

## Supplementary Materials

### Table of Contents

<b>A</b>	<b>Related work</b>	<b>2</b>
<b>B</b>	<b>Undiscounted infinite-horizon MDPs</b>	<b>3</b>
<b>C</b>	<b>Theorems, proofs and additional information for Section 3</b>	<b>3</b>
C.1	Contribution coefficients, hindsight distribution and graphical models . . . . .	3
C.2	Proof Theorem 1 . . . . .	4
C.3	Different policy gradient estimators leveraging contribution coefficients. . . . .	6
C.4	Proof of Proposition 2 . . . . .	7
C.5	Proof Theorem 3 . . . . .	7
C.6	Proof of Theorem 4 . . . . .	9
<b>D</b>	<b>Learning the contribution coefficients</b>	<b>9</b>
D.1	Proof proposition 5 . . . . .	9
D.2	Learning contribution coefficients via contrastive classification. . . . .	10
D.3	Successor representations. . . . .	11
<b>E</b>	<b>Experimental details and additional results</b>	<b>11</b>
E.1	Dynamic programming setup . . . . .	11
E.2	Bias, variance and SNR metrics . . . . .	13
E.3	Linear key-to-door environment setup . . . . .	14
E.4	Reward switching setup . . . . .	14
E.5	Tree environment setup . . . . .	14
E.6	Training details . . . . .	15
E.7	Additional results . . . . .	17
<b>F</b>	<b>Contribution analysis in continuous spaces and POMDPs</b>	<b>19</b>
F.1	Spurious contributions in continuous state spaces . . . . .	19
F.2	Deterministic continuous reward functions can lead to excessive variance . . . . .	19
F.3	Smoothing can alleviate excess variance trading variance for bias . . . . .	20
F.4	Continuous action spaces . . . . .	21
F.5	Partially Observable MDPs . . . . .	22
<b>G</b>	<b>Learning contribution coefficients from non-rewarding observations</b>	<b>23</b>
G.1	Latent learning . . . . .	23
G.2	Optimal rewarding outcome encodings for credit assignment. . . . .	23
G.3	Counterfactual reasoning on rewarding states . . . . .	23
G.4	Learning credit assignment representations with an information bottleneck . . . . .	24
<b>H</b>	<b>Contribution analysis and causality</b>	<b>25</b>
H.1	Causal interpretation of COCOA . . . . .	25
H.2	Extending COCOA to counterfactual interventions . . . . .	26
<b>I</b>	<b>Contribution analysis with temporal discounting</b>	<b>26</b>
<b>J</b>	<b>Bootstrapping with COCOA</b>	<b>27</b>
<b>K</b>	<b>Additional details</b>	<b>30</b>
K.1	Compute resources . . . . .	30
K.2	Software and libraries . . . . .	30

## A Related work

Our work builds upon Hindsight Credit Assignment (HCA) [1] which has sparked a number of follow-up studies. We generalize HCA towards estimating contributions upon rewarding outcomes instead of rewarding states and show through a detailed variance analysis that HCA suffers from spurious contributions leading to high variance, while using rewards or rewarding objects as rewarding outcome encodings leads to low-variance gradient estimators. Follow-up work on HCA reduces the variance of the HCA advantage estimates by combining it with Monte Carlo estimates [36] or using temporal difference errors instead of rewards [37]. Alipov et al. [27] leverages the latter approach to scale up HCA towards more complex environments. However, all of the above approaches still suffer from spurious contributions, a significant source of variance in the HCA gradient estimator. In addition, recent studies have theoretically reinterpreted the original HCA formulation from different angles: Ma and Pierre-Luc [38] create the link to Conditional Monte Carlo methods [39, 40] and Arumugam et al. [41] provide an information theoretic perspective on credit assignment. Moreover, Young [42] applies HCA to the problem of estimating gradients in neural networks with stochastic, discrete units.

Long-term credit assignment in RL is hard due to the high variance of the sum of future rewards for long trajectories. A common technique to reduce the resulting variance in the policy gradients is to subtract a baseline from the sum of future rewards [6, 7, 43, 44]. To further reduce the variance, a line of work introduced state-action-dependent baselines [45–48]. However, Tucker et al. [49] argues that these methods offer only a small benefit over the conventional state-dependent baselines, as the current action often only accounts for a minor fraction of the total variance. More recent work proposes improved baselines by incorporating hindsight information about the future trajectory into the baseline, accounting for a larger portion of the variance [1, 3, 50–52]. Mesnard et al. [3] learn a summary metric of the uncontrollable, external environment influences in the future trajectory, and provide this hindsight information as an extra input to the value baseline. Nota et al. [50] consider partially observable MDPs, and leverage the future trajectory to more accurately infer the current underlying Markov state, thereby providing a better value baseline. Harutyunyan et al. [1] propose return-HCA, a different variant of HCA that uses a return-conditioned hindsight distribution to construct a baseline, instead of using state-based hindsight distributions for estimating contributions. Finally, Guez et al. [51] and Venuto et al. [52] learn a summary representation of the full future trajectory and provide it as input to the value function, while imposing an information bottleneck to prevent the value function from overly relying on this hindsight information.

Environments with sparse rewards and long delays between actions and corresponding rewards put a high importance on long-term credit assignment. A popular strategy to circumvent the sparse and delayed reward setting is to introduce reward shaping [53–60]. These approaches add auxiliary rewards to the sparse reward function, aiming to guide the learning of the policy with dense rewards. A recent line of work introduces a reward shaping strategy specifically designed for long-term credit assignment, where rewards are decomposed and distributed over previous state-action pairs that were instrumental in achieving that reward [4, 5, 30, 32, 61–66]. To determine how to redistribute rewards, these approaches rely on heuristic contribution analyses, such as via the access of memory states [4], linear decompositions of rewards [30, 62–66] or learned sequence models [5, 32, 61]. Leveraging our unbiased contribution analysis framework to reach more optimal reward transport is a promising direction for future research.

When we have access to a (learned) differentiable world model of the environment, we can achieve precise credit assignment by leveraging path-wise derivatives, i.e. backpropagating value gradients through the world model [18–22, 67]. For stochastic world models, we need access to the noise variables to compute the path-wise derivatives. The Dreamer algorithms [19–21] approach this by computing the value gradients on *simulated* trajectories, where the noise is known. The Stochastic Value Gradient (SVG) method [18] instead *infers* the noise variables on real *observed* trajectories. To enable backpropagating gradients over long time spans, Ma et al. [68] equip the learned recurrent world models of SVG with an attention mechanism, allowing the authors to leverage Sparse Attentive Backtracking [69] to transmit gradients through skip connections. Buesing et al. [34] leverages the insights from SVG in partially observable MDPs, using the inferred noise variables to estimate the effect of counterfactual policies on the expected return. Importantly, the path-wise derivatives leveraged by the above model-based credit assignment methods are not compatible with discrete action spaces, as sensitivities w.r.t. discrete actions are undefined. In contrast, COCOA can leverage model-based information for credit assignment, while being compatible with discrete actions.

Incorporating hindsight information has a wide variety of applications in RL. Goal-conditioned policies [70] use a goal state as additional input to the policy network or value function, thereby generalizing it to arbitrary goals. Hindsight Experience Replay [71] leverages hindsight reasoning to learn almost as much from undesired outcomes as from desired ones, as at hindsight, we can consider every final, possibly undesired state as a ‘goal’ state and update the value functions or policy network [72] accordingly. Goyal et al. [73] train an inverse environment model to simulate alternative past trajectories leading to the same rewarding state, hence leveraging hindsight reasoning to create a variety of highly rewarding trajectories. A recent line of work frames RL as a supervised sequence prediction problem, learning a policy conditioned on goal states or future returns [74–77]. These models are trained on past trajectories or offline data, where we have in hindsight access to states and returns, considering them as targets for the learned policy.

Finally, Temporal Difference (TD) learning [78] has a rich history of leveraging proximity in time as a proxy for credit assignment [2].  $TD(\lambda)$  [78] considers eligibility traces to trade off bias and variance, crediting past state-action pairs in the trajectory for the current reward proportional to how close they are in time. van Hasselt et al. [79] estimate expected eligibility traces, taking into account that the same rewarding state can be reached from various previous states. Hence, not only the past state-actions on the trajectory are credited and updated, but also counterfactual ones that lead to the same rewarding state. Extending the insights from COCOA towards temporal difference methods is an exciting direction for future research.

## B Undiscounted infinite-horizon MDPs

In this work, we consider an undiscounted MDP with a finite state space  $\mathcal{S}$ , bounded rewards and an infinite horizon. To ensure that the expected return and value functions remain finite, we require some standard regularity conditions [2]. We assume the MDP contains an absorbing state  $s_\infty$  that transitions only to itself and has zero reward. Moreover, we assume *proper* transition dynamics, meaning that an agent following any policy will eventually end up in the absorbing state  $s_\infty$  with probability one as time goes to infinity.

The discounted infinite horizon MDP formulation is a special case of this setting, with a specific class of transition dynamics. An explicit discounting of future rewards  $\sum_{k \geq 0} \gamma^k R_k$  with discount factor  $\gamma \in [0, 1]$ , is equivalent to considering the above undiscounted MDP setting, but modifying the state transition probability function  $p(S_{t+1} | S_t, A_t)$  such that each state-action pair has a fixed probability  $(1 - \gamma)$  of transitioning to the absorbing state [2]. Hence, all the results considered in this work can be readily applied to the discounted MDP setting, by modifying the environment transitions as outlined above. We can also explicitly incorporate temporal discounting in the policy gradients and contribution coefficients, which we outline in Appendix I.

The undiscounted infinite-horizon MDP with an absorbing state can also model episodic RL with a fixed time horizon. We can include time as an extra feature in the states  $S$ , and then have a probability of 1 to transition to the absorbing state when the agent reaches the final time step.

## C Theorems, proofs and additional information for Section 3

### C.1 Contribution coefficients, hindsight distribution and graphical models

Here, we provide more information on the derivation of the contribution coefficients of Eq. 3.

**Abstracting time.** In an undiscounted environment, it does not matter at which point in time the agent achieves the rewarding outcome  $u'$ . Hence, the contribution coefficients (3) sum over all future time steps, to reflect that the rewarding outcome  $u'$  can be encountered at any future time, and to incorporate the possibility of encountering  $u'$  multiple times. Note that when using temporal discounting, we can adjust the contribution coefficients accordingly (c.f. App I).

**Hindsight distribution.** To obtain the hindsight distribution  $p(A_t = a | S_t = s, U' = u')$ , we convert the classical graphical model of an MDP (c.f. Fig. 5a) into a graphical model that incorporates the time  $k$  into a separate node (c.f. Fig. 5b). Here, we rewrite  $p(U_k = u' | S = s, A = a)$ , the probability distribution of a rewarding outcome  $U_k$   $k$  time steps later, as  $p(U' = u' | S = s, A = a, K = k)$ . By giving  $K$  the geometric distribution  $p(K = k) = (1 - \beta)\beta^{k-1}$  for some  $\beta$ , we can rewrite the infinite sums used in the time-independent contribution coefficients as a marginalization

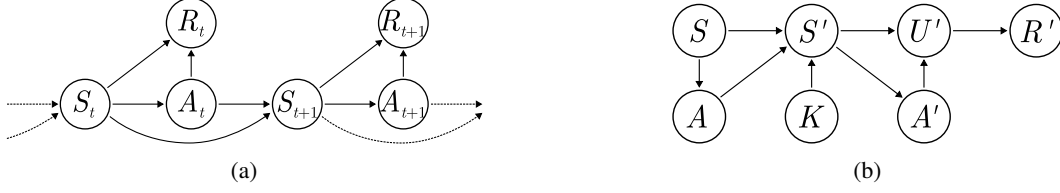


Figure 5: (a) Graphical model of the MDP. (b) Graphical model of the MDP, where we abstracted time.

534 over  $K$ :

$$\sum_{k \geq 1} p(U_k = s' \mid S = s, A = a) = \lim_{\beta \rightarrow 1} \sum_{k \geq 1} p(U_k = u' \mid S = s, A = a) \beta^{k-1} \quad (8)$$

$$= \lim_{\beta \rightarrow 1} \frac{1}{1 - \beta} \sum_{k \geq 1} p(U' = u' \mid S = s, A = a, K = k) p(K = k) \quad (9)$$

$$= \lim_{\beta \rightarrow 1} \frac{1}{1 - \beta} p(U' = u' \mid S = s, A = a) \quad (10)$$

535 Note that this limit is finite, as for  $k \rightarrow \infty$ , the probability of reaching an absorbing state  $s_\infty$  and  
 536 corresponding rewarding outcome  $u_\infty$  goes to 1 (and we take  $u' \neq u_\infty$ ). Via Bayes rule, we then  
 537 have that

$$w(s, a, u') = \frac{\sum_{k \geq 1} p(U_k = u' \mid S = s, A = a)}{\sum_{k \geq 1} p(U_k = u' \mid S = s)} - 1 \quad (11)$$

$$= \frac{p(U' = u' \mid S = s, A = a)}{p(U' = u \mid S = s)} - 1 = \frac{p(A = a \mid S = s, U' = u')}{\pi(a \mid s)} - 1 \quad (12)$$

538 where we assume  $\pi(a \mid s) > 0$  and where we dropped the limit of  $\beta \rightarrow 1$  in the notation, as we will  
 539 always take this limit henceforth. Note that when  $\pi(a \mid s) = 0$ , the hindsight distribution  $p^\pi(a \mid s, u')$   
 540 is also equal to zero, and hence the right-hand-side of the above equation is undefined. However,  
 541 the middle term using  $p^\pi(U_k = u' \mid S = s, A = a)$  is still well-defined, and hence we can use the  
 542 contribution coefficients even for actions where  $\pi(a \mid s) = 0$ .

## 543 C.2 Proof Theorem 1

544 We start with a lemma showing that the contribution coefficients can be used to estimate the advantage  
 545  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ , by generalizing Theorem 1 of Harutyunyan et al. [1] towards rewarding  
 546 outcomes.

547 **Definition 2** (Fully predictive, repeated from main text). A rewarding outcome encoding  $U$  is fully  
 548 predictive of the reward  $R$ , if the following conditional independence condition holds:  $p^\pi(R_k = r \mid$   
 549  $S_0 = s, A_0 = a, U_k = u) = p^\pi(R = r \mid U = u)$ , where the right-hand side does not depend on the  
 550 time  $k$ .

551 **Lemma 6.** For each state-action pair  $(s, a)$  with  $\pi(a \mid s) > 0$ , and assuming that  $u = f(s, a, r)$  is  
 552 fully predictive of the reward (c.f. Definition 2), we have that

$$A^\pi(s, a) = r(s, a) - \sum_{a' \in \mathcal{A}} \pi(a' \mid s) r(s, a') + \mathbb{E}_{T \sim \mathcal{T}(s, \pi)} \left[ \sum_{k \geq 1} w(s, a, U_k) R_k \right] \quad (13)$$

553 with the advantage function  $A^\pi$ , and the reward function  $r(s, a) \triangleq \mathbb{E}[R \mid s, a]$ .

554 *Proof.* We rewrite the undiscounted state-action value function in the limit of a discounting factor  
 555  $\beta \rightarrow 1_-$  (the minus sign indicating that we approach 1 from the left):

$$Q(s, a) = \lim_{\beta \rightarrow 1_-} \mathbb{E}_{T \sim \mathcal{T}(s, a, \pi)} \left[ \sum_{k \geq 1} \beta^k R_k \right] \quad (14)$$

$$= r(s, a) + \lim_{\beta \rightarrow 1_-} \sum_{r \in \mathcal{R}} \sum_{k \geq 1} \beta^k p^\pi(R_k = r \mid s, a) r \quad (15)$$

$$= r(s, a) + \lim_{\beta \rightarrow 1_-} \sum_{r \in \mathcal{R}} \sum_{u \in \mathcal{U}} \sum_{k \geq 1} \beta^k p^\pi(R_k = r, U_k = u \mid s, a) r \quad (16)$$

$$= r(s, a) + \lim_{\beta \rightarrow 1_-} \sum_{r \in \mathcal{R}} \sum_{u \in \mathcal{U}} \sum_{k \geq 1} \beta^k p^\pi(R = r \mid U = u) r p^\pi(U_k = u \mid s, a) \quad (17)$$

$$= r(s, a) + \lim_{\beta \rightarrow 1_-} \sum_{u \in \mathcal{U}} r(u) \sum_{k \geq 1} \beta^k p^\pi(U_k = u \mid s, a) \quad (18)$$

$$= r(s, a) + \lim_{\beta \rightarrow 1_-} \sum_{u \in \mathcal{U}} r(u) \sum_{k \geq 1} \beta^k p^\pi(U_k = u \mid s) \frac{\sum_{k' \geq 1} \beta^{k'} p^\pi(U_{k'} = u \mid s, a)}{\sum_{k' \geq 1} \beta^{k'} p^\pi(U_{k'} = u \mid s)} \quad (19)$$

556 where we use the property that  $u$  is *fully predictive* of the reward (c.f. Definition 1), and where we  
 557 define  $r(s, a) \triangleq \mathbb{E}[R \mid s, a]$  and  $r(u) \triangleq \mathbb{E}[R \mid u]$ . Using the graphical model of Fig. 5b where we  
 558 abstract time  $k$  (c.f. App. C.1), we have that  $\frac{1}{1-\beta} p_\beta^\pi(R = r' \mid s, a) = \sum_{k=0} \beta^k p^\pi(R_k = r \mid s, a)$ .  
 559 Leveraging Bayes rule, we get

$$p^\pi(a \mid s, U' = u) = \lim_{\beta \rightarrow 1_-} \frac{p_\beta^\pi(U' = u \mid s, a) \pi(a \mid s)}{p_\beta^\pi(U' = u \mid s)} = \lim_{\beta \rightarrow 1_-} \frac{\sum_{k \geq 1} \beta^k p^\pi(U_k = u \mid s, a)}{\sum_{k \geq 1} \beta^k p^\pi(U_k = u \mid s)} \pi(a \mid s) \quad (20)$$

560 Where we dropped the limit of  $\beta \rightarrow 1_-$  in the notation of  $p^\pi(a \mid s, R = r)$ , as we henceforth always  
 561 consider this limit. Taking everything together, we have that

$$Q^\pi(s, a) = r(s, a) + \sum_{u \in \mathcal{U}} \sum_{k \geq 1} p^\pi(U_k = u \mid s) \frac{p^\pi(a \mid s, U' = u)}{\pi(a \mid s)} r(u) \quad (21)$$

$$= r(s, a) + \sum_{u \in \mathcal{U}} \sum_{k \geq 1} p^\pi(U_k = u \mid s) (w(s, a, u) + 1) r(u) \quad (22)$$

$$= r(s, a) + \sum_{u \in \mathcal{U}} \sum_{k \geq 1} p^\pi(U_k = u \mid s) (w(s, a, u) + 1) \sum_{r \in \mathcal{R}} p(R_k = r \mid U_k = u) r \quad (23)$$

$$= r(s, a) + \mathbb{E}_{T \sim \mathcal{T}(s, \pi)} \left[ \sum_{k \geq 1} (w(s, a, U_k) + 1) R_k \right] \quad (24)$$

562 Subtracting the value function, we get

$$A^\pi(s, a) = r(s, a) - \sum_{a' \in \mathcal{A}} \pi(a' \mid s) r(s, a') + \mathbb{E}_{T \sim \mathcal{T}(s, \pi)} \left[ \sum_{k \geq 1} w(s, a, U_k) R_k \right] \quad (25)$$

563 □

564 Now we are ready to prove Theorem 1.

565 *Proof.* Using the policy gradient theorem [24], we have

$$\nabla_{\theta} V^{\pi}(s_0) = \mathbb{E}_{T \sim \mathcal{T}(s_0, \pi)} \left[ \sum_{t \geq 0} \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a | S_t) A^{\pi}(S_t, a) \right] \quad (26)$$

$$= \mathbb{E}_{T \sim \mathcal{T}(s_0, \pi)} \left[ \sum_{t \geq 0} \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a | S_t) \left[ r(S_t, a) + \sum_{k \geq 1} w(S_t, a, U_{t+k}) R_{t+k} \right] \right] \quad (27)$$

$$= \mathbb{E}_{T \sim \mathcal{T}(s_0, \pi)} \left[ \sum_{t \geq 0} \nabla_{\theta} \log \pi(A_t | S_t) R_t + \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a | S_t) \sum_{k \geq 1} w(S_t, a, U_{t+k}) R_{t+k} \right] \quad (28)$$

566 where we used that removing a baseline  $\sum_{a' \in \mathcal{A}} \pi(a' | s) r(s, a')$  independent from the actions  
 567 does not change the policy gradient, and replaced  $\mathbb{E}_{T \sim \mathcal{T}(s_0, \pi)} [\sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a | S_t) r(S_t, a)]$  by its  
 568 sampled version  $\mathbb{E}_{T \sim \mathcal{T}(s_0, \pi)} [\nabla_{\theta} \log \pi(A_t | S_t) R_t]$ . Hence the policy gradient estimator of Eq. 4 is  
 569 unbiased.  $\square$

### 570 C.3 Different policy gradient estimators leveraging contribution coefficients.

571 In this work, we use the COCOA estimator of Eq. 4 to estimate policy gradients leveraging the  
 572 contribution coefficients of Eq. 3, as this estimator does not need a separate reward model, and it  
 573 works well for the small action spaces we considered in our experiments. However, one can design  
 574 other unbiased policy gradient estimators compatible with the same contribution coefficients.

575 Harutyunyan et al. [1] introduced the HCA gradient estimator of Eq. 2, which can readily be extended  
 576 to rewarding outcomes  $U$ :

$$\hat{\nabla}_{\theta} V^{\pi} = \sum_{t \geq 0} \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi(a | S_t) \left( r(S_t, a) + \sum_{k \geq 1} w(S_t, a, U_{t+k}) R_{t+k} \right) \quad (29)$$

577 This gradient estimator uses a reward model  $r(s, a)$  to obtain the rewards corresponding to counter-  
 578 factual actions, whereas the COCOA estimator (4) only uses the observed rewards  $R_k$ .

579 For large action spaces, it might become computationally intractable to sum over all possible actions.  
 580 In this case, we can sample independent actions from the policy, instead of summing over all actions,  
 581 leading to the following policy gradient.

$$\hat{\nabla}_{\theta} V^{\pi} = \sum_{t \geq 0} \nabla_{\theta} \log \pi(A_t | S_t) R_t + \frac{1}{M} \sum_m \nabla_{\theta} \log \pi(a^m | S_t) \sum_{k=1}^{\infty} w(S_t, a^m, U_{t+k}) R_{t+k} \quad (30)$$

582 where we sample  $M$  actions from  $a^m \sim \pi(\cdot | S_t)$ . Importantly, for obtaining an unbiased policy  
 583 gradient estimate, the actions  $a^m$  should be sampled independently from the actions used in the  
 584 trajectory of the observed rewards  $R_{t+k}$ . Hence, we cannot take the observed  $A_t$  as a sample  $a^m$ , but  
 585 need to instead sample independently from the policy. This policy gradient estimator can be used for  
 586 continuous action spaces (c.f. App. F, and can also be combined with the reward model as in Eq.  
 587 29.

588 One can show that the above policy gradient estimators are unbiased with a proof akin to the one of  
 589 Theorem 1.

590 **Comparison of COCOA to using time as a credit assignment heuristic.** In Section 2 we discussed  
 591 briefly the following discounted policy gradient:

$$\hat{\nabla}_{\theta}^{\text{REINFORCE}, \gamma} V^{\pi}(s_0) = \sum_{t \geq 0} \nabla_{\theta} \log \pi(A_t | S_t) \sum_{k \geq t} \gamma^{k-t} R_k \quad (31)$$

592 with discount factor  $\gamma \in [0, 1]$ . Note that we specifically put no discounting  $\gamma^t$  in front of  
 593  $\nabla_{\theta} \log \pi(A_t | S_t)$ , which would be required for being an unbiased gradient estimate of the dis-  
 594 counted expected total return, as the above formulation is most used in practice [10, 25]. Rewriting  
 595 the summations reveals that this policy gradient uses time as a heuristic for credit assignment:

$$\hat{\nabla}_{\theta}^{\text{REINFORCE}, \gamma} V^{\pi}(s_0) = \sum_{t \geq 0} R_t \sum_{k \leq t} \gamma^{t-k} \nabla_{\theta} \log \pi(A_k | S_k). \quad (32)$$

596 We can rewrite the COCOA gradient estimator (4) with the same reordering of summation, showcasing  
 597 that it leverages the contribution coefficient for providing precise credit to past actions, instead of  
 598 using the time discounting heuristic:

$$\hat{\nabla}_\theta^U V^\pi(s_0) = \sum_{t \geq 0} R_t \left[ \nabla_\theta \log \pi(A_t | S_t) + \sum_{k < t} \sum_{a \in \mathcal{A}} w(S_k, a, U_t) \nabla_\theta \pi(a | S_k) \right] \quad (33)$$

#### 599 C.4 Proof of Proposition 2

600 *Proof.* Proposition 2 assumes that each action sequence  $\{A_m\}_{m=t}^{t+k}$  leads to a unique state  $s'$ . Hence,  
 601 all previous actions can be decoded perfectly from the state  $s'$ , leading to  $p^\pi(A_t = a | S_t, S' =$   
 602  $s') = \delta(a = a_t)$ , with  $\delta$  the indicator function and  $a_t$  the action taken in the trajectory that led to  $s'$ .  
 603 Filling this into the COCOA gradient estimator leads to

$$\hat{\nabla}_\theta^S V^\pi(s_0) = \sum_{t \geq 0} \nabla_\theta \log \pi(A_t | S_t) R_t + \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a | S_t) \sum_{k \geq 1} \left( \frac{\delta(a = A_{t+k})}{\pi(a | S_t)} - 1 \right) R_{t+k} \quad (34)$$

$$= \sum_{t \geq 0} \frac{\nabla_\theta \pi(A_t | S_t)}{\pi(A_t | S_t)} \sum_{k \geq 0} R_{t+k} \quad (35)$$

$$= \sum_{t \geq 0} \nabla_\theta \log \pi(A_t | S_t) \sum_{k \geq 0} R_{t+k} \quad (36)$$

604 where we used that  $\sum_a \nabla_\theta \pi(a | s) = 0$ . □

#### 605 C.5 Proof Theorem 3

606 *Proof.* Theorem 3 considers the case where the environment only contains a reward at the final time  
 607 step  $t = T$ , and where we optimize the policy only on a single (initial) time step  $t = 0$ . Then, the  
 608 policy gradients are given by

$$\hat{\nabla}_\theta^U V^\pi(U_T, R_T) = \sum_a \nabla_\theta \pi(a | s) w(s, a, U_T) R_T \quad (37)$$

$$\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi(A_0, R_T) = \nabla_\theta \log \pi(A_0 | s) R_T \quad (38)$$

609 With  $U$  either  $S$ ,  $R$  or  $Z$ , and  $s$  the state at  $t = 0$ . As only the last time step can have a reward  
 610 by assumption, and the encoding  $U$  needs to retain the predictive information of the reward, the  
 611 contribution coefficients for  $u$  corresponding to a nonzero reward are equal to

$$w(s, a, u) = \frac{p^\pi(A_0 = a | S_0 = s, U_T = u)}{\pi(a | s)} - 1 \quad (39)$$

612 The coefficients corresponding to zero-reward outcomes are multiplied with zero in the  
 613 gradient estimator, and can hence be ignored. Now, we proceed by showing that  
 614  $\mathbb{E}[\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi(A_0, R_T) | S_T, R_T] = \hat{\nabla}_\theta^S V^\pi(S_T, R_T)$ ,  $\mathbb{E}[\hat{\nabla}_\theta^S V^\pi(S_T, R_T) | U'_T, R_T] =$   
 615  $\hat{\nabla}_\theta^{U'} V^\pi(U'_T, R_T)$ ,  $\mathbb{E}[\hat{\nabla}_\theta^{U'} V^\pi(U'_T, R_T) | U_T, R_T] = \hat{\nabla}_\theta^U V^\pi(U_T, R_T)$ , and  $\mathbb{E}[\hat{\nabla}_\theta^U V^\pi(U_T, R_T) |$   
 616  $R_T] = \hat{\nabla}_\theta^R V^\pi(R_T)$ , after which we can use the law of total variance to prove our theorem.

617 As  $S_T$  is fully predictive of  $R_T$ , the following conditional independence holds  $p^\pi(A_0 | S_0 =$   
 618  $s, S_T, R_T) = p^\pi(A_0 | S_0 = s, S_T)$  (c.f. Definition 1). Hence, we have that

$$\mathbb{E}[\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi(A_0, R_T) | S_T, R_T] \quad (40)$$

$$= \sum_{a \in \mathcal{A}} p(A_0 = a | S_0 = s, S_T) \frac{\nabla_\theta \pi(a | s)}{\pi(a | s)} R_T = \hat{\nabla}_\theta^S V^\pi(S_T, R_T) \quad (41)$$

619 where we used that  $\sum_a \nabla_\theta \pi(a | s) = \nabla_\theta \sum_a \pi(a | s) = 0$ .

620 Similarly, as  $S_T$  is fully predictive of  $U'_T$  (c.f. Definition 1), we have that

$$p(A | S, U'_T) = \sum_{s_T \in \mathcal{S}} p(S_T = s_T | S, U'_T) p(A | S, S_T = s_T, U'_T) \quad (42)$$

$$= \sum_{s_T \in \mathcal{S}} p(S_T = s_T | S, U'_T) p(A | S, S_T = s_T) \quad (43)$$

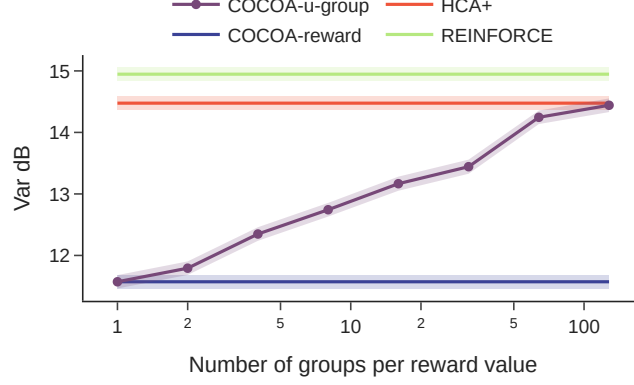


Figure 6: **Less informative rewarding outcome encodings lead to gradient estimators with lower variance.** Normalized variance in dB using ground-truth coefficients and a random uniform policy, for various gradient estimators on the tree environment (shaded region represents standard error over 10 random environments). To increase the information content of  $U$ , we increase the number of different encodings  $u$  corresponding to the same reward value (c.f.  $n_g$  in Section E.5), indicated by the x-axis. COCOA-u-group indicates the COCOA estimator with rewarding outcome encodings of increasing information content, whereas COCOA-reward and HCA+ have fixed rewarding outcome encodings of  $U = R$  and  $U = S$  respectively.

621 Using the conditional independence relation  $p^\pi(S_T | S_0, U'_T, R_T) = p^\pi(S_T | S_0, U'_T)$ , following  
 622 from d-separation in Fig. 5b, this leads us to

$$\mathbb{E}[\hat{\nabla}_\theta^S V^\pi(S_T, R_T) | U'_T, R_T] \quad (44)$$

$$= \sum_{a \in \mathcal{A}} \sum_{s_T \in \mathcal{S}} p(S_T = s_T | S_0 = s, U'_T) p(A_0 = a | S = s, S_T = s_T) \frac{\nabla_{\theta} \pi(a | s)}{\pi(a | s)} R_T \quad (45)$$

$$= \sum_{a \in \mathcal{A}} p(A_0 = a | S = s, U'_T) \frac{\nabla_{\theta} \pi(a | s)}{\pi(a | s)} R_T \quad (46)$$

$$= \hat{\nabla}_\theta^{U'} V^\pi(U'_T, R_T) \quad (47)$$

623 Using the same derivation leveraging the fully predictive properties (c.f. Definition 1), we get

$$\mathbb{E}[\hat{\nabla}_\theta^{U'} V^\pi(U'_T, R_T) | U_T, R_T] = \hat{\nabla}_\theta^U V^\pi(U_T, R_T) \quad (48)$$

$$\mathbb{E}[\hat{\nabla}_\theta^U V^\pi(U_T, R_T) | R_T] = \hat{\nabla}_\theta^R V^\pi(R_T) \quad (49)$$

624 Now we can use the law of total variance, which states that  $\mathbb{V}[X] = \mathbb{E}[\mathbb{V}[X | Y]] + \mathbb{V}[\mathbb{E}[X | Y]]$ .  
 625 Hence, we have that

$$\mathbb{V}[\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi] = \mathbb{V}[\mathbb{E}[\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi | S_T]] + \mathbb{E}[\mathbb{V}[\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi | S_T]] \quad (50)$$

$$= \mathbb{V}[\hat{\nabla}_\theta^S V^\pi] + \mathbb{E}[\mathbb{V}[\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi | S_T]] \succcurlyeq \mathbb{V}[\hat{\nabla}_\theta^S V^\pi] \quad (51)$$

626 as  $\mathbb{E}[\mathbb{V}[\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi | S_T]]$  is positive semi definite. Using the same construction for the other  
 627 pairs, we arrive at

$$\mathbb{V}[\hat{\nabla}_\theta^R V^\pi(s_0)] \preccurlyeq \mathbb{V}[\hat{\nabla}_\theta^U V^\pi(s_0)] \preccurlyeq \mathbb{V}[\hat{\nabla}_\theta^{U'} V^\pi(s_0)] \preccurlyeq \mathbb{V}[\hat{\nabla}_\theta^S V^\pi(s_0)] \preccurlyeq \mathbb{V}[\hat{\nabla}_\theta^{\text{REINFORCE}} V^\pi(s_0)] \quad (52)$$

628 thereby concluding the proof.  $\square$

629 **Additional empirical verification.** To get more insight into how the information content of the  
 630 rewarding-outcome encoding  $U$  relates to the variance of the COCOA gradient estimator, we repeat  
 631 the experiment of Fig. 2, but now plot the variance as a function of the amount of information in the  
 632 rewarding outcome encodings, for a fixed state overlap of 3. Fig. 6 shows that the variance of the



resulting COCOA gradient estimators interpolate between COCOA-reward and HCA+, with more informative encodings leading to a higher variance. These empirical results show that the insights of Theorem 3 hold in this more general setting of estimating full policy gradients in an MDP with random rewards.

## C.6 Proof of Theorem 4

*Proof.* This proof follows a similar technique to the policy gradient theorem [24]. Let us first define the expected number of occurrences of  $u$  starting from state  $s$  as  $O^\pi(u, s) = \sum_{k \geq 1} p^\pi(U_k = u \mid S_0 = s)$ , and its action equivalent  $O^\pi(u, s, a) = \sum_{k \geq 1} p^\pi(U_k = u \mid S_0 = s, A_0 = a)$ . Expanding  $O^\pi(u, s)$  leads to

$$\nabla_\theta O^\pi(u, s) = \nabla_\theta \left[ \sum_{a \in \mathcal{A}} \pi(a \mid s) \sum_{s' \in \mathcal{S}} p(s' \mid s, a) (\delta(s' = s) + O^\pi(u, s')) \right] \quad (53)$$

$$= \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a \mid s) O^\pi(u, s, a) + \sum_{s' \in \mathcal{S}} p(s' \mid s) \nabla_\theta O^\pi(u, s') \quad (54)$$

Now define  $\phi(u, s) = \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a \mid s) O^\pi(u, s, a)$ , and  $p^\pi(S_l = s' \mid S_0 = s)$  as the probability of reaching  $s'$  starting from  $s$  in  $l$  steps. We have that  $\sum_{s'' \in \mathcal{S}} p^\pi(S_l = s'' \mid S_0 = s) p^\pi(S_1 = s' \mid S_0 = s'') = p^\pi(S_{l+1} = s' \mid S_0 = s)$ . Leveraging this relation and recursively applying the above equation leads to

$$\nabla_\theta O^\pi(u, s) = \sum_{s' \in \mathcal{S}} \sum_{l=0}^{\infty} p^\pi(S_l = s' \mid S_0 = s) \phi(u, s') \quad (55)$$

$$= \sum_{s' \in \mathcal{S}} \sum_{l=0}^{\infty} p^\pi(S_l = s' \mid S_0 = s) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a \mid s') O^\pi(u, s', a) \quad (56)$$

$$\propto \mathbb{E}_{S \sim \mathcal{T}(s, \pi)} \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a \mid S) O^\pi(u, S, a) \quad (57)$$

where in the last step we normalized  $\sum_{l=0}^{\infty} p^\pi(S_l = s' \mid S_0 = s)$  with  $\sum_{s' \in \mathcal{S}} \sum_{l=0}^{\infty} p^\pi(S_l = s' \mid S_0 = s)$ , resulting in the state distribution  $S \sim \mathcal{T}(s, \pi)$  where  $S$  is sampled from trajectories starting from  $s$  and following policy  $\pi$ .

Finally, we can rewrite the contribution coefficients as

$$w(s, a, u) = \frac{O^\pi(u, s, a)}{O^\pi(u, s)} - 1 \quad (58)$$

which leads to

$$\nabla_\theta O^\pi(u, s) \propto \mathbb{E}_{S \sim \mathcal{T}(s, \pi)} \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a \mid S) w(S, a, u) O^\pi(u, S) \quad (59)$$

where we used that  $\sum_{a \in \mathcal{A}} \nabla_\theta \pi(a \mid s) = 0$ , thereby concluding the proof.  $\square$

## D Learning the contribution coefficients

### D.1 Proof proposition 5

*Proof.* In short, as the logits  $l$  are a deterministic function of the state  $s$ , they do not provide any further information on the hindsight distribution and we have that  $p^\pi(A_0 = a \mid S_0 = s, L_0 = l, U' = u') = p^\pi(A_0 = a \mid S_0 = s, U' = u')$ .

657 We can arrive at this result by observing that  $p(A_0 = a \mid S_0 = s, L_0 = l) = p(A_0 = a \mid S_0 = s) =$   
 658  $\pi(a \mid s)$ , as the policy logits encode the policy distribution. Using this, we have that

$$p^\pi(A_0 = a \mid S_0 = s, L_0 = l, U' = u') \quad (60)$$

$$= \frac{p^\pi(A_0 = a, U' = u' \mid S_0 = s, L_0 = l)}{p^\pi(U' = u' \mid S_0 = s, L_0 = l)} \quad (61)$$

$$= \frac{p^\pi(U' = u' \mid S_0 = s, A_0 = a, L_0 = l)p(A_0 = a \mid S_0 = s, L_0 = l)}{\sum_{a' \in \mathcal{A}} p^\pi(A_0 = a', U' = u' \mid S_0 = s, L_0 = l)} \quad (62)$$

$$= \frac{p^\pi(U' = u' \mid S_0 = s, A_0 = a)p(A_0 = a \mid S_0 = s)}{\sum_{a' \in \mathcal{A}} p^\pi(U' = u' \mid S_0 = s, A_0 = a')p(A_0 = a' \mid S_0 = s)} \quad (63)$$

$$= \frac{p^\pi(A_0 = a, U' = u' \mid S_0 = s)}{p^\pi(U' = u' \mid S_0 = s)} \quad (64)$$

$$= p^\pi(A_0 = a \mid S_0 = s, U' = u') \quad (65)$$

659  $\square$

## 660 D.2 Learning contribution coefficients via contrastive classification.

661 Instead of learning the hindsight distribution, we can estimate the probability ratio  $p^\pi(A_t = a \mid S_t =$   
 662  $s, U' = u')/\pi(a \mid s)$  directly, by leveraging contrastive classification. Proposition 7 shows that by  
 663 training a binary classifier  $D(a, s, u')$  to distinguish actions sampled from  $p^\pi(A_t = a \mid S_t = s, R' =$   
 664  $r')$  versus  $\pi(a \mid s)$ , we can directly estimate the probability ratio.

665 **Proposition 7.** *Consider the contrastive loss*

$$L = \mathbb{E}_{s, u' \sim \mathcal{T}(s_0, \pi)} [\mathbb{E}_{a \sim p^\pi(a \mid s, u')} [\log D(a, s, u')] + \mathbb{E}_{a \sim \pi(a \mid s)} [\log (1 - D(a, s, u'))]] , \quad (66)$$

666 *and  $D^*(s, a, u')$  its minimizer. Then the following holds*

$$w(s, a, u') = \frac{D^*(a, s, u')}{1 - D^*(a, s, u')} - 1. \quad (67)$$

667 *Proof.* Consider a fixed pair  $s, u'$ . We can obtain the discriminator  $D^*(a, s, u')$  that maximizes  
 668  $\mathbb{E}_{a \sim p^\pi(a \mid s, u')} [\log D(a, s, u')] + \mathbb{E}_{a \sim \pi(a \mid s)} [\log (1 - D(a, s, u'))]$  by taking the point-wise derivative  
 669 of this objective and equating it to zero:

$$p^\pi(a \mid s, u') \frac{1}{D^*(a, s, u')} - \pi(a \mid s) \frac{1}{1 - D^*(a, s, u')} = 0 \quad (68)$$

$$\Rightarrow \frac{p^\pi(a \mid s, u')}{\pi(a \mid s)} = \frac{D^*(a, s, u')}{1 - D^*(a, s, u')} \quad (69)$$

670 As for any  $(a, b) \in \mathbb{R}_{>0}^2$ , the function  $f(x) = a \log x + b \log(1 - x)$  achieves its global maximum on  
 671 support  $x \in [0, 1]$  at  $\frac{a}{a+b}$ , the above maximum is a global maximum. Repeating this argument for all  
 672  $(s, u')$  pairs concludes the proof.  $\square$

673 We can approximate the training objective of  $D$  by sampling  $a^{(m)}, s^{(m)}, u'^{(m)}$  along the observed  
 674 trajectories, while leveraging that we have access to the policy  $\pi$ , leading to the following loss

$$L = \sum_{m=1}^M \left[ -\log D(a^{(m)}, s^{(m)}, u'^{(m)}) - \sum_{a' \in \mathcal{A}} \pi(a' \mid s^{(m)}) \log (1 - D(a', s^{(m)}, u'^{(m)})) \right] \quad (70)$$

675 **Numerical stability.** Assuming  $D$  uses the sigmoid nonlinearity on its outputs, we can improve the  
 676 numerical stability by computing the logarithm and sigmoid jointly. This results in

$$\log \sigma(x) = -\log(1 + \exp(-x)) \quad (71)$$

$$\log(1 - \sigma(x)) = -x - \log(1 + \exp(-x)) \quad (72)$$

677 The probability ratio  $D/(1 - D)$  can then be computed with

$$\frac{\sigma(x)}{1 - \sigma(x)} = \exp[-\log(1 + \exp(-x)) + x + \log(1 + \exp(-x))] = \exp(x) \quad (73)$$

### 678 D.3 Successor representations.

679 If we take  $u'$  equal to the state  $s'$ , we can observe that the sum  $\sum_{k \geq 1} p^\pi(S_{t+k} = s' | s, a)$  used in  
 680 the contribution coefficients (3) is equal to the successor representation  $M(s, a, s')$  introduced by  
 681 Dayan [28]. Hence, we can leverage temporal differences to learn  $M(s, a, s')$ , either in a tabular  
 682 setting [28], or in a deep feature learning setting [29]. Using the successor representations, we  
 683 can construct the state-based contribution coefficients as  $w(s, a, s') = M(s, a, s') / [\sum_{\tilde{a} \in \mathcal{A}} \pi(\tilde{a} |$   
 684  $s)M(s, \tilde{a}, s')] - 1$ .

685 For a rewarding outcome encoding  $U$  different from  $S$ , we can recombine the state-based successor  
 686 representations to obtain the required contribution coefficients using the following derivation.

$$\sum_{k \geq 0} p^\pi(U_k = u' | S_0 = s, A_0 = a) = \sum_{k \geq 0} \sum_{s' \in \mathcal{S}} p^\pi(U_k = u', S_k = s' | S_0 = s, A_0 = a) \quad (74)$$

$$= \sum_{k \geq 0} \sum_{s' \in \mathcal{S}} p(U' = u' | S' = s') p^\pi(S_k = s' | S_0 = s, A_0 = a) \quad (75)$$

$$= \sum_{s' \in \mathcal{S}} p(U' = u' | S' = s') \sum_{k \geq 0} p^\pi(S_k = s' | S_0 = s, A_0 = a) \quad (76)$$

$$= \sum_{s' \in \mathcal{S}} p(U' = u' | S' = s') M(s, a, s') \quad (77)$$

687 where in the second line we used that  $S$  is fully predictive of  $U$  (c.f. Definition 1). Note that  
 688  $p(U' | S')$  is policy independent, and can hence be approximated efficiently using offline data. We  
 689 use this recombination of successor representations in our dynamic programming setup to compute  
 690 the ground-truth contribution coefficients (c.f. App. E).

## 691 E Experimental details and additional results

692 Here, we provide additional details to all experiments performed in the manuscript and present  
 693 additional results. The Python code to run our experiments can be found in the supplementary  
 694 materials zip file.

### 695 E.1 Dynamic programming setup

696 In order to delineate the different policy gradient methods considered in this work, we develop  
 697 a framework to compute ground-truth policy gradients and advantages as well as ground-truth  
 698 contribution coefficients. We will first show how we can recursively compute a quantity closely  
 699 related to the successor representation using dynamic programming which then allows us to obtain  
 700 expected advantages and finally expected policy gradients.

#### 701 E.1.1 Computing ground truth quantities

702 To compute ground truth quantities, we assume that the environment reaches an absorbing, terminal  
 703 state  $s_\infty$  after  $T$  steps for all states  $s$ , i.e.  $p^\pi(S_T = s_\infty | s_0 = s) = 1$ . This assumption is satisfied  
 704 by the linear key-to-door and tree environment we consider. As a helper quantity, we define the  
 705 successor representation:

$$M(s, a, s', T) = \sum_{k=1}^T p^\pi(S_k = s' | S_0 = s, A_0 = a) \quad (78)$$

706 which captures the cumulative probability over all time steps of reaching state  $s'$  when starting from  
 707 state  $S_0 = s$ , choosing initial action  $A_0 = a$  and following the policy  $\pi$  thereafter. We can compute  
 708  $M(s, a, s', T)$  recursively as

$$M(s, a, s', t) = \sum_{s'' \in \mathcal{S}} p(S_1 = s'' | S_0 = s, A_0 = a) \sum_{a'' \in \mathcal{A}} \pi(a'' | s'') (\mathbb{1}_{s''=s'} + M(s'', a'', s', t-1)) \quad (79)$$

709 where  $\mathbb{1}$  is the indicator function and  $M(s, a, s', 0)$  is initialized to 0 everywhere. Then, the various  
 710 quantities can be computed exactly as follows.

### Contribution coefficients

$$w(s, a, u') = \frac{M(s, a, u', T)}{\sum_{a' \in \mathcal{A}} \pi(a' | s) M(s, a', u', T)} - 1 \quad (80)$$

711 where

$$M(s, a, u', T) = \sum_{s' \in \mathcal{S}} M(s, a, s', T) \sum_{a' \in \mathcal{A}} \pi(a' | s') \sum_{r'} p(r' | s', a') \mathbb{1}_{f(s', a', r')=u'}, \quad (81)$$

712 similar to the successor representations detailed in Section D. To discover the full state space  $\mathcal{S}$  for  
 713 an arbitrary environment, we perform a depth-first search through the environment transitions (which  
 714 are deterministic in our setting), and record the encountered states.

### Value function

$$V(s) = \sum_{a' \in \mathcal{A}} \pi(a' | s) r(s, a') + \sum_{k=1}^T \sum_{s' \in \mathcal{S}} p^\pi(S_k = s' | S_0 = s, \pi) \sum_{a' \in \mathcal{A}} \pi(a' | s') r(s', a') \quad (82)$$

$$= \sum_{a' \in \mathcal{A}} \pi(a' | s) r(s, a') + \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} \pi(a | s) M(s, a, s', T) \sum_{a' \in \mathcal{A}} \pi(a' | s') r(s', a') \quad (83)$$

715 where  $r(s', a') = \sum_{r'} r' p(r' | s', a')$

### Action-value function

$$Q(s, a) = r(s, a) + \sum_{k=1}^T \sum_{s' \in \mathcal{S}} p^\pi(S_k = s' | S_0 = s, A_0 = a) \sum_{a' \in \mathcal{A}} \pi(a' | s') r(s', a') \quad (84)$$

$$= r(s, a) + \sum_{s' \in \mathcal{S}} M(s, a, s', T) \sum_{a' \in \mathcal{A}} \pi(a' | s') r(s', a') \quad (85)$$

716 where  $r(s', a') = \sum_{r'} r' p(r' | s', a')$

### 717 E.1.2 Computing expected advantages

718 As a next step we detail how to compute the true expected advantage given a potentially imperfect  
 719 estimator of the contribution coefficient  $\hat{w}$ , the value function  $\hat{V}$  or the action-value function  $\hat{Q}$ .

### COCO

$$\mathbb{E}_\pi[\hat{A}^\pi(s, a)] = r(s, a) - \sum_{a' \in \mathcal{A}} \pi(a' | s) r(s, a') + \mathbb{E}_{T \sim \mathcal{T}(s, \pi)} \left[ \sum_{k=1}^T \hat{w}(s, a, U_k) R_k \right] \quad (86)$$

$$= r(s, a) - \sum_{a' \in \mathcal{A}} \pi(a' | s) r(s, a') + \sum_{u' \in \mathcal{U}} \sum_{k=1}^T p^\pi(U_k = u' | S_0 = s) \hat{w}(s, a, u') r(u') \quad (87)$$

$$= r(s, a) - \sum_{a' \in \mathcal{A}} \pi(a' | s) r(s, a') + \sum_{u' \in \mathcal{U}} \hat{w}(s, a, u') r(u') \sum_{a' \in \mathcal{A}} \pi(a' | s) M(s, a', u', T) \quad (88)$$

720 where  $r(u') = \sum_{r'} r' p(r' | u')$  and  $M(s, a', u', T)$  is defined as in E.1.1.

### Advantage

$$\mathbb{E}_\pi[\hat{A}^\pi(s, a)] = Q(s, a) - \hat{V}(s) \quad (89)$$

721 where  $Q$  is the ground truth action-value function obtained following E.1.1.

### Q-critic

$$\mathbb{E}_\pi[\hat{A}^\pi(s, a)] = \hat{Q}(s, a) - \sum_{a' \in \mathcal{A}} \pi(a' | s) \hat{Q}(s, a') \quad (90)$$

### 722 E.1.3 Computing expected policy gradients

723 Finally, we show how we can compute the expected policy gradient,  $\mathbb{E}_\pi[\hat{\nabla}_\theta V^\pi]$  of a potentially  
724 biased gradient estimator, given the expected advantage  $\mathbb{E}_\pi[\hat{A}^\pi(s, a)]$ .

$$\mathbb{E}_\pi[\hat{\nabla}_\theta V^\pi] = \sum_{k=0}^T \mathbb{E}_{S_k \sim \mathcal{T}(s_0, \pi)} \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s) \mathbb{E}_\pi[\hat{A}^\pi(S_k, a)] \quad (91)$$

$$= \sum_{k=0}^T \sum_{s \in \mathcal{S}} p^\pi(S_k = s|s_0) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s) \mathbb{E}_\pi[\hat{A}^\pi(s, a)] \quad (92)$$

725 Using automatic differentiation to obtain  $\nabla_\theta \pi(a|s_0)$ , and the quantity  $M(s, a, s', T)$  we defined  
726 above, we can then compute the expected policy gradient.

$$\mathbb{E}_\pi[\hat{\nabla}_\theta V^\pi] = \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s_0) \mathbb{E}_\pi[\hat{A}^\pi(s_0, a)] + \quad (93)$$

$$\sum_{s \in \mathcal{S}} \sum_{a_0 \in \mathcal{A}} \pi(a_0 | s_0) M(s_0, a_0, s, T) \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|s) \mathbb{E}_\pi[\hat{A}^\pi(s, a)] \quad (94)$$

727 **Computing the ground-truth policy gradient.** To compute the ground-truth policy gradient, we  
728 can apply the same strategy as for Eq. 93 and replace the expected (possibly biased) advantage  
729  $\mathbb{E}_\pi[\hat{A}^\pi(s_0, a)]$  by the ground-truth advantage function, computed with the ground-truth action-value  
730 function (c.f. Section E.1.1).

### 731 E.2 Bias, variance and SNR metrics

732 To analyze the quality of the policy gradient estimators, we use the signal-to-noise ratio (SNR), which  
733 we further subdivide into variance and bias. A higher SNR indicates that we need fewer trajectories  
734 to estimate accurate policy gradients, hence reflecting better credit assignment. To obtain meaningful  
735 scales, we normalize the bias and variance by the norm of the ground-truth policy gradient.

$$SNR = \frac{\|\nabla_\theta V^\pi\|^2}{\mathbb{E}_\pi[\|\hat{\nabla}_\theta V^\pi - \nabla_\theta V^\pi\|^2]} \quad (95)$$

$$\text{Variance} = \frac{\mathbb{E}_\pi[\|\hat{\nabla}_\theta V^\pi - \mathbb{E}_\pi[\hat{\nabla}_\theta V^\pi]\|^2]}{\|\nabla_\theta V^\pi\|^2} \quad (96)$$

$$\text{Bias} = \frac{\|\mathbb{E}_\pi[\hat{\nabla}_\theta V^\pi] - \nabla_\theta V^\pi\|^2}{\|\nabla_\theta V^\pi\|^2} \quad (97)$$

736 We compute the full expectations  $\mathbb{E}_\pi[\hat{\nabla}_\theta V^\pi]$  and ground-truth policy gradient  $\nabla_\theta V^\pi$  by lever-  
737 aging our dynamic programming setup, while for the expectation of the squared differences  
738  $\mathbb{E}_\pi[\|\hat{\nabla}_\theta V^\pi - \nabla_\theta V^\pi\|^2]$  we use Monte Carlo sampling with a sample size of 512. We report  
739 the metrics in Decibels in all figures.

740 **Focusing on long-term credit assignment.** As we are primarily interested in assessing the long-  
741 term credit assignment capabilities of the gradient estimators, we report the statistics of the policy  
742 gradient estimator corresponding to learning to pick up the key or not. Hence, we compare the SNR,  
743 variance and bias of a partial policy gradient estimator considering only  $t = 0$  in the outer sum  
744 (corresponding to the state with the key) for all considered estimators (c.f. Table 1).

745 **Shadow training.** Policy gradients evaluated during training depend on the specific learning  
746 trajectory of the agent. Since all methods' policy gradient estimators contain noise, these trajectories  
747 are likely different for the different methods. As a result, it is difficult to directly compare the quality  
748 of the policy gradient estimators, since it depends on the specific data generated by intermediate  
749 policies during training. In order to allow for a controlled comparison between methods independent  
750 of the noise introduced by different trajectories, we consider a *shadow training* setup in which the  
751 policy is trained with the Q-critic method using ground-truth action-value functions. We can then  
752 compute the policy gradients for the various estimators on the same shared data along this learning  
753 trajectory without using it to actually train the policy. We use this strategy to generate the results  
754 shown in Fig. 1B (right), Fig. 2B and Fig. 3C-D.

### 755 E.3 Linear key-to-door environment setup

756 We simplify the key-to-door environment previously considered by various papers [e.g. 3, 4, 30], to a  
 757 one-dimensional, linear track instead of the original two-dimensional grid world. This version still  
 758 captures the difficulty of long-term credit assignment but reduces the computational burden allowing  
 759 us to thoroughly analyze different policy gradient estimators with the aforementioned dynamic  
 760 programming based setup. The environment is depicted in Fig. 3A. Here, the agent needs to pick  
 761 up a key in the first time step, after which it engages in a distractor task of picking up apples which  
 762 can either be to the left or the right of the agent and which can stochastically assume two different  
 763 reward values. Finally, the agent reaches a door which it can open with the key to collect a treasure  
 764 reward.

765 In our simulations we represent states using a nine-dimensional vector, encoding the relative position  
 766 on the track as a floating number, a boolean for each item that could be present at the current location  
 767 (empty, apple left, apple right, door, key, treasure) as well as boolean indicating whether the agent  
 768 has the key and a boolean for whether the agent does not have the key.

769 There are four discrete actions the agent can pick at every time step: pick the key, pick to the left,  
 770 pick to the right and open the door. Regardless of the chosen action, the agent will advance to the  
 771 next position in the next time step. If it has correctly picked up the key in the first time step and  
 772 opened the door in the penultimate time step, it will automatically pick up the treasure, not requiring  
 773 an additional action.

774 Hung et al. [4] showed that the signal-to-noise ratio (SNR) of the REINFORCE policy gradient [6]  
 775 for solving the main task of picking up the key can be approximated by

$$\text{SNR}^{\text{REINF}} \approx \frac{\|\mathbb{E}[\hat{\nabla}_{\theta}^{\text{REINF}} V^{\pi}]\|^2}{C(\theta)\mathbb{V}[\sum_{t \in T_2} R_t] + \text{Tr}[\mathbb{V}[\hat{\nabla}_{\theta}^{\text{REINF}} V^{\pi} \mid \text{no T2}]}]. \quad (98)$$

776 with  $\hat{\nabla}_{\theta}^{\text{REINF}} V^{\pi}$  the REINFORCE estimator (c.f. Table 1),  $C(\theta)$  a reward-independent constant,  $T_2$   
 777 the set of time steps corresponding to the distractor task, and  $\text{Tr}[\mathbb{V}[\hat{\nabla}_{\theta}^{\text{REINF}} V^{\pi} \mid \text{no T2}]]$  the trace of  
 778 the covariance matrix of the REINFORCE estimator in an equivalent task setup without distractor  
 779 rewards. Hence, we can adjust the difficulty of the task by increasing the number of distractor rewards  
 780 and their variance. We perform experiments with environments of length  $L \in \{20, 40, \dots, 100\}$   
 781 choosing the reward values such that the total distractor reward remains approximately constant.  
 782 Concretely, distractor rewards are sampled as  $r_{\text{distractor}} \sim \mathcal{U}(\{\frac{2}{L}, \frac{18}{L}\})$  and the treasure leads to a  
 783 deterministic reward of  $r_{\text{treasure}} = \frac{4}{L}$ .

### 784 E.4 Reward switching setup

785 While learning to get the treasure in the key-to-door environment requires long term credit assignment  
 786 as the agent needs to learn to 1) pickup the key and 2) open the door, learning to stop picking up  
 787 the treasure does not require long term credit assignment, since the agent can simply learn to stop  
 788 opening the door.

789 We therefore reuse the linear key-to-door environment of length  $L = 40$ , with the single difference  
 790 that we remove the requirement to open the door in order to get the treasure reward. The agent thus  
 791 needs to perform similar credit assignment to both get the treasure and stop getting the treasure.  
 792 When applying reward switching, we simply flip the sign of the treasure reward while keeping the  
 793 distractor reward unchanged.

### 794 E.5 Tree environment setup

795 We parameterize the tree environment by its depth  $d$ , the number of actions  $n_a$  determining the  
 796 branching factor, and the *state overlap*  $o_s$ , defined as the number of overlapping children from two  
 797 neighbouring nodes. The states are represented by two integers:  $i < d$  representing the current level  
 798 in the tree, and  $j$  the position of the state within that level. The root node has state  $(i, j) = (0, 0)$ ,  
 799 and the state transitions are deterministic and given by  $i \leftarrow i + 1$  and  $j \leftarrow j(n_a - o_s) + a$ , where  
 800 we represent the action by an integer  $0 \leq a < n_a$ . We assign each state-action pair a reward  
 801  $r \in \{-2, -1, 0, 1, 2\}$ , computed as

$$r(s, a) = ((\text{idx}(s) + ap + \text{seed}) \bmod n_r) - n_r / 2 \quad (99)$$

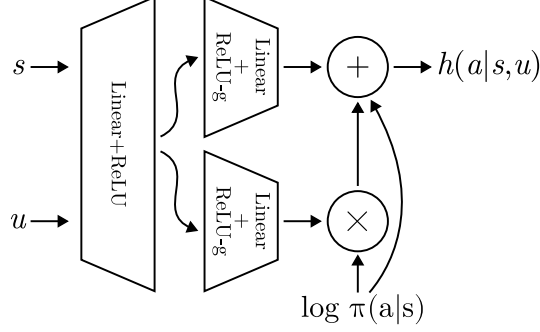


Figure 7: Schematic of the neural network architecture of the hindsight models

with the modulo operator  $\text{mod}$ , the number of reward values  $n_r = 5$ ,  $\text{idx}(s)$  a unique integer index for the state  $s = (i, j)$ , a large prime number  $p$ , and an environment seed.

To introduce rewarding outcome encodings with varying information content, we group state-action pairs corresponding to the same reward value in  $n_g$  groups:

$$u = (\text{idx}(s) + ap + \text{seed}) \bmod (n_r n_g) - n_r / 2 \quad (100)$$

**Environment parameters for Figure 2.** For the experiment of Fig. 2, we use 6 actions and a depth of 4. We plotted the variance of the COCOA estimator for encodings  $U$  corresponding to  $n_g = 4$  and  $n_g = 32$ .

## E.6 Training details

### E.6.1 Architecture

We use separate fully-connected ReLU networks to parameterize the policy, value function, and action-value function. For the policy we use two hidden layers of size 64, for the value function and action-value function respectively we use a single hidden layer of size 256. For the hindsight model of both HCA and COCOA we found a simple multilayer perceptron with the state  $s$  and rewarding outcome encoding  $u'$  as inputs to perform poorly. We hypothesize that this is due to the policy dependence of the hindsight distribution creating a moving target during learning as the policy changes. Leveraging Proposition 5, we therefore add the policy logits as an extra input to the hindsight network to ease tracking the changing policy. We found good performance using a simple hypernetwork, depicted in Fig. 7 that combines the policy logits with the state and hindsight object inputs through a multiplicative interaction, outputting a logit for each possible action. The multiplicative interaction denoted by  $\otimes$  consists of a matrix multiplication of a matrix output by the network with the policy logits, and can be interpreted as *selecting* a combination of policy logits to add to the output channel. In order to allow gating with both positive and negative values, we use a gated version of the ReLU nonlinearity in the second layer which computes the difference between the split, rectified input effectively halving the output dimension:

$$\text{ReLU-g}(x) = \text{ReLU}(x_{0:n/2}) - \text{ReLU}(x_{n/2:n}) \quad (101)$$

with  $n$  the dimension of  $x$ . Gating in combination with the multiplicative interaction is a useful inductive bias for the hindsight model, since for actions which have zero contribution towards the rewarding outcome  $u$  the hindsight distribution is equal to the policy. To increase the performance of our HCA+ baseline, we provide both the policy logits and one minus the policy logits to the multiplicative interaction.

### E.6.2 Optimization

For training of all models we use the AdamW optimizer with default parameters only adjusting the learning rates and clipping the global norm of the policy gradient. We use entropy regularization in combination with epsilon greedy to ensure sufficient exploration to discover the optimal policy. To estimate (action-) value functions we use  $\text{TD}(\lambda)$  treating each  $\lambda$  as a hyperparameter. For all linear layers we use the default initialization of Haiku [80] where biases are initialized as zero and weights are sample from a truncated Gaussian with standard deviation  $\frac{1}{\sqrt{n_{\text{input}}}}$  where  $n_{\text{input}}$  is the input dimension of the respective layer and.

Table 2: The range of values swept over for each hyperparameter in a grid search for the linear key-to-door environment.

Hyperparameter	Range
lr_agent	{0.003, 0.001, 0.0003, 0.0001}
lr_hindsight	{0.01, 0.003, 0.001}
lr_value	{0.01, 0.003, 0.001}
lr_qvalue	{0.01, 0.003, 0.001}
lr_features	{0.01, 0.003, 0.001}
td_lambda_value	{0., 0.5, 0.9, 1.}
td_lambda_qvalue	{0., 0.5, 0.9, 1.}
entropy_reg	{0.3, 0.1, 0.03, 0.01}

### E.6.3 Reward features

Learned reward features should both be fully predictive of the reward (c.f. Theorem 1), and contain as little additional information about the underlying state-action pair as possible (c.f. Theorem 3). We can achieve the former by training a network to predict the rewards given a state-action pair, and take the penultimate layer as the feature  $u = f(s, a)$ . For the latter, there exist multiple approaches. When using a deterministic encoding  $u = f(s, a)$ , we can bin the features such that similar features predicting the same reward are grouped together. When using a stochastic encoding  $p(U | S, A)$  we can impose an information bottleneck on the reward features  $U$ , enforcing the encoding to discard as much information about  $U$  as possible [81, 82]. We choose the deterministic encoding approach, as our Dynamic Programming routines require a deterministic encoding.<sup>1</sup> We group the deterministic rewarding outcome encodings by discretizing the reward prediction network up to the penultimate layer.

**Architecture.** For the neural architecture of the reward prediction network we choose a multilayer perceptron with a single hidden layer of size 128 per action that takes as input the state representation. We use the threshold function,  $\mathbb{1}_{x>0.05}$ , as nonlinearity and use straight-through estimation when computing gradients. We apply the same nonlinearity to the weights before the penultimate layer and initialize them with a Gaussian distribution of mean 0.05 and unit variance. Weights of the readout layer are initialized with a Gaussian distribution of mean 0 and std  $\frac{1}{\sqrt{128}}$ .

**Loss function.** We train the reward network on the mean squared error loss against the reward. To avoid spurious contributions (c.f. 3.3), we encourage the network to learn sparse features that discard information irrelevant to the prediction of the reward, by adding a L1 regularization term to all weights up to the penultimate layer with a strength of  $\eta_{L_1} = 0.001$ . The readout weights are trained with standard L2 regularization with strength  $\eta_{L_2} = 0.03$ . All weights regularization are treated as weight decay, with the decay applied after the gradient update.

**Pretraining.** To learn the reward features, we collect the first 30 mini-batches of episodes in a buffer using a frozen random policy. We then sample triplets  $(s, a, r)$  from the buffer and train with full-batch gradient descent using the Adam optimizer over 20000 steps with a learning rate of 0.003. Once trained, the reward network is frozen, and the features at the penultimate layer are used to train the contribution coefficient as in other COCOA methods. To ensure a fair comparison, other methods are already allowed to train the policy on the first 30 batches of episodes.

### E.6.4 Hyperparameters

**Linear key-to-door setup (performance)** For all our experiments on the linear key-to-door environment, we chose a batch size of 8, while using a batch size of 512 to compute the average performance, SNR, bias and variance metrics. We followed a 2-step selection procedure for selecting the hyperparameters: first, we retain the set of hyperparameters for which the environment can be solved for at least 90% of all seeds, given a large amount of training budget. An environment is considered to be solved for a given seed when the probability of picking up the treasure is above 90%. Then, out of all those hyperparameters, we select the one which maximizes the cumulative amount of treasures picked over 10000 training batches. We used 30 seeds for each set of hyperparameters to

<sup>1</sup>In principle, our dynamic programming routines can be extended to allow for probabilistic rewarding outcome encodings and probabilistic environment transitions, which we leave to future work.



Table 3: The value selected for each hyperparameter on the linear key-to-door environment. The best performing hyperparameters were identical accross all environment lenght.

Hyperparameter	COCOA-reward	COCOA-feature	HCA+	Q-critic	Advantage	REINFORCE
lr_agent	0.0003	0.0003	0.0003	0.0003	0.001	0.0003
lr_hindsight	0.003	0.003	0.003	-	-	-
lr_value	-	-	-	-	0.001	-
lr_qvalue	-	-	-	0.003	-	-
lr_features	-	0.003	-	-	-	-
td_lambda_value	-	-	-	-	1.	-
td_lambda_qvalue	-	-	-	0.9	-	-

Table 4: The entropy regularization value selected for each environment length of the linear key-to-door environment. The values were obtained by linearly interpolating in log-log space between the best performing entropy regularization strength between environment length 20 and 100.

Environment length	20	40	60	80	100	100 (reward-aliasing)
entropy_reg	0.03	0.0187	0.0142	0.0116	0.01	0.0062

identify the best performing ones, then drew 30 fresh seeds for our evaluation. The range considered for our hyperparameter search can be found in Table 2, and the selected hyperparameters in Table 3. Surprisingly, we found that for the environment length considered, the same hyperparameters were performing best, with the exception of `entropy_reg`. For the final values of `entropy_reg` for each environment length, we linearly interpolated in log-log space between the best performing values, 0.03 for length 20 and 0.01 for 100. The values can be found in Table 4.

**Linear key-to-door setup (shadow training)** For measuring the bias and variances of different methods in the shadow training setting, we used the best performing hyperparameters found in the performance setting. We kept a batch size of 8 for the behavior policy and shadow training, while using a batch size of 512 during evaluation.

**Reward switching setup** For the reward switching experiment, we chose hyperparameters following a similar selection procedure as in the linear key-to-door setup, but in the simplified door-less environment of length 40, without any reward switching. We found very similar hyperparameters to work well despite the absence of a door compared to the linear key-to-door setup. However, in order to ensure that noisy methods such as REINFORCE fully converged before the moment of switching the reward, we needed to train the models for 60000 training batches before the switch. To stabilize the hindsight model during this long training period, we added coarse gradient norm clipping by 1.. Furthermore, we found that a slightly decreased learning rate of 0.001 for the Q-critic performed best. Once the best hyperparameters were found, we applied the reward switching to record the speed of adaptation for each algorithm. We kept a batch size of 8 for the behavior policy and shadow training, while using a batch size of 512 during evaluation.

## E.7 Additional results

We perform additional experiments to corroborate the findings of the main text. Specifically, we investigate how *reward aliasing* affects COCOA-reward and COCOA-feature and show that there is no significant difference in performance between HCA and our simplified version, HCA+.

### E.7.1 Reward aliasing

For COCOA-reward we use the scalar reward value to identify rewarding outcomes in the hindsight distribution, i.e.  $U = R$ . In cases where multiple rewarding outcomes yield an identical scalar reward value, the hindsight distribution cannot distinguish between them and has to estimate a common hindsight probability, rendering learning of the contribution coefficients potentially more difficult. In contrast, COCOA-feature learns hindsight features of rewarding objects that are predictive of the reward. Even when multiple rewarding objects lead to an identical scalar reward value their corresponding features are likely different and hence the performance of COCOA-feature should be similar to the case where reward values differ.

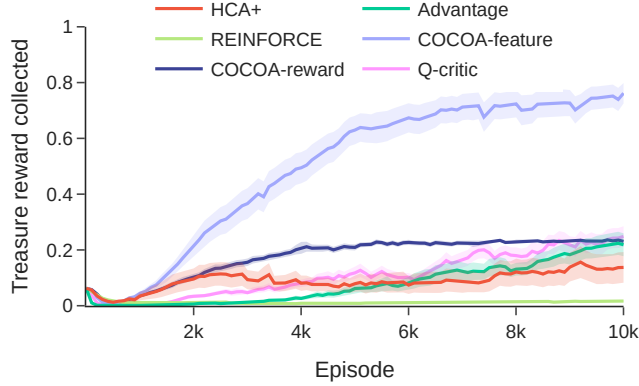


Figure 8: **COCOA-features is robust to reward aliasing.** On a version of the linear key-to-door environment where one of the distractor reward values has the same magnitude as the treasure reward, COCOA-reward can no longer distinguish between the distractor and treasure reward and as a result its performance decreases. COCOA-feature is robust to this manipulation since it relies on learned features to distinguish rewarding objects.

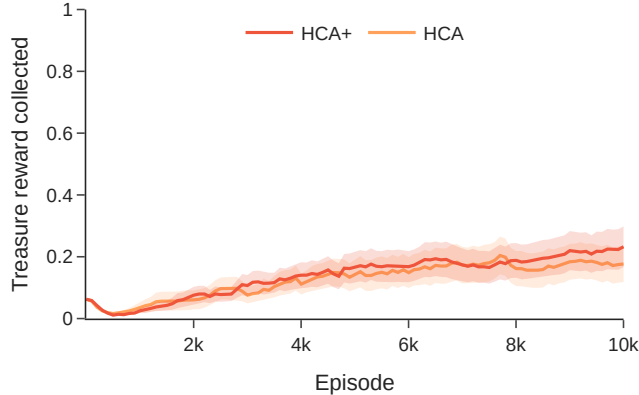


Figure 9: **No performance difference between the original HCA method and our modified variant HCA+.** In the linear key-to-door environment of length 103 both the original HCA method with an additional learned reward model and our simplified version HCA+ perform similarly in terms of performance measured as the percentage of treasure reward collected.

912 In Fig. 8, we test this *reward aliasing* setting experimentally and slightly modify the linear key-to-door  
 913 environment of length  $L = 100$  by giving the treasure reward the same value as one of the two  
 914 possible values of the stochastic distractor rewards. Concretely, we sample the distractor rewards  
 915 uniformly as  $r_{\text{distractor}} \sim \mathcal{U}(\{\frac{2}{L}, \frac{18}{L}\})$  but now set the deterministic treasure reward to  $r_{\text{treasure}} = \frac{2}{L}$ .  
 916 As expected, COCOA-feature is robust to reward aliasing, continuing to perform well on the task  
 917 of picking up the treasure while performance of COCOA-reward noticeably suffers. Note that the  
 918 performance of all methods has slightly decreased compared to Fig. 3, as the magnitude of the  
 919 treasure reward is now smaller relative to the variance of the distractor rewards, resulting in a worse  
 920 SNR for all methods.

#### 921 E.7.2 HCA vs HCA+

922 Our policy gradient estimator for  $U = S$  presented in Eq. 4 differs slightly from Eq. 2, the version  
 923 originally introduced by Harutyunyan et al. [1], as we remove the need for a learned reward model  
 924  $r(s, a)$ . We empirically verify that this simplification does not lead to a decrease in performance in  
 925 Fig. 9. We run the longest and hence most difficult version of the linear key-to-door environment  
 926 considered in our experiments and find no significant difference in performance between our simplified  
 927 version (HCA+) and the original variant (HCA).

## F Contribution analysis in continuous spaces and POMDPs

In Section 3.3 we showed that HCA can suffer from spurious contributions, as state representations need to contain detailed features to allow for a capable policy. The same level of detail however is detrimental when assigning credit to actions for reaching a particular state, since at some resolution almost every action will lead to a slightly different outcome. Measuring the contribution towards a specific state ignores that often the same reward could be obtained in a slightly different state, hence overvaluing the importance of past actions. Many commonly used environments, such as pixel-based environments, continuous environments, and partially observable MDPs exhibit this property to a large extent due to their fine-grained state representations. Here, we take a closer look at how spurious contributions arise in continuous environments and Partially Observable MDPs (POMDPs).

### F.1 Spurious contributions in continuous state spaces

When using continuous state spaces,  $p^\pi(S_k = s' | s, a)$  represents a probability density function (PDF) instead of a probability. A ratio of PDFs  $p(X = x)/p(X = x')$  can be interpreted as a likelihood ratio of how likely a sample  $X$  will be close to  $x$  versus  $x'$ . Using PDFs, the contribution coefficients with  $U = S$  used by HCA result in

$$w(s, a, s') = \frac{\sum_{k \geq 1} p^\pi(S_{t+k} = s' | S_t = s, A_t = a)}{\sum_{k \geq 1} p^\pi(S_{t+k} = s' | S_t = s)} - 1 = \frac{p^\pi(A_t = a | S_t = s, S' = s')}{\pi(a | s)} - 1 \quad (102)$$

and can be interpreted as a likelihood ratio of encountering a state ‘close’ to  $s'$ , starting from  $(s, a)$  versus starting from  $s$  and following the policy. The contribution coefficient will be high if  $p^\pi(S' = s' | s, a)$  has a high probability density at  $s'$  conditioned on action  $a$ , compared to the other actions  $a' \neq a$ . Hence, the less  $p^\pi(S_k = s' | s, a)$  and  $p^\pi(S_k = s' | s, a')$  overlap around  $s'$ , the higher the contribution coefficient for action  $a$ . The variance of the distribution  $p^\pi(S_k = s' | s, a)$ , and hence its room for overlap, is determined by how diffuse the environment transitions and policy are. For example, a very peaked policy and nearly deterministic environment transitions leads to a sharp distribution  $p^\pi(S_k = s' | s, a)$ .

The randomness of the policy and environment is a poor measure for the contribution of an action in obtaining a reward. For example, consider the musician’s motor system, that learns to both: (i) play the violin and (ii) play the keyboard. We assume a precise control policy and near deterministic environment dynamics, resulting in peaked distributions  $p^\pi(S_k = s' | s, a)$ . For playing the violin, a slight change in finger position significantly influences the pitch, hence the reward function sharply declines around the target state with perfect pitch. For playing the keyboard, the specific finger position matters to a lesser extent, as long as the correct key is pressed, resulting in a relatively flat reward function w.r.t. the precise finger positioning. The state-based contribution coefficients of Eq. 102 result in a high contribution for the action taken in the trajectory for both tasks. For playing violin, this could be a good approximation of what we intuitively think of as a ‘high contribution’, whereas for the keyboard, this overvalues the importance of the past action in many cases. From this example, it is clear that measuring contributions towards rewarding *states* in continuous state spaces can lead to spurious contributions, as the contributions mainly depend on how diffuse the policy and environment dynamics are, and ignore that different reward structures can require vastly different contribution measures.

### F.2 Deterministic continuous reward functions can lead to excessive variance

Proposition 2 shows that HCA degrades to REINFORCE in environments where each action sequence leads to distinct states. In Section F.1, we discussed that continuous state spaces exhibit this property to a large extent due to their fine-grained state representation. If each environment state has a unique reward value, COCOA can suffer from a similar degradation to HCA. As the rewarding outcome encoding  $U$  needs to be fully predictive of the reward (c.f. Theorem 1), it needs to have a distinct encoding for each unique reward, and hence for each state.

With a continuous reward function this can become a significant issue, if the reward function is deterministic. Due to the fine-grained state representation in continuous state spaces, almost every action will lead to a (slightly) different state. When we have a deterministic, continuous reward function, two nearby but distinct states will often lead to nearby but distinct rewards. Hence, to a large extent, the reward contains the information of the underlying state, encoded in its infinite precision of a real value, resulting in COCOA degrading towards the high-variance HCA estimator.

The above problem does not occur when the reward function  $p(R | S, A)$  is probabilistic, even for continuous rewards. If the variance of  $p(R | S, A)$  is bigger than zero, the variance of  $p(S, A | R)$  is also bigger than zero under mild assumptions. This means that it is not possible to perfectly decode the underlying state using a specific reward value  $r$ , and hence COCOA does not degrade towards HCA in this case. Intuitively, different nearby states can lead to the same sampled reward, removing the spurious contributions. In the following section, we will use this insight to alleviate spurious contributions, even for deterministic reward functions.

### F.3 Smoothing can alleviate excess variance trading variance for bias

When each state has a unique reward, an unbiased COCOA estimator degrades to the high-variance HCA estimator, since  $U$  needs to be fully predictive of the reward, and hence each state needs a unique rewarding outcome encoding. Here, we propose two ways forward to overcome the excess variance of the COCOA estimator in this extreme setting that trade variance for bias.

**Rewarding outcome binning.** One intuitive strategy to avoid that each state has a unique rewarding outcome encoding  $U$ , is to group rewarding outcome encodings corresponding to nearby rewards together, resulting in discrete bins. As the rewarding outcome encodings now contain less details, the resulting COCOA estimator will have lower variance. However, as now the rewarding outcome encoding is not anymore fully predictive of the reward, the COCOA estimator will be biased. An intimately connected strategy is to change the reward function in the environment to a discretized reward function with several bins. Policy gradients in the new discretized environments will not be exactly equal to the policy gradients in the original environment, however for fine discretizations, we would not expect much bias. Similarly for binning the rewarding outcome encodings, when we group few nearby rewarding outcome encodings together, we expect a low bias, but also a low variance reduction. Increasing the amount of rewarding outcomes we group together we further lower the bias, at a cost of an increasing bias, hence creating a bias-variance trade-off.

**Stochastic rewarding outcomes.** We can generalize the above binning technique towards rewarding outcome encodings that bin rewards stochastically. In Section F.2, we discussed that when the reward function is probabilistic, the excessive variance problem is less pronounced, as different states can lead to the same reward. When dealing with a deterministic reward function, we can introduce stochasticity in the rewarding outcome encoding to leverage the same principle and reduce variance at the cost of increasing bias. For example, we can introduce the rewarding outcome encoding  $U \sim \mathcal{N}(R, \sigma)$ , with  $\mathcal{N}$  corresponding to a Gaussian distribution. As this rewarding outcome encoding is not fully predictive of the reward, it will introduce bias. We can control this bias-variance trade-off with the variance  $\sigma$ : a small sigma corresponds to a sharp distribution, akin to a fine discretization in the above binning strategy, and hence a low bias. Increasing  $\sigma$  leads to more states that could lead to the same rewarding outcome encoding, hence lowering the variance at the cost of increasing the bias.

**Implicit smoothing by noise perturbations or limited capacity networks.** Interestingly, the above strategy of defining stochastic rewarding outcomes is equivalent to adjusting the training scheme of the hindsight model  $h(a | s, u')$  by adding noise to the input  $u'$ . Here, we take  $U$  equal to the (deterministic) reward  $R$ , but add noise to it while training the hindsight model. Due to the noise, the hindsight model cannot perfectly decode the action  $a$  from its input, resulting in the same effects as explicitly using stochastic rewarding outcomes. Adding noise to the input of a neural network is a frequently used regularization technique. Hence, an interesting route to investigate is whether other regularization techniques on the hindsight model, such as limiting its capacity, can result in a bias-variance trade-off for HCA and COCOA.

**Smoothing hindsight states for HCA.** We can apply the same (stochastic) binning technique to hindsight states for HCA, creating a more coarse-grained state representation for backward credit assignment. However, the bias-variance trade-off is more difficult to control for states compared to rewards. The sensitivity of the reward function to the underlying state can be big in some regions of state-space, whereas small in others. Hence, a uniform (stochastic) binning of the state-space will likely result in sub-optimal bias-variance trade-offs, as it will be too coarse-grained in some regions causing large bias, whereas too fine-grained in other regions causing a low variance-reduction. Binning rewards in contrast does not suffer from this issue, as binning is directly performed in a space relevant to the task.

Proposition 8 provides further insight on what the optimal smoothing or binning for HCA looks like. Consider the case where we have a discrete reward function with not too many distinct

values compared to the state space, such that COCOA-reward is a low-variance gradient estimator. Proposition 8 shows that we can recombine the state-based hindsight distribution  $p^\pi(A_0 = a \mid S_0 = s, S' = s')$  into the reward-based hindsight distribution  $p^\pi(A_0 = a \mid S_0 = s, R' = r')$ , by leveraging a smoothing distribution  $\sum_{s' \in \mathcal{S}} p^\pi(S' = s' \mid S_0 = s, R' = r')$ . When we consider a specific hindsight state  $s''$ , this means that we can obtain the low-variance COCOA-reward estimator, by considering all states  $S'$  that could have lead to the same reward  $r(s'')$ , instead of only  $s''$ . Hence, instead of using a uniform stochastic binning strategy with e.g.  $S' \sim \mathcal{N}(s'', \sigma)$ , a more optimal binning strategy is to take the reward structure into account through  $p^\pi(S' \mid S_0 = s, R' = r(s''))$

**Proposition 8.**

$$p^\pi(A_0 = a \mid S_0 = s, R' = r') = \sum_{s' \in \mathcal{S}} p^\pi(S' = s' \mid S_0 = s, R' = r') p^\pi(A_0 = a \mid S_0 = s, S' = s') \quad (103)$$

*Proof.* As  $A_0$  is d-separated from  $R'$  conditioned on  $S'$  and  $S_0$  (c.f. Fig. 5b) we can use the implied conditional independence to prove the proposition:

$$p^\pi(A_0 = a \mid S_0 = s, R' = r') \quad (104)$$

$$= \sum_{s' \in \mathcal{S}} p^\pi(S' = s', A_0 = a \mid S_0 = s, R' = r') \quad (105)$$

$$= \sum_{s' \in \mathcal{S}} p^\pi(S' = s' \mid S_0 = s, R' = r') p^\pi(A_0 = a \mid S_0 = s, S' = s', R' = r') \quad (106)$$

$$= \sum_{s' \in \mathcal{S}} p^\pi(S' = s' \mid S_0 = s, R' = r') p^\pi(A_0 = a \mid S_0 = s, S' = s') \quad (107)$$

□

#### F.4 Continuous action spaces

In the previous section, we discussed how continuous state spaces can lead to spurious contributions resulting high variance for the HCA estimator. Here, we briefly discuss how COCOA can be applied to continuous action spaces, and how Proposition 2 translates to this setting.

**Gradient estimator.** We can adjust the COCOA gradient estimator of Eq. 4 towards continuous action spaces by replacing the sum over  $a' \in \mathcal{A}$  by an integral over the action space  $A'$ .

$$\hat{\nabla}_\theta^U V^\pi(s_0) = \sum_{t \geq 0} \nabla_\theta \log \pi(A_t \mid S_t) R_t + \int_{\mathcal{A}} da \nabla_\theta \pi(a \mid S_t) \sum_{k \geq 1} w(S_t, a, U_{t+k}) R_{t+k} \quad (108)$$

In general, computing this integral is intractable. We can approximate the integral by standard numerical integration methods, introducing a bias due to the approximation. Alternatively, we introduced in App. C.3 another variant of the COCOA gradient estimator that samples independent actions  $A'$  from the policy instead of summing over the whole action space. This variant can readily be applied to continuous action spaces, resulting in

$$\hat{\nabla}_\theta V^\pi = \sum_{t \geq 0} \nabla_\theta \log \pi(A_t \mid S_t) R_t + \frac{1}{M} \sum_m \nabla_\theta \log \pi(a^m \mid S_t) \sum_{k=1}^{\infty} w(S_t, a^m, U_{t+k}) R_{t+k} \quad (109)$$

where we sample  $M$  actions independently from  $a^m \sim \pi(\cdot \mid S_t)$ . This gradient estimator is unbiased, but introduces extra variance through the sampling of actions.

**Spurious contributions.** Akin to discrete action spaces, HCA in continuous action spaces can suffer from spurious contributions when distinct action sequences lead to unique states. In this case, previous actions can be perfectly decoded from the hindsight state, and we have that the probability density function  $p^\pi(A_0 = a \mid S_0 = s, S' = s')$  is equal to the Dirac delta function  $\delta(a = a_0)$ , with  $a_0$  the action taken in the trajectory that led to  $s'$ . Substituting this Dirac delta function into the policy gradient estimator of Eq. 108 results in the high-variance REINFORCE estimator. When we use the COCOA estimator of Eq. 109 using action sampling, HCA will even have a higher variance compared to REINFORCE.

## 1067 F.5 Partially Observable MDPs

1068 In many environments, agents do not have access to the complete state information  $s$ , but instead  
 1069 get observations with incomplete information. Partially observable Markov decision processes  
 1070 (POMDPs) formalize this case by augmenting MDPs with an observation space  $\mathcal{O}$ . Instead of directly  
 1071 observing Markov states, the agent now acquires an observation  $o_t$  at each time step. The probability  
 1072 of observing  $o$  in state  $s$  after performing action  $a$  is given by  $p_o(o \mid s, a)$ .

1073 A popular strategy in deep RL methods to handle partial observability is learning an internal state  $x_t$   
 1074 that summarizes the observation history  $h_t = \{o_{t'+1}, a_{t'}, r_{t'}\}_{t'=0}^{t-1}$ , typically by leveraging recurrent  
 1075 neural networks [83–86]. This internal state is then used as input to a policy or value function. This  
 1076 strategy is intimately connected to estimating *belief states* [87]. The observation history  $h_t$  provides  
 1077 us with information on what the current underlying Markov state  $s_t$  is. We can formalize this by  
 1078 introducing the *belief state*  $b_t$ , which captures the sufficient statistics for the probability distribution  
 1079 over the Markov states  $s$  conditioned on the observation history  $h_t$ . Theoretically, such belief states  
 1080 can be computed by doing Bayesian probability calculus using the environment transition model,  
 1081 reward model and observation model:

$$p_B(s \mid b_t) = p(S_t = s \mid H_t = h_t) \quad (110)$$

1082 Seminal work has shown that a POMDP can be converted to an MDP, using the belief states  $b$  as new  
 1083 Markov states instead of the original Markov states  $s$  [88]. Now the conventional optimal control  
 1084 techniques can be used to solve the belief-state MDP, motivating the use of standard RL methods  
 1085 designed for MDPs, combined with learning an internal state  $x$ .

1086 **HCA suffers from spurious contributions in POMDPs.** As the internal state  $x$  summarizes  
 1087 the complete observation history  $h$ , past actions can be accurately decoded based on  $h$ , causing  
 1088 HCA to degrade to the high-variance REINFORCE estimator (c.f. Proposition 2). Here, the tension  
 1089 between forming good state representations for enabling capable policies and good representations for  
 1090 backward credit assignment is pronounced clearly. To enable optimal decision-making, a good internal  
 1091 state  $x$  needs to encode which underlying Markov states the agent most likely occupies, as well as the  
 1092 corresponding uncertainty. To this end, the internal state needs to incorporate information about the  
 1093 full history. However, when using the same internal state for backward credit assignment, this leads  
 1094 to spurious contributions, as previous actions are encoded directly into the internal state.

1095 To make the spurious contributions more tangible, let us consider a toy example where the state  
 1096 space  $\mathcal{S}$  consists of three states  $x, y$  and  $z$ . We assume the internal state represents a belief state  
 1097  $b = \{b^1, b^2\}$ , which is a sufficient statistic for the belief distribution:

$$p(S = x \mid b) = b^1, \quad p(S = y \mid b) = b^2, \quad p(S = z \mid b) = 1 - b^1 - b^2 \quad (111)$$

1098 We assume that  $b_t$  is deterministically computed from the history  $h_t$ , e.g. by an RNN. Now consider  
 1099 that at time step  $k$ , our belief state is  $b_k = \{0.5, 0.25\}$  and we get a reward that resulted from the  
 1100 Markov state  $x$ . As the belief states are deterministically computed, the distribution  $p^\pi(B' = b' \mid b, a)$   
 1101 is a Dirac delta distribution. Now consider that the action  $a$  does not influence the belief distribution  
 1102 over the rewarding state  $x$ , but only changes the belief distribution over the non-rewarding states  
 1103 (e.g.  $B_k = \{0.5, 0.23\}$  instead of  $B_k = \{0.5, 0.25\}$  when taking action  $a'$  instead of  $a$ ). As the  
 1104 Dirac delta distributions  $p^\pi(B' = b' \mid b, a)$  for different  $a$  do not overlap, we get a high contribution  
 1105 coefficient for the action  $a_0$  that was taken in the actual trajectory that led to the belief state  $b'$ , and  
 1106 a low contribution coefficient for all other actions, even though the actions did not influence the  
 1107 distribution over the rewarding state.

1108 The spurious contributions originate from measuring contributions towards reaching a certain internal  
 1109 belief state, while ignoring that the same reward could be obtained in different belief states as well.  
 1110 Adapting Proposition 8 towards these internal belief states provides further insight on the difference  
 1111 between HCA and COCOA-reward:

$$p^\pi(A_0 = a \mid X_0 = x, R' = r') = \sum_{X' \in \mathcal{X}} p^\pi(X' = x' \mid X_0 = x, R' = r') p^\pi(A_0 = a \mid X_0 = x, X' = x') \quad (112)$$

1112 Here, we see that the contribution coefficients of COCOA-reward take into account that the same  
 1113 reward can be obtained while being in different internal belief states  $x'$  by averaging over them, while  
 1114 HCA only considers a single internal state.

## G Learning contribution coefficients from non-rewarding observations

### G.1 Latent learning

Building upon HCA [1], we learn the contribution coefficients (3) by approximating the hindsight distribution  $p^\pi(a \mid s, u')$ . The quality of this model is crucial for obtaining low-bias gradient estimates. However, its training data is scarce, as it is restricted to learn from on-policy data, and rewarding observations in case the reward or rewarding object is used as encoding  $U$ . We can make the model that approximates the hindsight distribution less dependent on the policy by providing the policy logits as input (c.f. Proposition 5). Enabling COCOA-reward or COCOA-feature to learn from non-rewarding states is a more fundamental issue, as in general, the rewarding outcome encodings corresponding to zero rewards do not share any features with those corresponding to non-zero rewards.

Empowering COCOA with *latent learning* is an interesting direction to make the learning of contribution coefficients more sample efficient. We refer to *latent learning* as learning useful representations of the environment structure without requiring task information, which we can then leverage for learning new tasks more quickly [89]. In our setting, this implies learning useful representations in the hindsight model without requiring rewarding observations, such that when new rewards are encountered, we can quickly learn the corresponding contribution coefficients, leveraging the existing representations.

### G.2 Optimal rewarding outcome encodings for credit assignment.

Theorem 3 shows that the less information a rewarding outcome encoding  $U$  contains, the lower the variance of the corresponding COCOA gradient estimator (4). Latent learning on the other hand considers the sample-efficiency and corresponding bias of the learned contribution coefficients: when the hindsight model can learn useful representations with encodings  $U$  corresponding to zero rewards, it can leverage those representations to quickly learn the contribution coefficients for rewarding outcome encodings with non-zero rewards, requiring less training data to achieve a low bias.

These two requirements on the rewarding outcome encoding  $U$  are often in conflict. To obtain a low-variance gradient estimator,  $U$  should retain as little information as possible while still being predictive of the reward. To enable latent learning for sample-efficient learning of the contribution coefficients, the hindsight model needs to pick up on recurring structure in the environment, requiring keeping as much information as possible in  $U$  to uncover the structural regularities. Using the state as rewarding outcome encoding is beneficial for enabling latent learning, as it contains rich information on the environment structure, but results in spurious contributions and a resulting high variance. Using the rewards or rewarding object as rewarding outcome encoding removes spurious contributions resulting in low variance, but renders latent learning difficult.

A way out of these conflicting pressures for an optimal rewarding outcome encoding is to use rewards or rewarding objects as the optimal encoding for low variance, but extract these contribution coefficients from models that allow for sample-efficient, latent learning. One possible strategy is to learn probabilistic world models [18, 19, 86] which can be done using both non-rewarding and rewarding observations, and use those to approximate the contribution coefficients of Eq. 3. Another strategy that we will explore in more depth is to learn hindsight models based on the state as rewarding outcome encoding to enable latent learning, but then recombine those learned hindsight models to obtain contribution coefficients using the reward or rewarding object as  $U$ .

### G.3 Counterfactual reasoning on rewarding states

HCA results in spurious contributions because it computes contributions towards reaching a precise rewarding state, while ignoring that the same reward could be obtained in other (nearby) states. Proposition 9 (a generalization of Proposition 8), shows that we can reduce the spurious contributions of the state-based contribution coefficients by leveraging counterfactual reasoning on rewarding states. Here, we obtain contribution coefficients for a certain rewarding outcome encoding (e.g. the rewarding object) by considering which other states  $s'$  could lead to the same rewarding outcome, and averaging over the corresponding coefficients  $w(s, a, s')$ .

**Proposition 9.** *Assuming  $S'$  is fully predictive of  $U'$ , we have that*

$$w(s, a, u') = \sum_{s' \in \mathcal{S}} p^\pi(S' = s' \mid S_0 = s, U' = u') w(s, a, s') \quad (113)$$

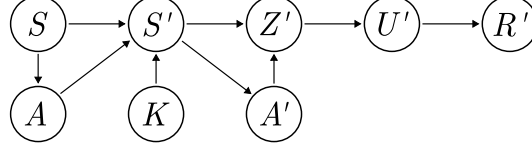


Figure 10: Schematic of the graphical model used in our variational information bottleneck approach.

The proof follows the same technique as that of Proposition 8. Proposition 8) shows that it is possible to learn state-based contribution coefficients, enabling latent learning, and obtain a low-variance COCOA estimator by recombining the state-based contribution coefficients into coefficients with less spurious contributions, if we have access to the generative model  $p^\pi(S' = s' \mid S_0 = s, U' = u')$ . This model embodies the counterfactual reasoning on rewarding states: ‘which other states  $s'$  are likely considering I am currently in  $u'$  and visited state  $s$  somewhere in the past’.

In general, learning this generative model is as difficult or more difficult than approximating the hindsight distribution  $p^\pi(a \mid s, u')$ , as it requires rewarding observations as training data. Hence, it is not possible to directly use Proposition 9 to combine latent learning with low-variance estimators. In the following section, we propose a possible way forward to circumvent this issue.

#### G.4 Learning credit assignment representations with an information bottleneck

Here, we outline a strategy where we learn a latent representation  $Z$  that retains useful information on the underlying states  $S$ , and crucially has a latent space structure such that  $p^\pi(Z' = z' \mid S = s, U' = u')$  is easy to approximate. Then, leveraging Proposition 9 (and replacing  $S'$  by  $Z'$ ) allows us to learn a hindsight representation based on  $Z$ , enabling latent learning, while reducing the spurious contributions by counterfactual reasoning with  $p^\pi(z' \mid s, u')$ .

To achieve this, we use the Variational Information Bottleneck approach [81, 90], closely related to the  $\beta$  Variational Autoencoder [91]. Fig. 10 shows the graphical model with the relations between the various variables and encoding: we learn a probabilistic encoding  $p(Z \mid S, A; \theta)$  parameterized by  $\theta$ , and we assume that the latent variable  $Z$  is fully predictive of the rewarding outcome encoding  $U$  and reward  $R$ . We aim to maximize the mutual information  $\mathcal{I}(Z'; S', A' \mid S, U')$  under some information bottleneck. We condition on  $S$  and  $U'$ , to end up later with a decoder and prior variational model that we can combine with Proposition 9.

Following rate-distortion theory, Alemi et al. [90] consider the following tractable variational bounds on the mutual information

$$H - D \leq \mathcal{I}(Z'; S', A' \mid S, U') \leq \text{Rate} \quad (114)$$

with entropy  $H$ , distortion  $D$  and rate defined as follows:

$$H = -\mathbb{E}_{S', A', S, U'} [\log p^\pi(S', A' \mid S, U')] \quad (115)$$

$$D = -\mathbb{E}_{S, U'} \left[ \mathbb{E}_{S', A' \mid S, U'} \left[ \int dz' p(z' \mid s', a'; \theta) \log q(s', a' \mid s, z'; \psi) \right] \right] \quad (116)$$

$$\text{Rate} = \mathbb{E}_{S, U'} \left[ \mathbb{E}_{S', A' \mid S, U'} [D_{\text{KL}}(p(z' \mid s', a'; \theta) \parallel q(z' \mid s, u'; \phi))] \right] \quad (117)$$

where the *decoder*  $q(s', a' \mid s, z'; \psi)$  is a variational approximation to  $p(s', a' \mid s, z')$ , and the *marginal*  $q(z' \mid s, u'; \phi)$  is a variational approximation to the true marginal  $p(z' \mid s, u')$ . The distortion  $D$  quantifies how well we can decode the state-action pair  $(s', a')$  from the encoding  $z'$ , by using a decoder  $q(s', a' \mid s, z'; \psi)$  parameterized by  $\psi$ . The distortion is reminiscent of an autoencoder loss, and hence encourages the encoding  $p(z' \mid s', a')$  to retain as much information as possible on the state-action pair  $(s', a')$ . The rate measures the average KL-divergence between the encoder and variational marginal. In information theory, this rate measures the extra number of bits (or nats) required to encode samples from  $Z'$  with an optimal code designed for the variational marginal  $q(z' \mid s, u'; \phi)$ .

We can use the rate to impose an information bottleneck on  $Z'$ . If we constrain the rate to  $\text{Rate} \leq a$  with  $a$  some positive constant, we restrict the amount of information  $Z'$  can encode about  $(S', A')$ , as  $p(z' \mid s', a'; \theta)$  needs to remain close to the marginal  $q(z' \mid s, u'; \phi)$ , quantified by the KL-divergence. We can maximize the mutual information  $\mathcal{I}(Z'; S', A' \mid S, U')$  under the information bottleneck by



1205 minimizing the distortion under this rate constraint. Instead of imposing a fixed constraint on the rate,  
 1206 we can combine the rate and distortion into an unconstrained optimization problem by leveraging the  
 1207 Lagrange formulation

$$\min_{\theta, \phi, \psi} D + \beta \text{Rate} \quad (118)$$

1208 Here, the  $\beta$  parameter determines the strength of the information bottleneck. This formulation  
 1209 is equivalent to the  $\beta$  Variational Autoencoder [91], and for  $\beta = 1$  we recover the Variational  
 1210 Autoencoder [92].

1211 To understand why this information bottleneck approach is useful to learn credit assignment represen-  
 1212 tations  $Z$ , we examine the rate in more detail. We can rewrite the rate as

$$\text{Rate} = \mathbb{E}_{S, U'} [D_{\text{KL}}(p(z' | s, u') \| q(z' | s, u'; \phi))] + \quad (119)$$

$$\mathbb{E}_{S, U', Z'} [D_{\text{KL}}(p(s', a' | s, z') \| p(s', a' | s, u'))] \quad (120)$$

1213 Hence, optimizing the information bottleneck objective (118) w.r.t.  $\phi$  fits the variational marginal  
 1214  $q(z' | s, u'; \phi)$  to the true marginal  $p(z' | s, u')$  induced by the encoder  $p(z' | s', a'; \theta)$ . Proposition 9  
 1215 uses this true marginal to recombine coefficients based on  $Z$  into coefficients based on  $U$ . Hence, by  
 1216 optimizing the information bottleneck objective, we learn a model  $q(z' | s, u'; \phi)$  that approximates  
 1217 the true marginal, which we then can use to obtain contribution coefficients with less spurious  
 1218 contributions by leveraging Proposition 9. Furthermore, minimizing the rate w.r.t.  $\theta$  shapes the latent  
 1219 space of  $Z'$  such that the true marginal  $p(z' | s, u')$  moves closer towards the variational marginal  
 1220  $q(z' | s, u'; \phi)$ .

1221 In summary, the outlined information bottleneck approach for learning a credit assignment representa-  
 1222 tion  $Z$  is a promising way forward to merge the powers of latent learning with a low-variance COCOA  
 1223 gradient estimator. The distortion and rate terms in the information bottleneck objective of Eq. 118  
 1224 represent a trade-off parameterized by  $\beta$ . Minimizing the distortion results in a detailed encoding  
 1225  $Z'$  with high mutual information with the state-action pair  $(S', A')$ , which can be leveraged for  
 1226 latent learning. Minimizing the rate shapes the latent space of  $Z'$  in such way that the true marginal  
 1227  $p(z' | s, u')$  can be accurately approximated within the variational family of  $q(z' | s, u'; \phi)$ , and  
 1228 fits the parameters  $\phi$  resulting in an accurate marginal model. We can then leverage the variational  
 1229 marginal  $q(z' | s, u'; \phi)$  to perform counterfactual reasoning on the rewarding state encodings  $Z'$   
 1230 according to Proposition 9, resulting in a low-variance COCOA-estimator.

## 1231 H Contribution analysis and causality

### 1232 H.1 Causal interpretation of COCOA

1233 COCOA is closely connected to causality theory [33], where the contribution coefficients (3) cor-  
 1234 respond to performing *Do* – *interventions* on the causal graph to estimate their effect on future  
 1235 rewards. To formalize this connection with causality, we need to use a new set of tools beyond  
 1236 conditional probabilities, as causation is in general not the same as correlation. We start with rep-  
 1237 resenting the MDP combined with the policy as a directed acyclic graphical model (c.f. Fig. 5a in  
 1238 App. C.1). In causal reasoning, we have two different ‘modes of operation’. On the one hand, we can  
 1239 use observational data, corresponding to ‘observing’ the states of the nodes in the graphical model,  
 1240 which is compatible with conditional probabilities measuring correlations. On the other hand, we can  
 1241 perform *interventions* on the graphical model, where we manually set a node, e.g.  $A_t$  to a specific  
 1242 value  $a$  independent from its parents  $S_t$ , and see what the influence is on the probability distributions  
 1243 of other nodes in the graph, e.g.  $R_{t+1}$ . These interventions are formalized with *do-calculus* [33],  
 1244 where we denote an intervention of putting a node  $V_i$  equal to  $v$  as  $\text{Do}(V_i = v)$ , and can be used to  
 1245 investigate causal relations.

1246 Using the graphical model of Fig. 5b that abstracts time, we can use do-interventions to quantify the  
 1247 causal contribution of an action  $A_t = a$  taken in state  $S_t = s$  upon reaching the rewarding outcome  
 1248  $U' = u'$  in the future as,

$$w_{\text{Do}}(s, a, u') = \frac{p^\pi(U' = u' | S_t = s, \text{Do}(A_t = a))}{\sum_{\tilde{a} \in \mathcal{A}} \pi(\tilde{a} | s) p^\pi(U' = u' | S_t = s, \text{Do}(A_t = \tilde{a}))} - 1. \quad (121)$$

1249 As conditioning on  $S$  satisfies the *backdoor criterion* [33] for  $U'$  w.r.t.  $A$ , the interventional  
 1250 distribution  $p^\pi(U' = u' | S_t = s, \text{Do}(A_t = a))$  is equal to the observational distribution  
 1251  $p^\pi(U' = u' | S_t = s, A_t = a)$ . Hence, the causal contribution coefficients of Eq. 121 are  
 1252 equal to the contribution coefficients of Eq. 3 used by COCOA.

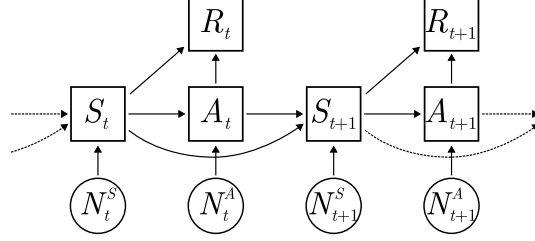


Figure 11: Structural causal model (SCM) of the MDP. Squares represent deterministic functions and circles random variables.

## H.2 Extending COCOA to counterfactual interventions

Within causality theory, *counterfactual reasoning* goes one step further than causal reasoning, by incorporating the hindsight knowledge of the external state of the world in its reasoning. Applied to COCOA, this more advanced counterfactual reasoning would evaluate the query: ‘How does taking action  $a$  influence the probability of reaching a rewarding outcome, compared to taking alternative actions  $a'$ , *given everything else remains the same*’. To formalize the difference between causal and counterfactual reasoning, we need to convert the causal DAG of Figure 5a into a structural causal model (SCM), as shown in Figure 11. The SCM expresses all conditional distributions as deterministic functions with independent noise variables  $N$ , akin to the reparameterization trick [92]. In causal reasoning, we perform do-interventions on nodes, which is equivalent to cutting all incoming edges to a node. To compute the resulting probabilities, we still use the prior distribution over the noise variables  $N$ . Counterfactual reasoning goes one step further. First, it infers the posterior probability  $p^\pi(\{N\} \mid T = \tau)$ , with  $\{N\}$  the set of all noise variables, given the observed trajectory  $\tau$ . Then it performs a Do-intervention as in causal reasoning. However now, to compute the resulting probabilities on the nodes, it uses the posterior noise distribution combined with the modified graph.

One possible strategy to estimate contributions using the above counterfactual reasoning is to explicitly estimate the posterior noise distribution and combine it with forward dynamics models to obtain the counterfactual probability of reaching specific rewarding outcomes. Leveraging the work of Buesing et al. [34] is a promising starting point for this direction of future work. Alternatively, we can avoid explicitly modeling the posterior noise distribution, by leveraging the hindsight distribution combined with the work of Mesnard et al. [3]. Here, we aim to learn a parameterized approximation  $h$  to the counterfactual hindsight distribution  $p^\pi_\tau(A_t = a \mid S_t = s, U' = u')$ , where the  $\tau$ -subscript indicates the counterfactual distribution incorporating the noise posterior. Building upon the approach of Mesnard et al. [3], we can learn  $h(a \mid s', s, \Phi_t(\tau))$  to approximate the counterfactual hindsight distribution, with  $\Phi_t(\tau)$  a summary statistic trained to encapsulate the information of the posterior noise distribution  $p(\{N\} \mid T = \tau)$ . Mesnard et al. [3] show that such a summary statistic  $\Phi_t(\tau)$  can be used to amortize the posterior noise estimation if it satisfies the following two conditions: (i) it needs to provide useful information for predicting the counterfactual hindsight distribution and (ii) it needs to be independent from the action  $A_t$ . We can achieve both characteristics by (i) training  $h(a \mid s', s, \Phi_t(\tau))$  on the hindsight action classification task and backpropagating the gradients to  $\Phi_t(\tau)$ , and (ii) training  $\Phi_t$  on an *independence maximization* loss  $\mathcal{L}_{\text{IM}}(s)$ , which is minimized iff  $A_t$  and  $\Phi_t$  are conditionally independent given  $S_t$ . An example is to minimize the KL divergence between  $\pi(a_t \mid s_t)$  and  $p(a_t \mid s_t, \Phi_t(\tau))$  where the latter can be approximated by training a classifier  $q(a_t \mid s_t, \Phi_t(\tau))$ . Leveraging this approach to extend COCOA towards counterfactual interventions is an exciting direction for future research.

## I Contribution analysis with temporal discounting

As discussed in Appendix B, we can implicitly incorporate discounting into the COCOA framework by adjusting the transition probabilities to have a fixed probability of  $(1 - \gamma)$  of transitioning to the absorbing state  $s_\infty$  at each time step.

We can also readily incorporate explicit time discounting into the COCOA framework, which we discuss here. We consider now a discounted MDP defined as the tuple  $(\mathcal{S}, \mathcal{A}, p, p_r, \gamma)$ , with discount factor  $\gamma \in [0, 1]$ , and  $(\mathcal{S}, \mathcal{A}, p, p_r)$  as defined in Section 2. The discounted value function  $V_\gamma^\pi(s) =$

Table 5: Comparison of discounted policy gradient estimators

Method	Policy gradient estimator ( $\hat{\nabla}_\theta V^\pi(s_0)$ )
REINFORCE	$\sum_{t \geq 0} \gamma^t \nabla_\theta \log \pi(A_t   S_t) \sum_{k \geq 0} \gamma^k R_{t+k}$
Advantage	$\sum_{t \geq 0} \gamma^t \nabla_\theta \log \pi(A_t   S_t) \left( \sum_{k \geq 0} \gamma^k R_{t+k} - V_\gamma^\pi(S_t) \right)$
Q-critic	$\sum_{t \geq 0} \gamma^t \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a   S_t) Q_\gamma^\pi(S_t, a)$
COCOA	$\sum_{t \geq 0} \gamma^t \nabla_\theta \log \pi(A_t   S_t) R_t + \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a   S_t) \sum_{k \geq 1} \gamma^k w_\gamma(S_t, a, U_{t+k}) R_{t+k}$
HCA+	$\sum_{t \geq 0} \gamma^t \nabla_\theta \log \pi(A_t   S_t) R_t + \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a   S_t) \sum_{k \geq 1} \gamma^k w_\gamma(S_t, a, S_{t+k}) R_{t+k}$

1296  $\mathbb{E}_{T \sim \mathcal{T}(s, \pi)} [\sum_{t=0}^{\infty} \gamma^t R_t]$  and action value function  $Q_\gamma^\pi(s, a) = \mathbb{E}_{T \sim \mathcal{T}(s, a, \pi)} [\sum_{t=0}^{\infty} \gamma^t R_t]$  are the  
 1297 expected discounted return when starting from state  $s$ , or state  $s$  and action  $a$  respectively. Table 5  
 1298 shows the policy gradient estimators of  $V_\gamma^\pi(s)$  for REINFORCE, Advantage and Q-critic.

1299 In a discounted MDP, it matters at which point in time we reach a rewarding outcome  $u$ , as the  
 1300 corresponding rewards are discounted. Hence, we adjust the contribution coefficients to

$$w_\gamma(s, a, u') = \frac{\sum_{k \geq 1} \gamma^k p^\pi(U_{t+k} = u' | S_t = s, A_t = a)}{\sum_{k \geq 1} \gamma^k p^\pi(U_{t+k} = u' | S_t = s)} - 1 \quad (122)$$

$$= \frac{p_\gamma^\pi(A_t = a | S_t = s, U' = u')}{\pi(a | s)} - 1 \quad (123)$$

1301 Here, we define the discounted hindsight distribution  $p_\gamma^\pi(A_t = a | S_t = s, U' = u')$  similarly to the  
 1302 undiscounted hindsight distribution explained in Appendix C.1, but now using a different probability  
 1303 distribution on the time steps  $k$ :  $p_\beta(K = k) = (1 - \beta)\beta^{k-1}$ , where we take  $\beta = \gamma$ . We can readily  
 1304 extend Theorems 1-4 to the explicit discounted setting, by taking  $\beta = \gamma$  instead of the limit of  $\beta \rightarrow 1$ ,  
 1305 and using the discounted COCOA policy gradient estimator shown in Table 5.

1306 To approximate the discounted hindsight distribution  $p_\gamma^\pi(A_t = a | S_t = s, U' = u')$ , we need to  
 1307 incorporate the temporal discounting into the classification cross-entropy loss:

$$L_\gamma = \mathbb{E}_\pi \left[ \sum_{t \geq 0} \sum_{k \geq 1} \gamma^k CE(h_\gamma(\cdot | S_t, U_{t+k}), \delta(a = A_t)) \right], \quad (124)$$

1308 with  $CE$  the cross-entropy loss,  $h_\gamma(\cdot | S_t, U_{t+k})$  the classification model that approximates the  
 1309 discounted hindsight distribution, and  $\delta(a = A_t)$  a one-hot encoding of  $A_t$ .

## 1310 J Bootstrapping with COCOA

1311 Here, we show how COCOA can be combined with  $n$ -step returns, and we make a correction to  
 1312 Theorem 7 of Harutyunyan et al. [1] which considers  $n$ -step returns for HCA.

1313 Consider the graphical model of Fig. 12a where we model time,  $K$ , as a separate node in the graphical  
 1314 model (c.f. App. C.1). To model  $n$ -step returns, we now define the following prior probability on  $K$ ,  
 1315 parameterized by  $\beta$ :

$$p_{n, \beta}(K = k) = \begin{cases} \frac{\beta^k}{z} & \text{if } 1 \leq k \leq n-1 \\ 0 & \text{else} \end{cases} \quad z = \beta \frac{1 - \beta^{n-1}}{1 - \beta} \quad (125)$$

1316 Using  $p_{n, \beta}^\pi$  as the probability distribution induced by this graphical model, we have that

$$p_{n, \beta}^\pi(U' = u | S = s, A = a) = \sum_k p_{n, \beta}^\pi(U' = u, K = k | S = s, A = a) \quad (126)$$

$$= \sum_k p_{n, \beta}(K = k) p_{n, \beta}^\pi(U' = u | S = s, A = a, K = k) \quad (127)$$

$$= \frac{1 - \beta}{\beta(1 - \beta^{n-1})} \sum_{k=1}^{n-1} \beta^k p^\pi(U_k = u | S_0 = s, A_0 = a) \quad (128)$$

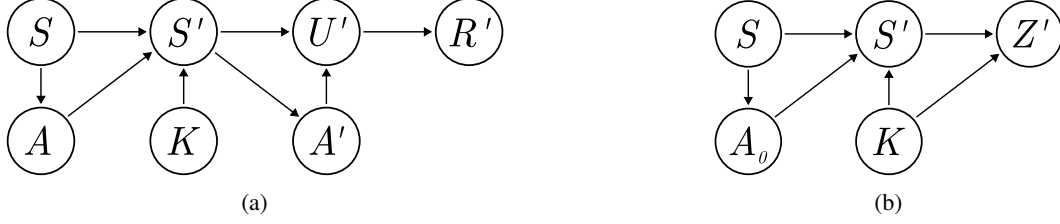


Figure 12: (a) Graphical model where we abstract time. (b) Graphical model implicitly used in the proof of Theorem 7 in Harutyunyan et al. [1].

1317 We introduce the following contribution coefficients that we will use for n-step returns

$$w_{n,\beta}(s, a, u') = \frac{\sum_{k=1}^{n-1} \beta^k p^\pi(U_k = u' | S_0 = s, A_0 = a)}{\sum_{k=1}^{n-1} \beta^k p^\pi(U_k = u' | S_0 = s)} \quad (129)$$

$$= \frac{p_{n,\beta}^\pi(U' = u' | S = s, A = a)}{p_{n,\beta}^\pi(U' = u' | S = s)} - 1 = \frac{p_{n,\beta}^\pi(A = a | S = s, U' = u')}{\pi(a | s)} - 1 \quad (130)$$

$$w_n(s, a, s') = \frac{p^\pi(S_n = s' | S_0 = s, A_0 = a)}{p^\pi(S_n = s' | S_0 = s)} - 1 = \frac{p^\pi(A = a | S = s, S_n = s')}{\pi(a | s)} - 1 \quad (131)$$

1318 Now we are ready to prove the n-step return theorem for the discounted MDP setting. We can recover  
1319 the undiscounted setting by taking the limit of  $\gamma \rightarrow 1_-$ .

1320 **Theorem 10.** Consider state  $s$  and action  $a$  for which it holds that  $\pi(a | s) > 0$  and take  $\beta$  equal to the  
1321 discount factor  $\gamma \in [0, 1]$ . Furthermore, assume that the rewarding outcome encoding  $u = f(s, a, r)$   
1322 is fully predictive of the reward (c.f. Definition 1). Then the advantage  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$   
1323 is equal to

$$A^\pi(s, a) = r(s, a) - r^\pi(s) + \mathbb{E}_{\mathcal{T}(s, \pi)} \left[ \sum_{k=1}^{n-1} \gamma^k w_{n,\beta}(s, a, U_k) R_k + \gamma^n w_n(s, a, S_n) V^\pi(S_n) \right] \quad (132)$$

1324 with  $r(s, a)$  the reward model and  $r^\pi(s) = \sum_a \pi(a | s) r(s, a)$ .

1325 *Proof.* We start with the action-value function  $Q^\pi$ , and will subtract the value  $V^\pi$  to obtain the result  
1326 on the advantage function.

$$Q(s, a) = \mathbb{E}_{T \sim \mathcal{T}(s, a, \pi)} \left[ \sum_{k \geq 1} \gamma^k R_k \right] \quad (133)$$

$$= r(s, a) + \sum_{r' \in \mathcal{R}} \sum_{k \geq 1} \gamma^k p^\pi(R_k = r' | s, a) r' \quad (134)$$

$$= r(s, a) + \sum_{r' \in \mathcal{R}} \sum_{u' \in \mathcal{U}} \sum_{k=1}^{n-1} \gamma^k p^\pi(R_k = r', U_k = u' | s, a) r' + \quad (135)$$

$$\gamma^n \sum_{s' \in \mathcal{S}} p^\pi(S_n = s' | s, a) V^\pi(s') \quad (136)$$

$$= r(s, a) + \sum_{r' \in \mathcal{R}} \sum_{u' \in \mathcal{U}} \sum_{k=1}^{n-1} \gamma^k p^\pi(R' = r' | U' = u') p^\pi(U_k = u' | s, a) r' + \quad (137)$$

$$\gamma^n \sum_{s' \in \mathcal{S}} p^\pi(S_n = s' | s, a) V^\pi(s') \quad (138)$$

$$= r(s, a) + \sum_{u' \in \mathcal{U}} r(u') \sum_{k=1}^{n-1} \gamma^k p^\pi(U_k = u' | s, a) + \gamma^n \sum_{s' \in \mathcal{S}} p^\pi(S_n = s' | s, a) V^\pi(s') \quad (139)$$

$$= r(s, a) + \sum_{u' \in \mathcal{U}} r(u') \sum_{k=1}^{n-1} \gamma^k p^\pi(U_k = u' | s) \frac{\sum_{k'=1}^{n-1} \gamma^{k'} p^\pi(U_{k'} = u | s, a)}{\sum_{k'=1}^{n-1} \gamma^{k'} p^\pi(U_{k'} = u | s)} + \quad (140)$$

$$\gamma^n \sum_{s' \in \mathcal{S}} p^\pi(S_n = s' | s) \frac{p^\pi(S_n = s' | s, a)}{p^\pi(S_n = s' | s)} V^\pi(s') \quad (141)$$

$$= r(s, a) + \sum_{u' \in \mathcal{U}} r(u') \sum_{k=1}^{n-1} \gamma^k p^\pi(U_k = u' | s) (w_{n,\beta}(s, a, u') + 1) \quad (142)$$

$$+ \gamma^n \sum_{s' \in \mathcal{S}} p^\pi(S_n = s' | s) (w_n(s, a, s') + 1) V^\pi(s') \quad (143)$$

$$(144)$$

1327 where we use that  $U'$  is fully predictive of the reward  $R'$ , and define  $r(u') = \sum_{r' \in \mathcal{R}} p(R' = r' |$   
 1328  $U' = u')r'$  By subtracting the value function, we get

$$A(s, a) = r(s, a) - r^\pi(s) + \sum_{u' \in \mathcal{U}} r(u') \sum_{k=1}^{n-1} \gamma^k p^\pi(U_k = u' | s) w_{n,\beta}(s, a, u') \quad (145)$$

$$+ \gamma^n \sum_{s' \in \mathcal{S}} p^\pi(S_n = s' | s) w_n(s, a, s') V^\pi(s') \quad (146)$$

$$= r(s, a) - r^\pi(s) + \mathbb{E}_{\mathcal{T}(s, \pi)} \left[ \sum_{k=1}^{n-1} \gamma^k w_{n,\beta}(s, a, U_k) R_k + \gamma^n w_n(s, a, S_n) V^\pi(S_n) \right] \quad (147)$$

1329

□

1330 Finally, we can sample from this advantage function to obtain an n-step COCOA gradient estimator,  
 1331 akin to Theorem 1.

1332 Note that we require to learn the state-based contribution coefficients  $w_n(s, a, s')$  to bootstrap the  
 1333 value function into the n-step return, as the value function requires a Markov state  $s'$  as input instead  
 1334 of a rewarding outcome encoding  $u'$ . Unfortunately, these state-based contribution coefficients will  
 1335 suffer from spurious contributions, akin to HCA, introducing a significant amount of variance into  
 1336 the n-step COCOA gradient estimator. We leave it to future research to investigate whether we  
 1337 can incorporate value functions into an n-step return, while using rewarding-outcome contribution  
 1338 coefficients  $w(s, a, u')$  instead of state-based contribution coefficients  $w_n(s, a, s')$ .

1339 **Learning the contribution coefficients.** We can learn the contribution coefficients  $w_{\beta,n}(s, a, u')$   
 1340 with the same strategies as described in Section 3, but now with training data from  $n$ -step trajectories  
 1341 instead of complete trajectories. If we use a discount  $\gamma \neq 1$ , we need to take this discount factor into  
 1342 account in the training distribution or loss function (c.f. App. I).

1343 **Correction to Theorem 7 of Harutyunyan et al. [1].** Harutyunyan et al. [1] propose a theorem  
 1344 similar to Theorem 10, with two important differences. The first one concerns the distribution on  $K$   
 1345 in the graphical model of Fig. 12a. Harutyunyan et al. [1] implicitly use this graphical model, but  
 1346 with a different prior probability distribution on  $K$ :

$$p_{n,\beta}^{HCA}(K = k) = \begin{cases} \beta^{k-1}(1 - \beta) & \text{if } 1 \leq k \leq n - 1 \\ \beta^{n-1} & \text{if } k = n \\ 0 & \text{else} \end{cases} \quad (148)$$

1347 The graphical model combined with the distribution on  $K$  defines the hindsight distribution  
 1348  $p_{n,\beta,HCA}^\pi(A = a | S = s, S' = s')$ . The second difference is the specific  $Q$ -value estimator Haru-  
 1349 tyunyan et al. [1] propose. They use the hindsight distribution  $p_{n,\beta,HCA}^\pi(A = a | S = s, S' = s')$  in  
 1350 front of the value function (c.f. Theorem 10), which considers that  $s'$  can be reached at any time step  
 1351  $k \sim p_{n,\beta}^{HCA}(k)$ , whereas Theorem 10 uses  $w_n(s, a, s')$  which considers that  $s'$  is reached exactly at  
 1352 time step  $k = n$ .

1353 To the best of our knowledge, there is an error in the proposed proof of Theorem 7 by Harutyunyan et al.  
 1354 [1] for which we could not find a simple fix. For the interested reader, we briefly explain the error. One

1355 indication of the problem is that for  $\beta \rightarrow 1$ , all the probability mass of  $p_{n,\beta}^{HCA}(K = k)$  is concentrated  
 1356 at  $k = n$ , hence the corresponding hindsight distribution  $p_{n,\beta,HCA}^\pi(A = a \mid S = s, S' = s')$   
 1357 considers only hindsight states  $s'$  encountered at time  $k = n$ . While this is not a mathematical error, it  
 1358 does not correspond to the intuition of a ‘time independent hindsight distribution’ the authors provide.  
 1359 In the proof itself, a conditional independence relation is assumed that does not hold. The authors  
 1360 introduce a helper variable  $Z$  defined on the state space  $\mathcal{S}$ , with a conditional distribution

$$\mu_k(Z = z \mid S' = s') = \begin{cases} \delta(z = s') & \text{if } 1 \leq k \leq n - 1 \\ \tilde{d}^\pi(z \mid s') & \text{if } k = n \end{cases} \quad (149)$$

1361 with the normalized discounted visit distribution  $\tilde{d}^\pi(z \mid s') = (1 - \gamma) \sum_k \gamma^k p^\pi(S_k = z \mid S_0 = s)$ .  
 1362 We can model this setting as the graphical model visualized in Fig. 12b. In the proof (last line on page  
 1363 15 in the supplementary materials of Harutyunyan et al. [1]), the following conditional independence  
 1364 is used:

$$p^\pi(A_0 = a \mid S_0 = s, S' = s', Z = z) = p^\pi(A_0 = a \mid S_0 = s, S' = s') \quad (150)$$

1365 However, Fig. 12b shows that  $S'$  is a collider on the path  $A_0 \rightarrow S' \leftarrow K \rightarrow Z$ . Hence, by  
 1366 conditioning on  $S'$  we open this collider path, making  $A_0$  dependent on  $Z$  conditioned on  $S_0$  and  $S'$ ,  
 1367 thereby invalidating the assumed conditional independence. For example, if  $Z$  is different from  $S'$ ,  
 1368 we know that  $K = n$  (c.f. Eq. 149), hence  $Z$  can contain information about action  $A_0$ , beyond  $S'$ , as  
 1369  $S'$  ignores at which point in time  $s'$  is encountered.

## 1370 **K Additional details**

### 1371 **K.1 Compute resources**

1372 We used Linux workstations with Nvidia RTX 2080 and Nvidia RTX 3090 GPUs for development  
 1373 and conducted hyperparameter searches and experiments using 5 TPUv2-8, 5 TPUv3-8 and 1 Linux  
 1374 server with 8 Nvidia RTX 3090 GPUs over the course of 9 months. All of the final experiments  
 1375 presented take less than a few hours to complete using a single Nvidia RTX 3090 GPU. In total, we  
 1376 spent an estimated amount of 2 GPU months.

### 1377 **K.2 Software and libraries**

1378 For the results produced in this paper we relied on free and open-source software. We implemented  
 1379 our experiments in Python using JAX [93, Apache License 2.0] and the Deepmind Jax Ecosystem  
 1380 [80, Apache License 2.0]. For experiment tracking we used wandb [94, MIT license] and for the  
 1381 generation of plots we used plotly [95, MIT license].

## References

- [1] Anna Harutyunyan, Will Dabney, Thomas Mesnard, Mohammad Gheshlaghi Azar, Bilal Piot, Nicolas Heess, Hado P van Hasselt, Gregory Wayne, Satinder Singh, Doina Precup, and Remi Munos. Hindsight Credit Assignment. In H. Wallach, H. Larochelle, A. Beygelzimer, F Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 12488–12497. Curran Associates, Inc., 2019.
- [2] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning, second edition: An Introduction*. MIT Press, November 2018.
- [3] Thomas Mesnard, Théophane Weber, Fabio Viola, Shantanu Thakoor, Alaa Saade, Anna Harutyunyan, Will Dabney, Tom Stepleton, Nicolas Heess, Arthur Guez, Marcus Hutter, Lars Buesing, and Rémi Munos. Counterfactual Credit Assignment in Model-Free Reinforcement Learning. In *Proceedings of the 38 th International Conference on Machine Learning, PMLR 139*, 2021.
- [4] Chia-Chun Hung, Timothy Lillicrap, Josh Abramson, Yan Wu, Mehdi Mirza, Federico Carnevale, Arun Ahuja, and Greg Wayne. Optimizing agent behavior over long time scales by transporting value. *Nature Communications*, 10(1):5223, November 2019. Number: 1 Publisher: Nature Publishing Group.
- [5] Jose A. Arjona-Medina, Michael Gillhofer, Michael Widrich, Thomas Unterthiner, Johannes Brandstetter, and Sepp Hochreiter. RUDDER: Return Decomposition for Delayed Rewards. In H. Wallach, H. Larochelle, A. Beygelzimer, F. Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 13566–13577. Curran Associates, Inc., 2019.
- [6] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3):229–256, May 1992.
- [7] Evan Greensmith, Peter L. Bartlett, and Jonathan Baxter. Variance Reduction Techniques for Gradient Estimates in Reinforcement Learning. *Journal of Machine Learning Research*, 5(Nov):1471–1530, 2004.
- [8] Philip Thomas. Bias in Natural Actor-Critic Algorithms. In *Proceedings of the 31st International Conference on Machine Learning*, pages 441–448. PMLR, January 2014. ISSN: 1938-7228.
- [9] Sham Kakade. Optimizing Average Reward Using Discounted Rewards. In David Helmbold and Bob Williamson, editors, *Computational Learning Theory*, Lecture Notes in Computer Science, pages 605–615, Berlin, Heidelberg, 2001. Springer.
- [10] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-Dimensional Continuous Control Using Generalized Advantage Estimation. In *International Conference on Learning Representations*. arXiv, 2016. arXiv:1506.02438 [cs].
- [11] Peter Marbach and John Tsitsiklis. Approximate Gradient Methods in Policy-Space Optimization of Markov Reward Processes. *Discrete Event Dynamic Systems*, page 38, 2003.
- [12] OpenAI, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębniak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique P. d O. Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with Large Scale Deep Reinforcement Learning, December 2019. arXiv:1912.06680 [cs, stat].
- [13] OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider, Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba, and Lei Zhang. Solving Rubik’s Cube with a Robot Hand, October 2019. arXiv:1910.07113 [cs, stat].
- [14] Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P.

- 1433 Agapiou, Max Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin  
1434 Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu  
1435 Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKin-  
1436 ney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris  
1437 Apps, and David Silver. Grandmaster level in StarCraft II using multi-agent reinforcement  
1438 learning. *Nature*, 575(7782):350–354, November 2019. Number: 7782 Publisher: Nature  
1439 Publishing Group.
- 1440 [15] Richard S. Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM*  
1441 *SIGART Bulletin*, 2(4):160–163, July 1991.
- 1442 [16] Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-  
1443 Efficient Reinforcement Learning with Stochastic Ensemble Value Expansion. In *32nd Confer-*  
1444 *ence on Neural Information Processing Systems*. arXiv, 2018. arXiv:1807.01675 [cs, stat].
- 1445 [17] Thanard Kurutach, Ignasi Clavera, Yan Duan, Aviv Tamar, and Pieter Abbeel. Model-Ensemble  
1446 Trust-Region Policy Optimization. 2018.
- 1447 [18] Nicolas Heess, Gregory Wayne, David Silver, Timothy Lillicrap, Tom Erez, and Yuval Tassa.  
1448 Learning Continuous Control Policies by Stochastic Value Gradients. In *Advances in Neural*  
1449 *Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- 1450 [19] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to Control:  
1451 Learning Behaviors by Latent Imagination. In *International Conference on Learning Represen-*  
1452 *tations*, March 2020. Number: arXiv:1912.01603 arXiv:1912.01603 [cs].
- 1453 [20] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering Atari  
1454 with Discrete World Models. In *International Conference on Learning Representations*, 2021.  
1455 Number: arXiv:2010.02193 arXiv:2010.02193 [cs, stat].
- 1456 [21] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering Diverse Domains  
1457 through World Models, January 2023. arXiv:2301.04104 [cs, stat].
- 1458 [22] Mikael Henaff, William F. Whitney, and Yann LeCun. Model-Based Planning with Discrete  
1459 and Continuous Actions, April 2018. arXiv:1705.07177 [cs].
- 1460 [23] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap,  
1461 Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous Methods for Deep Reinforce-  
1462 ment Learning. In *Proceedings of The 33rd International Conference on Machine Learning*,  
1463 pages 1928–1937. PMLR, June 2016. ISSN: 1938-7228.
- 1464 [24] Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy Gradient  
1465 Methods for Reinforcement Learning with Function Approximation. In *Advances in Neural*  
1466 *Information Processing Systems 12*, 1999.
- 1467 [25] Chris Nota and Philip S Thomas. Is the Policy Gradient a Gradient? In *Proc. of the 19th*  
1468 *International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2020)*,  
1469 page 9, 2020.
- 1470 [26] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-  
1471 Dimensional Continuous Control Using Generalized Advantage Estimation, October 2018.  
1472 arXiv:1506.02438 [cs].
- 1473 [27] Vyacheslav Alipov, Riley Simmons-Edler, Nikita Putintsev, Pavel Kalinin, and Dmitry Vetrov.  
1474 Towards Practical Credit Assignment for Deep Reinforcement Learning, February 2022.  
1475 arXiv:2106.04499 [cs].
- 1476 [28] Peter Dayan. Improving Generalization for Temporal Difference Learning: The Successor  
1477 Representation. *Neural Computation*, 5(4):613–624, July 1993. Conference Name: Neural  
1478 Computation.
- 1479 [29] Tejas D. Kulkarni, Ardavan Saeedi, Simanta Gautam, and Samuel J. Gershman. Deep Successor  
1480 Reinforcement Learning, June 2016. Number: arXiv:1606.02396 arXiv:1606.02396 [cs, stat].



- [30] David Raposo, Sam Ritter, Adam Santoro, Greg Wayne, Theophane Weber, Matt Botvinick, Hado van Hasselt, and Francis Song. Synthetic Returns for Long-Term Credit Assignment. *arXiv:2102.12425 [cs]*, February 2021. arXiv: 2102.12425.
- [31] John W. Roberts and Russ Tedrake. Signal-to-Noise Ratio Analysis of Policy Gradient Algorithms. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21*, pages 1361–1368. Curran Associates, Inc., 2009.
- [32] Vihang P. Patil, Markus Hofmarcher, Marius-Constantin Dinu, Matthias Dorfer, Patrick M. Blies, Johannes Brandstetter, Jose A. Arjona-Medina, and Sepp Hochreiter. Align-RUDDER: Learning From Few Demonstrations by Reward Redistribution. *Proceedings of Machine Learning Research*, 162, 2022. arXiv: 2009.14108.
- [33] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. Adaptive Computation and Machine Learning. November 2018.
- [34] Lars Buesing, Theophane Weber, Yori Zwols, Sebastien Racaniere, Arthur Guez, Jean-Baptiste Lespiau, and Nicolas Heess. Woulda, Coulda, Shoulda: Counterfactually-Guided Policy Search. In *International Conference on Learning Representations*, 2019. arXiv: 1811.06272.
- [35] Ioana Bica, Ahmed M. Alaa, James Jordon, and Mihaela van der Schaar. Estimating Counterfactual Treatment Outcomes over Time Through Adversarially Balanced Representations. In *International Conference on Learning Representations*. arXiv, February 2020. arXiv:2002.04083 [cs, stat].
- [36] Pushi Zhang, Li Zhao, Guoqing Liu, Jiang Bian, Minlie Huang, Tao Qin, and Tie-Yan Liu. Independence-aware Advantage Estimation. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*, pages 3349–3355, Montreal, Canada, August 2021. International Joint Conferences on Artificial Intelligence Organization.
- [37] Kenny Young. Variance Reduced Advantage Estimation with  $\Delta$  Hindsight Credit Assignment. *arXiv:1911.08362 [cs]*, September 2020. arXiv: 1911.08362.
- [38] Michel Ma and Bacon Pierre-Luc. Counterfactual Policy Evaluation and the Conditional Monte Carlo Method. In *Offline Reinforcement Learning Workshop, NeurIPS*, 2020.
- [39] Paul Bratley, Bennett L. Fox, and Linus E. Schrage. *A Guide to Simulation*. Springer, New York, NY, 1987.
- [40] J. M. Hammersley. Conditional Monte Carlo. *Journal of the ACM*, 3(2):73–76, April 1956.
- [41] Dilip Arumugam, Peter Henderson, and Pierre-Luc Bacon. An Information-Theoretic Perspective on Credit Assignment in Reinforcement Learning. *arXiv:2103.06224 [cs, math]*, March 2021. arXiv: 2103.06224.
- [42] Kenny Young. Hindsight Network Credit Assignment: Efficient Credit Assignment in Networks of Discrete Stochastic Units. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(8):8919–8926, June 2022. Number: 8.
- [43] Théophane Weber, Nicolas Heess, Lars Buesing, and David Silver. Credit Assignment Techniques in Stochastic Computation Graphs. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, pages 2650–2660. PMLR, April 2019. ISSN: 2640-3498.
- [44] Jonathan Baxter and Peter L. Bartlett. Infinite-Horizon Policy-Gradient Estimation. *Journal of Artificial Intelligence Research*, 15:319–350, November 2001.
- [45] Shixiang Gu, Timothy Lillicrap, Zoubin Ghahramani, Richard E. Turner, and Sergey Levine. Q-Prop: Sample-Efficient Policy Gradient with An Off-Policy Critic. 2017.
- [46] Philip S. Thomas and Emma Brunskill. Policy Gradient Methods for Reinforcement Learning with Function Approximation and Action-Dependent Baselines, June 2017. arXiv:1706.06643 [cs].
- [47] Hao Liu\*, Yihao Feng\*, Yi Mao, Dengyong Zhou, Jian Peng, and Qiang Liu. Action-dependent Control Variates for Policy Optimization via Stein Identity. February 2022.

- [48] Cathy Wu, Aravind Rajeswaran, Yan Duan, Vikash Kumar, Alexandre M. Bayen, Sham Kakade, Igor Mordatch, and Pieter Abbeel. Variance Reduction for Policy Gradient with Action-Dependent Factorized Baselines. February 2022.
- [49] George Tucker, Surya Bhupatiraju, Shixiang Gu, Richard Turner, Zoubin Ghahramani, and Sergey Levine. The Mirage of Action-Dependent Baselines in Reinforcement Learning. In *Proceedings of the 35th International Conference on Machine Learning*, pages 5015–5024. PMLR, July 2018. ISSN: 2640-3498.
- [50] Chris Nota, Philip Thomas, and Bruno C. Da Silva. Posterior Value Functions: Hindsight Baselines for Policy Gradient Methods. In *Proceedings of the 38th International Conference on Machine Learning*, pages 8238–8247. PMLR, July 2021. ISSN: 2640-3498.
- [51] Arthur Guez, Fabio Viola, Theophane Weber, Lars Buesing, Steven Kapturowski, Doina Precup, David Silver, and Nicolas Heess. Value-driven Hindsight Modelling. *Advances in Neural Information Processing Systems*, 33, 2020.
- [52] David Venuto, Elaine Lau, Doina Precup, and Ofir Nachum. Policy Gradients Incorporating the Future. January 2022.
- [53] Andrew Y. Ng, Daishi Harada, and Stuart J. Russell. Policy Invariance Under Reward Transformations: Theory and Application to Reward Shaping. In *Proceedings of the Sixteenth International Conference on Machine Learning*, ICML ’99, pages 278–287, San Francisco, CA, USA, June 1999. Morgan Kaufmann Publishers Inc.
- [54] Jürgen Schmidhuber. Formal Theory of Creativity, Fun, and Intrinsic Motivation (1990–2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, September 2010.
- [55] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying Count-Based Exploration and Intrinsic Motivation. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [56] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. December 2018.
- [57] Ofir Marom and Benjamin Rosman. Belief Reward Shaping in Reinforcement Learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), April 2018. Number: 1.
- [58] Farzan Memarian, Wonjoon Goo, Rudolf Lioutikov, Scott Niekum, and Ufuk Topcu. Self-Supervised Online Reward Shaping in Sparse-Reward Environments, July 2021. arXiv:2103.04529 [cs].
- [59] Halit Bener Suay, Tim Brys, Matthew E. Taylor, and Sonia Chernova. Learning from Demonstration for Shaping through Inverse Reinforcement Learning. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*, AAMAS ’16, pages 429–437, Richland, SC, May 2016. International Foundation for Autonomous Agents and Multiagent Systems.
- [60] Yuchen Wu, Melissa Mozifian, and Florian Shkurti. Shaping Rewards for Reinforcement Learning with Imperfect Demonstrations using Generative Models, November 2020. arXiv:2011.01298 [cs].
- [61] Johan Ferret, Raphaël Marinier, Matthieu Geist, and Olivier Pietquin. Self-Attentional Credit Assignment for Transfer in Reinforcement Learning. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, pages 2655–2661, July 2020. arXiv:1907.08027 [cs].
- [62] Zhizhou Ren, Ruihan Guo, Yuan Zhou, and Jian Peng. Learning Long-Term Reward Redistribution via Randomized Return Decomposition. January 2022.
- [63] Yonathan Efroni, Nadav Merlis, and Shie Mannor. Reinforcement Learning with Trajectory Feedback. In *The Thirty-Fifth AAAI Conference on Artificial Intelligence*. arXiv, March 2021.
- [64] Minah Seo, Luiz Felipe Vecchietti, Sangkeum Lee, and Dongsoo Har. Rewards Prediction-Based Credit Assignment for Reinforcement Learning With Sparse Binary Rewards. *IEEE*

- 1580 ACCESS, 7:118776–118791, 2019. Accepted: 2019-09-24T11:21:52Z Publisher: IEEE-INST  
1581 ELECTRICAL ELECTRONICS ENGINEERS INC.
- 1582 [65] Markel Sanz Ausin, Hamoon Azizsoltani, Song Ju, Yeo Jin Kim, and Min Chi. InferNet for  
1583 Delayed Reinforcement Tasks: Addressing the Temporal Credit Assignment Problem. In *2021*  
1584 *IEEE International Conference on Big Data (Big Data)*, pages 1337–1348, December 2021.
- 1585 [66] Hamoon Azizsoltani, Yeo Jin Kim, Markel Sanz Ausin, Tiffany Barnes, and Min Chi. Unob-  
1586 served Is Not Equal to Non-existent: Using Gaussian Processes to Infer Immediate Rewards  
1587 Across Contexts. In *Proceedings of the Twenty-Eighth International Joint Conference on*  
1588 *Artificial Intelligence*, pages 1974–1980, 2019.
- 1589 [67] John Schulman, Nicolas Heess, Theophane Weber, and Pieter Abbeel. Gradient Estimation  
1590 Using Stochastic Computation Graphs. In *Advances in Neural Information Processing Systems*,  
1591 volume 28. Curran Associates, Inc., 2015.
- 1592 [68] Michel Ma, Pierluca D’Oro, Yoshua Bengio, and Pierre-Luc Bacon. Long-Term Credit As-  
1593 signment via Model-based Temporal Shortcuts. In *Deep RL Workshop NeurIPS*, October  
1594 2021.
- 1595 [69] Nan Rosemary Ke, Anirudh Goyal ALIAS PARTH GOYAL, Olexa Bilaniuk, Jonathan Binas,  
1596 Michael C Mozer, Chris Pal, and Yoshua Bengio. Sparse Attentive Backtracking: Temporal  
1597 Credit Assignment Through Reminding. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,  
1598 N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*  
1599 *31*, pages 7640–7651. Curran Associates, Inc., 2018.
- 1600 [70] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal Value Function  
1601 Approximators. In *Proceedings of the 32nd International Conference on Machine Learning*,  
1602 pages 1312–1320. PMLR, June 2015. ISSN: 1938-7228.
- 1603 [71] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder,  
1604 Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight Experience  
1605 Replay. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,  
1606 Inc., 2017.
- 1607 [72] Paulo Rauber, Avinash Ummadisingu, Filipe Mutz, and Jürgen Schmidhuber. Hindsight policy  
1608 gradients. December 2018.
- 1609 [73] Anirudh Goyal, Philemon Brakel, William Fedus, Soumye Singhal, Timothy Lillicrap, Sergey  
1610 Levine, Hugo Larochelle, and Yoshua Bengio. Recall Traces: Backtracking Models for Efficient  
1611 Reinforcement Learning. In *International Conference on Learning Representations*. arXiv,  
1612 January 2019. arXiv:1804.00379 [cs, stat].
- 1613 [74] Juergen Schmidhuber. Reinforcement Learning Upside Down: Don’t Predict Rewards – Just  
1614 Map Them to Actions, June 2020. arXiv:1912.02875 [cs].
- 1615 [75] Rupesh Kumar Srivastava, Pranav Shyam, Filipe Mutz, Wojciech Jaśkowski, and Jürgen  
1616 Schmidhuber. Training Agents using Upside-Down Reinforcement Learning, September 2021.  
1617 arXiv:1912.02877 [cs].
- 1618 [76] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter  
1619 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision Transformer: Reinforcement Learning  
1620 via Sequence Modeling. In *Advances in Neural Information Processing Systems*, volume 34,  
1621 pages 15084–15097. Curran Associates, Inc., 2021.
- 1622 [77] Michael Janner, Qiyang Li, and Sergey Levine. Offline Reinforcement Learning as One  
1623 Big Sequence Modeling Problem. In *Advances in Neural Information Processing Systems*,  
1624 volume 34, pages 1273–1286. Curran Associates, Inc., 2021.
- 1625 [78] Richard S. Sutton. Learning to predict by the methods of temporal differences. *Machine*  
1626 *Learning*, 3(1):9–44, August 1988.
- 1627 [79] Hado van Hasselt, Sephora Madjiheurem, Matteo Hessel, David Silver, André Barreto, and  
1628 Diana Borsa. Expected Eligibility Traces. In *Association for the Advancement of Artificial*  
1629 *Intelligence*. arXiv, February 2021. arXiv:2007.01839 [cs, stat].

- [80] Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter Buchlovsky, David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Antoine Dedieu, Claudio Fantacci, Jonathan Godwin, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel, Shaobo Hou, Steven Kapturowski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch, Lena Martens, Hamza Merzic, Vladimir Mikulik, Tamara Norman, George Papamakarios, John Quan, Roman Ring, Francisco Ruiz, Alvaro Sanchez, Rosalia Schneider, Eren Sezener, Stephen Spencer, Srivatsan Srinivasan, Wojciech Stokowiec, Luyu Wang, Guangyao Zhou, and Fabio Viola. The DeepMind JAX Ecosystem, 2020.
- [81] Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep Variational Information Bottleneck. 2017.
- [82] Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method, April 2000. arXiv:physics/0004057.
- [83] Matthew Hausknecht and Peter Stone. Deep Recurrent Q-Learning for Partially Observable MDPs. In *Association for the Advancement of Artificial Intelligence*. arXiv, 2015. arXiv:1507.06527 [cs].
- [84] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning Latent Dynamics for Planning from Pixels. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2555–2565. PMLR, May 2019. ISSN: 2640-3498.
- [85] Karol Gregor, Danilo Jimenez Rezende, Frederic Besse, Yan Wu, Hamza Merzic, and Aaron van den Oord. Shaping Belief States with Generative Environment Models for RL. In *33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*. arXiv, June 2019. Number: arXiv:1906.09237 arXiv:1906.09237 [cs, stat].
- [86] Karol Gregor, George Papamakarios, Frederic Besse, Lars Buesing, and Theophane Weber. Temporal Difference Variational Auto-Encoder. 2019.
- [87] Matthijs T J Spaan. Partially Observable Markov Decision Processes. *Reinforcement Learning*, page 27.
- [88] K. J. Astrom. Optimal Control of Markov decision processes with incomplete state estimation. *J. Math. Anal. Applic.*, 10:174–205, 1965.
- [89] Edward C. Tolman. Cognitive maps in rats and men. *Psychological Review*, 55:189–208, 1948. Place: US Publisher: American Psychological Association.
- [90] Alexander A. Alemi, Ben Poole, Ian Fischer, Joshua V. Dillon, Rif A. Saurous, and Kevin Murphy. Fixing a Broken ELBO. In *Proceedings of the 35th International Conference on Machine Learning*,. arXiv, February 2018. arXiv:1711.00464 [cs, stat].
- [91] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. July 2022.
- [92] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. arXiv, May 2014. Number: arXiv:1312.6114 arXiv:1312.6114 [cs, stat].
- [93] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- [94] Lukas Biewald. Experiment Tracking with Weights and Biases, 2020.
- [95] Plotly Technologies Inc. Collaborative data science, 2015. Place: Montreal, QC Publisher: Plotly Technologies Inc.