

---

# Generative Actor-Critic

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Conventional Reinforcement Learning (RL) algorithms, typically focused on esti-  
2 mating or maximizing expected returns, face challenges when online finetuning  
3 offline pretrained models. This paper introduces Generative Actor Critic (GAC),  
4 a novel framework that decouples sequential decision-making into two stages:  
5 learning a generative model of the joint distribution of trajectories and their returns,  
6  $p(\tau, y)$ , and then performing decision-making via inference on this learned model.  
7 GAC offers a new perspective, framing *policy evaluation* as learning this compre-  
8 hensive distribution, and enabling versatile *policy improvement* strategies through  
9 inference. To operationalize GAC, we introduce a specific instantiation based on a  
10 latent variable model that features a continuous latent plan vector. We develop novel  
11 inference strategies for both exploitation (optimizing latent plans for expected re-  
12 turns) and exploration (sampling latent plans conditioned on dynamically adjusted  
13 target returns). Experiments on Gym-MuJoCo benchmarks demonstrate GAC’s  
14 strong offline performance and significantly enhanced offline-to-online adaptation  
15 compared to state-of-the-art methods, even in absence of stepwise rewards.

## 16 1 Introduction

17 A central objective in sequential decision-making is to *maximize the expected returns on trajec-*  
18 *tries* [1], denoted as  $\mathbb{E}_{p(\tau)}[Y(\tau)]$ , where each trajectory  $\tau$  consists of a sequence of states and actions,  
19 and  $Y$  is a utility function assigning return  $y$  to each trajectory. Conventional Reinforcement Learning  
20 (RL) algorithms are designed to estimate and optimize this expectation during their training phase.  
21 The widely-adopted Actor-Critic framework [2, 3], for instance, exemplifies this by learning a *critic*  
22 to evaluate expected returns and an *actor* to refine the policy towards maximizing these returns.  
23 This expectation-centric paradigm is particularly well-suited for online learning settings common  
24 in traditional RL [1], due to its amenability to efficient, iterative updates from agent-environment  
25 interactions. However, in the context of modern Generative AI (GenAI) pipelines—often charac-  
26 terized by an extensive offline pre-training stage on vast datasets preceding online fine-tuning—we  
27 propose to move beyond this expectation-centric approach. We advocate for decoupling the decision  
28 modeling process into two distinct stages: (1) train-time generative modeling of the joint distribution  
29 over trajectories  $\tau$  and their corresponding returns  $y$ :

$$p(\tau, y)$$

30 and (2) test-time decision-making, framed as an inference query based on this learned joint distribution.  
31 In this paper, we introduce this approach as the Generative Actor-Critic (GAC) framework.

32 The Generative Actor-Critic framework fundamentally shifts the paradigm from estimating trajec-  
33 tories’ expected returns to learning a comprehensive distribution of trajectories and their returns.  
34 This holistic approach offers several key advantages. In particular, modeling the distribution  $p(\tau, y)$   
35 naturally entails a generative critic  $p(y|\tau)$  and thus achieves *generalized policy evaluation*; this  
36 perspective underscores the ability of GAC to utilize data from various sources when modeling the

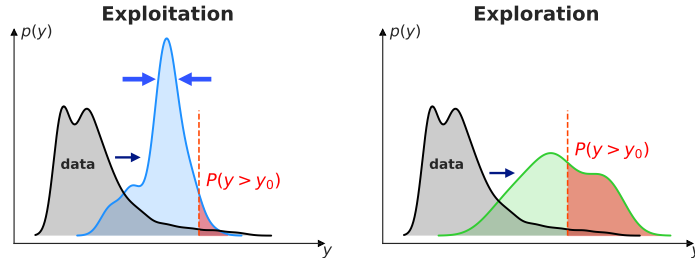


Figure 1: **Conceptual illustration of exploitation and exploration.** After modeling the data distribution (gray), GAC supports distinct test-time inference queries for various decision-making objectives. **(Left)** For *exploitation*, the objective is to maximize the expected return, leading to a focused, low-variance policy that targets high-certainty outcomes (blue). **(Right)** For *exploration*, the generative model is conditioned on a target distribution shifted towards higher returns, guiding the search for novel and potentially superior trajectories (green).

37 correlations between behaviors and outcomes, making GAC particularly well-suited for scenarios  
 38 involving offline pre-training followed by online fine-tuning. Furthermore, by explicitly modeling  
 39 the entire distribution, GAC can capture complex, multi-modal relationships between behaviors and  
 40 outcomes—a capacity often limited in methods focusing solely on expectation. While this advantage  
 41 is shared by distributional RL [4], the latter typically learns a return distribution conditioned on  
 42 state-action pairs (often via a distributional Bellman equation) and subsequently falls back to the  
 43 expectation of this learned return distribution for policy extraction [4–6]. In contrast, GAC directly  
 44 models the joint distribution of trajectories and returns, enabling a fundamental shift that replaces  
 45 expectation-based actor with a versatile inference process on  $p(\tau|y)p(y)$ . For instance, one could  
 46 query the learned model to identify trajectories that maximize expected future returns, similar to  
 47 policy optimization in distributional RL [5]. Alternatively, one could steer the generative process  
 48 by conditioning on desired high returns to sample corresponding trajectories, akin to mechanisms  
 49 in Decision Transformer (DT) [7, 8] and Diffuser [9, 10]. One could also sample diverse yet high-  
 50 performing trajectories [11] to achieve *generalized policy improvement*. As we will demonstrate, such  
 51 test-time flexibility can be leveraged, particularly in the offline-to-online setting, to address aspects of  
 52 the classical exploration-exploitation dilemma. Fig. 1 illustrates the intuition. Starting from the data  
 53 distribution, *exploitation* is to refine the policy to focus on a narrow distribution of high-certainty  
 54 returns, while *exploration* is to actively seek novel, potentially superior outcomes by shifting the  
 55 target return distribution towards uncharted, higher-value regions.

56 To operationalize the Generative Actor Critic framework and realize its potential, particularly for  
 57 complex decision-making tasks and effective offline-to-online adaptation, this paper introduces  
 58 several core methodological contributions: First, we instantiate GAC using a latent-variable model  
 59 structured as  $p(\tau, y, z) = p(\tau|z)p(y|z)p(z)$ . Inspired by recent advances like Latent Plan Transformer  
 60 (LPT) [12], a continuous latent plan  $z$  is introduced to effectively capture the correlation between  
 61 high-dimensional trajectories  $\tau$  and their corresponding low-dimensional returns  $y$ . Building on LPT’s  
 62 demonstration of efficient inference via such a latent space, we introduce non-trivial architectural  
 63 and algorithmic modifications to significantly boost offline performance. Second, leveraging this  
 64 latent structure, we design designated inference queries for exploitation and exploration. Exploitation  
 65 is formulated as maximizing the expected return  $\mathbb{E}[y|z]$  by performing gradient ascent directly in  
 66 the latent plan space. For exploration, we propose to sample  $z$  from the posterior  $p(z|y^+)$  of target  
 67 returns from  $p(y^+)$ . These targets are dynamically set by sampling high-performing base returns  
 68  $y$  from the empirical distribution  $p_{\text{data}}(y)$  (i.e., a replay buffer’s top-k) and adding a small target  
 69 improvement  $\Delta y$ , thereby guiding exploration towards potentially superior outcomes. The small  
 70 target improvement and the prioritized replay buffer [13, 14] offer a “trust region” [15] for online  
 71 fine-tuning. Empirically, GAC achieves competitive performance on Gym-MuJoCo benchmarks. Our  
 72 framework demonstrates strong offline learning capabilities and significantly enhanced adaptation in  
 73 offline-to-online scenarios compared to state-of-the-art baselines, even when step-wise rewards are  
 74 absent during parts of its operation.

## 75 2 Preliminaries

76 **Actor-Critic** The conventional decision-making process in RL is Markov Decision Process  
 77 (MDP)[1], represented by a tuple  $\mathcal{M} = \langle S, A, P, \pi, R, \rho, T \rangle$ , comprising a state space  $S$ , ac-  
 78 tion space  $A$ , transition function  $P : S \times A \mapsto \Pi(S)$ , policy  $\pi : S \mapsto \Pi(A)$ , reward function

79  $R : S \times A \mapsto \mathbb{R}$ , initial state distribution  $\rho : \Pi(S)$ , and horizon  $T$ . The decision-making objective  
80 is to maximize  $\mathbb{E}[y]$ , *i.e.*, the expectation of *return*  $y := \sum_{t=0}^T R(s_t, a_t)$  accumulated along the  
81 trajectories  $\tau = [s_0, a_0, s_1, a_1, \dots, a_{T-1}, s_T]$ . One of the most useful identities in this formulation  
82 is the Bellman equation: if we define  $Q(s_t, a_t) := R(s_t, a_t) + \mathbb{E}_{P, \pi}[\sum_{k=1}^{T-t} R(s_{t+k}, a_{t+k})]$ , we will  
83 have  $Q(s_t, a_t) = R(s_t, a_t) + \mathbb{E}_{P, \pi}[Q(s_{t+1}, a_{t+1})]$ . This Bellman equation, as well as its various  
84 variants [1, 16, 17, 3], provides a foundation for learning a value function  $Q$  for the expectation  
85 of *return-to-go* at each step  $RTG_t := \sum_{k=0}^{T-t} R(s_{t+k}, a_{t+k})$ . Subsequently, the policy  $\pi(a_t|s_t)$  can  
86 be updated along an approximated gradient direction  $\nabla_{a_t} Q(s_t, a_t)$ . That gives us the Actor-Critic  
87 framework [2, 18], where the *critic*  $Q$  estimates  $\mathbb{E}[RTG_t|s_t, a_t]$ , and the *actor*  $\pi$  optimizes the  
88 expectation from the critic.

89 **Distributional RL** Bellemare et al. [4] proposed to shift from the expectation-centric Bellman equa-  
90 tion to the distributional version,  $RTG_t|s_t, a_t = R(s_t, a_t) + RTG_{t+1}|s_{t+1}, a_{t+1}$ . The fundamental  
91 difference is that  $RTG_t$  is a random variable, while the previous object of interest,  $Q_t$ , is its mean.  
92 Formally speaking, in distributional RL we model  $p(RTG_t|s_t, a_t)$  when training the critic, and only  
93 calculate the expectation  $\mathbb{E}[RTG_t|s_t, a_t]$  when extracting policy for the actor [4–6]. Bellemare et al.  
94 [4] showed that such a generative critic could express the multi-modality of  $p(RTG_t|s_t, a_t)$  caused  
95 by *state aliasing* [19], *i.e.*, ambiguity of hidden states under insufficient memory context. However,  
96 Bellemare et al. [4] also formally derived that greedy policy selection (*i.e.* policy optimization toward  
97  $\nabla_{a_t} \mathbb{E}[RTG_t|s_t, a_t] = 0$ ) leads to unstable value iteration and does not guarantee contraction.

98 **RL via Sequence Modeling** Denoting the policy optimization result as  $p(a_t|s_t, Q(s_t, a_t) = Q^*)$ ,  
99 where  $Q^*$  is the optimal value, we obtain a generative perspective for the actor. This is what underlies  
100 the idea of return-conditioned behavior cloning [20, 7, 21]. Among them, Decision Transformer [7]  
101 proposed a generalized form of policy  $\pi(a_t|s_{\leq t}, a_{< t}, RTG_{\leq t})$  that is learned with offline data similar  
102 to the autoregressive training of language models. However, in test time, since the model does not  
103 have access to the  $RTG$ s anymore,  $RTG_0$  is set to be a particular high value (*e.g.*, the maximum  
104 or a quantile above the mean) in the offline data, and is updated as  $RTG_{t+1} = RTG_t - r_t$  after  
105 observing rewards. To set more plausible targets, Lee et al. [11] proposed to model not only the  
106 action, but also the rewards and  $RTG$ s. This connects the generative actor  $\pi(a_t|s_{\leq t}, a_{< t}, RTG_{\leq t})$   
107 with the generative critic  $p(RTG_t|s_{\leq t}, a_{\leq t}, RTG_{< t})$  in distributional RL and is closely related to  
108 our GAC framework. The difference is that GAC does not assume step-wise rewards, which is argued  
109 in Latent Plan Transformer [12] as a more naturalistic setup that results in decision-making models  
110 that can do *trajectory stitching* [22]. We will pinpoint a principled synergy between the generative  
111 critic and the generative actor for online policy improvement.

112 **Offline-to-Online RL** Offline-to-online RL [23–30] aims to enhance the sample efficiency of online  
113 fine-tuning by leveraging offline pre-training. However, conventional expectation-centric approaches  
114 often face challenges during this transition, as their critics can become unreliable due to the action  
115 distribution mismatch between the static offline dataset and the evolving online policy [31]. Efforts to  
116 mitigate this mismatch [32, 30] can, in turn, inadvertently curtail the exploratory capabilities of the  
117 learned actor. Generative actors, such as Decision Transformer, can be fine-tuned using a replay buffer  
118 that combines offline and online experiences [8]; yet, they often lack stable and scalable mechanisms  
119 for targeted exploration in the online phase. In contrast, GAC addresses this gap by providing a  
120 distributional perspective for targeted exploration:  $p(\tau|y^+)p(y^+)$ . This capability stems from GAC’s  
121 distinct approach of jointly modeling trajectories and their associated returns, which facilitates more  
122 principled and flexible exploration strategies during online interaction.

### 123 3 Generative Actor Critic (GAC)

124 GAC is a framework that separates generative decision-modeling and decision-making as inference.  
125 In Section 3.1, we introduce several probabilistic inference queries for different decision-making  
126 strategies, motivating the employment of latent-variable generative models for efficient inference. In  
127 Section 3.2, we introduce a learning algorithm for GAC based on a latent-variable model inspired by  
128 Latent Plan Transformer (LPT) [12]. In Section 3.3 we describe the interplay between the actor, the  
129 critic, and the replay buffer [13, 14] in realizing policy improvement in exploratory online finetuning.

130 **3.1 Decision-Making as Inference**

131 As GAC features a generative model of  $p(\tau, y)$ , decision-making can be framed as principled  
 132 inference over it. Previous works [16, 33–36] in decision-making as inference mainly take a stepwise  
 133 perspective, in which the policy can be viewed as amortized variational distribution  $q_{y^*}(a_t|s_t)$  for the  
 134 ground-truth posterior  $p(a_t|s_t, y = y^*)$ , where  $y^*$  denotes the optimal return. However, this posterior,  
 135 along with its generalized versions in DT [7, 8], Diffuser [9, 10], LPT [12], is believed to have an  
 136 *optimism bias* [16, 34]. Consider the decision of whether to buy a lottery, where  $y^* = \$10M$  but it is  
 137 almost impossible to get. The decision from  $p(a|s, y = y^*)$  would be “to buy”, while we all know  
 138 that the optimal decision is “not to buy”. Ziebart [16] attempted to bypass this issue by structuring  
 139 the variational distribution of  $p(\tau|y = y^*)$  as  $\prod_0^{T-1} q_{y^*}(a_t|s_t)P(s_{t+1}|s_t, a_t)$ , where  $P(s_{t+1}|s_t, a_t)$   
 140 is the ground-truth world dynamics.

141 We want to advocate for a more general solution to optimism bias. The issue is not in the inference  
 142 process, but in the inference objective. As we all know, the most intuitive decision-making objective  
 143 is to maximize the expected return  $\mathbb{E}[y|\tau]$ . And the actual optimal decision “not to buy” is optimal  
 144 under this objective. Obviously, to a generative model  $p(\tau, y)$ ,  $\max_{\tau} \mathbb{E}[y|\tau]$  is a inference query  
 145 distinct from  $p(\tau|y = y^*)$ . In fact, the logic underlying the decision policy  $p(\tau|y = y^*)$  in the  
 146 context of lottery is “If you want to win a lottery, you must first buy a lottery.”

147 Solving  $\max_{\tau} \mathbb{E}[y|\tau]$  directly would employ gradient ascent that requires costly backpropagation  
 148 through time. This is where latent modeling excels. Consider a generative model with continuous  
 149 latent variables  $z$

$$p_{\theta}(\tau, y, z) = p_{\alpha}(z)p_{\beta}(\tau|z)p_{\gamma}(y|z), \quad (1)$$

150 where conditional independence is assumed between the trajectory  $\tau$  and its return  $y$ , positioning  $z$   
 151 as an information bottleneck. Intuitively, the latent  $z$  are *plans* that abstract trajectories around their  
 152 returns. As long as  $p_{\gamma}(y|z)$  and  $\mathbb{E}[y|z]$  have analytical forms, gradient-based inference can be lifted  
 153 to the variational posterior  $q_{\phi}(z)$  in the continuous latent space:

$$\max_{\phi} \mathbb{E}_{q_{\phi}(z)} [\mathbb{E}[y|z]]. \quad (2)$$

154 Once  $z$  is inferred, the policy is entailed by  $p_{\beta}(\tau|z)$ . The efficiency gain over  $\max_{\tau} \mathbb{E}[y|\tau]$  is apparent.

155 Continuous latent modeling also enables other inference queries that would not be imaginable  
 156 otherwise. In fact, the optimization in Eq. (2) can be viewed as a special type of joint sampling over  
 157  $p(z, y)$ . Returning to the lottery example, a decision policy  $p(z|y)$  would not be unreasonable if  
 158 we had more nuanced control over optimism. In other words, if, instead of sampling from a fixed  
 159 target  $y^*$ , we sample from a marginal distribution  $p(y^+)$  that is slightly shifted right from  $p(y)$ ,  
 160 the *optimism* in  $p_{\theta}(z|y^+)$  is tamed with the awareness of  $p(y^+)$ . To sample from  $p_{\theta}(z|y^+)$ , we can  
 161 employ classical Variational Bayes, which maximizes the evidence lower bound (ELBO) [37, 38]:

$$\max_{\phi} \mathbb{E}_{q_{\phi}(z)} [\log p_{\gamma}(y^+|z)] - D_{\text{KL}}(q_{\phi}(z)||p_{\alpha}(z)). \quad (3)$$

162 This is equivalent to  $\min_{\phi} D_{\text{KL}}(q_{\phi}(z)||p_{\theta}(z|y^+))$ , where  $p_{\theta}(z|y) \propto p(z)p_{\gamma}(y|z)$ . We will introduce  
 163 how to sample from  $p(y^+)$  in Section 3.3.

164 **3.2 Generative Modeling**

165 Latent Plan Transformer [12] is an instantiation of the latent-variable generative model described in  
 166 Eq. (1). Building on their conceptualization, we made two modifications that align well with some  
 167 common sense about decision-making: (1) the initial state  $s_0$  is not generated from a plan, so we  
 168 move it into the conditioning; (2) the latent abstractions of the trajectories and their returns are better  
 169 to be distinct and associated instead of identical, so we reserve one vector  $z_y$  for returns, leave the  
 170 rest  $z_{\setminus y}$  for trajectories, and associate them in the prior  $p_{\alpha}(z)$ . This gives us the factorization:

$$p_{\theta}(\tau, y, z) = \rho(s_0)p_{\alpha}(z|s_0)p_{\beta}(\tau|s_0, z_{\setminus y})p_{\gamma}(y|z_y). \quad (4)$$

171 Following LPT, the prior  $p_{\alpha}(z)$  is implemented as an implicit generative model, but we replace the  
 172 UNet transformation [39] in LPT with a Transformer encoder with bidirectional mask. To sample  
 173  $z \sim p_{\alpha}(z)$ , we first sample a set of isotropic Gaussian  $\epsilon \sim \mathcal{N}(0, I)$  and transform them with the  
 174 Transformer  $z = f_{\alpha}(s_0, \epsilon)$ .

175 We also inherit LPT’s trajectory generator as a  $z$ -conditioned autoregressive model with a finite context  
 176 window of size  $M$ :  $p_\beta(\tau|s_0, z_{\setminus y}) = \prod_{t=0}^{T-1} p_\beta(a_t|s_{\leq t}, a_{< t}, z_{\setminus y})p_\beta(s_{t+1}|s_{t-M:t}, a_{t-M:t}, z_{\setminus y})$ . This  
 177 design forces the latent  $z$  to serve as global carriers of information, bridging temporal segments that  
 178 would otherwise be disconnected due to the limited context. The generator is parameterized by a  
 179 causal Transformer decoder that incorporates the latent  $z$  via cross-attention. The stepwise policy and  
 180 transition distributions are modeled as Gaussians with learnable or fixed variances.

181 The return predictor models  $p_\gamma(y|z_y)$  as a Gaussian  $\mathcal{N}(\mu_\gamma(z_y), \sigma_y^2)$ , where  $\mu_\gamma$  is an MLP that  
 182 predicts  $y$  from the special latent vector  $z_y$ . The variance  $\sigma_y^2$  is either treated as a hyperparameter or  
 183 learned jointly. Note that  $p_\gamma(y|z_y)$  being a Gaussian would not constrain the capability of modeling  
 184 multi-modal return distribution because  $p(y|\tau) = \int p(y|z_y)p(z_y|\tau)dz_y$  can be very expressive.

185 For a data point  $(\tau, y)$ , Maximum Likelihood Estimate (MLE) of latent-variable models is intractable.  
 186 Instead, we introduce a variational distribution  $q_\phi(\epsilon|\tau, y) = \mathcal{N}(\mu, \sigma^2)$ , parameterized by  $\phi = (\mu, \sigma)$ ,  
 187 *i.e.*, the mean vector  $\mu$  and a diagonal covariance matrix  $\sigma^2$  [40, 41], and maximize the evidence  
 188 lower bound (ELBO) [37, 38]:

$$ELBO(\theta, \phi) = \mathbb{E}_{q_\phi(\epsilon)}[\log p_\theta(\tau, y|s_0, \epsilon)] - D_{\text{KL}}(q_\phi(\epsilon)||p(\epsilon)). \quad (5)$$

189 We use the re-parametrization trick [42] in  $\mathbb{E}_q$ . The training employs alternating update of the  
 190 local parameters  $\phi$  specific to each  $(\tau, y)$  with classical Variational Bayes [38, 37] and the global  
 191 parameters  $\theta = (\alpha, \beta, \gamma)$  shared across all training data. Kong et al. [12] recently showed that  
 192 this training scheme can be scaled up to GPT-2 scale language models. To further accelerate and  
 193 stabilize training, we store the optimized  $\phi_{(\tau, y)}^*$  for each trajectory-return pair  $(\tau, y)$ , which are used  
 194 as initializations whenever the data point is redrawn to constitute the batch.

195 Inference queries in Eq. (2) and Eq. (3) can be adjusted accordingly to incorporate conditioning on  
 196  $s_0$  and reparametrization with  $\epsilon$ . Between them, Eq. (2) enables exploitation during test time, while  
 197 Eq. (3) facilitates exploratory online fine-tuning.

### 198 3.3 Exploratory Online Fine-tuning

199 A key advantage of the GAC framework is that the training objective—maximizing the ELBO in  
 200 Eq. (5)—remains consistent for both offline pre-training and online fine-tuning. The online phase  
 201 focuses on progressively improving the generative model by collecting higher-quality trajectories and  
 202 incorporating them into the training data. This is achieved through a principled exploration strategy  
 203 followed by staged model updates.

204 Our exploration mechanism is designed to sample trajectories from an improved distribution by  
 205 targeting returns slightly higher than what the model has previously achieved. Ideally, one might tilt  
 206 the model’s prior  $p(z)$  to generate a desired  $p(y^+)$ . However, we adopt a simpler and more direct  
 207 approach. We leverage a prioritized replay buffer, initialized with the offline dataset, which serves as  
 208 an empirical return distribution  $p_{\text{data}}(y)$ . To generate a target  $y^+$  from the target distribution  $p(y^+)$ ,  
 209 we first sample a high-performing return  $y$  from the top-k quantile of the replay buffer. We then add a  
 210 small, positive increment  $\Delta y$  to it, creating an optimistic target  $y^+ = y + \Delta y$ . This target is used to  
 211 condition the inference process as described in Eq. (3), yielding latent plans  $z \sim p(z|y^+)$  that guide  
 212 the agent to explore potentially out-of-distribution (OOD) yet superior trajectories.

213 While the returns collected by the actor  $p(z|y)p_{\text{data}}(y)$  from the environment are generally consistent  
 214 with the target  $y$  sampled from the empirical distribution, the introduction of  $\Delta y$  creates a mismatch  
 215 between the target distribution  $p(y^+)$  and the ground-truth return distribution. How to reliably set  
 216  $\Delta y$  to balance exploration with stability is an interesting research question that we leave for future  
 217 work. In this work, we treat  $\Delta y$  as a manually tuned hyperparameter. Empirically, we find that this  
 218 optimistic distributional targeting strategy effectively guides exploration; while it does not always lead  
 219 to trajectories that meet the optimistic  $y^+$ , it consistently outperforms fixed target  $y^*$  and gradually  
 220 shifts the distribution of collected returns toward higher values, as shown in Fig. 2.

221 The online finetuning process is structured in stages. In each stage, we first collect a set of new  
 222 trajectories using the exploration policy described above. These newly collected trajectory-return  
 223 pairs are then added to the replay buffer, enriching the dataset with recent, high-quality experiences.  
 224 Following data collection, we fine-tune the GAC model on the updated replay buffer for a set number  
 225 of steps, using the same ELBO maximization objective as in the pre-training phase. This cycle of

Table 1: **Offline results with and without Step-wise Rewards.** Comparison between the average normalized returns of our GAC method against several baselines, evaluated without any online finetuning. We assess the pure exploitation capability of our model, denoted as  $\text{GAC-}\mathbb{E}[y]$ , by optimizing over the entire latent space. For each of 5 independent seeds, we generate 100 evaluation trajectories from the pre-trained checkpoint and report the mean and standard deviation of their returns. The best result in each row is highlighted in bold.

Dataset	TD3+BC	IQL	CQL	DT	QDT	LPT	$\text{GAC-}\mathbb{E}[y]$
Results without step-wise rewards							
hopper-medium	27.6	35.1	23.3	57.3	50.7	58.5	<b>67.1</b> $\pm$ 0.8
hopper-medium-replay	10.6	13.9	7.7	50.8	38.7	71.2	<b>81.4</b> $\pm$ 1.1
walker2d-medium	67.8	49.1	0.0	69.9	63.7	77.8	<b>79.3</b> $\pm$ 0.9
walker2d-medium-replay	8.7	5.3	3.2	51.6	29.6	72.3	<b>78.9</b> $\pm$ 1.1
halfcheetah-medium	20.8	8.5	1.0	42.4	42.4	43.1	<b>43.3</b> $\pm$ 0.6
halfcheetah-medium-replay	12.8	5.2	7.8	33.3	32.8	39.6	<b>40.1</b> $\pm$ 0.4
Results with step-wise rewards							
hopper-medium	60.4	63.8	58.0	60.3	66.5	-	-
hopper-medium-replay	64.4	92.1	48.6	63.7	52.1	-	-
walker2d-medium	82.7	79.9	79.2	73.3	67.1	-	-
walker2d-medium-replay	85.6	73.7	74.1	60.2	58.2	-	-
halfcheetah-medium	48.1	47.4	44.4	42.1	42.3	-	-
halfcheetah-medium-replay	44.8	44.1	46.2	34.1	35.6	-	-

226 exploration, replay buffer update, and model fine-tuning is repeated, allowing GAC to progressively  
 227 adapt and improve its performance through online interaction.

## 228 4 Experiments

229 We evaluate GAC on standard offline and offline-to-online reinforcement learning benchmarks from  
 230 the D4RL Gym-MuJoCo suite (Halfcheetah, Hopper, and Walker2D). Our experiments are conducted  
 231 in a setting where only total trajectory returns are available. We compare GAC against state-of-the-  
 232 art methods and analyze the effectiveness of its different inference strategies for exploitation and  
 233 exploration.

234 We compare GAC against strong baselines across several paradigms: offline RL (IQL[26], CQL[32],  
 235 TD3+BC[24]), sequence modeling approaches (DT[7], ODT[8], QDT[43], LPT[12]), and classic  
 236 online RL (PPO[44], SAC[3]) for reference. Notably, except for LPT, these methods were designed  
 237 for step-wise rewards. We adapt them to our trajectory-return-only setting for a fair comparison,  
 238 highlighting the distinct advantage of models that do not rely on dense reward signals.

239 **Offline Pre-training.** We report the offline pre-training results in Table 1, evaluated using the  
 240 exploitative inference from Eq. (2) (termed  $\text{GAC-}\mathbb{E}[y]$ ). Our model consistently outperforms all  
 241 baselines across the tested environments. GAC’s advantage stems from its unique generative approach.  
 242 Unlike conventional methods rooted in temporal difference learning, GAC forgoes explicit per-step  
 243 credit assignment. Instead, it learns an implicit understanding of action-outcome relationships via  
 244 cross-attention across the latent token sequence. This allows the model to learn a smooth and  
 245 structured latent space of behaviors, enabling it to compose novel, high-quality trajectories by  
 246 interpolating between successful sub-sequences found in the training data. The robustness of this  
 247 approach is particularly evident in Mujoco environments. While TD-based methods like IQL, CQL,  
 248 and TD3+BC are heavily dependent on dense, step-wise rewards, GAC—much like other generative  
 249 models such as DT, QDT, and LPT—maintains high performance even with only trajectory-level  
 250 returns. This highlights a key strength of our framework: GAC achieves superior performance despite  
 251 operating without the granular, step-wise reward signals that most of these powerful baselines rely on  
 252 for effective training.

253 **Inference Strategies.** We analyze the practical implications of the different inference strategies  
 254 introduced in Section 3.1. We compare three approaches:  $\text{GAC-}\mathbb{E}[y]$  for pure exploitation via latent  
 255 space optimization;  $\text{GAC-}y^*$ , which conditions on a standard high-value target return similar to  
 256 DT[7]; and  $\text{GAC-}p(y^+)$ , our proposed exploration strategy that conditions on dynamically adjusted

Table 2: **Comparing inference objectives for pretrained GAC.** We report the mean and standard deviation of returns over 100 rollouts for exploitation (GAC- $\mathbb{E}[y]$ ), exploration (GAC- $p(y^+)$ ), fixed target conditioning (GAC- $y^*$ ), and sampling from the prior ( $p(y|z)p(z)$ ).

Dataset	GAC- $\mathbb{E}[y]$	GAC- $p(y^+)$	GAC- $y^*$	$p(y z)p(z)$
hopper-medium	<b>67.6</b> $\pm$ 5.2	63.1 $\pm$ <b>8.3</b>	60.6 $\pm$ 8.0	55.8 $\pm$ 5.2
hopper-medium-replay	<b>81.8</b> $\pm$ 9.9	66.1 $\pm$ 19.3	36.3 $\pm$ <b>23.8</b>	36.5 $\pm$ 23.4
walker2d-medium	<b>80.2</b> $\pm$ 4.6	78.9 $\pm$ <b>8.9</b>	78.5 $\pm$ 6.7	75.6 $\pm$ 10.1
walker2d-medium-replay	<b>79.3</b> $\pm$ 10.9	63.1 $\pm$ <b>18.6</b>	59.5 $\pm$ 14.1	48.9 $\pm$ 26.6
halfcheetah-medium	<b>42.6</b> $\pm$ 4.2	41.1 $\pm$ <b>6.0</b>	40.3 $\pm$ 4.2	41.2 $\pm$ 7.5
halfcheetah-medium-replay	<b>39.8</b> $\pm$ 7.8	38.8 $\pm$ <b>8.9</b>	36.7 $\pm$ 8.2	33.0 $\pm$ 12.0

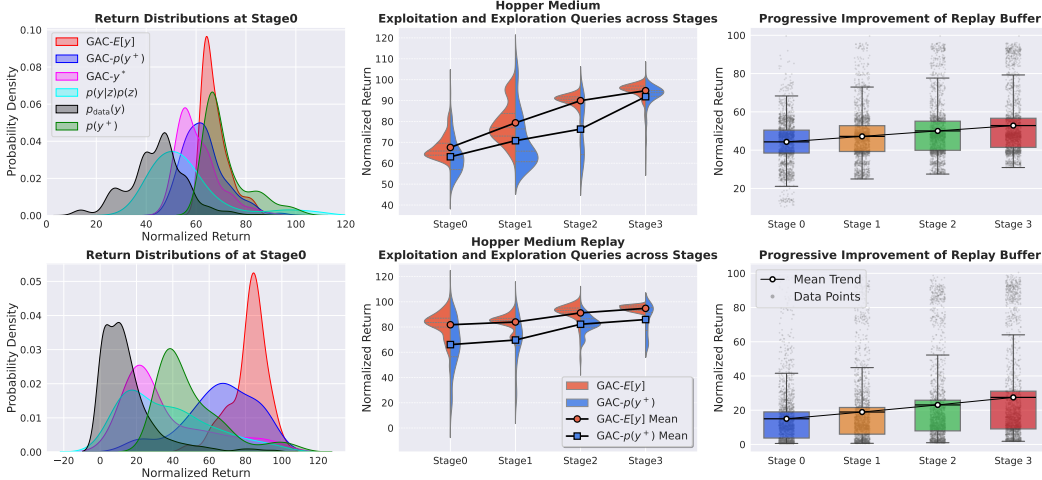


Figure 2: **Illustration of GAC’s return distributions.** The left column displays the return distributions achieved by various inference methods after offline pre-training, highlighting the explorative and exploitative behavior. The middle column presents a split violinplot for GAC- $\mathbb{E}[y]$  and GAC- $p(y^+)$  inference, showing the performance gains of checkpoints and the explorative and exploitative behavior at three sequential fine-tuning stages. The right column visualizes the progressive improvement in the complete dataset’s quality, as evidenced by the upward shift of the dataset’s return distribution over the fine-tuning stages.

257 targets sampled from the replay buffer. This comparison is designed to evaluate the trade-offs between  
 258 pure optimization, fixed-target conditioning, and our more adaptive, data-driven exploration method.

259 The return distributions collected by actors from different inference objectives, visualized in the  
 260 left column of Fig. 2 and summarized in Table 2, highlight their distinct behaviors. The return  
 261 distribution  $p(y|z)$ , generated from random prior latent plans from  $p(z)$ , is dispersed by principle  
 262 and effectively covers the returns seen in the training data, confirming that the model has learned a  
 263 well-structured latent space without collapsing to a single mode. In sharp contrast, the GAC- $\mathbb{E}[y]$   
 264 strategy produces a focused, low-variance distribution concentrated on high-certainty, high-return  
 265 outcomes, underscoring its efficacy as a pure exploitation method. To encourage exploration, we  
 266 sample target returns from  $p(y^+)$ , which apparently have better coverage at high-return and even  
 267 out-of-distribution regions than GAC- $\mathbb{E}[y]$ . Based on these targets, GAC- $p(y^+)$  achieves a high  
 268 mean return but with significantly greater variance than GAC- $\mathbb{E}[y]$ . While GAC- $p(y^+)$ ’s optimistic  
 269 conditioning on  $y^+$  leads to a deliberate mismatch between its resulting distribution and  $p(y^+)$ , it  
 270 reflects a realistic improvement over the  $p_{\text{data}}(y)$ , showcasing successful guided exploration toward  
 271 better outcomes. As GAC- $p(y^+)$ ’s target distribution is anchored at the data distribution, it also  
 272 outperforms GAC- $y^*$ , whose reliance on a single fixed target often leads to overestimation (*optimism*  
 273 *bias*) of the outcomes.

274 **Online Fine-tuning.** Table 3 reports the final returns compared to other methods in most environments  
 275 after fine-tuning with just 300 online trajectories. While the gains are relatively small in some tasks,  
 276 this is often a consequence of its higher initial offline performance. Notably, GAC remains highly  
 277 effective even without dense, step-wise rewards—a setting where pure online RL methods like  
 278 PPO and SAC typically fail. Unlike methods that rely on explicit optimism or conservatism, GAC

Table 3: **Online training results with and without step-wise rewards.** Comparison of normalized returns before and after online fine-tuning. We fine-tune GAC for 3 stages with data collected by GAC- $p(y^+)$  and report the final GAC- $\mathbb{E}[y]$  performance from five seeds. The best final result for each dataset is highlighted in bold.

Dataset	Online RL		ODT		CQL		IQL		LPT		GAC	
	PPO	SAC	online	$\delta$	online	$\delta$	online	$\delta$	online	$\delta$	online	$\delta$
Results without step-wise rewards												
hopper-m	11.5	9.8	57.6	0.3	29.6	6.3	25.0	-10.1	64.8	6.3	<b>94.1</b> $\pm$ 1.1	27.0
hopper-m-r	-	-	65.2	14.4	8.4	0.7	12.6	-1.3	72.4	1.2	<b>95.5</b> $\pm$ 0.8	14.1
walker2d-m	9.1	4.0	70.7	0.8	1.9	1.9	50.1	1.0	79.5	1.7	<b>84.7</b> $\pm$ 0.3	5.4
walker2d-m-r	-	-	57.3	5.7	0.5	-2.7	6.9	1.6	79.0	6.7	<b>84.6</b> $\pm$ 0.5	5.7
halfcheetah-m	1.2	1.9	40.7	-1.7	2.8	1.8	8.9	0.4	43.2	0.1	<b>44.3</b> $\pm$ 0.1	1.0
halfcheetah-m-r	-	-	24.4	-8.4	6.4	-1.4	7.4	2.2	40.6	1.0	<b>40.8</b> $\pm$ 0.1	0.7
Results with step-wise rewards												
hopper-m	77.8	75.0	97.5	30.6	60.4	2.4	66.8	3.0	-	-	-	-
hopper-m-r	-	-	88.9	2.3	56.3	7.7	96.2	4.1	-	-	-	-
walker2d-m	43.9	47.5	76.8	4.6	79.6	0.4	80.3	0.4	-	-	-	-
walker2d-m-r	-	-	76.9	4.0	75.8	1.7	70.6	-3.1	-	-	-	-
halfcheetah-m	25.0	24.4	42.4	-0.6	45.6	1.2	47.4	0.0	-	-	-	-
halfcheetah-m-r	-	-	40.2	0.4	43.0	-3.2	44.1	0.0	-	-	-	-

279 improves by leveraging its principled exploration strategy. This strategy systematically probes for  
 280 trajectories with higher returns by conditioning on targets set slightly beyond the empirical best  
 281 of the current dataset. When these exploratory rollouts yield higher-performing trajectories, they  
 282 are added to the replay buffer. This process gradually shifts the collected data distribution toward  
 283 superior outcomes, which is exemplified in the Hopper task where our method achieves its most  
 284 significant performance gain. To visualize this dynamic process, Fig. 2 plots the return distributions  
 285 from different inference objectives across several fine-tuning stages. The figure clearly illustrates the  
 286 interplay between exploration and exploitation, tracking the progressive performance improvement  
 287 via the continuous, positive shift in the collected data’s return distribution. The implementation details  
 288 and more results can be found in Appendix A.

## 289 5 Conclusion

290 In this paper, we introduced Generative Actor-Critic (GAC), a novel framework that decouples  
 291 decision-making from generative decision-modeling by learning the joint distribution of trajecto-  
 292 ries and their returns,  $p(\tau, y)$ . Our latent variable-based instantiation enables principled inference  
 293 strategies for both exploitation and exploration, leading to state-of-the-art offline performance and  
 294 highly efficient online adaptation on Gym-MuJoCo benchmarks in the absence of step-wise rewards.  
 295 Ultimately, GAC provides an effective framework that bridges generative modeling and reinforcement  
 296 learning, opening a promising path toward building more capable and adaptable autonomous agents.

## 297 References

- 298 [1] Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1.  
 299 MIT press Cambridge, 1998.
- 300 [2] Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information*  
 301 *processing systems*, 12, 1999.
- 302 [3] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-  
 303 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*  
 304 *conference on machine learning*, pages 1861–1870. Pmlr, 2018.
- 305 [4] Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforce-  
 306 ment learning. In *International conference on machine learning*, pages 449–458. PMLR,  
 307 2017.

- 308 [5] Gabriel Barth-Maroon, Matthew W Hoffman, David Budden, Will Dabney, Dan Horgan,  
309 TB Dhruva, Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributed distributional  
310 deterministic policy gradients. In *International Conference on Learning Representations*, 2018.
- 311 [6] Audrunas Gruslys, Will Dabney, Mohammad Gheshlaghi Azar, Bilal Piot, Marc Bellemare,  
312 and Remi Munos. The reactor: A fast and sample-efficient actor-critic agent for reinforcement  
313 learning. In *International Conference on Learning Representations*, 2018.
- 314 [7] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter  
315 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: reinforcement learning via  
316 sequence modeling, 2 june 2021. URL <http://arxiv.org/abs/2106.01345>, 2021.
- 317 [8] Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In *international*  
318 *conference on machine learning*, pages 27042–27059. PMLR, 2022.
- 319 [9] Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion  
320 for flexible behavior synthesis. In *International Conference on Machine Learning*, pages  
321 9902–9915. PMLR, 2022.
- 322 [10] Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B Tenenbaum, Tommi S Jaakkola, and Pulkit  
323 Agrawal. Is conditional generative modeling all you need for decision making? In *The Eleventh*  
324 *International Conference on Learning Representations*, 2023.
- 325 [11] Kuang-Huei Lee, Ofir Nachum, Mengjiao Sherry Yang, Lisa Lee, Daniel Freeman, Sergio  
326 Guadarrama, Ian Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, et al. Multi-game  
327 decision transformers. *Advances in Neural Information Processing Systems*, 35:27921–27936,  
328 2022.
- 329 [12] Deqian Kong, Dehong Xu, Minglu Zhao, Bo Pang, Jianwen Xie, Andrew Lizarraga, Yuhao  
330 Huang, Sirui Xie, and Ying Nian Wu. Latent plan transformer for trajectory abstraction:  
331 Planning as latent space inference. *Advances in Neural Information Processing Systems*, 37:  
332 123379–123401, 2024.
- 333 [13] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G  
334 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.  
335 Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 336 [14] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay.  
337 *arXiv preprint arXiv:1511.05952*, 2015.
- 338 [15] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust  
339 region policy optimization. In *International conference on machine learning*, pages 1889–1897.  
340 PMLR, 2015.
- 341 [16] Brian D Ziebart. *Modeling purposeful adaptive behavior with the principle of maximum causal*  
342 *entropy*. Carnegie Mellon University, 2010.
- 343 [17] Roy Fox, Ari Pakman, and Naftali Tishby. Taming the noise in reinforcement learning via soft  
344 updates. In *32nd Conference on Uncertainty in Artificial Intelligence 2016, UAI 2016*, pages  
345 202–211. Association For Uncertainty in Artificial Intelligence (AUAI), 2016.
- 346 [18] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.  
347 Deterministic policy gradient algorithms. In *International conference on machine learning*,  
348 pages 387–395. Pmlr, 2014.
- 349 [19] Andrew Kachites McCallum. *Reinforcement learning with selective perception and hidden*  
350 *state*. University of Rochester, 1996.
- 351 [20] Rupesh Kumar Srivastava, Pranav Shyam, Filipe Mutz, Wojciech Jaśkowski, and Jürgen  
352 Schmidhuber. Training agents using upside-down reinforcement learning. *arXiv preprint*  
353 *arXiv:1912.02877*, 2019.
- 354 [21] Scott Emmons, Benjamin Eysenbach, Ilya Kostrikov, and Sergey Levine. Rvs: What is essential  
355 for offline rl via supervised learning? *arXiv preprint arXiv:2112.10751*, 2021.

- 356 [22] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for  
357 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 358 [23] Jialong Wu, Haixu Wu, Zihan Qiu, Jianmin Wang, and Mingsheng Long. Supported policy  
359 optimization for offline reinforcement learning. *Advances in Neural Information Processing*  
360 *Systems*, 35:31278–31291, 2022.
- 361 [24] Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.  
362 *Advances in neural information processing systems*, 34:20132–20145, 2021.
- 363 [25] Jiafei Lyu, Xiaoteng Ma, Xiu Li, and Zongqing Lu. Mildly conservative q-learning for offline  
364 reinforcement learning. *Advances in Neural Information Processing Systems*, 35:1711–1724,  
365 2022.
- 366 [26] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit  
367 q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
- 368 [27] Jianxiong Li, Xiao Hu, Haoran Xu, Jingjing Liu, Xianyuan Zhan, and Ya-Qin Zhang.  
369 Proto: Iterative policy regularized offline-to-online reinforcement learning. *arXiv preprint*  
370 *arXiv:2305.15669*, 2023.
- 371 [28] Alex Beeson and Giovanni Montana. Improving td3-bc: Relaxed policy constraint for offline  
372 learning and stable online fine-tuning. *arXiv preprint arXiv:2211.11802*, 2022.
- 373 [29] Haichao Zhang, We Xu, and Haonan Yu. Policy expansion for bridging offline-to-online  
374 reinforcement learning. *arXiv preprint arXiv:2302.00935*, 2023.
- 375 [30] Mitsuhiro Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn,  
376 Aviral Kumar, and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online  
377 fine-tuning. *Advances in Neural Information Processing Systems*, 36:62244–62269, 2023.
- 378 [31] Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online  
379 reinforcement learning via balanced replay and pessimistic q-ensemble. In *Conference on Robot*  
380 *Learning*, pages 1702–1712. PMLR, 2022.
- 381 [32] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning  
382 for offline reinforcement learning. *Advances in neural information processing systems*, 33:  
383 1179–1191, 2020.
- 384 [33] Matthew Botvinick and Marc Toussaint. Planning as inference. *Trends in cognitive sciences*, 16  
385 (10):485–488, 2012.
- 386 [34] Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and  
387 review. *arXiv preprint arXiv:1805.00909*, 2018.
- 388 [35] Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and  
389 Martin Riedmiller. Maximum a posteriori policy optimisation. In *International Conference on*  
390 *Learning Representations*, 2018.
- 391 [36] Aoyang Qin, Feng Gao, Qing Li, Song-Chun Zhu, and Sirui Xie. Learning non-markovian  
392 decision-making from state-only sequences. *Advances in Neural Information Processing*  
393 *Systems*, 36:6596–6618, 2023.
- 394 [37] Kevin P Murphy. Machine learning: A probabilistic perspective (adaptive computation and  
395 machine learning series). *The MIT Press: London, UK*, 2018.
- 396 [38] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational  
397 inference. *the Journal of machine Learning research*, 14(1):1303–1347, 2013.
- 398 [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks  
399 for biomedical image segmentation. In *Medical image computing and computer-assisted*  
400 *intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5–9,*  
401 *2015, proceedings, part III 18*, pages 234–241. Springer, 2015.

- 402 [40] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for  
403 statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- 404 [41] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduc-  
405 tion to variational methods for graphical models. *Machine learning*, 37:183–233, 1999.
- 406 [42] Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013.
- 407 [43] Taku Yamagata, Ahmed Khalil, and Raul Santos-Rodriguez. Q-learning decision transformer:  
408 Leveraging dynamic programming for conditional sequence modelling in offline rl. In *Internat-  
409 ional Conference on Machine Learning*, pages 38989–39007. PMLR, 2023.
- 410 [44] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
411 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 412 [45] Ying Li, Lei Cheng, Feng Yin, Michael Minyi Zhang, and Sergios Theodoridis. Overcoming  
413 posterior collapse in variational autoencoders via em-type training. In *ICASSP 2023-2023 IEEE  
414 International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5.  
415 IEEE, 2023.
- 416 [46] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna:  
417 A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM  
418 SIGKDD international conference on knowledge discovery & data mining*, pages 2623–2631,  
419 2019.

Table 4: GAC Hyperparameters.

Architecture-Related	Range
Embedding dimension	choice(64, 128, 192, 256, 512)
Number of latent tokens	choice(1, 2, 4, 8)
Number of attention heads	choice(1, 2, 4, 8)
Number of decoder layers	choice(1, 2, 3, 4)
Number of encoder layers	choice(1, 2, 3, 4)
Context length	choice(1, 4, 16, 32, 64)
Architecture-Unrelated	Range
Outer-loop learning rate	interval(1e-4, 1e-2)
Outer-loop weight decay	interval(1e-4, 1e-2)
Outer-loop training steps	choice(1, 5, 10, 20)
Inner-loop learning rate	interval(0.01, 0.5)
Batch size	choice(800, 700, 600, 500)

Table 5: The best GAC Architecture-Related Hyperparameters.

Parameter	h-m	h-m-r	w-m	w-m-r	ha-m	ha-m-r
Embedding dimension	192	128	192	192	128	256
Number of latent tokens	4	8	4	4	8	4
Number of attention heads	4	4	2	8	1	4
Number of decoder layers	3	4	3	4	3	3
Number of encoder layers	4	1	3	3	2	2
Context length	4	4	1	1	4	1

## 420 A More Details on MuJoCo

421 Inspired by Li et al. [45], we employ a more sophisticated training methodology than that of LPT  
422 [12]. This approach aims to tighten the approximation between the Evidence Lower Bound (ELBO)  
423 and the log-likelihood, and mitigating the issue of posterior collapse. Specifically, for each data  
424 point, we conduct inner-loop optimization until convergence, which is defined by a maximum of  
425 100 training steps and an early-stopping threshold of 1e-4. In the outer loop, unlike the single-step  
426 approach in [12], we treat the number of training steps as a tunable hyperparameter. Given that  
427 GAC decouples decision-modeling from decision-making, we can directly assess the quality of  
428 our generative model by its loss value, rather than relying on immediate policy evaluation in the  
429 environment. Consequently, we leverage the Optuna framework [46] for hyperparameter optimization,  
430 with the objective of maximizing the ELBO. The complete set of hyperparameters subject to tuning  
431 is detailed in Table 4. Following the pre-training phase, architecture-related hyperparameters are  
432 fixed, while the remaining parameters are optimized during each subsequent fine-tuning stage. Using  
433 the optimal hyperparameters identified (see Table 5 and Table 6<sup>1</sup>), we then configure the inner-loop  
434 inference with more stringent convergence criteria: a maximum of 1000 steps and an early-stopping  
435 threshold of 1e-5.

436 For GAC- $p(y^+)$  to collect online data, we introduce another two parameters to obtain a target return  
437 set. We perform a data processing procedure on the current return dataset, which we denote as  $D$ .  
438 This process involves a quantile-based filtering and a subsequent transformation. The procedure is  
439 parameterized by two tunable hyperparameters: the quantile threshold,  $q$ , and an additive constant,  $\Delta y$ .  
440 First, we determine the  $q$ -th quantile of the dataset  $D$ , which we denote as  $y_q$ . All data points in  $D$  that  
441 are greater than  $y_q$  are selected to form a new subset,  $D_{\text{sub}} = \{y \in D \mid y > y_q\}$ . We finally randomly  
442 sample from  $D_{\text{sub}}$  and plus  $\Delta y$  to obtain the dynamic target return as  $\{y + \Delta y \mid y \in D_{\text{sub}}\}$ . This  
443 procedure allows us to isolate and transform the upper tail of the data distribution, with the specific  
444 threshold and transformation intensity controlled by the hyperparameters  $q$  and  $\Delta y$ , respectively. We  
445 present these additional parameters for inference and the standard target return for each environment

<sup>1</sup>The hyperparameters with interval range are fp64 after searching, only displayed four decimal places in Table 6. 'h', 'w', 'ha', 'm' and 'r' denotes 'hopper', 'walker2d', 'halfcheetah', 'medium' and 'replay' respectively.

Table 6: The best GAC Architecture-Unrelated Hyperparameters.

Parameter	h-m	h-m-r	w-m	w-m-r	ha-m	ha-m-r
Stage 0 (pre-training)						
Outer-loop learning rate	0.0010	0.0003	0.0023	0.0026	0.0026	0.0012
Outer-loop weight decay	0.0002	0.0056	0.0027	0.0011	0.0020	0.0071
Outer-loop training steps	10	1	5	1	10	10
Inner-loop learning rate	0.2365	0.2058	0.3980	0.0795	0.0501	0.0452
Batch size	700	500	600	500	500	700
Stage 1 (fine-tuning)						
Outer-loop learning rate	0.0002	0.0018	0.0008	0.0017	0.0062	0.0009
Outer-loop weight decay	0.0002	0.0026	0.0029	0.0062	0.0044	0.0060
Outer-loop training steps	1	20	1	20	5	20
Inner-loop learning rate	0.0458	0.2116	0.1223	0.0712	0.2363	0.2836
Batch size	800	700	800	500	700	500
Stage 2 (fine-tuning)						
Outer-loop learning rate	0.0004	0.0010	0.0013	0.0014	0.0093	0.0013
Outer-loop weight decay	0.0006	0.0007	0.0022	0.0050	0.0001	0.0002
Outer-loop training steps	1	1	5	20	5	5
Inner-loop learning rate	0.1018	0.0712	0.0124	0.1691	0.1435	0.2154
Batch size	700	500	800	500	600	700
Stage 3 (fine-tuning)						
Outer-loop learning rate	0.3755	0.0005	0.0007	0.0010	0.0039	0.0022
Outer-loop weight decay	0.0077	0.0001	0.0027	0.0007	0.0023	0.0015
Outer-loop training steps	5	5	20	1	20	20
Inner-loop learning rate	0.0023	0.010	0.0337	0.0712	0.1193	0.1984
Batch size	600	500	600	500	500	600

Table 7: The Inference Hyperparameters.

Parameter	h-m	h-m-r	w-m	w-m-r	ha-m	ha-m-r
Standard target return	3600	3600	6000	6000	5000	5000
Stage0 ( $q, \Delta y$ )	(0.8,10)	(0.8,10)	(0.8,10)	(0.8,10)	(0.6,3)	(0.6,3)
Stage1 ( $q, \Delta y$ )	(0.8,10)	(0.8,10)	(0.8,10)	(0.8,10)	(0.6,3)	(0.6,3)
Stage2 ( $q, \Delta y$ )	(0.8,10)	(0.8,10)	(0.6,3)	(0.6,5)	(0.6,3)	(0.6,1)
Stage3 ( $q, \Delta y$ )	(0.8,10)	(0.8,10)	(0.6,3)	(0.6,5)	(0.6,3)	(0.6,1)

446 in Table 7. We finally present all of the inference results for all Gym-MuJoCo tasks with GAC- $\mathbb{E}[y]$ ,  
447 GAC- $p(y^+)$ , GAC- $y^*$  and PRIOR inference strategies for one seed in Fig. 1.

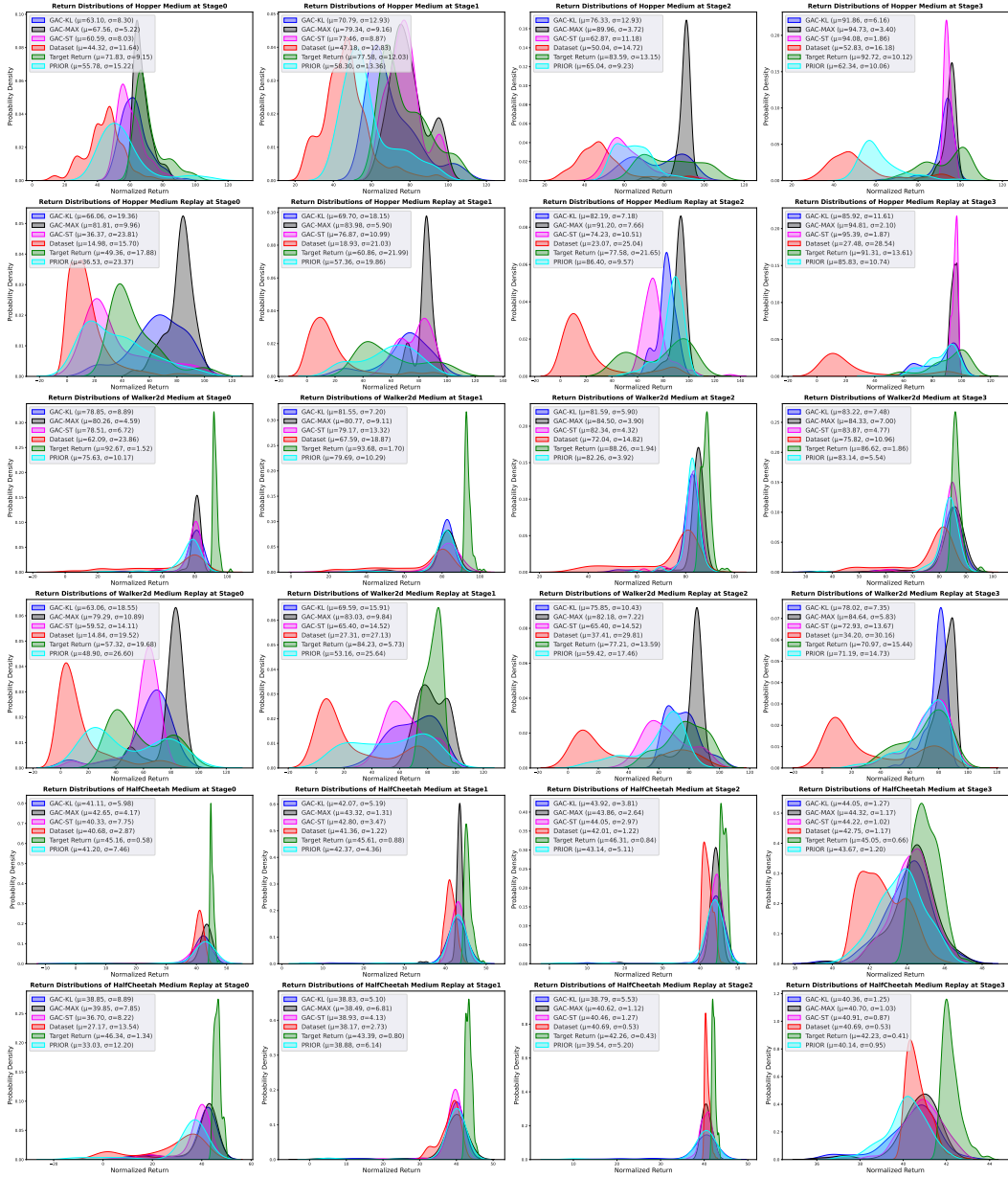


Figure 1: Complete MuJoCo Results.