

---

# A Variational Formulation of Reinforcement Learning in Infinite-Horizon Markov Decision Processes

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Reinforcement learning in infinite-horizon Markov decision processes (MDPs) is  
2 typically framed as expected discounted return maximization. In this paper, we  
3 formulate an alternative principle for optimal sequential decision-making in infinite-  
4 horizon MDPs: variational Bayesian inference in transdimensional probabilistic  
5 models. In particular, we specify a probabilistic model of random size and consider  
6 the variational problem of finding an approximation to the posterior distribution  
7 over state–action trajectories conditioned on state–action trajectories that reflect a  
8 desired behavior. We derive a tractable variational objective for infinite-horizon  
9 settings, prove a variational dynamic-discount policy iteration theorem, show that  
10 fixed discount factor KL-regularized reinforcement learning objectives are special  
11 cases of dynamic-discount variational objectives, and prove that learning dynamic  
12 discount factors is optimal.

## 13 1 Introduction

14 We provide a Bayesian framework for deriving behavior-driven optimal decision rules for sequential  
15 decision problems. In particular, we provide a mathematical justification for learned, dynamic dis-  
16 count factors in KL-regularized reinforcement learning, which have been proposed as an empirically  
17 useful tool in recently developed reinforcement learning algorithms [6, 8, 4], and establish a rigorous  
18 foundation for framing modern reinforcement learning methods as probabilistic inference. Although  
19 control as inference has gained in popularity, the treatment of infinite-horizon settings in previous  
20 works is ad-hoc and not probabilistically well-motivated. With this work, we hope to address this  
21 shortcoming and provide a clear formulation of control as inference that carefully disambiguates  
22 modeling and inference assumptions.

23 Levine [3] and Haarnoja et al. [1] presented a framework for framing maximum-entropy reinforcement  
24 learning as Bayesian inference in probabilistic models over finite-horizon state–action trajectories.  
25 However, most modern reinforcement learning problems are not formulated as finite but as *infinite-*  
26 *horizon* problems [5, 9]. To apply their probabilistic formulation of reinforcement  
27 learning to infinite-horizon problems, Levine [3] and Haarnoja et al. [1] introduce a fixed discount-  
28 factor into their formulation post-hoc and without providing a probabilistic justification for doing  
29 so. In this paper, we show that including a (fixed) discount factor as proposed by Levine [3] and  
30 Haarnoja et al. [1] is a special case of a more general probabilistic framing of the problem, leads to  
31 a variational formulation with a loose evidence lower bound, and can provably be improved upon  
32 by framing Bayesian variational inference in infinite-horizon MDPs as variational inference in a  
33 *transdimensional* probabilistic model.

34 To derive a learning algorithm that allows us to infer a policy that reflects the behavior encoded  
35 in desired state trajectories, we frame the problem of finding an optimal policy as computing  
36 an approximation to the conditional distribution over state–action trajectories given state–action

37 trajectories that reflect a desired behavior. We formulate a corresponding probabilistic model and  
 38 derive tractable variational objectives for finite- and infinite-horizon settings. Based on these results,  
 39 we define a novel Bellman backup operator and show that for tabular settings, the repeated application  
 40 of the operator converges to an optimal policy and an optimal dynamic discount factor. Building  
 41 on this result, we show that fixed discount factor KL-regularized reinforcement learning objectives  
 42 are special cases of the dynamic-discount objectives derived here and demonstrate that variationally  
 43 learned, dynamic discount factors are optimal in KL-regularized RL.

## 44 2 Preliminaries

45 Standard reinforcement learning (RL) addresses reward maximization in a Markov decision pro-  
 46 cess (MDP) defined by the tuple  $(\mathcal{S}, \mathcal{A}, p_{\mathbf{S}_0}, p_d, r, \gamma)$  [10, 11], where  $\mathcal{S}$  and  $\mathcal{A}$  denote the state  
 47 and action space, respectively,  $p_0$  denotes the initial state distribution,  $p_d$  is a state transition dis-  
 48 tribution,  $r$  is an immediate reward function, and  $\gamma$  is a discount factor. To sample trajectories,  
 49 an initial state is sampled according to  $p_{\mathbf{S}_0}$ , and successive states are sampled from the state tran-  
 50 sition distribution  $\mathbf{S}_{t+1} \sim p_d(\cdot | \mathbf{s}_t, \mathbf{a}_t)$  and actions from a policy  $\mathbf{A}_t \sim \pi(\cdot | \mathbf{s}_t)$ . We will write  
 51  $\mathcal{T}_{0:t} = \{\mathbf{S}_0, \mathbf{A}_0, \mathbf{S}_1, \dots, \mathbf{S}_t, \mathbf{A}_t\}$  to represent a finite-horizon and  $\mathcal{T}_0 \doteq \{\mathbf{S}_t, \mathbf{A}_t\}_{t=0}^{\infty}$  to represent  
 52 an infinite-horizon stochastic state–action trajectory, and write  $\tau_{0:t} = \{\mathbf{s}_0, \mathbf{a}_0, \mathbf{s}_1, \dots, \mathbf{s}_t, \mathbf{a}_t\}$  and  
 53  $\tau_0 \doteq \{\mathbf{s}_t, \mathbf{a}_t\}_{t=0}^{\infty}$  for the respective trajectory realizations. Given a reward function  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$   
 54 and discount factor  $\gamma \in (0, 1)$ , the objective in reinforcement learning is to find a policy  $\pi$  that  
 55 maximizes the returns, defined as  $\mathbb{E}_{p_\pi} [\sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t)]$ , where  $p_\pi$  denotes the distribution of states  
 56 induced by a policy  $\pi$ .

## 57 3 A Variational Formulation of RL in Infinite-Horizon MDPs

58 Desired behaviors for artificial agents are often abstract and hard to encode into reward functions.  
 59 However, in practice, it is often easy to represent desired behaviors via demonstrations. Such  
 60 demonstrations can be thought of as sample state–action trajectories from a distribution over optimal  
 61 state–action trajectories. In the remainder of this paper, we will demonstrate how to use variational  
 62 Bayesian inference to infer an optimal policy from a set of optimal state–action demonstrations.

63 To start the exposition, we note that for every state in the environment, there exists a desired, or  
 64 optimal, behavior that an agent *could* take. We denote this optimal behavior for any given state as  
 65 the set of state–action trajectories by  $\tau^\Omega$ . Throughout, we will use the index  $\Omega$  to denote optimality.  
 66 Hence, for any state  $\mathbf{s} \in \mathcal{S}$ , assuming the MDP is ergodic and transition dynamics are deterministic,  
 67 there exists a set of actions that will set an agent on an optimal state–action trajectory, that is,  
 68  $\mathcal{A}^\Omega \doteq \{\mathbf{a} \in \mathcal{A} | \mathbf{s}' \sim p_d(\mathbf{s}' | \mathbf{s}, \mathbf{a}) : \mathbf{s}' \in \tau^\Omega\}$ , meaning there exists a set of actions that will set an  
 69 agent on the optimal state–action trajectory with probability one.

70 If the transition dynamics are stochastic, each state–action pair will have some probability less than  
 71 one of transitioning the agent onto an optimal state–action trajectory. Denoting the event of a state  
 72 being in the optimal state–action trajectory by  $\mathbf{s} \in \mathcal{S}^\Omega$ , where  $\mathcal{S}^\Omega \doteq \{\mathbf{s} \in \tau^\Omega\}$ , we can define a  
 73 random variable  $\xi(\mathcal{S}^\Omega) \doteq \mathbb{I}\{\mathbf{s}' \in \mathcal{S}^\Omega\}$ . We then have that  $\xi = 1$  if the state  $\mathbf{s}'$  into which an agent  
 74 transitioned after taking action  $\mathbf{a}$  in state  $\mathbf{s}$  is in the optimal trajectory and  $\xi(\mathcal{S}^\Omega) = 0$  otherwise. The  
 75 probability of transitioning into a state on the optimal state–action trajectory at time step  $t + 1$  is then  
 76 given by

$$\mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) = \int_{\mathcal{S}^\Omega} p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) d\mathbf{s}_{t+1} = \int_{\mathcal{S}} p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \mathbb{I}\{\mathbf{s}_{t+1} \in \mathcal{S}^\Omega\} d\mathbf{s}_{t+1}. \quad (1)$$

77 In other words, the probability of transitioning into a state on the optimal state–action trajectory  
 78 corresponds to marginalization over the set of optimal states  $\mathcal{S}^\Omega$ . Equation (1) is a likelihood function.

79 Similarly, by the Markov property, the joint probability of transitioning into a state on the optimal  
 80 state–action trajectory and staying on it from time step 1 to time step  $t^* \doteq t + 1$ , given a state–action  
 81 trajectory, factorizes and is given by

$$\begin{aligned} \mathbb{P}(\xi_{1:t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_t, \mathbf{a}_t) \\ = \prod_{t'=0}^t \int_{\mathcal{S}^\Omega} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) d\mathbf{s}_{t'+1} = \prod_{t'=0}^t \int_{\mathcal{S}} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \mathbb{I}\{\mathbf{s}_{t'+1} \in \mathcal{S}^\Omega\} d\mathbf{s}_{t'+1}, \end{aligned} \quad (2)$$

82 where we start from  $t' = 0$  without loss of generality.

### 83 3.1 Warm-Up: Finite-Horizon Reinforcement Learning as Variational Inference

84 First, we consider the finite-horizon setting. This formulation only diverges slightly from prior work  
85 but will help us transition to the transdimensional model formulation for the infinite-horizon setting.

86 With the notion of trajectory-dependent optimality described in the previous section, we can now  
87 specify a model over finite-horizon state–action trajectories and  $\xi_{1:t+1}(\mathcal{S}^\Omega)$ ,

$$p(\tau_{0:t}, \xi_{1:t+1}^*(\mathcal{S}^\Omega)) \doteq p_{\mathbf{S}_0}(\mathbf{s}_0) \prod_{t'=0}^t \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_t \mid \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} \mid \mathbf{s}_{t'}),$$

88 where  $\tilde{\tau}_{0:t}$  is a state–action trajectory starting at state  $\mathbf{S}_0$  and ending at state  $\mathbf{S}_t$ ,  $p(\mathbf{a}_t \mid \mathbf{s}_t)$  is  
89 a conditional action prior,  $p_d(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)$  is the environment’s state transition distribution, and  
90  $\xi_{1:t+1}^*(\mathcal{S}^\Omega) \doteq \{\xi_{t'}(\mathcal{S}^\Omega) = 1\}_{t'=1}^{t+1}$  is the set of events corresponding to transitioning onto an optimal  
91 trajectory. By extension, the probability of transitioning onto an optimal state–action trajectory and  
92 remaining on it for  $t^*$  time steps given a state and a prior policy is given by the marginal likelihood

$$\mathbb{P}(\xi_{1:t+1}(\mathcal{S}^\Omega) = 1 \mid \mathbf{s}_0) = \iint_{\mathcal{A}^{t+1} \mathcal{S}^t} p_{\tilde{\tau}_{0:t} \mid \mathbf{S}_0}(\tilde{\tau}_{0:t} \mid \mathbf{s}_0) \left( \prod_{t'=0}^t \int_{\mathcal{S}^\Omega} p_d(\mathbf{s}_{t'+1} \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) d\mathbf{s}_{t'+1} \right) d\mathbf{s}_{1:t} d\mathbf{a}_{0:t} \quad (3)$$

93 where  $p_{\tilde{\tau}_{0:t} \mid \mathbf{S}_0}(\tilde{\tau}_{0:t} \mid \mathbf{s}_0) \doteq p(\mathbf{a}_t \mid \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} \mid \mathbf{s}_{t'}) \quad (4)$

94 is a prior distribution over state–action trajectories. Using an indicator function  $\mathbb{I}\{\mathbf{s}_{t+1} \in \mathcal{S}^\Omega\}$   
95 denoting whether the next state is on the desired state–action trajectory, the marginal likelihood in  
96 Equation (3) can equivalently be expressed as

$$\begin{aligned} & \mathbb{P}(\xi_{1:t+1}(\mathcal{S}^\Omega) = 1 \mid \mathbf{s}_0) \\ &= \iint_{\mathcal{A}^{t+1} \mathcal{S}^t} p_{\tilde{\tau}_{0:t} \mid \mathbf{S}_0}(\tilde{\tau}_{0:t} \mid \mathbf{s}_0) \left( \prod_{t'=0}^t \int_{\mathcal{S}} p_d(\mathbf{s}_{t'+1} \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) \mathbb{I}\{\mathbf{s}_{t'+1} \in \mathcal{S}^\Omega\} d\mathbf{s}_{t'+1} \right) d\mathbf{s}_{1:t} d\mathbf{a}_{0:t}. \end{aligned} \quad (5)$$

97 This marginalization establishes the connection between the full joint distribution in Equation (3) and  
98 the likelihood of remaining on an optimal state–action trajectory under a state–action trajectory prior  
99 and the likelihood function defined in Equation (1).

100  $\mathbb{P}(\xi_{1:t+1}(\mathcal{S}^\Omega) = 1 \mid \mathbf{s}_0)$  is the marginal likelihood of remaining on the optimal state trajectory from  
101 time step 1 to time step  $t+1$  under the prior policy  $p(\mathbf{a}_t \mid \mathbf{s}_t)$  and the dynamics model  $p_d(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)$ .  
102 Using Bayes’ Theorem, we could use the marginal likelihood to compute the posterior distribution  
103 over state–action trajectories,  $p_{\tilde{\tau}_{0:t} \mid \xi_{1:t+1}^*}(\cdot \mid \xi_{1:t+1}^*(\mathcal{S}^\Omega))$ . Unfortunately, the marginal likelihood  
104 in Equation (5) is intractable for all but the simplest probabilistic models.

105 To infer an approximate posterior distribution over state–action trajectories instead, we express  
106 posterior inference as the variational minimization problem

$$\min_{q_{\tilde{\tau}_{0:t}} \in \hat{\mathcal{Q}}} \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:t}}(\cdot) \parallel p_{\tilde{\tau}_{0:t} \mid \xi_{1:t+1}^*}(\cdot \mid \xi_{1:t+1}^*(\mathcal{S}^\Omega))), \quad (6)$$

107 where  $\mathbb{D}_{\text{KL}}(\cdot \parallel \cdot)$  is the KL divergence, and  $\hat{\mathcal{Q}}$  denotes the variational family over which to optimize.  
108 We consider a family of distributions parameterized by a policy  $\pi$  and defined by

$$q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t}) \doteq p_{\mathbf{S}_0}(\mathbf{s}_0) \pi(\mathbf{a}_t \mid \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) \pi(\mathbf{a}_{t'} \mid \mathbf{s}_{t'}), \quad (7)$$

109 where  $\pi \in \Pi$ , a family of policy distributions, and where  $p_{\mathbf{S}_0}(\mathbf{s}_0) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} \mid \mathbf{s}_{t'}, \mathbf{a}_{t'})$  is the  
110 true state transition distribution up to and including the state transition at  $t$ . In Proposition 1 (Fixed-  
111 Time Variational Objective), we show that under this variational family, the inference problem  
112 in Equation (6) can be equivalently stated as the problem of maximizing an entropy-regularized  
113 expected reward function at every time step, where the reward function is given by the log-likelihood  
114 of transitioning onto an optimal state–action trajectory given a state–action pair. This is effectively  
115 the result obtained by Ziebart et al. [13], Levine [3], and Haarnoja et al. [1].

### 116 3.2 Infinite-Horizon Reinforcement Learning as Variational Bayesian Inference

117 To derive an infinite-horizon objective, we modify the probabilistic model used above. To represent  
 118 the possibility that an agent may stay on the optimal state trajectory for *any* number of time steps,  
 119 that is, for state–action trajectories of varying lengths, we treat the length of the trajectory itself as a  
 120 random variable,  $T$ , and define the model

$$p(\tilde{\tau}_{0:t}, \xi_{1:t+1}^*(\mathcal{S}^\Omega), t) \doteq p_T(t) p_{\mathbf{s}_0}(\mathbf{s}_0) \prod_{t'=0}^t \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) p_d(\mathbf{s}_{t'+1} \mid \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} \mid \mathbf{s}_{t'}),$$

121 where  $p_T(t)$  is the probability of remaining on the optimal state trajectory for  $t + 1$  time steps. Since  
 122 the trajectory length is itself a random variable, the joint distribution is a *transdimensional* distribution  
 123 defined on  $\bigsqcup_{t=0}^\infty \{t\} \times \mathcal{S}^t \times \mathcal{A}^t$  [2].

124 Unlike in the fixed-horizon setting, the variational Bayesian inference problem in the infinite-horizon  
 125 setting corresponds to finding the posterior distribution over both state–action trajectories *and* the  
 126 length of the optimal state trajectory  $T$  conditioned on the desired behavior  $\xi_{1:t+1}^*(\mathcal{S}^\Omega)$ . Analogously  
 127 to the steps above, we can express this inference problem variationally as

$$\min_{q_{\tilde{\tau}_{0:T}, T} \in \mathcal{Q}} \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \parallel p_{\tilde{\tau}_{0:T}, T \mid \xi_{1:T+1}}(\cdot \mid \xi_{1:t+1}^*(\mathcal{S}^\Omega))), \quad (8)$$

128 where  $t$  denotes the time step immediately *before* the outcome is achieved,  $\mathcal{Q}$  denotes the variational  
 129 family. Under this variational distribution, we can obtain an unfactorized variational objective that  
 130 does in general not lend itself to stochastic gradient-based optimization (and off-policy reinforcement  
 131 learning). The variational objective is given in Proposition 3, but we omit it here for brevity.

132 To obtain a variational objective amenable to stochastic variational inference and off-policy reinforcement  
 133 learning, we define the variational family as follows:  $q_{\tilde{\tau}_{0:T}, T}(\tilde{\tau}_{0:t}, t) = q_{\tilde{\tau}_{0:t} \mid T}(\tilde{\tau}_{0:t} \mid t) q_T(t)$ ,  
 134 where  $q_T$  is a distribution over  $T$  in some variational family  $\mathcal{Q}_T$  parameterized by

$$q_T(t) = q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0), \quad (9)$$

135 with Bernoulli random variables  $\Delta_t$  denoting the event of “remaining on the optimal state trajectory  
 136 from time step 1 to time step  $t+1$ ,” we can equivalently express the variational problem in Equation (8)  
 137 recursively in a way that is tractable and amenable to off-policy optimization:

138 **Theorem 1** (Dynamic-Discount Behavior-Driven RL as Variational Inference). *Let  $q_T(t)$  and*  
 139  *$q_{\tilde{\tau}_{0:t} \mid T}(\tilde{\tau}_{0:t} \mid t, \mathbf{s}_0)$  be as defined in Equation (7) and Equation (9), and define a behavior-driven state*  
 140 *value function,*

$$V^\pi(\mathbf{s}_t, \mathcal{S}^\Omega; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t \mid \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T)] - \mathbb{D}_{\text{KL}}(\pi(\cdot \mid \mathbf{s}_t) \parallel p(\cdot \mid \mathbf{s}_t)), \quad (10)$$

141 *a behavior-driven state–action value function*

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) + q(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; \pi, q_T)], \quad (11)$$

142 *and a behavior-driven reward-like function*

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) \doteq \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 \mid \mathbf{s}_t, \mathbf{a}_t) - q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \parallel p_{\Delta_{t+1}}). \quad (12)$$

143 *Then given an optimal state trajectory  $\mathcal{S}^\Omega$ ,*

$$\mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \parallel p_{\tilde{\tau}_{0:T}, T \mid \xi_{1:T+1}}(\cdot \mid \xi_{1:T+1}^*(\mathcal{S}^\Omega))) = -\mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) + C = -V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T) + C,$$

144 *where  $C \doteq \log p(\xi_{1:T+1}^*)$  is independent of  $\pi$  and  $q_T$ , and hence maximizing  $V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; \pi, q_T)$  is*  
 145 *equivalent to minimizing Equation (8) and hence, the following holds:*

$$\arg \min_{\pi \in \Pi, q_T \in \mathcal{Q}_T} \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \parallel p_{\tilde{\tau}_{0:T}, T \mid \xi_{1:T+1}}(\cdot \mid \xi_{1:T+1}^*(\mathcal{S}^\Omega))) = \arg \max_{\pi \in \Pi, q_T \in \mathcal{Q}_T} V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T).$$

146 Theorem 1 tells us that the solution to the variational problem we started out with (Equation (8)),  
 147 is in fact the solution to an infinite-horizon reinforcement learning problem with a reward function  
 148 determined by the likelihood of transitioning onto an optimal trajectory, a learned, dynamic discount  
 149 factor, and KL divergence regularization. In Appendix A, we prove that dynamic-discount factor RL  
 150 is optimal and preferred over fixed discount factors. For detailed proofs, see the appendix.

## 151 4 Conclusion

152 Using a variational framing of the inference problem, we showed that optimized, dynamic discount  
 153 factors are optimal in KL-regularized RL and that fixed discount factor methods are a special (less  
 154 optimal) case of this formulation. We hope that this work contributes to bridging the gap between  
 155 reinforcement learning and probabilistic inference research and helps establish a mutual reference  
 156 point from which to derive new insights and methods.

157 **References**

- 158 [1] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-  
159 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*  
160 *Conference on Machine Learning*, 2018.
- 161 [2] Matthew Hoffman, Nando Freitas, Arnaud Doucet, and Jan Peters. An expectation maximiza-  
162 tion algorithm for continuous markov decision processes with arbitrary reward. In *Artificial*  
163 *intelligence and statistics*, pages 232–239, 2009.
- 164 [3] Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and  
165 review. 2018.
- 166 [4] Cong Lu, Philip J. Ball, Tim G. J. Rudner, Jack Parker-Holder, Michael A. Osborne, and  
167 Yee Whye Teh. Challenges and Opportunities in Offline Reinforcement Learning from Visual  
168 Observations. 2022.
- 169 [5] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
170 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. In *NIPS*  
171 *Workshop on Deep Learning*, pages 1–9, 2013.
- 172 [6] Rafael Rafailov, Tianhe Yu, Aravind Rajeswaran, and Chelsea Finn. Offline reinforcement  
173 learning from images with latent space models. In Ali Jadbabaie, John Lygeros, George J.  
174 Pappas, Pablo A. Parrilo, Benjamin Recht, Claire J. Tomlin, and Melanie N. Zeilinger, editors,  
175 *Proceedings of the 3rd Conference on Learning for Dynamics and Control*, volume 144 of  
176 *Proceedings of Machine Learning Research*, pages 1154–1168. PMLR, 07 – 08 June 2021.
- 177 [7] Tim G. J. Rudner, Cong Lu, Michael A. Osborne, Yarin Gal, and Yee Whye Teh. On pathologies  
178 in KL-regularized reinforcement learning from expert demonstrations. In *Advances in Neural*  
179 *Information Processing Systems 34*. 2021.
- 180 [8] Tim G. J. Rudner, Vitchyr H. Pong, Rowan McAllister, Yarin Gal, and Sergey Levine. Outcome-  
181 driven reinforcement learning via variational inference. In *Advances in Neural Information*  
182 *Processing Systems 34*. 2021.
- 183 [9] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driess-  
184 che, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander  
185 Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap,  
186 Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the  
187 game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, jan 2016.  
188 ISSN 0028-0836.
- 189 [10] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. 1998.
- 190 [11] Csaba Szepesvári. *Algorithms for Reinforcement Learning*, volume 4. 2010.
- 191 [12] Martha White. Unifying task specification in reinforcement learning. In *International Confer-*  
192 *ence on Machine Learning*, pages 3742–3750. PMLR, 2017.
- 193 [13] Brian D Ziebart, Andrew Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy  
194 inverse reinforcement learning. In *AAAI Conference on Artificial Intelligence*, pages 1433–1438,  
195 2008.

196 **Supplementary Material**

197 **Table of Contents**

198 **A Dynamic-Discout Behavior-Driven Reinforcement Learning** 7

199 **B Proofs** 8

200 B.1 Finite- and Infinite-Horizon Variational Objectives . . . . . 8

201 B.2 Recursive Variational Objective & Bellman Backup Operator . . . . . 10

202 B.3 Optimal Variational Posterior over  $T$  . . . . . 13

203 B.4 Dynamic-Discout Behavior-Driven Policy Iteration . . . . . 15

204 B.5 Lemmas . . . . . 17

## 205 Appendix A Dynamic-Discount Behavior-Driven Reinforcement Learning

206 Building on Theorem 1, we will now define a dynamic-discount behavior-driven Bellman backup operator and  
 207 use it to derive a policy iteration theorem for variational, dynamic-discount reinforcement learning. In particular,  
 208 we define:

209 **Definition 1** (Dynamic-Discount Behavior-Driven Bellman Backup Operator). *Given a function  $Q : \mathcal{S} \times \mathcal{A} \times$*   
 210  *$\mathcal{S} \rightarrow \mathbb{R}$ , define the operator  $\mathcal{T}^\pi$  as*

$$\mathcal{T}^\pi Q(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; q_T)], \quad (\text{A.1})$$

211 where  $r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta)$  is from Theorem 1 (Dynamic-Discount Behavior-Driven RL as Variational Inference)  
 212 and

$$V(\mathbf{s}_t, \mathcal{S}^\Omega; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T)] + \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)). \quad (\text{A.2})$$

213 This dynamic-discount, behavior-driven Bellman backup operator is identical to the Bellman backup operator for  
 214 KL-regularized reinforcement learning [7] except for the learned, dynamic discount factor,  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0)$ .

215 In tabular settings, repeated application of this Bellman operator will result in an optimal policy and an optimal  
 216 dynamic discount factor. More specifically, alternating between policy evaluation and optimization of the  
 217 variational distribution over the state–action trajectory and the trajectory length converges to an optimal policy.

218 **Theorem 2** (Variational Dynamic-Discount Behavior-Driven Policy Iteration). *Assume  $|\mathcal{A}| < \infty$  and that the*  
 219 *MDP is ergodic.*

220 1. *Dynamic-Discount Behavior-Driven Policy Evaluation (D2BD-PE): Given policy  $\pi$  and a function  $Q^0 :$*   
 221  *$\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ , define  $Q^{i+1} = \mathcal{T}^\pi Q^i$ . Then the sequence  $Q^i$  converges to the lower bound in Theorem 1.*

222 2. *Dynamic-Discount Behavior-Driven Policy Improvement (D2BD-PI): The policy*

$$\pi^+ = \arg \max_{\pi' \in \Pi} \left\{ \mathbb{E}_{\pi'(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T)] - \mathbb{D}_{\text{KL}}(\pi'(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \right\} \quad (\text{A.3})$$

223 and the variational distribution over  $T$  recursively defined in terms of

$$\begin{aligned} q^+(\Delta_{t+1} = 0 | \mathbf{s}_0; \pi, Q^\pi) \\ = \sigma \left( \mathbb{E}_{\pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1}) p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)] + \sigma^{-1} (p_{\Delta_{t+1}}(\Delta_{t+1} = 0)) \right) \end{aligned} \quad (\text{A.4})$$

224 improve the variational objective. In other words,  $V^{\pi^+}(\mathbf{s}_0, \mathcal{S}^\Omega; q_T) \geq V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T)$  and  
 225  $V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T^+) \geq V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T)$  for all  $\mathbf{s}_0 \in \mathcal{S}$ .

226 3. *Alternating between D2BD-PE and D2BD-PI converges to a policy  $\pi^*$  and a variational distribution over  $T$ ,*  
 227  *$q_T^*$ , such that  $Q^{\pi^*}(\mathbf{s}, \mathbf{a}, \mathcal{S}^\Omega; q_T^*) \geq Q^\pi(\mathbf{s}, \mathbf{a}, \mathcal{S}^\Omega; q_T)$  for all  $(\pi, q_T) \in \Pi \times \mathcal{Q}_T$  and any  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ .*

228 An implication of this result is that an optimal policy found via dynamic-discount behavior-driven policy iteration  
 229 has at least as high a state value at  $\mathbf{S}_0 = \mathbf{s}_0$  as it would under a fixed discount factor. That is, for  $p_T$  given by a  
 230 fixed geometric distribution with parameter  $\gamma$ , the state–action value function simplifies to the standard Bellman  
 231 backup operator,

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; p_T) \doteq \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; \pi, p_T)], \quad (\text{A.5})$$

232 and

$$Q^\pi(\mathbf{s}_0, \mathbf{a}_0, \mathcal{S}^\Omega; q_T^*) \geq Q^\pi(\mathbf{s}_0, \mathbf{a}_0, \mathcal{S}^\Omega; p_T). \quad (\text{A.6})$$

233 In other words, dynamic discount factors are optimal in KL-regularized reinforcement learning and can be  
 234 justified using the variational Bayesian inference formulation described here.

## 235 Appendix B Proofs

### 236 B.1 Finite- and Infinite-Horizon Variational Objectives

237 In this section, we present detailed derivations and proofs for the results in the main text.

238 **Proposition 1** (Fixed-Time Variational Objective). *Let the variational distribution  $q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t})$  be as defined*  
 239 *in Equation (7). Then, given a horizon length  $t^*$  and optimal state trajectory  $\mathcal{S}^\Omega$ ,*

$$\mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:t}}(\cdot) \parallel p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*)) = \log p(\Delta_{1:t+1}^*) - \bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega), \quad (\text{B.7})$$

240 where

$$\bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega) \doteq \mathbb{E}_{q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t})} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \parallel p(\cdot | \mathbf{s}_{t'})) \right], \quad (\text{B.8})$$

241 and since  $\log p(\Delta_{1:t+1}^*)$  is constant in  $\pi$ ,

$$\arg \min_{\pi \in \Pi} \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:t}}(\cdot) \parallel p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*)) = \arg \max_{\pi \in \Pi} \bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega). \quad (\text{B.9})$$

242 *Proof.* To find an approximation to the posterior  $p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*)$ , we can use variational inference.

243 To do so, we consider the trajectory distribution under  $p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*)$ , which by Bayes' Theorem is  
 244 given by

$$p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*) = \frac{p_{\mathbf{S}_0}(\mathbf{s}_0) \prod_{t'=0}^t \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) p(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{p(\Delta_{1:t+1}^*)}. \quad (\text{B.10})$$

245 Inferring an approximation to the posterior distribution  $p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*)$  then becomes equivalent to

246 finding a variational distribution  $q_{\tilde{\tau}_{0:t}|\mathbf{S}_0}(\cdot | \mathbf{s}_0)$ , which induces a trajectory distribution  $q_{\tilde{\tau}_{0:t}}(\cdot)$  that minimizes  
 247 the KL divergence from  $q_{\tilde{\tau}_{0:t}}(\cdot)$  to  $p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*)$ :

$$\min_{q_{\tilde{\tau}_{0:t}} \in \hat{\mathcal{Q}}} \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:t}}(\cdot) \parallel p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*)). \quad (\text{B.11})$$

248 If we find a distribution  $q_{\tilde{\tau}_{0:t}}(\cdot)$  for which the resulting KL divergence is zero, then  $q_{\tilde{\tau}_{0:t}}(\cdot)$  is the exact  
 249 posterior. If the KL divergence is positive, then  $q_{\tilde{\tau}_{0:t}}(\cdot)$  is an approximate posterior. To solve the variational  
 250 problem in Equation (B.11), we can define a factorized variational family

$$q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t}) \doteq p_{\mathbf{S}_0}(\mathbf{s}_0) \pi(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} q_{\mathbf{S}_{t'+1}|\mathbf{S}_{t'}, \mathbf{A}_{t'}}(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}), \quad (\text{B.12})$$

251 where  $\mathbf{A}_{0:t}$  and  $\mathbf{S}_{0:t}$  are latent variables over which to infer an approximate posterior distribution. Returning to  
 252 the variational problem in Equation (B.11), we can now write

$$\begin{aligned} \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:t}}(\cdot) \parallel p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\cdot | \Delta_{1:t+1}^*)) &= \int_{\mathcal{A}^{t+1}} \int_{\mathcal{S}^{t+1}} q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t}) \log \frac{q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t})}{p_{\tilde{\tau}_{0:t}|\Delta_{1:t+1}}(\tilde{\tau}_{0:t} | \Delta_{1:t+1}^*)} d\mathbf{s}_{0:t} d\mathbf{a}_{0:t} \\ &= -\bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega) + \log p(\Delta_{1:t+1}^*), \end{aligned} \quad (\text{B.13})$$

253 where

$$\begin{aligned} \bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega) &\doteq \mathbb{E}_{q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t})} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \right. \\ &\quad \left. + \log p(\mathbf{a}_t | \mathbf{s}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t) + \sum_{t'=0}^{t-1} \log p(\mathbf{a}_{t'} | \mathbf{s}_{t'}) - \log \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \right. \\ &\quad \left. + \sum_{t'=0}^{t-1} \log p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \log q_{\mathbf{S}_{t'+1}|\mathbf{S}_{t'}, \mathbf{A}_{t'}}(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \end{aligned} \quad (\text{B.14})$$

254 and

$$\log p(\Delta_{1:t+1}^*) = \log \int_{\mathcal{A}^{t+1}} \int_{\mathcal{S}^{t+1}} \mathbb{P}(\Delta_{1:t+1} = 1 | \mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_t, \mathbf{a}_t) p_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t}) d\mathbf{s}_{0:t} d\mathbf{a}_{0:t} \quad (\text{B.15})$$

$$= \log \int_{\mathcal{A}^{t+1}} \int_{\mathcal{S}^{t+1}} \left( \prod_{t'=0}^t \int_{\mathcal{S}} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \mathbb{I}\{\mathbf{s}_{t'+1} \in \mathcal{S}^\Omega\} d\mathbf{s}_{t'+1} \right) p_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t}) d\mathbf{s}_{0:t} d\mathbf{a}_{0:t} \quad (\text{B.16})$$



255 is a log-marginal likelihood. Following Haarnoja et al. [1], we define the variational distribution over next states  
 256 as the true transition dynamics, that is,

$$q_{\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t}(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) = p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t), \quad (\text{B.17})$$

257 so that

$$q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t}) \doteq p_{\mathbf{s}_0}(\mathbf{s}_0)\pi(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'})\pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}). \quad (\text{B.18})$$

258 We can then simplify  $\bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega)$  to

$$\bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega) = \mathbb{E}_{q_{\tilde{\tau}_{0:t}}}(\tilde{\tau}_{0:t}) \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \| p(\cdot | \mathbf{s}_{t'})) \right]. \quad (\text{B.19})$$

259 Since  $\log p(\Delta_{1:t+1}^*)$  is constant in  $\pi$ , solving the variational optimization problem in Equation (B.11) is  
 260 equivalent to maximizing the variational objective with respect to  $\pi \in \Pi$ , where  $\Pi$  is a family of policy  
 261 distributions.  $\square$

262 **Corollary 1.** *The objective in Equation (B.19) corresponds to KL-regularized reinforcement learning with a*  
 263 *time-varying reward function given by*

$$r(\mathbf{s}_{t'}, \mathbf{a}_{t'}, \Delta_{t'+1}) \doteq \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}).$$

264 *Proof.* Let

$$r(\mathbf{s}_{t'}, \mathbf{a}_{t'}, \Delta_{t'+1}) \doteq \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}). \quad (\text{B.20})$$

265 and note that the objective

$$\bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega) = \mathbb{E}_{q_{\tilde{\tau}_{0:t}}}(\tilde{\tau}_{0:t}) \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \| p(\cdot | \mathbf{s}_{t'})) \right]. \quad (\text{B.21})$$

266 can equivalently written as

$$\bar{\mathcal{F}}(\pi, \mathcal{S}^\Omega) = \mathbb{E}_{q_{\tilde{\tau}_{0:t}}}(\tilde{\tau}_{0:t}) \left[ \sum_{t'=0}^t r(\mathbf{s}_{t'}, \mathbf{a}_{t'}, \Delta_{t'+1}) + \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \| p(\cdot | \mathbf{s}_{t'})) \right], \quad (\text{B.22})$$

267 which, as shown in Haarnoja et al. [1], can be written in the form of Equation (11).  $\square$

268 **Proposition 3** (Unfactorized Dynamic-Discount Behavior-Driven RL as Variational Inference). *Let*

$$q_{\tilde{\tau}_{0:T}, T}(\tilde{\tau}_{0:t}, t) = q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t} | t)q_T(t), \quad (\text{B.23})$$

269 *let*  $q_T(t)$  *be a variational distribution defined on*  $t \in \mathbb{N}_0$ , *and let*  $q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t} | t)$  *be as defined in Equation (7).*

270 *Then, given an optimal state trajectory*  $\mathcal{S}^\Omega$ , *we have that*

$$\mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \| p_{\tilde{\tau}_{0:T}, T|\Delta_{1:T+1}}(\cdot | \Delta_{1:T+1}^*)) = \log p(\Delta_{1:T+1}^*) - \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega), \quad (\text{B.24})$$

271 *where*

$$\mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) \doteq \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t} | t)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \| p_{\tilde{\tau}_{0:T}, T}(\cdot)) \right] \quad (\text{B.25})$$

272 *and*  $\log p(\Delta_{1:T+1}^*)$  *is constant in*  $\pi$  *and*  $q_T$ .

273 *Proof.* In general, solving the variational problem

$$\min_{q \in \mathcal{Q}} \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \| p_{\tilde{\tau}_{0:T}, T|\Delta_{1:T+1}}(\cdot | \Delta_{1:T+1}^*)) \quad (\text{B.26})$$

274 is challenging, but as in the fixed-time setting, we can take advantage of the fact that, by choosing a variational  
 275 family parameterized by

$$q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t} | t) \doteq \pi(\mathbf{a}_t | \mathbf{s}_t) \prod_{t'=0}^{t-1} p_d(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}), \quad (\text{B.27})$$

276 with  $\pi \in \Pi$ , we can follow the same steps as in the proof for Proposition 1 (Fixed-Time Variational Objective)  
 277 and show that given an optimal state trajectory  $\mathcal{S}^\Omega$ ,

$$\mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \| p_{\tilde{\tau}_{0:T}, T|\Delta_{1:T+1}}(\cdot | \Delta_{1:T+1}^*)) = \log p(\Delta_{1:T+1}^*) - \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega), \quad (\text{B.28})$$

278 where

$$\begin{aligned} & \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) \\ & \doteq \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}|T|\mathbf{s}_0}(\cdot | \mathbf{s}_0) \| p_{\tilde{\tau}_{0:T}|T|\mathbf{s}_0}(\cdot | \mathbf{s}_0)) \right], \end{aligned} \quad (\text{B.29})$$

279 where  $q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}, t) \doteq q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)q_T(t)$ , and hence, solving the variational problem in Equation (8) is  
280 equivalent to maximizing  $\mathcal{F}(\pi, q_T, \mathcal{S}^\Omega)$  with respect to  $\pi$  and  $q_T$ .  $\square$

## 281 B.2 Recursive Variational Objective & Bellman Backup Operator

282 **Proposition 4** (Factorized Dynamic-Discount Behavior-Driven RL as Variational Inference). *Let the variational*  
283 *distribution factorize as*

$$q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}, t) = q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)q_T(t), \quad (\text{B.30})$$

284 let

$$q_T(t) = q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) \quad (\text{B.31})$$

285 be a variational distribution defined on  $t \in \mathbb{N}_0$ , and let  $q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)$  be as defined in Equation (7). Then,  
286 given an optimal state trajectory  $\mathcal{S}^\Omega$ , Equation (B.25) can be rewritten as

$$\mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) = \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) \right) \left( r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \right) \right] \quad (\text{B.32})$$

287 where

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) \doteq \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}), \quad (\text{B.33})$$

288 *Proof.* Consider the variational objective  $\mathcal{F}(\pi, q_T, \mathcal{S}^\Omega)$  in Equation (B.25):

$$\begin{aligned} & \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) \\ & = \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}|T}(\cdot) \| p_{\tilde{\tau}_{0:T}|T}(\cdot)) \right] \end{aligned} \quad (\text{B.34})$$

$$= \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \log \frac{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)q_T(t)}{p_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)p_T(t)} d\tilde{\tau}_{0:t} \right] \quad (\text{B.35})$$

$$= \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \log \frac{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)}{p_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \right] - \sum_{t=0}^{\infty} q_T(t) \log \frac{q_T(t)}{q_T(t)}. \quad (\text{B.36})$$

289 Noting that  $\sum_{t=0}^{\infty} q_T(t) \log \frac{q_T(t)}{q_T(t)} = \mathbb{D}_{\text{KL}}(q_T \| p_T)$ , we can write

$$\begin{aligned} & \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) \\ & = \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \log \frac{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)}{p_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \right] - \mathbb{D}_{\text{KL}}(q_T \| p_T) \end{aligned} \quad (\text{B.37})$$

$$\begin{aligned} & = \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \\ & \quad - \sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \left[ \log \frac{q_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)}{p_{\tilde{\tau}_{0:T}|T}(\tilde{\tau}_{0:t}|t)} \right] - \mathbb{D}_{\text{KL}}(q_T \| p_T). \end{aligned} \quad (\text{B.38})$$

290 Further noting that for an infinite-horizon trajectory distribution

$$q_{\tilde{\tau}_{t'} | \mathbf{s}_{t'}}(\tilde{\tau}_{t'} | \mathbf{s}_{t'}) \doteq \prod_{t=t'}^{\infty} p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi(\mathbf{a}_t | \mathbf{s}_t), \quad (\text{B.39})$$

291 trajectory realization  $\tilde{\tau}_{t+1} \doteq \{\tau_{t'}\}_{t'=t+1}^{\infty}$ , and any joint probability density  $f(\mathbf{s}_t, \mathbf{a}_t)$ ,

$$\sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_{0:T} | T}(\tilde{\tau}_{0:t} | t)} \left[ f(\mathbf{s}_t, \mathbf{a}_t) \right] \quad (\text{B.40})$$

$$= \sum_{t=0}^{\infty} \left( \int q_{\tilde{\tau}_{T+1}}(\tilde{\tau}_{t+1}) \left( \int_{\mathcal{S}^t \times \mathcal{A}^t} q_{\tilde{\tau}_{0:t}}(\tilde{\tau}_{0:t}) q_T(t) f(\mathbf{s}_t, \mathbf{a}_t) d\tilde{\tau}_{0:t} \right) d\tilde{\tau}_{t+1} \right) \quad (\text{B.41})$$

$$= \sum_{t=0}^{\infty} \left( \mathbb{E}_{q_{\tilde{\tau}_{0:T} | T}(\tilde{\tau}_{0:t} | t)} \left[ q_T(t) f(\mathbf{s}_t, \mathbf{a}_t) \right] \cdot \underbrace{\left( \int q_{\tilde{\tau}_{T+1}}(\tilde{\tau}_{t+1}) d\tilde{\tau}_{t+1} \right)}_{=1} \right) \quad (\text{B.42})$$

$$= \sum_{t=0}^{\infty} \left( \left( \int_{\mathcal{S}^t \times \mathcal{A}^t} q(\tilde{\tau}_{0:t}) q_T(t) f(\mathbf{s}_t, \mathbf{a}_t) d\tilde{\tau}_{0:t} \right) \cdot \underbrace{\left( \int q_{\tilde{\tau}_{T+1}}(\tilde{\tau}_{t+1}) d\tilde{\tau}_{t+1} \right)}_{=1} \right) \quad (\text{B.42})$$

$$= \sum_{t=0}^{\infty} \int q_{\tilde{\tau}_0}(\tilde{\tau}_0) q_T(t) f(\mathbf{s}_t, \mathbf{a}_t) d\tilde{\tau}_0 \quad (\text{B.43})$$

$$= \int q_{\tilde{\tau}_0}(\tilde{\tau}_0) \sum_{t=0}^{\infty} q_T(t) f(\mathbf{s}_t, \mathbf{a}_t) d\tilde{\tau}_0, \quad (\text{B.44})$$

292 we can express Equation (B.38) in terms of the infinite-horizon state–action trajectory

$$q_{\tilde{\tau}_0}(\tilde{\tau}_0) \doteq \prod_{t=0}^{\infty} p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi(\mathbf{a}_t | \mathbf{s}_t) \quad (\text{B.45})$$

293 as

$$\begin{aligned} \mathcal{F}(\pi, q_T, \mathcal{S}^{\Omega}) &= \int q_{\tilde{\tau}_0}(\tilde{\tau}_0) \sum_{t=0}^{\infty} q_T(t) \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^{\Omega}) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) d\tilde{\tau} \\ &\quad - \sum_{t=0}^{\infty} q_T(t) \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T} | T}(\cdot | t) \| p_{\tilde{\tau}_{0:T} | T}(\cdot | t)) - \mathbb{D}_{\text{KL}}(q_T \| p_T) \end{aligned} \quad (\text{B.46})$$

$$\begin{aligned} &= \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t=0}^{\infty} q_T(t) \left( \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^{\Omega}) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \right. \right. \\ &\quad \left. \left. - \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T} | T}(\cdot | t) \| p_{\tilde{\tau}_{0:T} | T}(\cdot | t)) \right) \right] - \mathbb{D}_{\text{KL}}(q_T \| p_T). \end{aligned} \quad (\text{B.47})$$

294 Using Lemma 5 and the definition of  $q_T(t)$  in Equation (9), we can rewrite this objective as

$$\begin{aligned} &\mathcal{F}(\pi, q_T, \mathcal{S}^{\Omega}) \\ &= \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) \right) q_{\Delta_{t'}}(\Delta_{t'} = 1) \left( \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^{\Omega}) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \right. \right. \\ &\quad \left. \left. - \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T} | T}(\cdot | t) \| p_{\tilde{\tau}_{0:T} | T}(\cdot | t)) \right) \right] - \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) \right) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}) \\ &= \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q(\Delta_{t'} = 0) \right) \right. \\ &\quad \cdot \left( q(\Delta_{t+1} = 1) \left( \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^{\Omega}) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T} | T}(\cdot | t) \| p_{\tilde{\tau}_{0:T} | T}(\cdot | t)) \right) \right. \\ &\quad \left. \left. - \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}) \right) \right], \end{aligned} \quad (\text{B.49})$$

295 with

$$\begin{aligned} & \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \parallel p_{\Delta_{t+1}}) \\ &= q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \log \frac{q_{\Delta_{t+1}}(\Delta_{t+1} = 0)}{p_{\Delta_{t+1}}(\Delta_{t+1} = 0)} + (1 - q_{\Delta_{t+1}}(\Delta_{t+1} = 0)) \log \frac{1 - q_{\Delta_{t+1}}(\Delta_{t+1} = 0)}{1 - p_{\Delta_{t+1}}(\Delta_{t+1} = 0)}. \end{aligned} \quad (\text{B.50})$$

296 Next, to re-express  $\mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}|\mathcal{T}}(\cdot | t) \parallel p_{\tilde{\tau}_{0:T}|\mathcal{T}}(\cdot | t))$  as a sum over Kullback-Leibler divergences between  
297 distributions over single action random variables, we note that

$$\mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}|\mathcal{T}}(\cdot | t) \parallel p_{\tilde{\tau}_{0:T}|\mathcal{T}}(\cdot | t)) = \int_{\mathcal{S}^t \times \mathcal{A}^t} q_{\tilde{\tau}_{0:T}|\mathcal{T}}(\tilde{\tau}_{0:t} | t) \log \frac{q_{\tilde{\tau}_{0:T}|\mathcal{T}}(\tilde{\tau}_{0:t} | t)}{p_{\tilde{\tau}_{0:T}|\mathcal{T}}(\tilde{\tau}_{0:t} | t)} d\tilde{\tau}_{0:t} \quad (\text{B.51})$$

$$= \int_{\mathcal{S}^t \times \mathcal{A}^t} q_{\tilde{\tau}_{0:T}|\mathcal{T}}(\tilde{\tau}_{0:t} | t) \log \frac{\prod_{t'=1}^t \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{\prod_{t'=1}^t p(\mathbf{a}_{t'} | \mathbf{s}_{t'})} d\tilde{\tau}_{0:t} \quad (\text{B.52})$$

$$= \int_{\mathcal{S}^t \times \mathcal{A}^t} q_{\tilde{\tau}_{0:T}|\mathcal{T}}(\tilde{\tau}_{0:t} | t) \sum_{t'=0}^t \log \frac{\pi(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{p(\mathbf{a}_{t'} | \mathbf{s}_{t'})} d\tilde{\tau}_{0:t} \quad (\text{B.53})$$

$$= \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t'=0}^t \int_{\mathcal{A}} \pi(\mathbf{a}_{t'} | \mathbf{s}_{t'}) \log \frac{\pi(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{p(\mathbf{a}_{t'} | \mathbf{s}_{t'})} d\mathbf{a}_{t'} \right] \quad (\text{B.54})$$

$$= \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t'=0}^t \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \parallel p(\cdot | \mathbf{s}_{t'})) \right], \quad (\text{B.55})$$

298 where we have used the same marginalization trick as above to express the expression in terms of an infinite-  
299 horizon trajectory distribution, which allows us to express Equation (B.49) as

$$\begin{aligned} & \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) \\ &= \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q(\Delta_{t'} = 0) \right) \cdot \left( q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \left( \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \right. \right. \right. \\ & \quad \left. \left. - \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t'=0}^t \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \parallel p(\cdot | \mathbf{s}_{t'})) \right] - \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \parallel p_{\Delta_{t+1}}) \right) \right] \\ &= \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q(\Delta_{t'} = 0) \right) \cdot \left( q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \left( \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \right. \right. \right. \right. \\ & \quad \left. \left. - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \parallel p(\cdot | \mathbf{s}_{t'})) \right) \right) - \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \parallel p_{\Delta_{t+1}}) \right]. \end{aligned} \quad (\text{B.56})$$

300 Rearranging and dropping redundant expectation operators, we can now express the objective as

$$\begin{aligned} & \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) \\ &= \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ - \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) \right) \left( q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \parallel p_{\Delta_{t+1}}) \right) \right] \\ & \quad + \underbrace{\sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \right)}_{=q_T(t)} \\ & \quad \cdot \mathbb{E}_{q_{\tilde{\tau}_0}(\tilde{\tau}_0)} \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \parallel p(\cdot | \mathbf{s}_{t'})) \right], \end{aligned} \quad (\text{B.57})$$

301 whereupon we note that the last term can be expressed as

$$\sum_{t=0}^{\infty} q_T(t) \mathbb{E}_{q_{\tilde{\tau}_0}}(\tilde{\tau}_0) \left[ \sum_{t'=0}^t \log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \| p(\cdot | \mathbf{s}_{t'})) \right] \quad (\text{B.58})$$

$$= \mathbb{E}_{q_{\tilde{\tau}_0}}(\tilde{\tau}_0) \left[ \sum_{t=0}^{\infty} \sum_{t'=0}^t q_T(t) (\log \mathbb{P}(\xi_{t'+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_{t'}, \mathbf{a}_{t'}) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t'}) \| p(\cdot | \mathbf{s}_{t'}))) \right]$$

$$= \mathbb{E}_{q_{\tilde{\tau}_0}}(\tilde{\tau}_0) \left[ \sum_{t=0}^{\infty} q(T \geq t) (\log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t))) \right] \quad (\text{B.59})$$

$$= \mathbb{E}_{q_{\tilde{\tau}_0}}(\tilde{\tau}_0) \left[ \sum_{t=0}^{\infty} \underbrace{\left( \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) \right)}_{(\text{by Lemma 2})} (\log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t))) \right], \quad (\text{B.60})$$

302 where the second line follows from expanding the sums and regrouping terms. By substituting the expression  
 303 in Equation (B.60) into Equation (B.57), we obtain an objective expressed entirely in terms of distributions over  
 304 single-index random variables:

$$\mathcal{F}(\pi, q_T, \mathcal{S}^\Omega)$$

$$= \mathbb{E}_{q_{\tilde{\tau}_0}}(\tilde{\tau}_0) \left[ \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) \right) \cdot \left( \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \right) \right] \quad (\text{B.61})$$

$$= \mathbb{E}_{q_{\tilde{\tau}_0}}(\tilde{\tau}_0) \left[ \sum_{t=0}^{\infty} \left( \prod_{t'=1}^t q_{\Delta_{t'}}(\Delta_{t'} = 0) \right) \left( r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \right) \right], \quad (\text{B.62})$$

305 where we defined

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) \doteq \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}), \quad (\text{B.63})$$

306 which concludes the proof.  $\square$

307 **Theorem 1** (Recursive Dynamic-Discount Behavior-Driven RL as Variational Inference). *Let  $q_T(t)$  and*  
 308  *$q_{\tilde{\tau}_{0:t}|T}(\tilde{\tau}_{0:t} | t)$  be as defined in Equation (7) and Equation (9), and define*

$$V^\pi(\mathbf{s}_t, \mathcal{S}^\Omega; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} \left[ Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \right] - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)), \quad (\text{B.64})$$

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) + q(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} \left[ V^\pi(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; \pi, q_T) \right], \quad (\text{B.65})$$

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) \doteq \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}). \quad (\text{B.66})$$

309 Then given an optimal state trajectory  $\mathcal{S}^\Omega$ ,

$$\mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \| p_{\tilde{\tau}_{0:T}, T | \Delta_{1:T+1}}(\mathcal{S}^\Omega)(\cdot | \Delta_{1:T+1}^*(\mathcal{S}^\Omega))) = -\mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) + C = -V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T) + C,$$

310 where  $C \doteq \log p(\Delta_{1:T+1}^*)$  is independent of  $\pi$  and  $q_T$ , and hence maximizing  $V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; \pi, q_T)$  is equivalent  
 311 to minimizing Equation (8). In other words,

$$\begin{aligned} & \arg \min_{\pi \in \Pi, q_T \in \mathcal{Q}_T} \{ \mathbb{D}_{\text{KL}}(q_{\tilde{\tau}_{0:T}, T}(\cdot) \| p_{\tilde{\tau}_{0:T}, T | \mathbf{s}_0, \mathbf{s}_{T^*}}(\cdot | \Delta_{1:T+1}^*)) \} \\ &= \arg \max_{\pi \in \Pi, q_T \in \mathcal{Q}_T} \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) = \arg \max_{\pi \in \Pi, q_T \in \mathcal{Q}_T} V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T). \end{aligned}$$

312 *Proof.* The proof follows directly from the proof of Theorem 1 in Rudner et al. [8].  $\square$

### 313 B.3 Optimal Variational Posterior over $T$

314 **Proposition 2** (Optimal Variational Distribution over  $T$ ). *The optimal variational distribution  $q_T^*$  with respect*  
 315 *to Equation (10) is defined recursively in terms of  $q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0) \forall t \in \mathbb{N}_0$  by*

$$q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0; \pi, Q^\pi) = \sigma(\Lambda(\mathbf{s}_t, \pi, q_T, Q^\pi) + \sigma^{-1}(p_{\Delta_{t+1}}(\Delta_{t+1} = 0))), \quad (\text{B.67})$$

316 where

$$\Lambda(\mathbf{s}_t, \pi, q_T, Q^\pi) \doteq \mathbb{E}_{\pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1}) p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T) - \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t)]$$

317 and  $\sigma(\cdot)$  is the sigmoid function, that is,  $\sigma(x) = \frac{1}{e^{-x} + 1}$  and  $\sigma^{-1}(x) = \log \frac{x}{1-x}$ .

318 *Proof.* Consider  $\mathcal{F}(\pi, q_T, \mathcal{S}^\Omega)$ :

$$\begin{aligned} \mathcal{F}(\pi, q_T, \mathbf{s}_t, \mathcal{S}^\Omega) &= \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T)] \\ &= \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E} [V(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; q_T)]]. \end{aligned} \quad (\text{B.68})$$

319 Since the variational objective  $\mathcal{F}(\pi, q_T, \mathbf{s}_t, \mathcal{S}^\Omega)$  can be expressed recursively as

$$V^\pi(\mathbf{s}_t, \mathcal{S}^\Omega; q_T) \doteq \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T)] - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)),$$

320 with

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \doteq r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; q_T)],$$

$$r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) \doteq \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}),$$

321 and since  $\mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}})$  is strictly convex in  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0)$ , we can find the globally optimal  
322 Bernoulli distribution parameters  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0)$  for all  $t \in \mathbb{N}_0$  recursively. That is, it is sufficient to solve  
323 the problem

$$q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0) \doteq \arg \max_{q_{\Delta_{t+1}}(\Delta_{t+1}=0)} \left\{ \mathcal{F}(\pi, q_T, \mathcal{S}^\Omega) \right\} = \arg \max_{q_{\Delta_{t+1}}(\Delta_{t+1}=0)} \left\{ \mathcal{F}(\pi, q_{\Delta_1}, \dots, q_{\Delta_{t+1}}, \dots, \mathbf{s}_0, \mathcal{S}^\Omega) \right\} \quad (\text{B.69})$$

324 for a fixed  $t + 1$ . To do so, we take the derivative of  $\mathcal{F}(\pi, q_{\Delta_1}, \dots, q_{\Delta_{t+1}}, \dots, \mathbf{s}_0, \mathcal{S}^\Omega)$ , which—defined  
325 recursively—is given by

$$\begin{aligned} &\mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T)] - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \\ &= \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} \left[ r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; q_T)] \right] - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \end{aligned} \quad (\text{B.70})$$

$$\begin{aligned} &= \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} \left[ \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}) \right. \\ &\quad \left. + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; q_T)] \right] - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \end{aligned} \quad (\text{B.71})$$

$$\begin{aligned} &= \mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} \left[ \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}) \right. \\ &\quad \left. + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; q_T)] \right] - \mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)), \end{aligned} \quad (\text{B.72})$$

326 with respect to  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0)$  and set it to zero, which yields

$$\begin{aligned} 0 &= -\mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)} \left[ \mathbb{E}_{\pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1}) p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)] \right] \\ &\quad + \log \frac{1 - q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0)}{1 - p_{\Delta_{t+1}}(\Delta_{t+1} = 0)} - \log \frac{q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0)}{p_{\Delta_{t+1}}(\Delta_{t+1} = 0)}. \end{aligned} \quad (\text{B.73})$$

327 Rearranging, we get

$$\frac{q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0)}{1 - q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0)} = \exp \left( \mathbb{E} [Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)] + \log \frac{p_{\Delta_{t+1}}(\Delta_{t+1} = 0)}{1 - p_{\Delta_{t+1}}(\Delta_{t+1} = 0)} \right), \quad (\text{B.74})$$

328 where the expectation is taken with respect to  $\pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1}) p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \pi(\mathbf{a}_t | \mathbf{s}_t)$  and the  $Q$ -function  
329 depends on  $q(\Delta_{t'})$  with  $t' > t$ , but not on  $q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0)$ . Solving for  $q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0)$ . Solving for  
330  $q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0)$ , we obtain

$$q_{\Delta_{t+1}}^*(\Delta_{t+1} = 0) = \frac{\exp(\mathbb{E}_{p_{\pi p_d} \pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)] + \log \frac{p_{\Delta_{t+1}}(\Delta_{t+1}=0)}{1 - p_{\Delta_{t+1}}(\Delta_{t+1}=0)})}{1 + \exp(\mathbb{E}_{p_{\pi p_d} \pi(\mathbf{a}_t | \mathbf{s}_t)} [Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)] + \log \frac{p_{\Delta_{t+1}}(\Delta_{t+1}=0)}{1 - p_{\Delta_{t+1}}(\Delta_{t+1}=0)})} \quad (\text{B.75})$$

$$= \sigma \left( \mathbb{E}_{p_{\pi p_d}} [Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)] + \sigma^{-1}(p_{\Delta_{t+1}}(\Delta_{t+1} = 0)) \right), \quad (\text{B.76})$$

331 where  $p_{\pi p_d} \doteq \pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1}) p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$ ,  $\sigma(\cdot)$  is the sigmoid function with  $\sigma(x) = \frac{1}{e^{-x} + 1}$  and  
332  $\sigma^{-1}(x) = \log \frac{x}{1-x}$ . This concludes the proof.  $\square$

333 **Remark 1.** As can be seen from Proposition 2 (Optimal Variational Distribution over  $T$ ), the optimal approxima-  
334 tion to the posterior over  $T$  trades off short-term rewards via  $\mathbb{E}_{\pi(\mathbf{a}_t | \mathbf{s}_t)}[r(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta)]$ , long-term rewards  
335 via  $\mathbb{E}_{\pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1})p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)}[Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)]$ , and the prior log-odds of not achieving the outcome  
336 at a given point in time conditioned on the outcome not having been achieved yet,  $\frac{p_{\Delta_{t+1}}(\Delta_{t+1}=0)}{1-p_{\Delta_{t+1}}(\Delta_{t+1}=0)}$ .

#### 337 B.4 Dynamic-Discount Behavior-Driven Policy Iteration

338 **Theorem 2** (Variational Dynamic-Discount Behavior-Driven Policy Iteration). Assume  $|\mathcal{A}| < \infty$  and that the  
339 MDP is ergodic.

340 1. *Dynamic-Discount Behavior-Driven Policy Evaluation (D2BD-PE):* Given policy  $\pi$  and a function  $Q^0 : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ , define  $Q^{i+1} = \mathcal{T}^\pi Q^i$ . Then the sequence  $Q^i$  converges to the lower bound in Theorem 1  
341 (Dynamic-Discount Behavior-Driven RL as Variational Inference).  
342

343 2. *Dynamic-Discount Behavior-Driven Policy Improvement (D2BD-PI):* The policy

$$\pi^+ = \arg \max_{\pi' \in \Pi} \left\{ \mathbb{E}_{\pi'(\mathbf{a}_t | \mathbf{s}_t)} \left[ Q^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \right] - \mathbb{D}_{\text{KL}}(\pi'(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \right\} \quad (\text{B.77})$$

344 and the variational distribution over  $T$  recursively defined in terms of

$$\begin{aligned} q^+(\Delta_{t+1} = 0; \pi, Q^\pi) \\ = \sigma \left( \mathbb{E}_{\pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1})p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [Q^\pi(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)] + \sigma^{-1}(p_{\Delta_{t+1}}(\Delta_{t+1} = 0)) \right) \end{aligned} \quad (\text{B.78})$$

345 improve the variational objective. In other words, we have that  $V^{\pi^+}(\mathbf{s}_0, \mathcal{S}^\Omega; q_T) \geq V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T)$  and  
346  $V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T^+) \geq V^\pi(\mathbf{s}_0, \mathcal{S}^\Omega; q_T)$  for all  $\mathbf{s}_0 \in \mathcal{S}$ .

347 3. *Alternating between D2BD-PE and D2BD-PI converges to a policy  $\pi^*$  and a variational distribution over  $T$ ,  
348  $q_T^*$ , such that  $Q^{\pi^*}(\mathbf{s}, \mathbf{a}, \mathcal{S}^\Omega; q_T^*) \geq Q^\pi(\mathbf{s}, \mathbf{a}, \mathcal{S}^\Omega; q_T)$  for all  $(\pi, q_T) \in \Pi \times \mathcal{Q}_T$  and any  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ .*

349 *Proof.* Parts of this proof are adapted from the proof given in Haarnoja et al. [1], modified for the Bellman  
350 operator proposed in Definition 1.

351 1. *Dynamic-Discount Behavior-Driven Policy Evaluation (D2BD-PE):* Instead of absorbing the entropy term  
352 into the  $Q$ -function, we can define an entropy-augmented reward as

$$\begin{aligned} r^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) \doteq \log \mathbb{P}(\xi_{t+1}(\mathcal{S}^\Omega) = 1 | \mathbf{s}_t, \mathbf{a}_t) - \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} \| p_{\Delta_{t+1}}) \\ + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [\mathbb{D}_{\text{KL}}(\pi(\cdot | \mathbf{s}_{t+1}) \| p(\cdot | \mathbf{s}_{t+1}))]. \end{aligned} \quad (\text{B.79})$$

353 We can then write an update rule according to Definition 1 as

$$\tilde{Q}(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \leftarrow r^\pi(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_\Delta) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{\pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1})p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [\tilde{Q}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T)], \quad (\text{B.80})$$

354 where  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \leq 1$ . This update is similar to a Bellman update [10], but with a discount  
355 factor given by  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0)$ . In general, this discount factor  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0)$  can be computed  
356 dynamically based on the current state and action, such as in Equation (B.67). As discussed in White [12],  
357 this Bellman operator is still a contraction mapping so long as the Markov chain induced by the current  
358 policy is ergodic and there exists a state such that  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0) < 1$ . The first condition is true by  
359 assumption. The second condition is true since  $q_{\Delta_{t+1}}(\Delta_{t+1} = 0)$  is given by Equation (B.67), which  
360 is always strictly between 0 and 1. Therefore, we apply convergence results for policy evaluation with  
361 transition-dependent discount factors [12] to this contraction mapping, and the result immediately follows.

362 2. *Dynamic-Discount Behavior-Driven Policy Improvement (D2BD-PI):* Let  $\pi_{\text{old}} \in \Pi$  and let  $Q^{\pi_{\text{old}}}$  and  $V^{\pi_{\text{old}}}$   
363 be the behavior-driven state and state-action value functions from Definition 1, let  $q_T$  be some variational  
364 distribution over  $T$ , and let  $\pi_{\text{new}}$  be given by

$$\pi_{\text{new}}(\mathbf{a}_t | \mathbf{s}_t) = \arg \max_{\pi' \in \Pi} \left\{ \mathbb{E}_{\pi'(\mathbf{a}_t | \mathbf{s}_t)} \left[ Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \right] - \mathbb{D}_{\text{KL}}(\pi'(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \right\} \quad (\text{B.81})$$

$$= \arg \max_{\pi' \in \Pi} \mathcal{J}_{\pi_{\text{old}}}(\pi'(\mathbf{a}_t, \mathbf{s}_t), q_T). \quad (\text{B.82})$$

365 Then, it must be true that  $\mathcal{J}_{\pi_{\text{old}}}(\pi_{\text{old}}(\mathbf{a}_t | \mathbf{s}_t); q_T) \leq \mathcal{J}_{\pi_{\text{old}}}(\pi_{\text{new}}(\mathbf{a}_t | \mathbf{s}_t); q_T)$ , since one could set  
366  $\pi_{\text{new}} = \pi_{\text{old}} \in \Pi$ . Thus,

$$\begin{aligned} \mathbb{E}_{\pi_{\text{new}}(\mathbf{a}_t | \mathbf{s}_t)} \left[ Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \right] - \mathbb{D}_{\text{KL}}(\pi_{\text{new}}(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \\ \geq \mathbb{E}_{\pi_{\text{old}}(\mathbf{a}_t | \mathbf{s}_t)} \left[ Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \right] - \mathbb{D}_{\text{KL}}(\pi_{\text{old}}(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)), \end{aligned} \quad (\text{B.83})$$

367 and since

$$V^{\pi_{\text{old}}}(\mathbf{s}_t, \mathcal{S}^\Omega; q_T) = \mathbb{E}_{\pi_{\text{old}}(\mathbf{a}_t | \mathbf{s}_t)} \left[ Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \right] - \mathbb{D}_{\text{KL}}(\pi_{\text{old}}(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)), \quad (\text{B.84})$$

368 we get

$$\mathbb{E}_{\pi_{\text{new}}(\mathbf{a}_t | \mathbf{s}_t)} \left[ Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \right] - \mathbb{D}_{\text{KL}}(\pi_{\text{new}}(\cdot | \mathbf{s}_t) \| p(\cdot | \mathbf{s}_t)) \geq V^{\pi_{\text{old}}}(\mathbf{s}_t, \mathcal{S}^\Omega; q_T). \quad (\text{B.85})$$

369 We can now write the Bellman equation as

$$\begin{aligned} & Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \\ &= q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V^{\pi_{\text{old}}}(\mathbf{s}_{t+1}, \mathcal{S}^\Omega; q_T)] \\ &\leq q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \log p_d(\mathbf{g} | \mathbf{s}_t, \mathbf{a}_t) + q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \mathbb{E}_{p(\mathbf{s}_{t'} | \mathbf{s}_t, \mathbf{a}_t)} [\mathbb{E}_{\pi_{\text{new}}(\mathbf{a}_{t'} | \mathbf{s}_{t'})} [Q^{\pi_{\text{old}}}(\mathbf{s}_{t'}, \mathbf{a}_{t'}, \mathcal{S}^\Omega; q_T)] \\ &\quad - \mathbb{D}_{\text{KL}}(\pi_{\text{new}}(\cdot | \mathbf{s}_{t'}) \| p(\cdot | \mathbf{s}_{t'}))], \end{aligned} \quad (\text{B.86})$$

370

$$\leq Q^{\pi_{\text{new}}}(\mathbf{s}_t, \mathbf{a}_t, \mathcal{S}^\Omega; q_T) \quad (\text{B.88})$$

371 where we defined  $t' \doteq t + 1$ , repeatedly applied the Bellman backup operator defined in [Definition 1](#) and  
372 used the bound in [Equation \(B.85\)](#). Convergence follows from Dynamic-Discount Behavior-Driven Policy  
373 Evaluation above.

374 3. Locally Optimal Variational Dynamic-Discount Behavior-Driven Policy Iteration: Define  $\pi^i$  to be a policy  
375 at iteration  $i$ . By ODPI for a given  $q_T$ , the sequence of state-action value functions  $\{Q^{\pi^i}(q_T)\}_{i=1}^\infty$  is  
376 monotonically increasing in  $i$ . Since the reward is finite and the negative KL divergence is upper bounded by  
377 zero,  $Q^\pi(q_T)$  is upper bounded for  $\pi \in \Pi$  and the sequence  $\{\pi^i\}_{i=1}^\infty$  converges to some  $\pi^*$ . To see that  $\pi^*$   
378 is an optimal policy, note that it must be the case that  $\mathcal{J}_{\pi^*}(\pi^*(\mathbf{a}_t | \mathbf{s}_t); q_T) > \mathcal{J}_{\pi^*}(\pi(\mathbf{a}_t | \mathbf{s}_t); q_T)$  for any  
379  $\pi \in \Pi$  with  $\pi \neq \pi^*$ . By the argument used in ODPI above, it must be the case that the behavior-driven  
380 state-action value of the converged policy is higher than that of any other non-converged policy in  $\Pi$ , that is,  
381  $Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t; q_T) > Q^\pi(\mathbf{s}_t, \mathbf{a}_t; q_T)$  for all  $\pi \in \Pi$  and any  $q_T^i \in \mathcal{Q}_T$  and  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ . Therefore, given  
382  $q_T$ ,  $\pi^*$  must be optimal in  $\Pi$ , which concludes the proof.

383 4. Globally Optimal Variational Dynamic-Discount Behavior-Driven Policy Iteration: Let  $\pi^i$  be a policy and  
384 let  $q_T^i$  be variational distributions over  $T$  at iteration  $i$ . By Locally Optimal Variational Dynamic-Discount  
385 Behavior-Driven Policy Iteration, for a fixed  $q_T^i$  with  $q_T^i = q_T^j, \forall i, j \in \mathbb{N}_0$ , the sequence of  $\{(\pi^i, q_T^i)\}_{i=1}^\infty$   
386 increases the objective [Equation \(B.24\)](#) at each iteration and converges to a stationary point in  $\pi^i$ , where  
387  $Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t; q_T^i) > Q^\pi(\mathbf{s}_t, \mathbf{a}_t; q_T^i)$  for all  $\pi \in \Pi$  and any  $q_T^i \in \mathcal{Q}_T$  and  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ . Since the  
388 objective in [Equation \(B.24\)](#) is concave in  $q_T$ , it must be the case that for,  $q_T^{*i} \in \mathcal{Q}_T$ , the optimal variational  
distribution over  $T$  at iteration  $i$ , defined recursively by

$$\begin{aligned} & q^{*i}(\Delta_{t+1} = 0; \pi^i, Q^{\pi^i}) \\ &= \sigma \left( \mathbb{E}_{\pi(\mathbf{a}_{t+1} | \mathbf{s}_{t+1})} p_d(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) [Q^{\pi^i}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}, \mathcal{S}^\Omega; q_T(\pi^i, Q^{\pi^i}))] + \sigma^{-1}(p_{\Delta_{t+1}}(\Delta_{t+1} = 0)) \right), \end{aligned}$$

389 for  $t \in \mathbb{N}_0$ ,  $Q^\pi(\mathbf{s}_t, \mathbf{a}_t; q_T^i) > Q^{\pi^i}(\mathbf{s}_t, \mathbf{a}_t; q_T^i)$  for all  $\pi \in \Pi$  and any  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ . Note that  $q_T$  is  
390 defined implicitly in terms of  $\pi^i$  and  $Q^{\pi^i}$ , that is, the optimal variational distribution over  $T$  at iteration  
391  $i$  is defined as a function of the policy and  $Q$ -function at iteration  $i$ . Hence, it must then be true that for  
392  $Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t; q_T^*) > Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t; q_T)$  for all  $q_T^*(\pi^*, Q^{\pi^*}) \in \mathcal{Q}_T$  and for any  $\pi^* \in \Pi$  and  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ .  
393 In other words, for an optimal policy and corresponding  $Q$ -function, there exists an optimal variational  
394 distribution over  $T$  that maximizes the  $Q$ -function, given the optimal policy. Repeating locally optimal  
395 variational behavior-driven policy iteration under the new variational distribution  $q_T^*(\pi^*, Q^{\pi^*})$  will yield an  
396 optimal policy  $\pi^{**}$  and computing the corresponding optimal variational distribution,  $q_T^{**}(\pi^{**}, Q^{\pi^{**}})$  will  
397 further increase the variational objective such that for  $\pi^{**} \in \Pi$  and  $q_T^{**}(\pi^{**}, Q^{\pi^{**}}) \in \mathcal{Q}_T$ , we have that

$$Q^{\pi^{**}}(\mathbf{s}_t, \mathbf{a}_t; q_T^{**}) > Q^{\pi^{**}}(\mathbf{s}_t, \mathbf{a}_t; q_T^*) > Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t; q_T^*) > Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t; q_T) \quad (\text{B.89})$$

398 for any  $\pi^* \in \Pi$  and  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ . Hence, global optimal variational behavior-driven policy iteration  
399 increases the variational objective at every step. Since the objective is upper bounded (by virtue of the  
400 rewards being finite and the negative KL divergence being upper bounded by zero) and the sequence of  
401  $\{(\pi^i, q_T^i)\}_{i=1}^\infty$  increases the objective [Equation \(B.24\)](#) at each iteration, by the monotone convergence  
402 theorem, the objective value converges to a supremum and since the objective function is concave the



403 supremum is unique. Hence, since the supremum is unique and obtained via global optimal variational  
 404 outcome-driven policy iteration on  $(\pi, q_T) \in \Pi \times \mathcal{Q}_T$ , the sequence of  $\{(\pi^i, q_T^i)\}_{i=1}^\infty$  converges to a  
 405 unique stationary point  $(\pi^*, q_T^*) \in \Pi \times \mathcal{Q}_T$ , where  $Q^{\pi^*}(s_t, \mathbf{a}_t; q_T^*) > Q^\pi(s_t, \mathbf{a}_t; q_T^i)$  for all  $\pi \in \Pi$  and  
 406 any  $q_T^i \in \mathcal{Q}_T$  and  $(s, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$ .

407 □

408 **Corollary 3** (Optimality of Variational Outcome Driven Policy Iteration). *Variational Dynamic-*  
 409 *Discount Behavior-Driven Policy Iteration on  $(\pi, q_T) \in \Pi \times \mathcal{Q}_T$  results in an optimal policy at least as good*  
 410 *or better than any optimal policy attainable from policy iteration on  $\pi \in \Pi$  alone.*

411 **Remark 2.** *The convergence proof of ODPE assumes a transition-dependent discount factor [12], because the*  
 412 *variational distribution used in Equation (B.67) depends on the next state and action as well as on the desired*  
 413 *outcome.*

## 414 B.5 Lemmas

415 **Lemma 1.** *Let  $q(T = t) \doteq q(T = t | T \geq t) \prod_{i=1}^t q(T \neq i - 1 | T \geq i - 1)$  be a discrete probability*  
 416 *distribution with support  $\mathbb{N}_0$ . Then for any  $t \in \mathbb{N}_0$ , we have that*

$$q(T \geq t) = \sum_{i=t}^{\infty} q(T = i | T \geq i) \prod_{j=1}^i q(T \neq j - 1 | T \geq j - 1) = \prod_{i=1}^t q(T \neq i - 1 | T \geq i - 1). \quad (\text{B.90})$$

417 *Proof.* We proof the statement by induction on  $t$ .

418 Base case: For  $t = 0$ ,  $q(T \geq 0) = 1$  by definition of the empty product.

419 Inductive case: Note that  $q(T \leq t) = \prod_{i=1}^t q(T = i - 1 | T \geq i - 1)$ . Show that

$$q(T \geq t) = \prod_{i=1}^t q(T \neq i - 1 | T \geq i - 1) \implies q(T \geq t + 1) = \prod_{i=1}^{t+1} q(T \neq i - 1 | T \geq i - 1). \quad (\text{B.91})$$

420 Consider  $q(T \geq t + 1) = \sum_{i=t+1}^{\infty} q(T = i | T \geq i) \prod_{j=1}^i q(T \neq j - 1 | T \geq j - 1)$ . To proof the inductive  
 421 hypothesis, we need to show that the following equality is true:

$$\begin{aligned} & \sum_{i=t+1}^{\infty} q(T = i | T \geq i) \prod_{j=1}^i q(T \neq j - 1 | T \geq j - 1) = \prod_{i=1}^{t+1} q(T \neq i - 1 | T \geq i - 1) \quad (\text{B.92}) \\ \iff & \sum_{i=t}^{\infty} q(T = i | T \geq i) \prod_{j=1}^i q(T \neq j - 1 | T \geq j - 1) - q(T = t | T \geq t) \prod_{j=1}^t q(T \neq j - 1 | T \geq j - 1) \\ & = q(T \neq t | T \geq t) \prod_{i=1}^t q(T \neq i - 1 | T \geq i - 1). \end{aligned} \quad (\text{B.93})$$

422 By the inductive hypothesis,

$$q(T \geq t) = \sum_{i=t}^{\infty} q(T = i | T \geq i) \prod_{j=1}^i q(T \neq j - 1 | T \geq j - 1) = \prod_{i=1}^t q(T \neq i - 1 | T \geq i - 1), \quad (\text{B.94})$$

423 and so

$$\text{Equation (B.93)} \iff \prod_{j=1}^t q(T \neq j | T \geq j) - q(T \neq t + 1 | T \geq t + 1) \quad (\text{B.95})$$

$$\cdot \prod_{j=1}^t q(T = j | T \geq j) = q(T \neq t | T \geq t) \prod_{i=1}^t q(T \neq i - 1 | T \geq i - 1). \quad (\text{B.96})$$

424 Factoring out  $\prod_{i=1}^t q(T \neq i-1 | T \geq i-1)$ , we get

$$\Leftrightarrow \prod_{j=1}^t q(T \neq j-1 | T \geq j-1) \underbrace{(1 - q(T = t | T \geq t))}_{=q(T \neq t | T \geq t)} = q(T \neq t | T \geq t) \prod_{j=1}^t q(T = j-1 | T \geq j-1) \quad (\text{B.97})$$

$$\Leftrightarrow q(T \neq t | T \geq t) \prod_{j=1}^t q(T \neq j-1 | T \geq j-1) = q(T \neq t | T \geq t) \prod_{j=1}^t q(T \neq j-1 | T \geq j-1), \quad (\text{B.98})$$

425 which proves the inductive hypothesis.  $\square$

426 **Lemma 2.** Let  $q_T(t)$  and  $p_T(t)$  be discrete probability distributions with support  $\mathbb{N}_0$ , let  $\Delta_t$  be a Bernoulli  
427 random variable, with success defined as  $T = t+1$  given that  $T \geq t$ , and let  $q_{\Delta_t}$  be a discrete probability  
428 distribution over  $\Delta_t$  for  $t \in \mathbb{N} \setminus \{0\}$ , so that

$$\begin{aligned} q_{\Delta_{t+1}}(\Delta_{t+1} = 0) &\doteq q(T \neq t | T \geq t) \\ q_{\Delta_{t+1}}(\Delta_{t+1} = 1) &\doteq q(T = t | T \geq t). \end{aligned} \quad (\text{B.99})$$

429 Then we can write  $q(T = t) = q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \prod_{i=1}^t q_{\Delta_i}(\Delta_i = 0)$  for any  $t \in \mathbb{N}_0$  and have that

$$q(T \geq t) = \sum_{i=t}^{\infty} q_{\Delta_{i+1}}(\Delta_{i+1} = 1) \prod_{j=1}^i q_{\Delta_j}(\Delta_j = 0) = \prod_{i=1}^t q_{\Delta_i}(\Delta_i = 0). \quad (\text{B.100})$$

430 *Proof.* By Lemma 1, we have that for any  $t \in \mathbb{N}_0$

$$q(T \geq t) = \sum_{i=t}^{\infty} q(T = i | T \geq i) \prod_{j=1}^i q(T \neq j-1 | T \geq j-1) = \prod_{i=1}^t q(T \neq i-1 | T \geq i-1). \quad (\text{B.101})$$

431 The result follows by replacing  $q(T = i | T \geq i)$  by  $q_{\Delta_{i+1}}(\Delta_{i+1} = 1)$ ,  $q(T \neq j-1 | T \geq j-1)$  by  
432  $q_{\Delta_j}(\Delta_j = 0)$ , and  $q(T \neq i-1 | T \geq i-1)$  by  $q_{\Delta_i}(\Delta_i = 0)$ .  $\square$

433 **Lemma 3.** Let  $q_T(t)$  and  $p_T(t)$  be discrete probability distributions with support  $\mathbb{N}_0$ . Then for any  $k \in \mathbb{N}_0$ ,

$$\begin{aligned} &\mathbb{E}_{t \sim q(T | T \geq k)} \left[ \log \frac{q(T = t | T \geq k)}{p(T = t | T \geq k)} \right] \\ &= f(q, p, k) + q(T \neq k | T \geq k) \mathbb{E}_{t \sim q(T | T \geq k+1)} \left[ \log \frac{q(T = t | T \geq k+1)}{p(T = t | T \geq k+1)} \right]. \end{aligned} \quad (\text{B.102})$$

434 *Proof.* Consider  $\mathbb{E}_{t \sim q(T | T \geq k)} \left[ \log \frac{q(T=t | T \geq k)}{p(T=t | T \geq k)} \right]$  and note that by the law of total expectation we can rewrite  
435 it as

$$\begin{aligned} &\mathbb{E}_{t \sim q(T | T \geq k)} \left[ \log \frac{q(T = t | T \geq k)}{p(T = t | T \geq k)} \right] \\ &= q(T = k | T \geq k) \mathbb{E}_{t \sim q(T | T = k)} \left[ \log \frac{q(T = t | T \geq k)}{p(T = t | T \geq k)} \right] \\ &\quad + q(T \neq k | T \geq k) \mathbb{E}_{t \sim q(T | T \geq k+1)} \left[ \log \frac{q(T = t | T \geq k)}{p(T = t | T \geq k)} \right] \\ &= q(T = k | T \geq k) \log \frac{q(T = k | T \geq k)}{p(T = k | T \geq k)} + q(T \neq k | T \geq k) \mathbb{E}_{t \sim q(T | T \geq k+1)} \left[ \log \frac{q(T = t | T \geq k)}{p(T = t | T \geq k)} \right]. \end{aligned} \quad (\text{B.104})$$

436 For all values of  $T \geq k+1$ , we have that

$$q(T = t | T \geq k) = q(T = t | T \geq k+1) q(T \neq k | T \geq k) \quad (\text{B.105})$$

$$p(T = t | T \geq k) = p(T = t | T \geq k+1) p(T \neq k | T \geq k) \quad (\text{B.106})$$

437 and so we can rewrite the expectation in Equation (B.104) as

$$\mathbb{E}_{t \sim q(T | T \geq k+1)} \left[ \log \frac{q(T = t | T \geq k)}{p(T = t | T \geq k)} \right] = \mathbb{E}_{t \sim q(T | T \geq k+1)} \left[ \log \frac{q(T = t | T \geq k)}{p(T = t | T \geq k)} + \log \frac{q(T \neq k | T \geq k)}{p(T \neq k | T \geq k)} \right] \quad (\text{B.107})$$

$$= \mathbb{E}_{t \sim q(T | T \geq k+1)} \left[ \log \frac{q(T = t | T \geq k)}{p(T = t | T \geq k)} \right] + \log \frac{q(T \neq k | T \geq k)}{p(T \neq k | T \geq k)} \quad (\text{B.108})$$

438 Combining Equation (B.108) with Equation (B.104), we have

$$\begin{aligned} & \mathbb{E}_{t \sim q(T|T \geq k)} \left[ \log \frac{q(T=t|T \geq k)}{p(T=t|T \geq k)} \right] \\ &= \underbrace{q(T=k|T \geq k) \log \frac{q(T=k|T \geq k)}{p(T=k|T \geq k)} + q(T \neq k|T \geq k) \log \frac{q(T \neq k|T \geq k)}{p(T \neq k|T \geq k)}}_{\doteq f(q,p,k)} \\ & \quad + q(T \neq k|T \geq k) \mathbb{E}_{t \sim q(T|T \geq k+1)} \left[ \log \frac{q(T=t|T \geq k+1)}{p(T=t|T \geq k+1)} \right], \end{aligned} \quad (\text{B.109})$$

439 which concludes the proof.  $\square$

440 **Lemma 4.** Let  $q_T(t)$  and  $p_T(t)$  be discrete probability distributions with support  $\mathbb{N}_0$ . Then the KL divergence  
441 from  $q_T$  to  $p_T$  can be written as

$$\mathbb{D}_{\text{KL}}(q_T \| p_T) = \sum_{t=0}^{\infty} q(T \geq t) f(q_T, p_T, t) \quad (\text{B.110})$$

442 where  $f(q_T, p_T, t)$  is shorthand for

$$f(q_T, p_T, t) = q(T=t|T \geq t) \log \frac{q(T=t|T \geq t)}{p(T=t|T \geq t)} + q(T \neq t|T \geq t) \log \frac{q(T \neq t|T \geq t)}{p(T \neq t|T \geq t)}. \quad (\text{B.111})$$

443 *Proof.* Note that  $q(T=k)$  denotes the probability that the distribution  $q$  assigns to the event  $T=k$  and  $q(T \geq$   
444  $m)$  denotes the tail probability, that is,  $q(T \geq m) = \sum_{t=m}^{\infty} q(T=t)$ . We will write  $q(T|T \geq m)$  to denote  
445 the conditional distribution of  $q$  given  $T \geq m$ , that is,  $q(T=k|T \geq m) = \mathbb{1}[k \geq m]q(T=k)/q(T \geq m)$ .  
446 We will use analogous notation for  $p$ .

447 By the definition of the KL divergence and using the fact that, since the support is lowerbounded by  $T=0$ ,  
448  $q(T=0) = q(T=0|T \geq 0)$ , we have

$$\mathbb{D}_{\text{KL}}(q_T \| p_T) = \mathbb{E}_{t \sim q(T)} \left[ \log \frac{q(T=t)}{p(T=t)} \right] = \mathbb{E}_{t \sim q(T|T \geq 0)} \left[ \log \frac{q(T=t|T \geq 0)}{p(T=t|T \geq 0)} \right]. \quad (\text{B.112})$$

449 Using Lemma 3 with  $k=0, 1, 2, 3, \dots$ , we can expand the above expression to get

$$\mathbb{D}_{\text{KL}}(q_T \| p_T) = f(q_T, p_T, 0) + q(T \neq 0|T \geq 0) \mathbb{E}_{t \sim q(T|T \geq 1)} \left[ \log \frac{q(T=t|T \geq 1)}{p(T=t|T \geq 1)} \right] \quad (\text{B.113})$$

$$\begin{aligned} &= f(q, p, 0) + q(T \neq 0|T \geq 1) f(q_T, p_T, 1) \\ & \quad + q(T \neq 0|T \geq 0) q(T \neq 1|T \geq 1) \mathbb{E}_{t \sim q(T|T \geq 2)} \left[ \log \frac{q(T=t|T \geq 2)}{p(T=t|T \geq 2)} \right] \end{aligned} \quad (\text{B.114})$$

$$\begin{aligned} &= \underbrace{1}_{=q(T \geq 0)} \cdot f(q, p, 0) \\ & \quad + \underbrace{q(T \neq 0|T \geq 0)}_{=q(T \geq 1)} f(q, p, 1) \\ & \quad + \underbrace{q(T \neq 0|T \geq 0) q(T \neq 1|T \geq 1)}_{=q(T \geq 2)} f(q_T, p_T, 2) \\ & \quad + \underbrace{q(T \neq 0|T \geq 0) q(T \neq 1|T \geq 1) q(T \neq 2|T \geq 2)}_{=q(T \geq 3)} \mathbb{E}_{t \sim q(T|T \geq 3)} \left[ \log \frac{q(T=t|T \geq 3)}{p(T=t|T \geq 3)} \right] \end{aligned} \quad (\text{B.115})$$

$$= \sum_{t=0}^{\infty} q(T \geq t) f(q_T, p_T, t), \quad (\text{B.116})$$

450 where  $f(q_T, p_T, t)$  is shorthand for

$$f(q_T, p_T, t) = q(T=t|T \geq t) \log \frac{q(T=t|T \geq t)}{p(T=t|T \geq t)} + q(T \neq t|T \geq t) \log \frac{q(T \neq t|T \geq t)}{p(T \neq t|T \geq t)}. \quad (\text{B.117})$$

451 and we used the fact that, by Lemma 1,

$$q(T \geq t) = \prod_{k=1}^t q(T \neq k-1|T \geq k-1). \quad (\text{B.118})$$

452 This completes the proof.  $\square$

453 **Lemma 5.** Let  $q_T(t)$  and  $p_T(t)$  be discrete probability distributions with support  $\mathbb{N}_0$ , let  $\Delta_t$  be a Bernoulli  
454 random variable, with success defined as  $T = t$  given that  $T \geq t$ , and let  $q_{\Delta_t}$  and  $p_{\Delta_t}$  be discrete probability  
455 distributions over  $\Delta_t$  for  $t \in \mathbb{N}_0 \setminus \{0\}$ , so that

$$q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \doteq q(T \neq t | T \geq t) \quad q_{\Delta_{t+1}}(\Delta_{t+1} = 1) \doteq q(T = t | T \geq t) \quad (\text{B.119})$$

$$p_{\Delta_{t+1}}(\Delta_{t+1} = 0) \doteq p(T \neq t | T \geq t) \quad p_{\Delta_{t+1}}(\Delta_{t+1} = 1) \doteq p(T = t | T \geq t). \quad (\text{B.120})$$

456 Then the KL divergence from  $q_T$  to  $p_T$  can be written as

$$\mathbb{D}_{\text{KL}}(q_T || p_T) = \sum_{t=0}^{\infty} \left( \prod_{k=1}^t q_{\Delta_k}(\Delta_k = 0) \right) \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} || p_{\Delta_{t+1}}) \quad (\text{B.121})$$

457 *Proof.* The result follows from Lemma 4, Equation (B.118), Equation (B.119), and the definition of  $f$ . In detail,  
458 from Lemma 1, and Equation (B.119) we have that

$$q(T \geq t) = \prod_{k=1}^t q(T \neq k - 1 | T \geq k - 1) = \prod_{k=1}^t q_{\Delta_k}(\Delta_k = 0). \quad (\text{B.122})$$

459 From the definition of  $f(q_T, p_T, t)$ , we have

$$f(q_T, p_T, t) = q(T = t | T \geq t) \log \frac{q(T = t | T \geq t)}{p(T = t | T \geq t)} + q(T \neq t | T \geq t) \log \frac{q(T \neq t | T \geq t)}{p(T \neq t | T \geq t)} \quad (\text{B.123})$$

$$= q_{\Delta_{t+1}}(\Delta_{t+1} = 0) \log \frac{q_{\Delta_{t+1}}(\Delta_{t+1} = 0)}{p_{\Delta_{t+1}}(\Delta_{t+1} = 0)} + q(\Delta_{t+1} = 1) \log \frac{q_{\Delta_{t+1}}(\Delta_{t+1} = 1)}{p_{\Delta_{t+1}}(\Delta_{t+1} = 1)} \quad (\text{B.124})$$

$$= \mathbb{D}_{\text{KL}}(q_{\Delta_{t+1}} || p_{\Delta_{t+1}}). \quad (\text{B.125})$$

460 Combining Equation (B.122), Equation (B.125), and Equation (B.110) completes the proof.  $\square$