Benchmarking Massively Parallelized Multi-Task Reinforcement Learning for Robotics Tasks

Viraj Joshi, Zifan Xu, Bo Liu, Peter Stone, Amy Zhang

Keywords: Multi-Task Learning, Reinforcement Learning, Robotics.

Summary

Multi-task Reinforcement Learning (MTRL) has emerged as a critical training paradigm for applying reinforcement learning (RL) to a set of complex real-world robotic tasks, which demands a generalizable and robust policy. At the same time, massively parallelized training has gained popularity, not only for significantly accelerating data collection through GPUaccelerated simulation but also for enabling diverse data collection across multiple tasks by simulating heterogeneous scenes in parallel. However, existing MTRL research has largely been limited to off-policy methods like SAC in the low-parallelization regime. MTRL could capitalize on the higher asymptotic performance of on-policy algorithms, whose batches require data from current policy, and as a result, take advantage of massive parallelization offered by GPU-accelerated simulation. To bridge this gap, we introduce a massively parallelized Multi-Task Benchmark for robotics (MTBench), an open-sourced benchmark featuring a broad distribution of 50 manipulation tasks and 20 locomotion tasks, implemented using the GPU-accelerated simulator IsaacGym. MTBench also includes four base RL algorithms combined with seven state-of-the-art MTRL algorithms and architectures, providing a unified framework for evaluating their performance. Our extensive experiments highlight the superior speed of evaluating MTRL approaches using MTBench, while also uncovering unique challenges that arise from combining massive parallelism with MTRL.

Contribution(s)

- This paper introduces MTBench, a unified GPU-accelerated benchmark for massively parallelized multi-task reinforcement learning (MTRL) in two robotics settings, manipulation and locomotion.
 - **Context:** Existing robotics MTRL benchmarks, such as Meta-World (Yu et al., 2021), have impractically long experimental runtimes, hindering the development and reproducibility of MTRL research. Other GPU-accelerated benchmarks for robotics do not support MTRL out of the box. We address both of these concerns with our end-to-end MTRL benchmark.
- This paper conducts comprehensive experiments to evaluate all aspects of MTRL, including base RL algorithms, gradient manipulation methods, and neural network architectures.
 Context: We confirm whether the reliance on off-policy methods in the MTRL literature holds in the massively parallel regime, and then evaluate a suite of MTRL schemes using on-policy methods across our evaluation settings.
- This paper presents four key observations on applying existing MTRL schemes to massively parallelized training in robotics. These insights guide the selection of MTRL schemes and inform future research directions.
 - **Context:** Massively parallelized training is emerging as a popular paradigm, introducing unique challenges for existing RL methods (D'Oro et al., 2022; Li et al., 2023; Gallici et al., 2024; Singla et al., 2024). However, MTRL development has yet to leverage this paradigm.

Benchmarking Massively Parallelized Multi-Task Reinforcement Learning for Robotics Tasks

Viraj Joshi^{1,†}, Zifan Xu^{1,†}, Bo Liu¹, Peter Stone^{1,2}, Amy Zhang¹

{viraj_joshi,zfxu}@utexas.edu

Abstract

Multi-task Reinforcement Learning (MTRL) has emerged as a critical training paradigm for applying reinforcement learning (RL) to a set of complex real-world robotic tasks, which demands a generalizable and robust policy. At the same time, massively parallelized training has gained popularity, not only for significantly accelerating data collection through GPU-accelerated simulation but also for enabling diverse data collection across multiple tasks by simulating heterogeneous scenes in parallel. However, existing MTRL research has largely been limited to off-policy methods like SAC in the low-parallelization regime. MTRL could capitalize on the higher asymptotic performance of on-policy algorithms, whose batches require data from the current policy, and as a result, take advantage of massive parallelization offered by GPU-accelerated simulation. To bridge this gap, we introduce a massively parallelized Multi-Task Benchmark for robotics (MTBench), an open-sourced benchmark featuring a broad distribution of 50 manipulation tasks and 20 locomotion tasks, implemented using the GPU-accelerated simulator IsaacGym. MTBench also includes four base RL algorithms combined with seven state-of-the-art MTRL algorithms and architectures, providing a unified framework for evaluating their performance. Our extensive experiments highlight the superior speed of evaluating MTRL approaches using MTBench, while also uncovering unique challenges that arise from combining massive parallelism with MTRL. Code is available at https://github.com/Viraj-Joshi/MTBench

1 Introduction

Deep reinforcement learning has been successfully applied to a wide range of decision-making tasks, including Atari games (Mnih et al., 2013), the game of Go (Silver et al., 2016), and continuous control tasks (Hwangbo et al., 2019; Wurman et al., 2022). While these applications have achieved remarkable task-specific performance, recent research trends have shifted towards developing general-purpose agents capable of solving multiple tasks or adapting to diverse environments (Cobbe et al., 2020; Kirk et al., 2023; Park et al., 2024). This transition is partly motivated by the demands of real-world robotics applications, where versatility and robustness are essential. For example, tabletop manipulation often requires acquiring multiple skills to accomplish complex tasks (Pinto & Gupta, 2016; Yu et al., 2021) and legged locomotion demands adaptability to traverse challenging terrains (Lee et al., 2020; Liang et al., 2024).

To facilitate the learning of a general-purpose robotic agent, massively parallelized training (**1000 simulations) has gained popularity with the advancement of GPU-accelerated simulators (Liang et al., 2018b; Freeman et al., 2021; Makoviychuk et al., 2021; Mittal et al., 2023; Tao et al., 2024; Zakka et al., 2025). These simulators have significantly mitigated hardware and runtime constraints for learning *single tasks*, reducing experiment durations from days to minutes (Liang et al., 2018b;

¹The University of Texas at Austin

²Sonv AI

[†] equal contribution

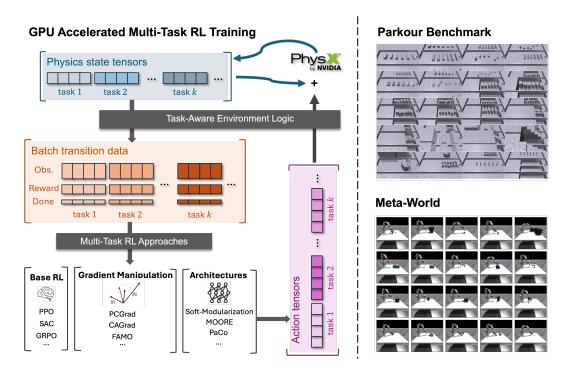


Figure 1: MTBench is a benchmark that leverages massive parallelism for MTRL in two robotics domains, Parkour and Meta-World, and provides MTRL implementations developed over the years. On the left, we see that IsaacGym's Tensor API enables us to assign blocks of environments to a desired task within the domain of interest, allowing for the setting and getting of the required information for RL training.

Rudin et al., 2022). However, in the multi-task setting, no out-of-the-box solution exists to allocate a fixed number of environments per task on a single GPU, allowing for simultaneous diverse data collection and end-to-end MTRL training. Additionally, massively parallelized online batched RL introduces new, non-trivial algorithmic challenges. For example, on-policy methods like PPO reach a saturation point beyond which additional parallelization no longer improves performance (Singla et al., 2024). Meanwhile, off-policy methods such as SAC and Q-Learning become unstable, losing their sample efficiency compared to on-policy methods as interaction with parallel environments unbalances the replay ratio (D'Oro et al., 2022; Li et al., 2023; Gallici et al., 2024).

On the other hand, learning general-purpose robotic agents has also motivated multi-task RL (MTRL), which aims to learn a single policy that maximizes average performance across multiple tasks. By leveraging task similarities (Pinto & Gupta, 2016), MTRL often enhances sample efficiency, requiring fewer transitions to match the performance of single-task counterparts. Prior research has primarily focused on addressing optimization challenges introduced by multiple learning signals, either from a gradient-based perspective (Yu et al., 2020; Liu et al., 2024; 2023a) or through neural architecture design (Yang et al., 2020; Sodhani et al., 2021; Sun et al., 2022; Hendawy et al., 2024). However, these MTRL approaches have been limited to using off-policy methods in low-parallelization settings using libraries like Ray (Liang et al., 2018a). With massive parallelization applied to MTRL, we no longer need to deal with how to distribute experience collection and learning, instead utilizing on-policy algorithms, whose batches require data from current experience and as a result, take advantage of the parallelization offered by GPU-based simulators.

To support large-scale MTRL experiments and advance the development of general-purpose robotic agents, we introduce a massively parallelized Multi-Task Benchmark for robotics (MTBench). This open-source benchmark includes a diverse set of 50 manipulation tasks and 20 locomotion tasks (right side of Figure 1), implemented using the GPU-accelerated simulator IsaacGym. Each task

allows for procedurally generating infinitely many variations by modifying factors such as initial states and terrain configurations. Additionally, MTBench integrates four base RL algorithms with seven state-of-the-art MTRL algorithms and architectures, providing a unified framework to evaluate their performance.

Based on our experiments, we highlight the following major observations:

- **(O1) On-Policy > Off-Policy:** Choosing between on-policy RL methods or off-policy methods affects performance more than the MTRL scheme applied in massively parallel training. Off-policy RL's asymptotic performance struggles to match on-policy RL in this regime.
- (O2) Prioritize Wall-Clock Time over Sample Efficiency: In the massively parallel regime, wall-clock efficiency is more critical than sample efficiency, as experience collection scales easily with more GPUs.
- (O3) Value Learning is the Key Bottleneck in MTRL: Multi-task RL struggles primarily with value estimation rather than policy learning, as gradient conflicts mostly impact the critic function.
- **(O4)** Curriculum Learning is Crucial for Sparse-Reward Tasks: MTRL alone does not help exploration in sparse-reward tasks; curriculum learning is essential for overcoming early stagnation.

2 Background

2.1 GPU Accelerated Simulation

Traditionally, simulators used for online RL rely on the coordination between CPU and GPU where the CPU handles physics simulation and observation/reward calculations while the GPU handles neural network training and inference, leading to frequent slow memory transfers between the two many times during the RL training process. Now, GPU-accelerated simulators provide access to the results of physics simulation on the GPU, and as a result, we have all relevant data - observations, actions, and rewards - remaining on the GPU throughout the learning process. This development allows for massive parallelization and as a result, dramatically reduces MTRL training time from days or weeks on thousands of CPU cores to just hours on a single GPU.

Specifically, NVIDIA IsaacGym offers a Tensor API that directly exposes the physics state of the world in Python, so we can directly populate and manage massively parallelized heterogeneous scenes for all tasks (Figure 1), avoiding the communication overhead of synchronizing experience collection and neural network training across distributed systems (Nair et al., 2015; Espeholt et al., 2018).

2.2 Multi-Task Reinforcement Learning

RL is formalized as a finite horizon, discrete-time MDP, which is represented by a tuple $\mathcal{M}=(\mathcal{S},\mathcal{A},\mathcal{P},r,\mu,\gamma)$, where $\mathcal{S}\in\mathbb{R}^n$ denotes the continuous state space, $\mathcal{A}\in\mathbb{R}^m$ denotes the continuous action space, $\mathcal{P}:\mathcal{S}\times\mathcal{A}\to\Delta(\mathcal{S})$ denotes the stochastic transition dynamics, $r:\mathcal{S}\times\mathcal{A}\to\mathbb{R}$ denotes the reward function, $\mu:\mathcal{S}\to\Delta(\mathcal{S})$ denotes the initial state distribution, and $\gamma\in[0,1)$ is the discount factor. A policy parameterized by $\theta,\pi_{\theta}(a_t|s_t):\mathcal{S}\to\Delta(\mathcal{A})$, is a probability distribution over actions conditioned on the current state. RL learns a policy π_{θ} such that it maximizes the expected cumulative discounted return $J(\theta)=\mathbb{E}_{s_0\sim\mu,\pi_{\theta}}[\sum_{t=0}^T \gamma^t r(s_t,a_t)]$ where $a_t\sim\pi_{\theta}$.

Problem statement Each task τ is sampled the task distribution $p(\mathcal{T})$ is a different MDP $\mathcal{M}^{\tau} = (\mathcal{S}^{\tau}, \mathcal{A}^{\tau}, \mathcal{P}^{\tau}, r^{\tau}, \mu^{\tau}, \gamma^{\tau})$. MTRL learns a single policy π_{θ} that maximizes the expected cumulative discounted return averaged across all tasks $J(\theta) = \sum_{\tau \in \mathcal{T}} J_{\tau}(\theta)$. The only restriction we place upon \mathcal{M}^{τ} is that their union shares a universal state space \mathcal{S} and by appending a task embedding z to the state, we give the policy the ability to distinguish what task each observation belongs to.

A change in any part of a \mathcal{M}^{τ} constitutes what it means to define a new task. In locomotion, each task from $p(\mathcal{T})$ would be associated with a different goal to reach in the same control setting, so

only r^{τ} would differ across tasks. In tabletop manipulation like Meta-World, the tasks range from basic skills like pushing and grasping to more advanced skills combining these basic skills, so the goals (r^{τ}) and state spaces (S^{τ}) vary across tasks but the action spaces A^{τ} are identical.

3 Benchmark

MTBench provides a unified framework for simulating two key robotics task categories: manipulation and locomotion, within the IsaacGym simulator. For manipulation, we incorporate 50 tasks from Meta-World (Yu et al., 2021), chosen for their simplicity, task diversity, and well-designed, shaped rewards. The locomotion domain includes 20 diverse quadrupedal Parkour tasks from Eurekaverse (Liang et al., 2024), the most comprehensive Parkour benchmark, encompassing a wide range of established locomotion challenges. As Figure 1 demonstrates, MTBench supports defining any custom subset of tasks and their associated number of environments, enabling researchers to craft different task sets of varying difficulty. This section provides a detailed overview of these task domains and the evaluation protocols.

3.1 Meta-World

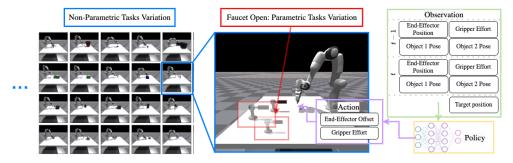


Figure 2: Illustrations of non-parametric tasks variation, parametric tasks variation of Faucet Open, and the observation and action space of the RL agents in the Meta-World benchmark.

Task Descriptions: Meta-World consists of 50 tabletop manipulation tasks that require a simulated one-armed robot (Franka Robotics, 2017) to interact with one or two objects in various ways, such as pushing, picking, and placing. Within each task, Meta-World provides parametric variation over the initial object position and target position. Each task has a pre-defined success criterion. Our re-implementation of Meta-World makes necessary changes by updating Sawyer to Franka Emika Panda and tuning the reward function of each task to ensure that the tasks are individually solvable.

Observation and Action Spaces: Despite sharing a common state space dimensionality, the semantic meaning of certain dimensions varies across tasks. The state representation comprises the end-effector's 3D position in \mathbb{R}^3 , the normalized gripper effort in \mathbb{R}^1 , the object 3D positions from two objects in \mathbb{R}^6 , and the quaternion representation of the two objects' orientation in \mathbb{R}^8 . For tasks involving a single object, the state dimensions corresponding to a second object are set to zero. To account for temporal dependencies, the observation space concatenates the state representations from two consecutive time steps and appends the 3D position of the target goal. This results in a final observation vector of 39 dimensions. The action space is also consistent across the tasks, comprising of the displacement of the end-effector in \mathbb{R}^3 and the normalized gripper effort in \mathbb{R}^1 . An overview of the observation and action can be seen in Figure 2.

Evaluation Settings: Following Yu et al. (2021), we explictly provide two evaluation settings: multi-task 10 (MT10) and multi-task 50 (MT50), where MT10 consists of 10 selected tasks and MT50 consists of all 50 tasks. During evaluation, we measure the *success rates* (SR) (Appendix

B.1) and the *cumulative reward* (R). When each *environment* has its parametric parameters randomly varied every reset, the evaluation is referred to as MT10-rand and MT50-rand.

3.2 Parkour Benchmark

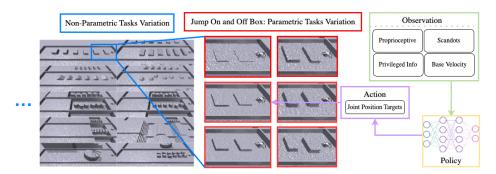


Figure 3: Illustrations of non-parametric tasks variation, parametric tasks variation of Jump On and Off Box, and the observation and action space of the RL agents in the Parkour benchmark.

Task Descriptions: In the Parkour tasks from Eurekaverse, the agent controls a quadrupedal Unitree Go1 robot (Unitree Robotics, 2021) to track predefined waypoints while traversing one of 20 different terrain categories (Liang et al., 2024). These tasks challenge various motor skills, including climbing boxes, walking on slopes, jumping, navigating stepping stones, ascending stairs, maneuvering through narrow hallways, weaving through agility poles in a zig-zag pattern, and maintaining balance. The left side of Figure 3 shows bird-eye views of these terrain categories.

Each task also provides parametric terrain variations defined by a set of terrain parameters, whose definitions and valid ranges are detailed in the supplementary materials. Additionally, each task introduces a one-dimensional continuous variable, termed *difficulty*, and a predefined mapping from the *difficulty* to a set of terrain parameters. This *difficulty* measure aligns with human intuition; for instance, high boxes present a greater challenge than lower boxes for a quadrupedal robot to jump on and off.

Observation and Action Spaces: The observation of the agent is slightly simplified for more efficient benchmarking compared to Eurekaverse. The observation is compromised by proprioceptive observation in \mathbb{R}^{48} , scandots of the terrain environments in \mathbb{R}^{132} , base linear velocity in \mathbb{R}^3 , and privileged information in \mathbb{R}^{29} . The action assigns joint position targets at a frequency of 50 Hz for a Proportional-Derivative (PD) controller. An overview of the observation and actions is shown in Figure 3. The reward function resembles Fu et al. (2023), which encourages positive linear and angular velocities that point to the next waypoint, while minimizing energy consumption.

Evaluation Settings: We define two evaluation settings: Parkour-easy and Parkour-hard. Parkour-easy consists of 200 terrains, with each of the 20 tasks assigned 10 terrains generated at the lowest difficulty level. In contrast, Parkour-hard also includes 200 terrains but distributes difficulty levels uniformly across the 10 terrains per task, providing a more diverse and challenging evaluation setting. Before the training, all the evaluated methods are pre-trained on flat ground to acquire the basic walking gait. Such a pre-training phase is typical in the literature (Zhuang et al., 2023; Cheng et al., 2024).

During evaluation, we measure *progress* (P) as the ratio of the current waypoint index to the total number of waypoints at the time of episode termination. An agent that successfully traverses the entire terrain achieves a *progress* score of 100%. The overall *progress* is computed as the average over 200 terrains, with each terrain evaluated across 10 independent runs.

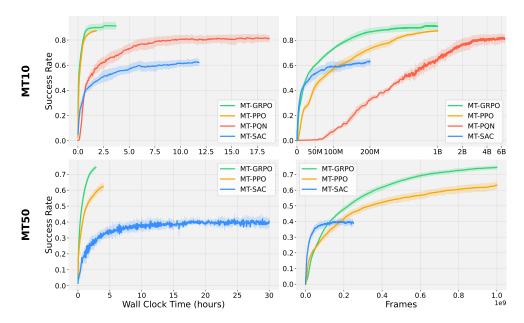


Figure 4: **Vanilla MTRL performance in Meta-World.** We report the pointwise 95% percentile bootstrap CIs of the average success rates using 10 seeds for each RL algorithm in the MT10-rand and MT50-rand evaluation settings. On-policy methods (MT-PPO, MT-GRPO) continue to improve with more experience, achieving a substantially higher success rate than the traditional off-policy method, MT-SAC, in substantially less time.

3.3 Algorithms

We re-implement a suite of algorithms and MTRL approaches using a popular learning library RL-Games (Makoviichuk & Makoviychuk, 2021), providing a unified benchmark for end-to-end vectorized MTRL training across many seeds and hyperparameters on a single GPU. Our benchmark is highly extensible towards new RL algorithms as well as approaches within the two axes of MTRL research, gradient manipulation, and neural architectures. There is a brief overview in Appendix C.

Base MTRL Algorithms We implement four RL algorithms: MT-PPO, a multi-task version of Proximal Policy Optimization (Schulman et al., 2017); MT-GRPO (Shao et al., 2024), a variant of PPO introduced for language modeling but adapted here for control; MT-SAC, a multi-task version of Soft Actor-Critic (Haarnoja et al., 2018); and MT-PQN, a novel multi-task extension to Parallel Q-learning (Gallici et al., 2024) to handle continuous control problems. All algorithms are multi-task versions of their single-task counterparts, simply by augmenting the observation space with one-hot task embeddings.

MTRL Schemes We implement two categories of MTRL schemes that can be easily combined with any of our base algorithms. The first category consists of gradient manipulation methods: PCGrad (Yu et al., 2020), CAGrad (Liu et al., 2024), and FAMO (Liu et al., 2023a). The second category consists of multi-task architectures: Soft-Modularization (Yang et al., 2020), CARE (Sodhani et al., 2021), PaCo (Sun et al., 2022), and MOORE (Hendawy et al., 2024). The prefix "MH" (multi-head) is prepended to name of the MTRL approach to denote one output head per task, and otherwise "SH" (single-head) to denote tasks sharing one head.

Curriculum Learning Unlike Meta-World, where reward functions are carefully designed with dense rewards, locomotion tasks often rely on sparse reward signals (e.g., moving forward to the next waypoints). As a result, strategies like curriculum learning have been widely adopted to facilitate learning in challenging tasks, such as running (Margolis et al., 2024) and jumping onto high

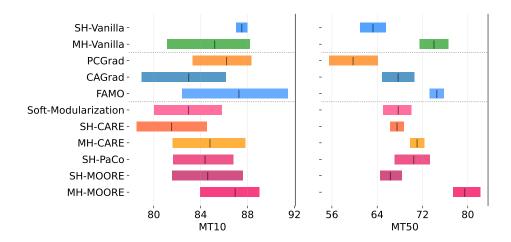


Figure 5: We compare the 95% bootstrapped confidence intervals of the average success rate of all MTRL approaches using MT-PPO for the MT10-rand and MT50-rand evaluation settings of Meta-World. Each approach uses 1B frames per run over 10 seeds. Exact numbers are in Table 2.

platforms (Liang et al., 2024). Inspired by this, we incorporate a simple curriculum strategy to train Parkour-hard tasks. In Parkour-hard, each task consists of ten terrains with varying levels of *difficulty*. Agents always begin on the easiest terrain and progress to more challenging ones if they achieve a *progress* of at least 80% in their current terrain. We refer to the Parkour-hard training with curriculum learning by Parkour-hard-cl.

4 Results

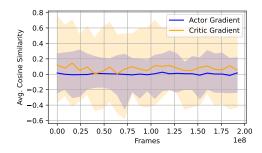
In this section, we present the results of our benchmark across our evaluation settings and empirically justify the aforementioned four major observations.

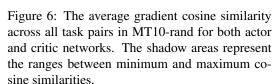
4.1 Choosing the MTRL Algorithm (O1, O2)

To illustrate how on-policy MTRL methods leverage massive parallelism, we first evaluate two on-policy methods, MT-PPO and MT-GRPO, alongside two off-policy methods, MT-SAC and MT-PQN, in Meta-World. Figure 4 presents the learning curves with respect to both wall-clock time and the number of environment interactions. Since this observation is concerned with answering what the best base MTRL algorithm is, we tune all aspects of each method to achieve its highest success rate, including using different network architectures. The full hyperparameter and model details are in the supplementary materials.

On-policy methods outperform traditional off-policy methods. Using MT-SAC as a representative of traditional off-policy algorithms used for MTRL, Figure 4 shows there is a substantial performance gap in success rates between MT-PPO and MT-SAC in both evaluation settings and, more importantly, a substantial wall-clock time difference as well (roughly 22 minutes and 12 hours after 200M frames of collected experience in MT10-rand). While MT-SAC can match MT-PPO's runtime by simply matching the gradient steps per epoch that MT-PPO takes, this results in a near-zero success rate. Furthermore, as the number of tasks increases to the MT50-rand setting, these gaps increase.

Traditional off-policy methods in the massively parallelized regime cannot effectively leverage increased environment interaction, as their stability, performance, and runtime greatly rely on the ratio of gradient updates to environment steps, i.e, update-to-data (UTD) ratio (D'Oro et al., 2022) be-





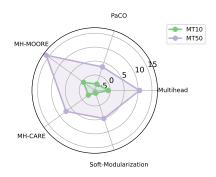


Figure 7: Success rate (SR) differences of five neural network architectures relative to the Vanilla baseline in MT10-rand and MT50-rand.

ing greater than or equal to 1. Research into leveraging massive parallelism in off-policy methods is gaining popularity but is either not yet adapted for the continuous control setting (Gallici et al., 2024) or requires distributed asynchronous processes spread across GPUs (Li et al., 2023).

Off-Policy methods can be designed for the massively parallelized regime. In Figure 4, we also included our adaptation of PQN (Gallici et al., 2024) to the multi-task continuous control setting. The details of our implementation are in Appendix A. Surprisingly, applying these simple changes to an originally discrete control algorithm and left to run long enough, MT-PQN can roughly match the performance of MT-PPO in MT10-rand. Considering PQN's performance and stability, similar simulation throughput to PPO, and lack of a replay buffer, suggests that smartly adapting PQN to multi-task continuous control tasks could be a promising research direction compared to actor-critic algorithms.

4.2 MTRL Approaches (O2, O3)

Figure 5 reports the 95% bootstrap confidence intervals of the mean success rate following Agarwal et al. (2021) of all MTRL approaches using MT-PPO. All of the gradient manipulation methods use the same three-layer MLP neural networks.

Multi-task architectures show greater performance gains with larger task sets. As shown in Figure 7, the benefits of multi-task architectures become more pronounced as the number of tasks increases. In MT10-rand, vanilla PPO asymptotically outperforms advanced multi-task architectures. However, in MT50-rand, the best-performing multi-task architecture, MH-MOORE, surpasses the vanilla approach by roughly 16% in success rate. This improvement is likely due to enhanced knowledge sharing that only manifests in training diverse enough tasks, such as MT50. However, similar performance gains are not observed in the Parkour benchmark, likely due to the insufficient task diversity in Parkour tasks.

Resolving gradient conflict consistently improves the performance. Gradient manipulation can outperform or match vanilla MT-PPO across all evaluation settings (middle section of Figure 5 and Table 1). This suggests that gradient conflicts are still a common optimization challenge in multitask RL problems. Among these methods, FAMO shows superior scalability with respect to an increasing number of tasks in its success rate as well as wall-clock training time, likely due to its simple strategy of adaptive task weighting, which eliminates the need for backpropagating through each task's loss.

Tasi	ks Parkour-easy	Parkour-hard	Parkour-hard-CL
Methods	P(%)↑	P(%) ↑	P(%) ↑
Vanilla	80.39 ± 0.43	54.51 ± 0.82	68.12 ± 0.43
Multihead	73.17 ± 2.53	49.65 ± 1.27	62.17 ± 1.05
PCGrad	$\textbf{80.61} \pm \textbf{0.78}$	$\textbf{55.98} \pm \textbf{0.24}$	68.37 ± 0.65
CAGrad	80.25 ± 0.32	55.88 ± 0.91	67.75 ± 0.43
FAMO	79.79 ± 0.29	55.56 ± 0.97	$\textbf{68.57} \pm \textbf{0.76}$
PaCo	78.63 ± 0.70	$\textbf{58.65} \pm \textbf{0.77}$	64.15 ± 0.98
SH-MOORE	64.61 ± 1.69	46.78 ± 0.65	49.53 ± 0.52
Soft-Modularizatio	n 69.29 ± 3.81	47.28 ± 0.23	51.43 ± 0.35

Table 1: The average success rate and standard deviation of MTRL approaches using MT-PPO in all Parkour evaluation settings. Each approach uses 250M frames per run over 10 seeds.

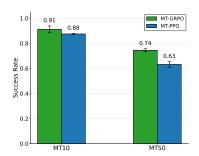


Figure 8: Eliminating the difficulty of critic estimation consistently improves performance over most MTRL approaches using MT-PPO when comparing the 95% bootstrapped CI of average success rates in Meta-World. Each approach uses 1B frames per run over 10 seeds.

Value learning is the key bottleneck in MTRL. Prior research in MTRL has shown that addressing gradient conflicts improves performance in off-policy actor-critic RL algorithms like MT-SAC. Our benchmarking results extend this observation to on-policy actor-critic algorithms, demonstrating that gradient conflicts also arise when learning the critic network in MT-PPO. However, we do not observe similar conflicts in policy optimization. This observation aligns with prior work using an actor-critic algorithm for large-scale multi-task learning (Hessel et al., 2019). Figure 6 shows the average cosine similarity across all task gradient pairs for both actor and critic networks, where critic gradients manifest lower minimum similarities.

4.3 Reward Sparsity (O4)

Although tasks in Meta-World and the Parkour Benchmark are defined independently of their reward functions, training performance is significantly influenced by reward design. We adopt commonly used reward formulations in both domains. In Meta-World, tasks utilize dense rewards, which provide continuous feedback to guide specific interactions between the robotic arm and objects. In contrast, the Parkour Benchmark employs a sparse reward scheme, where the agent is rewarded solely for maintaining forward velocity toward waypoints, without receiving additional signals for intermediate behaviors.

Dense rewards increase the complexity of multi-task critic learning. In multi-task RL, dense reward functions introduce challenges for critic learning, as different tasks exhibit varying reward distributions and gradient magnitudes. We can see in Figure 5 that addressing these conflicts in dense-reward multi-task settings such as MT10-rand and MT50-rand can improve performance. However, the performance gains are relatively marginal in a sparse-reward multi-task setting like the Parkour benchmark.

Curriculum learning is crucial for sparse-reward tasks. In environments with sparse rewards, standard MTRL methods do not inherently enhance exploration, as agents receive limited feedback in each task. This challenge is particularly evident in the Parkour Benchmark, where agents tend to adopt overly conservative behaviors in more difficult tasks. Curriculum learning addresses this issue by structuring task progression, enabling agents to first master simpler behaviors before tackling more complex ones. By gradually increasing task difficulty, curriculum learning improves exploration efficiency and yields a 10% performance gain in *progress*, as observed when comparing Parkour-hard and Parkour-hard-cl (columns 2 and 3 in Table 1).

4.4 Learning without a Critic (O3)

To further investigate the impact of gradient conflict in the critic on MTRL, we can eliminate the critic by increasing the horizon length in MT-PPO to be equal to the length of the episode.

MTRL can benefit from eliminating gradient conflict in the critic. In fact, eliminating the critic from MT-PPO is equivalent to implementing MT-GRPO (Shao et al., 2024) without the KL term. We use the Monte Carlo estimate of the return as the reward in the advantage calculation. In the dense reward setting, Figure 8 indicates that MT-GRPO is a simple baseline that nearly outperforms every MTRL approach (except MH-MOORE and FAMO in MT-50) using the same hyperparameters as MT-PPO and no baked-in MTRL design.

Massive Parallelism is well suited for reducing bias from an imperfect critic. By directly using Monte Carlo returns instead of bootstrapping, we effectively eliminate the bias introduced by imperfect critic estimation. This approach represents a clear bias-variance tradeoff: while removing the critic increases the variance of our gradient estimates, this increased variance can be effectively mitigated through large batches (of size episode length times the number of parallel environments) made possible by massive parallelization (Sutton et al., 1999).

5 Related

5.1 Parallelizing RL

As deep online RL relies on training neural networks (learners) and collecting experience (actors), many methods have explored how to parallelize both aspects to speed up training over the years. Early works leveraged low levels of parallelization without hardware accelerators mainly for Atari either in a distributed compute cluster of hundreds (in some cases thousands) of CPU cores (Nair et al., 2015) or a single machine using a multi-threaded approach (Mnih et al., 2016). Hybrid CPU-GPU distributed frameworks introduce accelerating learners with GPUs (Babaeizadeh et al., 2016; Espeholt et al., 2018; Horgan et al., 2018; Petrenko et al., 2020) along with actors collecting experience across CPUs.

Unlike Atari, robotic control tasks rely on physics simulators (Todorov et al., 2012), where distributed RL methods using CPU-based simulators would demand even more intense hardware requirements (Liang et al., 2018a; OpenAI et al., 2019) due to running multiple simulator instances in parallel. A wave of recent GPU-accelerated simulators (Liang et al., 2018b; Freeman et al., 2021; Makoviychuk et al., 2021; Mittal et al., 2023; Tao et al., 2024; Zakka et al., 2025) has essentially alleviated the experience collection constraint and shown success in rapidly learning single-task robotic control tasks (Allshire et al., 2021; Rudin et al., 2022) with the modest hardware requirement of 1 GPU.

5.2 GPU-Accelerated Benchmarks

Several RL benchmarks have arisen as a result of GPU-accelerated simulation, mainly in JAX-based game environments (Cobbe et al., 2020; Lange, 2022; Morad et al., 2023; Bonnet et al., 2023; Koyamada et al., 2023; Rutherford et al., 2024; Matthews et al., 2024). In contrast, a relatively small number of rigid-body robotic tasks are bundled with GPU-accelerated simulators or soft-body robotic tasks with other simulation platforms (Chen et al., 2022; Xing et al., 2024). To truly represent the multi-task challenge, MTBench precludes adapting popular, small task sets, e.g, robot tasks from DMControl (Tassa et al., 2018) or robosuite (Zhu et al., 2020), or combining them since their tasks significantly overlap, resulting in low diversity. For manipulation, Meta-World resolves both of these concerns and maintains continuity of MTRL research over other large task set alternatives like RLBench (James et al., 2020) or LIBERO (Liu et al., 2023b).

In the domain of locomotion, massively parallelized training has become the standard approach due to its simplicity and increased robustness (Hwangbo et al., 2019; Lee et al., 2020). Nevertheless, parkour-style locomotion—which requires qualitatively different motor skills across terrains—remains a challenging setting for multi-task learning, which recent work addresses by learning specialized policies for individual motor skills and subsequently distilling them into a unified policy (Zhuang et al., 2023). Emerging locomotion benchmarks such as HumanoidBench (Sferrazza et al., 2024) and the Parkour Benchmark (Liang et al., 2024) feature a broad variety of Parkour tasks, but have not yet been adopted for evaluating multi-task RL methods. Another important class of locomotion tasks involves humanoid motion imitation, such as those in PHC (Luo et al., 2023) and LocoMuJoCo (Al-Hafez et al., 2023), which exhibit inherently diverse task distributions due to the complex and high-dimensional nature of human motion.

6 Conclusions

We present MTBench, a highly extensible MTRL benchmark that includes a GPU-accelerated implementation of Meta-World and Parkour tasks, extensive gradient manipulation and neural architecture baselines, and an initial study on the current state as well as future directions of MTRL in the massively parallel regime. However, MTBench is limited to state-based MTRL to retain high simulation throughput, which we hope to resolve with pixel-based MTRL using NVIDIA IsaacLab in a future release of MTBench.

Future work can use MTBench beyond online MTRL methods. One can explore offline RL, imitation learning, or distillation methods by writing additional code to rapidly collect transitions from expert single-task agents. Another application of our benchmark could be as part of the 'finetune' step in the 'pretrain, then finetune' paradigm where one pre-trains on a diverse set of tasks using offline data and rapidly finetunes an agent online using our environments.

A PQN

Parallel Q-learning (Gallici et al., 2024) is a recent off-policy TD method designed for discrete action spaces and massively parallelized GPU-based simulators that casts aside the tricks introduced over the years to stabilize deep Q learning such as replay buffers (Mnih et al., 2013), target networks (Mnih et al., 2015) and double Q-networks (Wang et al., 2016) by simply introducing regularization in the function approximator like LayerNorm (Ba et al., 2016) or BatchNorm (Ioffe & Szegedy, 2015). Coupled with this architectural change, PQN exploits vectorized environments by collecting experience in parallel for T steps.

As our action space is continuous, we modify PQN through bang-off-bang control, treating continuous control as a multi-agent problem where each of the M actuators is an agent in a cooperative game following Seyde et al. (2023). Then, the state-action function $Q_{\theta}(\mathbf{s_t}, \mathbf{a_t})$ is factorized as the average of M different state-action functions $Q_{\theta}^i(\mathbf{s_t}, a_t^i)$, where the ith state-action function predicts the value of the bang-off-bang actions in ith action dimension following Sunehag et al. (2017).

$$Q_{\theta}(\mathbf{s_t}, \mathbf{a_t})) = \frac{1}{M} \sum_{i=1}^{M} Q_{\theta}^i(\mathbf{s_t}, a_t^i)$$
(1)

In code for the Meta-World setting, the output of the state-action function is of size (B, M, n_b) where B is the batch size, m is the action dimension/number of actuators (4) and n_b is the number of bins per dimension (3). The action value is recovered by first taking the max over the bin dimension and then the mean over the action dimension. By taking the max over the bin dimension, Seyde et al. (2023) sidestepped taking a max over the continuous action space. Now, we can compute the Bellman target and in the case of PQN, n-step returns.

$$y_t = r(\mathbf{s_t}, \mathbf{a_t}) + \gamma \frac{1}{M} \sum_{i=1}^{M} \max_{a_{t+1}^i} Q_{\theta}^i(\mathbf{s_{t+1}}, a_{t+1}^i)$$
 (2)

B Meta-World

B.1 Success

We report two evaluation metrics, the overall success rate averaged across tasks and the cumulative reward achieved by the multi-task policy. Following the original Meta-World, success is a boolean indicating whether the robot brings the object within an ϵ distance of the goal position at *any* point during the episode, which is less restrictive than works qualifying a success only if it occurs at the *end* of an episode. Mathematically, success occurs if $\|o - g\|_2 < \epsilon$ is satisfied at least once, where o is the object position and g is the goal position.

Rather than defining the success rate as the maximum success rate over some evaluation rollouts as some previous work did, the success rate is defined as the proportion of success in the large number of environments that terminate their episodes every step. The reported success rate is this success rate averaged over the last 5 epochs of training. Due to massive parallelization, there is no need to separately roll out the learned policy in a separate process.

B.2 More Results

Tasks		MT10-rand		MT50-rand	
Methods		SR ↑	R ↑	SR ↑	R ↑
Vanilla Multihead GRPO-Vanilla	85.19	[86.99, 87.97] [81.42, 88.21] [88.32, 93.96]	1032.99 [1016.82, 1045.94] 1005.69 [980.87, 1027.55] 916.32 [899.83, 933.60]	63.26 [60.91, 65.37] 74.03 [71.47, 76.58] 74.48 [73.31, 75.64]	817.77 [789.13, 842.97] 954.97 [939.72, 962.84] 916.83 [898.69, 935.66]
PCGrad CAGrad FAMO	82.98	[83.19, 88.32] [79.23, 86.27] [82.53, 91.57]	1038.27 [1022.88, 1050.59] 938.43 [896.83, 972.29] 1016.11 [964.30, 1053.88]	59.74 [55.52, 64.12] 67.70 [64.76, 70.53] 74.52 [73.25, 75.75]	760.99 [739.05, 772.13] 874.45 [845.62, 903.10] 961.03 [946.64, 976.15]
PaCo SH-MOORE MH-MOORE SH-CARE MH-CARE Soft-Modulariza	84.60 86.94 81.51 84.79	[81.61, 86.61] [81.55, 87.59] [83.91, 89.01] [78.52, 84.49] [81.34, 87.32] [80.15, 85.66]	995.39 [970.20, 1017.21] 1022.64 [1006.23, 1037.54] 1044.85 [1029.76, 1056.96] 964.28 [948.59, 979.89] 990.03 [972.35, 1006.34] 994.29 [980.24, 1009.03]	70.46 [67.01, 73.32] 66.33 [64.56, 68.29] 79.46 [77.40, 82.24] 67.51 [66.33, 68.72] 71.05 [69.88, 72.30] 67.72 [65.06, 69.93]	917.84 [881.61, 953.14] 837.70 [815.00, 860.89] 1019.59 [999.24, 1048.88] 842.04 [822.31, 864.66] 863.88 [850.43, 878.51] 860.41 [832.44, 883.77]

Table 2: 95% bootstrapped confidence intervals of the Meta-World evaluation metrics used to generate Figure 5 and Figure 8

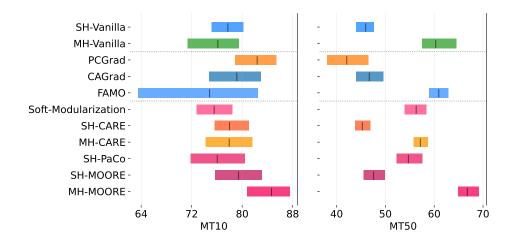


Figure 9: 95% bootstrapped CIs of the average success rate of all MT-PPO MTRL approaches using 250M frames per run over 10 seeds in Meta-World.

Acknowledgments

This work has taken place in the Learning Agents Research Group (LARG) at the Artificial Intelligence Laboratory, The University of Texas at Austin. LARG research is supported in part by the National Science Foundation (FAIN-2019844, NRT-2125858), the Office of Naval Research (N00014-18-2243), Army Research Office (W911NF-23-2-0004, W911NF-17-2-0181), Lockheed Martin, and Good Systems, a research grand challenge at the University of Texas at Austin. The views and conclusions contained in this document are those of the authors alone. Peter Stone serves as the Executive Director of Sony AI America and receives financial compensation for this work. The terms of this arrangement have been reviewed and approved by the University of Texas at Austin in accordance with its policy on objectivity in research.

References

- Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320, 2021.
- Firas Al-Hafez, Guoping Zhao, Jan Peters, and Davide Tateo. Locomujoco: A comprehensive imitation learning benchmark for locomotion. In 6th Robot Learning Workshop, NeurIPS, 2023.
- Arthur Allshire, Mayank Mittal, Varun Lodaya, Viktor Makoviychuk, Denys Makoviichuk, Felix Widmaier, Manuel Wüthrich, Stefan Bauer, Ankur Handa, and Animesh Garg. Transferring dexterous manipulation from gpu simulation to a remote real-world trifinger. *arXiv* preprint arXiv:2108.09779, 2021.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL https://arxiv.org/abs/1607.06450.
- Mohammad Babaeizadeh, Iuri Frosio, Stephen Tyree, Jason Clemons, and Jan Kautz. Reinforcement learning through asynchronous advantage actor-critic on a gpu. *arXiv preprint arXiv:1611.06256*, 2016.
- Clément Bonnet, Daniel Luo, Donal Byrne, Shikha Surana, Sasha Abramowitz, Paul Duckworth, Vincent Coyette, Laurence I Midgley, Elshadai Tegegn, Tristan Kalloniatis, et al. Jumanji: a diverse suite of scalable reinforcement learning environments in jax. arXiv preprint arXiv:2306.09884, 2023.
- Siwei Chen, Yiqing Xu, Cunjun Yu, Linfeng Li, Xiao Ma, Zhongwen Xu, and David Hsu. Daxbench: Benchmarking deformable object manipulation with differentiable physics. *arXiv* preprint arXiv:2210.13066, 2022.
- Xuxin Cheng, Kexin Shi, Ananye Agarwal, and Deepak Pathak. Extreme parkour with legged robots. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pp. 11443–11450. IEEE, 2024.
- Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In *International conference on machine learning*, pp. 2048– 2056. PMLR, 2020.
- Pierluca D'Oro, Max Schwarzer, Evgenii Nikishin, Pierre-Luc Bacon, Marc G Bellemare, and Aaron Courville. Sample-efficient reinforcement learning by breaking the replay ratio barrier. In *Deep Reinforcement Learning Workshop NeurIPS* 2022, 2022.
- Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *International conference on machine learning*, pp. 1407–1416. PMLR, 2018.

- Franka Robotics. Franka emika panda robot, 2017. URL https://www.franka.de. Accessed: 2025-02-17.
- C Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. Brax—a differentiable physics engine for large scale rigid body simulation. *arXiv* preprint arXiv:2106.13281, 2021.
- Zipeng Fu, Xuxin Cheng, and Deepak Pathak. Deep whole-body control: learning a unified policy for manipulation and locomotion. In *Conference on Robot Learning*, pp. 138–149. PMLR, 2023.
- Matteo Gallici, Mattie Fellows, Benjamin Ellis, Bartomeu Pou, Ivan Masmitja, Jakob Nicolaus Foerster, and Mario Martin. Simplifying deep temporal difference learning, 2024. URL https://arxiv.org/abs/2407.04811.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor, 2018. URL https://arxiv.org/abs/1801.01290.
- Ahmed Hendawy, Jan Peters, and Carlo D'Eramo. Multi-task reinforcement learning with mixture of orthogonal experts, 2024. URL https://arxiv.org/abs/2311.11385.
- Matteo Hessel, Hubert Soyer, Lasse Espeholt, Wojciech Czarnecki, Simon Schmitt, and Hado Van Hasselt. Multi-task deep reinforcement learning with popart. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3796–3803, 2019.
- Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado van Hasselt, and David Silver. Distributed prioritized experience replay, 2018. URL https://arxiv.org/abs/1803.00933.
- Jemin Hwangbo, Joonho Lee, Alexey Dosovitskiy, Dario Bellicoso, Vassilios Tsounis, Vladlen Koltun, and Marco Hutter. Learning agile and dynamic motor skills for legged robots. *Science Robotics*, 4(26):eaau5872, 2019.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift, 2015. URL https://arxiv.org/abs/1502.03167.
- Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 2020.
- Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of zero-shot generalisation in deep reinforcement learning. *Journal of Artificial Intelligence Research*, 76:201–264, 2023.
- Sotetsu Koyamada, Shinri Okano, Soichiro Nishimori, Yu Murata, Keigo Habara, Haruka Kita, and Shin Ishii. Pgx: Hardware-accelerated parallel game simulators for reinforcement learning. *Advances in Neural Information Processing Systems*, 36:45716–45743, 2023.
- Robert Tjarko Lange. gymnax: A JAX-based reinforcement learning environment library, 2022. URL http://github.com/RobertTLange/gymnax.
- Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learning quadrupedal locomotion over challenging terrain. *Science robotics*, 5(47):eabc5986, 2020.
- Zechu Li, Tao Chen, Zhang-Wei Hong, Anurag Ajay, and Pulkit Agrawal. Parallel *q*-learning: Scaling off-policy reinforcement learning under massively parallel simulation. In *International Conference on Machine Learning*. PMLR, 2023.
- Eric Liang, Richard Liaw, Robert Nishihara, Philipp Moritz, Roy Fox, Ken Goldberg, Joseph Gonzalez, Michael Jordan, and Ion Stoica. Rllib: Abstractions for distributed reinforcement learning. In *International conference on machine learning*, pp. 3053–3062. PMLR, 2018a.

- Jacky Liang, Viktor Makoviychuk, Ankur Handa, Nuttapong Chentanez, Miles Macklin, and Dieter Fox. Gpu-accelerated robotic simulation for distributed reinforcement learning. In *Conference on Robot Learning*, pp. 270–282. PMLR, 2018b.
- William Liang, Sam Wang, Hung-Ju Wang, Osbert Bastani, Dinesh Jayaraman, and Yecheng Jason Ma. Eurekaverse: Environment curriculum generation via large language models, 2024. URL https://arxiv.org/abs/2411.01775.
- Bo Liu, Yihao Feng, Peter Stone, and Qiang Liu. Famo: Fast adaptive multitask optimization, 2023a. URL https://arxiv.org/abs/2306.03792.
- Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning, 2023b. URL https://arxiv.org/abs/2306.03310.
- Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning, 2024. URL https://arxiv.org/abs/2110.14048.
- Zhengyi Luo, Jinkun Cao, Kris Kitani, Weipeng Xu, et al. Perpetual humanoid control for real-time simulated avatars. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10895–10904, 2023.
- Denys Makoviichuk and Viktor Makoviychuk. rl-games: A high-performance framework for reinforcement learning. https://github.com/Denys88/rl_games, May 2021.
- Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High performance gpu-based physics simulation for robot learning, 2021. URL https://arxiv.org/abs/2108.10470.
- Gabriel B Margolis, Ge Yang, Kartik Paigwar, Tao Chen, and Pulkit Agrawal. Rapid locomotion via reinforcement learning. *The International Journal of Robotics Research*, 43(4):572–587, 2024.
- Michael Matthews, Michael Beukman, Benjamin Ellis, Mikayel Samvelyan, Matthew Jackson, Samuel Coward, and Jakob Foerster. Craftax: A lightning-fast benchmark for open-ended reinforcement learning, 2024. URL https://arxiv.org/abs/2402.16801.
- Mayank Mittal, Calvin Yu, Qinxi Yu, Jingzhou Liu, Nikita Rudin, David Hoeller, Jia Lin Yuan, Ritvik Singh, Yunrong Guo, Hammad Mazhar, Ajay Mandlekar, Buck Babich, Gavriel State, Marco Hutter, and Animesh Garg. Orbit: A unified simulation framework for interactive robot learning environments. *IEEE Robotics and Automation Letters*, 8(6):3740–3747, 2023. DOI: 10.1109/LRA.2023.3270034.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013. URL https://arxiv.org/abs/1312.5602.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas Kirkeby Fidjeland, Georg Ostrovski, Stig Petersen, Charlie Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518:529–533, 2015. URL https://api.semanticscholar.org/CorpusID:205242740.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937. PmLR, 2016.

- Steven Morad, Ryan Kortvelesy, Matteo Bettini, Stephan Liwicki, and Amanda Prorok. Popgym: Benchmarking partially observable reinforcement learning. *arXiv preprint arXiv:2303.01859*, 2023.
- Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, et al. Massively parallel methods for deep reinforcement learning. *arXiv* preprint arXiv:1507.04296, 2015.
- OpenAI, Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, Jonas Schneider, Szymon Sidor, Josh Tobin, Peter Welinder, Lilian Weng, and Wojciech Zaremba. Learning dexterous in-hand manipulation, 2019. URL https://arxiv.org/abs/1808.00177.
- Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking offline goal-conditioned rl. *arXiv preprint arXiv:2410.20092*, 2024.
- Aleksei Petrenko, Zhehui Huang, Tushar Kumar, Gaurav Sukhatme, and Vladlen Koltun. Sample factory: Egocentric 3d control from pixels at 100000 fps with asynchronous reinforcement learning. In *International Conference on Machine Learning*, pp. 7652–7662. PMLR, 2020.
- Lerrel Pinto and Abhinav Gupta. Learning to push by grasping: Using multiple tasks for effective learning, 2016. URL https://arxiv.org/abs/1609.09025.
- Nikita Rudin, David Hoeller, Philipp Reist, and Marco Hutter. Learning to walk in minutes using massively parallel deep reinforcement learning. In *Conference on Robot Learning*, pp. 91–100. PMLR, 2022.
- Alexander Rutherford, Benjamin Ellis, Matteo Gallici, Jonathan Cook, Andrei Lupu, Gardar Ingvarsson, Timon Willi, Ravi Hammond, Akbir Khan, Christian Schroeder de Witt, Alexandra Souly, Saptarashmi Bandyopadhyay, Mikayel Samvelyan, Minqi Jiang, Robert Tjarko Lange, Shimon Whiteson, Bruno Lacerda, Nick Hawes, Tim Rocktaschel, Chris Lu, and Jakob Nicolaus Foerster. Jaxmarl: Multi-agent rl environments and algorithms in jax, 2024. URL https://arxiv.org/abs/2311.10090.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
- Tim Seyde, Peter Werner, Wilko Schwarting, Igor Gilitschenski, Martin Riedmiller, Daniela Rus, and Markus Wulfmeier. Solving continuous control via q-learning, 2023. URL https://arxiv.org/abs/2210.12566.
- Carmelo Sferrazza, Dun-Ming Huang, Xingyu Lin, Youngwoon Lee, and Pieter Abbeel. Humanoid-bench: Simulated humanoid benchmark for whole-body locomotion and manipulation. *arXiv* preprint arXiv:2403.10506, 2024.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/2402.03300.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Jayesh Singla, Ananye Agarwal, and Deepak Pathak. Sapg: Split and aggregate policy gradients. In *Proceedings of the 41st International Conference on Machine Learning (ICML 2024)*, Proceedings of Machine Learning Research, Vienna, Austria, July 2024. PMLR.

- Shagun Sodhani and Amy Zhang. Mtrl multi task rl algorithms. Github, 2021. URL https://github.com/facebookresearch/mtrl.
- Shagun Sodhani, Amy Zhang, and Joelle Pineau. Multi-task reinforcement learning with context-based representations. In *International Conference on Machine Learning*, 2021. URL https://api.semanticscholar.org/CorpusID:231879645.
- Lingfeng Sun, Haichao Zhang, Wei Xu, and Masayoshi Tomizuka. Paco: Parameter-compositional multi-task reinforcement learning. *ArXiv*, abs/2210.11653, 2022. URL https://api.semanticscholar.org/CorpusID:253080666.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-decomposition networks for cooperative multi-agent learning, 2017. URL https://arxiv.org/abs/1706.05296.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. Advances in neural information processing systems, 12, 1999.
- Stone Tao, Fanbo Xiang, Arth Shukla, Yuzhe Qin, Xander Hinrichsen, Xiaodi Yuan, Chen Bao, Xinsong Lin, Yulin Liu, Tse kai Chan, Yuan Gao, Xuanlin Li, Tongzhou Mu, Nan Xiao, Arnav Gurha, Zhiao Huang, Roberto Calandra, Rui Chen, Shan Luo, and Hao Su. Maniskill3: Gpu parallelized robotics simulation and rendering for generalizable embodied ai. arXiv preprint arXiv:2410.00425, 2024.
- Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Riedmiller. Deepmind control suite, 2018. URL https://arxiv.org/abs/1801.00690.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 5026–5033, 2012. DOI: 10.1109/IROS.2012.6386109.
- Unitree Robotics. *Go1 User Manual*. Unitree Robotics, 2021. Available at https://www.unitree.com/go1.
- Ziyu Wang, Tom Schaul, Matteo Hessel, Hado van Hasselt, Marc Lanctot, and Nando de Freitas. Dueling network architectures for deep reinforcement learning, 2016. URL https://arxiv.org/abs/1511.06581.
- Peter R Wurman, Samuel Barrett, Kenta Kawamoto, James MacGlashan, Kaushik Subramanian, Thomas J Walsh, Roberto Capobianco, Alisa Devlic, Franziska Eckert, Florian Fuchs, et al. Outracing champion gran turismo drivers with deep reinforcement learning. *Nature*, 602(7896):223–228, 2022.
- Eliot Xing, Vernon Luk, and Jean Oh. Stabilizing reinforcement learning in differentiable multiphysics simulation. *arXiv preprint arXiv:2412.12089*, 2024.
- Ruihan Yang, Huazhe Xu, Yi Wu, and Xiaolong Wang. Multi-task reinforcement learning with soft modularization, 2020. URL https://arxiv.org/abs/2003.13661.
- Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning, 2020. URL https://arxiv.org/abs/2001.06782.
- Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Avnish Narayan, Hayden Shively, Adithya Bellathur, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning, 2021. URL https://arxiv.org/abs/1910.10897.

Kevin Zakka, Baruch Tabanpour, Qiayuan Liao, Mustafa Haiderbhai, Samuel Holt, Jing Yuan Luo, Arthur Allshire, Erik Frey, Koushil Sreenath, Lueder A. Kahrs, Carlo Sferrazza, Yuval Tassa, and Pieter Abbeel. Mujoco playground: An open-source framework for gpu-accelerated robot learning and sim-to-real transfer., 2025. URL https://github.com/google-deepmind/mujoco_playground.

Yuke Zhu, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Abhishek Joshi, Soroush Nasiriany, and Yifeng Zhu. robosuite: A modular simulation framework and benchmark for robot learning. arXiv preprint arXiv:2009.12293, 2020.

Ziwen Zhuang, Zipeng Fu, Jianren Wang, Christopher Atkeson, Soeren Schwertfeger, Chelsea Finn, and Hang Zhao. Robot parkour learning. *arXiv preprint arXiv:2309.05665*, 2023.

Supplementary Materials

The following content was not necessarily subject to peer review.

C MTRL Approaches

Here, we present an overview of each state-of-the-art MTRL baseline in MTBench.

C.1 Gradient manipulation methods

Gradient manipulation methods compute a new gradient of the multi-task objective, incurring the overhead of solving an optimization problem per iteration as well as storing and computing K task gradients.

PCGrad: Projecting Conflicting Gradients (Yu et al., 2020) observe when the gradients of any two task objectives l_i conflict (defined as having negative cosine similarity) and when their magnitudes are sufficiently different, optimization using the average gradient will cause negative transfer. It attempts to resolve gradient confliction by a simple procedure manipulating each task gradient ∇l_i to be the result of iteratively removing the conflict with each task gradient ∇l_i , $\forall j \in [K], j \neq i$.

$$\nabla l_i' \leftarrow \nabla l_i - \frac{\nabla l_i^T \nabla l_j}{\|\nabla l_j\|^2} \nabla l_j \quad \text{if} \quad \nabla l_i^T \nabla l_j < 0$$
(3)

CAGrad: Conflict-Averse Gradient descent (Liu et al., 2024) resolves the gradient conflict by finding an update vector $d \in \mathbb{R}^m$ that minimizes the worst-case gradient conflict across all the tasks. More specifically, let g_i be the gradient of task $i \in [K]$, and g_0 be the gradient computed from the average loss, CAGrad seeks to solve such an optimization problem:

$$\max_{d \in \mathbb{R}^m} \min_{i \in [K]} \langle g_i, d \rangle \quad \text{s.t.} \quad \|d - g_0\| \le c \|g_0\| \tag{4}$$

Here, $c \in [0,1)$ is a pre-specified hyper-parameter that controls the convergence rate. The optimization problem looks for the best update vector within a local ball centered at the averaged gradient g_0 , which also minimizes the conflict in losses $\langle g_i, d \rangle$.

FAMO: Fast Adaptive Multitask Optimization (Liu et al., 2023a) addresses the under-optimization of certain tasks when using standard gradient descent on averaged losses without incurring the O(K) cost to compute and store all task gradients, which can be significant, especially as the number of tasks increases. FAMO leverages loss history to adaptively adjust task weights, ensuring balanced optimization across tasks while maintaining O(1) space and time complexity per iteration.

C.2 Neural Architectures

Neural Architecture methods seek to avoid task interference by learning shared representations, which are fed to the prediction head. Such representations accelerate MTRL.

CARE: Contextual Attention-based Representation learning (Sodhani & Zhang, 2021) utilizes metadata associated with the set of tasks to weight the representations learned by a mixture of encoders through the attention mechanism.

MOORE: Mixture Of Orthogonal Experts (Hendawy et al., 2024) uses a mixture of experts to encode the state and orthogonalizes those representations to encourage diversity, weighting these representations from a task encoder.

PaCo: Parameter Compositional (Sun et al., 2022) learns a base parameter set $\phi = [\phi_1 \cdots \phi_k]$ and task-specific compositional vector w_k such that multiplying ϕ and w_k represents the task parameters θ_k .

Soft-Modularization: Yang et al. (2020) also uses a mixture of experts to encode the state but also uses a routing network to softly combine the outputs at each layer based on the task.

D Hyperparameter Details

In this section, we provide hyperparameter values for each MTRL approach.

Description	value	variable_name
Number of environments	24576 / 24576	num_envs
Network hidden sizes	[256,128,64]	network.mlp.units
Minibatch size	16384 / 32768	minibatch_size
Horizon length	32	horizon
Mini-epochs	5	mini_epochs
Number of epochs	1272 / 1272	max_epochs
Episode length	150	episodeLength
Discount factor	0.99	gamma
Clip ratio	0.2	e_clip
Policy entropy coefficient	.005	entropy_coef
Optimizer learning rate	5e-4	learning_rate
Optimizer learning schedule	fixed	lr_schedule
Advantage estimation tau	0.95	tau
Value Normalization by task	True	normalize_value
Input Normalization by task	True	normalize_input
Separate critic and policy networks	True	network.separate
CARE-Specific Hyperparameters		
Network hidden sizes	[400,400,400]	care.units
Mixture of Encoders experts	6	encoder.num_experts
Mixture of Encoders layers	2	encoder.num_layers
Mixture of Encoders hidden dim	50	encoder.D
Attention temperature	1.0	encoder.temperature
Post-Attention MLP hidden sizes	[50,50]	attention.units
Context encoder hidden sizes	[50,50]	context_encoder.units
Context encoder bias	True	context_encoder.bias
MOORE-Specific Hyperparameters		
MoE experts	4/6	moore.num_experts
MoE layers	3	moore.num_layers
MoE hidden dim	400	moore.D
Activation before/after task encoding weighting	[Linear, Tanh]	moore.agg_activation
Task encoder hidden sizes	[256]	task_encoder.units
Task encoder bias	False	task_encoder.bias
PaCo-Specific Hyperparameters		
Number of Compositional Vectors	5 / 20	paco.K
Network hidden dim	400	paco.D
Network layers	3	paco.num_layers
Task encoder bias	False	task_encoder.bias
Task encoder init	orthogonal	task_encoder.compositional_initializer
Task encoder activation	softmax	task_encoder.activation
Soft-Modularization-Specific Hyperparameters	sortmax	task_cheoder.activation
MoE experts	2	soft_network.num_experts
MoE layers	4	soft_network.num_layer
State encoder hidden sizes	[256,256]	state_encoder.units
Task encoder hidden sizes	[256]	task_encoder.units
PCGrad Hyperparameters		
Number of environments	24576 / 8192	num_envs
Project actor gradient	False	project_actor_gradient
,		1 J — C
3 6	True	project_critic_gradient
Project critic gradient CAGrad Hyperparameters		project_critic_gradient
Project critic gradient CAGrad Hyperparameters	True	
Project critic gradient CAGrad Hyperparameters Number of environments	True 24576 / 6144	num_envs
Project critic gradient CAGrad Hyperparameters Number of environments Project actor gradient	True 24576 / 6144 False	num_envs project_actor_gradient
Project critic gradient CAGrad Hyperparameters Number of environments Project actor gradient Project critic gradient	True 24576 / 6144 False True	num_envs project_actor_gradient project_critic_gradient
Project critic gradient CAGrad Hyperparameters Number of environments Project actor gradient Project critic gradient Local ball radius for searching update vector	True 24576 / 6144 False	num_envs project_actor_gradient
Project critic gradient CAGrad Hyperparameters Number of environments Project actor gradient Project critic gradient Local ball radius for searching update vector FAMO Hyperparameters	True 24576 / 6144 False True 0.4	num_envs project_actor_gradient project_critic_gradient c
Project critic gradient CAGrad Hyperparameters Number of environments Project actor gradient Project critic gradient Local ball radius for searching update vector FAMO Hyperparameters Regularization coefficient	True 24576 / 6144 False True 0.4	num_envs project_actor_gradient project_critic_gradient c
Project critic gradient CAGrad Hyperparameters Number of environments Project actor gradient Project critic gradient Local ball radius for searching update vector FAMO Hyperparameters Regularization coefficient Learning rate of the task logits	True 24576 / 6144 False True 0.4 1e-3 1e-3	num_envs project_actor_gradient project_critic_gradient c gamma w_lr
Project critic gradient CAGrad Hyperparameters Number of environments Project actor gradient Project critic gradient Local ball radius for searching update vector FAMO Hyperparameters Regularization coefficient	True 24576 / 6144 False True 0.4	num_envs project_actor_gradient project_critic_gradient c

Table 3: Hyperparameters used for MTPPO. A '/' indicates the value used for Meta-World's MT10/MT50 respectively, and otherwise is identical for each setting.

Description	value	variable_name
Number of environments	4096 / 24576	num_envs
Minibatch size	16384 / 76800	minibatch_size
Episode length	150	episodeLength
Horizon length	150	horizon
Mini-epochs	5	mini_epochs
Number of epochs	1908 / 1272	max_epochs
Discount factor	0.99	gamma
Clip ratio	0.2	e_clip
Policy entropy coefficient	.005	entropy_coef
Optimizer learning rate	5e-4	learning_rate
Optimizer learning schedule	fixed	lr_schedule
Advantage estimation tau	0.95	tau
Value Normalization by task	True	normalize_value
Input Normalization by task	True	normalize_input
Separate critic and policy networks	True	network.separate

Table 4: Hyperparameters used for MT-GRPO in MT10 / MT50. A '/' indicates the value used for MT10/MT50 respectively and otherwise is identical for each setting.

Description	value	variable_name
Number of environments	8192	num_envs
Gamma	.99	gamma
Peng's Q(lambda)	.5	q_lambda
Number of minibatches	4	num_minibatches
Episode length	500	episodeLength
Bang-off-Bang	3	binsPerDim
Action Scale	.005	actionScale
Mini epochs	8	mini_epochs
Max grad norm	10.0	max_grad_norm
Horizon	16	horizon
Start epsilon	1.0	start_e
End epsilon	0.005	end_e
Decay epsilon	True	decay_epsilon
Fraction of exploration steps	.005	exploration_fraction
Critic learning rate	3e-4	critic_lr
Anneal learning rate	True	anneal_lr
Value Normalization by task	False	normalize_value
Input Normalization by task	False	normalize_input
Use residual connections	True	q.residual_network
Number of LayerNormAndResidualMLPs	2	q.num_blocks
Network hidden dim	256	q.D
Batch norm input	False	q.norm_first_layer

Table 5: Hyperparameters used for MT-PQN in MT10.

Description	value	variable_name
Number of environments	4096	num_envs
Network hidden sizes	[512,256,128]	network.mlp.units
Gamma	.99	gamma
Separate critic and policy networks	True	network.separate
Number of Gradient steps per epoch	32	gradient_steps_per_itr
Learnable temperature	True	learnable_ temperature
Use distangeled alpha	True	use_disentangled_alpha
Initial alpha	1	init_alpha
Alpha learning rate	5e-3	alpha_lr
Critic learning rate	5e-4	critic_lr
Critic tau	.01	critic_tau
Batch size	8192	batch_size
N-step reward	16	nstep
Grad norm	.5	grad_norm
Horizon	1	horizon
Value Normalization by task	True	normalize_value
Input Normalization by task	True	normalize_input
Replay Buffer Size	5000000	replay_buffer_size
Target entropy coef	1.0	target_entropy_coef

Table 6: Hyperparameters used for MT-SAC in MT10/MT50. A '/' indicates the value used for MT10/MT50 respectively and otherwise is identical for each setting. MT-SAC is very sensitive to the number of environments and replay ratio in the massively parallel regime.

Description	value	variable_name
Minibatch size	16384	minibatch_size
Horizon length	32	horizon
Mini-epochs	5	mini_epochs
Number of epochs	2000 / 4000	max_epochs
Episode length	800	•
Discount factor	0.99	gamma
Clip ratio	0.2	e_clip
Policy entropy coefficient	.005	entropy_coef
Optimizer learning rate	5e-4	learning_rate
Optimizer learning schedule	adaptive	lr_schedule
Advantage estimation tau	0.95	tau
Value Normalization by task	False	normalize_value
Input Normalization by task	False	normalize_input
Separate critic and policy networks	True	network.separate
MOORE-Specific Hyperparameters		
MoE experts	2	moore.num_experts
MoE layers	2	moore.num_layers
MoE hidden dim	256	moore.D
Activation before/after task encoding weighting	[Linear, Linear]	moore.agg_activation
Task encoder hidden sizes	[128]	
Task encoder bias	False	task_encoder.bias
Multihead	False	multi_head
PaCo-Specific Hyperparameters		
Number of Compositional Vectors	5	paco.K
Network hidden dim	400	paco.D
Network layers	3	paco.num_layers
Task encoder bias	False	task_encoder.bias
Task encoder init	orthogonal	task_encoder.compositional_initializer
Task encoder activation	softmax	task_encoder.activation
Soft-Modularization-Specific Hyperparameters		
MoE experts	2	soft_network.num_experts
MoE layers	2	soft_network.num_layer
State encoder hidden sizes	[256,256]	state_encoder.units
Task encoder hidden sizes	[128]	task_encoder.units
PCGrad Hyperparameters		
Project actor gradient	False	project_actor_gradient
Project critic gradient	True	project_critic_gradient
CAGrad Hyperparameters		
Project actor gradient	False	project_actor_gradient
Project critic gradient	True	project_critic_gradient
Local ball radius for searching update vector	0.4	C
FAMO Hyperparameters		
Regularization coefficient	1e-4	gamma
Learning rate of the task logits	5e-3	w_lr
Small value for the clipping of the task logits	1e-3	epsilon
Normalize the task logits gradients	True	norm_w_grad
Transmittee the mark region granteins	1140	

Table 7: Hyperparameters used for MT-PPO in Parkour Benchmark. A '/' indicates the value used for Parkour-easy/Parkour-hard respectively and otherwise is identical for each setting.