# Towards a practical measure of interference for reinforcement learning

**Vincent Liu[1], Adam White[12], Hengshuai Yao[3], Martha White[1]**
[1]University of Alberta
[2]DeepMind
[3]Huawei Technologies
{vliu1, amw8, whitem}@ualberta.ca, hengshuai.yao@huawei.com

## Abstract

Catastrophic interference is common in many network-based learning systems, and many proposals exist for mitigating it. But, before we overcome interference we must understand it better. In this work, we provide a definition of interference for control in reinforcement learning. We systematically evaluate our new measures, by assessing correlation with several measures of learning performance, including stability, sample efficiency, and online and offline control performance across a variety of learning architectures. Our new interference measure allows us to ask novel scientific questions about commonly used deep learning architectures. In particular we show that target network frequency is a dominating factor for interference, and that updates on the last layer result in significantly higher interference than updates internal to the network. This new measure can be expensive to compute; we conclude with motivation for an efficient proxy measure and empirically demonstrate it is correlated with our definition of interference.

## 1 Introduction

Generalization is a key property of reinforcement learning algorithms with function approximation. It is important for an agent to generalize from previous encountered samples to a larger subset of samples which have not been seen. Generalization has been extensively studied in supervised learning, where we normally assume that we can sample iid inputs from a fixed input distribution and the targets are sampled from a fixed conditional distribution.

The assumption of iid inputs, however, does not hold in general. When learning on a correlated stream of data, as in RL, the learner might fit the learned function to recent data and potentially overwrite or forget previously learned information. This issue is called *catastrophic interference*. Interference occurs even in the iid prediction setting: an update on some set of states is said to interfere with predictions in another state when that update decreases accuracy for that state. This interference is catastrophic if it causes significant forgetting, which is typically only observed with temporally correlated data, such as in RL [4, 14, 21] or in the sequential multi-task learning setting [19, 28]. The conventional wisdom is that catastrophic interference is particularly problematic in the control setting in RL, even single-task RL, because (a) when an agent explores, it receives a sequence of observations, which are likely to be temporally correlated; (b) the agent is changing its policy while learning, making the sequence of observations non-stationary; and (c) the agent uses its own estimates as targets (as in temporal difference learning), which makes the target outputs non-stationary.

It is as yet difficult to verify this conventional wisdom, as we do not have effective means to measure interference. It is commonly held that replay, target networks and the choice of representation [21] all mitigate interference, and so improve performance. But, without a clear definition and way to measure interference in RL, it is hard to test these hypotheses. There has been work quantifying

interference for supervised learning [6, 13, 17, 28], with some empirical work even correlating catastrophic forgetting and properties of task sequences in supervised learning [25]. In prediction, however, the definition of interference is relatively straightforward: interference corresponds to decreases in prediction accuracy, which can be measured using a stored test set. This definition, unfortunately, does not extend to the control setting: if we use value function accuracy, then we have a changing performance measure as the policy changes. Several papers have investigated generalization and transfer in RL [7, 10, 26, 27], demonstrating that learning on new environments results in drops in performance on previously learned environments [7], or re-initialization can help a plateaued agent make further progress [11]. These works, however, do not directly measure levels of interference, and instead focus on test performance on new environments or new segments of environments.

In this paper, we propose a definition of interference for control in RL using an existing performance measure, called the Optimality Residual (OR). The interference is defined as the change in OR, with two statistics to reflect the presence of catastrophic interference. We evaluate of our interference measures by computing the correlation with several performance metrics, including sample efficiency and stability. We also use these measures to investigate the role of common deep RL techniques, including target networks, experience replay buffer size, mini-batch size, network size, and interference in different layers. It is difficult—or in some cases impossible—to estimate this exact interference measure. We provide an approximation, by deriving an upper bound on the OR, and demonstrate empirically that the approximation is strongly correlated with the exact interference.

## 2  Background

In reinforcement learning (RL), an agent interacts with its environment, receiving observations and selecting actions to maximize a reward signal. We assume the environment can be formalized as a Markov decision process (MDP). An MDP is a tuple $(\mathcal{S}, \mathcal{A}, \mathrm{Pr}, R, \gamma)$ where $\mathcal{S}$ is a set of states, $\mathcal{A}$ is an set of actions, $\mathrm{Pr} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to [0, 1]$ is the transition probability, $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ is the reward function, and $\gamma \in [0, 1]$ a discount factor. The goal of the agent is to find a policy $\pi : \mathcal{S} \times \mathcal{A} \to [0, 1]$ to maximize the expected discounted sum of rewards.

Given a fixed policy $\pi$, the action-value function $Q^\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is defined as $Q^\pi(s, a) := \mathbb{E}[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a]$, where $R_{t+1}$ denotes the reward at time $t + 1$, i.e. $R_{t+1} = R(S_t, A_t, S_{t+1})$, $S_{t+1} \sim \mathrm{Pr}(\cdot | S_t, A_t)$, and actions are taken according to policy $\pi$: $A_t \sim \pi(\cdot | S_t)$. The optimal value function $Q^*$ is defined as $Q^*(s, a) := \sup_\pi Q(s, a)$, with $\pi^*$ the policy that is greedy w.r.t. $Q^*$. The optimal value function can be obtained using the Bellman optimality operator for action values $\mathcal{T} : \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|} \to \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$:

$$(\mathcal{T}Q)(s, a) := \sum_{s' \in \mathcal{S}} \mathrm{Pr}(s'|s, a) \left[ R(s, a, s') + \gamma \max_{a' \in \mathcal{A}} Q(s', a') \right]$$

$Q^*$ is the unique solution of the Bellman equation $\mathcal{T}Q = Q$. Q-learning is built on this operator, iteratively updating to find the fixed point of the Bellman optimality operator.

We can use neural networks to learn an approximation to the optimal action-value. For $Q_{\boldsymbol{\theta}}$ the approximation, with parameters $\boldsymbol{\theta}$, the online update for non-linear semi-gradient Q learning is

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \alpha \delta_t \nabla_{\boldsymbol{\theta}_t} Q_{\boldsymbol{\theta}_t}(S_t, A_t) \qquad \text{where } \delta_t := R_{t+1} + \gamma \left[ \max_{a' \in \mathcal{A}} Q_{\boldsymbol{\theta}_t}(S_{t+1}, a') - Q_{\boldsymbol{\theta}_t}(S_t, A_t) \right].$$

This update with NNs typically leads to unstable performance, so is often augmented with experience replay [20] and target networks, introduced in DQN [23]. Replay consists of storing transitions in a buffer $D$, and performing mini-batch updates sampled from this buffer, per step. Target networks use an older set of parameters $\bar{\boldsymbol{\theta}}$ for $\max_{a' \in \mathcal{A}} Q_{\bar{\boldsymbol{\theta}}_t}(s', a')$, to make the update target more stationary.

## 3  A Simple Example Relating Interference and Control Performance

Before discussing our definition and measure of interference, it is useful to use a controlled setting to illustrate how algorithmic choices impact interference and performance. For example, we expect agents with poor representations to suffer from more interference. If we have a very good hand-designed, sparse representation—such as tile-coding—we expect much less interference than a neural

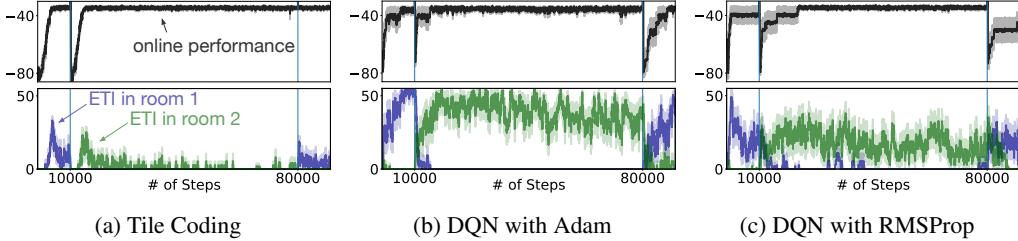|                | (a) Tile Coding | (b) DQN with Adam | (c) DQN with RMSProp |

Figure 1: The Two-Room example. We plot the learning curve of Q-learning with different architecture choices. The three stages are indicated by the two vertical lines. ETI is a measure of interference, which is defined in a latter section. The curves are averaged over 10 runs with one standard error.

network representation that generalizes aggressively. We use three such agents for demonstration: Q-learning with tile-coding, DQN with the Adam optimizer and DQN with the RMSprop optimizer.

The controlled environment, called Two-Rooms, consists of two open rooms with different start and goal states. The trick is that in the first room the agent should navigate up and to the right, and in the second room down and to the left. The input state contains the xy position of the agent, and which room the agent is in. The tile coding agent represents each room independently, whereas DQN is free to generalize across rooms. The agent begins life in one room and trains just long enough (10k steps) to learn a near optimal policy. Then the agent is placed in the second room and trained much longer (70k steps) than required—to the point that over specialization is possible. Finally, the agent is placed back to learn in the first room, to evaluate the impact of extended training in the second room.

In Figure 1, we show online learning curves and the corresponding interference (defined in Section 4.3) in each room separately. Generally, we can see that when the agent is learning, there is interference; the key issue is whether learning in one room interfere with the other. The tile-coding representation—with no features shared between rooms—has no interference in one room, while training in the other. The performance of the DQN agents drops when transfering from room 2 back to room 1. The interference is catastrophic: the agent using RMSProp does not recover the optimal policy, and the agent using Adam learns more slowly than starting from scratch.

# 4 Measuring Interference in RL

In this section, we define interference for control in RL. We start by discussing the definition of interference in RL for the prediction setting, where we learn $Q^\pi$; we do this for clarity and to provide a contrast to the control setting. We highlight that to define whether an update causes interference requires an answer to the question: interference according to what objective? We propose a natural choice for control: the distance to the optimal action-value function. We discuss two ways to summarize interference over time, to gauge whether an agent has high or low interference.

## 4.1 Interference in Prediction

In the prediction setting, we estimate $Q^\pi$ for a fixed $\pi$. A typical measure of prediction error is the mean-squared value error (MSVE), with state-action weighting $d : \mathcal{S} \times \mathcal{A} \to [0, \infty)$

$$\text{MSVE}(\boldsymbol{\theta}) := \|Q^\pi - Q_{\boldsymbol{\theta}}\|_d^2 = \sum_{s \in \mathcal{S}, a \in \mathcal{A}} d(s, a)(Q^\pi(s, a) - Q_{\boldsymbol{\theta}}(s, a))^2$$

To quantify *expected interference*, we can look at the difference in MSVE before and after an update: $\text{MSVE}(\boldsymbol{\theta}_{t+1}) - \text{MSVE}(\boldsymbol{\theta}_t)$. If this value is positive, the update generally degraded performance and there was more interference on average than positive generalization. If this value is negative, the update generally improved performance and there was more positive generalization than interference.

There are existing interference measures based on gradient similarity that could be used for the prediction setting. To see why, assume we can directly minimize the MSVE and so have loss $L(\boldsymbol{\theta}, s, a) = \frac{1}{2}(Q^\pi(s, a) - Q_{\boldsymbol{\theta}}(s, a))^2$. If we perform an update using $(s_t, a_t)$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t, s_t, a_t) = \boldsymbol{\theta}_t + \alpha(Q^\pi(s_t, a_t) - Q_{\boldsymbol{\theta}_t}(s_t, a_t)) \nabla_{\boldsymbol{\theta}} Q_{\boldsymbol{\theta}_t}(s_t, a_t)$$

then the interference of that update to one point $(s, a)$ is $L(\boldsymbol{\theta}_{t+1}, s, a) - L(\boldsymbol{\theta}_t, s, a)$. Using a Taylor series expansion, we get the following approximation assuming we have a small step-size $\alpha$:

$$L(\boldsymbol{\theta}_{t+1}; s, a) - L(\boldsymbol{\theta}_t; s, a) \approx \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t; s_t, a_t)^\top (\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t) = -\alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t; s_t, a_t)^\top \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t; s, a)$$

3

This approximation corresponds to *gradient alignment*, which has been used to learn neural networks that are more robust to interference [22, 28]. They measure if $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t; s_t, a_t)^\top \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t; s, a) > 0$, to determine if there is positive generalization between two samples; they generally encourage these dot-products to be positive. Other work used gradient cosine similarity, to measure the level of transferability between tasks [8], and to measure the level of interference between objectives [30]. A somewhat similar measure was used to measure generalization in reinforcement learning [1, 4], using the dot product of the gradients of Q functions $\nabla_{\boldsymbol{\theta}} Q_{\boldsymbol{\theta}_t}(s_t, a_t)^\top \nabla_{\boldsymbol{\theta}} Q_{\boldsymbol{\theta}_t}(s, a)$. This is related in the sense that, for the MSVE with $\delta_t = Q^\pi(s_t, a_t) - Q_{\boldsymbol{\theta}_t}(s_t, a_t)$, $\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t; s_t, a_t)^\top \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_t; s, a) = \delta_t \delta_i \nabla_{\boldsymbol{\theta}} Q_{\boldsymbol{\theta}_t}(s_t, a_t)^\top \nabla_{\boldsymbol{\theta}} Q_{\boldsymbol{\theta}_t}(s, a)$. This measure neglects the direction of the gradients, and so measures both positive generalization as well as interference.

In all the above, interference is measured relative to a chosen performance objective. This performance objective could even be different than the objective directly optimized by the agent. For example, the agent could optimize the MSPBE, as is done by TD-learning, and performance measured with MSVE. We could also have chosen to define the interference using the MSPBE as the performance objective. This is all to say that defining interference is relative to many givens: we need to clearly specify our performance objective, the update for the weights and what samples are used in that update. The same nuance arises in the control setting, which we discuss next.

## 4.2 Interference in Control

Given a value estimation $Q_{\boldsymbol{\theta}}$, let $\pi_{\boldsymbol{\theta}}$ be the policy with respect to the current estimation $Q_{\boldsymbol{\theta}}$. For example, $\pi_{\boldsymbol{\theta}}$ can be the greedy policy w.r.t. $Q_{\boldsymbol{\theta}}$. A previously proposed measure [9, 33] for the quality of a policy is the distance between the action-values for that policy and the optimal action-values

$$\text{OR}(\boldsymbol{\theta}) := \sum_{s \in \mathcal{S}, a \in \mathcal{A}} d(s, a) |Q^*(s, a) - Q^{\pi_{\boldsymbol{\theta}}}(s, a)| = \mathbb{E}_d[Q^*(S, A) - Q^{\pi_{\boldsymbol{\theta}}}(S, A)].$$

We call this the *Optimality Residual* (OR). The distribution $d$ specifies the importance of a state-action pair in the OR. Often, it corresponds to the sampling distribution. For example, $d(s, a) = \nu(s) u(a|s)$ where $\nu$ is a start-state distribution and $u$ is a behavior policy. Notice that the absolute value is not included in the second line, because $Q^*(s, a) \geq Q^{\pi_{\boldsymbol{\theta}}}(s, a)$ for all policies. This objective is one appropriate choice, because the target $Q^*(s, a)$ does not change as the policy changes.

Once we have this objective, the definition for expected interference parallels the prediction setting

$$\text{EI}(\boldsymbol{\theta}_t, B_t) := \mathbb{E}_d[\text{OR}(\boldsymbol{\theta}_{t+1}, (s, a)) - \text{OR}(\boldsymbol{\theta}_t, (s, a))] = \mathbb{E}_d[Q^{\pi_{\boldsymbol{\theta}_t}}(S, A) - Q^{\pi_{\boldsymbol{\theta}_{t+1}}}(S, A)]. \quad (1)$$

where $B_t$ is the mini-batch of data used to update $\theta$ and $\text{OR}(\boldsymbol{\theta}, (s, a)) := Q^*(s, a) - Q^{\pi_{\boldsymbol{\theta}}}(s, a)$.

When running experiments in reinforcement learning, where we have a simulator, it is in fact possible to estimate this quantity. One of the primary motivations for measuring interference is to facilitate investigation by researchers. The OR can be estimated simply by using rollouts from a given $(s, a)$. The policy $\pi_{\boldsymbol{\theta}_t}$ can be started from $(s, a)$ multiple times, generating multiple trajectories. These can be used to get a sample average estimate of the expected return from $(s, a)$ under $\pi_{\boldsymbol{\theta}_t}$. This can then be repeated for $\pi_{\boldsymbol{\theta}_{t+1}}$. The EI is the average OR across $(s, a) \sim d$. In general, though, estimating the EI can be very expensive, because a large number of rollouts may be needed to get accurate estimates [29]. In RL experiments without simulators, it is generally not feasible. In Section 6, we discuss a more practical approach to approximate the EI. First, though, we validate the utility of this true EI.

## 4.3 Summarizing Interference over Time

To determine the impact of interference on agent performance, we need to be provide summary statistics of interference over time. The above are instaneous interference measures, which can tell us how much interference occurred after an update. However, this interference might have long range impacts, and so performance changes on this step might be impacted by interference many steps ago.

A simple choice is to use an average EI over the last window of time. Unfortunately, this choice is problematic because the EI is signed. A negative EI actually indicates improvement—good generalization. An agent could oscillate between positive and negative EIs, with the average appearing to be near zero. The mean of skewed, potentially multi-modal distributions is not a particularly suitable choice, and we can consider other statistics.

To be more systematic about the choice, let $X$ be the random variable corresponding to EI over the desired window of time. For example, if the agent has been learning for 1000 steps, and the desired

window of time is all learning, then $X$ is a scalar RV with a density over the possible instantaneous EIs over this window of 1000 steps. The empirical distribution is the 1000 values of EI.

We consider two statistics, one to measure if the agent had large interference values and the other if interference was highly variable. Catastrophic interference may occur even with only a few steps of very large interference; when reported as an average over time, these large values might be dominated by many small ones. Instead, we can look at the average of the top 10% of interference values—the largest interference—over the window of time. If it is large, then at least 10% of the time the agent had large interference. This type of measure has been used to measure risk, and termed Conditional Value at Risk or sometimes Expected Tail Loss. Correspondingly, we call this the Expected Tail Interference (ETI), defined as

$$\text{ETI}_\alpha(X) = \mathbb{E}[X|X \geq \text{Percentile}_{1-\alpha}(X)] \tag{2}$$

$\text{Percentile}_{1-\alpha}(X)$ is the $(1-\alpha)$-percentile of the distribution of $X$. In our experiments, we set $\alpha = 0.1$.

Finally, we can also provide a more accurate measure of variance by considering the interquartile range: the difference between the 75th and 25th percentiles. We call this the Interference Dispersion

$$\text{Interference Dispersion}(X) = \text{Percentile}_{0.75}(X) - \text{Percentile}_{0.25}(X). \tag{3}$$

Previous work [5] has also considered using conditional value at risk and interquartile range to measure the reliability of reinforcement learning algorithms.

## 5   Empirical Evaluation: Correlation between Interference and Performance

In the section, we evaluate the utility of the interference measures by computing the correlation with several performance measures, including efficiency, stability and episodic return. The goal is both to validate the utility of these measures of interference—as they would not be useful if uncorrelated with performance—as well as to investigate the impact of common deep RL techniques on interference and control performance.

**Environments**   We use Two-Rooms, designed to induce interference across the rooms, and Cart-pole, in which interference has previously been shown to be problematic [14]. Two-Rooms is designed so that the agent has sufficient information to learn optimal policies for each room, but the overlap in inputs for the two rooms is likely to cause interference for standard neural network architectures. Cart-pole involves balancing a pole [3]. Though a simple environment, deep RL agents fail in this domain, or learn unstable policies, as we show in our experiments, and so it provides a useful setting to understand the role of interference on performance. The agent is run a maximal number of steps: 90k for Two-Rooms and 20k for Cart-pole. We run for a fixed number of steps, rather than episode, because otherwise some agents get more environment interactions if they have long episodes. All experiments are averaged over 10 runs.

**Agents**   We investigate well-known deep RL techniques to improving learning, including experience replay, mini-batch updating, Adam optimization (particularly the addition of momentum), and target networks. We consider networks of two hidden layers, with various number of nodes in each layer, batch sizes, buffer sizes and target network update delay. The set of each hyperparameter and other experiment details are in Appendix B.

**Performance Metrics**   We consider four performance measures: *average episodic return* (AER), *consecutive stable performance*, *stable AER* and *sample efficiency*. The AER reflects accumulated reward by the agent, across all steps of learning. It is computed as follows. For each step $i$ during learning, the agent has an associated expected return $\bar{G}_i$: how much reward it currently gets within an episode, in expectation. This can be estimated using multiple runs, or using a recent window of returns, to get estimate $\bar{G}_i$. The AER is the average of these across the last 50% of steps: $\text{AER} := \frac{2}{n} \sum_{i=n/2}^{n} \bar{G}_i$. The AER reflects the agents performance, on average, across the second half of its lifespan. We use the second half to gauge performance, because we are interested in assessing the impact of interference in what the agent has learned.

The AER can be measured using online or offline return. An *offline* $\bar{G}_i$ is an estimate of the expected return $\mathbb{E}_{sa \sim d}[Q^{\pi_{\theta_t}}(s,a)]$ for the policy at time step $i$, measured by averaging over Monte Carlo rollouts. It asks how well the agent would perform if it freezes its policy, and no longer performs updates. An *online* $\bar{G}_i$ is an average over the most recent episodic returns obtained by the agent online, computed using an exponential average with weighting 0.1.
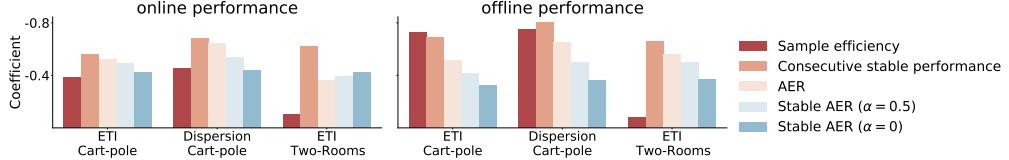
Figure 2: Kendall's rank coefficient on Cart-pole and Two-Rooms. Interference Dispersion is close to zero in Two-Rooms and the correlation is small, so we do not report the numbers in the figure.

We define consecutive stable performance as the maximum number of consecutive steps above a performance threshold (60 step for Two-Rooms, 200 for Cart-pole), divided by the total number of steps. If that number is 1, the agent's threshold performance is maximally stable; if it is zero, it is maximally unstable. Sample complexity corresponds to the first step $i$ that the agent reaches a performance threshold for $k$ consecutive steps (we use $k = 500$), divided by the total number of steps. Sample efficiency is $1-$ sample complexity. If the agent has less interference, we expect the agents to learn a good policy faster, though an agent that generalizes aggressively—and has high interference—might have good efficiency, but may not stably remain at this performance. Finally, stable AER is defined as $\beta \text{AER} + (1 - \beta)E[\bar{G}|\bar{G} \leq \text{Percentile}_{0.1}(\bar{G})]$ where $\beta$ represents the risk profile of the algorithm designer. If the agent has high AER but is unstable, then it will have lower stable AER under a small risk-tolerance $\beta$.

We measure Kendall's Rank-Correlation Coefficient, as in [15], which reflects if two different measures rank agents similarly. It is agnostic to magnitude or precise numbers: if the interference and performance measure both say agent 1 is better than agent 2, then they are reporting similar outcomes. See Appendix B.2 for the formula.

**Results** We show the correlation coefficients between the two interference measures, ETI and Interference Dispersion, and the above four performance measures, in Figure 2. We expect negative correlations, since high interference should correspond to low performance. The overall conclusion is that ETI and Interference Dispersion are both negatively correlated with all performance measures, providing some evidence for the validity of these interference measures.

Next, we look at correlations between performance and interference, at a more fine-grained algorithmic level. To do so, we use a scatter plot for each agent, labeled based on the choice of mini-batch size, buffer size and target network update frequency. The y-axis is performance, and the x-axis interference, allowing a visual inspection of correlation between the two as well as general trends for each algorithm choice. We create one scatter plot per environment, per performance measure, and per interference measure; we include only a subset in Figure 3 and the remainder in Appendix C.1. We find several conclusions. 1) The batch size, buffer size and network size did not seem to have a large impact on either interference or performance; instead, target network frequency was the dominating factor. 2) The target network frequency had opposite performance in the two environments: it increases interference in Cart-pole and reduced it in Two-Rooms. In Two-Rooms, target networks improve stable performance at the cost of reducing efficiency.

Besides optimization, another important component of deep reinforcement learning is the function approximator. Therefore, we conduct an experiment to measure interference within a network, in Appendix C.2. We find that updates on the last layer result in significantly higher interference than updates in the internal layers. The result motivates future research directions to mitigate interference: (1) strategies to mitigate interference in the last layer, and (2) algorithms to learn representation such that updating the last layer on top of these representation is robust to interference.
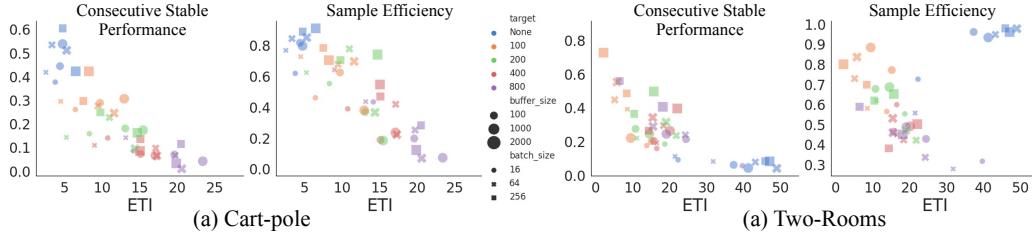


Figure 3: ETI vs sample efficiency and consecutive stable performance in Cart-pole and Two-Rooms, for a variety of Deep RL agents with the smallest network size of 128x128.

6

# 6 Approximating the Expected Interference with TD Errors

It can be impractical to compute the EI, and instead we will need to approximate it. One obvious strategy is simply to estimate $Q^{\pi_t}$ from sampled data, and use estimates $\hat{Q}^{\pi_{t-1}} - \hat{Q}^{\pi_t}$ from a set of sampled states, such as sampled start states. The estimate $\hat{Q}^{\pi_{t-1}}$ could be used to initialize $\hat{Q}^{\pi_t}$, so that fewer updates are needed, as likely $\pi_t$ and $\pi_{t-1}$ are not too different. Unfortunately, such a simple strategy, and ideas related to directly estimating this difference, perform poorly (see Appendix D.1). The issue is that approximation of EI with these estimates seems highly sensitive to accuracy, and it is expensive—or impossible if there is insufficient data—to get highly accurate estimates.

Instead, we want a proxy measure that is more likely to maintain the same sign as EI: reflect performance improvements if the agent got better, and performance degradation otherwise. A natural proxy measure is the Bellman error. The Bellman error reflects if the agent has gotten closer to a fixed point; if it reduced between steps, then this suggests the agent is closer to the fixed point and likely that there is a performance improvement. Fortunately, there is quite a lot of theory relating the Bellman error to $V^\pi$. We extend previous results—namely Lemma 4.3 and Theorem 5.3 in Munos [24]—to the action-value setting. Though relatively straightforward, modifications were needed to allow for differences in distribution over action selection from start states, particularly in the redefinition of concentration coefficients used below. We first present a lemma that upper bounds the EI in terms of the Bellman error. All proofs are in Appendix A.

**Lemma 1.** *Let $Q \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\pi$ be a greedy policy with respect to $Q$. Then*

$$(Q^* - Q^\pi) \le A|\mathcal{T}Q - Q| \qquad and \qquad d(Q^* - Q^\pi) \le dA|\mathcal{T}Q - Q| \qquad (4)$$

*where $A := [(\mathbf{I} - \gamma P\Pi^{\pi^*})^{-1} + (\mathbf{I} - \gamma P\Pi^\pi)^{-1}]$, with $\frac{1-\gamma}{2}A$ a stochastic matrix.*

This bound tells us that we can sample the state-action pairs proportionally to $dA$ to upper bound the OR. Sampling according to $dA$, however, is typically infeasible and here again we need some approximation. We can usually only expect to have a sampled set of transitions, under some behavior policy, resulting in states $s$ in each transition sampled according to some $\mu : \mathcal{S} \to [0, \infty)$. We can additionally bound this sampling error, by using concentration coefficients. Assume $d(s, a) = \nu(s)/|\mathcal{A}|$, where implicitly actions are sample uniformly from $s$. We show the result for any $p \ge 1$ and any policies with non-zero support on all actions in Theorem 1, with the informal result written here with $p = 1$ and uniform policies for simplicity.

**Theorem 1.** *[Informal] Let $\nu$ and $\mu$ be probability measures on $\mathcal{S}$. For $\pi$ greedy w.r.t. $Q \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$*

$$\sum_{s,a} d(s,a)(Q^*(s,a) - Q^{\pi_t}(s,a)) \le \frac{2}{1-\gamma}[C(\nu,\mu)]\sum_{s,a} \frac{\mu(s)}{|\mathcal{A}|}|(\mathcal{T}Q_t)(s,a) - Q_t(s,a)|.$$

The concentration coefficient $C(\nu, \mu)$ reflects differences in state visitation, starting from $\nu$ versus $\mu$, defined precisely in Appendix A. We test three practical choices of $\mu$, with $\hat{d}(s, a) = \mu(s)/|\mathcal{A}|$.

If this approximation is relatively good, then $\mathbb{E}_{(s,a)\sim d}[Q^* - Q^{\pi_{\boldsymbol{\theta}_t}}]$ is approximately proportional to $\mathbb{E}_{(s,a)\sim\hat{d}}[|\mathcal{T}Q_{\boldsymbol{\theta}_t} - Q_{\boldsymbol{\theta}_t}|]$. Recall that EI is $\mathbb{E}_{(s,a)\sim d}[(Q^* - Q^{\pi_{\boldsymbol{\theta}_t}}) - (Q^* - Q^{\pi_{\boldsymbol{\theta}_{t-1}}})]$. Therefore, a potentially reasonable approximation of EI using the Bellman error is $\mathbb{E}_{(s,a)\sim\hat{d}}[|\mathcal{T}Q_{\boldsymbol{\theta}_t} - Q_{\boldsymbol{\theta}_t}| - |\mathcal{T}Q_{\boldsymbol{\theta}_{t-1}} - Q_{\boldsymbol{\theta}_{t-1}}|]$. Even this approximation remains difficult to sample, due to the double sampling problem for Bellman error. Fortunately, we only need to approximate the difference rather than each term. This can be reasonably well approximated uses differences in TD error. Let $\delta(\boldsymbol{\theta}; s, a, r, s') := r + \gamma \max_{a' \in \mathcal{A}} Q_{\boldsymbol{\theta}}(s', a') - Q_{\boldsymbol{\theta}}(s, a)$. By the bias-variance decomposition [2], we can show that

$$\mathbb{E}[\delta(\boldsymbol{\theta}; s, a, r, s')^2] = \mathbb{E}[|\mathcal{T}Q_{\boldsymbol{\theta}}(s,a) - Q_{\boldsymbol{\theta}}(s,a)|^2] + \mathbb{E}[|r + \max_{a'} Q_{\boldsymbol{\theta}}(s', a') - \mathcal{T}Q_{\boldsymbol{\theta}}(s,a)|^2].$$

The first term is the desired Bellman error, and the second term the variance of the targets. If the environment is deterministic, then this variance is zero. More generally, the *Approximate EI*, using TD errors, satisfies

$$\text{AEI} := \mathbb{E}_{(s,a)\sim\hat{d}}[\delta(\boldsymbol{\theta}_t; s, a, r, s')^2 - \delta(\boldsymbol{\theta}_{t-1}; s, a, r, s')^2]$$

$$= \mathbb{E}_{(s,a)\sim\hat{d}}[|\mathcal{T}Q_{\boldsymbol{\theta}_t}(s,a) - Q_{\boldsymbol{\theta}_t}(s,a)|^2 - |\mathcal{T}Q_{\boldsymbol{\theta}_{t-1}}(s,a) - Q_{\boldsymbol{\theta}_{t-1}}(s,a)|^2]$$

$$+ \mathbb{E}_{(s,a)\sim\hat{d}}[|r + \max_{a'} Q_{\boldsymbol{\theta}_t}(s', a') - \mathcal{T}Q_{\boldsymbol{\theta}_t}(s,a)|^2 - |r + \max_{a'} Q_{\boldsymbol{\theta}_{t-1}}(s', a') - \mathcal{T}Q_{\boldsymbol{\theta}_{t-1}}(s,a)|^2].$$

The second expectation is likely to be small, because the two parameters likely have similar variances.

## 6.1 Choosing a Measure $\mu$ to Approximate the Expected Interference

The quality of the approximation is heavily based on the sampling distribution $\mu$. Ideally, we want a measure $\mu$ such that the concentration coefficient $C(\nu, \mu)$ is small, though this is difficult to ascertain. We only have a stream of observations of the agent interacting with the environment, and further can likely only keep a subset of those in a buffer. Sampling from such a buffer is implicitly sampling from a measure $\mu$, where the data acts like a non-parametric sampling distribution. We can consider multiple strategies for adjusting this sampling distribution, both by choosing what to store in the buffer and by re-weighting samples obtained from the buffer, similarly to importance sampling.

We consider three practical choices. The first, which we call *buffer*, involves simply sampling from the most recent transitions. The AEI is then approximated by averaging the differences in TD errors from uniformly sampled transitions from this buffer. The second strategy, which we call *reservoir*, approximates uniform sampling from all the past transition, by maintaining a reservoir buffer. The third strategy, which we call *discounted*, involves reweighting transitions in the reservoir buffer. To approximately sampling from the discounted future state distribution, we re-weight each transition by $(1 - \gamma)\gamma^t$ where t is the number of steps in that episode. We use re-weighting instead of sampling since we would like the measure to have smaller variance.

## 6.2 Empirical Correlation between EI and AEI

We empirically demonstrate that the approximations of interference are correlated with true interference in Two-Rooms and Cart-pole. We sample 1000 transitions, which is a relatively small number compared to the state space, and so more reflective of realistic limitations. We measure Pearson correlation in Figure 4, between EI and AEI per step as well as ETI and Approximate ETI, for two agents. We provide the details in Appendix D.

Though there are several approximation steps above, we find that AEI correlate highly with EI, most clearly in Cart-pole but also in Two-Rooms. The sampling strategies are similarly effective, though reservoir sampling seems to be most effective. We also conduct the same experiments for AEI as in Section 5 are in Appendix D, with similar conclusions, though with slightly reduced correlations to performance measures.
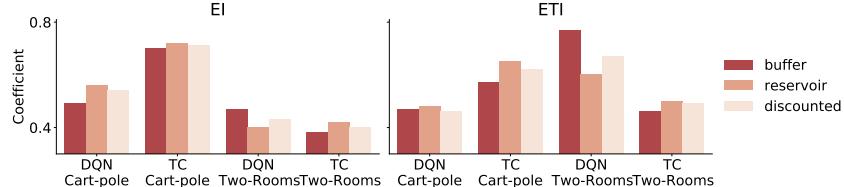


Figure 4: Correlation coefficients with EI and ETI in two domains. All result have p-value $< 0.01$.

## 7 Conclusion

In this paper, we propose a definition of interference for control in RL, and provide a practical approximation using TD errors. We validate the utility of the interference measures by computing the correlation with several performance metric. Using the proposed measures, we provide some insights into interference in deep reinforcement learning algorithms. We highlighted the role of the target network, which we found significantly increased interference and decreased performance in a setting where it was not needed. In another setting, however, the lack of a target network resulted in fast but unstable learning, and we found the opposite conclusion. In both cases, the correlation to interference was clear, for both the true and approximate measures.

This is one of the first papers specifically attempting to define interference for control, and naturally has limitations. One important next step is to expand the set of environments, and agents. In this first small-scale study, we developed a methodology for such experiments, which can be leveraged to extend to new settings. Another important step is to further explore approximations to the true interference, as well as find more clear theoretical reasons why we see that the change in TD-errors performs so well as a proxy. Finally, this paper focuses on deterministic, greedy policies with learned action-values. There is some evidence that a mixture of policies might be more robust to interference [16, 31]. Stochastic policies naturally fit in our definition of EI, but our approximation may not be as suitable.

## Broader Impact

This work focuses on characterizing and understanding an RL agent's behavior. It is unlikely to have a direct impact on society although it may guide future research with such an impact. For example, future research following this work may involve developing stable and practical RL algorithms applied to real world problems.

## References

[1] Joshua Achiam, Ethan Knight, and Pieter Abbeel. Towards characterizing divergence in deep q-learning. *arXiv:1903.08894*, 2019.

[2] András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.

[3] Andrew G Barto, Richard S Sutton, and Charles W Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE transactions on systems, man, and cybernetics*, (5):834–846, 1983.

[4] Emmanuel Bengio, Joelle Pineau, and Doina Precup. Interference and generalization in temporal difference learning. *arXiv preprint arXiv:2003.06350*, 2020.

[5] Stephanie C.Y. Chan, Anoop Korattikara, Sam Fishman, John Canny, and Sergio Guadarrama. Measuring the reliability of reinforcement learning algorithms. In *International Conference on Learning Representations*, 2020.

[6] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 532–547, 2018.

[7] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. *arXiv preprint arXiv:1812.02341*, 2018.

[8] Yunshu Du, Wojciech M Czarnecki, Siddhant M Jayakumar, Razvan Pascanu, and Balaji Lakshminarayanan. Adapting auxiliary losses using gradient similarity. *arXiv preprint arXiv:1812.02224*, 2018.

[9] Amir-massoud Farahmand. Regularization in reinforcement learning. 2011.

[10] Jesse Farebrother, Marlos C Machado, and Michael Bowling. Generalization and regularization in dqn. *arXiv preprint arXiv:1810.00123*, 2018.

[11] William Fedus, Dibya Ghosh, John D. Martin, Marc G. Bellemare, Yoshua Bengio, and Hugo Larochelle. On catastrophic interference in atari 2600 games, 2020.

[12] Yihao Feng, Lihong Li, and Qiang Liu. A kernel loss for solving the bellman equation. In *Advances in Neural Information Processing Systems*, pages 15430–15441, 2019.

[13] Stanislav Fort, Paweł Krzysztof Nowak, and Srini Narayanan. Stiffness: A new perspective on generalization in neural networks. *arXiv preprint arXiv:1901.09491*, 2019.

[14] Benjamin Frederick Goodrich. Neuron clustering for mitigating catastrophic forgetting in supervised and reinforcement learning. 2015.

[15] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019.

[16] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. 2002.

[17] Ronald Kemker, Marc McClure, Angelina Abitino, Tyler L Hayes, and Christopher Kanan. Measuring catastrophic forgetting in neural networks. In *Thirty-second AAAI conference on artificial intelligence*, 2018.

[18] Maurice G Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938.

[19] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 2017.

[20] Long-Ji Lin. Reinforcement learning for robots using neural networks. Technical report, Carnegie-Mellon Univ Pittsburgh PA School of Computer Science, 1993.

[21] Vincent Liu, Raksha Kumaraswamy, Lei Le, and Martha White. The utility of sparse representations for control in reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4384–4391, 2019.

[22] David Lopez-Paz et al. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems*, 2017.

[23] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, and others. Human-level control through deep reinforcement learning. *Nature*, 2015.

[24] Rémi Munos. Performance bounds in l_p-norm for approximate value iteration. *SIAM journal on control and optimization*, 2007.

[25] Cuong V Nguyen, Alessandro Achille, Michael Lam, Tal Hassner, Vijay Mahadevan, and Stefano Soatto. Toward understanding catastrophic forgetting in continual learning. *arXiv preprint arXiv:1908.01091*, 2019.

[26] Charles Packer, Katelyn Gao, Jernej Kos, Philipp Krähenbühl, Vladlen Koltun, and Dawn Song. Assessing generalization in deep reinforcement learning. *arXiv preprint arXiv:1810.12282*, 2018.

[27] Aravind Rajeswaran, Kendall Lowrey, Emanuel V Todorov, and Sham M Kakade. Towards generalization and simplicity in continuous control. In *Advances in Neural Information Processing Systems*, pages 6550–6561, 2017.

[28] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv:1810.11910*, 2018.

[29] Touqir Sajed, Wesley Chung, and Martha White. High-confidence error estimates for learned value functions. *arXiv preprint arXiv:1808.09127*, 2018.

[30] Tom Schaul, Diana Borsa, Joseph Modayil, and Razvan Pascanu. Ray interference: a source of plateaus in deep reinforcement learning. *arXiv:1904.11455*, 2019.

[31] Nino Vieillard, Olivier Pietquin, and Matthieu Geist. Deep conservative policy iteration. *arXiv preprint arXiv:1906.09784*, 2019.

[32] Tao Wang, Michael Bowling, and Dale Schuurmans. Dual representations for dynamic programming and reinforcement learning. In *2007 IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning*, 2007.

[33] Ronald J Williams. Tight performance bounds on greedy policies based on imperfect value functions. Technical report, Citeseer, 1993.

# A Proofs and Technical Details

The proof of Lemma 1 and Theorem 1 are modified from Munos [24]. To begin with, we introduce notation in a matrix form. Define the transition matrix $P \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ where $P(s, a, s') = \Pr(s, a, s')$. Given a policy $\pi$, we define $\Pi^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}||\mathcal{A}|}$ as

$$\Pi^\pi = \begin{pmatrix} \pi(\mathbf{s_1}) & & & \\ & \pi(\mathbf{s_2}) & & \\ & & \ddots & \\ & & & \pi(\mathbf{s_{|\mathcal{S}|}}) \end{pmatrix}$$

where $\pi(\mathbf{s_i}) = \begin{bmatrix} \pi(a_1|s_i) \dots \pi(a_{|\mathcal{A}|}|s_i) \end{bmatrix}$ and all other components are zeros. This notation of $\Pi^\pi$, first introduced in Wang et al. [32], is convenient to use since $\Pi^\pi P$ gives the state to state transition and $P\Pi^\pi$ gives the state-action to state-action transition.

Given an action-value function $Q$, we define the Bellman operator w.r.t. a policy $\pi$ by.

$$\mathcal{T}^\pi Q = \mathbf{r} + \gamma P\Pi^\pi Q.$$

where $\mathbf{r}(s, a) := \sum_{s' \in \mathcal{S}} \Pr(s, a, s') R(s, a, s')$ is the expected immediate reward from state $s$ after taking action $a$. Let $\pi_Q$ denote the greedy policy w.r.t. $Q$, the Bellman optimality operator is defined by

$$\mathcal{T}Q = \mathbf{r} + \gamma P\Pi^{\pi_Q} Q.$$

Since $\pi_Q$ is the greedy policy, we can show that, for any policy $\pi$,

$$\mathcal{T}Q \geq \mathcal{T}^\pi Q.$$

Here $\geq$ denotes the component-wise inequality. Moreover, it is known that the $Q^{\pi_Q}$ is the fixed point of the operator, that is,

$$\mathcal{T}Q^{\pi_Q} = Q^{\pi_Q}.$$

**Lemma 1.** *Let $Q \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\pi$ be a greedy policy with respect to Q. Then*

$$d(Q^* - Q^\pi) \leq dA|\mathcal{T}Q - Q| \tag{5}$$

*where $A := [(\mathbf{I} - \gamma P\Pi^{\pi^*})^{-1} + (\mathbf{I} - \gamma P\Pi^\pi)^{-1}]$, with $\frac{1-\gamma}{2} A$ a stochastic matrix.*

*Proof of Lemma 1.* Using the fact that $Q^* = \mathcal{T}^{\pi^*} Q^*, \mathcal{T}Q \geq \mathcal{T}^{\pi^*} Q$ and $\mathcal{T}Q = \mathcal{T}^\pi Q$, we can show

$$\begin{aligned} Q^* - Q^\pi &= \mathcal{T}^{\pi^*} Q^* - \mathcal{T}^{\pi^*} Q^\pi + \mathcal{T}^{\pi^*} Q^\pi - \mathcal{T}Q + \mathcal{T}Q - \mathcal{T}^\pi Q^\pi \\ &\leq (\mathcal{T}^{\pi^*} Q^* - \mathcal{T}^{\pi^*} Q^\pi + \mathcal{T}^{\pi^*} Q^\pi - \mathcal{T}^{\pi^*} Q) + (\mathcal{T}^\pi Q - \mathcal{T}^\pi Q^\pi) \\ &\leq \gamma P\Pi^{\pi^*} (Q^* - Q^\pi + Q^\pi - Q) + \gamma P\Pi^\pi (Q - Q^\pi) \\ &= \gamma P\Pi^{\pi^*} (Q^* - Q^\pi) + (\gamma P\Pi^{\pi^*} - \gamma P\Pi^\pi)(Q^\pi - Q). \end{aligned}$$

Note that $(\mathbf{I} - \gamma P\Pi^*)$ is invertible, so we have

$$Q^* - Q^\pi \leq (\mathbf{I} - \gamma P\Pi^{\pi^*})^{-1} (\gamma P\Pi^{\pi^*} - \gamma P\Pi^\pi)(Q^\pi - Q).$$

Moreover, we can derive a component-wise equality between the Bellman residual and $(Q^\pi - Q)$:

$$\begin{aligned} (\mathbf{I} - \gamma P\Pi^\pi)(Q^\pi - Q) &= Q^\pi - Q - \gamma P\Pi^\pi Q^\pi + \gamma P\Pi^\pi Q \\ &= Q^\pi - Q + \mathbf{r} + \gamma P\Pi^\pi Q - (\mathbf{r} + P\Pi^\pi Q^\pi) \\ &= Q^\pi - Q + \mathcal{T}^\pi Q - \mathcal{T}^\pi Q^\pi \\ &= \mathcal{T}^\pi Q - Q = \mathcal{T}Q - Q. \end{aligned}$$

Therefore,

$$\begin{aligned} Q^* - Q^\pi &\leq (\mathbf{I} - \gamma P\Pi^{\pi^*})^{-1} (\gamma P\Pi^{\pi^*} - \gamma P\Pi^\pi)(\mathbf{I} - \gamma P\Pi^\pi)^{-1}(\mathcal{T}Q - Q) \\ &= (\mathbf{I} - \gamma P\Pi^{\pi^*})^{-1} [(\mathbf{I} - \gamma P\Pi^\pi) - (\mathbf{I} - \gamma P\Pi^{\pi^*})](\mathbf{I} - \gamma P\Pi^\pi)^{-1}(\mathcal{T}Q - Q) \\ &= [(\mathbf{I} - \gamma P\Pi^*)^{-1} - (\mathbf{I} - \gamma P\Pi^{\pi^*})^{-1}](\mathcal{T}Q - Q) \\ &\leq [(\mathbf{I} - \gamma P\Pi^\pi)^{-1} + (\mathbf{I} - \gamma P\Pi^{\pi^*})^{-1}]|\mathcal{T}Q - Q|. \end{aligned}$$

$\square$

**Definition 1.** *Let $b$ be a policy such that $b(\cdot|s)$ has full support over the action space for all states, $\pi_1, ..., \pi_m$ be a sequence of policies, and $\nu$ and $\mu$ be two measures on $\mathcal{S}$. For any integer $m \geq 1$, we define*

$$c(m) := \sup_{\pi_1, ... \pi_m, s \in \mathcal{S}, a \in \mathcal{A}} \frac{(\nu \Pi^{\pi_1} P \ldots P \Pi^{\pi_m})(s, a)}{\mu \Pi^b(s, a)}.$$

*Let $c(0) := 1$ and $c(m) := \infty$ if $\nu \Pi^{\pi_1} P \ldots \Pi^{\pi_m}$ is not absolutely continuous w.r.t. $\mu \Pi^u$. We define the discounted future state distribution concentration coefficients as*

$$C(\nu, \mu, b) := (1 - \gamma) \sum_{m=0}^{\infty} \gamma^m c(m).$$

**Remark.** In practice, we could choose the behavior policy $b$ as an uniform random policy or a $\epsilon$-greedy policy.

**Theorem 1.** *Let $Q \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$, $\pi$ be a greedy policy with respect to $Q$, $u$ be an uniform policy and $b$ be a behavior policy. Let $\nu$ and $\mu$ be two probability measures on $\mathcal{S}$, and $d = \nu \Pi^u$. Then,*

$$\sum_{s \in \mathcal{S}, a \in \mathcal{A}} d(s, a)|Q^*(s, a) - Q^\pi(s, a)|^p \leq \left[\frac{2}{1 - \gamma}\right]^p C(\nu, \mu, b) \sum_{s,a} \mu(s)b(a|s)|(\mathcal{T}Q)(s, a) - Q(s, a)|^p.$$

*Proof of Theorem 1.* We can write

$$Q^* - Q^\pi \leq A|\mathcal{T}Q - Q|$$

where $A = \left[(\mathbf{I} - \gamma P \Pi^{\pi^*})^{-1} + (\mathbf{I} - \gamma P \Pi^\pi)^{-1}\right]$ and $\frac{1-\gamma}{2} A$ is a stochastic matrix. Then,

$$\sum_{s \in \mathcal{S}, a \in \mathcal{A}} d(s, a)|Q^*(s, a) - Q^\pi(s, a)|^p \leq \left[\frac{2}{1 - \gamma}\right]^p \sum_{s,a} d(s, a) \left[\frac{1 - \gamma}{2} A|\mathcal{T}Q - Q|\right]^p (s, a)$$

$$\leq \left[\frac{2}{1 - \gamma}\right]^p \sum_{s,a} d(s, a) \left[\frac{1 - \gamma}{2} A|\mathcal{T}Q - Q|^p\right](s, a)$$

$$\leq \left[\frac{2}{1 - \gamma}\right]^p C(\nu, \mu, b) \sum_{s,a} \mu(s)b(a|s)|\mathcal{T}Q - Q|^p(s, a)$$

The second inequality follows from Jensen's inequality. The third inequality follows from $d(\frac{1-\gamma}{2} A) \leq (1 - \gamma) \sum_{m=0}^{\infty} \gamma^m c(m) \mu \Pi^b = C(\nu, \mu, b) \mu \Pi^b$. $\quad\square$

## B  Experimental Details

### B.1  Experiment set-up

We experiment with two environments: (1) Two-Rooms: we set the maximum steps per episode to 200, and the number of training steps to 90k, and (2) Cart-pole from OpenAI gym (https://gym.openai.com/): We set the maximum steps per episode to 500, and the number of training steps to 20k. We use a discounting factor $\gamma = 0.99$ in both environments.

The environment Two-Rooms consists of two rooms with different start and goal states. In the first room the agent should navigate up and to the right, and in the second room down and to the left. The input state contains the xy position of the agent, which is in $[0, 1]^2$ for both rooms, and which room the agent is in, which is in $\{0, 1\}$.

For all experiments, we use a two-layer neural network with ReLU activation, and use He intialization to initialize the neural networks.

For the experiments in Section 5, we generate a set of hyper-parameter $\Theta$ by choosing each parameter in the set:

- buffer size $\in \{100, 1000, 2000\}$

- batch size $\in \{16, 64, 256\}$

- Hidden size $\in \{128, 256, 512\}$

- Target network update frequency $\in \{0, 100, 200, 400, 800\}$ where zero means no target network is used

For tile coding, we use 4 tiles and 16 tilings with a constant step size. The step size are searched in the set $\{0.2, 0.1, 0.05, 0.025\}$ by the best online AER. For SR-NN in Appendix C.2, we fix $\beta = 0.1$ and use a grid search for the key parameter: $\lambda_{SKL} \in \{0.01, 0.001, 0.0001\}$. For Section 6.2 and Appendix D.1, we choose a standard neural network with hidden size of 128, batch size of 64, buffer size of 1000 and no target network.

## B.2 Kendall's rank-correlation coefficient

Inspired from [15], we use Kendall's rank-correlation coefficient [18] to check the correlation between a performance metric and a statistics of our interference measures. Let $\Theta$ be a set of hyperparameters and $K := \cup_{\boldsymbol{\theta} \in \Theta} \{(g(\boldsymbol{\theta}), s(\boldsymbol{\theta}))\}$ where $s(\theta)$ is a statistics of our interference measures and $g(\theta)$ is a performance measures corresponding to a hyperparameter configuration $\theta$. Kendall's rank coefficient $\tau$ is defined as

$$\tau := \frac{1}{|K||K-1|} \sum_{(g_1,s_1) \in K} \sum_{(g_2,s_2) \in K/(g_1,s_1)} \operatorname{sign}(g_1 - g_2)\operatorname{sign}(s_1 - s_2).$$

The coefficient varies between $-1$ and $1$.

# C   Additional Experiments of Section 5

## C.1   Correlation between interference and performance

We show the scatter plots for Cart-pole in Figure 7 and 8, and for Two-Rooms in Figure 9. In Two-Rooms, we are interested in the performance when the agent has trained on room 2 for a long time. Therefore, we measure interference and performance for the second half of training on room 2.

The results show show relatively consistent correlation between ETI and performance measures. The notable exception is in Two-Rooms, when there are no target networks. The agents have high interference, but also high AER, for all three variants of AER. The consecutive stable performance and sample efficiency plots sheds some light on why this occurs. Target networks slow learning in this environment, but then maintain stable performance above the threshold for consecutive stable performance. These same agents, though, look worse in terms of AER, than the No Target Network agents, which oscillate more but manage to get to higher performance. The plots are skewed by the fact that, with Target Networks, learning is not quite done when we start measuring interference, in that second half of Room 2. Consequently, though the agent is above the threshold of acceptable performance, it is still on the rise. The lower 10% of the returns is much lower for some of the agents with target networks, than those without, because of this fact. If we allowed the agents to learn for even longer, this point that drop low on the AER plots would likely move up higher, and we would see a clear trend from the cluster of points near zero interference and high performance, the cluster of points with high interference (those without target networks).

## C.2   Measuring interference within a network

In deep reinforcement learning, neural networks are used as the function approximatior. We want to understand how much interference is due to the internal layers and the last layer. The last layer typically does not have an activation; hence we can view a value function as a *two-part approximation* with a representation function and a linear weight $Q_{\mathbf{w},\beta}(s,a) := \boldsymbol{\phi}_\beta(s,a)^\top \mathbf{w}$, where $\mathbf{w} \in \mathbb{R}^d$ is the weights in the last layer and $\boldsymbol{\phi}_\beta : \mathcal{S} \times \mathcal{A} \to \mathbb{R}^d$ is the *representation* learned by the network with weights $\beta$, composed of all the hidden layers in the network. The function $\boldsymbol{\phi}_\beta(s,a)$ corresponds to the last layer in the network, with $\beta$ the weights of the network.

To study interference separately within the network, we use the stochastic block coordinate descent (SBCD) Q-learning to update $\beta$ and $\mathbf{w}$ seperately:

$$\text{Representation Learning Network (RLN) updates: } \beta_{t+1} = \beta_t - \alpha_1 \sum_{i=1}^{B} \nabla_{\beta_t} L(\beta_t, \mathbf{w}_t; s_i, a_i)$$

$$\text{Value Learning Network updates (VLN) updates: } \mathbf{w}_{t+1} = \mathbf{w}_t - \alpha_2 \sum_{i=1}^{B} \nabla_{\mathbf{w}_t} L(\beta_{t+1}, \mathbf{w}_t; s_i, a_i).$$

where $B$ is the mini-batch size, and $\alpha_1$ and $\alpha_2$ are learning rate.

We measure interference for RLN and VLN updates separately at every step, and report the ETI and ETI for approximations in Table 1. We can observe that VLN has much higher ETI than RLN. The result suggests that updates on the last layer result in significantly higher interference than updates on the internal layers, even when we decrease the learning rate for VLN.

We include a baseline SR-NN [21], which learns a sparse representation $\phi_\beta$, to see how representation learning can reduce interference in VLN. SR-NN uses the distributional regularizers to learn sparse representation in neural networks:

$$\min_{\beta} \sum_{i=1}^{B} L(\beta, \mathbf{w}; s_i, a_i) + \lambda_{SKL} \sum_{j=1}^{d} SKL(\bar{\phi}_{\beta,j})$$

where $SKL$ is a regularization on the expected activation, i.e., $\bar{\phi}_{\beta,j} = \sum_{i=1}^{B} \phi_{\beta,j}(s_i, a_i)$ and $\phi_{\beta,j}$ denote the $j$-th component of $\phi_\beta$. Table 1 shows that SBCD with SR-NN has a lower ETI for VLN.

Table 1: ETI for updating VLN and RLN on Cart-pole. We report the control performance as a baseline to see the magnitude of interference. Bold numbers show that ETI for VLN is significantly larger than ETI for RLN. The number are averaged over 10 runs with one standard error.

| | ETI for RLN | ETI for VLN | Control Performance |
|---|---|---|---|
| SBCD Q-learning | $5.05 \pm 0.27$ | $\mathbf{18.23} \pm 3.58$ | $84.45 \pm 0.76$ |
| SBCD Q-learning (smaller $\alpha_2$) | $3.83 \pm 0.26$ | $\mathbf{14.05} \pm 1.59$ | $86.85 \pm 0.39$ |
| SBCD with SR-NN | $3.48 \pm 0.21$ | $\mathbf{5.09} \pm 0.32$ | $89.52 \pm 0.41$ |

# D  Additional Experiments of Section 6

## D.1  Empirical comparison of approximation strategies

Besides approximate EI using TD errors, we test two approximation baselines. First, we could in fact directly approximate the change in Bellman error using recent insights on Kernel Bellman Errors [12], though the approximation is still quite expensive to compute. For example, if we use $M$ transitions to evaluate TD errors, which requires $O(M)$ computation, evaluating Kernel loss requires $O(M^2)$ computation. Hence, we use only 100 transitions (from a reservoir buffer) to evaluate the approximation. Second, we can estimate $\hat{Q}^{\pi_{\theta_t}} \approx Q^{\pi_{\theta_t}}$ from sampled data in the buffer using off-policy policy evaluation (OPE), and directly approximate $EI \approx \mathbb{E}_d[\hat{Q}^{\pi_{\theta_t}}(S, A) - \hat{Q}^{\pi_{\theta_{t+1}}}(S, A)]$ from a set of sampled state-action pairs from $d$. At each evaluation step, we run off-policy SARSA algorithm for 10 epochs over the data stored in a vanilla replay buffer. We call the first baseline *kernel*, and the second baseline *OPE*.

During training, we compute $AEI_i$ and the true measure $EI_i$ every $k$ steps. We collect all data points $(X_i, Y_i)$ from the second half of training steps, over 10 runs, and report Pearson correlation coefficient between AEI and EI. Formally, Pearson correlation coefficient between two sets of measures $X$ and $Y$ is defined as

$$r_{X,Y} = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_i (X_i - \bar{X})^2} \sqrt{\sum_i (Y_i - \bar{Y})^2}}.$$

We show the results in Figure 5. The results suggest that change in TD errors has higher correlation coefficients than other approximation baselines.
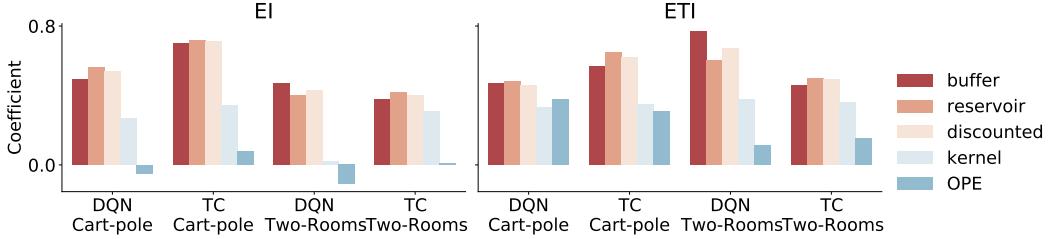
Figure 5: Correlation coefficients with AEI and Approximate ETI in two domains.

## D.2 Results using AEI for Deep RL

In this section, we present the same experiments as in Section 5 and Appendix C.2, with the Approximate EI. We can draw similar conclusions, though with slightly reduced correlations to performance measures. Figure 6 shows that Approximate ETI and ID are negatively correlated with several performance measures. Table 2 shows that VLN has higher Approximate ETI than RLN.
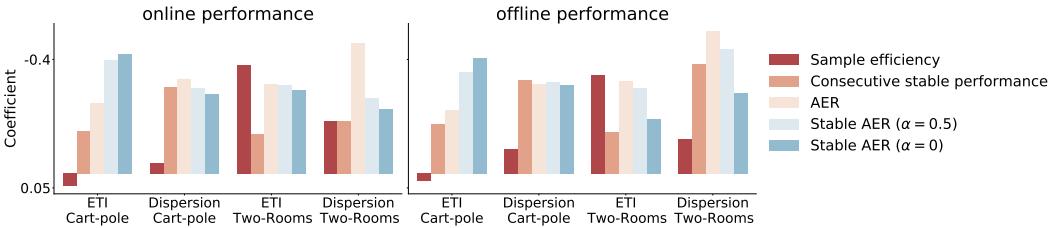


Figure 6: Kendall's rank coefficient on Cart-pole and Two-Rooms.

|  | ETI for RLN | ETI for VLN | Performance |
|---|---|---|---|
| SBCD Q-learning | $0.22 \pm 0.02$ | $\mathbf{1.02} \pm 0.29$ | $84.45 \pm 0.76$ |
| SBCD Q-learning (smaller $\alpha_2$) | $0.08 \pm 0.01$ | $\mathbf{0.29} \pm 0.04$ | $86.85 \pm 0.39$ |
| SBCD with SR-NN | $\mathbf{0.12} \pm 0.01$ | $0.08 \pm 0.01$ | $89.52 \pm 0.41$ |

Table 2: Approximate ETI for updating VLN and RLN separately. Bold numbers show that ETI for VLN is significantly larger than ETI for RLN for SBCD Q-learning. For SBCD with SR-NN, ETI for RLN is larger than ETI for VLN. The number are averaged over 10 runs with one standard error.

(a) Online sample efficiency.

(b) Offline sample efficiency.

(c) Online consecutive stable performance.

(d) Offline consecutive stable performance.

(e) Online AER.

(f) Offline AER.

(g) Online stable AER ($\beta = 0.5$).

(h) Offline stable AER ($\beta = 0.5$).

(i) Online stable AER ($\beta = 0$).

(j) Offline stable AER ($\beta = 0$).

Figure 7: ETI vs performance measures in Cart-pole, for a variety of Deep RL agents.

(a) Online sample efficiency.

(b) Offline sample efficiency.

(c) Online consecutive stable performance.

(d) Offline consecutive stable performance.

(e) Online AER.

(f) Offline AER.

(g) Online stable AER ($\beta = 0.5$).

(h) Offline stable AER ($\beta = 0.5$).

(i) Online stable AER ($\beta = 0$).

(j) Offline stable AER ($\beta = 0$).

Figure 8: Interference Dispersion vs performance measures in Cart-pole, for a variety of Deep RL agents.

(a) Online sample efficiency.

(b) Offline sample efficiency.

(c) Online consecutive stable performance.

(d) Offline consecutive stable performance.

(e) Online AER.

(f) Offline AER.

(g) Online stable AER ($\beta = 0.5$).

(h) Offline stable AER ($\beta = 0.5$).

(i) Online AER ($\beta = 0$).
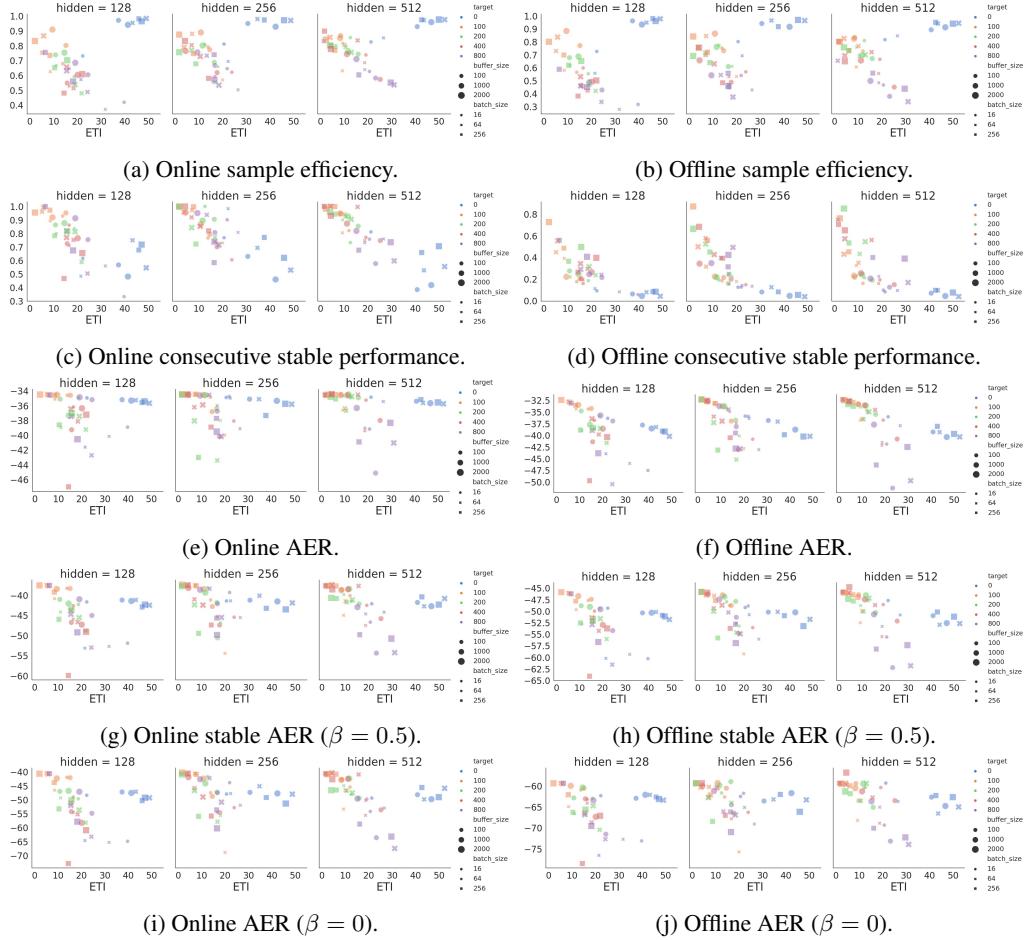
(j) Offline AER ($\beta = 0$).

Figure 9: ETI vs performance measures in Two-Rooms, for a variety of Deep RL agents.