

Demystifying Reproducibility in Meta- and Multi-Task Reinforcement Learning

Anonymous Authors¹

Abstract

Establishing the significance of experimental results in reinforcement learning (RL) is difficult. This is compounded by the additional complexity of meta- and multi-task RL (MTRL), a rapidly-growing research area which lacks well-defined baselines. We analyze several design decisions each author must make when they implement a meta-RL or MTRL algorithm, and use over 500 experiments to show that these seemingly-small details can create statistically-significant variations in a single algorithm’s performance that exceed the reported performance differences between algorithms themselves. Informed by this analysis, we precisely define several important hyperparameters, design decisions, and evaluation metrics for meta-RL and MTRL methods, so that we can compare these methods reproducibly. We then provide multi-seed benchmark results for seven of the most popular meta-RL and MTRL algorithms on the most challenging benchmarks currently available. Finally, we share with the community an open source package of these algorithm reference implementations, which use our consistent definitions, achieve state-of-the-art-performance, and seek to follow the original works introducing these algorithms as closely as possible.

1. Introduction

Deep RL methods have excelled at creating agents that perform well on single, narrowly-defined tasks, such as playing a video game which always has the same rules (Berner et al., 2019; Mnih et al., 2013), or controlling a robot to do the same task over-and-over in a world which never changes (Levine et al., 2018). However, to thrive in the real world, an autonomous agent must be able to adapt

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

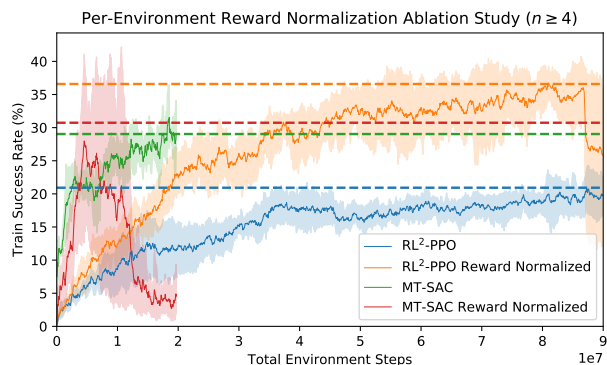


Figure 1. Reward normalization ablation study of RL²-PPO and MT-SAC success rates on ML10/MT10, in which one experiment uses per-task reward normalization and the other does not. Vertical lines indicate best epoch overall for each experimental variant. Mean and 95% confidence interval for five seeds. We find a statistically significant result that reward scaling improves the peak performance of RL²-PPO ($p = 0.04\%$).

quickly to ever-changing environments, objectives, and circumstances (Thrun, 1996). Meta- and multi-task reinforcement learning methods (meta-RL and MTRL, respectively) seek to achieve these capabilities by learning shared structure among tasks, and leveraging it to succeed at new tasks. Meta-RL methods focus on rapidly learning new skills at test time given a small amount of labeled experience, while MTRL methods aim for zero-shot transfer.

Consistently implementing and measuring progress for both approaches presents numerous challenges and complexities beyond those that affect single-task RL: they require defining a training and evaluation pipeline over multiple tasks, which often have disparate specifications (such as reward magnitudes); they are multi-objective optimization problems, and therefore not all metrics used in single-task RL can be naïvely applied to meta- and MTRL (as we explain in detail in Section 3); and these complex algorithms have complex implementations with many design decisions which are unique to meta- and MTRL. We show in Section 4 that **these design decisions matter**, and can statistically explain **between 30 and 100%** of the observed performance difference between popular algorithms.

Concurrent with the rapidly-rising amount of RL research,

the community has taken a renewed interest in reproducibility. This topic was summarized by Henderson et al. (2018a) by asking the question, “[How can we] ensure continued progress in the field by minimizing wasted effort stemming from results that are non-reproducible and easily misinterpreted?”

We believe a measurement-driven strategy is the best path forward for meta- and MTRL, because the fields of supervised learning and single-task RL have achieved great success and real-world impact by adopting three key pillars of reproducible research: standardized challenges (benchmarks) (Russakovsky et al., 2015; Brockman et al., 2016), well-defined performance measurements on those challenges (Papineni et al., 2002), and widely-shared reference implementations, which exhibit repeatable performance measurements on these benchmarks. In this paper, we aim to demystify the reproducibility of popular meta- and MTRL algorithms by contributing to all three of these areas.

Our principal contributions are:

- (1) A set of statistically-tested ablation experiments showing that meta- and MTRL algorithms have many “*small details with large effects*,” which are distinct from their underlying RL components, and can be attributed to differences of between 30-100% of baseline performance.
- (2) An analysis of the random seed sensitivity of 7 popular meta- and MTRL algorithms (MAML, RL², PEARL, MT-PPO, MT-TRPO, MT-SAC, and Task Embeddings (Finn et al., 2017; Duan et al., 2016b; Rakelly et al., 2019; Hausman et al., 2018)) using 9 or more different random seeds on the Meta-World ML45 and MT50 benchmarks (Yu et al., 2019), and for the meta-RL algorithms, a comparison of the number of samples needed to successfully adapt to a new task at test time.
- (3) A set of rigorously-tested, open source reference implementations of all the above algorithms, which were used to perform this analysis, and a well-defined community process for contributing new reference implementations to the set. For comprehensive experimental results and code, see the our supplementary material and the anonymous project website.¹

2. Background and Related Work

Meta- and multi-task RL: The meta- and multi-task RL research community has identified many promising approaches for creating adaptive agents which combine RL with techniques such as recurrent neural networks (Duan et al., 2016b; Wang et al., 2016), attention mechanisms (Mishra et al., 2017), variational inference (Rakelly et al., 2019; Hausman et al., 2018; Sæmunds-

son et al., 2018; Xu et al., 2018; Gupta et al., 2018a), second-order and constrained optimization (Finn et al., 2017; Rothfuss et al., 2018), adversarial methods (Sung et al., 2017), and transfer learning and distillation (Teh et al., 2017).

Reproducibility of deep reinforcement learning: While some works address the reproducibility problem directly (Henderson et al., 2018a), others study different aspects of the problem, such as determining appropriate statistical significance tests and norms for presenting RL results (Colas et al., 2019), or quantifying the sensitivity of RL algorithms to small implementation decisions, hyperparameters, and random seeds (Nagarajan et al., 2018; Liu et al., 2019b; Antoniou et al., 2018; Henderson et al., 2018b; Islam et al., 2017). These decisions are often ambiguous given the formal definition of the algorithm, and are not uniformly applied across implementations. Henderson et al. (2018a) and Engstrom et al. (2020) show that the performance differences created by these decisions can exceed the measured performance margins between single-task RL algorithms, and Deleu et al. (2018) showed similar results on small challenges for a single meta-RL algorithm (MAML). Our work studies differences in implementation which are unique to meta- and MTRL by studying 7 algorithms, and shows that we can explain between 30% and 100% of the measured performance differences between these algorithms, *using implementation decisions alone*.

Benchmarks for RL algorithms: The first pillar of reproducible RL is standardized challenges, and in RL that corresponds to environments. Following from the success of single-task RL benchmarks in many domains (Bellemare et al., 2013; Beattie et al., 2016; Kempka et al., 2016; Tassa et al., 2018; Brockman et al., 2016; Dosovitskiy et al., 2017; Nichol et al., 2018b; Kolve et al., 2017; Brodeur et al., 2017; Plappert et al., 2018; Yang et al., 2019; Stone et al., 2005), researchers have begun to propose a new generation of benchmarks which are specifically designed to evaluate meta- and multi-task RL methods. In addition to introducing the notion of a train-test split among the tasks, tactics used by these new benchmarks include challenging agents to adapt to new environments created with random variations on existing continuous control challenges (Zhang et al., 2018), diverse multi-task robotic environment (James et al., 2019), procedural generation of game levels (Cobbe et al., 2019; Juliani et al., 2019), and principled variations on small tasks which are specifically designed to assess RL agent generalization in several ways (Osband et al., 2020).

Empirical performance studies of RL Empirical studies that seek to characterize the baseline performance of model-free (Duan et al., 2016a) and model-based (Wang et al., 2019) deep RL algorithms on a suite of simulated continuous control tasks (Brockman et al., 2016), and of model-free RL algorithms on real robots (Mahmood et al., 2018),

¹For comprehensive results and code, see <https://icml2020submission6857.github.io/>

provide important goal posts for future proposed methods. We provide a multi-seed of the performance comparison of meta- and MTRL algorithms on a standard benchmark, in hopes of establishing reproducible performance standards for the community.

Metrics and Evaluations for RL: The second pillar of reproducible RL is standardized performance evaluation procedures and metrics. Most works proposing new RL algorithms evaluate them by comparing their performance on benchmark tasks (*i.e.* average discounted return) as a function of total samples collected, and the choice of hyperparameters and random seeds for both the proposed and baseline methods are not specified by any shared standard (Khetarpal et al., 2018). Some authors have proposed new metrics and evaluation protocols in an effort to make measurements are more consistent (Toromanoff et al., 2019), or more predictive of other important aspects of performance, such as generalization and adaptation (Osband et al., 2020), safety and reliability (Chan et al., 2019), or performance in the real world (Dulac-Arnold et al., 2019). A contribution of this work is precise definitions of important hyperparameters, and an evaluation procedure which is relevant to any meta- or multi-task RL method.

Reference implementations for RL: We believe the third and final pillar for reproducible research is widely-disseminated reference implementations. The recent popularity of neural networks and deep RL set off a rapid proliferation of open source libraries for single-task RL (Dhariwal et al., 2017; Liang et al., 2018; Guadarrama et al., 2018; Fan et al., 2018; Schaarschmidt et al., 2018; Kolesnikov & Hrinchuk, 2019; Loon et al., 2019; Fujita et al., 2019). A few of these specifically state adherence to the original algorithm definitions an explicit goal of their implementations (Hill et al., 2018; Castro et al., 2018; Achiam, 2018), which we will show in Section 4 is an important property of reference implementations. While there exist open source implementations for individual meta-RL algorithms (Rothfuss et al., 2018; Finn et al., 2017; Rakelly et al., 2019; Deleu et al., 2019; Deleu, 2018), no repository implements more than 3 algorithms, and *each uses different design decisions and definitions of key hyperparameters*. We show that this makes it challenging to obtain benchmark results which researchers can be confident are comparable, (even if the individual implementations are strong). To address this, we contribute an open source repository of 7 reference implementations with shared hyperparameter definitions and design decisions.

The most closely related work to this one is Yu et al. (2019), which introduces the Meta-World benchmark for meta-RL and multi-task RL algorithms, which we use extensively in this work. Whereas Yu et al. (2019) present single-seed performance results to demonstrate the usefulness of the

proposed benchmark, our work provides results 9 or more random seeds to establish a reproducible standard of performance on that benchmark. Furthermore Yu et al. (2019) were forced to use algorithm implementations from several disparate sources; our work provides a unified set of open source reference implementations for 7 multi-task and meta-RL algorithms, so that the community can benefit from a shared definition of each algorithm and its hyperparameters. Additionally, our work precisely defines—both mathematically and in source code—an evaluation procedure for calculating the performance of any meta-RL or multi-task RL algorithm on any adaptation benchmark such as Meta-World.

3. Defining Consistent Metrics and Hyperparameters for Meta- and Multi-Task Reinforcement Learning

3.1. Preliminaries

We presume some prior task distribution \mathcal{T} , which may be discrete or continuous, and from which the algorithm can sample new tasks. We further presume two disjoint sets have been drawn from this distribution, $T_{\text{train}}, T_{\text{test}} \in \mathcal{T}$ and sample meta-train and meta-test batches from them, respectively. In the case of multi-task RL, $T_{\text{train}} = T_{\text{test}}$ and we may neglect the distinction between meta-train and meta-test evaluation.

3.2. Definitions and Hyperparameters for Meta- and Multi-Task RL

We precisely define important hyperparameters for meta- and multi-task RL, so that we can consistently evaluate and compare algorithms of from these families.

Meta-batch: A batch of some number of *tasks* sampled from the task distribution \mathcal{T} .

Meta-train / meta-test batch size: Size of a meta-train or meta-test batch, measured in *number of tasks*.

Adaptation batch size: Number of samples used by a meta-RL algorithm to adapt to a specific task.

Exploration policy: The policy used by the meta-RL algorithm to acquire adaptation samples on a new task.

Adapted policy: A policy expected to perform well on a single task, computed by a meta-RL algorithm’s adaptation procedure.

Evaluation sample size: The number of trajectories used to measure an adapted policy’s success rate.

Task sampling strategy: A rule used to choose the members of each meta-batch, *e.g.* round-robin, uniform sampling with replacement, k-held out sampling, etc.

Algorithm 1 Meta-RL Evaluation Procedure

Input: Meta-RL algorithm \mathcal{A} , with exploration subroutine $\pi_{\text{explore}} \leftarrow \mathcal{A}_{\text{explore}}()$, and adaptation subroutine $\pi_t \leftarrow \mathcal{A}_{\text{adapt}}(\mathcal{D}_t)$, adaptation batch size N_{adapt} , task rollout function \mathcal{D} , success $\leftarrow \text{rollout}(\pi, t, N)$, evaluation sample size N_{eval} , maximum trajectory length N_{traj}

$n_{\text{success}} \leftarrow 0$

for all $t \in T_{\text{test}}$ **do**

$\pi_{\text{explore}} \leftarrow \mathcal{A}_{\text{explore}}()$

$\mathcal{D}_t, \cdot \leftarrow \text{rollout}(\pi_{\text{explore}}, t, N_{\text{adapt}})$

$\pi_t \leftarrow \mathcal{A}_{\text{adapt}}(\mathcal{D}_t)$

 Clone $\pi_{t,k} \leftarrow \text{clone}(\pi_t)$ N_{eval} times

for 1 to N_{eval} **do**

$\cdot, s \leftarrow \text{rollout}(\pi_{t,k}, t, N_{\text{traj}})$

$n_{\text{success}} \leftarrow n_{\text{success}} + s$

end for

end for

Output: Average Success Rate = $\frac{n_{\text{success}}}{|T_{\text{test}}| \times N_{\text{eval}}}$

3.3. Evaluation Procedure for Meta-RL

We want to fairly compare meta-RL algorithms fairly (so that we do not give too much preference to any algorithm) and consistently (so that we can draw meaningful conclusions across algorithms). To that end, we define a fairly straightforward, but precise, method of defining meta-RL algorithms, and consistently use it in our experiments.

The meta-evaluation procedure allows the meta-RL algorithm to explore for a fixed number of time steps to produce an adapted policy. The adapted policy is then cloned to produce as many adapted policies as the evaluation batch size. Then, each adapted policy is rolled out in the environment for **one** trajectory, and its success or failure is recorded.

See Algorithm 1 for pseudocode describing the procedure in more detail.

When finding the success rate for a single task, we still perform this evaluation procedure as described, with the meta-test task sampling strategy simply repeating the sample task multiple times. We do **not** collect multiple test trajectories per adapted policy, because this could overrepresent the performance of methods that can continue to adapt during the testing paths (in particular, RL²).

3.4. Evaluation Metrics for Meta- and Multi-Task RL

Average task success rate: Mean of the success rates for each individual task in a meta-batch.

Train / test success rate: Average task success rate for the tasks in a meta-train/meta-test batch, respectively.

Single-task sample complexity: The number of total environment steps \mathcal{S}_t needed for an RL algorithm to achieve a success rate (usually 100%) on a given task t .

Multi-task sample complexity (MTRL): The number of

total environment steps (across all tasks) \mathcal{S}_{MT} needed for a multi-task algorithm to achieve a success rate (usually 100%) *on average* across a distribution of tasks \mathcal{T} .

Adaptation sample efficiency (meta-RL): Number of adaptation samples necessary to achieve a particular test success rate. Only applies to meta-RL. May be presented as a curve or summary statistics.

3.5. Statistical Comparisons for Meta- and MTRL Algorithms

The methods and metrics above allow us to precisely measure the performance of multi-task and meta-RL algorithms. However, they do not describe how to evaluate differences between algorithms (or variants on the same algorithm) given these measurements. Consistently analyzing performance requires we define a statistical framework for deriving high-confidence results from our measurements.

Colas et al. (2018) describe a methodology for evaluating RL algorithms which we follow as closely as possible. That methodology requires deriving a scalar *performance measure* from each *run* (which typically corresponds to a training experiment on a single environment with a single random seed) of an algorithm. Although neither Colas et al. (2018) nor Colas et al. (2019) advocate for a particular performance measure, they use the average discounted return of the last 100 iterations of the run in their examples.

Average discounted return is misleading in meta- and MTRL: Meta- and MTRL are multi-objective optimization problems which seek to maximize the average discounted return over a *distribution* of tasks. However, the reward magnitude of any single task is arbitrary, and the reward magnitudes of tasks can vary wildly. For instance, the dense reward magnitudes of successful trajectories for two different tasks in the Meta-World benchmark can range between 10^2 and 10^5 . This means the predictive power of any average return metric calculated across tasks in meta- or MTRL is easily corrupted by tasks with very large or very small reward magnitudes. Because the return between tasks on multi-task benchmarks can be misleading (see Section 4.1.1), we follow Yu et al. (2019) and use average success rate, a sparse binary success signal provided by each environment, for all algorithm comparisons in this work.

Contending with asymptotic instability: Many meta- and MTRL algorithm implementations we evaluate, despite still being useful, are not always asymptotically stable. This means that always measuring on the last 100 iterations of the run, as in Colas et al. (2019), would not accurately represent these algorithms' performance. To minimize the impact of asymptotic instability, we instead use the average success rate of the best 100 sequential iterations as the performance measure of a run.

We use this performance measure throughout Section 4, and

find that this it has relatively low variance across runs and is consistent with our empirical observations of performance, as well the practices of prior works in reporting meta- and MTRL performance results. This measure allows us to make meaningful statistical statements even for highly asymptotically unstable algorithm variants. When we use it, we also typically plot the mean of the performance values of each variant’s runs as a dashed line across the plot.

As recommended in Colas et al. (2018), we apply Welch’s t-test to the performance measures of runs in order to measure the probability, p , of the null hypothesis (*i.e.* that both sets of runs are drawn from the same distribution).

4. Experiments

In this section we evaluate meta- and MTRL algorithms using the Meta-World benchmark (Yu et al., 2019), which is a suite of 50 simulated robotic arm tasks with shared problem structure and internal variations, arranged into three levels of difficulty, and designed to assess the ability of meta-RL and multi-task RL algorithms to exploit shared structure between tasks.

4.1. Implementation Sensitivity in Meta- and Multi-Task RL: *Small Details Make a Big Difference*

In this section we use ablation experiments to investigate a series of implementation choices which authors make when implementing meta- and multi-task RL algorithms.

We show that several minor details, most of which are not specified in the original works introducing these algorithms, can have a large impact on their relative performance. In particular, we show that implementation decisions can create performance differences *for a single algorithm* which can statistically explain absolute performance differences of 10-40% with 95% confidence on the Meta-World benchmarks. This is larger than most performance gaps we observe *between different algorithms*, which range between 0-25%.

4.1.1. REWARD NORMALIZATION

In Figure 1, we compare the performance of RL²-PPO and MT-SAC with and without naive per-task reward normalization. We find a 75% increase in our performance measure, a definitely statistically significant difference ($p = 0.04\%$).

Most recent implementations of reinforcement learning algorithms make use of generalized advantage estimation (GAE) (Schulman et al., 2015b) to normalize rewards. This typically ensures that the RL algorithm can learn regardless of the scale of the reward it receives from the environment. It is common for RL environments within the same benchmark suite to have vastly different reward scales (van Hasselt et al., 2016). Because single-task RL algorithms learn on

a single one of these environments at a time, this difference in reward scales typically does not have a significant effect on learning, and there are methods that have been shown to handle it even in extreme cases (van Hasselt et al., 2016). However, in multi-task RL and meta-RL, differences in reward scale between tasks can cause the algorithm to fixate on a subset of the tasks to the detriment of others. In the Meta-World benchmarks (Yu et al., 2019), rewards on different tasks can differ by multiple orders of magnitude during training. This is likely to cause meta-RL and multi-task RL algorithms to fixate on a subset of tasks. Although a change in the behavior of MT-SAC is clearly visible in the plot, we do not find evidence that it has an effect on performance. Examining individual runs of MT-SAC, we find that some of them reach a very high success rate, but all of them are very unstable. We speculate that applying naive reward scaling in this case makes the learning objective of the critic non-stationary, leading to instability in MT-SAC.

4.1.2. TIME LIMITS IN OFF-POLICY META- AND MULTI-TASK RL ALGORITHMS

In Figure 3 we show the effect of removing the time limit signal from two Meta-World benchmarks. In all Meta-World benchmarks, the time limit signal is conflated with the successful termination signal. Because off-policy algorithms rely on this signal in order to bootstrap correctly terminal states correctly (Mnih et al., 2013), this conflation makes it unnecessarily difficult for off-policy algorithms, such as MT-SAC and PEARL, to perform well in Meta-World benchmarks. We find that removing this conflation does not effect PEARL, but does produce a statistically significant increase in the performance of MT-SAC ($p = 1.5\%$) by 110%.

4.1.3. MAXIMUM ENTROPY RL

In Figure 4, we evaluate the effect of maximum-entropy RL (**not** entropy regularization) on the performance of MT-PPO and RL²-PPO. Applying our statistical method, we find that the added maximum entropy term impedes RL²-PPO from reaching the same peak performance on this benchmark ($p = 0.6\%$). However, we find no evidence that the max entropy term has affected MT-PPO on this benchmark. Unlike single-task RL algorithms, multi-task RL algorithms must simultaneously explore unsolved (low-reward) tasks while maintaining performance on solved (high-reward tasks). This is challenging. Williams & Peng (1991) introduce maximum entropy (MaxEnt), a technique for improving the exploration of policy gradient based RL methods, which adds a term to the reward that encourages high entropy action distributions. Likewise Teh et al. (2017) propose using related but distinct entropy-regularization methods for improving multi-task RL performance.

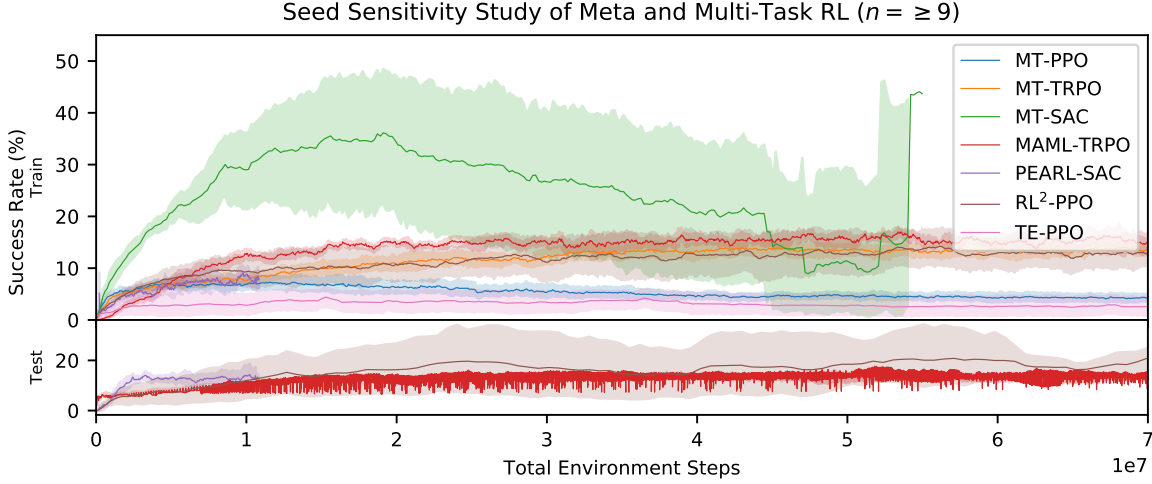


Figure 2. Top: Training success rate of 7 meta- and multi-task RL algorithms on the ML45 (RL²-PPO, MAML-TRPO, PEARL-SAC) and MT50 (MT-PPO, MT-TRPO, TE-PPO) benchmarks. Bottom: Test success rate of the 3 meta-RL algorithms on the ML45 benchmark. Mean and 95% confidence interval over 9 random seeds.

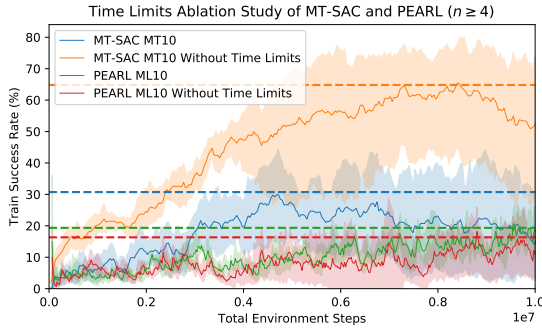


Figure 3. Ablation study of environment time limits on PEARL and MT-SAC success rates on ML10/MT10. By removing the terminal signal at the time limit, MT-SAC performance improves ($p = 1.5\%$) by 110%. Mean and 95% confidence interval using 4 seeds.

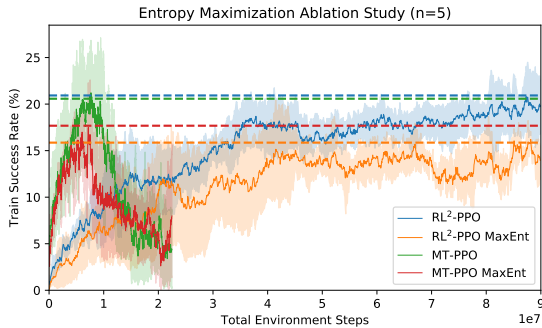


Figure 4. Ablation study of MaxEnt on RL²-PPO and MT-PPO. We find statistically significant evidence that MaxEnt prevents RL²-PPO from reaching its peak performance ($p = 0.6\%$). Horizontal lines indicate the mean performance measure for each experiment variant. Mean and 95% confidence interval using 5 seeds.

4.1.4. EFFECT OF THE “SUBROUTINE ALGORITHM”

In Figure 5 we compare the performance of MAML and RL² using different “subroutine algorithms.” Most meta-RL and multi-task RL algorithms use a single-task RL as some form of “subroutine” (often for solving some meta-objective as a variant of single-task RL) (Finn et al., 2017; Duan et al., 2016b). In some cases, the meta-RL algorithms doesn’t require a specific single-task algorithm, allowing implementers to choose from among one of a family of single-task algorithms. For example, the original description of MAML (Finn et al., 2017) uses the classical policy gradient method to optimize a meta-objective, but in practice it is typically replaced with a different algorithm from the family of policy gradient algorithms. Typically, it is replaced with TRPO (Schulman et al., 2015a), but in some cases (such as some comparisons in Liu et al. (2019a)), it is replaced with PPO. Recent work has indicated that PPO and TRPO can perform very similarly (Engstrom et al., 2020), which might indicate that the choice between them will not significantly effect the behavior of meta-RL algorithms that use them as subroutines. In this experiment, we test the effects of using either PPO or TRPO in two meta-RL algorithms: MAML and RL² on the ML10 benchmark. We find that MAML performs significantly worse when using PPO ($p = 0.016\%$). This confirms the findings of Rothfuss et al. (2018), who cite the low performance of MAML-PPO as one of the primary motivations for the creation of ProMP, which performs MAML-style meta updates without requiring constrained optimization. We also find that RL² performs much better with PPO ($p = 6.3 \times 10^{-6}$), which led us to use it for experiments throughout this work.

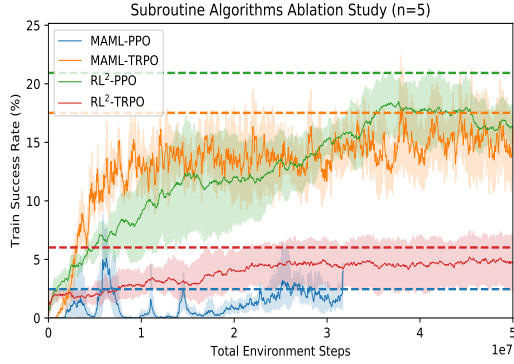


Figure 5. Comparison of using TRPO or PPO as the “subroutine algorithm” for RL² and MAML. As can clearly be seen in the graph, the selection of subroutine algorithm has a significant effect on performance ($p = 6.3 \times 10^{-6}$ for RL² and $p = 0.016\%$ for MAML). Mean and 95% confidence interval for 5 seeds.

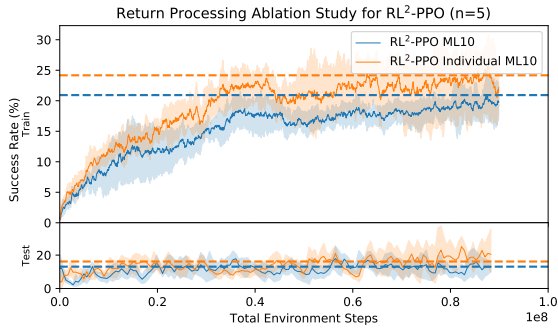


Figure 6. Return processing ablation study for RL²-PPO, with a single line-difference on how the algorithm processes returns. Mean and 95% confidence interval across 5 random seeds.

4.1.5. RETURNS IN RL²

In Figure 6 we compare two slightly different variants of RL². The algorithm RL² (Duan et al., 2016b) implements meta-learning by training a recurrent neural network (RNN) to maximize return over a series of consecutive episodes. However, the most commonly used implementation of RL², which was originally implemented for comparison in (Rothfuss et al., 2018), instead computes return separately for each episode within the trial, then optimizes the sum over those returns. Although there is a slight trend indicating that this modification has improved RL² performance, it is not statistically significant.

4.2. Evaluating Meta- and Multi-Task RL Algorithms on Consistent Metrics

In this section, we evaluate the properties of the selected meta-RL and multi-task RL algorithms. Our general approach to measure a property of an algorithm, such as its consistency across randomness in the environment, is to define a single hyper-parameter that, when varied, should measure that property. Notably, several of these hyper-

parameters would not be practical to control when deploying a meta-RL algorithm. For example, if the agent the algorithm is training is operating in a physical environment, it is not possible to meaningfully control the random seed.

4.2.1. RANDOM SEED SENSITIVITY

The first property of an algorithm we are interested in is the consistency of its performance across random variations in how it is run (both within its environment and the training process itself). In order to evaluate this, we vary the global random seed used by the training process, and measure the resulting success rate on the ML45 and MT50 benchmarks (Yu et al., 2019). The results are presented in Figure 2. As can be seen from the plot, these are very challenging meta-RL and multi-task RL benchmarks. Because of the number of seeds we have run, we are able to perform meaningful statistical comparisons of these algorithms on these benchmarks.

The primary feature of these results, which is clearly visible in the figure, is the high variability in performance across seeds among all algorithms. This indicates that these algorithms have highly variable performance with respect to random variations in the training process.

We apply our statistical method (as described in Section 3.5) to compare the performance of different algorithms on this benchmark. Although we find statistically significant differences between MT-SAC and the other multi-task algorithms ($p = 0.001$ for MT-PPO and $p = 0.01$ for MT-TRPO), we do not find any other differences.

4.2.2. TRANSFER AND ADAPTATION SAMPLE EFFICIENCY

The second property we would like to measure is how the meta-RL algorithms adapt to unseen tasks. In particular, we are interested in how many samples it takes from the new task to adapt. We evaluate this by varying the number of meta-test exploration trajectories. We evaluate every power of two up to 128 trajectories, as well as the performance without any adaptation. The resulting success rates with respect to adaptation experience is shown in Figure 7. We find no evidence that any Meta-RL algorithm successfully adapts to tasks in this benchmark over any of a wide range of adaptation batch sizes. However, given that MAML’s inner adaptation loop is itself an RL algorithm (either PPO or TRPO), which has been independently verified to be able to solve these tasks (Yu et al., 2019), it should eventually reach near 100% success on the new task.

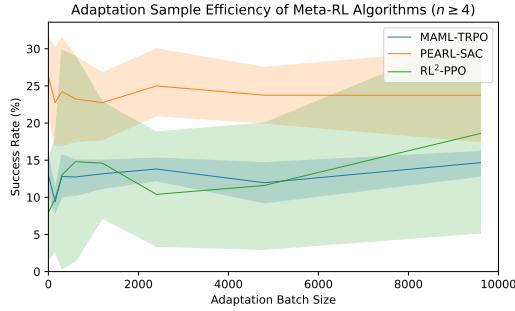


Figure 7. Adaptation sample efficiency of MAML-TRPO, RL²-PPO, and PEARL on the ML45 benchmark. Mean and 95% confidence interval for at least 4 seeds.

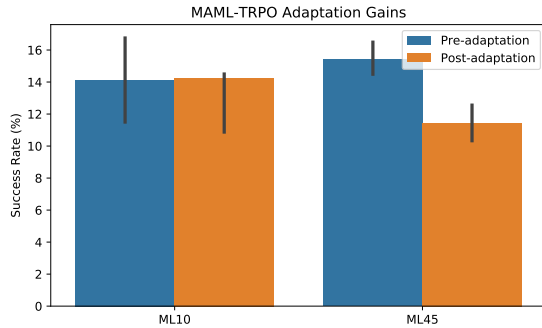


Figure 8. Study of MAML-TRPO gains to adaptation on ML10 and ML45 using an adaptation batch size of 4800 samples. Mean and 95% confidence interval for 5 seeds.

5. Analysis

5.1. Importance of Implementation Decisions in Meta- and MTRL

Figure 2 shows both the train and (where applicable) test performances of several Meta- and MTRL algorithms across 10 random seeds. It is apparent from looking at the graph that many very different algorithms have similar performance (within 25% of each other, with largely overlapping confidence intervals) on the ML45/MT50 benchmarks. However, several of our experiments in Section 4 showed small, often one-line changes that can cause performance changes from 30% to 110% of the baseline on these benchmarks. This implies that making comparisons between these algorithms requires great care in implementation as well as experimentation to discover potential changes caused by implementation differences.

5.2. Does MAML-TRPO actually adapt to shared structure tasks?

Figure 8 shows the test-time performance of MAML-TRPO on ML10 and ML45 over 5 random seeds, before and after the adaptation step. That is, the “Pre-adaptation” performance rates are for the un-adapted exploration policy. Our results suggest with high confidence that on the Meta-World

benchmarks, MAML-TRPO does not actually adapt to new tasks, and in fact its performance *degrades* with adaptation in these environments. Rather, it learns an exploration policy which has high zero-shot performance on the task distribution. We believe these results confirm some of the findings of (Raghu et al., 2019), who find that MAML’s effectiveness is largely attributable to learning representations which are reusable between tasks in the distribution (useful even for zero-shot performance), and not necessarily from rapidly adapting between members of the task distribution (few-shot adaptation).

5.3. Insights on improving meta-RL and multi-task RL benchmarks

The results of our environment time limit ablation study (Fig. 3) suggest that benchmark authors must be careful to not conflate desired environment time limits and the MDP termination signal. This can make benchmarks less-representative of true performance for off-policy algorithm such as MT-SAC, which rely on an accurate signal for MDP termination.

6. Conclusion

In this work, we have shown that in meta- and MTRL research, small implementation decisions make a *big* difference. Our results from Section 4.1 show that implementation decisions *alone* in meta- and MTRL algorithms can create absolute performance differences of 10-35% on the Meta-World benchmarks *for a single algorithm*. Combined with our results from Section 4.2, which show that statistically-significant performance differences between these algorithms range from 0-30%, we can conclude that it is possible that implementation decisions can confound performance improvement claims in meta- and MTRL algorithms research. What’s more, this work only examines a small fraction of possible design decisions and algorithms in meta- and MTRL, and it is very likely there exist many other subtle design decisions which create significant performance differences in these algorithms. This underscores the need for widely-disseminated reference implementations and consistent evaluation protocols, so that we can be confident in algorithmic performance comparisons in meta- and MTRL research.

In the future, we plan to expand our reference implementation library to include more meta- and multi-task RL methods such as PCGrad (Yu et al., 2020), MAESN (Gupta et al., 2018b), ProMP (Rothfuss et al., 2018), REPTILE (Nichol et al., 2018a), SNAIL (Mishra et al., 2017), and QT-Opt (Kalashnikov et al., 2018). We also look forward to welcoming contributions of reference implementations from throughout the research community, and have created a well-defined process for accepting new reference implementations into the library.

References

- Achiam, J. Spinning Up in Deep Reinforcement Learning. 2018.
- Antoniou, A., Edwards, H., and Storkey, A. How to train your maml. *arXiv preprint arXiv:1810.09502*, 2018.
- Beattie, C., Leibo, J. Z., Teplyaev, D., Ward, T., Wainwright, M., Küttler, H., Lefrancq, A., Green, S., Valdés, V., Sadik, A., Schrittwieser, J., Anderson, K., York, S., Cant, M., Cain, A., Bolton, A., Gaffney, S., King, H., Hassabis, D., Legg, S., and Petersen, S. Deepmind lab, 2016.
- 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, 2013.
- Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C., et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym, 2016.
- Brodeur, S., Perez, E., Anand, A., Golemo, F., Celotti, L., Strub, F., Rouat, J., Larochelle, H., and Courville, A. Home: A household multimodal environment. *arXiv preprint arXiv:1711.11017*, 2017.
- Castro, P. S., Moitra, S., Gelada, C., Kumar, S., and Bellemare, M. G. Dopamine: A Research Framework for Deep Reinforcement Learning. 2018. URL <http://arxiv.org/abs/1812.06110>.
- Chan, S. C., Fishman, S., Canny, J., Korattikara, A., and Guadarrama, S. Measuring the reliability of reinforcement learning algorithms. *arXiv preprint arXiv:1912.05663*, 2019.
- Cobbe, K., Hesse, C., Hilton, J., and Schulman, J. Leveraging procedural generation to benchmark reinforcement learning. *arXiv preprint arXiv:1912.01588*, 2019.
- Colas, C., Sigaud, O., and Oudeyer, P.-Y. How many random seeds? statistical power analysis in deep reinforcement learning experiments, 2018.
- Colas, C., Sigaud, O., and Oudeyer, P.-Y. A hitchhiker’s guide to statistical comparisons of reinforcement learning algorithms. *arXiv preprint arXiv:1904.06979*, 2019.
- Deleu, T. Model-Agnostic Meta-Learning for Reinforcement Learning in PyTorch, 2018. Available at: <https://github.com/tristandeleu/pytorch-maml-rl>.
- Deleu, T., Guioy, S., and Hosseini, S. On the reproducibility of gradient-based meta-reinforcement learning baselines. 2018.
- Deleu, T., Würfl, T., Samiei, M., Cohen, J. P., and Bengio, Y. Torchmeta: A Meta-Learning library for PyTorch, 2019. URL <https://arxiv.org/abs/1909.06576>. Available at: <https://github.com/tristandeleu/pytorch-meta>.
- Dhariwal, P., Hesse, C., Klimov, O., Nichol, A., Plappert, M., Radford, A., Schulman, J., Sidor, S., Wu, Y., and Zhokhov, P. Openai baselines. <https://github.com/openai/baselines>, 2017.
- Dosovitskiy, A., Ros, G., Codevilla, F., Lopez, A., and Koltun, V. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pp. 1–16, 2017.
- Duan, Y., Chen, X., Houthooft, R., Schulman, J., and Abbeel, P. Benchmarking deep reinforcement learning for continuous control. In *International Conference on Machine Learning*, pp. 1329–1338, 2016a.
- Duan, Y., Schulman, J., Chen, X., Bartlett, P. L., Sutskever, I., and Abbeel, P. RL²: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016b.
- Dulac-Arnold, G., Mankowitz, D., and Hester, T. Challenges of real-world reinforcement learning. *arXiv preprint arXiv:1904.12901*, 2019.
- Engstrom, L., Ilyas, A., Santurkar, S., Tsipras, D., Janoos, F., Rudolph, L., and Madry, A. Implementation matters in deep {rl}: A case study on {ppo} and {trpo}. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rletNlrtPB>.
- Fan, L., Zhu, Y., Zhu, J., Liu, Z., Zeng, O., Gupta, A., Creus-Costa, J., Savarese, S., and Fei-Fei, L. Surreal: Open-source reinforcement learning framework and robot manipulation benchmark. In *Conference on Robot Learning*, pp. 767–782, 2018.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1126–1135. JMLR. org, 2017.
- Fujita, Y., Kataoka, T., Nagarajan, P., and Ishikawa, T. Chainerrl: A deep reinforcement learning library. *arXiv preprint arXiv:1912.03905*, 2019.
- Guadarrama, S., Korattikara, A., Ramirez, O., Castro, P., Ethan, H., Fishman, S., Wang, K., Gonina, E., Wu,

- N., Kokiopoulou, E., Sbaiz, L., Smith, J., Bartók, G., Berent, J., Harris, C., Vanhoucke, V., and Brevdo, E. TF-Agents: A library for reinforcement learning in tensorflow. <https://github.com/tensorflow/agents>, 2018. URL <https://github.com/tensorflow/agents>. [Online; accessed 25-June-2019].
- Gupta, A., Mendonca, R., Liu, Y., Abbeel, P., and Levine, S. Meta-reinforcement learning of structured exploration strategies. *ArXiv*, abs/1802.07245, 2018a.
- Gupta, A., Mendonca, R., Liu, Y., Abbeel, P., and Levine, S. Meta-reinforcement learning of structured exploration strategies. *CoRR*, abs/1802.07245, 2018b. URL <http://arxiv.org/abs/1802.07245>.
- Hausman, K., Springenberg, J. T., Wang, Z., Heess, N., and Riedmiller, M. Learning an embedding space for transferable robot skills. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rk07ZXZRb>.
- Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., and Meger, D. Deep reinforcement learning that matters. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018a.
- Henderson, P., Romoff, J., and Pineau, J. Where did my optimum go?: An empirical analysis of gradient descent optimization in policy gradient methods. *arXiv preprint arXiv:1810.02525*, 2018b.
- Hill, A., Raffin, A., Ernestus, M., Gleave, A., Kanervisto, A., Traore, R., Dhariwal, P., Hesse, C., Klimov, O., Nichol, A., Plappert, M., Radford, A., Schulman, J., Sidor, S., and Wu, Y. Stable baselines. <https://github.com/hill-a/stable-baselines>, 2018.
- Islam, R., Henderson, P., Gomrokchi, M., and Precup, D. Reproducibility of benchmarked deep reinforcement learning tasks for continuous control. *arXiv preprint arXiv:1708.04133*, 2017.
- James, S., Ma, Z., Arrojo, D. R., and Davison, A. J. Rlbench: The robot learning benchmark & learning environment. *arXiv preprint arXiv:1909.12271*, 2019.
- Juliani, A., Khalifa, A., Berges, V. P., Harper, J., Teng, E., Henry, H., Crespi, A., Togelius, J., and Lange, D. Obstacle tower: A generalization challenge in vision, control, and planning. In *28th International Joint Conference on Artificial Intelligence, IJCAI 2019*, pp. 2684–2691. International Joint Conferences on Artificial Intelligence, 2019.
- Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pp. 651–673, 2018.
- Kempka, M., Wydmuch, M., Runc, G., Toczek, J., and Jaśkowski, W. ViZDoom: A Doom-based AI research platform for visual reinforcement learning. In *IEEE Conference on Computational Intelligence and Games*, pp. 341–348, Santorini, Greece, Sep 2016. IEEE. URL <http://arxiv.org/abs/1605.02097>. The best paper award.
- Khetarpal, K., Ahmed, Z., Cianflone, A., Islam, R., and Pineau, J. Re-evaluate: Reproducibility in evaluating reinforcement learning algorithms. 2018.
- Kolesnikov, S. and Hrinchuk, O. Catalyst. rl: A distributed framework for reproducible rl research. *arXiv preprint arXiv:1903.00027*, 2019.
- Kolve, E., Mottaghi, R., Han, W., VanderBilt, E., Weihs, L., Herrasti, A., Gordon, D., Zhu, Y., Gupta, A., and Farhadi, A. AI2-THOR: An Interactive 3D Environment for Visual AI. *arXiv*, 2017.
- Levine, S., Pastor, P., Krizhevsky, A., Ibarz, J., and Quillen, D. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International Journal of Robotics Research*, 37(4-5):421–436, 2018.
- Liang, E., Liaw, R., Nishihara, R., Moritz, P., Fox, R., Goldberg, K., Gonzalez, J., Jordan, M., and Stoica, I. Rllib: Abstractions for distributed reinforcement learning. In *International Conference on Machine Learning*, pp. 3053–3062, 2018.
- Liu, H., Socher, R., and Xiong, C. Taming MAML: Efficient unbiased meta-reinforcement learning. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 4061–4071, Long Beach, California, USA, 09–15 Jun 2019a. PMLR. URL <http://proceedings.mlr.press/v97/liu19g.html>.
- Liu, Z., Li, X., Kang, B., and Darrell, T. Regularization matters in policy optimization. *arXiv preprint arXiv:1910.09191*, 2019b.
- Loon, K. W., Graesser, L., and Cvitkovic, M. Slm lab: A comprehensive benchmark and modular software framework for reproducible deep reinforcement learning. *arXiv preprint arXiv:1912.12482*, 2019.

- Mahmood, A. R., Korenkevych, D., Vasani, G., Ma, W., and Bergstra, J. Benchmarking reinforcement learning algorithms on real-world robots. In *Conference on Robot Learning*, pp. 561–591, 2018.
- Mishra, N., Rohaninejad, M., Chen, X., and Abbeel, P. Meta-learning with temporal convolutions. *CoRR*, abs/1707.03141, 2017. URL <http://arxiv.org/abs/1707.03141>.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Nagarajan, P., Warnell, G., and Stone, P. The impact of nondeterminism on reproducibility in deep reinforcement learning. 2018.
- Nichol, A., Achiam, J., and Schulman, J. On first-order meta-learning algorithms. *CoRR*, abs/1803.02999, 2018a. URL <http://arxiv.org/abs/1803.02999>.
- Nichol, A., Pfau, V., Hesse, C., Klimov, O., and Schulman, J. Gotta learn fast: A new benchmark for generalization in rl, 2018b.
- Osband, I., Doron, Y., Hessel, M., Aslanides, J., Sezener, E., Saraiva, A., McKinney, K., Lattimore, T., Szepevari, C., Singh, S., Roy, B. V., Sutton, R., Silver, D., and Hasselt, H. V. Behaviour suite for reinforcement learning. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rygf-kSYwH>.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pp. 311–318. Association for Computational Linguistics, 2002.
- Plappert, M., Andrychowicz, M., Ray, A., McGrew, B., Baker, B., Powell, G., Schneider, J., Tobin, J., Chociej, M., Welinder, P., et al. Multi-goal reinforcement learning: Challenging robotics environments and request for research. *arXiv preprint arXiv:1802.09464*, 2018.
- Raghu, A., Raghu, M., Bengio, S., and Vinyals, O. Rapid learning or feature reuse? towards understanding the effectiveness of maml. *arXiv preprint arXiv:1909.09157*, 2019.
- Rakelly, K., Zhou, A., Finn, C., Levine, S., and Quillen, D. Efficient off-policy meta-reinforcement learning via probabilistic context variables. In *International Conference on Machine Learning*, pp. 5331–5340, 2019.
- Rothfuss, J., Lee, D., Clavera, I., Asfour, T., and Abbeel, P. Prompt: Proximal meta-policy search. *arXiv preprint arXiv:1810.06784*, 2018.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3): 211–252, 2015.
- Sæmundsson, S., Hofmann, K., and Deisenroth, M. P. Meta reinforcement learning with latent variable gaussian processes. In *UAI*, 2018.
- Schaarschmidt, M., Mika, S., Fricke, K., and Yoneki, E. Rlgraph: Modular computation graphs for deep reinforcement learning. *arXiv preprint arXiv:1810.09028*, 2018.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust region policy optimization. In *International conference on machine learning*, pp. 1889–1897, 2015a.
- Schulman, J., Moritz, P., Levine, S., Jordan, M. I., and Abbeel, P. High-dimensional continuous control using generalized advantage estimation. *CoRR*, abs/1506.02438, 2015b.
- Stone, P., Kuhlmann, G., Taylor, M. E., and Liu, Y. Keep-away soccer: From machine learning testbed to benchmark. In *Robot Soccer World Cup*, pp. 93–105. Springer, 2005.
- Sung, F., Zhang, L., Xiang, T., Hospedales, T. M., and Yang, Y. Learning to learn: Meta-critic networks for sample efficient learning. *ArXiv*, abs/1706.09529, 2017.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., de Las Casas, D., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., Lillicrap, T., and Riedmiller, M. Deepmind control suite, 2018.
- Teh, Y., Bapst, V., Czarnecki, W. M., Quan, J., Kirkpatrick, J., Hadsell, R., Heess, N., and Pascanu, R. Distral: Robust multitask reinforcement learning. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30*, pp. 4496–4506. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7036-distral-robust-multitask-reinforcement-learning.pdf>.
- Thrun, S. *Explanation-Based Neural Network Learning: A Lifelong Learning Approach*. Kluwer Academic Publishers, Boston, MA, 1996.
- Toromanoff, M., Wirbel, E., and Moutarde, F. Is deep reinforcement learning really superhuman on atari? *arXiv preprint arXiv:1908.04683*, 2019.

- van Hasselt, H. P., Guez, A., Guez, A., Hessel, M., Mnih, V., and Silver, D. Learning values across many orders of magnitude. In Lee, D. D., Sugiyama, M., Luxburg, U. V., Guyon, I., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems* 29, pp. 4287–4295. Curran Associates, Inc., 2016. URL <http://papers.nips.cc/paper/6076-learning-values-across-many-orders-of-magnitude.pdf>.
- Wang, J. X., Kurth-Nelson, Z., Tirumala, D., Soyer, H., Leibo, J. Z., Munos, R., Blundell, C., Kumaran, D., and Botvinick, M. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.
- Wang, T., Bao, X., Clavera, I., Hoang, J., Wen, Y., Langlois, E., Zhang, S., Zhang, G., Abbeel, P., and Ba, J. Benchmarking model-based reinforcement learning. *arXiv preprint arXiv:1907.02057*, 2019.
- Williams, R. J. and Peng, J. Function optimization using connectionist reinforcement learning algorithms. 1991.
- Xu, Z., van Hasselt, H., and Silver, D. Meta-gradient reinforcement learning. In *NeurIPS*, 2018.
- Yang, B., Zhang, J., Pong, V., Levine, S., and Jayaraman, D. Replab: A reproducible low-cost arm benchmark platform for robotic learning. *arXiv preprint arXiv:1905.07447*, 2019.
- Yu, T., Quillen, D., He, Z., Julian, R., Hausman, K., Finn, C., and Levine, S. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. *arXiv preprint arXiv:1910.10897*, 2019.
- Yu, T., Kumar, S., Gupta, A., Hausman, K., Levine, S., and Finn, C. Gradient surgery for multi-task learning, 2020. URL <https://openreview.net/forum?id=HJewiCVFPB>.
- Zhang, A., Wu, Y., and Pineau, J. Natural environment benchmarks for reinforcement learning. *arXiv preprint arXiv:1811.06032*, 2018.