

---

# Meta-GC-TTT: Training Offline Goal-Conditioned Policies for Test-Time Adaptation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In offline goal-conditioned reinforcement learning, Test-Time Training (TTT) can  
2 specialize a pre-trained policy to the current state and goal at deployment. This  
3 turns a broad goal-conditioned policy into a query-specific expert. Yet, standard  
4 offline pre-training optimizes the policy *before* TTT, not *after* it. As a result, the  
5 policy is not trained for the gradient dynamics it will face at test time. We introduce  
6 META-GC-TTT, a framework for learning test-time-trainable goal-conditioned  
7 policies. META-GC-TTT samples state-goal tasks from the offline dataset, adapts  
8 the policy with TTT, and updates the base policy for post-TTT performance. Our  
9 evaluation on the OGBench loco-navigation suite demonstrates that meta-learned  
10 initializations significantly improve zero-shot performance and achieve a  $(3 - 5\times)$   
11 increase in adaptation efficiency. Overall, offline goal-conditioned policies should  
12 not only be trained not only to act, but also to adapt.

## 13 1 Introduction

14 The standard paradigm of modern machine learning has shifted from training task-specific models  
15 from scratch to a two-phase process: *large-scale pre-training* followed by *specialized adaptation*.  
16 This evolution, exemplified by Foundation Models in language and vision [Radford et al., 2018,  
17 Touvron et al., 2023, Ouyang et al., 2022, Yang et al., 2025, Team et al., 2025, Radford et al., 2021,  
18 Dosovitskiy et al., 2021, Beyrer et al., 2024], leverages vast unlabeled datasets and self-supervised  
19 objectives to learn underlying data structures.

20 However, this paradigm has only recently begun to transition to Reinforcement Learning (RL) for  
21 robotics. Emerging research proposed Vision-Language-Action (VLA) models [Zitkovich et al.,  
22 2023, Kim et al., 2024, Black et al., 2026, Intelligence et al., 2025] and *offline goal-conditioned RL*  
23 (*GCRL*) [Liu et al., 2022, Eysenbach et al., 2022b] as promising candidates for robotic foundation  
24 models. In particular, GCRL offers two pivotal advantages: (i) *offline scalability*, which leverages  
25 massive, diverse datasets while bypassing the safety and cost constraints of online interaction; and  
26 (ii) *self-supervised objectives*, which provide a reward-free training signal analogous to next-token  
27 prediction in LLMs. By treating reached states as goals, this objective enables the extraction of  
28 general-purpose skills directly from unlabeled trajectories.

29 Yet, scaling alone may be insufficient for high-performance agents. While pre-training provides  
30 a broad “common sense,” complex tasks require a final specialization phase. This paradigm has  
31 recently been pushed to its logical extreme with *Test-Time Training (TTT)* [Sun et al., 2020, Hardt  
32 and Sun, 2024], where a model is adapted on-the-fly for each specific test instance. By performing  
33 gradient descent on data relevant to the current query at the moment of inference, TTT transforms a  
34 broad generalist model into a localized expert tailored to the immediate task.

35 At its core, this work connects goal-conditioned pre-training with test-time adaptation through *meta-*  
36 *learning*. If a model is intended to be adaptive at test-time, the training phase should explicitly

37 optimize for this capability. Meta-learning algorithms [Finn et al., 2017, Nichol et al., 2018] provide  
38 the framework to ensure a model may truly *learn to learn*, optimizing the initialization specifically  
39 for its ability to specialize via local gradients.

40 The main contributions of this paper are as follows:

- 41 • We introduce META-GC-TTT, a unified post-training framework that leverages meta-  
42 learning to optimize policies for test-time adaptation, instantiating three distinct algorithms:  
43 CONTINUAL-TTT, TT-FOMAML and TT-REPTILE.
- 44 • We propose a novel meta-update mechanism using an early-stopping approach based on  
45 goal-conditioned test loss. This dynamically determines the inner-loop horizon to prevent  
46 overfitting and stabilize adaptation.
- 47 • We provide an extensive evaluation across the loco-navigation suite of OGBench. Our results  
48 demonstrate that meta-learned initializations significantly improve zero-shot performance  
49 and achieve up to a  $5\times$  increase in adaptation efficiency compared to standard TTT.

## 50 2 Related Work

51 **Offline Reinforcement Learning** Offline RL focuses on learning optimal policies from static  
52 datasets without environment interaction [Ernst et al., 2005, Levine et al., 2020], and represents  
53 a promising paradigm for training robotics foundation models [Amin et al., 2025]. The primary  
54 challenge is the distributional shift between the behavior policy  $\pi_\beta$  and the learned policy  $\pi_\theta$ , which  
55 causes catastrophic value overestimation on out-of-distribution (OOD) actions. Conservative methods  
56 like CQL [Kumar et al., 2020] and in-sample methods like IQL [Kostrikov et al., 2021] have been  
57 developed to mitigate these issues. Such techniques are in practice deployed for the underlying policy  
58 through this work.

59 **Goal-Conditioned RL** GCRL enables agents to reach diverse targets by training goal-conditioned  
60 policies [Liu et al., 2022, Eysenbach et al., 2022b, Ma et al., 2022], and can be applied naturally to  
61 offline settings [Agarwal et al., 2023], as most objectives are self-supervised [Andrychowicz et al.,  
62 2017, Tian et al., 2021], and do not require reward labels. Despite the success of self-supervised  
63 objectives, evidence suggests that effective reasoning horizon remains a primary bottleneck for  
64 scaling these models to complex tasks [Park et al., 2025b]. We utilize the OGBench suite [Park et al.,  
65 2025a] to evaluate how META-GC-TTT reduces this effective horizon through localized adaptation.

66 **Test-Time Training (TTT)** The TTT paradigm involves adapting model parameters on-the-fly  
67 for each specific test instance [Sun et al., 2020]. In decision-making, TTT transforms a generalist  
68 model into a localized expert [Park et al., 2024, Bagatella et al., 2025]. Our work extends the GC-  
69 TTT framework [Bagatella et al., 2025] by aligning the initialization specifically for inference-time  
70 gradient dynamics.

71 **Optimization-based Meta-Learning** Meta-learning aims to learn representations that facilitate rapid  
72 adaptation [Finn et al., 2017]. Algorithms like MAML optimize a model’s initialization specifically  
73 for post-adaptation performance, while Reptile merges expert parameters with a base model via  
74 interpolation [Nichol et al., 2018]. Multiple variations of these methods have been proposed [Ra-  
75 jeswaran et al., 2019, Antoniou et al., 2018, Finn et al., 2018, Liu et al., 2019]. Our framework,  
76 META-GC-TTT, differentiates itself by integrating meta-learning with an efficient TTT inner loop.

## 77 3 Background

### 78 3.1 Goal-Conditioned Offline Reinforcement Learning

79 Goal-Conditioned Reinforcement Learning (GCRL) may be formalized through a goal-augmented  
80 Markov decision process (G-MDP) defined by the tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{G}, P, R, \gamma)$ , where  $\mathcal{S}$  is the state space,  
81  $\mathcal{A}$  is the action space,  $\mathcal{G}$  is the goal space (often  $\mathcal{G} \subseteq \mathcal{S}$ ),  $P(s'|s, a)$  represents goal-independent  
82 transition dynamics,  $R(s, g)$  is the goal-conditioned reward function (typically sparse:  $R(s, g) =$   
83  $\mathbf{1}_{d(s,g) \leq \epsilon}$ ), and  $\gamma \in [0, 1)$  is the discount factor. The objective is to learn a policy  $\pi(a|s, g)$  maximizing

84 an expectation over the *value function*:

$$V^\pi(s|g) = \mathbb{E}_{P,\pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, g) \mid s_0 = s \right]. \quad (1)$$

85 An state-action value function may be defined similarly:  $Q^\pi(s, a|g) =$   
 86  $\mathbb{E}_{P,\pi} [\sum_{t=0}^{\infty} \gamma^t R(s_t, g) \mid s_0 = s, a_0 = a]$ .

87 In the offline setting,  $\pi$  needs to be estimated from a static dataset  $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$  collected  
 88 by an unknown behavior policy  $\pi_\beta$ . A multitude of algorithms has been proposed for this setting  
 89 [Park et al., 2023, Zeng et al., 2023, Zheng et al., 2023, Eysenbach et al., 2022a, Zhou and Kao,  
 90 2026], among which we will consider two representative instances. Among supervised approaches,  
 91 Goal-Conditioned Behavioral Cloning (GCBC) [Yang et al., 2022] trains a policy parameterized  
 92 by  $\theta$  by maximizing the likelihood of actions taken in the dataset:

$$\mathcal{L}_{BC}(\theta) = \mathbb{E}_{(s,a,g) \sim \mathcal{D}_g} [-\log \pi_\theta(a|s, g)] \quad (2)$$

93 While stable, GCBC is limited by the quality of  $\pi_\beta$ . Reinforcement learning algorithms instead  
 94 rely on the estimation of  $V^\pi$  for improving the policy. Among those, Goal-Conditioned Implicit  
 95 Q-Learning (GCIQL) [Kostrikov et al., 2021, Zeng et al., 2023] trains value estimates through  
 96 expert regression and TD learning, and updates the policy  $\pi_\theta(a|s, g)$  via either Advantage  
 97 Weighted Regression (AWR) [Peng et al., 1910]:

$$\mathcal{L}_{AWR}(\theta) = \mathbb{E}_{(s,a,g) \sim \mathcal{D}_g} [-\exp(\beta A(s, a, g)) \log \pi_\theta(a|s, g)] \quad (3)$$

98 or DDPG-BC [Lillicrap et al., 2019, Fujimoto and Gu, 2021], which regularizes a deterministic  
 99 policy gradient with a behavioral cloning term:

$$\mathcal{L}_{DDPG-BC}(\theta) = \mathbb{E}_{(s,a,g) \sim \mathcal{D}_g} [-Q(s, \hat{a}, g) + \alpha \|\hat{a} - a\|^2] \text{ where } \hat{a} \sim \pi_\theta(\cdot|s, g). \quad (4)$$

100 In either case, the policy is in practice a neural network updated through first-order methods, making  
 101 it suitable for gradient-based test-time adaptation.

### 102 3.2 Goal-Conditioned Test-time Training (GC-TTT)

103 Test-Time Training (TTT) is a model adaptation paradigm where a base model is fine-tuned for each  
 104 specific test instance  $\tau$  during deployment [Sun et al., 2020, Hardt and Sun, 2024]. Given base model  
 105 parameters  $\theta_0$  and a task-specific support set  $\mathcal{D}_\tau$ , TTT utilizes a transformation operator  $U_\tau^K$  to  
 106 produce specialized *expert* parameters (or “fast weights”)  $\theta^{TTT} = U_\tau^K(\theta_0)$ . This operator typically  
 107 consists of  $K$  steps of stochastic gradient descent on some task-specific objective  $\mathcal{L}_\tau$ :

$$\theta_k \leftarrow \theta_{k-1} - \alpha \nabla_{\theta_{k-1}} \mathcal{L}_\tau(\theta_{k-1}, \mathcal{D}_\tau) \quad (5)$$

108 where  $\alpha$  is a learning rate. In the context of GCRL, Bagatella et al. [2025] recently demonstrated  
 109 that TTT can specialize goal-conditioned policies to specific test instances (in this case, state-goal  
 110 pairs  $(s, g) = \tau$ ). Given a base policy  $\pi_{\theta_0}$  and a state-goal pair  $(s, g)$ , GC-TTT (i) describes a  
 111 self-supervised procedure `filter`( $\mathcal{D}, s, g$ ) to extract an instance-specific dataset  $\mathcal{D}_{(s,g)} = \mathcal{D}_\tau$  from  
 112 an unlabeled pre-training dataset  $\mathcal{D}$  and (ii) simply applies the transformation operator  $U_{(s,g)}^K$  on top  
 113 of some underlying GCRL objective  $\mathcal{L}_{(s,g)}$ . We refer to Appendix A for more details.

### 114 3.3 Optimization-Based Meta-Learning

115 Optimization-based meta-learning seeks to find an initialization  $\theta$  that minimizes the “learning  
 116 distance” to future tasks  $\tau$ , ensuring the model is easy to fine-tune. Model-Agnostic Meta-Learning  
 117 (MAML) [Finn et al., 2017] explicitly optimizes for the performance of the model *after* adaptation.  
 118 The meta-objective finds an initialization  $\theta^\circ$  such that the parameters resulting from the inner-loop  
 119 adaptation  $U_g^K(\theta)$  perform optimally on a held-out query set  $\mathcal{D}_\tau^{test}$ :

$$\theta^\circ = \arg \min_{\theta \in \Theta} \mathbb{E}_\tau [\mathcal{L}_\tau(U_\tau^K(\theta), \mathcal{D}_\tau^{test})] \quad (6)$$

120 This induces a bilevel optimization where the outer loop updates  $\theta$  by back-propagating through  
 121 the gradient descent path of the inner loop. To avoid the computational cost of the second-order

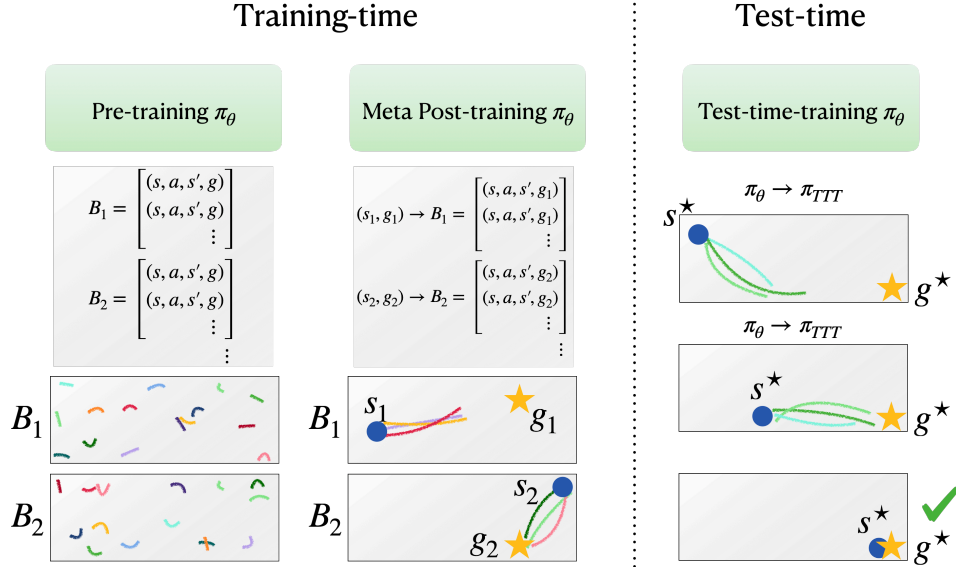


Figure 1: META-GC-TTT is an offline post-training scheme, based on meta-learning, to adapt a pre-trained policy and prepare it for test-time training. During post-training, we perform a meta-optimization loop by sampling starting states and goals as task instances, fine-tuning local copies of the policy on relevant batches of data and using a meta-optimization rule for back-propagating the knowledge of fine-tuned models into the base policy.

122 derivatives required by MAML, First-Order MAML (FOMAML) [Finn et al., 2017] approximates  
 123 the meta-gradient by treating the inner-loop update as a constant:

$$\theta \leftarrow \theta - \beta \cdot \nabla_{U_\tau^K(\theta)} \mathcal{L}_\tau(U_\tau^K(\theta)), \quad (7)$$

124 where  $\beta$  is a meta-learning rate. This assumes a locally linear loss landscape, making it memory-  
 125 efficient while remaining effective for small inner-loop learning rates. Reptile [Nichol et al., 2018] is  
 126 a first-order algorithm that simplifies the meta-update rule by performing model merging with the  
 127 adapted weights:

$$\theta \leftarrow \theta + \epsilon(U_{\mathcal{D}_\tau^K}(\theta) - \theta) \quad (8)$$

128 where  $\epsilon$  is a step size. Unlike MAML, it does not require an explicit query/support split for the  
 129 meta-gradient and is claimed to find an initialization that is geometrically close to the manifolds of  
 130 optimal solutions for the task distribution.

#### 131 4 META-GC-TTT

132 We propose META-GC-TTT, a post-training framework designed to align pre-trained offline GCRL  
 133 policies for efficient test-time adaptation, schematized in Figure 1. Through meta-learning, we trans-  
 134 form the policy into an *adaptive process*. Rather than only maximizing zero-shot performance with  
 135 standard pre-training objectives while remaining agnostic of test-time processes, we further optimize  
 136 the policy parameters  $\theta$  to be “ready for TT gradients,” ensuring that a small number of steps  $K$  on  
 137 a local goal-specific dataset  $\mathcal{D}_{(s,g)}$  can shift the behavior from a general prior to a specialized expert.

138 Given a goal-conditioned policy  $\pi_{\theta_0}$  pre-trained through an underlying RL algorithm (e.g. GCBC or  
 139 GCIQL) and an offline dataset  $\mathcal{D}$ , we can instantiate META-GC-TTT as a bilevel optimization routine  
 140 (Algorithm 1). In the inner loop, we leverage GC-TTT [Bagatella et al., 2025]: given a test-time  
 141 instance as a state-goal tuple  $(s, g)$ , GC-TTT filters the dataset  $\mathcal{D}$  to extract task-specific batches  
 142  $\mathcal{D}_{(s,g)}$ , and quickly adapts the policy through a few gradient descent steps. As the adaptation procedure  
 143 is gradient-based, we *amortize* it through an outer meta-learning loop, as described in Algorithm 1.

144 Algorithm 1 involves several routines, namely:

- 145 1. *Task sampling*: A meta-batch of  $B$  tasks, each composed of a state-goal pair  $(s_i, g_i)$ , is  
 146 sampled uniformly from the dataset  $\mathcal{D}$ .

---

**Algorithm 1** META-GC-TTT Post-training

---

**Require:** Pre-trained policy  $\pi_{\theta_0}$ , Dataset  $\mathcal{D}$ , rates  $\alpha, \beta$ , outer batch size  $B$ , outer steps  $N$  and inner steps  $K$

- 1:  $\theta \leftarrow \theta_0$
- 2: **for**  $n = 1$  **to**  $N$  **do** ▷ Outer loop
- 3:   Sample meta-batch of tasks  $(s_i, g_i)_{i=1}^B$  ▷ Task sampling
- 4:   **for**  $i = 1$  **to**  $B$  **do**
- 5:      $\mathcal{D}_{(s_i, g_i)} \leftarrow \text{filter}(\mathcal{D}, s_i, g_i)$  ▷ Data filtering
- 6:      $\mathcal{D}_{(s_i, g_i)}^{train}, \mathcal{D}_{(s_i, g_i)}^{test} \leftarrow \text{split}(\mathcal{D}_{(s_i, g_i)})$
- 7:      $\theta'_i \leftarrow \theta$
- 8:     **for**  $k = 1$  **to**  $K$  **do** ▷ Inner Loop: Adaptation with  $U_{(s_i, g_i)}^K$
- 9:        $\theta'_i \leftarrow \theta'_i - \alpha \nabla \mathcal{L}_{(s_i, g_i)}(\theta'_i, \text{sample}(\mathcal{D}_{(s_i, g_i)}^{train}))$
- 10:     **end for**
- 11:   **end for**
- 12:    $\theta \leftarrow \text{meta\_update}(\theta, \{\theta'_i\}, \{\mathcal{D}_i^{test}\}, \beta)$  ▷ Meta-update
- 13: **end for**

---

- 147   2. *Data filtering*: We follow the same procedure as GC-TTT (see Appendix A), and obtain a  
148   task-specific dataset  $\mathcal{D}_{(s_i, g_i)}$ .
- 149   3. The *inner loop* consists of  $K$  gradient steps of an offline GCRL loss on batches sampled  
150   from the training split, using SGD or other optimizers [Kingma and Ba, 2017], resulting in a  
151   specialized set of parameters  $\theta'_i$ .
- 152   4. The *meta-update* combines the specialized parameters  $\theta'_i$  with the base model  $\theta$ . We  
153   instantiate three algorithms based on the meta-update rule:
- 154   • **CONTINUAL-TTT**:  $\theta \leftarrow \frac{1}{B} \sum_{i=1}^B \theta'_i$ , updating the base parameters to the average of  
155   specialized parameters.
  - 156   • **TT-FOMAML**:  $\theta \leftarrow \theta - \beta \frac{1}{B} \sum_{i=1}^B \nabla_{\theta'_i} \mathcal{L}(\theta'_i, \mathcal{D}_i^{test})$ , using first-order gradients of  
157   the test loss [Finn et al., 2017].
  - 158   • **TT-REPTILE**:  $\theta \leftarrow \theta + \beta (\frac{1}{B} \sum_{i=1}^B (\theta'_i - \theta))$ , merging specialized parameters with  
159   base parameters via interpolation [Nichol et al., 2018].

160 **Test-Aware Meta-Update (Early-Stopping)** Standard meta-update procedures assume that  $K$   
161 gradient steps consistently yield an expert model superior to the base initialization. In practice,  
162 however, fixed-length adaptation risks inner-loop overfitting, vanishing gradients, or the accumulation  
163 of approximation errors [Antoniou et al., 2018]. To mitigate this, we additionally introduce a *test-*  
164 *aware* meta-update strategy that utilizes the test set  $\mathcal{D}_{(s_i, g_i)}^{test}$  as an inner-loop validation mechanism.  
165 We monitor the test loss at each step  $k \in [0, K]$  and perform the meta-update using the optimal  
166 checkpoint:

$$k^* = \arg \min_{k \in \{0, \dots, K\}} \mathcal{L}_{\mathcal{D}_{(s, g)}^{test}}(U_{(s, g)}^k(\theta)). \quad (9)$$

167 This approach provides two primary advantages. First, it *serves as a diagnostic tool* for the meta-  
168 learning process:  $k^* = 0$  prevents updates from poor specialization by rejecting the step,  $k^* = 1$  re-  
169 covers standard stochastic gradient descent (SGD) for TT-REPTILE and CONTINUAL-TTT, and  $k^* >$   
170 1 confirms that the initialization has acquired a structure capable of capturing complex dependencies  
171 through extended optimization. Second, it *enhances hyperparameter robustness* by eliminating the  
172 need to tune  $K$  as a fixed parameter. By setting  $K$  to a large upper bound, the mechanism dynamically  
173 determines the optimal adaptation length per task, facilitating autonomous deployment. While test loss  
174 remains a proxy for final performance, we evaluate the effectiveness of early stopping in Section 5.2.

## 175 5 Experiments

176 Having described a framework for post-training goal-conditioned policies, and aligning them to  
177 test-time adaptation, we now turn to a broad empirical evaluation. Following [Bagatella et al.,  
178 2025], we will rely on a suite of goal-conditioned loco-navigation tasks from OGBench [Park et al.,  
179 2025a], each of which defines five goal-conditioned tasks as shown in Figure 2. We will cover

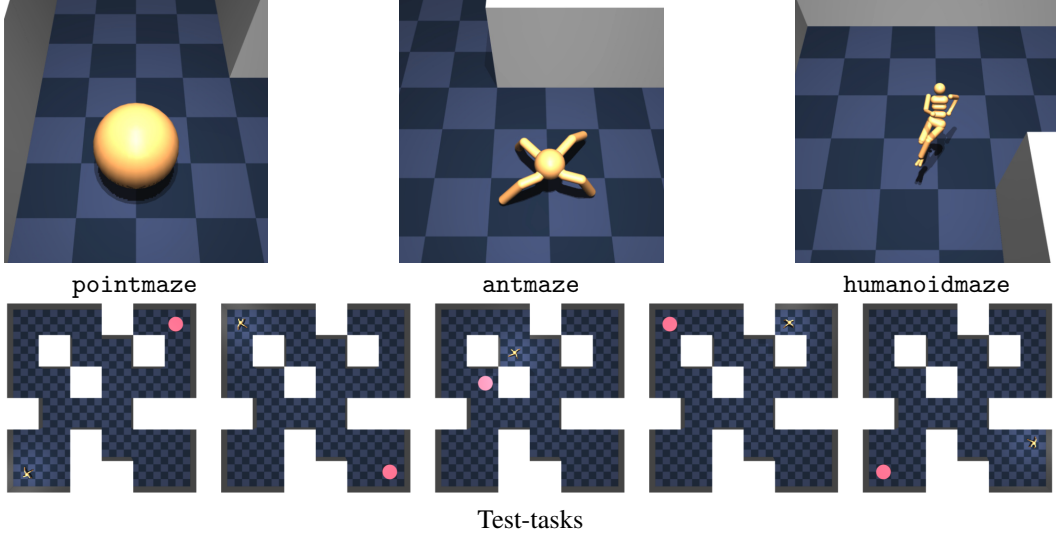


Figure 2: (Top) Maze environments from OGBench [Park et al., 2025a] representing increasing locomotion difficulty. (Bottom) The specific test-tasks used for evaluation.

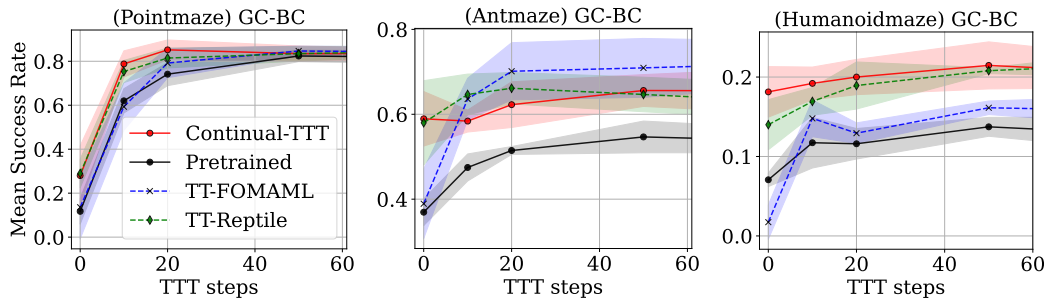


Figure 3: Meta-post-training improves zero-shot initialization and adaptation on high-quality datasets. Meta-learning variants consistently outperform standard GC-TTT baselines. TT-REPTILE and CONTINUAL-TTT provide superior starting success while TT-FOMAML optimizes for steeper recovery and higher peak performance despite initial zero-shot "dips."

180 three embodiments (pointmaze, antmaze and humanoidmaze) and two data sources: long, expert  
 181 demonstrations (navigate) as well as short, disjoint trajectory segments (stitch).

182 In each environment, we apply META-GC-TTT on top of two representative offline GCRL algorithms,  
 183 namely GCBC [Yang et al., 2022] (which is applied to expert navigate data) and GCIQL [Kostrikov  
 184 et al., 2021], which we evaluate on stitch data due to its off-policy nature. We pre-train each  
 185 algorithm for 400k steps (i.e., until convergence) to provide stable initializations for a subsequent 50k-  
 186 step meta-post-training phase using TT-REPTILE, TT-FOMAML, or CONTINUAL-TTT. After post-  
 187 training, we evaluate each policy in a zero-shot evaluation, as well as in combination with GC-TTT  
 188 [Bagatella et al., 2025] across a range of test-time compute budgets. We report success rates averaged  
 189 across the five standard evaluation tasks and three seeds. A detailed account of our checkpoint  
 190 selection criteria and hyperparameter configurations is provided in Section 5.2 and Appendix F.

191 Our main experimental results are organized around three questions, which define the following  
 192 subsections.

### 193 5.1 Does META-GC-TTT improve adaptation efficiency and performance?

194 **High-quality Data and Supervised Learning** Figure 3 reports success rates as a function of the  
 195 number of test-time training (TTT) steps  $K \in \{0 \text{ (zero-shot)}, 10, 20, 50\}$  when using GCBC as the  
 196 underlying algorithm, and training on navigate data. Across the environments, we observe that

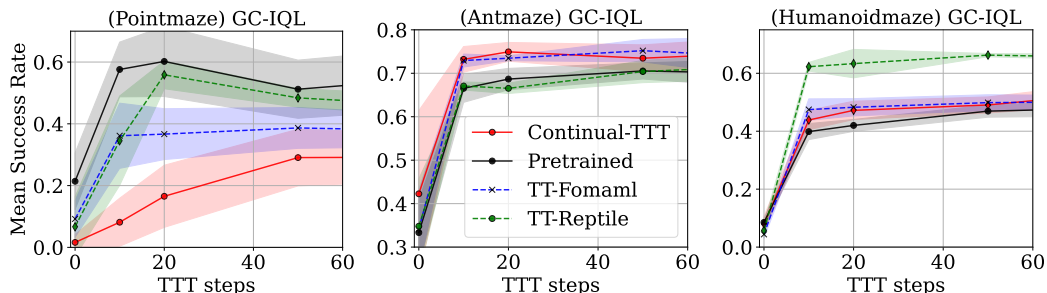


Figure 4: On mixed-quality datasets, meta-post-training positively impacts test-time training on two environments, while leading to instabilities on pointmaze.

197 meta-learned policies significantly improve both zero-shot performance and overall efficiency of  
 198 the adaptation process.

199 Regarding zero-shot performance, we observe substantial improvements, particularly with TT-  
 200 REPTILE and CONTINUAL-TTT. These meta-update rules implicitly optimize the base model as well  
 201 as the adapted model for performance on training tasks; assuming the task distribution is not disjoint,  
 202 this optimization generalizes effectively to similar test-time goals.

203 In contrast, TT-FOMAML specifically optimizes for adaptation potential rather than zero-shot  
 204 success, typically providing the steepest relative increase in performance within the first few adap-  
 205 tation steps. While this objective can occasionally degrade zero-shot performance—as seen in  
 206 humanoidmaze, where the initialization likely shifts toward a high-sensitivity region of the loss  
 207 landscape - it facilitates recovery. On antmaze, TT-FOMAML achieves the largest absolute gains  
 208 among all evaluated methods, reaching a peak success rate of 0.72 compared to the 0.55 baseline.

209 All three meta-update results translate into significant adaptation efficiency gains during deployment.  
 210 In pointmaze, meta-post-trained models achieve an 80% success rate in as few as 10 to 20 steps,  
 211 whereas the baseline requires 50 steps. This represents an approximate 3–5 $\times$  gain in computational  
 212 efficiency, demonstrating that meta-learning effectively aligns pre-trained policies for rapid online  
 213 specialization.

214 **Mixed Data and Off-policy Learning** While results on GCBC demonstrate the potential of META-  
 215 GC-TTT with expert data, standard offline RL assumes a more challenging regime: learning from  
 216 suboptimal, mixed-quality datasets (e.g. *stitch*), which in turn demand algorithms capable of  
 217 policy evaluation and improvement. To investigate this setting, we apply META-GC-TTT while  
 218 using GCIQL [Kostrikov et al., 2021] as the underlying algorithm. Unlike GCBC, GCIQL requires  
 219 goal-conditioned critics to model the optimality of actions; these estimators are trained with expectile  
 220 regression or TD learning, introducing significant optimization instabilities. In Figure 4 we observe  
 221 that meta-updates may fail to converge when applied to the highly filtered, off-policy batches  
 222 encountered during the meta-training phase, as *pointmaze* shows. We further study this phenomenon  
 223 in Appendix D, which shows that policy-only updates may circumvent these issues, and successfully  
 224 meta-learn.

225 Beyond *pointmaze*, we observe that META-GC-TTT performs well out-of-the-box. On *antmaze*,  
 226 CONTINUAL-TTT and TT-FOMAML consistently outperform the baseline across nearly all adap-  
 227 tation steps. In the high-dimensional *humanoidmaze*, TT-REPTILE emerges as the most effective  
 228 method, significantly exceeding the performance of its counterparts. These results suggest that no  
 229 single meta-learning algorithm is globally superior; rather, their efficacy depends on the specific  
 230 dynamics and state-space complexity of the environment.

231 Notably, when trained on *stitch* datasets, meta-learned policies do not exhibit the same significant  
 232 zero-shot improvements observed in the BC expert setting. We hypothesize that this may be related  
 233 to the additional complexity introduced by learning a goal-conditioned critic from fragmented data.

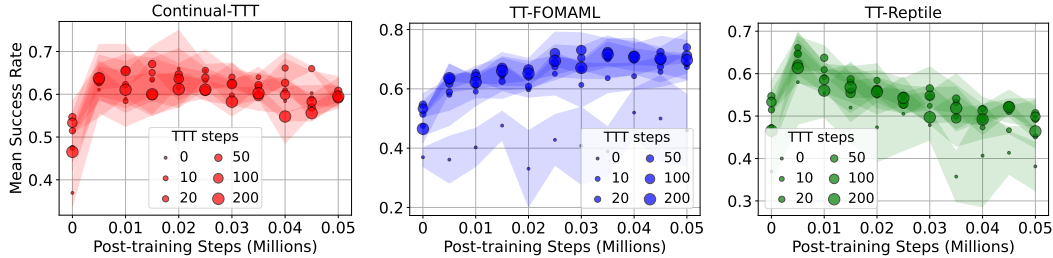


Figure 5: The test-time success rate evolution over post-training steps for different meta-update rules on `antmaze`. Performance curves are heavily dependent on the meta-algorithm used.

## 234 5.2 How much compute should be allocated for post-training?

235 While we have shown that meta-learning is beneficial for test-time performance and practical because  
 236 it can be performed entirely offline, it remains computationally expensive. Therefore, this section  
 237 studies trends in downstream performance as more compute is allocated to post-training. Figure 5 plots  
 238 success rates for an increasing number of post-training iterations with a representative environment-  
 239 algorithm pair (`antmaze`, GCBC) and each of the three meta-updates. In all cases, we observe that as  
 240 little as 5k outer loops are sufficient to significantly boost zero-shot and TTT performance. At the  
 241 same time, we also report three distinct behavioral patterns:

- 242 • CONTINUAL-TTT quickly drives zero-shot and TTT performance to a similar range, and  
 243 generally maintains performance through post-training;
- 244 • TT-FOMAML exhibits gradual but slower improvement, potentially due to an ill-  
 245 conditioned loss landscape;
- 246 • TT-REPTILE shows performance degradation over meta-learning steps, maintaining a gap  
 247 between zero-shot and TTT success rates.

248 These results indicate that the relationship between optimization objectives and actual task success  
 249 is non-trivial. The lack of a direct correlation suggests that minimizing surrogate loss does not  
 250 consistently yield improvements in loco-navigation performance.

## 251 5.3 Does early-stopping ease the burden of hyperparameter tuning?

252 While the previous sections demonstrate the effectiveness of META-GC-TTT, the framework intro-  
 253 duces several critical hyperparameters that govern the stability of meta-post-training and the efficiency  
 254 of test-time adaptation. Appendix F provides a comprehensive report of these parameters, including  
 255 inner/outer batch sizes and learning rates.

256 The number of inner-loop steps ( $K$ ) is particularly crucial from a practical standpoint. A significant  
 257 limitation in the applicability of standard optimization-based meta-learning is that  $K$  is typically fixed  
 258 a priori rather than tuned per task. Our findings indicate that the choice of  $K$  during meta-post-training  
 259 significantly impacts the Area Under the Success Curve (AUC, Appendix E).

260 Table 1 illustrates this trade-off: simply increasing  $K$  does not guarantee superior performance. No-  
 261 tably, TT-FOMAML exhibits a "collapse" as  $K$  increases, with success rates dropping to zero. This  
 262 is likely because the first-order approximation becomes increasingly inaccurate as error accumulates  
 263 over more inner steps, making the objective function overly complex on long horizons.

264 To lift the burden of choosing  $K$ , we test our proposed Early-Stopping strategy (ref. Section 4),  
 265 which selects the optimal checkpoint within a 200-step window based on test loss. As reported in  
 266 Table 1, this strategy consistently recovers performance lost to overfitting for both CONTINUAL-TTT  
 267 and TT-REPTILE. Interestingly, TT-FOMAML benefits less from early stopping. In this case, the  
 268 zero-shot model begins with a higher initial loss; while inner-loop optimization successfully reduces  
 269 this loss, the resulting meta-update remains suboptimal. This suggests that the performance bottleneck  
 270 for TT-FOMAML resides within the inherent limitations of the MAML meta-update approximation  
 271 rule. Ultimately, these results indicate that the ideal TTT horizon is highly task-dependent: while

Table 1: Inner-loop steps comparison. All numbers reported are Area Under Success Curve for TTT scaling curves. For TT-REPTILE and CONTINUAL-TTT, early-stopping retains performance by avoiding overfitting and keeping expert models that minimize the test loss.

| Method        | Inner-loop Steps |       |       |       |       |       | 200 Early-stopping |
|---------------|------------------|-------|-------|-------|-------|-------|--------------------|
|               | 1                | 10    | 20    | 50    | 100   | 200   |                    |
| antmaze       |                  |       |       |       |       |       |                    |
| TT-REPTILE    | 0.707            | 0.690 | 0.683 | 0.651 | 0.659 | 0.652 | 0.685              |
| TT-FOMAML     | 0.602            | 0.560 | 0.449 | 0.417 | 0.358 | 0.000 | 0.404              |
| CONTINUAL-TTT | 0.732            | 0.688 | 0.658 | 0.636 | 0.529 | 0.423 | 0.607              |
| humanoidmaze  |                  |       |       |       |       |       |                    |
| TT-REPTILE    | 0.623            | 0.630 | 0.634 | 0.641 | 0.636 | 0.533 | 0.624              |
| TT-FOMAML     | 0.305            | 0.186 | 0.066 | 0.005 | 0.045 | 0.000 | 0.000              |
| CONTINUAL-TTT | 0.466            | 0.483 | 0.510 | 0.511 | 0.408 | 0.248 | 0.518              |

272 complex goals require sustained refinement, simpler tasks risk performance degradation through  
 273 behavioral drift if optimized for too many steps.

## 274 6 Conclusion

275 This work investigates the synergy between meta-learning and test-time training (TTT) in the  
 276 context of offline goal-conditioned reinforcement learning to bridge the gap between large-scale  
 277 pre-training and task-specific specialization. Our empirical evaluation provides strong evidence  
 278 that meta-learning-based post-training, formalized through the META-GC-TTT framework, serves  
 279 as a powerful alignment tool that prepares pre-trained models for rapid adaptation and enhances  
 280 overall adaptation efficiency. Additionally, we proposed an early-stopping mechanism that maintains  
 281 robustness by preventing the model from overfitting to suboptimal data batches, and eliminates the  
 282 need to tune the number of inner steps as a fixed hyperparameter.

283 **Limitations and Future Work** Several promising avenues remain for the development of adaptive  
 284 agents. A deeper investigation into the optimization dynamics of specific meta-algorithms could  
 285 yield more robust hyperparameter selection strategies, particularly regarding the trade-off between  
 286 inner-loop learning rates and outer-loop stability. Furthermore, combining META-GC-TTT with hier-  
 287 archical algorithms (e.g., SAW [Zhou and Kao, 2026]) could mitigate scaling bottlenecks associated  
 288 with long horizons [Park et al., 2025b] by adapting high-level planners alongside low-level controllers.

## 289 References

- 290 Siddhant Agarwal, Ishan Durugkar, Peter Stone, and Amy Zhang. f-policy gradients: A general  
 291 framework for goal-conditioned rl using f-divergences. In *NeurIPS*, 2023.
- 292 Ali Amin, Raichelle Aniceto, Ashwin Balakrishna, Kevin Black, Ken Conley, Grace Connors, James  
 293 Darpinian, Karan Dhabalia, Jared DiCarlo, Danny Driess, et al.  $\pi^*$  0.6: a vla that learns from  
 294 experience. *arXiv preprint arXiv:2511.14759*, 1, 2025.
- 295 Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob  
 296 McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In  
 297 *NeurIPS*, 2017.
- 298 Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *arXiv preprint*  
 299 *arXiv:1810.09502*, 2018.
- 300 Marco Bagatella, Mert Albaba, Jonas Hübötter, Georg Martius, and Andreas Krause. Test-time offline  
 301 reinforcement learning on goal-related experience, 2025. URL [https://arxiv.org/abs/2507.](https://arxiv.org/abs/2507.18809)  
 302 18809.
- 303 Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz,  
 304 Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, Thomas

- 305 Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey Gritsenko,  
306 Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer,  
307 Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan  
308 Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen, and  
309 Xiaohua Zhai. Paligemma: A versatile 3b vlm for transfer, 2024. URL [https://arxiv.org/  
310 abs/2407.07726](https://arxiv.org/abs/2407.07726).
- 311 Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai,  
312 Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Sergey  
313 Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi, James  
314 Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky.  $\pi_0$ : A vision-language-  
315 action flow model for general robot control, 2026. URL <https://arxiv.org/abs/2410.24164>.
- 316 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
317 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,  
318 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale,  
319 2021. URL <https://arxiv.org/abs/2010.11929>.
- 320 Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning.  
321 *Journal of Machine Learning Research*, 6, 2005.
- 322 Benjamin Eysenbach, Tianjun Zhang, Sergey Levine, and Russ R Salakhutdinov. Contrastive learning  
323 as goal-conditioned reinforcement learning. *Advances in Neural Information Processing Systems*,  
324 35:35603–35620, 2022a.
- 325 Benjamin Eysenbach, Tianjun Zhang, Ruslan Salakhutdinov, and Sergey Levine. Contrastive learning  
326 as goal-conditioned reinforcement learning. In *NeurIPS*, 2022b.
- 327 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of  
328 deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- 329 Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. *Advances  
330 in neural information processing systems*, 31, 2018.
- 331 Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.  
332 *Advances in neural information processing systems*, 34:20132–20145, 2021.
- 333 Moritz Hardt and Yu Sun. Test-time training on nearest neighbors for large language models.  
334 In *The Twelfth International Conference on Learning Representations*, 2024. URL [https://  
335 openreview.net/forum?id=CNL2bku4ra](https://openreview.net/forum?id=CNL2bku4ra).
- 336 Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess,  
337 Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya Ghosh,  
338 Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin  
339 LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z.  
340 Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, James Tanner,  
341 Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury Zhilinsky.  $\pi_{0.5}$ : a  
342 vision-language-action model with open-world generalization, 2025. URL [https://arxiv.org/  
343 abs/2504.16054](https://arxiv.org/abs/2504.16054).
- 344 Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael  
345 Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin  
346 Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla:  
347 An open-source vision-language-action model, 2024. URL [https://arxiv.org/abs/2406.  
348 09246](https://arxiv.org/abs/2406.09246).
- 349 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL  
350 <https://arxiv.org/abs/1412.6980>.
- 351 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit  
352 q-learning, 2021. URL <https://arxiv.org/abs/2110.06169>.
- 353 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline  
354 reinforcement learning, 2020. URL <https://arxiv.org/abs/2006.04779>.

- 355 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,  
356 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 357 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,  
358 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning, 2019.  
359 URL <https://arxiv.org/abs/1509.02971>.
- 360 Hao Liu, Richard Socher, and Caiming Xiong. Taming maml: Efficient unbiased meta-reinforcement  
361 learning. In *International conference on machine learning*, pages 4061–4071. PMLR, 2019.
- 362 Minghuan Liu, Menghui Zhu, and Weinan Zhang. Goal-conditioned reinforcement learning: Problems  
363 and solutions, 2022. URL <https://arxiv.org/abs/2201.08299>.
- 364 Yecheng Jason Ma, Jason Yan, Dinesh Jayaraman, and Osbert Bastani. Offline goal-conditioned  
365 reinforcement learning via  $\beta$ -advantage regression. In *NeurIPS*, 2022.
- 366 Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv  
367 preprint arXiv:1803.02999*, 2018.
- 368 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
369 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow  
370 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–  
371 27744, 2022.
- 372 Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. Higl: Offline goal-conditioned  
373 rl with latent states as actions. *Advances in Neural Information Processing Systems*, 36:34866–  
374 34891, 2023.
- 375 Seohong Park, Kevin Frans, Sergey Levine, and Aviral Kumar. Is value learning really the main  
376 bottleneck in offline rl? *Advances in Neural Information Processing Systems*, 37:79029–79056,  
377 2024.
- 378 Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking  
379 offline goal-conditioned rl, 2025a. URL <https://arxiv.org/abs/2410.20092>.
- 380 Seohong Park, Kevin Frans, Deepinder Mann, Benjamin Eysenbach, Aviral Kumar, and Sergey  
381 Levine. Horizon reduction makes rl scalable. *arXiv preprint arXiv:2506.04168*, 2025b.
- 382 Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression:  
383 Simple and scalable off-policy reinforcement learning, 2019. URL <https://arxiv.org/abs/1910>.
- 384 Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language  
385 understanding by generative pre-training. 2018.
- 386 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
387 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
388 models from natural language supervision. In *International conference on machine learning*, pages  
389 8748–8763. PmLR, 2021.
- 390 Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit  
391 gradients. *Advances in neural information processing systems*, 32, 2019.
- 392 Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training  
393 with self-supervision for generalization under distribution shifts. In *International conference on  
394 machine learning*, pages 9229–9248. PMLR, 2020.
- 395 Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej,  
396 Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas  
397 Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon,  
398 Etienne Pot, Ivo Penchev, Gaël Liu, Francesco Visin, Kathleen Kenealy, Lucas Beyer, Xiaohai  
399 Zhai, Anton Tsitsulin, Robert Busa-Fekete, Alex Feng, Noveen Sachdeva, Benjamin Coleman,  
400 Yi Gao, Basil Mustafa, Iain Barr, Emilio Parisotto, David Tian, Matan Eyal, Colin Cherry, Jan-  
401 Thorsten Peter, Danila Sinopalnikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran Kazemi,  
402 Dan Malkin, Ravin Kumar, David Vilar, Idan Brusilovsky, Jiaming Luo, Andreas Steiner, Abe

403 Friesen, Abhanshu Sharma, Abheesht Sharma, Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa  
404 Saade, Alex Feng, Alexander Kolesnikov, Alexei Bendebury, Alvin Abdagic, Amit Vadi, András  
405 György, André Susano Pinto, Anil Das, Ankur Bapna, Antoine Miech, Antoine Yang, Antonia  
406 Paterson, Ashish Shenoy, Ayan Chakrabarti, Bilal Piot, Bo Wu, Bobak Shahriari, Bryce Petri  
407 Charlie Chen, Charline Le Lan, Christopher A. Choquette-Choo, CJ Carey, Cormac Brick, Daniel  
408 Deutsch, Danielle Eisenbud, Dee Cattle, Derek Cheng, Dimitris Paparas, Divyashree Shivakumar  
409 Sreepathihalli, Doug Reid, Dustin Tran, Dustin Zelle, Eric Noland, Erwin Huizenga, Eugene  
410 Kharitonov, Frederick Liu, Gagik Amirkhanyan, Glenn Cameron, Hadi Hashemi, Hanna Klimczak-  
411 Plucińska, Harman Singh, Harsh Mehta, Harshal Tushar Lehri, Hussein Hazimeh, Ian Ballantyne,  
412 Idan Szpektor, Ivan Nardini, Jean Pouget-Abadie, Jetha Chan, Joe Stanton, John Wieting, Jonathan  
413 Lai, Jordi Orbay, Joseph Fernandez, Josh Newlan, Ju yeong Ji, Jyotinder Singh, Kat Black, Kathy  
414 Yu, Kevin Hui, Kiran Vodrahalli, Klaus Greff, Linhai Qiu, Marcella Valentine, Marina Coelho,  
415 Marvin Ritter, Matt Hoffman, Matthew Watson, Mayank Chaturvedi, Michael Moynihan, Min Ma,  
416 Nabila Babar, Natasha Noy, Nathan Byrd, Nick Roy, Nikola Momchev, Nilay Chauhan, Noveen  
417 Sachdeva, Oskar Bunyan, Pankil Botarda, Paul Caron, Paul Kishan Rubenstein, Phil Culliton,  
418 Philipp Schmid, Pier Giuseppe Sessa, Pingmei Xu, Piotr Stanczyk, Pouya Tafti, Rakesh Shivanna,  
419 Renjie Wu, Renke Pan, Reza Rokni, Rob Willoughby, Rohith Vallu, Ryan Mullins, Sammy Jerome,  
420 Sara Smoot, Sertan Girgin, Shariq Iqbal, Shashir Reddy, Shruti Sheth, Siim Pöder, Sijal Bhatnagar,  
421 Sindhu Raghuram Panyam, Sivan Eiger, Susan Zhang, Tianqi Liu, Trevor Yacovone, Tyler Liechty,  
422 Uday Kalra, Utku Evcı, Vedant Misra, Vincent Roseberry, Vlad Feinberg, Vlad Kolesnikov,  
423 Woohyun Han, Woosuk Kwon, Xi Chen, Yinlam Chow, Yuvein Zhu, Zichuan Wei, Zoltan Egyed,  
424 Victor Cotruta, Minh Giang, Phoebe Kirk, Anand Rao, Kat Black, Nabila Babar, Jessica Lo,  
425 Erica Moreira, Luiz Gustavo Martins, Omar Sanseviero, Lucas Gonzalez, Zach Gleicher, Tris  
426 Warkentin, Vahab Mirrokni, Evan Senter, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia  
427 Hadsell, Yossi Matias, D. Sculley, Slav Petrov, Noah Fiedel, Noam Shazeer, Oriol Vinyals, Jeff  
428 Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Jean-Baptiste  
429 Alayrac, Rohan Anil, Dmitry, Lepikhin, Sebastian Borgeaud, Olivier Bachem, Armand Joulin,  
430 Alek Andreev, Cassidy Hardin, Robert Dadashi, and Léonard Hussenot. Gemma 3 technical report,  
431 2025. URL <https://arxiv.org/abs/2503.19786>.

432 Stephen Tian, Suraj Nair, Frederik Ebert, Sudeep Dasari, Benjamin Eysenbach, Chelsea Finn, and  
433 Sergey Levine. Model-based visual planning with self-supervised functional distances. In *ICLR*,  
434 2021.

435 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
436 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
437 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

438 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang  
439 Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,  
440 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin  
441 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,  
442 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui  
443 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang  
444 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger  
445 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan  
446 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.

447 Rui Yang, Yiming Lu, Wenzhe Li, Hao Sun, Meng Fang, Yali Du, Xiu Li, Lei Han, and Chongjie  
448 Zhang. Rethinking goal-conditioned supervised learning and its connection to offline rl. In *10th  
449 International Conference on Learning Representations, ICLR 2022*, 2022.

450 Zilai Zeng, Ce Zhang, Shijie Wang, and Chen Sun. Goal-conditioned predictive coding for offline  
451 reinforcement learning. *Advances in Neural Information Processing Systems*, 36:25528–25548,  
452 2023.

453 Chongyi Zheng, Ruslan Salakhutdinov, and Benjamin Eysenbach. Contrastive difference predictive  
454 coding. *arXiv preprint arXiv:2310.20141*, 2023.

455 John L. Zhou and Jonathan C. Kao. Flattening hierarchies with policy bootstrapping, 2026. URL  
456 <https://arxiv.org/abs/2505.14975>.

457 Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Ste-  
458 fan Welker, Ayzaan Wahid, Quan Vuong, Vincent Vanhoucke, Huong Tran, Radu Soricut, Anikait  
459 Singh, Jaspiar Singh, Pierre Sermanet, Pannag R Sanketi, Grecia Salazar, Michael S Ryoo, Krista  
460 Reymann, Kanishka Rao, Karl Pertsch, Igor Mordatch, Henryk Michalewski, Yao Lu, Sergey  
461 Levine, Lisa Lee, Tsang-Wei Edward Lee, Isabel Leal, Yuheng Kuang, Dmitry Kalashnikov, Ryan  
462 Julian, Nikhil J Joshi, Alex Irpan, brian ichter, Jasmine Hsu, Alexander Herzog, Karol Hausman,  
463 Keerthana Gopalakrishnan, Chuyuan Fu, Pete Florence, Chelsea Finn, Kumar Avinava Dubey,  
464 Danny Driess, Tianli Ding, Krzysztof Marcin Choromanski, Xi Chen, Yevgen Chebotar, Justice  
465 Carbajal, Noah Brown, Anthony Brohan, Montserrat Gonzalez Arenas, and Kehang Han. RT-2:  
466 Vision-language-action models transfer web knowledge to robotic control. In *7th Annual Confer-  
467 ence on Robot Learning*, 2023. URL <https://openreview.net/forum?id=XMQgwiJ7KSX>.

## 468 **A Goal-Conditioned Test-Time Training (GC-TTT)**

469 This section provides a brief introduction to Goal-Conditioned Test-time Training (GC-TTT)  
470 [Bagatella et al., 2025]. GC-TTT involves two distinct operations: dataset filtering and policy  
471 updates.

472 The filtering procedure actively condition the data distribution of an offline dataset  $\mathcal{D}$  on both the  
473 target goal  $g$  and the current state of the agent  $s$ , constructing  $\mathcal{D}_{(s,g)}$ . Trajectories  $\tau$  in  $\mathcal{D}$  are first  
474 filtered based on their proximity to the current agent state  $s$ , retaining only trajectories containing at  
475 least one state that is  $\epsilon$ -close to  $s$ :

$$\text{is\_relevant}(\tau, s, \epsilon) \iff \min_{s' \in \tau} d(s, s') \leq \epsilon \quad (10)$$

476 Trajectories that pass the relevance filter are subsequently ranked according to an optimality score,  
477 selecting the *top-k* (or *top-p%*) trajectories. In the absence of a learned value function, trajectories  
478 that do not reach the goal state (i.e. they have no states  $\epsilon$ -close to  $g$ ) are simply filtered out, and  
479 remaining trajectories are ranked based on a proxy to the goal distance (e.g. L2):

$$\text{critic\_free\_score}(\tau, g) = \max_{s' \in \tau} -d(s', g) \quad (11)$$

480 While computationally simple, this method requires the dataset to contain trajectories that successfully  
481 reach the goal. It is not applicable when optimal paths must be composed from many sub-optimal,  
482 fragmented segments. When a learned goal-conditioned value function  $V(\cdot|g): \mathcal{S} \rightarrow \mathbb{R}$  is available,  
483 every state in the relevant trajectories can be evaluated, and the trajectory score can be defined instead  
484 by the maximum value attained:

$$\text{critic\_based\_score}(\tau, g) = \max_{s' \in \tau} V(s'|g) \quad (12)$$

485 Among all trajectories surviving the relevance filter, we select the top- $p\%$  based on their optimality  
486 ranking. Crucially, we extract *sub-trajectories* from these samples, beginning at the state closest to  $s$   
487 and ending at the state closest to  $g$  (or the state with the highest  $V$ -value).

488 Once  $\mathcal{D}_{s,g}$  is constructed, it can be used to optimize near-arbitrary offline RL objectives [Yang et al.,  
489 2022, Kostrikov et al., 2021, Zhou and Kao, 2026].

## 490 **B Test-time Loss Visualization**

491 To understand the dynamics of TTT, we analyze how the loss varies w.r.t. test-time gradient steps  
492 of the actor and critic components of the Goal-Conditioned Implicit Q-Learning (GCIQL) training.  
493 Figure 6 illustrates these trends for a representative `antmaze stitch` task.

494 Across all evaluated learning rates, we observe a distinct discrepancy in optimization stability. The  
495 actor loss typically exhibits a smooth decrease, suggesting that the policy effectively refines its local  
496 behavior using the retrieved transitions. In contrast, the critic loss (and by extension, the total loss) is  
497 significantly more erratic. As the learning rate increases (e.g.,  $3 \times 10^{-4}$ ), the actor loss begins to  
498 lose its smoothness, yet it remains significantly more stable than the critic. This pattern suggests that  
499 while the agent can reliably improve its action selection at test time, the underlying value function is  
500 prone to noise that can potentially mislead the adaptation process.

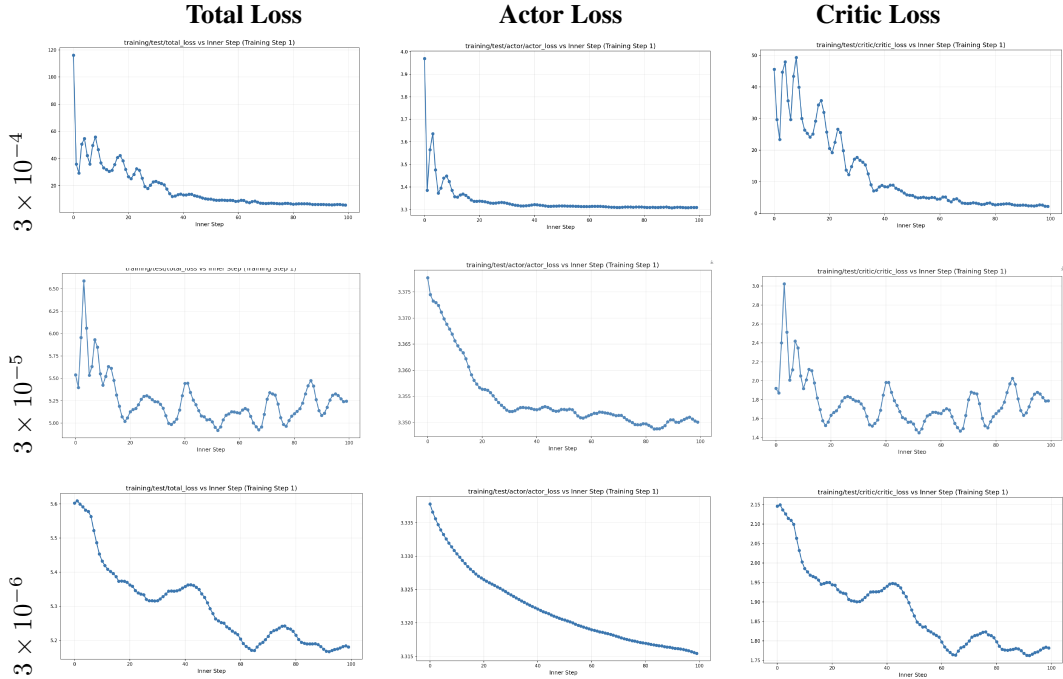


Figure 6: Comparison of total loss and its components, namely actor and critic losses, across multiple learning rates on a representative antmaze stitch task. The actor loss is typically smoother than the critic loss, which causes instabilities during META-GC-TTT.

## 501 C Actor-Only Test-time Training

502 The observed instability of the critic during TTT (Appendix B) suggests that updating the critic at test-  
 503 time may be unnecessary or even detrimental. To test this hypothesis, we perform an ablation study  
 504 comparing full adaptation (policy + critic) against policy-only updates. As shown in Table 2, updating  
 505 the actor only retains the majority of performance gains, particularly in complex environments like  
 506 antmaze and humanoidmaze.

Table 2: Area Under Success Curve (AUC) comparison: Full TTT (policy + critic) vs. actor-only TTT. Values represent the mean across 3 independent seeds.

| Setting / Learning Rate | $3 \times 10^{-4}$ | $5 \times 10^{-4}$ | $1 \times 10^{-3}$ | $3 \times 10^{-3}$ | $5 \times 10^{-3}$ |
|-------------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| <b>Pointmaze</b>        |                    |                    |                    |                    |                    |
| TTT (Full Update)       | 0.21               | 0.28               | 0.33               | 0.41               | 0.53               |
| TTT (Actor-Only)        | 0.22               | 0.24               | 0.27               | 0.32               | 0.35               |
| <b>Antmaze</b>          |                    |                    |                    |                    |                    |
| TTT (Full Update)       | 0.54               | 0.62               | 0.67               | 0.70               | 0.66               |
| TTT (Actor-Only)        | 0.55               | 0.61               | 0.65               | 0.68               | 0.65               |
| <b>Humanoidmaze</b>     |                    |                    |                    |                    |                    |
| TTT (Full Update)       | 0.42               | 0.39               | 0.35               | 0.03               | 0.00               |
| TTT (Actor-Only)        | 0.38               | 0.38               | 0.34               | 0.02               | 0.00               |

507 This indicates that the effectiveness of GC-TTT is largely driven by policy refinement and, as Park  
 508 et al. [2024] claim, the policy extraction is the bottleneck in offline RL. This finding is of significant  
 509 practical importance for the design of meta-learning pipelines: *by meta-post-training only the actor,*  
 510 *we may circumvent convergence issues.*

## 511 D Addressing Pathologies in Off-Policy META-GC-TTT

512 An unexpected result in our initial evaluation was the  
 513 failure of meta-learning algorithms to maintain baseline  
 514 performance in pointmaze stitch (see Figure 4).

515 In accordance to our findings in Sections B and C, we  
 516 investigate the significance of this degradation through  
 517 a slight change of the inner-loop objective, specifically  
 518 replacing the DDPG-BC loss with standard BC. As illus-  
 519 trated in Figure 7, META-GC-TTT shines in this setting,  
 520 with TT-REPTILE and TT-FOMAML achieving near-  
 521 ceiling success rates ( $> 95\%$ ) within a handful of TTT  
 522 steps, and maintaining this performance up to 200 steps.  
 523 This further supports our original conjecture: the initial  
 524 failures observed in Figure 4 likely stem from off-policy  
 525 learning in specific environment-hyperparameter config-  
 526 urations, and on-policy corrections are likely to mitigate  
 527 these issues.

528 More broadly, these substantial performance gains align with the assertions of Park et al. [2024],  
 529 suggesting that the primary bottleneck in offline RL is often policy extraction rather than value  
 530 learning. META-GC-TTT serves as a robust mechanism for this extraction, leveraging pre-trained  
 531 representations and goal-conditioned value functions to maximize task success.

## 532 E Success Metrics

533 To rigorously evaluate the adaptation performance of our models, we employ two primary statistical  
 534 tools: standard error estimation for visualizing uncertainty and the normalized area under the success  
 535 curve for summarizing adaptation efficiency.

536 **Uncertainty Estimation and Error Bars** To provide a robust estimate of performance across  
 537 different runs, we aggregate results over multiple random seeds at each test-time training (TTT)  
 538 budget. For a given configuration with  $n$  valid seeds, let  $s_i$  denote the success rate for seed  $i$ . The  
 539 mean success  $\hat{\mu}$  is calculated as:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n s_i.$$

540 The uncertainty displayed in our plots represents the seed-wise dispersion, estimated via the standard  
 541 error (SE):

$$\text{SE} = \frac{\hat{\sigma}}{\sqrt{n}}, \quad \text{where} \quad \hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (s_i - \hat{\mu})^2}.$$

542 To visualize the 95% confidence interval, we report the approximate confidence half-width as  
 543  $1.96 \times \text{SE}$ . This interval provides a measure of the reliability of the mean success rate, ensuring that  
 544 the observed improvements are statistically consistent across independent experimental runs.

545 **Area Under the Success Curve (AUC)** While success rates at specific TTT budgets are informative,  
 546 they do not fully capture the trade-off between computational cost and agent efficacy. We therefore  
 547 utilize the *normalized area under the success curve* (AUC) to provide a single-scalar summary  
 548 of adaptation performance. Given a set of points  $(x_k, y_k)$ , where  $x_k$  represents the number of  
 549 TTT\_steps and  $y_k$  represents the corresponding mean\_success, the AUC is computed using the  
 550 trapezoidal rule:

$$\text{AUC} = \frac{1}{x_{\max}} \sum_{k=1}^{K-1} \frac{y_k + y_{k+1}}{2} (x_{k+1} - x_k),$$

551 where  $x_{\max}$  is the maximum number of adaptation steps. This normalization yields a unitless score  
 552 bounded between 0 and 1. A higher AUC reflects a dual benefit: it indicates both how quickly the

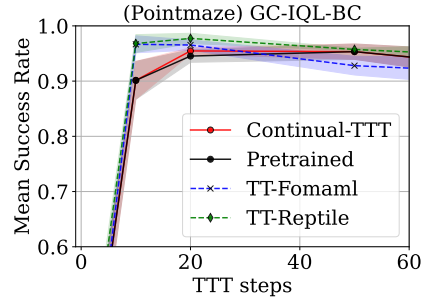


Figure 7: Switching the TTT objective to Behavioral Cloning (BC) on pointmaze stitch restores the effectiveness of meta-learning on top of a GCIQL backbone.

553 model specializes (adaptation efficiency) and the peak performance level achieved within the given  
 554 budget.

## 555 F Hyperparameter Choice and Ablations

556 **Hyperparameter Choice** For all hyperparameters not related to meta-learning, we use the default  
 557 settings as reported by [Park et al., 2025a] and [Bagatella et al., 2025]. META-GC-TTT introduces  
 558 several hyperparameters:

- 559 • **Learning Rates ( $\alpha, \beta$ ):** The inner learning rate ( $\alpha$ ) is the primary driver of test-time  
 560 improvement. We find that the optimal  $\alpha$  for meta-learning might differ from the one  
 561 used at test-time as it will affect post-training dynamics (see Figure 8). The meta-learning  
 562 rate ( $\beta$ ) is algorithm-specific: CONTINUAL-TTT requires none, TT-FOMAML uses an  
 563 meta-gradient learning rate (swept in  $\{3 \times 10^{-4}, 3 \times 10^{-5}\}$ ), and TT-REPTILE uses  $\beta \in$   
 564  $\{0.01, 0.1, 0.2, 0.5\}$  as a weight interpolation coefficient.
- 565 • **Batch Sizes ( $B_{inner}, B$ ):** To ensure a fair comparison and minimize the search space, we fix  
 566 the inner batch size to  $B_{inner} = 1024$ , following Park et al. [2025a]. To optimize memory  
 567 overhead we fix the meta-batch size  $B = 1$ . This is crucial for scaling as avoids storing  
 568 multiple copies of the model during the outer loop.
- 569 • **Number of Inner Steps ( $K$ ):** Determining the optimal  $K$  is challenging: a high  $K$  risks  
 570 overfitting on simple tasks, while a low  $K$  may be ineffective on complex navigation  
 571 goals. We conduct an extensive sweep  $K \in \{1, 10, 20, 50, 100, 200\}$  and evaluate our  
 572 **Early-Stopping** strategy fixing  $K = 100$  for expert datasets and  $K = 200$  for stitch  
 573 datasets.
- 574 • **Meta-Steps ( $N$ ):** We fix the post-training duration to  $N = 50k$  steps, which we empirically  
 575 found to be effective to improve performance with the tested datasets.

576 Table 3 reports all settings-specific hyper-parameters.

### 577 F.1 Inner Learning Rate Ablation

578 Our results indicate that the inner learning rate ( $\alpha$ ) used during both test-time training and meta-post-  
 579 training is the most critical factor for stability. While a high  $\alpha$  is beneficial to increase success rate at  
 580 test time (see Table 4), it frequently induces instability during the meta-learning process, ultimately  
 581 resulting in a degraded initialization.

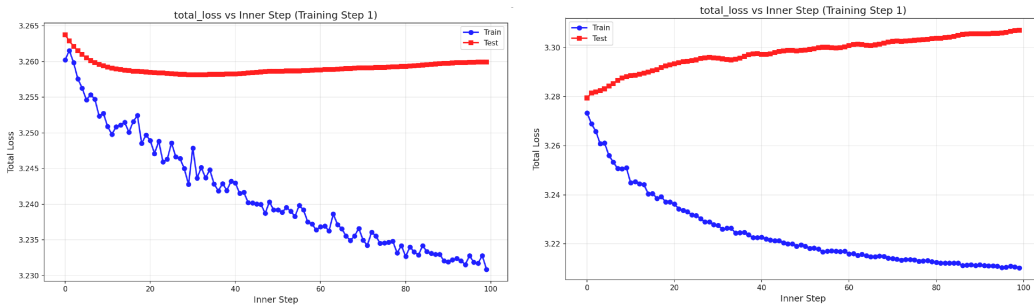


Figure 8: Visualization of the train and test loss of a task during post-training with various learning rates: bigger on the left ( $3 \times 10^{-4}$ ), smaller on the right ( $3 \times 10^{-5}$ )

582 As illustrated in Figure 8, high  $\alpha$  causes the model to overfit rapidly to the train batches encountered.  
 583 While this over-specialization is beneficial test-time performance, it is detrimental for meta-updating  
 584 the model. We empirically observe that models trained with higher learning rates exhibit lower  
 585 training losses but higher test losses throughout the post-training phase. Even when utilizing an  
 586 early-stopping strategy to mitigate this overfitting, the resulting updates are often rejected, thereby  
 587 stalling meta-learning progress.

| Algorithm                           | Actor loss | Inner LR           | Meta LR            | Inner steps ( $K$ ) | TTT LR             | Checkpoint | Early stop |
|-------------------------------------|------------|--------------------|--------------------|---------------------|--------------------|------------|------------|
| <b>Dataset: pointmaze_expert</b>    |            |                    |                    |                     |                    |            |            |
| CONTINUAL-TTT                       | BC         | $3 \times 10^{-5}$ | $3 \times 10^{-4}$ | 100                 | $3 \times 10^{-4}$ | 30k        | Yes        |
| TT-FOMAML                           | BC         | $3 \times 10^{-5}$ | $3 \times 10^{-4}$ | 100                 | $3 \times 10^{-4}$ | 35k        | Yes        |
| TT-REPTILE                          | BC         | $3 \times 10^{-5}$ | 0.01               | 100                 | $3 \times 10^{-4}$ | 45k        | Yes        |
| <b>Dataset: antmaze_expert</b>      |            |                    |                    |                     |                    |            |            |
| CONTINUAL-TTT                       | BC         | $3 \times 10^{-5}$ | $3 \times 10^{-4}$ | 100                 | $3 \times 10^{-4}$ | 10k        | Yes        |
| TT-FOMAML                           | BC         | $3 \times 10^{-5}$ | $3 \times 10^{-4}$ | 100                 | $3 \times 10^{-4}$ | 35k        | Yes        |
| TT-REPTILE                          | BC         | $3 \times 10^{-4}$ | 0.5                | 100                 | $3 \times 10^{-4}$ | 5k         | Yes        |
| <b>Dataset: humanoidmaze_expert</b> |            |                    |                    |                     |                    |            |            |
| CONTINUAL-TTT                       | BC         | $3 \times 10^{-5}$ | $3 \times 10^{-4}$ | 100                 | $3 \times 10^{-5}$ | 25k        | Yes        |
| TT-FOMAML                           | BC         | $3 \times 10^{-5}$ | $3 \times 10^{-4}$ | 100                 | $3 \times 10^{-4}$ | 40k        | Yes        |
| TT-REPTILE                          | BC         | $3 \times 10^{-5}$ | 0.5                | 100                 | $3 \times 10^{-5}$ | 50k        | Yes        |
| <b>Dataset: pointmaze_stitch</b>    |            |                    |                    |                     |                    |            |            |
| CONTINUAL-TTT                       | BC         | $3 \times 10^{-3}$ | $3 \times 10^{-4}$ | 200                 | $3 \times 10^{-3}$ | 30k        | Yes        |
| TT-FOMAML                           | BC         | $3 \times 10^{-3}$ | $3 \times 10^{-4}$ | 1                   | $3 \times 10^{-3}$ | 5k         | No         |
| TT-REPTILE                          | BC         | $3 \times 10^{-3}$ | $1 \times 10^{-4}$ | 10                  | $3 \times 10^{-3}$ | 40k        | No         |
| CONTINUAL-TTT                       | DDPG-BC    | $5 \times 10^{-3}$ | $3 \times 10^{-4}$ | 50                  | $5 \times 10^{-3}$ | 35k        | No         |
| TT-FOMAML                           | DDPG-BC    | $5 \times 10^{-3}$ | $3 \times 10^{-4}$ | 10                  | $5 \times 10^{-3}$ | 5k         | No         |
| TT-REPTILE                          | DDPG-BC    | $5 \times 10^{-3}$ | 0.1                | 1                   | $5 \times 10^{-3}$ | 20k        | No         |
| <b>Dataset: antmaze_stitch</b>      |            |                    |                    |                     |                    |            |            |
| CONTINUAL-TTT                       | DDPG-BC    | $3 \times 10^{-5}$ | $3 \times 10^{-4}$ | 100                 | $1 \times 10^{-3}$ | 25k        | Yes        |
| TT-FOMAML                           | DDPG-BC    | $3 \times 10^{-5}$ | $3 \times 10^{-4}$ | 100                 | $1 \times 10^{-3}$ | 15k        | Yes        |
| TT-REPTILE                          | DDPG-BC    | $3 \times 10^{-5}$ | 0.1                | 100                 | $1 \times 10^{-3}$ | 5k         | Yes        |
| <b>Dataset: humanoidmaze_stitch</b> |            |                    |                    |                     |                    |            |            |
| CONTINUAL-TTT                       | DDPG-BC    | $3 \times 10^{-4}$ | $3 \times 10^{-4}$ | 20                  | $3 \times 10^{-4}$ | 10k        | No         |
| TT-FOMAML                           | DDPG-BC    | $3 \times 10^{-4}$ | $3 \times 10^{-4}$ | 1                   | $3 \times 10^{-4}$ | 50k        | No         |
| TT-REPTILE                          | DDPG-BC    | $3 \times 10^{-4}$ | 0.1                | 200                 | $3 \times 10^{-4}$ | 10k        | Yes        |

Table 3: Best Hyper-parameters for each dataset.

Table 4: A higher learning rate  $\alpha$  is generally beneficial for GC-TTT. We report Area Under the Success Curve (AUC) over 3 seeds.

| Environment         | Inner Learning Rate ( $\alpha$ ) |                    |                    |                    |                    |
|---------------------|----------------------------------|--------------------|--------------------|--------------------|--------------------|
|                     | $3 \times 10^{-5}$               | $1 \times 10^{-4}$ | $3 \times 10^{-4}$ | $1 \times 10^{-3}$ | $3 \times 10^{-3}$ |
| pointmaze/stitch    | 0.21                             | 0.28               | 0.33               | 0.41               | <b>0.53</b>        |
| antmaze/stitch      | 0.54                             | 0.62               | 0.67               | <b>0.70</b>        | 0.66               |
| humanoidmaze/stitch | 0.44                             | 0.48               | <b>0.48</b>        | 0.40               | 0.03               |

## 588 G Pre-training Success Curves

589 Departing from the meta-learning analysis, Figure 9 illustrates test-time success rates across varying  
590 pre-training horizons. We observe diverse, environment-specific scaling behaviors that provide  
591 broader insights into offline GCRL performance. Notably, extending the pre-training budget beyond  
592 the 400k-step baseline used for our meta-learning initializations does not yield further gains.

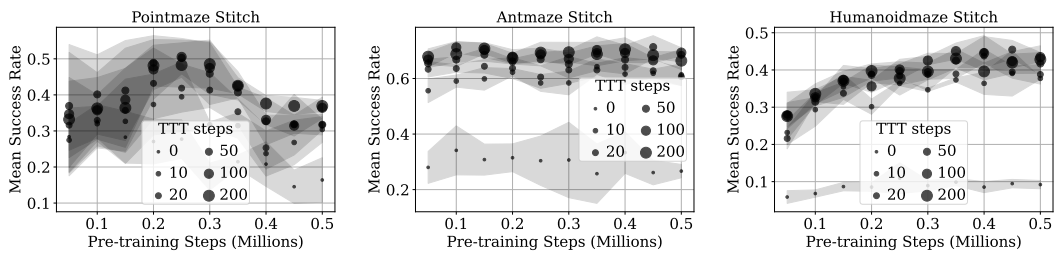


Figure 9: The test-time success rate evolution over pre-training steps is environment-dependent. On `pointmaze` stitch it starts degrading after 250k steps, on `antmaze` stitch it plateaus after 100k steps, while on `humanoidmaze` stitch converges at 350k steps. In all cases, 400k steps are sufficient for convergence.