

---

# Evaluating the Performance of Reinforcement Learning Algorithms

---

Scott M. Jordan<sup>1</sup> Yash Chandak<sup>1</sup> Daniel Cohen<sup>1</sup> Mengxue Zhang<sup>1</sup> Philip S. Thomas<sup>1</sup>

## Abstract

Performance evaluations are critical for quantifying algorithmic advances in reinforcement learning. Recent reproducibility analyses have shown that reported performance results are often inconsistent and difficult to replicate. In this work, we argue that the inconsistency of performance stems from the use of flawed evaluation metrics. Taking a step towards ensuring that reported results are consistent, we propose a new comprehensive evaluation methodology for reinforcement learning algorithms that produces reliable measurements of performance both on a single environment and when aggregated across environments. We demonstrate this method by evaluating a broad class of reinforcement learning algorithms on standard benchmark tasks.

## 1. Introduction

When applying reinforcement learning (RL), particularly to real-world applications, it is desirable to have algorithms that reliably achieve high levels of performance without requiring expert knowledge or significant human intervention. For researchers, having algorithms of this type would mean spending less time tuning algorithms to solve benchmark tasks and more time developing solutions to harder problems. Current evaluation practices do not properly account for the uncertainty in the results (Henderson et al., 2018) and neglect the difficulty of applying RL algorithms to a given problem. Consequently, existing RL algorithms are difficult to apply to real-world applications (Dulac-Arnold et al., 2019). To both make and track progress towards developing reliable and easy-to-use algorithms, we propose a principled evaluation procedure that quantifies the difficulty of using an algorithm.

For an evaluation procedure to be useful for measuring the

---

<sup>1</sup>College of Information and Computer Sciences, University of Massachusetts, MA, USA. Correspondence to: Scott Jordan <sjordan@cs.umass.edu>.

usability of RL algorithms, we suggest that it should have four properties. First, to ensure accuracy and reliability, an evaluation procedure should be *scientific*, such that it provides information to answer a research question, tests a specific hypothesis, and quantifies any uncertainty in the results. Second, the performance metric captures the *usability* of the algorithm over a wide variety of environments. For a performance metric to capture the usability of an algorithm, it should include the time and effort spent tuning the algorithm’s hyperparameters (e.g., step-size and policy structure). Third, the evaluation procedure should be *nonexploitative* (Balduzzi et al., 2018), meaning no algorithm should be favored by performing well on an over-represented subset of environments or by abusing a particular score normalization method. Fourth, an evaluation procedure should be *computationally tractable*, meaning that a typical researcher should be able to run the procedure and repeat experiments found in the literature.

As an evaluation procedure requires a question to answer, we pose the following to use throughout the paper: *which algorithm(s) perform well across a wide variety of environments with little or no environment-specific tuning?* Throughout this work, we refer to this question as the *general evaluation question*. This question is different from the one commonly asked in articles proposing a new algorithm, e.g., the *common question* is, can algorithm X outperform other algorithms on tasks A, B, and C? In contrast to the common question, the expected outcome for the general evaluation question is not to find methods that maximize performance with optimal hyperparameters but to identify algorithms that do not require extensive hyperparameter tuning and thus are easy to apply to new problems.

In this paper, we contend that the standard evaluation approaches do not satisfy the above properties, and are not able to answer the general evaluation question. Thus, we develop a new procedure for evaluating RL algorithms that overcomes these limitations and can accurately quantify the uncertainty of performance. The main ideas in our approach are as follows. We present an alternative view of an algorithm such that sampling its performance can be used to answer the general evaluation question. We define a new normalized performance measure, *performance percentiles*, which uses a relative measure of performance to compare algorithms across environments. We show how to use a

game-theoretic approach to construct an aggregate measure of performance that permits quantifying uncertainty. Lastly, we develop a technique, *performance bound propagation* (PBP), to quantify and account for uncertainty throughout the entire evaluation procedure. We provide source code so others may easily apply the methods we develop here.<sup>1</sup>

## 2. Notation and Preliminaries

In this section, we give notation used in this paper along with an overview of an evaluation procedure. In addition to this section, a list of symbols used in this paper is presented in Appendix C. We represent a performance metric of an algorithm,  $i \in \mathcal{A}$ , on an environment,  $j \in \mathcal{M}$ , as a random variable  $X_{i,j}$ . This representation captures the variability of results due to the choice of the random seed controlling the stochastic processes in the algorithm and the environment. The choice of the metric depends on the property being studied and is up to the experiment designer. The performance metric used in this paper is the average of the observed returns from one execution of the entire training procedure, which we refer to as the average return. The cumulative distribution function (CDF),  $F_{X_{i,j}}: \mathbb{R} \rightarrow [0, 1]$ , describes the performance distribution of algorithm  $i$  on environment  $j$  such that  $F_{X_{i,j}}(x) := \Pr(X_{i,j} \leq x)$ . The quantile function,  $Q_{X_{i,j}}(\alpha) := \inf_x \{x \in \mathbb{R} | F_{X_{i,j}}(x) \geq \alpha\}$ , maps a cumulative probability,  $\alpha \in (0, 1)$ , to a score such that  $\alpha$  proportion of samples of  $X_{i,j}$  are less than or equal to  $Q_{X_{i,j}}(\alpha)$ . A normalization function,  $g: \mathbb{R} \times \mathcal{M} \rightarrow \mathbb{R}$ , maps a score,  $x$ , an algorithm receives on an environment,  $j$ , to a normalized score,  $g(x, j)$ , which has a common scale for all environments. In this work, we seek an aggregate performance measure,  $y_i \in \mathbb{R}$ , for an algorithm,  $i$ , such that  $y_i := \sum_{j=1}^{|\mathcal{M}|} q_j \mathbf{E}[g(X_{i,j}, j)]$ , where  $q_j \geq 0$  for all  $j \in \{1, \dots, |\mathcal{M}|\}$  and  $\sum_{j=1}^{|\mathcal{M}|} q_j = 1$ . In Section 4, we discuss choices for the normalizing function  $g$  and weightings  $q$  that satisfy the properties specified in the introduction.

The primary quantities of interest in this paper are the aggregate performance measures for each algorithm and confidence intervals on that measure. Let  $y \in \mathbb{R}^{|\mathcal{A}|}$  be a vector representing the aggregate performance for each algorithm. We desire confidence intervals,  $Y^-, Y^+ \in \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|}$ , such that, for a confidence level  $\delta \in (0, 0.5]$ ,

$$\Pr(\forall i \in \{1, 2, \dots, |\mathcal{A}|\}, y_i \in [Y_i^-, Y_i^+]) \geq 1 - \delta.$$

To compute an aggregate performance measure and its confidence intervals that meet the criteria laid out in the introduction, one must consider the entire evaluation procedure. We view an evaluation procedure to have three main components: data collection, data aggregation, and reporting of the

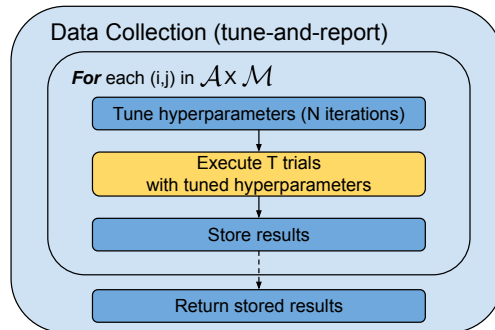


Figure 1. Data collection process of the tune-and-report method. The yellow box indicates trials using different random seeds.

results. During the data collection phase, samples are collected of the performance metric  $X_{i,j}$  for each combination  $(i, j)$  of an algorithm  $i \in \mathcal{A}$  and environment  $j \in \mathcal{M}$ . In the data aggregation phase, all samples of performance are normalized so the metric on each environment is on a similar scale, then they are aggregated to provide a summary of each algorithm’s performance across all environments. Lastly, the uncertainty of the results is quantified and reported.

## 3. Data Collection

In this section, we discuss how common data collection methods are unable to answer the general evaluation question and then present a new method that can. We first highlight the core difference in our approach to previous methods.

The main difference between data collection methods is in how the samples of performance are collected for each algorithm on each environment. Standard approaches rely on first tuning an algorithm’s hyperparameters, i.e., any input to an algorithm that is not the environment, and then generating samples of performance. Our method instead relies on having a definition of an algorithm that can automatically select, sample, or adapt hyperparameters. This method can be used to answer the general evaluation question because its performance measure represents the knowledge required to use the algorithm. We discuss these approaches below.

### 3.1. Current Approaches

A typical evaluation procedure used in RL research is the *tune-and-report* method. As depicted in Figure 1, the tune-and-report method has two phases: a tuning phase and a testing phase. In the tuning phase, hyperparameters are optimized either manually or via a hyperparameter optimization algorithm. Then after tuning, only the best hyperparameters are selected and executed for  $T$  trials using different random seeds to provide an estimate of performance.

<sup>1</sup>Source code for this paper can be found at <https://github.com/ScottJordan/EvaluationOfRLAlgs>.

The tune-and-report data collection method does not satisfy the usability requirement or the scientific requirement. Recall that our objective is to capture the difficulty of using a particular algorithm. Because the tune-and-report method ignores the amount of data used to tune the hyperparameter, an algorithm that only works well after significant tuning could be favored over one that works well without environment-specific tuning, thus, violating the requirements.

Consider an extreme example of an RL algorithm that includes all policy parameters as hyperparameters. This algorithm would then likely be optimal after any iteration of hyperparameter tuning that finds the optimal policy. This effect is more subtle in standard algorithms, where hyperparameter tuning infers problem-specific information about how to search for the optimal policy, (e.g., how much exploration is needed, or how aggressive policy updates can be). Furthermore, this demotivates the creation of algorithms that are easier to use but do not improve performance after finding optimal hyperparameters.

The tune-and-report method violates the scientific property by not accurately capturing the uncertainty of performance. Multiple i.i.d. samples of performance are taken after hyperparameter tuning and used to compute a bound on the mean performance. However, these samples of performance do not account for the randomness due to hyperparameter tuning. As a result, any statistical claim would be inconsistent with repeated evaluations of this method. This has been observed in several studies where further hyperparameter tuning has shown no difference in performance relative to baseline methods (Lucic et al., 2018; Melis et al., 2018).

The evaluation procedure proposed by Dabney (2014) addresses issues with uncertainty due to hyperparameter tuning and performance not capturing the usability of algorithms. Dabney’s evaluation procedure computes performance as a weighted average over all  $N$  iterations of hyperparameter tuning, and the entire tuning process repeats for  $T$  trials. Even though this evaluation procedure fixes the problems with the tune-and-report approach, it violates our computationally tractable property by requiring  $TN$  executions of the algorithm to produce just  $T$  samples of performance. In the case where  $N = 1$  it is not clear how hyperparameters should be set. Furthermore, this style of evaluation does not cover the case where it is prohibitive to perform hyperparameter tuning, e.g., slow simulations, long agent lifetimes, lack of a simulator, and situations where it is dangerous or costly to deploy a bad policy. In these situations, it is desirable for algorithms to be insensitive to the choice of hyperparameters or able to adapt them during a single execution. It is in this setting that the general evaluation question can be answered.

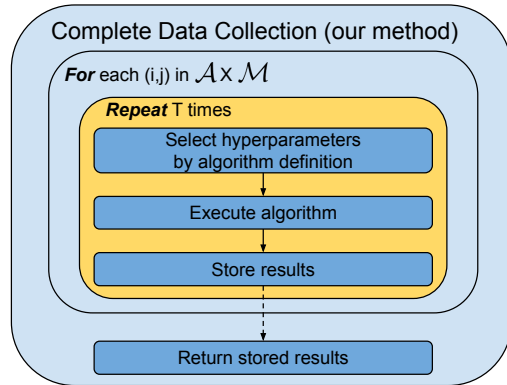


Figure 2. Data collection process using complete algorithm definitions. The yellow box indicates using different random seeds.

### 3.2. Our Approach

In this section, we outline our method, *complete data collection*, that does not rely on hyperparameter tuning. If there were no hyperparameters to tune, evaluating algorithms would be simpler. Unfortunately, how to automatically set hyperparameters has been an understudied area. Thus, we introduce the notion of a complete algorithm definition.

**Definition 1** (Algorithm Completeness). *An algorithm is complete on an environment  $j$ , when defined such that the only required input to the algorithm is meta-information about environment  $j$ , e.g., the number of state features and actions.*

Algorithms with a complete definition can be used on an environment and without specifying any hyperparameters. Note that this does not say that an algorithm cannot receive forms of problem specific knowledge, only that it is not required. A well-defined algorithm will be able to infer effective combinations of hyperparameters or adapt them during learning. There are many ways to make an existing algorithm complete. In this work, algorithms are made complete by defining a distribution from which to randomly sample hyperparameters. Random sampling may produce poor or divergent behavior in the algorithm, but this only indicates that it is not yet known how to set the hyperparameters of the algorithm automatically. Thus, when faced with a new problem, finding decent hyperparameters will be challenging. One way to make an adaptive complete algorithm is to include a hyperparameter optimization method in the algorithm. However, all tuning must be done within the same fixed amount of time and cannot propagate information over trials used to obtain statistical significance.

Figure 2 shows the complete data collection method. For this method we limit the scope of algorithms to only include ones with complete definitions; thus, it does not violate any of the properties specified. This method satisfies the scien-

tific requirement since it is designed to answer the general evaluation question, and the uncertainty of performance can be estimated using all of the trials. Again, this data collection method captures the difficulty of using an algorithm since the complete definition encodes the knowledge necessary for the algorithm to work effectively. The compute time of this method is tractable, since  $T$  executions of the algorithm produces  $T$  independent samples of performance.

The practical effects of using the complete data collection method are as follows. Researchers do not have to spend time tuning each algorithm to try and maximize performance. Fewer algorithm executions are required to obtain a statistically meaningful result. With this data collection method, improving upon algorithm definitions will become significant research contributions and lead to algorithms that are easy to apply to many problems.

## 4. Data Aggregation

Answering the general evaluation question requires a ranking of algorithms according to their performance on all environments  $\mathcal{M}$ . The aggregation step accomplishes this task by combining the performance data generated in the collection phase and summarizing it across all environments. However, data aggregation introduces several challenges. First, each environment has a different range of scores that need to be normalized to a common scale. Second, a uniform weighting of environments can introduce bias. For example, the set of environments might include many slight variants of one domain, giving that domain a larger weight than a single environment coming from a different domain.

### 4.1. Normalization

The goal in score normalization is to project scores from each environment onto the same scale while not being exploitable by the environment weighting. In this section, we first show how existing normalization techniques are exploitable or do not capture the properties of interest. Then we present our normalization technique: performance percentiles.

#### 4.1.1. CURRENT APPROACHES

We examine two normalization techniques: performance ratios and policy percentiles. We discuss other normalization methods in Appendix A. The performance ratio is commonly used with the Arcade Learning Environment to compare the performance of algorithms relative to human performance (Mnih et al., 2015; Machado et al., 2018). The performance ratio of two algorithms  $i$  and  $k$  on an environment  $j$  is  $\mathbf{E}[X_{i,j}]/\mathbf{E}[X_{k,j}]$ . This ratio is sensitive to the location and scale of the performance metric on each environment, such that an environment with scores in the range

$[0, 1]$  will produce larger differences than those on the range  $[1000, 1001]$ . Furthermore, all changes in performance are assumed to be equally challenging, i.e., going from a score of 0.8 to 0.89 is the same difficulty as 0.9 to 0.99. This assumption of linearity of difficulty is not reflected on environments with nonlinear changes in the score as an agent improves, e.g., completing levels in Super Mario.

A critical flaw in the performance ratio is that it can produce an arbitrary ordering of algorithms when combined with the arithmetic mean,  $\sum_j q_j \mathbf{E}[X_{i,j}]/\mathbf{E}[X_{k,j}]$  (Fleming & Wallace, 1986), meaning a different algorithm in the denominator could change the relative rankings. Using the geometric mean can address this weakness of performance ratios, but does not resolve the other issues.

Another normalization technique is *policy percentiles*, a method that projects the score of an algorithm through the performance CDF of random policy search (Dabney, 2014). The normalized score for an algorithm,  $i$ , is  $F_{X_{\Pi,j}}(X_{i,j})$ , where  $F_{X_{\Pi,j}}$  is the performance CDF when a policy is sampled uniformly from a set of policies,  $\Pi$ , on an environment  $j$ , i.e.,  $\pi \sim U(\Pi)$ . Policy percentiles have a unique advantage in that performance is scaled according to how difficult it is to achieve that level of performance relative to random policy search. Unfortunately, policy percentiles rely on specifying  $\Pi$ , which often has a large search space. As a result, most policies will perform poorly, making all scores approach 1.0. It is also infeasible to use when random policy search is unlikely to achieve high levels of performance. Despite these drawbacks, the scaling of scores according to a notion of difficulty is desirable, so we adapt this idea to use any algorithm’s performance as a reference distribution.

#### 4.1.2. OUR APPROACH

An algorithm’s performance distribution can have an interesting shape with large changes in performance that are due to divergence, lucky runs, or simply that small changes to a policy can result in large changes in performance (Jordan et al., 2018). These effects can be seen in Figure 3, where there is a quick rise in cumulative probability for a small increase in performance. Inspired by Dabney (2014)’s policy percentiles, we propose *performance percentiles*, a score normalization technique that can represent these intricacies.

The probability integral transform shows that projecting a random variable through its CDF transforms the variable to be uniform on  $[0, 1]$  (Dodge & Commenges, 2006). Thus, normalizing an algorithm’s performance by its CDF will equally distribute and represent a linear scaling of difficulty across  $[0, 1]$ . When normalizing performance against another algorithm’s performance distribution, the normalized score distribution will shift towards zero when the algorithm is worse than the normalizing distribution and shift towards one when it is superior. As seen in Figure 3, the CDF can

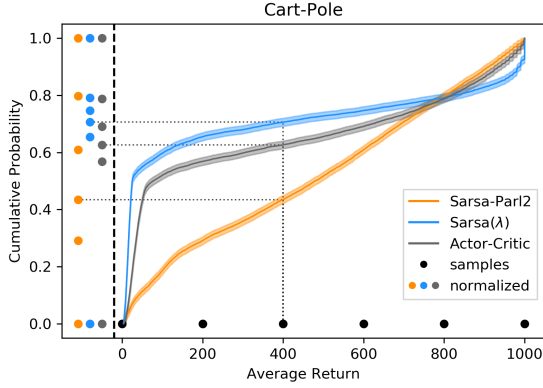


Figure 3. This plot shows the CDF of average returns for the Sarsa-Parl2, Sarsa( $\lambda$ ), and Actor-Critic algorithms on the Cart-Pole environment. Each line represents the empirical CDF using 10,000 trials and the shaded regions represent the 95% confidence intervals. To illustrate how the performance percentiles work, this plot shows how samples of performance (black dots) are normalized by each CDF, producing the normalized scores (colored dots). The correspondence between a single sample and its normalized score is shown by the dotted line.

be seen as encoding the relative difficulty of achieving a given level of performance, where large changes in an algorithm’s CDF output indicate a high degree of difficulty for that algorithm to make an improvement and similarly small changes in output correspond to low change in difficulty. In this context difficulty refers to the amount of random chance (luck) needed to achieve a given level of performance.

To leverage these properties of the CDF, we define performance percentiles, that use a weighted average of each algorithm’s CDF to normalize scores for each environment.

**Definition 2 (Performance Percentile).** *In an evaluation of algorithms,  $\mathcal{A}$ , the performance percentile for a score  $x$  on an environment,  $j$ , is  $F_{\bar{X}_j}(x, w_j)$ , where  $F_{\bar{X}_j}$  is the mixture of CDFs  $F_{\bar{X}_j}(x, w_j) := \sum_{i=1}^{|\mathcal{A}|} w_{j,i} F_{X_{i,j}}(x)$ , with weights  $w_j \in \mathbb{R}^{|\mathcal{A}|}$ ,  $\sum_{i=1}^{|\mathcal{A}|} w_{j,i} = 1$ , and  $\forall i w_{j,i} \geq 0$ .*

So we can say that performance percentiles capture the performance characteristic of an environment relative to some averaged algorithm. We discuss how to set the weights  $w_j$  in the next section.

Performance percentiles are closely related to the concept of (probabilistic) performance profiles (Dolan & Moré, 2002; Barreto et al., 2010). The difference being that performance profiles report the cumulative distribution of normalized performance metrics over a set of tasks (environments), whereas performance percentiles are a technique for normalizing scores on each task (environment).

## 4.2. Summarization

A weighting over environments is needed to form an aggregate measure. We desire a weighting over environments such that no algorithm can exploit the weightings to increase its ranking. Additionally, for the performance percentiles, we need to determine the weighting of algorithms to use as the reference distribution. Inspired by the work of Balduzzi et al. (2018), we propose a weighting of algorithms and environments, using the equilibrium of a two-player game.

In this game, one player,  $p$ , will try to select an algorithm to maximize the aggregate performance, while a second player,  $q$ , chooses the environment and reference algorithm to minimize  $p$ ’s score. Player  $p$ ’s pure strategy space,  $\mathcal{S}_1$ , is the set of algorithms  $\mathcal{A}$ , i.e.,  $p$  plays a strategy  $s_1 = i$  corresponding to an algorithm  $i$ . Player  $q$ ’s pure strategy space,  $\mathcal{S}_2$ , is the cross product of a set of environments,  $\mathcal{M}$ , and algorithms,  $\mathcal{A}$ , i.e., player  $q$  plays a strategy  $s_2 = (j, k)$  corresponding to a choice of environment  $j$  and normalization algorithm  $k$ . We denote the pure strategy space of the game by  $\mathcal{S} := \mathcal{S}_1 \times \mathcal{S}_2$ . A strategy,  $s \in \mathcal{S}$ , can be represented by a tuple  $s = (s_1, s_2) = (i, (j, k))$ .

The utility of strategy  $s$  is measured by a payoff function  $u_p: \mathcal{S} \rightarrow \mathbb{R}$  and  $u_q: \mathcal{S} \rightarrow \mathbb{R}$  for players  $p$  and  $q$  respectively. The game is defined to be zero sum, i.e.,  $u_q(s) = -u_p(s)$ . We define the payoff function to be  $u_p(s) := \mathbf{E}[F_{X_{k,j}}(X_{i,j})]$ . Both players  $p$  and  $q$  sample strategies from probability distributions  $p \in \Delta(\mathcal{S}_1)$  and  $q \in \Delta(\mathcal{S}_2)$ , where  $\Delta(\mathcal{X})$  is the set of all probability distributions over  $\mathcal{X}$ .

The equilibrium solution of this game naturally balances the normalization and environment weightings to counter each algorithm’s strengths without conferring an advantage to a particular algorithm. Thus, the aggregate measure will be useful in answering the general evaluation question.

After finding a solution  $(p^*, q^*)$ , the aggregate performance measure  $y_i$  for an algorithm  $i$  defined as

$$y_i := \sum_{j=1}^{|\mathcal{M}|} \sum_{k=1}^{|\mathcal{A}|} q_{j,k}^* \mathbf{E}[F_{X_{k,j}}(X_{i,j})]. \quad (1)$$

To find a solution  $(p^*, q^*)$ , we employ the  $\alpha$ -Rank technique (Omidshafiei et al., 2019), which returns a stationary distribution over the pure strategy space  $\mathcal{S}$ .  $\alpha$ -Rank allows for efficient computation of both the equilibrium and confidence intervals on the aggregate performance (Rowland et al., 2019). We detail this method and details of our implementation in Appendix B.

## 5. Reporting Results

As it is crucial to quantify the uncertainty of all claimed performance measures, we first discuss how to compute confidence intervals for both single environment and aggregate measures, then give details on displaying the results.

### 5.1. Quantifying Uncertainty

In keeping with our objective to have a scientific evaluation, we require our evaluation procedure to quantify any uncertainty in the results. When concerned with only a single environment, standard concentration inequalities can compute confidence intervals on the mean performance. Similarly, when displaying the distribution of performance, one can apply standard techniques for bounding the empirical distribution of performance. However, computing confidence intervals on the aggregate has additional challenges.

Notice that in (1) computing the aggregate performance requires two unknown values:  $q^*$  and the mean normalized performance,  $\mathbf{E}[F_{X_{k,j}}(X_{i,j})]$ . Since  $q^*$  depends on mean normalized performance, any uncertainty in the mean normalized performance results in uncertainty in  $q^*$ . To compute valid confidence intervals on the aggregate performance, the uncertainty through the entire process must be considered.

We introduce a process to compute the confidence intervals, which we refer to as *performance bound propagation* (PBP). We represent PBP as a function  $\text{PBP}: \mathcal{D} \times \mathbb{R} \rightarrow \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|}$ , which maps a dataset  $D \in \mathcal{D}$  containing all samples of performance and a confidence level  $\delta \in (0, 0.5]$ , to vectors  $Y^-$  and  $Y^+$  representing the lower and upper confidence intervals, i.e.,  $(Y^-, Y^+) = \text{PBP}(D, \delta)$ .

The overall procedure for PBP is as follows, first compute confidence intervals for each  $F_{X_{i,j}}$ , then using these intervals compute confidence intervals on each mean normalized performance, next determine an uncertainty set  $\mathcal{Q}$  for  $q^*$  that results from uncertainty in the mean normalized performance, finally for each algorithm find the minimum and maximum aggregate performance over the uncertainty in the mean normalized performances and  $\mathcal{Q}$ . We provide pseudocode in Appendix C and source code in the repository.

We prove that PBP produces valid confidence intervals for a confidence level  $\delta \in (0, 0.5]$  and a dataset  $D$  containing  $T_{i,j} > 1$  samples of performance for all algorithms  $i \in \mathcal{A}$  and environments  $j \in \mathcal{M}$ .

**Theorem 1.** *If  $(Y^-, Y^+) = \text{PBP}(D, \delta)$ , then*

$$\Pr(\forall i \in 1, 2, \dots, |\mathcal{A}|, y_i \in [Y_i^-, Y_i^+]) \geq 1 - \delta.$$

*Proof.* Although the creation of valid confidence intervals is critical to this contribution, due to space restrictions it is presented in Appendix C.  $\square$

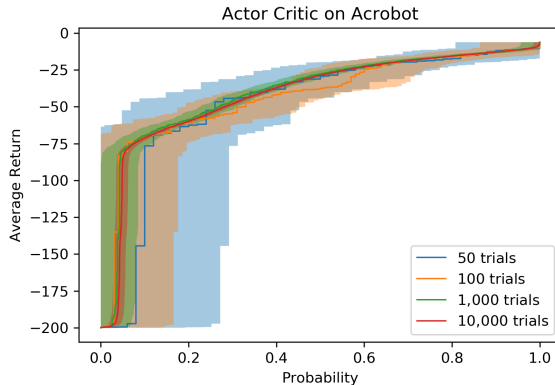


Figure 4. This plot shows the distribution of average returns for the Actor-Critic algorithm on the Acrobot environment. The  $x$ -axis represents a probability and the  $y$ -axis represents the average return such that the proportion of trials that have a value less than or equal to  $y$  is  $x$ , e.g., at  $x = 0.5$ ,  $y$  is the median. Each line represents the empirical quantile function using a different number of trials and the shaded regions represent the 95% confidence intervals computed using concentration inequalities. In this plot, the larger the area under the curve, the better the performance. This plot highlights the large amount of uncertainty when using small sample sizes and how much it decreases with more samples.

### 5.2. Displaying Results

In this section, we describe our method for reporting the results. There are three parts to our method: answering the stated hypothesis, providing tables and plots showing the performance and ranking of algorithms for all environments, and the aggregate score, then for each performance measure, provide confidence intervals to convey uncertainty.

The learning curve plot is a standard in RL and displays a performance metric (often the return) over regular intervals during learning. While this type of plot might be informative for describing some aspects of the algorithm’s performance, it does not directly show the performance metric used to compare algorithms, making visual comparisons less obvious. Therefore, to provide the most information to the reader, we suggest plotting the distribution of performance for each algorithm on each environment. Plotting the distribution of performance has been suggested in many fields as a means to convey more information, (Dolan & Moré, 2002; Farahmand et al., 2010; Reimers & Gurevych, 2017; Cohen et al., 2018). Often in RL, the object is to maximize a metric, so we suggest showing the quantile function over the CDF as it allows for a more natural interpretation of the performance, i.e., the higher the curve, the better the performance (Bellemare et al., 2013). Figure 4 show the performance distribution with 95% confidence intervals for different sample sizes. It is worth noting that when tuning

hyperparameters the data needed to compute these distributions is already being collected, but only the results from the tuned runs are being reported. By only reporting only the tuned performance it shows what an algorithm can achieve not what it is likely to achieve.

## 6. Experimental Results

In this section, we describe and report the results of experiments to illustrate how this evaluation procedure can answer the general evaluation question and identify when a modification to an algorithm or its definition improves performance. We also investigate the reliability of different bounding techniques on the aggregate performance measure.

### 6.1. Experiment Description

To demonstrate the evaluation procedure we compare the algorithms: Actor-Critic with eligibility traces (AC) (Sutton & Barto, 2018),  $Q(\lambda)$ , Sarsa( $\lambda$ ), (Sutton & Barto, 1998), NAC-TD (Morimura et al., 2005; Degris et al., 2012; Thomas, 2014), and proximal policy optimization (PPO) (Schulman et al., 2017). The learning rate is often the most sensitive hyperparameter in RL algorithms. So, we include three versions of Sarsa( $\lambda$ ),  $Q(\lambda)$ , and AC: a base version, a version that scales the step-size with the number of parameters (e.g., Sarsa( $\lambda$ )-s), and an adaptive step-size method, Parl2 (Dabney, 2014), that does not require specifying the step size. Since none of these algorithms have an existing complete definition, we create one by randomly sampling hyperparameters from fixed ranges. We consider all parameters necessary to construct each algorithm, e.g., step-size, function approximator, discount factor, eligibility trace decay. For the continuous state environments, each algorithm employs linear function approximation using the Fourier basis (Konidaris et al., 2011) with a randomly sampled order. See Appendix E for full details of each algorithm.

These algorithms are evaluated on 15 environments, eight discrete MDPs, half with stochastic transition dynamics, and seven continuous state environments: Cart-Pole (Florian, 2007), Mountain Car (Sutton & Barto, 1998), Acrobot (Sutton, 1995), and four variations of the pinball environment (Konidaris & Barto, 2009; Geramifard et al., 2015). For each independent trial, the environments have their dynamics randomly perturbed to help mitigate environment overfitting (Whiteson et al., 2011); see code for details. For further details about the experiment see Appendix F.

While these environments have simple features compared to the Arcade Learning Environment (Bellemare et al., 2013), they remain useful in evaluating RL algorithms for three reasons. First is that experiments finish quickly. Second, the environments provide interesting insights into an algorithm’s behavior. Third, as our results will show, there is not

Aggregate Performance		
Algorithms	Score	Rank
Sarsa-Parl2	0.4623( <b>0.3904</b> , 0.5537)	1 (2,1)
Q-Parl2	0.4366( <b>0.3782</b> , 0.5632)	2 (2,1)
AC-Parl2	0.1578(0.0765, <b>0.3129</b> )	3 (11,3)
Sarsa( $\lambda$ )-s	0.0930(0.0337, 0.2276)	4 (11,3)
AC-s	0.0851(0.0305, 0.2146)	5 (11,3)
Sarsa( $\lambda$ )	0.0831(0.0290, 0.2019)	6 (11,3)
AC	0.0785(0.0275, 0.2033)	7 (11,3)
$Q(\lambda)$ -s	0.0689(0.0237, 0.1973)	8 (11,3)
$Q(\lambda)$	0.0640(0.0214, 0.1780)	9 (11,3)
NAC-TD	0.0516(0.0180, 0.1636)	10 (11,3)
PPO	0.0508(0.0169, 0.1749)	11 (11,3)

Table 1. Aggregate performance measures for each algorithm and their rank. The parentheses contain the intervals computed using PBP and together all hold with 95% confidence. The bolded numbers identify the best ranked statistically significant differences.

yet a complete algorithm that can reliably solve each one.

We execute each algorithm on each environment for 10,000 trials. While this number of trials may seem excessive, our goal is to detect a statistically meaningful result. Detecting such a result is challenging because the variance of RL algorithms performance is high; we are comparing  $|\mathcal{A}| \times |\mathcal{M}| = 165$  random variables, and we do not assume the performances are normally distributed. Computationally, executing ten thousand trials is not burdensome if one uses an efficient programming language such as Julia (Bezanson et al., 2017) or C++, where we have noticed approximately two orders of magnitude faster execution than similar Python implementations. We investigate using smaller sample sizes at the end of this section.

### 6.2. Algorithm Comparison

The aggregate performance measures and confidence intervals are illustrated in Figure 5 and given in Table 1. Appendix I lists the performance tables and distribution plots for each environment. Examining the empirical performances in these figures, we notice two trends. The first is that our evaluation procedure can identify differences that are not noticeable in standard evaluations. For example, all algorithms perform near optimally when tuned properly (indicated by the high end of the performance distribution). The primary differences between algorithms are in the frequency of high performance and divergence (indicated by low end of the performance distribution). Parl2 methods rarely diverge, giving a large boost in performance relative to the standard methods.

The second trend is that our evaluation procedure can identify when theoretical properties do or do not make an al-

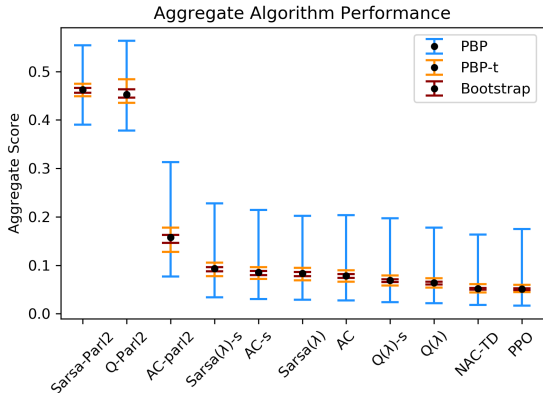


Figure 5. The aggregate performance for each algorithm with confidence intervals using PBP, PBP-t, and bootstrap. The width of each interval is scaled so all intervals hold with 95% confidence.

gorithm more usable. For example, Sarsa( $\lambda$ ) algorithms outperform their Q( $\lambda$ ) counterparts. This result *might* stem from the fact that Sarsa( $\lambda$ ) is known to converge with linear function approximation (Perkins & Precup, 2002) while Q( $\lambda$ ) is known to diverge (Baird, 1995; Wiering, 2004). Additionally, NAC-TD performs worse than AC despite that natural gradients are a superior ascent direction. This result is due in part because it is unknown how to set the three step-sizes in NAC-TD, making it more difficult to use than AC. Together these observations point out the deficiency in the way new algorithms have been evaluated. That is, tuning hyperparameters hides the lack of knowledge required to use the algorithm, introducing bias that favors the new algorithm. In contrast, our method forces this knowledge to be encoded into the algorithm, leading to a more fair and reliable comparison.

### 6.3. Experiment Uncertainty

While the trends discussed above might hold true in general, we must quantify our uncertainty. Based on the confidence intervals given using PBP, we claim with 95% confidence that on these environments and according to our algorithm definitions, Sarsa-Parl2 and Q-Parl2 have a higher aggregate performance of average returns than all other algorithms in the experiment. It is clear that 10,000 trials per algorithm per environment is not enough to detect a unique ranking of algorithms using the nonparametric confidence intervals in PBP. We now consider alternative methods, PBP-t, and the percentile bootstrap. PBP-t replaces the nonparametric intervals in PBP with ones based on the Student’s t-distribution. We detail these methods in Appendix G. From Figure 5, it is clear that both alternative bounds are tighter and thus useful in detecting differences. Since assumptions of these bounds are different and not typically satisfied, it is

	PBP		PBP-t		Bootstrap	
Samples	FR	SIG	FR	SIG	FR	SIG
10	0.0	0.0	1.000	0.00	0.112	0.11
30	0.0	0.0	0.000	0.00	0.092	0.37
100	0.0	0.0	0.000	0.02	0.084	0.74
1,000	0.0	0.0	0.000	0.34	0.057	0.83
10,000	0.0	0.33	0.003	0.83	0.069	0.83

Table 2. Table showing the failure rate (FR) and proportion of significant pairwise comparison (SIG) identified for  $\delta = 0.05$  using different bounding techniques and sample sizes. The first column represents the sample size. The second, third, and fourth columns represent the results for PBP, PBP-t, and bootstrap bound methods respectively. For each sample size, 1,000 experiments were conducted.

unclear if they are valid.

To test the different bounding techniques, we estimate the failure rate of each confidence interval technique at different sample sizes. For this experiment we execute 1,000 trials of the evaluation procedure using sample sizes (trials per algorithm per environment) of 10, 30, 100, 1,000, and 10,000. There are a total of 11.14 million samples per algorithm per environment. To reduce computation costs, we limit this experiment to only include Sarsa( $\lambda$ )-Parl2, Q( $\lambda$ )-Parl2, AC-Parl2, and Sarsa( $\lambda$ )-s. Additionally, we reduce the environment set to be the discrete environments and Mountain Car. We compute the failure rate of the confidence intervals, where a valid confidence interval will have a failure rate less than or equal to  $\delta$ , e.g., for  $\delta = 0.05$  failure rate should be less than  $\leq 5\%$ . We report the failure rate and the proportion of statistically significant pairwise comparisons in Table 2. All methods use the same data, so the results are not independent.

The PBP method has zero failures indicating it is overly conservative. The failure rate of PBP-t is expected to converge to zero as the number of samples increase due to the central limit theorem. PBP-t begins to identify significant results at a sample size of 1,000, but it is only at 10,000 that it can identify all pairwise differences.<sup>2</sup> The bootstrap technique has the tightest intervals, but has a high failure rate.

These results are stochastic and will not necessarily hold with different numbers of algorithms and environments. So, one should use caution in making claims that rely on either PBP-t or bootstrap. Nevertheless, to detect statistically significant results, we recommend running between 1,000, and 10,000 samples, and using the PBP-t over bootstrap.

While this number of trials seems, high it is a necessity as comparison of multiple algorithms over many environments

<sup>2</sup>Sarsa-Parl2 and Q-Parl2 have similar performance on discrete environments so we consider detecting 83% of results optimal.



is a challenging statistical problem with many sources of uncertainty. Thus, one should be skeptical of results that use substantially fewer trials. Additionally, researchers are already conducting many trials that go unreported when tuning hyperparameters. Since our method requires no hyperparameter tuning, researchers can instead spend the same amount of time collecting more trials that can be used to quantify uncertainty.

There are a few ways that the number of trials needed can be reduced. The first is to think carefully about what question one should answer so that only a few algorithms and environments are required. The second is to use active sampling techniques to determine when to stop generating samples of performance for each algorithm environment pair (Rowland et al., 2019). It is important to caution the reader that this process can bias the results if the sequential tests are not accounted for (Howard et al., 2018).

Summarizing our experiments, we make the following observations. Our experiments with complete algorithms show that there is still more work required to make standard RL algorithms work reliably on even extremely simple benchmark problems. As a result of our evaluation procedure, we were able to identify performance differences in algorithms that are not noticeable under standard evaluation procedures. The tests of the confidence intervals suggest that both PBP and PBP-t provide reliable estimates of uncertainty. These outcomes suggest that this evaluation procedure will be useful in comparing the performance of RL algorithms.

## 7. Related Work

This paper is not the first to investigate and address issues in empirically evaluating algorithms. The evaluation of algorithms has become a significant enough topic to spawn its own field of study, known as *experimental algorithmics* (Fleischer et al., 2002; McGeoch, 2012).

In RL, there have been significant efforts to discuss and improve the evaluation of algorithms (Whiteson & Littman, 2011). One common theme has been to produce shared benchmark environments, such as those in the annual reinforcement learning competitions (Whiteson et al., 2010; Dimitrakakis et al., 2014), the Arcade Learning Environment (Bellemare et al., 2013), and numerous others which are too long to list here. Recently, there has been a trend of explicit investigations into the reproducibility of reported results (Henderson et al., 2018; Islam et al., 2017; Khetarpal et al., 2018; Colas et al., 2018). These efforts are in part due to the inadequate experimental practices and reporting in RL and general machine learning (Pineau et al., 2020; Lipton & Steinhardt, 2018). Similar to these studies, this work has been motivated by the need for a more reliable evaluation procedure to compare algorithms. The primary difference

in our work to these is that the knowledge required to use an algorithm gets included in the performance metric.

An important aspect of evaluation not discussed so far in this paper is competitive versus scientific testing (Hooker, 1995). Competitive testing is the practice of having algorithms compete for top performance on benchmark tasks. Scientific testing is the careful experimentation of algorithms to gain insight into how an algorithm works. The main difference in these two approaches is that competitive testing only says *which* algorithms worked well but not *why*, whereas scientific testing directly investigates the what, when, how, or why better performance can be achieved.

There are several examples of recent works using scientific testing to expand our understanding of commonly used methods. Lyle et al. (2019) compares distributional RL approaches using different function approximation schemes showing that distributional approaches are only effective when nonlinear function approximation is used. Tucker et al. (2018) explore the sources of variance reduction in action dependent control variates showing that improvement was small or due to additional bias. Witty et al. (2018) and Atrey et al. (2020) investigate learned behaviors of an agent playing Atari 2600 games using ToyBox (Foley et al., 2018), a tool designed explicitly to enable carefully controlled experimentation of RL agents. While, at first glance the techniques developed here seem to be only compatible with competitive testing, this is only because we specified question with a competitive answer. The techniques developed here, particularly complete algorithm definitions, can be used to accurately evaluate the impact of various algorithmic choices. This allows for the careful experimentation to determine what components are essential to an algorithm.

## 8. Conclusion

The evaluation framework that we propose provides a principled method for evaluating RL algorithms. This approach facilitates fair comparisons of algorithms by removing unintentional biases common in the research setting. By developing a method to establish high-confidence bounds over this approach, we provide the framework necessary for reliable comparisons. We hope that our provided implementations will allow other researchers to easily leverage this approach to report the performances of the algorithms they create.

## Acknowledgements

The authors would like to thank Kaleigh Clary, Emma Tosch, and members of the Autonomous Learning Laboratory: Blossom Metevier, James Kostas, and Chris Nota, for discussion and feedback on various versions of this manuscript. Additionally, we would like to thank the reviewers and meta-reviewers for their comments, which helped improved this

paper.

This work was performed in part using high performance computing equipment obtained under a grant from the Collaborative R&D Fund managed by the Massachusetts Technology Collaborative. This work was supported in part by a gift from Adobe. This work was supported in part by the Center for Intelligent Information Retrieval. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect those of the sponsor. Research reported in this paper was sponsored in part by the CCDC Army Research Laboratory under Cooperative Agreement W911NF-17-2-0196 (ARL IoBT CRA). The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

## References

- Anderson, T. W. Confidence limits for the value of an arbitrary bounded random variable with a continuous distribution function. *Bulletin of The International and Statistical Institute*, 43:249–251, 1969.
- Atrey, A., Clary, K., and Jensen, D. D. Exploratory not explanatory: Counterfactual analysis of saliency maps for deep reinforcement learning. In *8th International Conference on Learning Representations, ICLR*. OpenReview.net, 2020.
- Baird, L. C. Residual algorithms: Reinforcement learning with function approximation. In Prieditis, A. and Russell, S. J. (eds.), *Machine Learning, Proceedings of the Twelfth International Conference on Machine Learning*, pp. 30–37. Morgan Kaufmann, 1995.
- Balduzzi, D., Tuyls, K., Pérolat, J., and Graepel, T. Re-evaluating evaluation. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems, NeurIPS.*, pp. 3272–3283, 2018.
- Barreto, A. M. S., Bernardino, H. S., and Barbosa, H. J. C. Probabilistic performance profiles for the experimental evaluation of stochastic algorithms. In Pelikan, M. and Branke, J. (eds.), *Genetic and Evolutionary Computation Conference, GECCO*, pp. 751–758. ACM, 2010.
- Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, jun 2013.
- Bezanson, J., Edelman, A., Karpinski, S., and Shah, V. B. Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1):65–98, 2017.
- Cohen, D., Jordan, S. M., and Croft, W. B. Distributed evaluations: Ending neural point metrics. *CoRR*, abs/1806.03790, 2018.
- Colas, C., Sigaud, O., and Oudeyer, P. How many random seeds? Statistical power analysis in deep reinforcement learning experiments. *CoRR*, abs/1806.08295, 2018.
- Csáji, B. C., Jungers, R. M., and Blondel, V. D. Pagerank optimization by edge selection. *Discrete Applied Mathematics*, 169:73–87, 2014.
- Dabney, W. C. *Adaptive step-sizes for reinforcement learning*. PhD thesis, University of Massachusetts Amherst, 2014.
- de Kerchove, C., Ninove, L., and Dooren, P. V. Maximizing pagerank via outlinks. *CoRR*, abs/0711.2867, 2007.
- Degrís, T., Pilarski, P. M., and Sutton, R. S. Model-free reinforcement learning with continuous action in practice. In *American Control Conference, ACC*, pp. 2177–2182, 2012.
- Dimitrakakis, C., Li, G., and Tziortziotis, N. The reinforcement learning competition 2014. *AI Magazine*, 35(3): 61–65, 2014.
- Dodge, Y. and Commenges, D. *The Oxford dictionary of statistical terms*. Oxford University Press on Demand, 2006.
- Dolan, E. D. and Moré, J. J. Benchmarking optimization software with performance profiles. *Math. Program.*, 91(2):201–213, 2002.
- Dulac-Arnold, G., Mankowitz, D. J., and Hester, T. Challenges of real-world reinforcement learning. *CoRR*, abs/1904.12901, 2019.
- Dvoretzky, A., Kiefer, J., and Wolfowitz, J. Asymptotic minimax character of a sample distribution function and of the classical multinomial estimator. *Annals of Mathematical Statistics*, 27:642–669, 1956.
- Farahmand, A. M., Ahmadabadi, M. N., Lucas, C., and Araabi, B. N. Interaction of culture-based learning and cooperative co-evolution and its application to automatic behavior-based system design. *IEEE Trans. Evolutionary Computation*, 14(1):23–57, 2010.
- Fercoq, O., Akian, M., Bouhtou, M., and Gaubert, S. Ergodic control and polyhedral approaches to pagerank optimization. *IEEE Trans. Automat. Contr.*, 58(1):134–148, 2013.

- Fleischer, R., Moret, B. M. E., and Schmidt, E. M. (eds.). *Experimental Algorithmics, From Algorithm Design to Robust and Efficient Software [Dagstuhl seminar, September 2000]*, volume 2547 of *Lecture Notes in Computer Science*, 2002. Springer.
- Fleming, P. J. and Wallace, J. J. How not to lie with statistics: The correct way to summarize benchmark results. *Commun. ACM*, 29(3):218–221, 1986.
- Florian, R. V. Correct equations for the dynamics of the cart-pole system. *Center for Cognitive and Neural Studies (Coneural)*, Romania, 2007.
- Foley, J., Tosch, E., Clary, K., and Jensen, D. Toybox: Better Atari Environments for Testing Reinforcement Learning Agents. In *NeurIPS 2018 Workshop on Systems for ML*, 2018.
- Geramifard, A., Dann, C., Klein, R. H., Dabney, W., and How, J. P. RLPy: A value-function-based reinforcement learning framework for education and research. *Journal of Machine Learning Research*, 16:1573–1578, 2015.
- Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., and Meger, D. Deep reinforcement learning that matters. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18)*, pp. 3207–3214, 2018.
- Hooker, J. N. Testing heuristics: We have it all wrong. *Journal of Heuristics*, 1(1):33–42, 1995.
- Howard, S. R., Ramdas, A., McAuliffe, J., and Sekhon, J. S. Uniform, nonparametric, non-asymptotic confidence sequences. *arXiv: Statistics Theory*, 2018.
- Islam, R., Henderson, P., Gomrokchi, M., and Precup, D. Reproducibility of benchmarked deep reinforcement learning tasks for continuous control. *CoRR*, abs/1708.04133, 2017.
- Jordan, S. M., Cohen, D., and Thomas, P. S. Using cumulative distribution based performance analysis to benchmark models. In *Critiquing and Correcting Trends in Machine Learning Workshop at Neural Information Processing Systems*, 2018.
- Khetarpal, K., Ahmed, Z., Cianflone, A., Islam, R., and Pineau, J. Re-evaluate: Reproducibility in evaluating reinforcement learning algorithms. 2018.
- Konidaris, G., Osentoski, S., and Thomas, P. S. Value function approximation in reinforcement learning using the fourier basis. In *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence, AAAI*, 2011.
- Konidaris, G. D. and Barto, A. G. Skill discovery in continuous reinforcement learning domains using skill chaining. In *Advances in Neural Information Processing Systems 22.*, pp. 1015–1023. Curran Associates, Inc., 2009.
- Lipton, Z. C. and Steinhardt, J. Troubling trends in machine learning scholarship. *CoRR*, abs/1807.03341, 2018.
- Lucic, M., Kurach, K., Michalski, M., Gelly, S., and Bousquet, O. Are gans created equal? A large-scale study. In *Advances in Neural Information Processing Systems 31.*, pp. 698–707, 2018.
- Lyle, C., Bellemare, M. G., and Castro, P. S. A comparative analysis of expected and distributional reinforcement learning. In *The Thirty-Third AAAI Conference on Artificial Intelligence*, pp. 4504–4511. AAAI Press, 2019.
- Machado, M. C., Bellemare, M. G., Talvitie, E., Veness, J., Hausknecht, M. J., and Bowling, M. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *J. Artif. Intell. Res.*, 61:523–562, 2018.
- Massart, P. The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *The Annals of Probability*, 18(3): 1269–1283, 1990.
- McGeoch, C. C. *A Guide to Experimental Algorithmics*. Cambridge University Press, 2012.
- Melis, G., Dyer, C., and Blunsom, P. On the state of the art of evaluation in neural language models. In *6th International Conference on Learning Representations, ICLR*. OpenReview.net, 2018.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M. A., Fidjeland, A., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 2015.
- Morimura, T., Uchibe, E. . i. e. j., and Doya, K. Utilizing the natural gradient in temporal difference reinforcement learning with eligibility traces. In *International Symposium on Information Geometry and Its Applications*, pp. 256–263, 2005.
- Omidshafiei, S., Papadimitriou, C., Piliouras, G., Tuyls, K., Rowland, M., Lespiau, J.-B., Czarnecki, W. M., Lanctot, M., Perolat, J., and Munos, R.  $\alpha$ -rank: Multi-agent evaluation by evolution. *Scientific reports*, 9(1):1–29, 2019.
- Page, L., Brin, S., Motwani, R., and Winograd, T. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.

- Perkins, T. J. and Precup, D. A convergent form of approximate policy iteration. In *Advances in Neural Information Processing Systems 15*, pp. 1595–1602. MIT Press, 2002.
- Pineau, J., Vincent-Lamarre, P., Sinha, K., Larivière, V., Beygelzimer, A., d’Alché-Buc, F., Fox, E. B., and Larochelle, H. Improving reproducibility in machine learning research (A report from the NeurIPS 2019 reproducibility program). *CoRR*, abs/2003.12206, 2020.
- Puterman, M. L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. Wiley, 1994.
- Reimers, N. and Gurevych, I. Reporting score distributions makes a difference: Performance study of LSTM-networks for sequence tagging. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP*, pp. 338–348, 2017.
- Rowland, M., Omidshafiei, S., Tuyls, K., Pérolat, J., Valko, M., Piliouras, G., and Munos, R. Multiagent evaluation under incomplete information. In *Advances in Neural Information Processing Systems 3, NeurIPS*, pp. 12270–12282, 2019.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- Sutton, R. S. Generalization in reinforcement learning: Successful examples using sparse coarse coding. In *Advances in Neural Information Processing Systems 8*, pp. 1038–1044, 1995.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning - an introduction*. Adaptive computation and machine learning. MIT Press, 1998.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Thomas, P. Bias in natural actor-critic algorithms. In *Proceedings of the 31th International Conference on Machine Learning, ICML*, pp. 441–448, 2014.
- Tucker, G., Bhupatiraju, S., Gu, S., Turner, R. E., Ghahramani, Z., and Levine, S. The mirage of action-dependent baselines in reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning, ICML*, pp. 5022–5031, 2018.
- Whiteson, S. and Littman, M. L. Introduction to the special issue on empirical evaluations in reinforcement learning. *Mach. Learn.*, 84(1-2):1–6, 2011.
- Whiteson, S., Tanner, B., and White, A. M. Report on the 2008 reinforcement learning competition. *AI Magazine*, 31(2):81–94, 2010.
- Whiteson, S., Tanner, B., Taylor, M. E., and Stone, P. Protecting against evaluation overfitting in empirical reinforcement learning. In *2011 IEEE Symposium on Adaptive Dynamic Programming And Reinforcement Learning, ADPRL*, pp. 120–127. IEEE, 2011.
- Wiering, M. Convergence and divergence in standard and averaging reinforcement learning. In *Machine Learning: ECML 2004, 15th European Conference on Machine Learning*, volume 3201 of *Lecture Notes in Computer Science*, pp. 477–488. Springer, 2004.
- Williams, R. J. and Baird, L. C. Tight performance bounds on greedy policies based on imperfect value functions. 1993.
- Witty, S., Lee, J. K., Tosch, E., Atrey, A., Littman, M., and Jensen, D. Measuring and characterizing generalization in deep reinforcement learning. *arXiv preprint arXiv:1812.02868*, 2018.