

Reinforcement Learning in Sparse-Reward Environments with Hindsight Policy Gradients

Paulo Rauber^{1, 2, 3, 4}

Avinash Ummadisingu²

Filipe Mutz^{5, 6}

Jürgen Schmidhuber^{1, 2, 3, 7}

¹IDSIA. Lugano, Switzerland.

²USI. Lugano, Switzerland.

³SUPSI. Lugano, Switzerland.

⁴QMUL. London, United Kingdom.

⁵IFES. Serra, Brazil.

⁶UFES. Serra, Brazil.

⁷NNAISENSE. Lugano, Switzerland

Keywords: Reinforcement learning, policy gradients, policy search, multi-goal reinforcement learning, goal-conditional policy, hierarchical reinforcement learning, sparse rewards, hindsight, counterfactual, generalization, sample efficiency.

Abstract

A reinforcement learning agent that needs to pursue different goals across episodes requires a goal-conditional policy. In addition to their potential to generalize desirable behavior to unseen goals, such policies may also enable higher-level planning based on subgoals. In sparse-reward environments, the capacity to exploit information about the degree to which an arbitrary goal has been achieved while another goal was intended appears crucial to enable sample efficient learning. However, reinforcement learning agents have only recently been endowed with such capacity for hindsight. In this paper, we demonstrate how hindsight can be introduced to policy gradient methods, generalizing this idea to a broad class of successful algorithms. Our experiments on a diverse selection of sparse-reward environments show that hindsight leads to a remarkable increase in sample efficiency.

1 Introduction

In a traditional reinforcement learning setting, an agent interacts with an environment in a sequence of episodes, observing states and acting according to a policy that ideally maximizes expected cumulative reward. If an agent is required to pursue different *goals* across episodes, its *goal-conditional policy* may be represented by a probability distribution over actions for every combination of state and goal. This distinction between states and goals is particularly useful when the probability of a state transition given an action is independent of the goal pursued by the agent.

Learning such goal-conditional behavior has received significant attention in ma-

chine learning and robotics, especially because a goal-conditional policy may generalize desirable behavior to goals that were never encountered by the agent (Schmidhuber and Huber, 1990; Da Silva et al., 2012; Kupcsik et al., 2013; Deisenroth et al., 2014; Schaul et al., 2015; Zhu et al., 2017; Kober et al., 2012; Ghosh et al., 2018; Mankowitz et al., 2018; Pathak et al., 2018; Schmidhuber, 2019). Consequently, developing goal-based curricula to facilitate learning has also attracted considerable interest (Fabisch and Metzen, 2014; Florensa et al., 2017; Sukhbaatar et al., 2018; Srivastava et al., 2013; Schmidhuber, 2013). In hierarchical reinforcement learning, goal-conditional policies may enable agents to plan using subgoals, which abstracts the details involved in lower-level decisions (Oh et al., 2017; Vezhnevets et al., 2017; Kulkarni et al., 2016; Levy et al., 2019).

In a typical *sparse-reward environment*, an agent receives a non-zero reward only upon reaching a *goal state*. Besides being natural, this task formulation avoids the potentially difficult problem of *reward shaping*, which often biases the learning process towards suboptimal behavior (Ng et al., 1999). Unfortunately, sparse-reward environments remain particularly challenging for traditional reinforcement learning algorithms (Andrychowicz et al., 2017; Florensa et al., 2017). For example, consider an agent tasked with traveling between cities. In a sparse-reward formulation, if reaching a desired destination by chance is unlikely, a learning agent will rarely obtain reward signals. At the same time, it seems natural to expect that an agent will learn how to reach the cities it visited regardless of its desired destinations.

In this context, the capacity to exploit information about the degree to which an arbitrary goal has been achieved while another goal was intended is called *hindsight*.

This capacity was recently introduced by Andrychowicz et al. (2017) to off-policy reinforcement learning algorithms that rely on experience replay (Lin, 1992). In earlier work, Karkus et al. (2016) introduced hindsight to policy search based on Bayesian optimization (Metzen et al., 2015). This work was recently extended by Pinsler et al. (2019).

In this paper, we demonstrate how hindsight can be introduced to policy gradient methods (Williams, 1986, 1992; Sutton et al., 1999a), generalizing this idea to a successful class of reinforcement learning algorithms (Peters and Schaal, 2008; Duan et al., 2016).

In contrast to previous work on hindsight, our approach relies on *importance sampling* (Bishop, 2013). In reinforcement learning, importance sampling has been traditionally employed in order to efficiently reuse information obtained by earlier policies during learning (Precup et al., 2000; Peshkin and Shelton, 2002; Jie and Abbeel, 2010; Thomas et al., 2015; Munos et al., 2016). In comparison, our approach attempts to efficiently learn about different goals using information obtained by the current policy for a specific goal. This approach leads to multiple formulations of a *hindsight policy gradient* that relate to well-known policy gradient results.

In comparison to conventional (goal-conditional) policy gradient estimators, our proposed estimators lead to remarkable sample efficiency on a diverse selection of sparse-reward environments.

2 Preliminaries

We denote random variables by upper case letters and assignments to these variables by corresponding lower case letters. We let $\text{Val}(X)$ denote the set of valid assignments to a random variable X . We also omit the subscript that typically relates a probability function to random variables when there is no risk of ambiguity. For instance, we may use $p(x)$ to denote $p_X(x)$ and $p(y)$ to denote $p_Y(y)$.

Consider an agent that interacts with its environment in a sequence of episodes, each of which lasts for exactly T time steps. The agent receives a goal $g \in \text{Val}(G)$ at the beginning of each episode. At every time step t , the agent observes a state $s_t \in \text{Val}(S_t)$, receives a reward $r(s_t, g) \in \mathbb{R}$, and chooses an action $a_t \in \text{Val}(A_t)$. For simplicity of notation, suppose that $\text{Val}(G)$, $\text{Val}(S_t)$, and $\text{Val}(A_t)$ are finite for every t .

In our setting, a goal-conditional policy defines a probability distribution over actions for every combination of state and goal. The same policy is used to make decisions at every time step.

Let $\tau = s_1, a_1, s_2, a_2, \dots, s_{T-1}, a_{T-1}, s_T$ denote a trajectory. We assume that the probability $p(\tau \mid g, \theta)$ of trajectory τ given goal g and a policy parameterized by $\theta \in \text{Val}(\Theta)$ is given by

$$p(\tau \mid g, \theta) = p(s_1) \prod_{t=1}^{T-1} p(a_t \mid s_t, g, \theta) p(s_{t+1} \mid s_t, a_t). \quad (1)$$

In contrast to a Markov decision process, this formulation allows the probability of a state transition given an action to change across time steps within an episode. More importantly, it implicitly states that the probability of a state transition given an action is independent of the goal pursued by the agent, which we denote by $S_{t+1} \perp\!\!\!\perp G \mid S_t, A_t$.

For every τ, g , and θ , we also assume that $p(\tau | g, \theta)$ is non-zero and differentiable with respect to θ .

Assuming that $G \perp \Theta$, the expected return $\eta(\theta)$ of a policy parameterized by θ is given by

$$\eta(\theta) = \mathbb{E} \left[\sum_{t=1}^T r(S_t, G) \mid \theta \right] = \sum_g p(g) \sum_{\tau} p(\tau | g, \theta) \sum_{t=1}^T r(s_t, g). \quad (2)$$

The action-value function is given by $Q_t^\theta(s, a, g) = \mathbb{E} \left[\sum_{t'=t+1}^T r(S_{t'}, g) \mid S_t = s, A_t = a, g, \theta \right]$, the value function by $V_t^\theta(s, g) = \mathbb{E} [Q_t^\theta(s, A_t, g) \mid S_t = s, g, \theta]$, and the advantage function by $A_t^\theta(s, a, g) = Q_t^\theta(s, a, g) - V_t^\theta(s, g)$.

3 Goal-conditional policy gradients

This section presents results for goal-conditional policies that are analogous to well-known results for conventional policies (Peters and Schaal, 2008). They establish the foundation for the results presented in the next section. Additional proofs are included in Appendix A for completeness.

The objective of policy gradient methods is finding policy parameters that achieve maximum expected return. When combined with Monte Carlo techniques (Bishop, 2013), the following result allows pursuing this objective using gradient-based optimization.

Theorem 3.1 (Intermediary goal-conditional policy gradient). *The gradient $\nabla \eta(\theta)$ of the expected return with respect to θ is given by*

$$\nabla \eta(\theta) = \sum_g p(g) \sum_{\tau} p(\tau | g, \theta) \left[\sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \theta) \right] \left[\sum_{t=1}^T r(s_t, g) \right]. \quad (3)$$

Proof. The partial derivative $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ of the expected return $\eta(\boldsymbol{\theta})$ with respect to θ_j is given by

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} \frac{\partial}{\partial\theta_j} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^T r(s_t, g). \quad (4)$$

The *likelihood-ratio trick* allows rewriting the previous equation as

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \frac{\partial}{\partial\theta_j} \log p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^T r(s_t, g). \quad (5)$$

Note that

$$\log p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) = \log p(s_1) + \sum_{t=1}^{T-1} \log p(a_t | s_t, g, \boldsymbol{\theta}) + \sum_{t=1}^{T-1} \log p(s_{t+1} | s_t, a_t). \quad (6)$$

Therefore,

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \left[\sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \right] \left[\sum_{t=1}^T r(s_t, g) \right]. \quad (7)$$

□

More conveniently, the following result can be obtained by noting that an action is independent of any previous state given the current state, the goal, and the policy parameters (see App. A.2).

Theorem 3.2 (Goal-conditional policy gradient). *The gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(s_{t'}, g). \quad (8)$$

In order to reduce the variance of the gradient estimator, the following result allows employing a so-called *baseline* (see App. A.4).

Theorem 3.3 (Goal-conditional policy gradient, baseline formulation). *For every $t, \boldsymbol{\theta}$, and associated real-valued (baseline) function b_t^θ , the gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \left[\left[\sum_{t'=t+1}^T r(s_{t'}, g) \right] - b_t^\theta(s_t, g) \right]. \quad (9)$$

Appendix A.7 presents the constant baselines that minimize the (elementwise) variance of the corresponding estimator. However, such baselines are usually impractical to compute (or estimate), and the variance of the estimator may be reduced further by a baseline function that depends on state and goal. Although generally suboptimal, it is typical to let the baseline function b_t^θ approximate the value function V_t^θ (Greensmith et al., 2004).

The action-value function is related to the goal-conditional policy gradient by the following result (see App. A.5).

Lemma 3.1 (Goal-conditional policy gradient, action-value formulation). *The gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) Q_t^\theta(s_t, a_t, g). \quad (10)$$

Lastly, actor-critic methods may rely on the following result for goal-conditional policies (see App. A.6).

Theorem 3.4 (Goal-conditional policy gradient, advantage formulation). *The gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) A_t^{\boldsymbol{\theta}}(s_t, a_t, g). \quad (11)$$

4 Hindsight policy gradients

This section presents the novel ideas that introduce hindsight to policy gradient methods. Additional proofs can be found in Appendix B.

Importance sampling is a traditional technique used to obtain estimates related to a random variable $X \sim p$ using samples from an arbitrary positive distribution q . This technique relies on the following equalities:

$$\mathbb{E}_{p(X)} [f(X)] = \sum_x p(x) f(x) = \sum_x \frac{q(x)}{q(x)} p(x) f(x) = \mathbb{E}_{q(X)} \left[\frac{p(X)}{q(X)} f(X) \right]. \quad (12)$$

Suppose that the reward $r(s, g)$ is known for every combination of state s and goal g , as in previous work on hindsight (Andrychowicz et al., 2017; Karkus et al., 2016; Pinsler et al., 2019). In that case, it is possible to evaluate a trajectory obtained while trying to achieve an original goal g' for an alternative goal g . This information can be exploited using a central result based on importance sampling.

Theorem 4.1 (Every-decision hindsight policy gradient). *For an arbitrary (original) goal g' , the gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T \left[\prod_{k=1}^{T-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g). \quad (13)$$

Proof. Starting from Theorem 3.2, importance sampling allows rewriting the partial derivative $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ as

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} \frac{p(\boldsymbol{\tau} | g', \boldsymbol{\theta})}{p(\boldsymbol{\tau} | g, \boldsymbol{\theta})} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(s_{t'}, g). \quad (14)$$

Using Equation 1,

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \left[\prod_{k=1}^{T-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] \sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(s_{t'}, g). \quad (15)$$

□

In the formulation presented above, every reward is multiplied by the ratio between the likelihood of the corresponding trajectory under an alternative goal and the likelihood under the original goal (see Eq. 1). Intuitively, every reward should instead be multiplied by a *likelihood ratio* that only considers the corresponding trajectory up to the previous action. This intuition underlies the following important result, named after an analogous result for action-value functions by Precup et al. (2000).

Theorem 4.2 (Per-decision hindsight policy gradient). *For an arbitrary (original) goal g' , the gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g). \quad (16)$$

Proof. Starting from Eq. 15, the partial derivative $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ can be rewritten as

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{t=1}^{T-1} \sum_{t'=t+1}^T \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) r(s_{t'}, g). \quad (17)$$

If we split every trajectory into states and actions before and after t' , then $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ is given by

$$\sum_g p(g) \sum_{t=1}^{T-1} \sum_{t'=\tau+1}^T \sum_{s_{1:t'-1} a_{1:t'-1}} p(s_{1:t'-1}, a_{1:t'-1} | g', \boldsymbol{\theta}) \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) z, \quad (18)$$

where z is defined by

$$z = \sum_{s_{t':T}} \sum_{a_{t':T-1}} p(s_{t':T}, a_{t':T-1} | s_{1:t'-1}, a_{1:t'-1}, g', \boldsymbol{\theta}) \left[\prod_{k=t'}^{T-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g). \quad (19)$$

Using Lemma D.2 (see App. D.2) and canceling terms,

$$z = \sum_{s_{t':T}} \sum_{a_{t':T-1}} p(s_{t'} | s_{t'-1}, a_{t'-1}) \left[\prod_{k=t'}^{T-1} p(a_k | s_k, g, \boldsymbol{\theta}) p(s_{k+1} | s_k, a_k) \right] r(s_{t'}, g). \quad (20)$$

Using Lemma D.2 once again,

$$z = \sum_{s_{t':T}} \sum_{a_{t':T-1}} p(s_{t':T}, a_{t':T-1} | s_{1:t'-1}, a_{1:t'-1}, g, \boldsymbol{\theta}) r(s_{t'}, g). \quad (21)$$

Using the fact that $S_{t'} \perp\!\!\!\perp G | S_{1:t'-1}, A_{1:t'-1}, \boldsymbol{\Theta}$,

$$z = \sum_{s_{t'}} r(s_{t'}, g) p(s_{t'} | s_{1:t'-1}, a_{1:t'-1}, g, \boldsymbol{\theta}) = \sum_{s_{t'}} r(s_{t'}, g) p(s_{t'} | s_{1:t'-1}, a_{1:t'-1}, g', \boldsymbol{\theta}). \quad (22)$$

Substituting z into Expression 18 and returning to an expectation over trajectories,

$$\frac{\partial}{\partial\theta_j} \eta(\boldsymbol{\theta}) = \sum_{\tau} p(\tau | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=\tau+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g). \quad (23)$$

□

The following lemma allows introducing baselines to hindsight policy gradients (see App. B.4).

Lemma 4.1. For every $g', t, \boldsymbol{\theta}$, and associated real-valued (baseline) function b_t^θ ,

$$\sum_{\tau} p(\tau | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \left[\prod_{k=1}^t \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] b_t^\theta(s_t, g) = \mathbf{0}. \quad (24)$$

Appendix B.7 presents the constant baselines that minimize the (elementwise) variance of the corresponding gradient estimator. By analogy with the conventional practice, we suggest letting the baseline function b_t^θ approximate the value function V_t^θ instead.

The action-value function is related to the hindsight policy gradient by the following result (see App. B.5).

Lemma 4.2 (Hindsight policy gradient, action-value formulation). For an arbitrary goal g' , the gradient $\nabla \eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by

$$\nabla \eta(\boldsymbol{\theta}) = \sum_{\tau} p(\tau | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \left[\prod_{k=1}^t \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] Q_t^\theta(s_t, a_t, g). \quad (25)$$

Importantly, the choice of likelihood ratio in Lemma 4.1 is far from unique. However, besides leading to straightforward estimation, it also underlies the advantage formulation presented below.

Theorem 4.3 (Hindsight policy gradient, advantage formulation). For an arbitrary (original) goal g' , the gradient $\nabla \eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by

$$\nabla \eta(\boldsymbol{\theta}) = \sum_{\tau} p(\tau | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \left[\prod_{k=1}^t \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] A_t^\theta(s_t, a_t, g). \quad (26)$$

Fortunately, the following result allows approximating the advantage under a goal using a state transition collected while pursuing another goal (see App. D.4).

Theorem 4.4. *For every t and θ , the advantage function A_t^θ is given by*

$$A_t^\theta(s, a, g) = \mathbb{E} [r(S_{t+1}, g) + V_{t+1}^\theta(S_{t+1}, g) - V_t^\theta(s, g) \mid S_t = s, A_t = a]. \quad (27)$$

5 Hindsight gradient estimators

This section details gradient estimation based on the results presented in the previous section. The corresponding proofs can be found in Appendix C.

Consider a dataset (batch) $\mathcal{D} = \{(\tau^{(i)}, g^{(i)})\}_{i=1}^N$ where each trajectory $\tau^{(i)}$ is obtained using a policy parameterized by θ in an attempt to achieve a goal $g^{(i)}$ chosen by the environment.

The following result points to a straightforward estimator based on Theorem 4.2 (see App. C.1).

Theorem 5.1. *The per-decision hindsight policy gradient estimator, given by*

$$\frac{1}{N} \sum_{i=1}^N \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(A_t^{(i)} \mid S_t^{(i)}, G^{(i)} = g, \theta) \sum_{t'=t+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(A_k^{(i)} \mid S_k^{(i)}, G^{(i)} = g, \theta)}{p(A_k^{(i)} \mid S_k^{(i)}, G^{(i)}, \theta)} \right] r(S_{t'}^{(i)}, g), \quad (28)$$

is a consistent and unbiased estimator of the gradient $\nabla \eta(\theta)$ of the expected return.

In preliminary experiments, we found that this estimator leads to unstable learning progress, which is probably due to its potential high variance. The following result, inspired by weighted importance sampling (Bishop, 2013), represents our attempt to trade variance for bias (see App. C.2).

Theorem 5.2. *The weighted per-decision hindsight policy gradient estimator, given by*

$$\sum_{i=1}^N \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(A_t^{(i)} | S_t^{(i)}, G^{(i)} = g, \theta) \sum_{t'=t+1}^T \frac{\left[\prod_{k=1}^{t'-1} \frac{p(A_k^{(i)} | S_k^{(i)}, G^{(i)}=g, \theta)}{p(A_k^{(i)} | S_k^{(i)}, G^{(i)}, \theta)} \right] r(S_{t'}^{(i)}, g)}{\sum_{j=1}^N \left[\prod_{k=1}^{t'-1} \frac{p(A_k^{(j)} | S_k^{(j)}, G^{(j)}=g, \theta)}{p(A_k^{(j)} | S_k^{(j)}, G^{(j)}, \theta)} \right]}, \quad (29)$$

is a consistent estimator of the gradient $\nabla \eta(\theta)$ of the expected return.

In simple terms, the likelihood ratio for every combination of trajectory, (alternative) goal, and time step is normalized across trajectories by this estimator. In Appendix C.3, we present a result that enables the corresponding consistency-preserving *weighted baseline*.

Consider a set $\mathcal{G}^{(i)} = \{g \in \text{Val}(G) \mid \text{exists a } t \text{ such that } r(s_t^{(i)}, g) \neq 0\}$ composed of so-called *active goals* during the i -th episode. The feasibility of the proposed estimators relies on the fact that only active goals correspond to non-zero terms inside the expectation over goals in Expressions 28 and 29. In many natural sparse-reward environments, active goals will correspond directly to states visited during episodes (for instance, the cities visited while trying to reach other cities), which enables computing said expectation exactly when the goal distribution is known.

The proposed estimators have remarkable properties that differentiate them from previous (weighted) importance sampling estimators for off-policy learning. For instance, although a trajectory is often more likely under the original goal than under an alternative goal, in policies with strong optimal substructure, a high probability of a trajectory between the state a and the goal (state) c that goes through the state b may naturally allow for a high probability of the corresponding (sub)trajectory between the state a and the goal (state) b . In other cases, the (unnormalized) likelihood ratios may

become very small for some (alternative) goals after a few time steps across all trajectories. After normalization, in the worst case, this may even lead to equivalent ratios for such goals for a given time step across all trajectories. In any case, it is important to note that only likelihood ratios associated to active goals for a given episode will affect the gradient estimate. Additionally, an original goal will always have (unnormalized) likelihood ratios equal to one for the corresponding episode.

Under mild additional assumptions, the proposed estimators also allow using a dataset containing goals chosen arbitrarily (instead of goals drawn from the goal distribution). Although this feature is not required by our experiments, we believe that it may be useful to circumvent *catastrophic forgetting* during curriculum learning (McCloskey and Cohen, 1989; Kirkpatrick et al., 2017).

6 Experiments

This section reports results of an empirical comparison between goal-conditional policy gradient estimators and hindsight policy gradient estimators.¹ Because there are no well-established sparse-reward environments intended to test agents under multiple goals, our experiments focus on our own selection of environments, which is described in Section 6.1. Section 6.2 details the implementation of estimators, policies, and baselines. Section 6.3 documents our experimental protocol. Section 6.4 analyses the results of the corresponding experiments. Unabridged results are presented in Section 6.5. Section 6.6 provides a supplementary empirical study of likelihood ratios, and Section 6.7

¹An open-source implementation of these estimators is available on <http://paulorauber.com/hpg>.

contains an empirical comparison with hindsight experience replay.

6.1 Environments

The environments presented in this section are diverse in terms of stochasticity, state space dimensionality and size, relationship between goals and states, and number of actions. In every one of these environments, the agent receives the remaining number of time steps plus one as a reward for reaching the goal state, which also ends the episode. In every other situation, the agent receives no reward.

Bit flipping environment. The agent starts every episode in the same state (0, represented by k bits), and its goal is to reach a randomly chosen state. The actions allow the agent to toggle (flip) each bit individually. The maximum number of time steps is $k + 1$. Despite its apparent simplicity, this environment is an ideal testbed for reinforcement learning algorithms intended to deal with sparse rewards, since obtaining a reward by chance is unlikely even for a relatively small k . Andrychowicz et al. (2017) employed a similar environment to evaluate their hindsight approach.

Grid world environments. The agent starts every episode in a (possibly random) position on an 11×11 grid, and its goal is to reach a randomly chosen (non-initial) position. Some of the positions on the grid may contain impassable obstacles (walls). The actions allow the agent to move in the four cardinal directions. Moving towards walls causes the agent to remain in its current position. A state or goal is represented by a pair of integers between 0 and 10. The maximum number of time steps is 32. In the *empty room* environment, the agent starts every episode in the upper left corner of the

grid, and there are no walls. In the *four rooms* environment (Sutton et al., 1999b), the agent starts every episode in one of the four corners of the grid (see Fig. 1). There are walls that partition the grid into four rooms, such that each room provides access to two other rooms through single openings (doors). With probability 0.2, the action chosen by the agent is ignored and replaced by a random action.

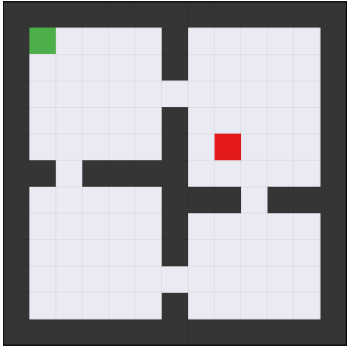


Figure 1: Four rooms.

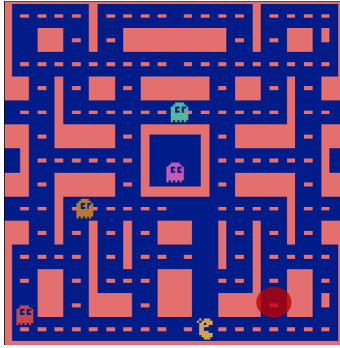


Figure 2: Ms. Pac-man.

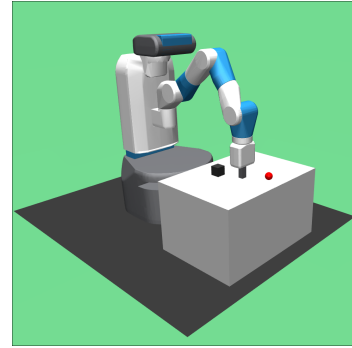


Figure 3: FetchPush.

***Ms. Pac-man* environment.** In this variant of the homonymous game for ATARI 2600 (see Fig. 2), the agent starts every episode close to the center of the map, and its goal is to reach a randomly chosen (non-initial) position on a 14×19 grid defined on the game screen. The actions allow the agent to move in the four cardinal directions for 13 game ticks. A state is represented by the result of preprocessing a sequence of game screens (images) as described in Section 6.2. A goal is represented by a pair of integers. The maximum number of time steps is 28, although an episode will also end if the agent is captured by an enemy. In comparison to the grid world environments considered in the previous section, this environment is additionally challenging due to its high-dimensional states and the presence of enemies.

FetchPush environment. This is a variant of the environment recently proposed by Plappert et al. (2018) to assess goal-conditional policy learning algorithms in a challenging task of practical interest (see Fig. 3). In a simulation, a robotic arm with seven degrees of freedom is required to push a randomly placed object (block) towards a randomly chosen position. The arm starts every episode in the same configuration. In contrast to the original environment, the actions in our variant allow increasing the desired velocity of the gripper along each of two orthogonal directions by ± 0.1 or ± 1 , leading to a total of eight actions. A state is represented by a 28-dimensional real vector that contains the following information: positions of the gripper and block; rotational and positional velocities of the gripper and block; relative position of the block with respect to the gripper; state of the gripper; and current desired velocity of the gripper along each direction. A goal is represented by three coordinates. The maximum number of time steps is 50.

6.2 Implementation

Importantly, the weighted per-decision hindsight policy gradient estimator used in our experiments (*HPG*) does not precisely correspond to Expression 29. Firstly, the original estimator requires a constant number of time steps T , which would often require the agent to act *beyond* the end of an episode in the environments that we consider. Secondly, although it is feasible to compute Expression 29 exactly when the goal distribution is known (as explained in Sec. 5), we sometimes subsample the sets of active goals per episode. Furthermore, when including a baseline that approximates the value function, we again consider only active goals, which by itself generally results in an in-

consistent estimator ($HPG+B$). As will become evident in the following sections, these *compromised* estimators still lead to remarkable sample efficiency.

In every experiment, a policy is represented by a feedforward neural network with a *softmax* output layer. The input to such a policy is a pair composed of state and goal. A baseline function is represented by a feedforward neural network with a single (linear) output neuron. The input to such a baseline function is a triple composed of state, goal, and time step. The baseline function is trained to approximate the value function using the mean squared (one-step) temporal difference error (Sutton and Barto, 1998). Parameters are updated using Adam (Kingma and Ba, 2014). The networks are given by the following.

Bit flipping environments and grid world environments. Both policy and baseline networks have two hidden layers, each with 256 hyperbolic tangent units. Every weight is initially drawn from a Gaussian distribution with mean 0 and standard deviation 0.01 (and redrawn if far from the mean by two standard deviations), and every bias is initially zero.

Ms. Pac-man environment. The policy network is represented by a convolutional neural network. The network architecture is given by a convolutional layer with 32 filters (8×8 , stride 4); convolutional layer with 64 filters (4×4 , stride 2); convolutional layer with 64 filters (3×3 , stride 1); and three fully-connected layers, each with 256 units. Every unit uses a hyperbolic tangent activation function. Every weight is initially set using variance scaling (Glorot and Bengio, 2010), and every bias is initially zero. These design decisions are similar to the ones made by Mnih et al. (2015).

A sequence of images obtained from the Arcade Learning Environment (Bellemare et al., 2013) is preprocessed as follows. Individually for each color channel, an elementwise maximum operation is employed between two consecutive images to reduce rendering artifacts. Such $210 \times 160 \times 3$ preprocessed image is converted to grayscale, cropped, and rescaled into an 84×84 image x_t . A sequence of images $x_{t-12}, x_{t-8}, x_{t-4}, x_t$ obtained in this way is *stacked* into an $84 \times 84 \times 4$ image, which is an input to the policy network (recall that each action is repeated for 13 game ticks). The goal information is concatenated with the *flattened* output of the last convolutional layer.

FetchPush environment. The policy network has three hidden layers, each with 256 hyperbolic tangent units. Every weight is initially set using variance scaling (Glorot and Bengio, 2010), and every bias is initially zero.

6.3 Evaluation

We assess sample efficiency through *learning curves* and *average performance* scores, which are obtained as follows. After collecting a number of batches (composed of trajectories and goals), each of which enables one step of gradient ascent, an agent undergoes *evaluation*. During evaluation, the agent interacts with the environment for a number of episodes, selecting actions with maximum probability according to its policy. A learning curve shows the average return obtained during each evaluation step, averaged across multiple *runs* (independent learning procedures). The curves presented in this text also include a 95% bootstrapped confidence interval. The average performance

is given by the average return across evaluation steps, averaged across runs. During both training and evaluation, goals are drawn uniformly at random. Note that there is no held-out set of goals for evaluation, since we are interested in evaluating sample efficiency instead of generalization.

For every combination of environment and batch size, grid search is used to select hyperparameters for each estimator according to average performance scores (after the corresponding standard deviation across runs is subtracted, as suggested by Duan et al. (2016)). *Definitive results* are obtained by using the best hyperparameters found for each estimator in additional runs. In most cases, we present definitive results for small (2) and medium (16) batch sizes.

Tables 1, 2, and 3 document the experimental settings. The number of runs, training batches, and batches between evaluations are reported separately for hyperparameter search and definitive runs. The number of training batches is adapted according to how soon each estimator leads to apparent convergence. Note that it is very difficult to establish this setting before hyperparameter search. The number of batches between evaluations is adapted so that there are 100 evaluation steps in total.

Other settings include the sets of policy and baseline learning rates under consideration for hyperparameter search, and the number of active goals subsampled per episode. In Tables 1, 2, and 3, $\mathcal{R}_1 = \{\alpha \times 10^{-k} \mid \alpha \in \{1, 5\} \text{ and } k \in \{2, 3, 4, 5\}\}$ and $\mathcal{R}_2 = \{\beta \times 10^{-5} \mid \beta \in \{1, 2.5, 5, 7.5, 10\}\}$.

As already mentioned, the definitive runs use the best combination of hyperparameters (learning rates) found for each estimator. Every setting was carefully chosen during preliminary experiments to ensure that the best result for each estimator is representa-

tive. In particular, the best performing learning rates rarely lie on the extrema of the corresponding search range. In the single case where the best performing learning rate found by hyperparameter search for a goal-conditional policy gradient estimator was such an extreme value (FetchPush, for a small batch size), evaluating one additional learning rate lead to decreased average performance.

Table 1: Experimental settings for the bit flipping environments

	Bit flipping (8 bits)			Bit flipping (16 bits)		
	Batch size 2	Batch size 16	Batch size 16	Batch size 2	Batch size 2	Batch size 16
Runs (definitive)	20	20	20	20	20	20
Training batches (definitive)	5000	1400	1400	15000	15000	1000
Batches between evaluations (definitive)	50	14	14	150	150	10
Runs (search)	10	10	10	10	10	10
Training batches (search)	4000	1400	1400	4000	4000	1000
Batches between evaluations (search)	40	14	14	40	40	10
<i>Policy learning rates</i>	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1
<i>Baseline learning rates</i>	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1
Episodes per evaluation	256	256	256	256	256	256
Maximum active goals per episode	∞	∞	∞	∞	∞	∞

Table 2: Experimental settings for the grid world environments

	Empty room			Four rooms		
	Batch size 2	Batch size 16	Batch size 16	Batch size 2	Batch size 2	Batch size 16
Runs (definitive)	20	20	20	20	20	20
Training batches (definitive)	2200	200	200	10000	10000	1700
Batches between evaluations (definitive)	22	2	2	100	100	17
Runs (search)	10	10	10	10	10	10
Training batches (search)	2500	800	800	10000	10000	3500
Batches between evaluations (search)	25	8	8	100	100	35
<i>Policy learning rates</i>	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1
<i>Baseline learning rates</i>	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1	\mathcal{R}_1
Episodes per evaluation	256	256	256	256	256	256
Maximum active goals per episode	∞	∞	∞	∞	∞	∞

Table 3: Experimental settings for the Ms. Pac-man and FetchPush environments

	Ms. Pac-man			FetchPush		
	Batch size 2	Batch size 16	Batch size 16	Batch size 2	Batch size 2	Batch size 16
Runs (definitive)	10	10	10	10	10	10
Training batches (definitive)	40000	12500	12500	40000	40000	12500
Batches between evaluations (definitive)	400	125	125	400	400	125
Runs (search)	5	5	5	5	5	5
Training batches (search)	40000	12000	12000	40000	40000	15000
Batches between evaluations (search)	800	120	120	800	800	300
<i>Policy learning rates</i>	\mathcal{R}_2	\mathcal{R}_2	\mathcal{R}_2	\mathcal{R}_2	\mathcal{R}_2	\mathcal{R}_2
Episodes per evaluation	240	240	240	512	512	512
Maximum active goals per episode	∞	3	3	∞	∞	3

6.4 Analysis

This section summarizes the unabridged results presented in Section 6.5.

Bit flipping environments. Figure 4 presents the learning curves for $k = 8$. Goal-conditional policy gradient estimators with and without an approximate value function baseline (*GCPG+B* and *GCPG*, respectively) obtain excellent policies and lead to comparable sample efficiency. *HPG+B* obtains excellent policies more than 400 batches earlier than these estimators, but its policies degrade upon additional training. Additional experiments strongly suggest that the main cause of this issue is the fact that the value function baseline is still very poorly fit by the time that the policy exhibits desirable behavior. In comparison, *HPG* obtains excellent policies as early as *HPG+B*, but its policies remain remarkably stable upon additional training.

The learning curves for $k = 16$ are presented in Figure 5. Clearly, both *GCPG* and *GCPG+B* are unable to obtain policies that perform better than chance, which is explained by the fact that they rarely incorporate reward signals during training. Confirming the importance of hindsight, *HPG* leads to stable and sample efficient learning. Although *HPG+B* also obtains excellent policies, they deteriorate upon additional training.

Similar results can be observed for a small batch size (see Sec. 6.5.3). The average performance results documented in Section 6.5.1 confirm that *HPG* leads to remarkable sample efficiency. Importantly, Sections 6.5.4 and 6.5.5 present hyperparameter sensitivity plots suggesting that *HPG* is less sensitive to hyperparameter settings than the other estimators. A hyperparameter sensitivity plot displays the average performance

achieved by each hyperparameter setting (sorted from best to worst along the horizontal axis). Section 6.5.5 also documents an ablation study where the likelihood ratios are removed from HPG, which notably promotes increased hyperparameter sensitivity. This study confirms the usefulness of the correction prescribed by importance sampling.

Grid world environments. Figure 6 shows the learning curves for the empty room environment. Clearly, every estimator obtains excellent policies, although HPG and HPG+B improve sample efficiency by at least 200 batches. The learning curves for the four rooms environment are presented in Figure 7. In this surprisingly challenging environment, every estimator obtains unsatisfactory policies. However, it is still clear that HPG and HPG+B improve sample efficiency. In contrast to the experiments presented in the previous section, HPG+B does not give rise to instability, which we attribute to easier value function estimation. Similar results can be observed for a small batch size (see Sec. 6.5.3). HPG achieves the best average performance in every grid world experiment except for a single case, where the best average performance is achieved by HPG+B (see Sec. 6.5.1). The hyperparameter sensitivity plots presented in Sections 6.5.4 and 6.5.5 once again suggest that HPG is less sensitive to hyperparameter choices, and that ignoring likelihood ratios promotes increased sensitivity (at least in the four rooms environment).

Ms. Pac-man environment. Figure 8 presents the learning curves for a medium batch size. Approximate value function baselines are excluded from this experiment due to the significant cost of systematic hyperparameter search. Although HPG obtains better policies during early training, GCPG obtains better final policies. However, for such

a medium batch size, only 3 active goals per episode (out of potentially 28) are subsampled for HPG. Although this harsh subsampling brings computational efficiency, it also appears to handicap the estimator. This hypothesis is supported by the fact that HPG outperforms GCPG for a small batch size, when all active goals are used (see Secs. 6.5.1 and 6.5.3). Policies obtained using each estimator are illustrated by videos included on the project website.

FetchPush environment. Figure 9 presents the learning curves for a medium batch size. HPG obtains good policies after a reasonable number of batches, in sharp contrast to GCPG. For such a medium batch size, only 3 active goals per episode (out of potentially 50) are subsampled for HPG, showing that subsampling is a viable alternative to reduce the computational cost of hindsight. Similar results are observed for a small batch size, when all active goals are used (see Secs. 6.5.1 and 6.5.3). Policies obtained using each estimator are illustrated by videos included on the project website.

6.5 Results

6.5.1 Average performance results

Table 4: Definitive average performance results

	Bit flipping (8 bits)		Bit flipping (16 bits)	
	Batch size 2	Batch size 16	Batch size 2	Batch size 16
	HPG	4.60 ± 0.06	4.72 ± 0.02	7.11 ± 0.12
GCPG	1.81 ± 0.61	3.44 ± 0.30	0.00 ± 0.00	0.00 ± 0.00
HPG+B	3.40 ± 0.46	4.04 ± 0.10	5.35 ± 0.40	6.09 ± 0.29
GCPG+B	0.64 ± 0.58	3.31 ± 0.58	0.00 ± 0.00	0.00 ± 0.00

	Empty room		Four rooms	
	Batch size 2	Batch size 16	Batch size 2	Batch size 16
	HPG	20.22 ± 0.37	16.83 ± 0.84	7.38 ± 0.16
GCPG	12.54 ± 1.01	10.96 ± 1.24	4.64 ± 0.57	6.12 ± 0.54
HPG+B	19.90 ± 0.29	17.12 ± 0.44	7.28 ± 1.28	8.08 ± 0.18
GCPG+B	12.69 ± 1.16	10.68 ± 1.36	4.26 ± 0.55	6.61 ± 0.49

	Ms. Pac-man		FetchPush	
	Batch size 2	Batch size 16	Batch size 2	Batch size 16
	HPG	6.58 ± 1.96	6.80 ± 0.64	6.10 ± 0.34
GCPG	5.29 ± 1.67	6.92 ± 0.58	3.48 ± 0.15	4.42 ± 0.28

6.5.2 Learning curves (batch size 16)

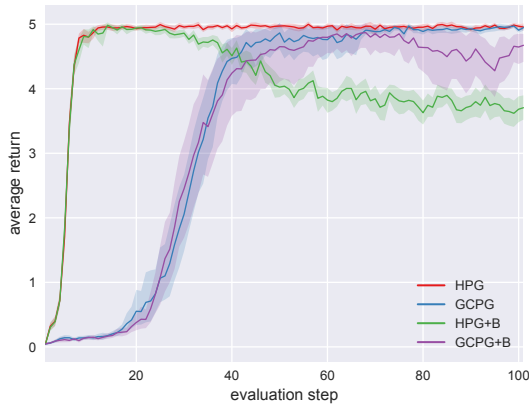


Figure 4: Bit flipping ($k = 8$).

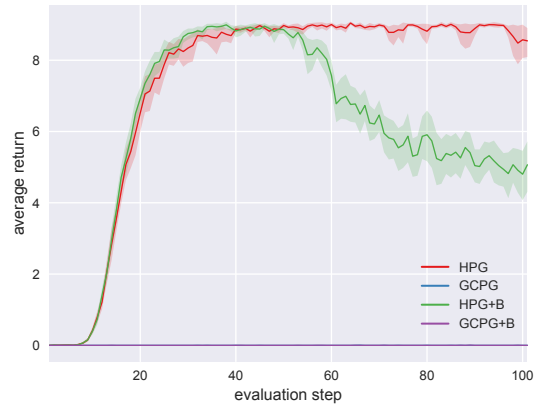


Figure 5: Bit flipping ($k = 16$).

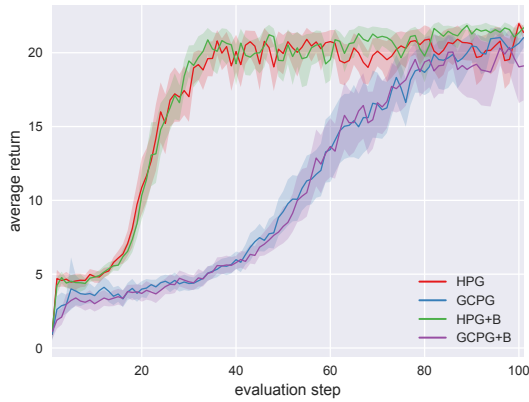


Figure 6: Empty room.

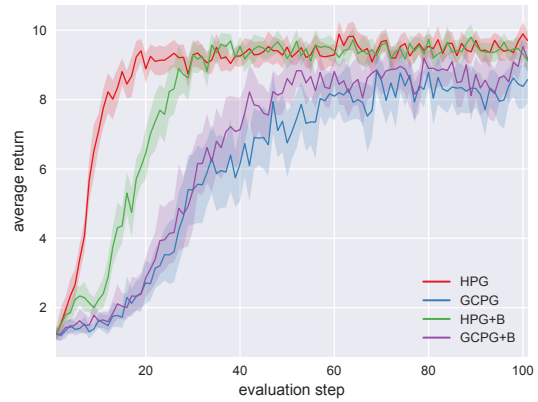


Figure 7: Four rooms.

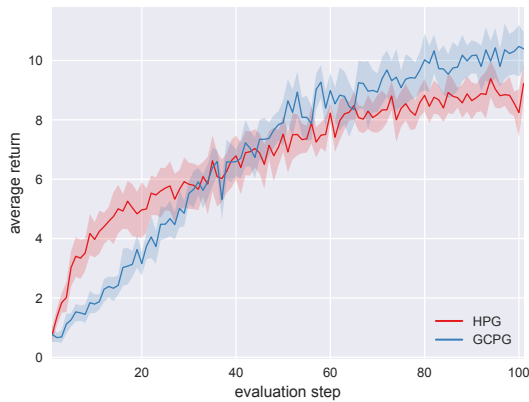


Figure 8: Ms. Pac-man.



Figure 9: FetchPush.

6.5.3 Learning curves (batch size 2)

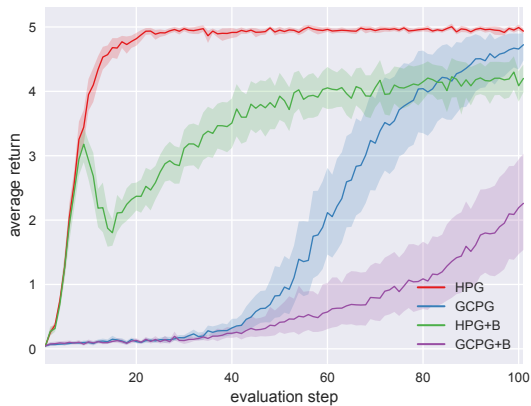


Figure 10: Bit flipping ($k = 8$).

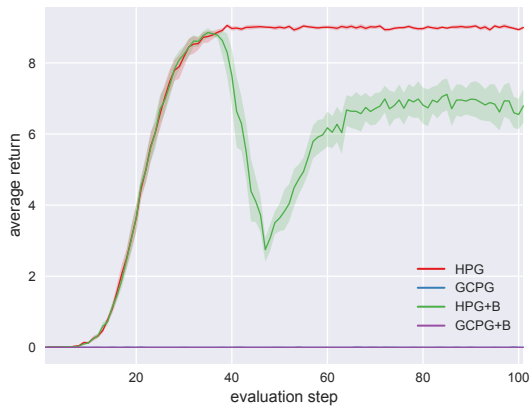


Figure 11: Bit flipping ($k = 16$).

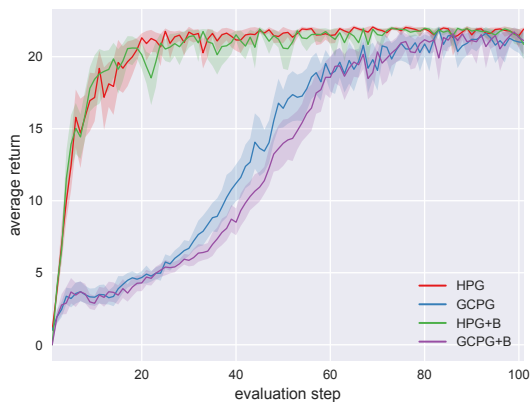


Figure 12: Empty room.

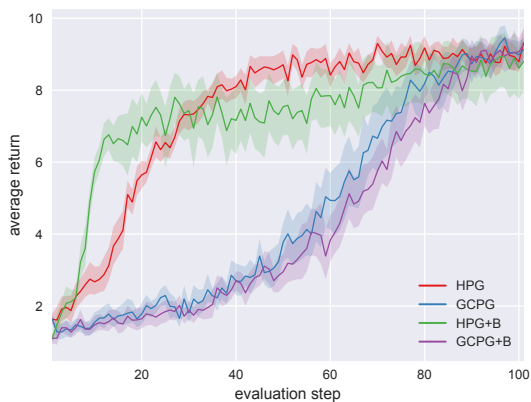


Figure 13: Four rooms.

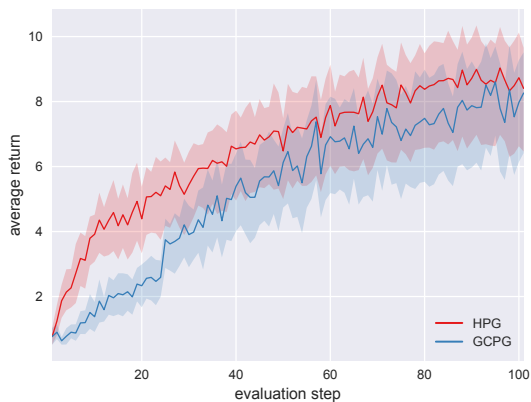


Figure 14: Ms. Pac-man.

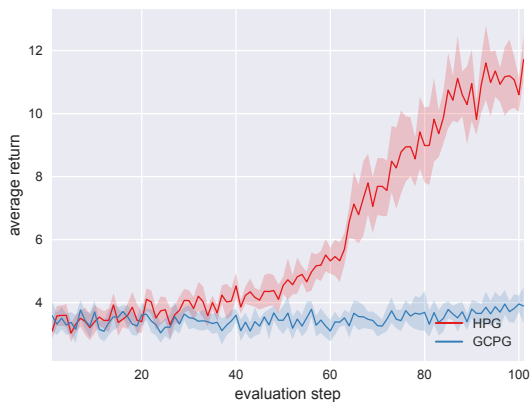


Figure 15: FetchPush.

6.5.4 Hyperparameter sensitivity plots (batch size 16)

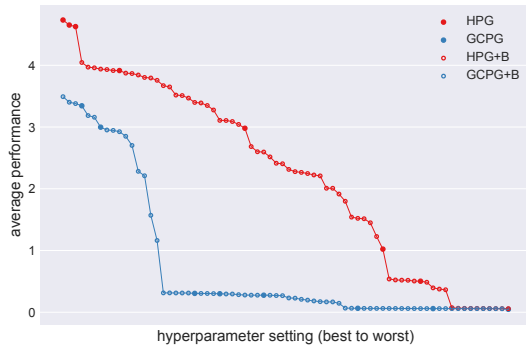


Figure 16: Bit flipping ($k = 8$).

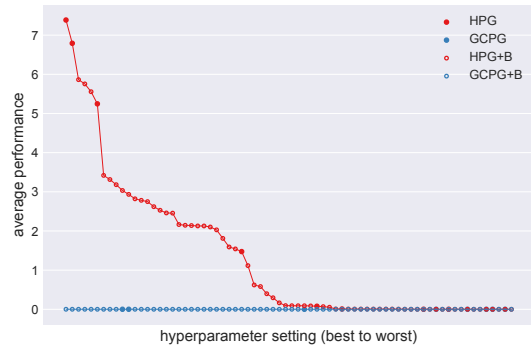


Figure 17: Bit flipping ($k = 16$).

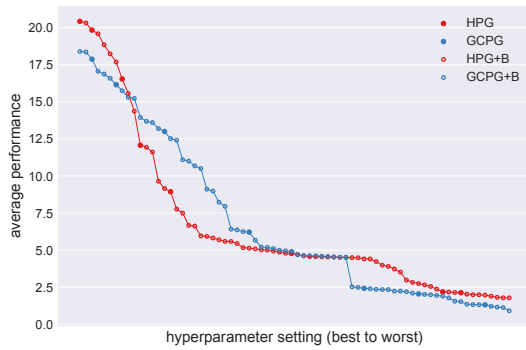


Figure 18: Empty room.



Figure 19: Four rooms.



Figure 20: Ms. Pac-man.

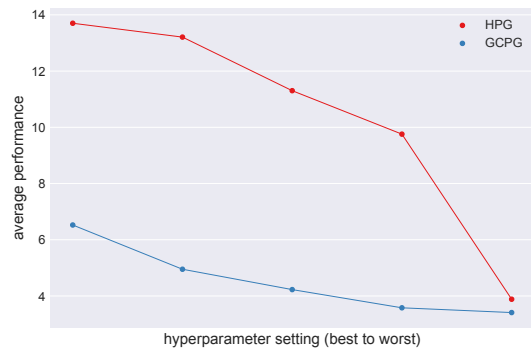


Figure 21: FetchPush.

6.5.5 Hyperparameter sensitivity plots (batch size 2)

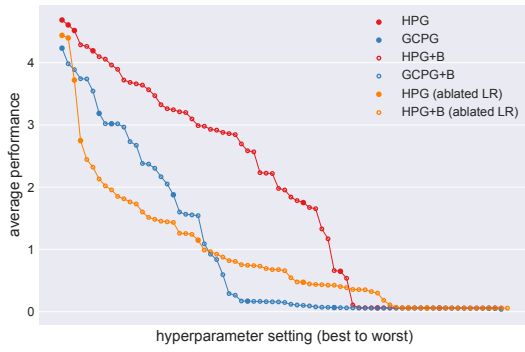


Figure 22: Bit flipping ($k = 8$).

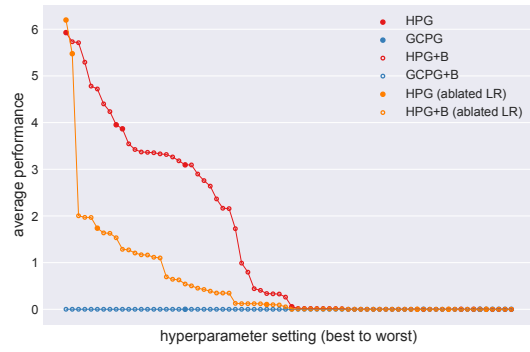


Figure 23: Bit flipping ($k = 16$).

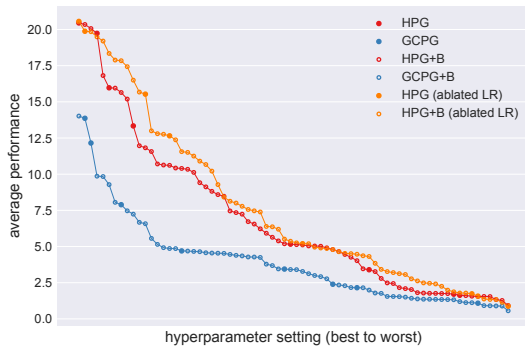


Figure 24: Empty room.

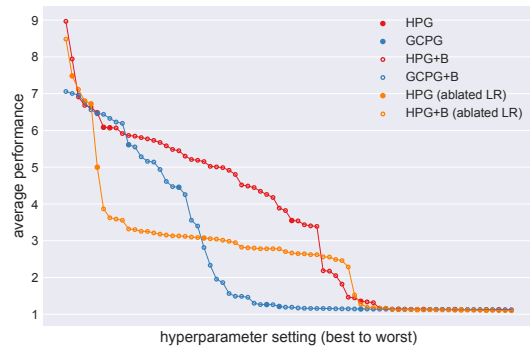


Figure 25: Four rooms.



Figure 26: Ms. Pac-man.



Figure 27: FetchPush.

6.6 Likelihood ratio study

This section presents a study of the *active* (normalized) likelihood ratios computed by agents during training. A likelihood ratio is considered active if and only if it multiplies a non-zero reward (see Expression 29). Note that only these likelihood ratios affect gradient estimates based on HPG.

This study is conveyed through plots that encode the distribution of active likelihood ratios computed during training, individually for each time step within an episode. Each plot corresponds to an agent that employs HPG and obtains the highest definitive average performance for a given environment (Figs. 28-33). Note that the length of the largest bar for a given time step is fixed to aid visualization.

The most important insight provided by these plots is that likelihood ratios behave very differently across environments, even for equivalent time steps (for instance, compare bit flipping environments to grid world environments). In contrast, after the first time step, the behavior of likelihood ratios changes slowly across time steps within the same environment. In any case, alternative goals have a significant effect on gradient estimates, which agrees with the results presented in the previous sections.

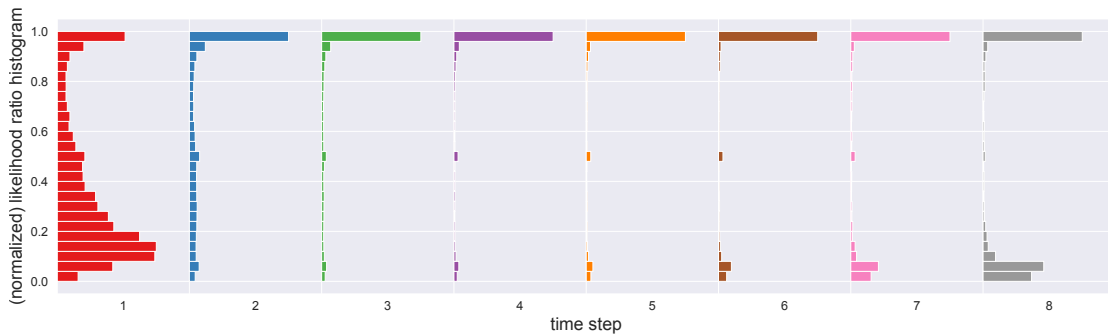


Figure 28: Bit flipping ($k = 8$, batch size 16).

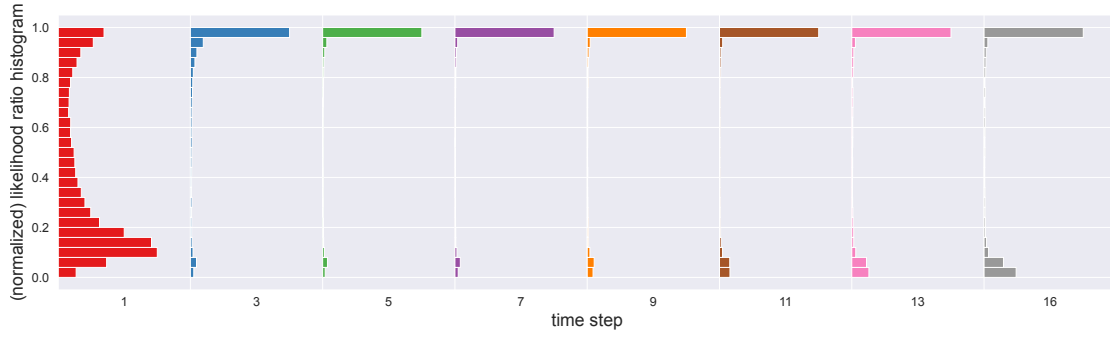


Figure 29: Bit flipping ($k = 16$, batch size 16).

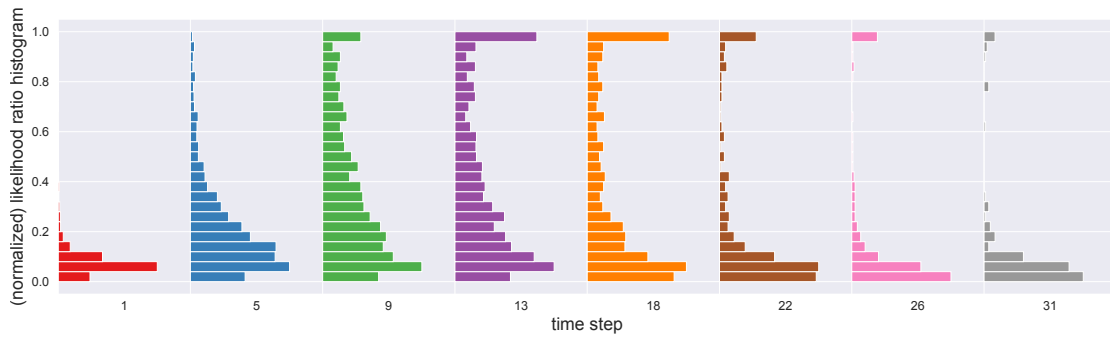


Figure 30: Empty room (batch size 16).

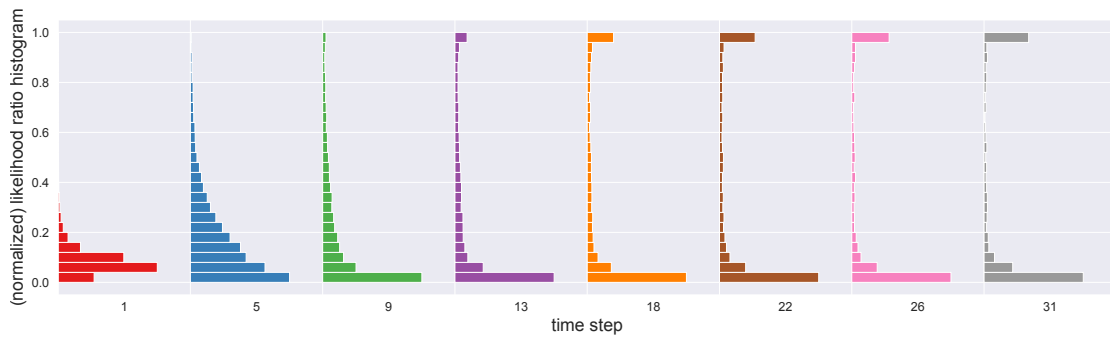


Figure 31: Four rooms (batch size 16).

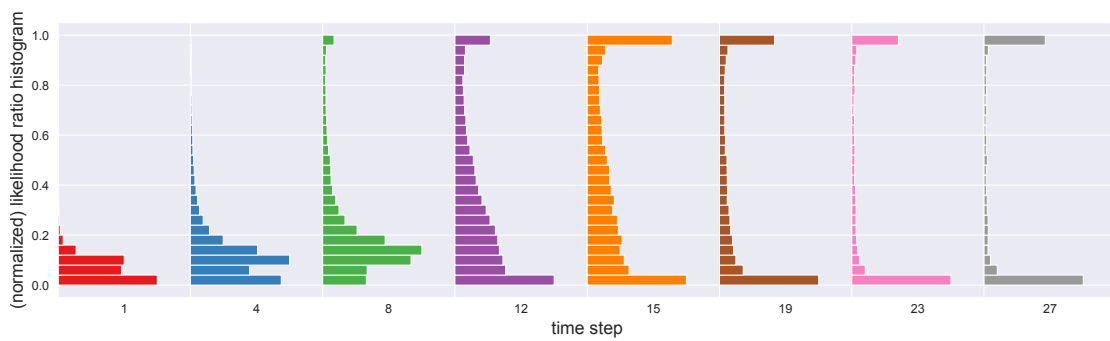


Figure 32: Ms. Pac-man (batch size 16).

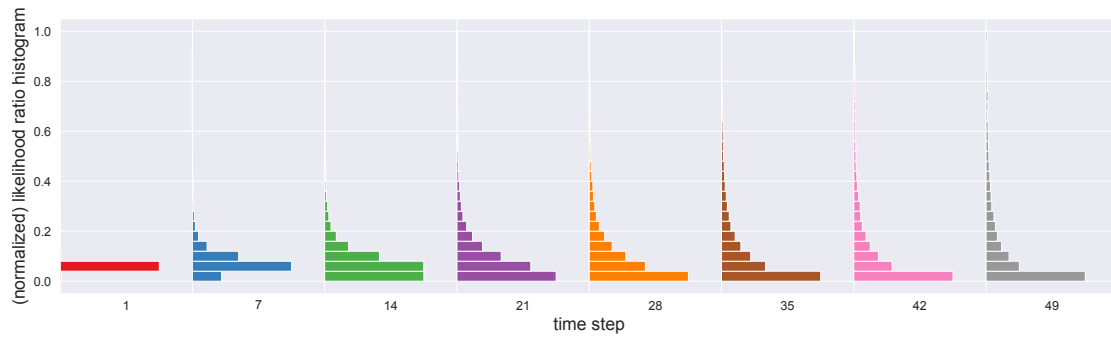


Figure 33: FetchPush (batch size 16).

6.7 Hindsight experience replay study

This section documents an empirical comparison between goal-conditional policy gradients (GCPG), hindsight policy gradients (HPG), deep Q-networks (Mnih et al., 2015, DQN), and a combination of DQN and hindsight experience replay (Andrychowicz et al., 2017, DQN+HER).

Experience replay. Our implementations of both DQN and DQN+HER are based on OpenAI Baselines (Dhariwal et al., 2017), and use mostly the same hyperparameters that Andrychowicz et al. (2017) used in their experiments on environments with discrete action spaces, all of which resemble our bit flipping environments. The only notable differences in our implementations are the lack of both Polyak-averaging and temporal difference target clipping.

Concretely, a *cycle* begins when an agent collects a number of episodes (16) by following an ϵ -greedy policy derived from its deep Q-network ($\epsilon = 0.2$). The corresponding transitions are included in a replay buffer, which contains at most 10^6 transitions. In the case of DQN+HER, hindsight transitions derived from a *final* strategy are also included in this replay buffer (Andrychowicz et al., 2017, Sec. 4.5). Completing the cycle, for a total of 40 different batches, a batch composed of 128 transitions chosen at random from the replay buffer is used to define a loss function and allow one step of gradient-based minimization. The targets required to define these loss functions are computed using a copy of the deep Q-network from the start of the corresponding cycle. Parameters are updated using Adam (Kingma and Ba, 2014). A discount factor of $\gamma = 0.98$ is used, and seems necessary to improve the stability of both DQN and

DQN+HER.

Network architectures. In every experiment, the deep Q-network is implemented by a feedforward neural network with a linear output neuron corresponding to each action. The input to such a network is a triple composed of state, goal, and time step. The network architectures are the same as those described in Section 6.2, except that every weight is initially set using variance scaling (Glorot and Bengio, 2010), and all hidden layers use rectified linear units (Nair and Hinton, 2010). For the Ms. Pac-man environment, the time step information is concatenated with the flattened output of the last convolutional layer (together with the goal information). In comparison to the architecture employed by Andrychowicz et al. (2017) for environments with discrete action spaces, our architectures have one or two additional hidden layers (besides the convolutional architecture used for Ms. Pac-man).

Experimental protocol. The experimental protocol employed in our comparison is very similar to the one described in Section 6.3. Each agent is evaluated periodically, after a number of cycles that depends on the environment. During this evaluation, the agent collects a number of episodes by following a greedy policy derived from its deep Q-network.

For each environment, grid search is used to select the learning rates for both DQN and DQN+HER according to average performance scores (after the corresponding standard deviation across runs is subtracted, as described in Section 6.3). The candidate sets of learning rates are the following. Bit flipping and grid world environments: $\{\alpha \times 10^{-k} \mid \alpha \in \{1, 5\} \text{ and } k \in \{2, 3, 4, 5\}\}$, FetchPush: $\{10^{-2}, 5 \times 10^{-3}, 10^{-3}, 5 \times$

$10^{-4}, 10^{-4}\}$, Ms. Pac-man: $\{10^{-3}, 5 \times 10^{-4}, 10^{-4}, 5 \times 10^{-5}, 10^{-5}\}$. These sets were carefully chosen such that the best performing learning rates do not lie on their extrema.

Definitive results for a given environment are obtained by using the best hyperparameters found for each method in additional runs. These definitive results are directly comparable to our previous results for GCPG and HPG (batch size 16), since every method will have interacted with the environment for the same number of episodes before each evaluation step. For each environment, the number of runs, the number of training batches (cycles), the number of batches (cycles) between evaluations, and the number of episodes per evaluation step are the same as those listed in Tables 1, 2, and 3.

Analysis. The definitive results for the different environments are represented by learning curves (Figs. 34-39, Pg. 42). In the bit flipping environment for $k = 8$ (Figure 34), HPG and DQN+HER lead to equivalent sample efficiency, while GCPG lags far behind and DQN is completely unable to learn. In the bit flipping environment for $k = 16$ (Figure 35), HPG surpasses DQN+HER in sample efficiency by a small margin, while both GCPG and DQN are completely unable to learn. In the empty room environment (Figure 36), HPG is arguably the most sample efficient method, although DQN+HER is more stable upon obtaining a good policy. GCPG eventually obtains a good policy, whereas DQN exhibits instability. In the four rooms environment (Figure 37), DQN+HER outperforms all other methods by a large margin. Although DQN takes much longer to obtain good policies, it would likely surpass both HPG and GCPG given additional training cycles. In the Ms. Pac-man environment (Figure 38), DQN+HER

once again outperforms all other methods, which achieve equivalent sample efficiency (although DQN appears unstable by the end of training). In the FetchPush environment (Figure 39), HPG dramatically outperforms all other methods. Both DQN+HER and DQN are completely unable to learn, while GCPG appears to start learning by the end of the training process. Note that active goals are harshly subsampled to increase the computational efficiency of HPG for both Ms. Pac-man and FetchPush (see Sec. 6.3 and Sec. 6.4).

Discussion. Our results suggest that the decision between applying HPG or DQN+HER in a particular sparse-reward environment requires experimentation. In contrast, the decision to apply hindsight was always successful.

Note that we have not employed heuristics that are known to sometimes increase the performance of policy gradient methods (such as entropy bonuses, reward scaling, learning rate annealing, and simple statistical baselines) to avoid introducing confounding factors. We believe that such heuristics would allow both GCPG and HPG to achieve good results in both the four rooms environment and Ms. Pac-man. Furthermore, whereas hindsight experience replay is directly applicable to state-of-the-art techniques, our work can probably benefit from being extended to state-of-the-art policy gradient approaches, which we intend to explore in future work. Similarly, we believe that additional heuristics and careful hyperparameter settings would allow DQN+HER to achieve good results in the FetchPush environment. This is evidenced by the fact that Andrychowicz et al. (2017) achieve good results using the deep deterministic policy gradient (Lillicrap et al., 2016, DDPG) in a similar environment (with a continuous ac-

tion space and a different reward function). The empirical comparisons between either GCPG and HPG or DQN and DQN+HER are comparatively more conclusive, since the similarities between the methods minimize confounding factors.

Regardless of these empirical results, policy gradient approaches constitute one of the most important classes of model-free reinforcement learning methods, which by itself warrants studying how they can benefit from hindsight. Our approach is also complementary to previous work, since it is entirely possible to combine a critic trained by hindsight experience replay with an actor that employs hindsight policy gradients. Although hindsight experience replay does not require a correction analogous to importance sampling, indiscriminately adding hindsight transitions to the replay buffer is problematic, which has mostly been tackled by heuristics (Andrychowicz et al., 2017, Sec. 4.5). In contrast, our approach seems to benefit from incorporating all available information about goals at every update, which also avoids the need for a replay buffer.

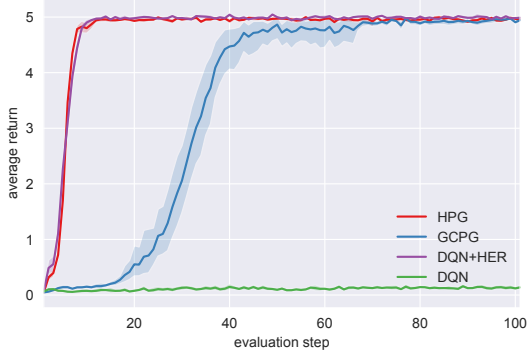


Figure 34: Bit flipping ($k = 8$).

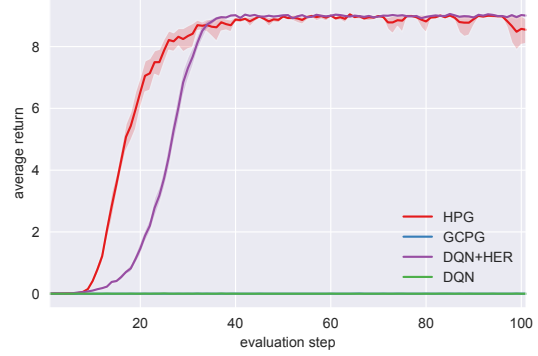


Figure 35: Bit flipping ($k = 16$).



Figure 36: Empty room.

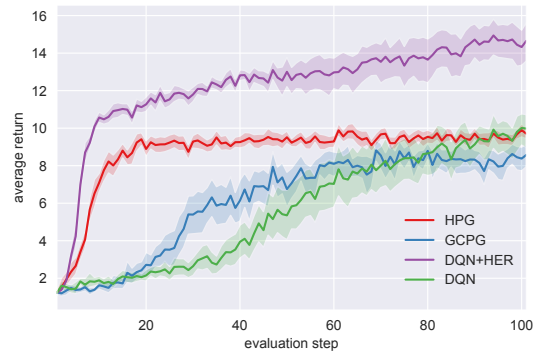


Figure 37: Four rooms.

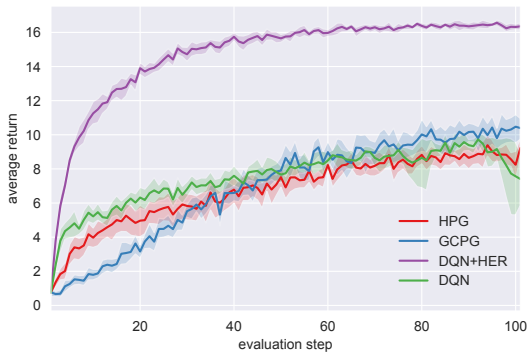


Figure 38: Ms. Pac-man.



Figure 39: FetchPush.

7 Conclusion

We introduced techniques that enable learning goal-conditional policies using hindsight. In this context, hindsight refers to the capacity to exploit information about the degree to which an arbitrary goal has been achieved while another goal was intended. Prior to our work, hindsight has been limited to off-policy reinforcement learning algorithms that rely on experience replay (Andrychowicz et al., 2017) and policy search based on Bayesian optimization (Karkus et al., 2016; Pinsler et al., 2019).

In addition to the fundamental hindsight policy gradient, our technical results include its baseline and advantage formulations. These results are based on a self-contained goal-conditional policy framework that is also introduced in this text. Besides the straightforward estimator built upon the per-decision hindsight policy gradient, we also presented a consistent estimator inspired by weighted importance sampling, together with the corresponding baseline formulation. A variant of this estimator leads to remarkable comparative sample efficiency on a diverse selection of sparse-reward environments, especially in cases where direct reward signals are extremely difficult to obtain. This crucial feature allows natural task formulations that require just trivial reward shaping.

The main drawback of hindsight policy gradient estimators appears to be their computational cost, which is directly related to the number of active goals in a batch. This issue may be mitigated by subsampling active goals, which generally leads to inconsistent estimators. Fortunately, our experiments suggest that this is a viable alternative. Note that the success of hindsight experience replay also depends on an active goal subsampling heuristic (Andrychowicz et al., 2017, Sec. 4.5). The inconsistent hindsight

policy gradient estimator with a value function baseline employed in our experiments sometimes leads to unstable learning, which is likely related to the difficulty of fitting such a value function without hindsight. This hypothesis is consistent with the fact that such instability is observed only in the most extreme examples of sparse-reward environments. Although our preliminary experiments in using hindsight to fit a value function baseline have been successful, this may be accomplished in several ways, and requires a careful study of its own. Further experiments are also required to evaluate hindsight on dense-reward environments.

There are many possibilities for future work besides integrating hindsight policy gradients into systems that rely on goal-conditional policies: deriving additional estimators; implementing and evaluating hindsight (advantage) actor-critic methods; assessing whether hindsight policy gradients can successfully circumvent catastrophic forgetting during curriculum learning of goal-conditional policies; approximating the reward function to reduce required supervision; analysing the variance of the proposed estimators; studying the impact of active goal subsampling; and evaluating every technique on continuous action spaces.

Acknowledgments

We thank Sjoerd van Steenkiste, Klaus Greff, Imanol Schlag, and the anonymous reviewers of previous versions of this work for their valuable feedback. This research was supported by the Swiss National Science Foundation (grant 200021_165675/1), the European Research Council (Advanced Grant 742870), and CAPES (Filipe Mutz, PDSE, 88881.133206/2016-01). We are grateful to Nvidia Corporation for donating a *DGX-1*

machine and to IBM for donating a *Minsky* machine.

A Goal-conditional policy gradients

This appendix contains proofs related to the results presented in Section 3: Theorem 3.1 (App. A.1), Theorem 3.2 (App. A.2), Theorem 3.3 (App. A.4), and Theorem 3.4 (App. A.6). Appendix A.7 presents optimal constant baselines for goal-conditional policies. The remaining subsections contain auxiliary results.

A.1 Theorem 3.1

Theorem 3.1 (Intermediary goal-conditional policy gradient). *The gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \left[\sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \right] \left[\sum_{t=1}^T r(s_t, g) \right]. \quad (3)$$

Proof. The partial derivative $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ of the expected return $\eta(\boldsymbol{\theta})$ with respect to θ_j is given by

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} \frac{\partial}{\partial\theta_j} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^T r(s_t, g). \quad (30)$$

The *likelihood-ratio trick* allows rewriting the previous equation as

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \frac{\partial}{\partial\theta_j} \log p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^T r(s_t, g). \quad (31)$$

Note that

$$\log p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) = \log p(s_1) + \sum_{t=1}^{T-1} \log p(a_t | s_t, g, \boldsymbol{\theta}) + \sum_{t=1}^{T-1} \log p(s_{t+1} | s_t, a_t). \quad (32)$$

Therefore,

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \left[\sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \right] \left[\sum_{t=1}^T r(s_t, g) \right]. \quad (33)$$

□

A.2 Theorem 3.2

Theorem 3.2 (Goal-conditional policy gradient). *The gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(s_{t'}, g). \quad (8)$$

Proof. Starting from Eq. 33, the partial derivative $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ of $\eta(\boldsymbol{\theta})$ with respect to θ_j is given by

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g | \boldsymbol{\theta}) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^T r(s_t, g) \sum_{t'=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(a_{t'} | s_{t'}, g, \boldsymbol{\theta}). \quad (34)$$

The previous equation can be rewritten as

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_{t=1}^T \sum_{t'=1}^{T-1} \mathbb{E} \left[r(S_t, G) \frac{\partial}{\partial\theta_j} \log p(A_{t'} | S_{t'}, G, \boldsymbol{\theta}) | \boldsymbol{\theta} \right]. \quad (35)$$

Let c denote an expectation inside Eq. 35 for $t' \geq t$. In that case, $A_{t'} \perp\!\!\!\perp S_t | S_{t'}, G, \boldsymbol{\Theta}$, and so

$$c = \sum_{s_t} \sum_{s_{t'}} \sum_g \sum_{a_{t'}} p(a_{t'} | s_{t'}, g, \boldsymbol{\theta}) p(s_t, s_{t'}, g | \boldsymbol{\theta}) r(s_t, g) \frac{\partial}{\partial\theta_j} \log p(a_{t'} | s_{t'}, g, \boldsymbol{\theta}). \quad (36)$$

Reversing the likelihood-ratio trick,

$$c = \sum_{s_t} \sum_{s_{t'}} \sum_g p(s_t, s_{t'}, g | \boldsymbol{\theta}) r(s_t, g) \frac{\partial}{\partial\theta_j} \sum_{a_{t'}} p(a_{t'} | s_{t'}, g, \boldsymbol{\theta}) = 0. \quad (37)$$

Therefore, the terms where $t' \geq t$ can be dismissed from Eq. 35, leading to

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=1}^T r(S_t, G) \sum_{t'=1}^{t-1} \frac{\partial}{\partial\theta_j} \log p(A_{t'} | S_{t'}, G, \boldsymbol{\theta}) | \boldsymbol{\theta} \right]. \quad (38)$$

The previous equation can be conveniently rewritten as

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(A_t | S_t, G, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(S_{t'}, G) | \boldsymbol{\theta} \right]. \quad (39)$$

□

A.3 Lemma A.1

Lemma A.1. For every $j, t, \boldsymbol{\theta}$, and associated real-valued (baseline) function b_t^θ ,

$$\sum_{t=1}^{T-1} \mathbb{E} \left[\frac{\partial}{\partial \theta_j} \log p(A_t | S_t, G, \boldsymbol{\theta}) b_t^\theta(S_t, G) | \boldsymbol{\theta} \right] = 0. \quad (40)$$

Proof. Letting c denote an expectation inside Eq. 40,

$$c = \sum_{s_t} \sum_g \sum_{a_t} p(a_t | s_t, g, \boldsymbol{\theta}) p(s_t, g | \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) b_t^\theta(s_t, g). \quad (41)$$

Reversing the likelihood-ratio trick,

$$c = \sum_{s_t} \sum_g p(s_t, g | \boldsymbol{\theta}) b_t^\theta(s_t, g) \frac{\partial}{\partial \theta_j} \sum_{a_t} p(a_t | s_t, g, \boldsymbol{\theta}) = 0. \quad (42)$$

□

A.4 Theorem 3.3

Theorem 3.3 (Goal-conditional policy gradient, baseline formulation). For every $t, \boldsymbol{\theta}$, and associated real-valued (baseline) function b_t^θ , the gradient $\nabla \eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by

$$\nabla \eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \left[\left[\sum_{t'=t+1}^T r(s_{t'}, g) \right] - b_t^\theta(s_t, g) \right]. \quad (9)$$

Proof. The result is obtained by subtracting Eq. 40 from Eq. 39. Importantly, for every combination of $\boldsymbol{\theta}$ and t , it would also be possible to have a distinct baseline function for each parameter in $\boldsymbol{\theta}$. □

A.5 Lemma 3.1

Lemma 3.1 (Goal-conditional policy gradient, action-value formulation). *The gradient*

$\nabla\eta(\boldsymbol{\theta})$ *of the expected return with respect to* $\boldsymbol{\theta}$ *is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) Q_t^\theta(s_t, a_t, g). \quad (10)$$

Proof. Starting from Eq. 39 and rearranging terms,

$$\frac{\partial}{\partial\theta_j} \eta(\boldsymbol{\theta}) = \sum_{t=1}^{T-1} \sum_g \sum_{s_t} \sum_{a_t} p(s_t, a_t, g | \boldsymbol{\theta}) \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{s_{t+1:T}} p(s_{t+1:T} | s_t, a_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(s_{t'}, g). \quad (43)$$

By the definition of action-value function,

$$\frac{\partial}{\partial\theta_j} \eta(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(A_t | S_t, G, \boldsymbol{\theta}) Q_t^\theta(S_t, A_t, G) \mid \boldsymbol{\theta} \right]. \quad (44)$$

□

A.6 Theorem 3.4

Theorem 3.4 (Goal-conditional policy gradient, advantage formulation). *The gradient*

$\nabla\eta(\boldsymbol{\theta})$ *of the expected return with respect to* $\boldsymbol{\theta}$ *is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) A_t^\theta(s_t, a_t, g). \quad (11)$$

Proof. The result is obtained by choosing $b_t^\theta = V_t^\theta$ and subtracting Eq. 40 from Eq.

44. □

A.7 Theorem A.1

For arbitrary j and $\boldsymbol{\theta}$, consider the following definitions of f and h .

$$f(\boldsymbol{\tau}, g) = \sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(s_{t'}, g), \quad (45)$$

$$h(\boldsymbol{\tau}, g) = \sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}). \quad (46)$$

For every $b_j \in \mathbb{R}$, using Theorem 3.2 and the fact that $\mathbb{E}[h(\mathcal{T}, G) \mid \boldsymbol{\theta}] = 0$ by Lemma A.1,

$$\frac{\partial}{\partial \theta_j} \eta(\boldsymbol{\theta}) = \mathbb{E}[f(\mathcal{T}, G) \mid \boldsymbol{\theta}] = \mathbb{E}[f(\mathcal{T}, G) - b_j h(\mathcal{T}, G) \mid \boldsymbol{\theta}]. \quad (47)$$

Theorem A.1. *Assuming $\text{Var}[h(\mathcal{T}, G) \mid \boldsymbol{\theta}] > 0$, the (optimal constant baseline) b_j that minimizes $\text{Var}[f(\mathcal{T}, G) - b_j h(\mathcal{T}, G) \mid \boldsymbol{\theta}]$ is given by*

$$b_j = \frac{\mathbb{E}[f(\mathcal{T}, G)h(\mathcal{T}, G) \mid \boldsymbol{\theta}]}{\mathbb{E}[h(\mathcal{T}, G)^2 \mid \boldsymbol{\theta}]}. \quad (48)$$

Proof. The result is an application of Lemma D.4. □

B Hindsight policy gradients

This appendix contains proofs related to the results presented in Section 4: Theorem 4.1 (App. B.1), Theorem 4.2 (App. B.2), Lemma 4.1 (App. B.3), Theorem B.1 (App. B.4), and Theorem 4.3 (App. B.6). Appendix B.7 presents optimal constant baselines for hindsight policy gradients. Appendix B.5 contains an auxiliary result.

B.1 Theorem 4.1

The following theorem relies on importance sampling, a traditional technique used to obtain estimates related to a random variable $X \sim p$ using samples from an arbitrary positive distribution q . This technique relies on the following equalities:

$$\mathbb{E}_{p(X)} [f(X)] = \sum_x p(x) f(x) = \sum_x \frac{q(x)}{q(x)} p(x) f(x) = \mathbb{E}_{q(X)} \left[\frac{p(X)}{q(X)} f(X) \right]. \quad (49)$$

Theorem 4.1 (Every-decision hindsight policy gradient). *For an arbitrary (original) goal g' , the gradient $\nabla \eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla \eta(\boldsymbol{\theta}) = \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T \left[\prod_{k=1}^{T-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g). \quad (13)$$

Proof. Starting from Theorem 3.2, importance sampling allows rewriting the partial derivative $\partial \eta(\boldsymbol{\theta}) / \partial \theta_j$ as

$$\frac{\partial}{\partial \theta_j} \eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} \frac{p(\boldsymbol{\tau} | g', \boldsymbol{\theta})}{p(\boldsymbol{\tau} | g', \boldsymbol{\theta})} p(\boldsymbol{\tau} | g, \boldsymbol{\theta}) \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(s_{t'}, g). \quad (50)$$

Using Equation 1,

$$\frac{\partial}{\partial \theta_j} \eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \left[\prod_{k=1}^{T-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T r(s_{t'}, g). \quad (51)$$

□

B.2 Theorem 4.2

Theorem 4.2 (Per-decision hindsight policy gradient). *For an arbitrary (original) goal g' , the gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g). \quad (16)$$

Proof. Starting from Eq. 51, the partial derivative $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ can be rewritten as

$$\frac{\partial}{\partial\theta_j} \eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{t=1}^{T-1} \sum_{t'=t+1}^T \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \left[\prod_{k=1}^{T-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) r(s_{t'}, g). \quad (52)$$

If we split every trajectory into states and actions before and after t' , then $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ is given by

$$\sum_g p(g) \sum_{t=1}^{T-1} \sum_{t'=t+1}^T \sum_{s_{1:t'-1}, a_{1:t'-1}} p(s_{1:t'-1}, a_{1:t'-1} | g', \boldsymbol{\theta}) \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) z, \quad (53)$$

where z is defined by

$$z = \sum_{s_{t':T}} \sum_{a_{t':T-1}} p(s_{t':T}, a_{t':T-1} | s_{1:t'-1}, a_{1:t'-1}, g', \boldsymbol{\theta}) \left[\prod_{k=t'}^{T-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g). \quad (54)$$

Using Lemma D.2 and canceling terms,

$$z = \sum_{s_{t':T}} \sum_{a_{t':T-1}} p(s_{t'} | s_{t'-1}, a_{t'-1}) \left[\prod_{k=t'}^{T-1} p(a_k | s_k, g, \boldsymbol{\theta}) p(s_{k+1} | s_k, a_k) \right] r(s_{t'}, g). \quad (55)$$

Using Lemma D.2 once again,

$$z = \sum_{s_{t':T}} \sum_{a_{t':T-1}} p(s_{t':T}, a_{t':T-1} | s_{1:t'-1}, a_{1:t'-1}, g, \boldsymbol{\theta}) r(s_{t'}, g). \quad (56)$$

Using the fact that $S_{t'} \perp\!\!\!\perp G \mid S_{1:t'-1}, A_{1:t'-1}, \Theta$,

$$z = \sum_{s_{t'}} r(s_{t'}, g) p(s_{t'} \mid s_{1:t'-1}, a_{1:t'-1}, g, \theta) = \sum_{s_{t'}} r(s_{t'}, g) p(s_{t'} \mid s_{1:t'-1}, a_{1:t'-1}, g', \theta). \quad (57)$$

Substituting z into Expression 53 and returning to an expectation over trajectories,

$$\frac{\partial}{\partial \theta_j} \eta(\theta) = \sum_{\tau} p(\tau \mid g', \theta) \sum_g p(g) \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_t \mid s_t, g, \theta) \sum_{t'=t+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(a_k \mid s_k, g, \theta)}{p(a_k \mid s_k, g', \theta)} \right] r(s_{t'}, g). \quad (58)$$

□

B.3 Lemma 4.1

Lemma 4.1. *For every g', t, θ , and associated real-valued (baseline) function b_t^θ ,*

$$\sum_{\tau} p(\tau \mid g', \theta) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t \mid s_t, g, \theta) \left[\prod_{k=1}^t \frac{p(a_k \mid s_k, g, \theta)}{p(a_k \mid s_k, g', \theta)} \right] b_t^\theta(s_t, g) = \mathbf{0}. \quad (24)$$

Proof. Let c denote the j -th element of the vector in the left-hand side of Eq. 24, such that

$$c = \sum_g p(g) \sum_{t=1}^{T-1} \mathbb{E} \left[\frac{\partial}{\partial \theta_j} \log p(A_t \mid S_t, g, \theta) \left[\prod_{k=1}^t \frac{p(A_k \mid S_k, g, \theta)}{p(A_k \mid S_k, g', \theta)} \right] b_t^\theta(S_t, g) \mid g', \theta \right]. \quad (59)$$

Using Lemma D.1 and writing the expectations explicitly,

$$c = \sum_g p(g) \sum_{t=1}^{T-1} \sum_{s_{1:t}} \sum_{a_{1:t}} p(s_{1:t}, a_{1:t} \mid g', \theta) \frac{\partial}{\partial \theta_j} \log p(a_t \mid s_t, g, \theta) \frac{p(s_{1:t}, a_{1:t} \mid g, \theta)}{p(s_{1:t}, a_{1:t} \mid g', \theta)} b_t^\theta(s_t, g). \quad (60)$$

Canceling terms, using Lemma D.1 once again, and reversing the likelihood-ratio

trick,

$$c = \sum_g p(g) \sum_{t=1}^{T-1} \sum_{s_{1:t}} \sum_{a_{1:t}} \frac{\partial}{\partial \theta_j} p(a_t | s_t, g, \boldsymbol{\theta}) \left[p(s_1) \prod_{k=1}^{t-1} p(a_k | s_k, g, \boldsymbol{\theta}) p(s_{k+1} | s_k, a_k) \right] b_t^\theta(s_t, g). \quad (61)$$

Pushing constants outside the summation over actions at time step t ,

$$c = \sum_g p(g) \sum_{t=1}^{T-1} \sum_{s_{1:t}} \sum_{a_{1:t-1}} \left[p(s_1) \prod_{k=1}^{t-1} p(a_k | s_k, g, \boldsymbol{\theta}) p(s_{k+1} | s_k, a_k) \right] b_t^\theta(s_t, g) \frac{\partial}{\partial \theta_j} \sum_{a_t} p(a_t | s_t, g, \boldsymbol{\theta}) = 0. \quad (62)$$

□

B.4 Theorem B.1

Theorem B.1 (Hindsight policy gradient, baseline formulation). *For every g' , t , $\boldsymbol{\theta}$, and associated real-valued (baseline) function b_t^θ , the gradient $\nabla \eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla \eta(\boldsymbol{\theta}) = \sum_{\boldsymbol{\tau}} p(\boldsymbol{\tau} | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) z, \quad (63)$$

where

$$z = \left[\sum_{t'=t+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g) \right] - \left[\prod_{k=1}^t \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] b_t^\theta(s_t, g). \quad (64)$$

Proof. The result is obtained by subtracting Eq. 24 from Eq. 16. Importantly, for every combination of $\boldsymbol{\theta}$ and t , it would also be possible to have a distinct baseline function for each parameter in $\boldsymbol{\theta}$. □

B.5 Lemma 4.2

Lemma 4.2 (Hindsight policy gradient, action-value formulation). *For an arbitrary goal g' , the gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_{\tau} p(\tau | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \left[\prod_{k=1}^t \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] Q_t^{\boldsymbol{\theta}}(s_t, a_t, g). \quad (25)$$

Proof. Starting from Eq. 44, the partial derivative $\partial\eta(\boldsymbol{\theta})/\partial\theta_j$ can be written as

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_{t=1}^{T-1} \sum_g p(g) \sum_{s_{1:t}} \sum_{a_{1:t}} p(s_{1:t}, a_{1:t} | g, \boldsymbol{\theta}) \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) Q_t^{\boldsymbol{\theta}}(s_t, a_t, g). \quad (65)$$

Using importance sampling, for an arbitrary goal g' ,

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \sum_{t=1}^{T-1} \sum_{s_{1:t}} \sum_{a_{1:t}} p(s_{1:t}, a_{1:t} | g', \boldsymbol{\theta}) \frac{p(s_{1:t}, a_{1:t} | g, \boldsymbol{\theta})}{p(s_{1:t}, a_{1:t} | g', \boldsymbol{\theta})} \frac{\partial}{\partial\theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) Q_t^{\boldsymbol{\theta}}(s_t, a_t, g). \quad (66)$$

Using Lemma D.1 and rewriting the previous equation using expectations,

$$\frac{\partial}{\partial\theta_j}\eta(\boldsymbol{\theta}) = \sum_g p(g) \mathbb{E} \left[\sum_{t=1}^{T-1} \frac{\partial}{\partial\theta_j} \log p(A_t | S_t, g, \boldsymbol{\theta}) \left[\prod_{k=1}^t \frac{p(A_k | S_k, g, \boldsymbol{\theta})}{p(A_k | S_k, g', \boldsymbol{\theta})} \right] Q_t^{\boldsymbol{\theta}}(S_t, A_t, g) | g', \boldsymbol{\theta} \right]. \quad (67)$$

□

B.6 Theorem 4.3

Theorem 4.3 (Hindsight policy gradient, advantage formulation). *For an arbitrary (original) goal g' , the gradient $\nabla\eta(\boldsymbol{\theta})$ of the expected return with respect to $\boldsymbol{\theta}$ is given by*

$$\nabla\eta(\boldsymbol{\theta}) = \sum_{\tau} p(\tau | g', \boldsymbol{\theta}) \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(a_t | s_t, g, \boldsymbol{\theta}) \left[\prod_{k=1}^t \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] A_t^{\boldsymbol{\theta}}(s_t, a_t, g). \quad (26)$$

Proof. The result is obtained by choosing $b_t^\theta = V_t^\theta$ and subtracting Eq. 59 from Eq. 67. □

B.7 Theorem B.2

For arbitrary g', j , and θ , consider the following definitions of f and h .

$$f(\tau) = \sum_g p(g) \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_t | s_t, g, \theta) \sum_{t'=t+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \theta)}{p(a_k | s_k, g', \theta)} \right] r(s_{t'}, g), \quad (68)$$

$$h(\tau) = \sum_g p(g) \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_t | s_t, g, \theta) \prod_{k=1}^t \frac{p(a_k | s_k, g, \theta)}{p(a_k | s_k, g', \theta)}. \quad (69)$$

For every $b_j \in \mathbb{R}$, using Theorem 4.2 and the fact that $\mathbb{E}[h(\mathcal{T}) | g', \theta] = 0$ by Lemma 4.1,

$$\frac{\partial}{\partial \theta_j} \eta(\theta) = \mathbb{E}[f(\mathcal{T}) | g', \theta] = \mathbb{E}[f(\mathcal{T}) - b_j h(\mathcal{T}) | g', \theta]. \quad (70)$$

Theorem B.2. *Assuming $\text{Var}[h(\mathcal{T}) | g', \theta] > 0$, the (optimal constant baseline) b_j that minimizes $\text{Var}[f(\mathcal{T}) - b_j h(\mathcal{T}) | g', \theta]$ is given by*

$$b_j = \frac{\mathbb{E}[f(\mathcal{T})h(\mathcal{T}) | g', \theta]}{\mathbb{E}[h(\mathcal{T})^2 | g', \theta]}. \quad (71)$$

Proof. The result is an application of Lemma D.4. □

C Hindsight gradient estimators

This appendix contains proofs related to the estimators presented in Section 5: Theorem 5.1 (App. C.1) and Theorem 5.2 (App. C.2). Appendix C.3 presents a result that enables a consistency-preserving *weighted baseline*.

In this appendix, we will consider a dataset $\mathcal{D} = \{(\boldsymbol{\tau}^{(i)}, g^{(i)})\}_{i=1}^N$ where each trajectory $\boldsymbol{\tau}^{(i)}$ is obtained using a policy parameterized by $\boldsymbol{\theta}$ in an attempt to achieve a goal $g^{(i)}$ chosen by the environment. Because \mathcal{D} is an *iid* dataset given Θ ,

$$p(\mathcal{D} \mid \boldsymbol{\theta}) = p(\boldsymbol{\tau}^{(1:N)}, g^{(1:N)} \mid \boldsymbol{\theta}) = \prod_{i=1}^N p(\boldsymbol{\tau}^{(i)}, g^{(i)} \mid \boldsymbol{\theta}) = \prod_{i=1}^N p(g^{(i)})p(\boldsymbol{\tau}^{(i)} \mid g^{(i)}, \boldsymbol{\theta}). \quad (72)$$

C.1 Theorem 5.1

Theorem 5.1. *The per-decision hindsight policy gradient estimator, given by*

$$\frac{1}{N} \sum_{i=1}^N \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(A_t^{(i)} \mid S_t^{(i)}, G^{(i)} = g, \boldsymbol{\theta}) \sum_{t'=t+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(A_k^{(i)} \mid S_k^{(i)}, G^{(i)} = g, \boldsymbol{\theta})}{p(A_k^{(i)} \mid S_k^{(i)}, G^{(i)}, \boldsymbol{\theta})} \right] r(S_{t'}^{(i)}, g), \quad (28)$$

is a consistent and unbiased estimator of the gradient $\nabla \eta(\boldsymbol{\theta})$ of the expected return.

Proof. Let $I_j^{(N)}$ denote the j -th element of the estimator, which can be written as

$$I_j^{(N)} = \frac{1}{N} \sum_{i=1}^N I(\boldsymbol{\tau}^{(i)}, G^{(i)}, \boldsymbol{\theta})_j, \quad (73)$$

where

$$I(\boldsymbol{\tau}, g', \boldsymbol{\theta})_j = \sum_g p(g) \sum_{t=1}^{T-1} \frac{\partial}{\partial \theta_j} \log p(a_t \mid s_t, g, \boldsymbol{\theta}) \sum_{t'=t+1}^T \left[\prod_{k=1}^{t'-1} \frac{p(a_k \mid s_k, g, \boldsymbol{\theta})}{p(a_k \mid s_k, g', \boldsymbol{\theta})} \right] r(s_{t'}, g). \quad (74)$$

Using Theorem 4.2, the expected value $\mathbb{E} \left[I_j^{(N)} \mid \boldsymbol{\theta} \right]$ is given by

$$\mathbb{E} \left[I_j^{(N)} \mid \boldsymbol{\theta} \right] = \frac{1}{N} \sum_{i=1}^N \sum_{g^{(i)}} p(g^{(i)}) \mathbb{E} \left[I(\mathcal{T}^{(i)}, g^{(i)}, \boldsymbol{\theta})_j \mid g^{(i)}, \boldsymbol{\theta} \right] = \frac{1}{N} \sum_{i=1}^N \sum_{g^{(i)}} p(g^{(i)}) \frac{\partial}{\partial \theta_j} \eta(\boldsymbol{\theta}) = \frac{\partial}{\partial \theta_j} \eta(\boldsymbol{\theta}). \quad (75)$$

Therefore, $I_j^{(N)}$ is an unbiased estimator of $\partial \eta(\boldsymbol{\theta}) / \partial \theta_j$.

Conditionally on Θ , the random variable $I_j^{(N)}$ is an average of iid random variables with expected value $\partial \eta(\boldsymbol{\theta}) / \partial \theta_j$ (see Eq. 75). By the strong law of large numbers (Sen and Singer, 1994, Theorem 2.3.13),

$$I_j^{(N)} \xrightarrow{\text{a.s.}} \frac{\partial}{\partial \theta_j} \eta(\boldsymbol{\theta}). \quad (76)$$

Therefore, $I_j^{(N)}$ is a consistent estimator of $\partial \eta(\boldsymbol{\theta}) / \partial \theta_j$. □

C.2 Theorem 5.2

Theorem 5.2. *The weighted per-decision hindsight policy gradient estimator, given by*

$$\sum_{i=1}^N \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(A_t^{(i)} \mid S_t^{(i)}, G^{(i)} = g, \boldsymbol{\theta}) \sum_{t'=t+1}^T \frac{\left[\prod_{k=1}^{t'-1} \frac{p(A_k^{(i)} \mid S_k^{(i)}, G^{(i)} = g, \boldsymbol{\theta})}{p(A_k^{(i)} \mid S_k^{(i)}, G^{(i)}, \boldsymbol{\theta})} \right] r(S_{t'}^{(i)}, g)}{\sum_{j=1}^N \left[\prod_{k=1}^{t'-1} \frac{p(A_k^{(j)} \mid S_k^{(j)}, G^{(j)} = g, \boldsymbol{\theta})}{p(A_k^{(j)} \mid S_k^{(j)}, G^{(j)}, \boldsymbol{\theta})} \right]}, \quad (29)$$

is a consistent estimator of the gradient $\nabla \eta(\boldsymbol{\theta})$ of the expected return.

Proof. Let $W_j^{(N)}$ denote the j -th element of the estimator, which can be written as

$$W_j^{(N)} = \sum_g p(g) \sum_{t=1}^{T-1} \sum_{t'=t+1}^T \frac{X(g, t, t')_j^{(N)}}{Y(g, t, t')_j^{(N)}}, \quad (77)$$

where

$$X(g, t, t')_j^{(N)} = \frac{1}{N} \sum_{i=1}^N X(\mathcal{T}^{(i)}, G^{(i)}, g, t, t', \boldsymbol{\theta})_j, \quad (78)$$

$$Y(g, t, t')_j^{(N)} = \frac{1}{N} \sum_{i=1}^N Y(\mathcal{T}^{(i)}, G^{(i)}, g, t, t', \boldsymbol{\theta})_j, \quad (79)$$

$$X(\boldsymbol{\tau}, g', g, t, t', \boldsymbol{\theta})_j = \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] \frac{\partial}{\partial \theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) r(s_{t'}, g), \quad (80)$$

$$Y(\boldsymbol{\tau}, g', g, t, t', \boldsymbol{\theta})_j = \left[\prod_{k=1}^{t'-1} \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right]. \quad (81)$$

Consider the expected value $E_{X_i} = \mathbb{E} \left[X(\mathcal{T}^{(i)}, G^{(i)}, g, t, t', \boldsymbol{\theta})_j \mid \boldsymbol{\theta} \right]$, which is given

by

$$E_{X_i} = \sum_{g^{(i)}} p(g^{(i)}) \mathbb{E} \left[\left[\prod_{k=1}^{t'-1} \frac{p(A_k | S_k, g, \boldsymbol{\theta})}{p(A_k | S_k, G = g^{(i)}, \boldsymbol{\theta})} \right] \frac{\partial}{\partial \theta_j} \log p(A_t | S_t, g, \boldsymbol{\theta}) r(S_{t'}, g) \mid G = g^{(i)}, \boldsymbol{\theta} \right]. \quad (82)$$

Using the fact that $t' > t$, Lemma D.1, and canceling terms, E_{X_i} can be written as

$$\sum_{g^{(i)}} p(g^{(i)}) \sum_{s_{1:t'}, a_{1:t'-1}} p(s_{t'} | s_{1:t'-1}, a_{1:t'-1}, G = g^{(i)}, \boldsymbol{\theta}) p(s_{1:t'-1}, a_{1:t'-1} | g, \boldsymbol{\theta}) \frac{\partial}{\partial \theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) r(s_{t'}, g). \quad (83)$$

Because $S_{t'} \perp\!\!\!\perp G \mid S_{1:t'-1}, A_{1:t'-1}, \boldsymbol{\Theta}$,

$$E_{X_i} = \mathbb{E} \left[\frac{\partial}{\partial \theta_j} \log p(A_t | S_t, g, \boldsymbol{\theta}) r(S_{t'}, g) \mid g, \boldsymbol{\theta} \right]. \quad (84)$$

Conditionally on $\boldsymbol{\Theta}$, the variable $X(g, t, t')_j^{(N)}$ is an average of iid random variables with expected value E_{X_i} . By the strong law of large numbers (Sen and Singer, 1994, Theorem 2.3.13), $X(g, t, t')_j^{(N)} \xrightarrow{\text{a.s.}} E_{X_i}$.

Using Lemma D.1, the expected value $E_{Y_i} = \mathbb{E} \left[Y(\mathcal{T}^{(i)}, G^{(i)}, g, t, t', \boldsymbol{\theta})_j \mid \boldsymbol{\theta} \right]$ is given by

$$E_{Y_i} = \sum_{g^{(i)}} p(g^{(i)}) \mathbb{E} \left[\frac{p(S_{1:t'-1}, A_{1:t'-1}^{(i)} | G^{(i)} = g, \boldsymbol{\theta})}{p(S_{1:t'-1}, A_{1:t'-1}^{(i)} | g^{(i)}, \boldsymbol{\theta})} \mid g^{(i)}, \boldsymbol{\theta} \right] = 1. \quad (85)$$

Conditionally on Θ , the variable $Y(g, t, t')_j^{(N)}$ is an average of iid random variables with expected value 1. By the strong law of large numbers, $Y(g, t, t')_j^{(N)} \xrightarrow{\text{a.s.}} 1$.

Because both $X(g, t, t')_j^{(N)}$ and $Y(g, t, t')_j^{(N)}$ converge almost surely to real numbers (Thomas, 2015, Ch. 3, Property 2),

$$\frac{X(g, t, t')_j^{(N)}}{Y(g, t, t')_j^{(N)}} \xrightarrow{\text{a.s.}} \mathbb{E} \left[\frac{\partial}{\partial \theta_j} \log p(A_t | S_t, g, \boldsymbol{\theta}) r(S_{t'}, g) | g, \boldsymbol{\theta} \right]. \quad (86)$$

By Theorem 3.2 and the fact that $W_j^{(N)}$ is a linear combination of terms $X(g, t, t')_j^{(N)} / Y(g, t, t')_j^{(N)}$,

$$W_j^{(N)} \xrightarrow{\text{a.s.}} \sum_g p(g) \sum_{t=1}^{T-1} \sum_{t'=t+1}^T \mathbb{E} \left[\frac{\partial}{\partial \theta_j} \log p(A_t | S_t, g, \boldsymbol{\theta}) r(S_{t'}, g) | g, \boldsymbol{\theta} \right] = \frac{\partial}{\partial \theta_j} \eta(\boldsymbol{\theta}). \quad (87)$$

□

C.3 Theorem C.1

Theorem C.1. *The weighted baseline estimator, given by*

$$\sum_{i=1}^N \sum_g p(g) \sum_{t=1}^{T-1} \nabla \log p(A_t^{(i)} | S_t^{(i)}, G^{(i)} = g, \boldsymbol{\theta}) \frac{\left[\prod_{k=1}^t \frac{p(A_k^{(i)} | S_k^{(i)}, G^{(i)} = g, \boldsymbol{\theta})}{p(A_k^{(i)} | S_k^{(i)}, G^{(i)}, \boldsymbol{\theta})} \right] b_t^\boldsymbol{\theta}(S_t^{(i)}, g)}{\sum_{j=1}^N \left[\prod_{k=1}^t \frac{p(A_k^{(j)} | S_k^{(j)}, G^{(j)} = g, \boldsymbol{\theta})}{p(A_k^{(j)} | S_k^{(j)}, G^{(j)}, \boldsymbol{\theta})} \right]}, \quad (88)$$

converges almost surely to zero.

Proof. Let $B_j^{(N)}$ denote the j -th element of the estimator, which can be written as

$$B_j^{(N)} = \sum_g p(g) \sum_{t=1}^{T-1} \frac{X(g, t)_j^{(N)}}{Y(g, t)_j^{(N)}}, \quad (89)$$

where

$$X(g, t)_j^{(N)} = \frac{1}{N} \sum_{i=1}^N X(\mathcal{T}^{(i)}, G^{(i)}, g, t, \boldsymbol{\theta})_j, \quad (90)$$

$$Y(g, t)_j^{(N)} = \frac{1}{N} \sum_{i=1}^N Y(\mathcal{T}^{(i)}, G^{(i)}, g, t, \boldsymbol{\theta})_j, \quad (91)$$

$$X(\boldsymbol{\tau}, g', g, t, \boldsymbol{\theta})_j = \left[\prod_{k=1}^t \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})} \right] \frac{\partial}{\partial \theta_j} \log p(a_t | s_t, g, \boldsymbol{\theta}) b_t^\theta(s_t, g), \quad (92)$$

$$Y(\boldsymbol{\tau}, g', g, t, \boldsymbol{\theta})_j = \prod_{k=1}^t \frac{p(a_k | s_k, g, \boldsymbol{\theta})}{p(a_k | s_k, g', \boldsymbol{\theta})}. \quad (93)$$

Using Eqs. 59 and 62, the expected value $E_{X_i} = \mathbb{E} \left[X(\mathcal{T}^{(i)}, G^{(i)}, g, t, \boldsymbol{\theta})_j \mid \boldsymbol{\theta} \right]$ is given by

$$E_{X_i} = \sum_{g^{(i)}} p(g^{(i)}) \mathbb{E} \left[X(\mathcal{T}^{(i)}, g^{(i)}, g, t, \boldsymbol{\theta})_j \mid g^{(i)}, \boldsymbol{\theta} \right] = 0. \quad (94)$$

Conditionally on Θ , the variable $X(g, t)_j^{(N)}$ is an average of iid random variables with expected value zero. By the strong law of large numbers (Sen and Singer, 1994, Theorem 2.3.13), $X(g, t)_j^{(N)} \xrightarrow{\text{a.s.}} 0$.

The fact that $Y(g, t)_j^{(N)} \xrightarrow{\text{a.s.}} 1$ is already established in the proof of Theorem 5.2. Because both $X(g, t)_j^{(N)}$ and $Y(g, t)_j^{(N)}$ converge almost surely to real numbers (Thomas, 2015, Ch. 3, Property 2),

$$\frac{X(g, t)_j^{(N)}}{Y(g, t)_j^{(N)}} \xrightarrow{\text{a.s.}} 0. \quad (95)$$

Because $B_j^{(N)}$ is a linear combination of terms $X(g, t)_j^{(N)} / Y(g, t)_j^{(N)}$, $B_j^{(N)} \xrightarrow{\text{a.s.}} 0$.

□

Clearly, if $E^{(N)}$ is a consistent estimator of a some quantity given $\boldsymbol{\theta}$, then so is $E^{(N)} - B_j^{(N)}$, which allows using this result in combination with Theorem 5.2.

D Fundamental results

This appendix presents results required by previous sections: Lemma D.1 (App. D.1), Lemma D.2 (App. D.2), Theorem 4.4 (App. D.4), and Lemma D.4 (App. D.5). Appendix D.3 contains an auxiliary result.

D.1 Lemma D.1

Lemma D.1. *For every $\tau, g, \boldsymbol{\theta}$, and $1 \leq t \leq T - 1$,*

$$p(s_{1:t}, a_{1:t} \mid g, \boldsymbol{\theta}) = p(s_1)p(a_t \mid s_t, g, \boldsymbol{\theta}) \prod_{k=1}^{t-1} p(a_k \mid s_k, g, \boldsymbol{\theta})p(s_{k+1} \mid s_k, a_k). \quad (96)$$

Proof. In order to employ backward induction, consider the case $t = T - 1$. By marginalization,

$$p(s_{1:T-1}, a_{1:T-1} \mid g, \boldsymbol{\theta}) = \sum_{s_T} p(\tau \mid g, \boldsymbol{\theta}) = \sum_{s_T} p(s_1) \prod_{k=1}^{T-1} p(a_k \mid s_k, g, \boldsymbol{\theta})p(s_{k+1} \mid s_k, a_k) \quad (97)$$

$$= p(s_1)p(a_{T-1} \mid s_{T-1}, g, \boldsymbol{\theta}) \prod_{k=1}^{T-2} p(a_k \mid s_k, g, \boldsymbol{\theta})p(s_{k+1} \mid s_k, a_k), \quad (98)$$

which completes the proof of the base case.

Assuming the inductive hypothesis is true for a given $2 \leq t \leq T - 1$ and considering the case $t - 1$,

$$p(s_{1:t-1}, a_{1:t-1} \mid g, \boldsymbol{\theta}) = \sum_{s_t} \sum_{a_t} p(s_1) p(a_t \mid s_t, g, \boldsymbol{\theta}) \prod_{k=1}^{t-1} p(a_k \mid s_k, g, \boldsymbol{\theta}) p(s_{k+1} \mid s_k, a_k) \quad (99)$$

$$= p(s_1) p(a_{t-1} \mid s_{t-1}, g, \boldsymbol{\theta}) \prod_{k=1}^{t-2} p(a_k \mid s_k, g, \boldsymbol{\theta}) p(s_{k+1} \mid s_k, a_k). \quad (100)$$

□

D.2 Lemma D.2

Lemma D.2. For every $\tau, g, \boldsymbol{\theta}$, and $1 \leq t \leq T$,

$$p(s_{t:T}, a_{t:T-1} \mid s_{1:t-1}, a_{1:t-1}, g, \boldsymbol{\theta}) = p(s_t \mid s_{t-1}, a_{t-1}) \prod_{k=t}^{T-1} p(a_k \mid s_k, g, \boldsymbol{\theta}) p(s_{k+1} \mid s_k, a_k). \quad (101)$$

Proof. The case $t = 1$ can be inspected easily. Consider $2 \leq t \leq T$. By definition,

$$p(s_{t:T}, a_{t:T-1} \mid s_{1:t-1}, a_{1:t-1}, g, \boldsymbol{\theta}) = \frac{p(s_{1:T}, a_{1:T-1} \mid g, \boldsymbol{\theta})}{p(s_{1:t-1}, a_{1:t-1} \mid g, \boldsymbol{\theta})}. \quad (102)$$

Using Lemma D.1,

$$p(s_{t:T}, a_{t:T-1} \mid s_{1:t-1}, a_{1:t-1}, g, \boldsymbol{\theta}) = \frac{p(s_1) \prod_{k=1}^{T-1} p(a_k \mid s_k, g, \boldsymbol{\theta}) p(s_{k+1} \mid s_k, a_k)}{p(s_1) p(a_{t-1} \mid s_{t-1}, g, \boldsymbol{\theta}) \prod_{k=1}^{t-2} p(a_k \mid s_k, g, \boldsymbol{\theta}) p(s_{k+1} \mid s_k, a_k)} \quad (103)$$

$$= \frac{\prod_{k=t-1}^{T-1} p(a_k \mid s_k, g, \boldsymbol{\theta}) p(s_{k+1} \mid s_k, a_k)}{p(a_{t-1} \mid s_{t-1}, g, \boldsymbol{\theta})}. \quad (104)$$

□

D.3 Lemma D.3

Lemma D.3. For every t and θ , the action-value function Q_t^θ is given by

$$Q_t^\theta(s, a, g) = \mathbb{E} [r(S_{t+1}, g) + V_{t+1}^\theta(S_{t+1}, g) \mid S_t = s, A_t = a]. \quad (105)$$

Proof. From the definition of action-value function and using the fact that $S_{t+1} \perp\!\!\!\perp G, \Theta \mid S_t, A_t$,

$$Q_t^\theta(s, a, g) = \mathbb{E} [r(S_{t+1}, g) \mid S_t = s, A_t = a] + \mathbb{E} \left[\sum_{t'=t+2}^T r(S_{t'}, g) \mid S_t = s, A_t = a, g, \theta \right]. \quad (106)$$

Let z denote the second term in the right-hand side of the previous equation, which can also be written as

$$z = \sum_{s_{t+1}} \sum_{a_{t+1}} \sum_{s_{t+2:T}} p(s_{t+1}, a_{t+1}, s_{t+2:T} \mid S_t = s, A_t = a, g, \theta) \sum_{t'=t+2}^T r(s_{t'}, g). \quad (107)$$

Consider the following three independence properties:

$$S_{t+1} \perp\!\!\!\perp G, \Theta \mid S_t, A_t, \quad (108)$$

$$A_{t+1} \perp\!\!\!\perp S_t, A_t \mid S_{t+1}, G, \Theta, \quad (109)$$

$$S_{t+2:T} \perp\!\!\!\perp S_t, A_t \mid S_{t+1}, A_{t+1}, G, \Theta. \quad (110)$$

Together, these properties can be used to demonstrate that

$$z = \sum_{s_{t+1}} p(s_{t+1} \mid S_t = s, A_t = a) \sum_{a_{t+1}} p(a_{t+1} \mid s_{t+1}, g, \theta) \sum_{s_{t+2:T}} p(s_{t+2:T} \mid s_{t+1}, a_{t+1}, g, \theta) \sum_{t'=t+2}^T r(s_{t'}, g). \quad (111)$$

From the definition of value function, $z = \mathbb{E} [V_{t+1}^\theta(S_{t+1}, g) \mid S_t = s, A_t = a]$.

□

D.4 Theorem 4.4

Theorem 4.4. For every t and θ , the advantage function A_t^θ is given by

$$A_t^\theta(s, a, g) = \mathbb{E} [r(S_{t+1}, g) + V_{t+1}^\theta(S_{t+1}, g) - V_t^\theta(s, g) \mid S_t = s, A_t = a]. \quad (27)$$

Proof. The result follows from the definition of advantage function and Lemma D.3. □

D.5 Lemma D.4

Consider a discrete random variable X and real-valued functions f and h . Suppose also that $\mathbb{E}[h(X)] = 0$ and $\text{Var}[h(X)] > 0$. Clearly, for every $b \in \mathbb{R}$, we have $\mathbb{E}[f(X) - bh(X)] = \mathbb{E}[f(X)]$.

Lemma D.4. The constant $b \in \mathbb{R}$ that minimizes $\text{Var}[f(X) - bh(X)]$ is given by

$$b = \frac{\mathbb{E}[f(X)h(X)]}{\mathbb{E}[h(X)^2]}. \quad (112)$$

Proof. Let $v = \text{Var}[f(X) - bh(X)]$. Using our assumptions and the definition of variance,

$$v = \mathbb{E}[(f(X) - bh(X))^2] - \mathbb{E}[f(X) - bh(X)]^2 = \mathbb{E}[(f(X) - bh(X))^2] - \mathbb{E}[f(X)]^2 \quad (113)$$

$$= \mathbb{E}[f(X)^2] - 2b\mathbb{E}[f(X)h(X)] + b^2\mathbb{E}[h(X)^2] - \mathbb{E}[f(X)]^2. \quad (114)$$

The first and second derivatives of v with respect to b are given by $dv/db = -2\mathbb{E}[f(X)h(X)] + 2b\mathbb{E}[h(X)^2]$ and $d^2v/db^2 = 2\mathbb{E}[h(X)^2]$. Our assumptions guarantee that $\mathbb{E}[h(X)^2] >$

0. Therefore, by Fermat's theorem, if b is a local minimum, then $dv/db = 0$, leading to the desired equality. By the second derivative test, b must be a local minimum.

□

References

- Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Abbeel, P., and Zaremba, W. (2017). Hindsight experience replay. In *Advances in Neural Information Processing Systems*, pages 5048–5058.
- Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. (2013). The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279.
- Bishop, C. M. (2013). *Pattern Recognition and Machine Learning*. Information science and statistics. Springer.
- Da Silva, B. C., Konidaris, G., and Barto, A. G. (2012). Learning parameterized skills. In *Proceedings of International Conference of Machine Learning*.
- Deisenroth, M. P., Englert, P., Peters, J., and Fox, D. (2014). Multi-task policy search for robotics. In *IEEE International Conference on Robotics and Automation, 2014*, pages 3876–3881.
- Dhariwal, P., Hesse, C., Klimov, O., Nichol, A., Plappert, M., Radford, A., Schulman, J., Sidor, S., Wu, Y., and Zhokhov, P. (2017). Openai baselines. <https://github.com/openai/baselines>.
- Duan, Y., Chen, X., Houthoofd, R., Schulman, J., and Abbeel, P. (2016). Benchmarking deep reinforcement learning for continuous control. In *Proceedings of International Conference on Machine Learning*, pages 1329–1338.

- Fabisch, A. and Metzen, J. H. (2014). Active contextual policy search. *The Journal of Machine Learning Research*, 15(1):3371–3399.
- Florensa, C., Held, D., Wulfmeier, M., Zhang, M., and Abbeel, P. (2017). Reverse curriculum generation for reinforcement learning. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 482–495.
- Ghosh, D., Singh, A., Rajeswaran, A., Kumar, V., and Levine, S. (2018). Divide-and-conquer reinforcement learning. In *International Conference on Learning Representations*.
- Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feed-forward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256.
- Greensmith, E., Bartlett, P. L., and Baxter, J. (2004). Variance reduction techniques for gradient estimates in reinforcement learning. *Journal of Machine Learning Research*, 5(Nov):1471–1530.
- Jie, T. and Abbeel, P. (2010). On a connection between importance sampling and the likelihood ratio policy gradient. In *Advances in Neural Information Processing Systems*, pages 1000–1008.
- Karkus, P., Kupcsik, A., Hsu, D., and Lee, W. S. (2016). Factored contextual policy search with bayesian optimization. *arXiv preprint arXiv:1612.01746*.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations*.

- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al. (2017). Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526.
- Kober, J., Wilhelm, A., Oztop, E., and Peters, J. (2012). Reinforcement learning to adjust parametrized motor primitives to new situations. *Autonomous Robots*, 33(4):361–379.
- Kulkarni, T. D., Narasimhan, K., Saeedi, A., and Tenenbaum, J. (2016). Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pages 3675–3683.
- Kupcsik, A. G., Deisenroth, M. P., Peters, J., and Neumann, G. (2013). Data-efficient generalization of robot skills with contextual policy search. In *Proceedings of the 27th AAAI Conference on Artificial Intelligence, AAAI 2013*, pages 1401–1407.
- Levy, A., Platt, R., and Saenko, K. (2019). Hierarchical reinforcement learning with hindsight. In *International Conference on Learning Representations*.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. (2016). Continuous control with deep reinforcement learning. *ICLR*.
- Lin, L. (1992). Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3/4):69–97.
- Mankowitz, D. J., Židek, A., Barreto, A., Horgan, D., Hessel, M., Quan, J., Oh, J., van

- Hasselt, H., Silver, D., and Schaul, T. (2018). Unicorn: Continual learning with a universal, off-policy agent. *arXiv preprint arXiv:1802.08294*.
- McCloskey, M. and Cohen, N. J. (1989). Catastrophic interference in connectionist networks: The sequential learning problem. *Psychology of Learning and Motivation-Advances in Research and Theory*, 24(C):109–165.
- Metzen, J. H., Fabisch, A., and Hansen, J. (2015). Bayesian optimization for contextual policy search. In *Proceedings of the Second Machine Learning in Planning and Control of Robot Motion Workshop., Hamburg*.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540):529.
- Munos, R., Stepleton, T., Harutyunyan, A., and Bellemare, M. (2016). Safe and efficient off-policy reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 1054–1062.
- Nair, V. and Hinton, G. E. (2010). Rectified linear units improve restricted boltzmann machines. In *ICML*.
- Ng, A. Y., Harada, D., and Russell, S. (1999). Policy invariance under reward transformations: Theory and application to reward shaping. In *International Conference on Machine Learning*, volume 99, pages 278–287.
- Oh, J., Singh, S., Lee, H., and Kohli, P. (2017). Zero-shot task generalization with

- multi-task deep reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning*.
- Pathak, D., Mahmoudieh, P., Luo, M., Agrawal, P., Chen, D., Shentu, F., Shelhamer, E., Malik, J., Efros, A. A., and Darrell, T. (2018). Zero-shot visual imitation. In *International Conference on Learning Representations*.
- Peshkin, L. and Shelton, C. R. (2002). Learning from scarce experience. In *Proceedings of the Nineteenth International Conference on Machine Learning*, pages 498–505.
- Peters, J. and Schaal, S. (2008). Reinforcement learning of motor skills with policy gradients. *Neural networks*, 21(4):682–697.
- Pinsler, R., Karkus, P., Kupcsik, A., Hsu, D., and Lee, W. S. (2019). Factored contextual policy search with bayesian optimization. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 7242–7248. IEEE.
- Plappert, M., Andrychowicz, M., Ray, A., McGrew, B., Baker, B., Powell, G., Schneider, J., Tobin, J., Chociej, M., Welinder, P., et al. (2018). Multi-goal reinforcement learning: Challenging robotics environments and request for research. *arXiv preprint arXiv:1802.09464*.
- Precup, D., Sutton, R. S., and Singh, S. P. (2000). Eligibility traces for off-policy policy evaluation. In *International Conference on Machine Learning*, pages 759–766.
- Schaul, T., Horgan, D., Gregor, K., and Silver, D. (2015). Universal value function approximators. In *Proceedings of the International Conference on Machine Learning*, pages 1312–1320.

- Schmidhuber, J. (2013). POWERPLAY: Training an increasingly general problem solver by continually searching for the simplest still unsolvable problem. *Frontiers in Psychology*. (Based on arXiv:1112.5309v1 [cs.AI], 2011).
- Schmidhuber, J. (2019). Reinforcement learning upside down: Don't predict rewards—just map them to actions. *arXiv preprint arXiv:1912.02875*.
- Schmidhuber, J. and Huber, R. (1990). *Learning to Generate Focus Trajectories for Attentive Vision*. Institut für Informatik.
- Sen, P. and Singer, J. (1994). *Large Sample Methods in Statistics: An Introduction with Applications*. Chapman & Hall/CRC Texts in Statistical Science. Taylor & Francis.
- Srivastava, R. K., Steunebrink, B. R., and Schmidhuber, J. (2013). First experiments with PowerPlay. *Neural Networks*, 41(0):130 – 136. Special Issue on Autonomous Learning.
- Sukhbaatar, S., Lin, Z., Kostrikov, I., Synnaeve, G., Szlam, A., and Fergus, R. (2018). Intrinsic motivation and automatic curricula via asymmetric self-play. In *International Conference on Learning Representations*.
- Sutton, R. S. and Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. Bradford Book.
- Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. (1999a). Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems 12*, pages 1057–1063.

- Sutton, R. S., Precup, D., and Singh, S. (1999b). Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211.
- Thomas, P. (2015). *Safe reinforcement learning*. PhD thesis, University of Massachusetts Amherst.
- Thomas, P., Theocharous, G., and Ghavamzadeh, M. (2015). High confidence policy improvement. In *International Conference on Machine Learning*, pages 2380–2388.
- Vezhnevets, A. S., Osindero, S., Schaul, T., Heess, N., Jaderberg, M., Silver, D., and Kavukcuoglu, K. (2017). FeUdal networks for hierarchical reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3540–3549.
- Williams, R. J. (1986). Reinforcement-learning in connectionist networks: A mathematical analysis. Technical Report 8605, Institute for Cognitive Science, University of California, San Diego.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256.
- Zhu, Y., Mottaghi, R., Kolve, E., Lim, J. J., Gupta, A., Fei-Fei, L., and Farhadi, A. (2017). Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *IEEE International Conference on Robotics and Automation*, pages 3357–3364.