et al., 2025a). In addition, we also describe the common hyperparameters used for BridgeData setup  
 778 at Table 4.

779 C.2 DEFINING TASK SUCCESSES  
780

781 Table 5 records the success rate of each task across all six tasks that we evaluate. Unlike Section 5.2,  
 782 we only measure whether each task has been completed to its fullest. While it does give out a stronger  
 783 signal on whether a task displays nonzero success rate, it does not provide as much information on  
 784 how the task was progressing overall.

785  
786 Table 5: Binary success counts for each task.  
787

Task	MQE (Ours)	GCBC	GCIQL	TRA
Single PnP	12/15	5/15	10/15	9/15
Double PnP	10/15	1/15	4/15	6/15
Triple PnP	4/15	0/15	0/15	1/15
Quadruple PnP	2/15	0/15	0/15	0/15
Open Drawer	12/15	10/15	5/15	11/15
Open Drawer & Place	5/15	0/15	0/15	1/15

797 D POLICY ROLLOUT IN BRIDGEDATA  
798

800 We provide the rollout of triple PnP. We especially consider TRA-g because it also exhibits  
 801 compositional generalization, yet there is no explicit policy improvement as compared to offline RL  
 802 methods such as MQE.

803 E ADDITIONAL ABLATIONS  
804

805 **What is the best distribution to sample  $s_t^w$ ?** We consider three different ways of sampling  
 806 waypoint  $s_t^w$  in our approach: (1) uniformly sampling between state and goal. (2) uniformly sampling  
 807 from Unif[1, ..., 50], and (3) we set  $k' = 50$ . We note that method (3) directly transforms our method  
 808 into performing n-step return within the quasimetric architecture, and that it also reduces the value  
 809 horizon due to sampling a fixed number of states ahead.

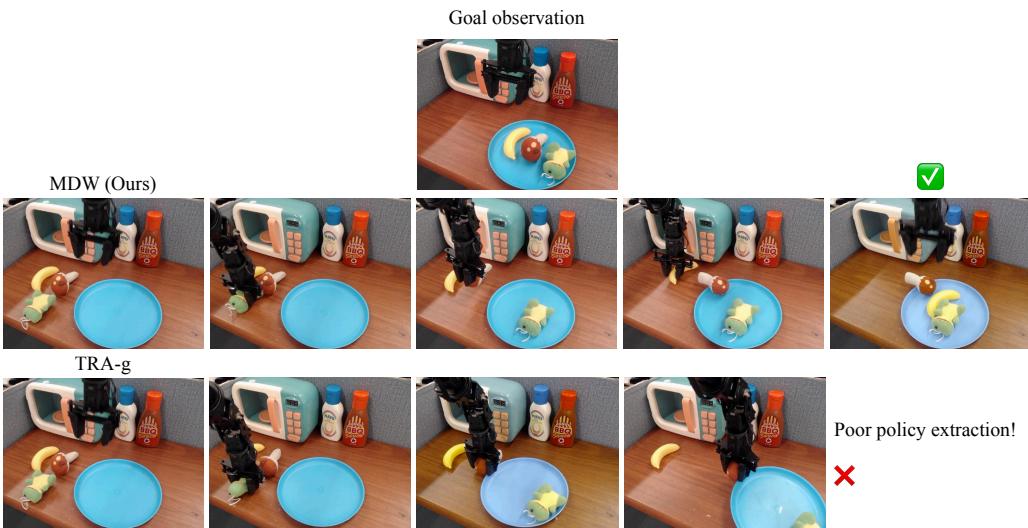


Figure 7: Inference from triple pnp task. We note that due to poor policy extraction and generalization, TRA-g is not able to complete the task.

Table 6 demonstrates that using a geometric distribution is the optimal way to sample such waypoints, and explicit horizon reduction techniques actually *hurt* the performance of MQE, as it learns a much less robust policy compared to random sampling.

Table 6:  $s_t^w$  Sampling Ablation

Configuration	Success Rate (%)
$k' \sim \text{Geom}(1 - \lambda)$	<b>26.5</b>
$k' \sim \text{Unif}[1, K]$	18.9
$k' \sim \text{Unif}[1, 50]$	17.8
$k' = 50$	1.7

**What role does  $\mathcal{L}_{\mathcal{I}}$  play?** We now ablate away  $\mathcal{L}_{\mathcal{I}}$  by removing line 6 of Algorithm 1, we examine whether the multistep value propagation is enough to learn a good policy.

Table 7:  $\mathcal{L}_{\mathcal{I}}$  Ablation

Configuration	Success Rate (%)
With $\mathcal{L}_{\mathcal{I}}$	<b>26.5</b>
Without $\mathcal{L}_{\mathcal{I}}$	7.9

Table 7 shows that  $\mathcal{L}_{\mathcal{I}}$  is in fact necessary, as without this invariance loss, the performance of MQE degrades by a large margin.