
Counterfactual Credit Assignment in Model-Free Reinforcement Learning

Thomas Mesnard^{*1} Théophile Weber^{*1} Fabio Viola¹ Shantanu Thakoor¹ Alaa Saade¹
Anna Harutyunyan¹ Will Dabney¹ Tom Stepleton¹ Nicolas Heess¹ Arthur Guez¹ Éric Moulines²
Marcus Hutter¹ Lars Buesing¹ Rémi Munos¹

Abstract

Credit assignment in reinforcement learning is the problem of measuring an action’s influence on future rewards. In particular, this requires separating *skill* from *luck*, i.e. disentangling the effect of an action on rewards from that of external factors and subsequent actions. To achieve this, we adapt the notion of counterfactuals from causality theory to a model-free RL setup. The key idea is to condition value functions on *future* events, by learning to extract relevant information from a trajectory. We formulate a family of policy gradient algorithms that use these future-conditional value functions as baselines or critics, and show that they are provably low variance. To avoid the potential bias from conditioning on future information, we constrain the hindsight information to not contain information about the agent’s actions. We demonstrate the efficacy and validity of our algorithm on a number of illustrative and challenging problems.

1. Introduction

Reinforcement learning (RL) agents act in their environments and learn to achieve desirable outcomes by maximizing a reward signal. A key difficulty is the problem of *credit assignment* (Minsky, 1961), i.e. to understand the relation between actions and outcomes, and to determine to what extent an outcome was caused by external, uncontrollable factors. In doing so we aim to disentangle the relative aspects of ‘skill’ and ‘luck’ in an agent’s performance. One possible solution to this problem is for the agent to build a model of the environment, and use it to obtain a more fine-

grained understanding of the effects of an action. While this topic has recently generated a lot of interest (Heess et al., 2015; Ha & Schmidhuber, 2018; Hamrick, 2019; Kaiser et al., 2019; Schrittwieser et al., 2019), it remains difficult to model complex, partially observed environments.

In contrast, model-free reinforcement learning algorithms such as policy gradient methods (Williams, 1992; Sutton et al., 2000) perform simple time-based credit assignment, where events and rewards happening after an action are credited to that action, *post hoc ergo propter hoc*. While unbiased in expectation, this coarse-grained credit assignment typically has high variance, and the agent will require a large amount of experience to learn the correct relation between actions and rewards. Another issue is that existing model-free methods are not capable of *counterfactual reasoning*, i.e. reasoning about what would have happened had different actions been taken *with everything else remaining the same*. Given a trajectory, model-free methods can in fact only learn about the actions that were actually taken to produce the data, and this limits the ability of the agent to learn efficiently.

As environments grow in complexity due to partial observability, scale, long time horizons, and increasing number of agents, actions taken by an agent will only affect a vanishing part of the outcome, making it increasingly difficult to learn from classical reinforcement learning algorithms. We need better credit assignment techniques.

In this paper, we investigate a new method of credit assignment for model-free reinforcement learning which we call *Counterfactual Credit Assignment* (CCA). CCA leverages *hindsight* information to implicitly perform counterfactual evaluation—an estimate of the return for actions other than the ones which were chosen. These counterfactual returns can be used to form unbiased and lower variance estimates of the policy gradient by building future-conditional baselines. Unlike classical Q functions, which also provide an estimate of the return for all actions but do so by averaging over all possible futures, our methods provide trajectory-specific counterfactual estimates, i.e. an estimate of the return for different actions, but keeping as many of the ex-

^{*}Equal contribution ¹DeepMind ²INRIA XPOP, CMAP, École Polytechnique, Palaiseau, France. Correspondence to: Théophile Weber <theophile@deepmind.com>, Thomas Mesnard <mesnard@deepmind.com>.

ternal factors constant between the return and its counterfactual estimate¹. Such a method would perform finer-grained credit assignment and could greatly improve data efficiency in environments with complex credit assignment structures. Our method is inspired by ideas from causality theory, but does not require learning a model of the environment.

Our main contributions are: a) introducing a family of novel policy gradient estimators that leverage hindsight information and generalizes previous approaches, b) proposing a practical instantiation of this algorithm with sufficiency conditions for unbiasedness and guarantees for lower variance, c) introducing a set of environments which further our understanding of when credit assignment is made difficult due to exogenous noise, long-term effects and task interleaving, and thus leads to poor policy learning, d) demonstrating the improved performance of our algorithm on these environments, e) formally connecting our results to notions of counterfactuals in causality theory, further linking the causal inference and reinforcement learning literatures.

2. Counterfactual Credit Assignment

2.1. Notation

We use capital letters for random variables and lower-case for the value they take. Consider a generic MDP $(\mathcal{X}, \mathcal{A}, p, r, \gamma)$. Given a current state $x \in \mathcal{X}$ and assuming an agent takes action $a \in \mathcal{A}$, the agent receives reward $r(x, a)$ and transitions to a state $y \sim p(\cdot|x, a)$. The state (resp. action, reward) of the agent at step t is denoted X_t (resp. A_t, R_t). The initial state of the agent X_0 is a fixed x_0 . The agent acts according to a policy π , i.e. action A_t is sampled from the policy $\pi_\theta(\cdot|X_t)$ where θ are the policy parameters, and aims to optimize the expected discounted return $\mathbb{E}[G] = \mathbb{E}[\sum_t \gamma^t R_t]$. The return G_t from step t is $G_t = \sum_{t' \geq t} \gamma^{t'-t} R_{t'}$. Note $G = G_0$. Finally, we define the score function $s_\theta(\pi_\theta, a, x) = \nabla_\theta \log \pi_\theta(a|x)$; the score function at time t is denoted $S_t = \nabla_\theta \log \pi_\theta(A_t|X_t)$. In the case of a partially observed environment, we assume the agent receives an observation E_t at every time step, and simply define X_t to be the set of all previous observations, actions and rewards $X_t = (O_{\leq t})$, with $O_t = (E_t, A_{t-1}, R_{t-1})$.² $\mathbb{P}(X)$ will denote the probability distribution of a random variable X .

2.2. Policy gradient algorithms

We begin by recalling two forms of policy gradient algorithms and the credit assignment assumptions they make. The first is the REINFORCE algorithm introduced by

¹From from a causality standpoint, one-step action-value functions are interventional concepts (“What would happen if”) instead of counterfactuals (“What would have happened if”).

²Previous actions and rewards are provided as part of the observation as it is generally beneficial to do so in partially observable Markov decision processes.

Williams (1992), which we will also call the single-action policy gradient estimator. The gradient of $\mathbb{E}[G]$ is given by:

$$\nabla_\theta \mathbb{E}[G] = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t S_t (G_t - V(X_t)) \right], \quad (1)$$

where $V(X_t) = \mathbb{E}[G_t|X_t]$. Let’s note here that $V(X_t)$ (resp. $Q(X_t, A_t) = \mathbb{E}[G_t|X_t, A_t]$) is the value function (resp. Q-function) for the policy π_θ but for notation simplicity the dependence on the policy will be implicit through the rest of this paper.

The appeal of this estimator lies in its simplicity and generality: to evaluate it, the only requirement is the ability to simulate trajectories, and compute both the score function and the return.

Note that subtracting the value function $V(X_t)$ from the return G_t does not bias the estimator and typically reduces variance, since the resulting estimate makes an action A_t more likely proportionally not to the return, but to which extent the return was higher than what was expected before the action was taken (Williams, 1992). Such a function will be called a *baseline* in the following. In theory, the baseline can be any function of X_t . It is however typically assumed that it does not depend on any variable ‘from the future’ (including the action about to be taken, A_t), i.e. with time index greater than t , since including variables which are (causally) affected by the action generally results in a biased estimator (Weber et al., 2019).

This estimator updates the policy through the score term; note however the learning signal only updates the policy $\pi_\theta(a|X_t)$ for the taken action $A_t = a$ (other actions are only updated through normalization of action probabilities). The policy gradient theorem from Sutton et al. (2000), which we will also call all-action policy gradient, shows it is possible to provide learning signal to all actions, given we have access to a Q-function, $Q(x, a) = \mathbb{E}[G_t|X_t = x, A_t = a]$, which we will call a *critic* in the following. The gradient of $\mathbb{E}[G]$ is given by:

$$\nabla_\theta \mathbb{E}[G] = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t \sum_{a \in \mathcal{A}} \nabla_\theta \pi(a|X_t) Q(X_t, a) \right]. \quad (2)$$

2.3. Intuitive example on hindsight reasoning and skill versus luck

Imagine a scenario in which Alice just moved to a new city, is learning to play soccer, and goes to the local soccer field to play a friendly game with a group of other kids she has never met. As the game goes on, Alice does not seem to play at her best and makes some mistakes. It turns out however her partner Megan is a strong player, and eventually scores the goal that makes the game a victory. What should Alice learn from this game?

When using the single-action policy gradient estimate, the

outcome of the game being a victory, and assuming a ± 1 reward scheme, all her actions taken during the game are made more likely; this is in spite of the fact that during this particular game she may not have played well and that the victory is actually due to her strong teammate. From an RL point of view, her actions are wrongly credited for the victory and positively reinforced as a result; effectively, Alice was lucky rather than skillful. Regular baselines do not mitigate this issue, as Alice did not a priori know the skill of Megan, resulting in an assumption that Megan was of average strength and therefore a guess that their team had a 50% chance of winning. This could be fixed by understanding that Megan’s strong play were not a consequence of Alice’s play, that her skill was a priori unknown but known in hindsight, and that it is therefore valid to retroactively include her skill level in the baseline. A hindsight baseline, conditioned on Megan’s estimated skill level, would therefore be closer to 1, driving the advantage estimate (and corresponding learning signal) close to 0.

As pointed out by Buesing et al. (2019), situations in which hindsight information is helpful in understanding a trajectory are frequent. In that work, the authors adopt a model-based framework, where hindsight information is used to ground counterfactual trajectories (i.e. trajectories under *different actions, but same randomness*). Our proposed approach follows a similar intuition, but is model-free: we attempt to *measure*—instead of *model*—information known in hindsight to compute a *future-conditional baseline*, but in a way that maintains unbiasedness. As we will see later, this corresponds to a constraint that the captured information must not have been caused by the agent.

2.4. Future-conditional (FC-PG) and Counterfactual (CCA-PG) Policy Gradient Estimators

Intuitively, our approach for assigning proper credit to action A_t relies on measuring statistics Φ_t that capture relevant information from the trajectory (e.g. including observations $O_{t'}$ at times t' greater than t). We then learn value functions or critics which are conditioned on the additional hindsight information contained in Φ_t . In general, these future-conditional values and critics would be biased for use in a policy gradient algorithm; we therefore use an importance correction term to eliminate this bias.

Theorem 1 (Future-Conditional Policy Gradient (FC-PG) estimators). *Let Φ_t be an arbitrary random variable. Assuming that $\frac{\pi(a|X_t)}{\mathbb{P}(a|X_t, \Phi_t)} < \infty$ for all a , the following is the single-action unbiased estimator of the gradient of $\mathbb{E}[G]$:*

$$\nabla_{\theta} \mathbb{E}[G] = \mathbb{E} \left[\sum_t \gamma^t S_t \left(G_t - \frac{\pi(A_t|X_t)}{\mathbb{P}(A_t|X_t, \Phi_t)} V(X_t, \Phi_t) \right) \right] \quad (3)$$

where $V(x, \phi) = \mathbb{E}[G_t | X_t = x, \Phi_t = \phi]$ is the future Φ -conditional value function.

With no requirements on Φ_t , we also have an all-action unbiased estimator:

$$\nabla_{\theta} \mathbb{E}[G] = \mathbb{E} \left[\sum_{t,a} \gamma^t \nabla_{\theta} \log \pi(a|X_t) \mathbb{P}(a|X_t, \Phi_t) Q(X_t, \Phi_t, a) \right]$$

where $Q(x, \phi, a) = \mathbb{E}[G_t | X_t = x, \Phi_t = \phi, A_t = a]$ is the future-conditional Q function (critic). Furthermore, we have $Q(X_t, a) = \mathbb{E} \left[Q(X_t, \Phi_t, a) \frac{\mathbb{P}(a|X_t, \Phi_t)}{\pi(a|X_t)} \right]$.

Intuitively, the $\frac{\pi(a|X_t)}{\mathbb{P}(a|X_t, \Phi_t)} < \infty$ condition means that knowing Φ_t should not preclude any action a which was possible for π from having potentially produced Φ_t . A counterexample is $\Phi_t = A_t$; knowing Φ_t precludes any action $a \neq A_t$ from having produced Φ_t . Typically, Φ_t will be chosen to a function of the present and future trajectory $(X_s, A_s, R_s)_{s>t}$. The estimators above are very general and generalize similar estimators (HCA) introduced by Harutyunyan et al. (2019) (see App. C for a discussion of how HCA can be rederived from FC-PG) and different choices of Φ will have varying properties. Φ may be hand-crafted using domain knowledge, or, as we will see later, learned using appropriate objectives. Note that in general an FC-PG estimator doesn’t necessarily have lower variance (a good proxy for fine-grained credit assignment) than the classical policy gradient estimator; this is due to the variance introduced by the importance weighting scheme. It would be natural to study an estimator where this effect is nullified through independence of the action and statistics Φ (resulting in a ratio of 1).

The resulting advantage estimate could thus be interpreted not just as an estimate of ‘what outcome should I expect’, but also a measure of ‘how (un)lucky did I get?’ and ‘what other outcomes might have been possible in this precise situation, had I acted differently’. It will in turn provide finer-grained credit for action A_t in a sense to be made precise below.

Corollary 1 (Counterfactual Policy Gradient (CCA-PG)). *If A_t is independent from Φ_t given X_t , the following is an unbiased single-action estimator of the gradient of $\mathbb{E}[G]$:*

$$\nabla_{\theta} \mathbb{E}[G] = \mathbb{E} \left[\sum_t \gamma^t S_t (G_t - V(X_t, \Phi_t)) \right]. \quad (4)$$

Furthermore, the hindsight advantage estimate has no higher variance than the forward one:

$$\mathbb{E} \left[(G_t - V(X_t, \Phi_t))^2 \right] \leq \mathbb{E} \left[(G_t - V(X_t))^2 \right].$$

Similarly, for the all-action estimator:

$$\nabla_{\theta} \mathbb{E}[G] = \mathbb{E} \left[\sum_t \gamma^t \sum_a \nabla_{\theta} \pi(a|X_t) Q(X_t, \Phi_t, a) \right]. \quad (5)$$

Also, we have for all a ,

$$Q(X_t, a) = \mathbb{E}[Q(X_t, \Phi_t, a) | X_t, A_t = a]$$

The benefit of the first estimator (equation 4) is clear: under the specified condition, and compared to the regular policy gradient estimator, the CCA estimator also has no bias, but the variance of its advantage estimate $G_t - V(X_t, \Phi_t)$ (the critical component behind variance of the overall estimator) is no higher.

For the all-action estimator, the benefits of CCA (equation 5) are less self-evident, since this estimator has *higher* variance than the regular all action estimator (which has variance 0). The interest here lies in bias due to learning imperfect Q functions. Both estimators require learning a Q function from data; any error in Q leads to a bias in π . Learning $Q(X_t, a)$ requires averaging over all possible trajectories initialized with state X_t and action a : in high variance situations, this will require a lot of data. In contrast, $Q(X_t, \Phi_t, a)$ predicts the average of the return G_t *conditional* on (X_t, Φ_t, a) ; if Φ_t has a high impact on G_t , the variance of that conditional return will be lower, and learning its average will in turn be far easier and data efficient.

2.5. Learning the relevant statistics: practical implementation of CCA-PG

The previous section proposes a sufficient condition on Φ for useful estimators to be derived. A question remains - how to compute such a Φ from the trajectory? While useful Φ could be handcrafted using expert knowledge, we propose to *learn* to extract Φ from the trajectory. The learning signal will be guided by two objectives: first, we will encourage Φ_t to be conditionally independent from A_t , as it is required for the estimator to be valid. Second, corollary 1 highlights that hindsight features which are predictive of the return lead to a decreased variance of the advantage estimate. To summarize, we want Φ to be predictive of the return while being independent of the action being currently credited. The corresponding hindsight conditional baseline would capture the ‘luck’ part of the outcome while the advantage estimate would capture the ‘skill’ aspect of it. We detail our agent components and losses below. See also Fig. 1 for a depiction of the resulting architecture and Appendix A for more details.

Agent components:

- **Agent network:** Our algorithm can generally be applied to arbitrary environments (e.g. POMDPs), so we assume the agent constructs an internal state X_t from past observations $(O_{t'})_{t' \leq t}$ using an arbitrary network, for instance an RNN, i.e. $X_t = \text{RNN}_{\theta_{\text{fs}}}(O_t, X_{t-1})^3$. From X_t the agent computes a policy $\pi_{\theta_{\text{fs}}}(a|X_t)$, where θ_{fs} denotes the parameters of the representation network and policy.
- **Hindsight network:** Additionally, we assume the

³Obviously, if the environment is fully observed, a feed-forward network suffices.

agent uses a hindsight network φ with parameters θ_{hs} which computes a hindsight statistic $\Phi_t = \varphi((X, A, R))$ which may depend arbitrarily on the vectors of observations, agent states and actions (in particular, it may depend on observations from timesteps $t' \geq t$).

- **Value network:** The third component is a future-conditional value network $V_{\theta_{\text{v}}}(X_t, \Phi_t)$, with parameters θ_{v} .
- **Hindsight predictor:** The last component is a probabilistic predictor h_{ω} with parameters ω that takes X_t, Φ_t as input and outputs a distribution over A_t which is used to enforce the independence condition.

Learning objectives:

- The first loss is the hindsight baseline loss $\mathcal{L}_{\text{hs}} = \sum_t (G_t - V_{\theta_{\text{v}}}(X_t, \Phi_t))^2$.
- The second loss is the independence loss, which ensures the conditional independence between A_t and Φ_t . There exists multiple ways to measure dependence between random variables; we assume a surrogate *independence maximization* (IM) loss $\mathcal{L}_{\text{IM}}(X_t)$ which is non-negative and zero if and only if A_t and Φ_t are conditionally independent given X_t . An example is to choose the Kullback-Leibler divergence between the distributions $\mathbb{P}(A_t|X_t)$ and $\mathbb{P}(A_t|X_t, \Phi_t)$. In this case, the KL can be estimated by $\sum_a \mathbb{P}(a|X_t) (\log \mathbb{P}(a|X_t) - \log \mathbb{P}(a|X_t, \Phi_t))$; $\log \mathbb{P}(a|X_t)$ is simply the policy $\pi(a|X_t)$; the posterior $\mathbb{P}(a|X_t, \Phi_t)$ is generally not known exactly, but we estimate it with the probabilistic predictor $h_{\omega}(A_t|X_t, \Phi_t)$, which we train with the next loss.
- The third loss is the hindsight predictor loss, which we train by minimizing the supervised learning loss $\mathcal{L}_{\text{sup}} = -\sum_t \mathbb{E}[\log h_{\omega}(A_t|X_t, \Phi_t)]$ on samples (X_t, A_t, Φ_t) from the trajectory (note that this is a proper scoring rule, i.e. the optimal solution to the loss is the true probability $\mathbb{P}(a|X_t, \Phi_t)$).
- The last loss is the policy gradients surrogate objective, implemented as $\mathcal{L}_{\text{PG}} = \sum_t \log \pi_{\theta}(A_t|X_t) (G_t - V(X_t, \Phi_t))$, where the bar notation indicates that the quantity is treated as a constant from the point of view of gradient computation, as is standard.

The overall loss is therefore $\mathcal{L} = \mathcal{L}_{\text{PG}} + \lambda_{\text{hs}} \mathcal{L}_{\text{hs}} + \lambda_{\text{sup}} \mathcal{L}_{\text{sup}} + \lambda_{\text{IM}} \mathcal{L}_{\text{IM}}$. We again want to highlight the very special role played by ω here: only \mathcal{L}_{sup} is optimized with respect to ω (the parameters of the probabilistic predictor), while all the other losses are optimized treating ω as a constant.

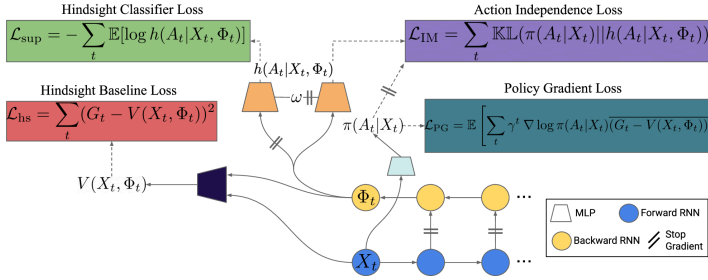


Figure 1: Counterfactual Credit Assignment in a nutshell: (1) The backward RNN which in this example computes the hindsight features is shaped by the hindsight baseline loss. This ensures that it is predictive of the return. (2) However, to have an unbiased baseline, this hindsight feature Φ_t needs to be independent from the action A_t . To that end, we first train a hindsight predictor that tries to predict what action has been taken a time t from X_t and Φ_t . (3) Then the action independence loss helps removing any information about A_t from the hindsight feature Φ_t (This only enforces that the output of the backward RNN Φ_t is independent from the action A_t . However, this could potentially translate in Φ_t being independent from further actions). This loss only impacts the backward RNN and no gradient is being applied to the hindsight predictor MLP. (4) Finally, the policy gradient loss helps improving the policy while no gradient is being sent to the hindsight baseline (i.e. as expressed by the bar notation).

3. Connections to causality

In this section we provide a formal connection between the CCA-PG estimator and counterfactuals in causality theory (this connection is investigated in greater depth in appendices F and G).

To this end, we assume that the MDP $(\mathcal{X}, \mathcal{A}, p, r, \gamma)$ in question is generated by an underlying structural causal models (SCM) analogous to (Buesing et al., 2019; Zhang, 2020). In this setting the trajectory $(X_s, A_s, R_s)_{s \geq t}$ and return G_t resulting from the agent-environment interaction is represented as the output of a deterministic function f taking as input the current state X_t , the action A_t , and a set of exogenous random variables \mathcal{E} which do not have any causal ancestors (in the graph). The latter represent the randomness required for sampling all future actions, transitions, and rewards. Such a "reparametrization" of trajectories and return is always possible, i.e. there is always an SCM (possibly non unique) that induces the same joint distribution \mathbb{P} as the original MDP. Intuitively, \mathcal{E} represent all factors external to A_t which affect the outcome⁴.

SCMs allow to formally define the notion of counterfactual. Given an observed trajectory $\tau = (X_s, A_s, R_s)_{s \geq t}$, we define the counterfactual trajectory τ' for an alternative action $A'_t = a'_t$ as the output of the following procedure:

- **Abduction:** infer the exogenous noise variables ϵ under the factual observation: $\epsilon \sim \mathbb{P}(\mathcal{E}|\tau)$.
- **Intervention:** Fix the value of A'_t to a'_t (mutilating incoming causal arrows).
- **Prediction:** Evaluate the counterfactual outcome τ' conditional on the fixed values \mathcal{E} and $A_t = a'_t$ yielding $\tau' = f(x_t, a'_t, \epsilon)$

The counterfactual distribution will be denoted $P(\tau'|\text{observe}(\tau), \text{do}(A'_t = a'_t))$. Note that it typically requires knowledge of the model (SCM) to be computed; samples from the models which do not expose

⁴Note that from this point of view, actions at future time-step are effectively 'chance' from the point of view of computing credit for action A_t

the exogenous variables \mathcal{E} are not typically not sufficient to identify the SCM, as several SCMs may correspond to the same distribution. However, under the CCA assumptions and an additional faithfulness assumption, we can show that the counterfactual return is indeed identifiable and is equal to the future conditional state-action value function:

Theorem 2. *Assume the causal model is faithful (i.e. that conditional independence assumptions are reflected in the graph structure and not only in the parameters). If Φ_t is conditionally independent from A_t given X_t , then the counterfactual distribution, having observed only Φ_t , is identifiable from samples of (X_t, Φ_t, A_t) , and we have*

$$\mathbb{E}[G(\tau')|\tau' \sim P(\tau'|X_t = x, \text{observe}(\Phi_t = \phi), \text{do}(A'_t = a))] = Q(X_t = x, A_t = a, \Phi_t = \phi) \quad (6)$$

4. Numerical experiments

Given its guarantees on lower variance and unbiasedness, we run all our experiments on the single action version of CCA-PG and leave the all-action version for future work. We first investigate a bandit with feedback task, then a task that requires short and long-term credit assignment (i.e. Key-to-Door), and finally an interleaved multi-task setup where each episode is composed of randomly sampled and interleaved tasks. All results for Key-to-Door and interleaved multi-task are reported as median performances over 10 seeds with quartiles represented by a shaded area.

4.1. Bandit with Feedback

We first demonstrate the benefits of hindsight value functions in a toy problem designed to highlight these. We consider a contextual bandit problem with feedback. At each time step, the agent receives a context $-N \leq C \leq N$ (where N is an environment parameter), and based on the context, chooses an action $-N \leq A \leq N$. The agent receives a reward $R = -(C - A)^2 + \epsilon_r$, where the exogenous noise ϵ_r is sampled from $\mathcal{N}(0, \sigma_r)$, as well as a feedback vector F which is a function of C, A and ϵ_r . More details about this problem as well as variants are presented in Appendix B.1.

For this problem, the optimal policy is to choose $A = C$, resulting in average reward of 0. However, the reward signal R is corrupted by the exogenous noise ϵ_r , uncorrelated to the action. The higher the standard deviation, the more difficult proper credit assignment becomes, as high rewards are more likely due to a high value of ϵ_r than an appropriate choice of action. On the other hand, the feedback F contains information about C , A and ϵ_r . If the agent can extract information Φ from F in order to capture information about ϵ_r and use it to compute a hindsight value function, the effect of the perturbation ϵ_r may be removed from the advantage estimate, resulting in a significantly lower variance estimator. However, if the agent blindly uses F to compute the hindsight value information, information about the action will ‘leak’ into the hindsight value, leading to an advantage estimate of 0 and no learning.

We investigate the proposed algorithm with $N = 10$. As can be seen on Fig. 2, increasing the variance of the exogenous noise leads to dramatic decrease of performance for the vanilla PG estimator without the hindsight baseline; in contrast, the CCA-PG estimator is generally unaffected by the exogenous noise. For very low level of exogenous noise however, CCA-PG suffers from a decrease in performance. This is due to the agent computing a hindsight statistic Φ which is not perfectly independent from A , leading to bias in the policy gradient update. To demonstrate this effect, and evaluate the importance of the independence constraint on performance, we run an ablation where we test lower values of the weight λ_{IM} of the independence maximization loss (leading to a larger mutual information between Φ and A) and indeed observed that the performance is dramatically degraded, as seen in Fig. 2.

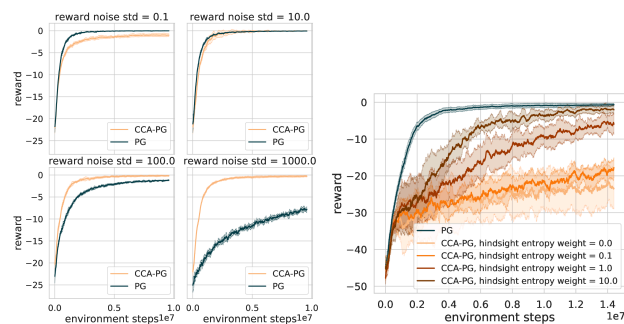


Figure 2: Top: Comparison of CCA-PG and PG in contextual bandits with feedback, for various levels of reward noise σ_r . Results are averaged over 6 independent runs with standard deviation represented by a shaded area. **Bottom:** Performance of CCA-PG on the bandit task, for different values of λ_{IM} . Properly enforcing the independence constraint prevents the degradation of performance.

4.2. Key-to-Door environments

Task Description. We investigate new versions of the Key-To-Door family of environments, initially proposed by Hung et al. (2019), as a testbed of tasks where credit

assignment is hard and is necessary for success. In this partially observable grid-world environment (cf. Fig. 7 in the appendix), the agent has to pick up a key in the first room, for which it has *no immediate reward*. In the second room, the agent can pick up 10 apples, that each give immediate rewards. In the final room, the agent may open a door (only if it is carrying a key), and receive a small reward for doing so. In this task, a single action (i.e picking up the key) has a direct impact on the reward it receives in the final room, however this signal is hard to detect as the episode return is largely driven by its performance in the second room (i.e picking up apples).

We now consider two instances of the Key-To-Door family that illustrate the difficulty of credit assignment in the presence of extrinsic variance. In the Low-Variance-Key-To-Door environment, each apple is worth a reward of 1 and opening the final door also gets a reward of 1. Thus, an agent that solves the apple phase perfectly sees very little variance in its episode return and the learning signal for picking up the key and opening the door is relatively strong.

High-Variance-Key-To-Door keeps the overall structure of the Key-To-Door task. The door keeps giving a deterministic reward of 1 when the key was grabbed but now the reward for each apple is randomly sampled to be either 1 or 10, and fixed within the episode. In this setting, even an agent that is skilled at picking up apples sees a large variance in episode returns, and thus the learning signal for picking up the key and opening the door is comparatively weaker. Appendix B.2.1 has some additional discussion illustrating the difficulty of learning in such a setting.

Results We test CCA-PG on these environments, and compare it against Actor-Critic (Williams, 1992), as well as State-conditional HCA and Return-conditional HCA (Harutyunyan et al., 2019) as baselines. An analysis of the relation between HCA and CCA is described in Appendix C. We test using both a backward-LSTM (referred to as CCA-PG RNN) or an attention model (referred to as CCA-PG Attn) for the hindsight function. Details for experimental setup are provided in Appendix B.2.2.

We evaluate agents both on their ability to maximize total reward, as well as to solve the specific credit assignment problem of picking up the key and opening the door. Fig. 3 compares CCA-PG with the baselines on the High-Variance-Key-To-Door task. Both CCA-PG architectures outperform the baselines in terms of total reward, as well as probability of picking up the key and opening the door.

This experiment highlights the capacity of CCA-PG to learn and incorporate trajectory-specific external factors into its baseline, resulting in lower variance estimators. Despite being a difficult task for credit assignment, CCA-PG is capable of solving it quickly and consistently. On the other

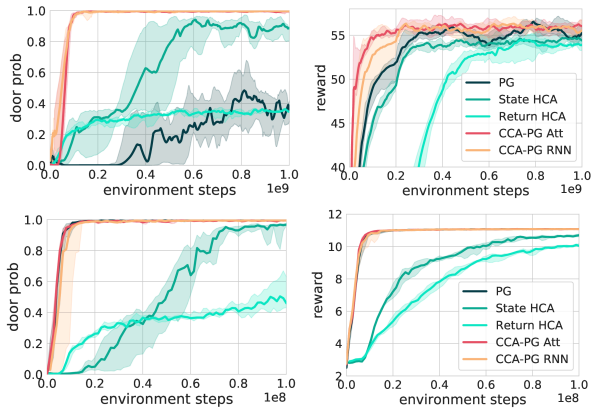


Figure 3: Probability of opening the door and total reward obtained on the **High-Variance-Key-To-Door** task (top two) and the **Low-Variance-Key-To-Door** task (bottom two).

hand, vanilla actor-critic is greatly impacted by this external variance, and needs around $3 \cdot 10^9$ environment steps to have an 80% probability of opening the door. CCA-PG also outperforms State- and Return- Conditional HCA, which do use hindsight information but in a more limited way than CCA-PG.

On the Low-Variance-Key-To-Door task (Fig. 3), due to the lack of extrinsic variance, standard actor-critic is able to perfectly solve the environment. However, it is interesting to note that CCA-PG still matches this perfect performance. On the other hand, the other hindsight methods struggle with both door-opening and apple-gathering. This might be explained by the fact that both these techniques do not guarantee lower variance, and rely strongly on their learned hindsight classifiers for their policy gradient estimators, which can be harmful when these quantities are not perfectly learned. See Appendix B.2.3 for additional experiments and ablations on these environments.

These experiments demonstrate that CCA-PG is capable of efficiently leveraging hindsight information to mitigate the challenge of external variance and learn strong policies that outperform baselines. At the same time, it suffers no drop in performance when used in cases where external variance is minimal.

4.3. Task Interleaving

Motivation. In the real world, human activity can be seen as solving a large number of loosely related problems. At an abstract level, one could see this lifelong learning process as solving problems not in a sequential, but an interleaved fashion instead. These problems are not solved sequentially, as one may temporarily engage with a problem and only continue engaging with it or receive feedback from its earlier actions significantly later. The structure of this interleaving will also typically vary over time.

To better understand the effects of interleaving on agent learning, we introduce a new class of environments capturing the structural properties mentioned above. In contrast to most work on multi-task learning, we do not assume a clear delineation between subtasks, nor focus on skill retention. The agent will encounter multiple tasks in a single episode in an interleaved fashion (switching between tasks will occur before a task gets completed), and will have to detect the implicitly boundaries between them.

Task Description. This task consists of pairs of query-answer rooms with different visual contexts that each indicates a different subtask. In the query room, the agent gets to pick between two colored boxes (out of 10 possible colors). Later, in the answer room, the agents gets to observe which of the two boxes was rewarding in the first room, and receives a reward if it picked the correct box (there is always exactly one rewarding color in the query room). The mapping of colors to whether it is rewarding or not is specific to each subtask and fixed across training. Each subtask would be relatively easy to solve if encountered in an isolated fashion. However, each episode is composed of *randomly sampled subtasks* and color pairs within those subtasks. Furthermore, query rooms and answer rooms of the sampled subtasks are presented in a random (interleaved) order which differs from one episode to another. Each episode are 140 steps long and it takes at least 9 steps for the agent to reach one colored square from its initial position. A visual example of what an episode looks like can be seen in Fig. 4.

There are six tasks, each classified as ‘easy’ or ‘hard’; easy tasks have high reward signals (i.e. easier for agents to pick up on), while hard tasks have low rewards. In the 2 tasks setup (resp. 4 tasks and 6 tasks), there is one (resp. two and two) ‘easy’ and one (resp. two and four) ‘hard’ task. More details about the experimental setup can be found in B.3.

In addition to the total reward, we record the probability of picking up the correct square for the easy and hard tasks separately. Performance in the hard tasks will indicate ability to do fine-grained credit assignment.

Results. While CCA-PG is able to perfectly solve both the ‘easy’ and ‘hard’ tasks in the three setups in less than $5 \cdot 10^8$ environment steps (Fig. 5), actor-critic is only capable to solve the ‘easy’ tasks for which the associated rewards are large. Even after $2 \cdot 10^9$ environment steps, actor-critic is still greatly impacted by the variance and remains incapable of solving ‘hard’ tasks in any of the three settings. CCA-PG also outperforms actor-critic in terms of the total reward obtained in each setting. State-conditional and Return-conditional HCA were also evaluated on this task but results are not reported as almost no learning was taking place on the ‘hard’ tasks. More results along with an ablation study can be found in Appendix B.3.

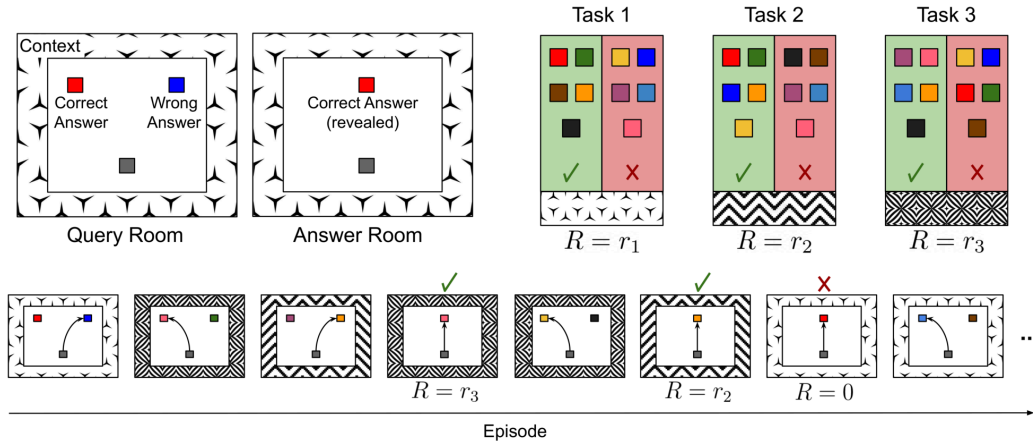


Figure 4: Task Interleaving Description. **Top left:** Delayed feedback contextual bandit problem. Given a context shown as a surrounding visual pattern, the agent has to decide to pick up one of the two colored squares where only one will be rewarding. The agent is later teleported to the second room where it is provided with the reward associated with its previous choice and a visual cue about which colored square it should have picked up. **Top right:** Different tasks with each a different color mapping, visual context and associated reward. **Bottom:** Example of a generated episode, composed of randomly sampled tasks and color pairs.

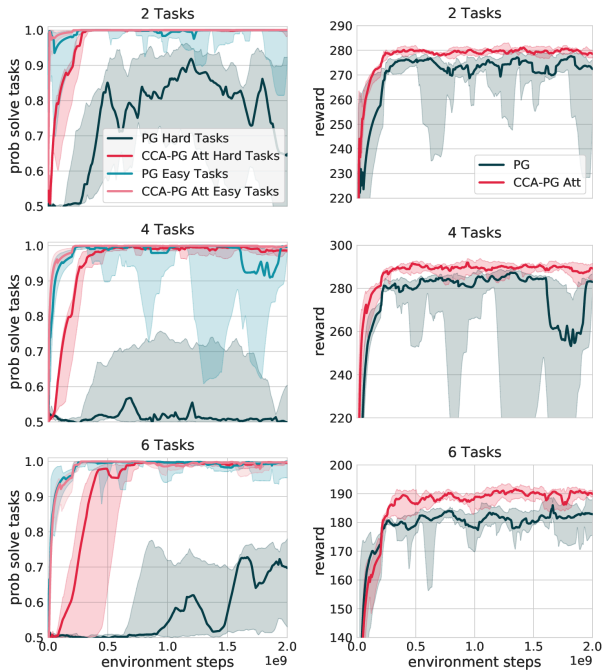


Figure 5: Probability of solving ‘easy’ and ‘hard’ tasks (left) and total reward (right) obtained for the **Multi Task Interleaving**. Left plots: Median over 10 seeds after doing a mean over the performances in ‘easy’ or ‘hard’ tasks.

Through efficient use of hindsight, CCA-PG is able to take into account trajectory-specific factors such as the kinds of rooms encountered in the episode and their associated rewards.

In the case of the Multi-Task Interleaving environment, an informative hindsight function would capture the reward for

different contexts and exposes as Φ_t all rewards obtained in the episode except those associated with the current context. This experiment again highlights the capacity of CCA-PG to solve hard credit assignment problems in a context where the return is affected by multiple distractors, while PG remains highly sensitive to them.

5. Related work

This paper builds on work from Buesing et al. (2019) which shows how causal models and real data can be combined to generate counterfactual trajectories and perform off-policy evaluation for RL. Their results however require an explicit model of the environment. In contrast, our work proposes a model-free approach, and focuses on policy improvement. Oberst & Sontag (2019) also investigate counterfactuals in reinforcement learning, point out the issue of non-identifiability of the correct SCM, and suggest a sufficient condition for identifiability; we discuss this issue in appendix G. Closely related to our work is Hindsight Credit Assignment, a concurrent approach from Harutyunyan et al. (2019). In this paper, the authors also investigate value functions and critics that depend on future information. However, the information the estimators depend on is fixed (future state or return) instead of being an arbitrary functions of the trajectory. Our FC estimators generalizes both the HCA and CCA estimators while CCA further characterizes which statistics of the future provide a useful estimator. Relations between HCA, CCA and FC are discussed in appendix C. The HCA approach is further extended by Young (2019), and Zhang et al. (2019) who minimize a surrogate for the variance of the estimator, but that surrogate cannot be guaranteed to actually lower the variance. Similarly to state-HCA, it treats each reward separately instead of taking

a trajectory-centric view as CCA. Guez et al. (2019) also investigate future-conditional value functions. Similarly to us, they learn statistics of the future Φ from which returns can be accurately predicted, and show that doing so leads to learning better representations (but use regular policy gradient estimators otherwise). Instead of enforcing an information-theoretic constraint, they bottleneck information through the size of the encoding Φ . In domain adaptation (Ganin et al., 2016; Tzeng et al., 2017), robustness to the training domain can be achieved by constraining the agent representation not to be able to discriminate between source and target domains, a mechanism similar to the one constraining hindsight features not being able to discriminate the agent’s actions. Also closely related to our paper, Bica et al. (2020) also leverages a similar mechanism to compute counterfactuals, for a different purpose than ours (computing treatment effects vs. policy improvement operators).

Both Andrychowicz et al. (2017) and Rauber et al. (2017) leverage the idea of using hindsight information to learn goal-conditioned policies. Hung et al. (2019) leverages attention-based systems and episode memory to perform long term credit assignment; however, their estimator will in general be biased. Ferret et al. (2019) looks at the question of transfer learning in RL and leverages transformers to derive a heuristic to perform reward shaping. Arjona-Medina et al. (2019) also addresses the problem of long-term credit assignment by redistributing delayed rewards earlier in the episode but their approach still fundamentally uses time as a proxy for credit.

Previous research also leverages the fact that baselines can include information unknown to the agent at time t (but potentially revealed in hindsight) but not affected by action A_t : for instance, when using independent multi-dimensional actions, the baseline for one dimension of the action vector can include the actions in other dimensions (Wu et al., 2018); or when the dynamic of the environment is partially driven by an exogenous and stochastic factor, independent of the agent’s actions, which can be included in the baseline (Mao et al., 2018). Similarly, in multi-agent environments, actions of other agents at the same time step (Foerster et al., 2018) can be used; and so can the full state of the simulator when learning control from pixels (Andrychowicz et al., 2020), or the use of opponent observations in Starcraft II (Vinyals et al., 2019). Note however that all of these require privileged information, both in the form of feeding information to the baseline inaccessible to the agent, and in knowing that this information is independent from the agent’s action A_t and therefore won’t bias the baseline. Our approach seeks to replicate a similar effect, but in a more general fashion and from an agent-centric point of view, where the agent *learns itself* which information from the future can be used to improve its baseline at time t .

6. Conclusion

In this paper we have considered the problem of credit assignment in RL. Building on insights from causality theory and structural causal models, we have investigated the concept of future-conditional value functions. Contrary to common practice these allow baselines and critics to condition on future events thus separating the influence of an agent’s actions on future rewards from the effects of other random events thus reducing the variance of policy gradient estimators. A key difficulty lies in the fact that unbiasedness relies on accurate estimation and minimization of mutual information. Learning inaccurate hindsight classifiers will result in miscalibrated estimation of luck, leading to bias in learning. Future research will investigate how to scale these algorithms to more complex environments, and the benefits of the more general FC-PG and all-actions estimators.

References

- Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Abbeel, O. P., and Zaremba, W. Hindsight experience replay. In *Advances in neural information processing systems*, pp. 5048–5058, 2017.
- Andrychowicz, O. M., Baker, B., Chociej, M., Jozefowicz, R., McGrew, B., Pachocki, J., Petron, A., Plappert, M., Powell, G., Ray, A., et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.
- Arjona-Medina, J. A., Gillhofer, M., Widrich, M., Unterthiner, T., Brandstetter, J., and Hochreiter, S. Rudder: Return decomposition for delayed rewards. In *Advances in Neural Information Processing Systems*, pp. 13544–13555, 2019.
- Bica, I., Alaa, A. M., Jordon, J., and van der Schaar, M. Estimating counterfactual treatment outcomes over time through adversarially balanced representations. *arXiv preprint arXiv:2002.04083*, 2020.
- Buesing, L., Weber, T., Zwols, Y., Racaniere, S., Guez, A., Lespiau, J.-B., and Heess, N. Woulda, coulda, shoulda: Counterfactually-guided policy search. *2019 International Conference for Learning Representations (ICLR)*, 2019.
- Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. Gated feedback recurrent neural networks. In *International conference on machine learning*, pp. 2067–2075, 2015.
- Ferret, J., Marinier, R., Geist, M., and Pietquin, O. Credit assignment as a proxy for transfer in reinforcement learning. *arXiv preprint arXiv:1907.08027*, 2019.
- Foerster, J. N., Farquhar, G., Afouras, T., Nardelli, N., and Whiteson, S. Counterfactual multi-agent policy gradients. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016.
- Glasserman, P. and Yao, D. D. Some guidelines and guarantees for common random numbers. *Management Science*, 38(6):884–908, 1992.
- Goyal, A., Lamb, A., Hoffmann, J., Sodhani, S., Levine, S., Bengio, Y., and Schölkopf, B. Recurrent independent mechanisms. *arXiv preprint arXiv:1909.10893*, 2019.
- Guez, A., Viola, F., Weber, T., Buesing, L., Kapturowski, S., Precup, D., Silver, D., and Heess, N. Value-driven hindsight modelling. <https://openreview.net/forum?id=rJxBa1HFvS>, 2019.
- Ha, D. and Schmidhuber, J. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- Hamrick, J. B. Analogues of mental simulation and imagination in deep learning. *Current Opinion in Behavioral Sciences*, 29:8–16, 2019.
- Harutyunyan, A., Dabney, W., Mesnard, T., Azar, M. G., Piot, B., Heess, N., van Hasselt, H. P., Wayne, G., Singh, S., Precup, D., et al. Hindsight credit assignment. In *Advances in Neural Information Processing Systems*, pp. 12467–12476, 2019.
- Heess, N., Wayne, G., Silver, D., Lillicrap, T., Erez, T., and Tassa, Y. Learning continuous control policies by stochastic value gradients. In *Advances in Neural Information Processing Systems*, pp. 2944–2952, 2015.
- Hinton, G., Srivastava, N., and Swersky, K. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Cited on*, 14(8), 2012.
- Hung, C.-C., Lillicrap, T., Abramson, J., Wu, Y., Mirza, M., Carnevale, F., Ahuja, A., and Wayne, G. Optimizing agent behavior over long time scales by transporting value. *Nature communications*, 10(1):1–12, 2019.
- Kaiser, L., Babaeizadeh, M., Milos, P., Osinski, B., Campbell, R. H., Czechowski, K., Erhan, D., Finn, C., Koza-kowski, P., Levine, S., et al. Model-based reinforcement learning for atari. *arXiv preprint arXiv:1903.00374*, 2019.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Mao, H., Venkatakrisnan, S. B., Schwarzkopf, M., and Alizadeh, M. Variance reduction for reinforcement learning in input-driven environments. In *International Conference on Learning Representations*, 2018.
- Minsky, M. Steps toward artificial intelligence. *Proceedings of the IRE*, 49(1):8–30, 1961.
- Oberst, M. and Sontag, D. Counterfactual off-policy evaluation with gumbel-max structural causal models. *arXiv preprint arXiv:1905.05824*, 2019.
- Parisotto, E., Song, H. F., Rae, J. W., Pascanu, R., Gulcehre, C., Jayakumar, S. M., Jaderberg, M., Kaufman, R. L., Clark, A., Noury, S., et al. Stabilizing transformers for reinforcement learning. *arXiv preprint arXiv:1910.06764*, 2019.

- Pearl, J. *Causality*. Cambridge university press, 2009a.
- Pearl, J. *Causality: Models, reasoning, and inference*. 2009b.
- Rauber, P., Ummadisingu, A., Mutz, F., and Schmidhuber, J. Hindsight policy gradients. *arXiv preprint arXiv:1711.06006*, 2017.
- Rezende, D. J. and Viola, F. Taming VAEs. *arXiv preprint arXiv:1810.00597*, 2018.
- Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., et al. Mastering atari, go, chess and shogi by planning with a learned model. *arXiv preprint arXiv:1911.08265*, 2019.
- Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pp. 1057–1063, 2000.
- Tzeng, E., Hoffman, J., Saenko, K., and Darrell, T. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7167–7176, 2017.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. Grandmaster level in starcraft II using multi-agent reinforcement learning. *Nature*, 575 (7782):350–354, 2019.
- Weber, T., Heess, N., Buesing, L., and Silver, D. Credit assignment techniques in stochastic computation graphs. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2650–2660, 2019.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- Wu, C., Rajeswaran, A., Duan, Y., Kumar, V., Bayen, A. M., Kakade, S., Mordatch, I., and Abbeel, P. Variance reduction for policy gradient with action-dependent factorized baselines. *2018 International Conference for Learning Representations (ICLR)*, 2018.
- Young, K. Variance reduced advantage estimation with δ -hindsight credit assignment. *arXiv preprint arXiv:1911.08362*, 2019.
- Zhang, J. Designing optimal dynamic treatment regimes: A causal reinforcement learning approach. In *International Conference on Machine Learning*, pp. 11012–11022. PMLR, 2020.
- Zhang, P., Zhao, L., Liu, G., Bian, J., Huang, M., Qin, T., and Tie-Yan, L. Independence-aware advantage estimation. <https://openreview.net/forum?id=B1eP504YDr>, 2019.