
Task-Agnostic Continual Reinforcement Learning: In Praise of a Simple Baseline

Massimo Caccia*
Amazon Web Services
Mila - Quebec AI Institute
Université de Montréal

Jonas Mueller†
Cleanlab

Taesup Kim†
Seoul National University

Laurent Charlin
Mila - Quebec AI Institute
HEC Montréal
Canada CIFAR AI Chair

Rasool Fakoor
Amazon Web Services

Abstract

We study task-agnostic continual reinforcement learning (TACRL) in which standard RL challenges are compounded with *partial observability* stemming from task agnosticism, as well as additional difficulties of continual learning (CL), i.e., learning on a non-stationary sequence of tasks. Here we compare TACRL methods with their soft upper bounds prescribed by previous literature: multi-task learning (MTL) methods which do not have to deal with non-stationary data distributions, as well as task-aware methods, which are allowed to operate under *full observability*. We consider a previously unexplored and straightforward baseline for TACRL, replay-based recurrent RL (3RL), in which we augment an RL algorithm with recurrent mechanisms to address partial observability and experience replay mechanisms to address catastrophic forgetting in CL.

Studying empirical performance in a sequence of RL tasks, we find surprising occurrences of 3RL matching and overcoming the MTL and task-aware soft upper bounds. We lay out hypotheses that could explain this inflection point of continual and task-agnostic learning research. Our hypotheses are empirically tested in continuous control tasks via a large-scale study of the popular multi-task and continual learning benchmark Meta-World. By analyzing different training statistics including gradient conflict, we find evidence that 3RL’s outperformance stems from its ability to quickly infer how new tasks relate with the previous ones, enabling forward transfer.

1 Introduction

Continual learning (CL) creates models and agents that can learn from a sequence of tasks. Continual learning agents promise to solve multiple tasks and adapt to new tasks without forgetting the previous one(s), a major limitation of standard learning agents [18, 52, 39, 32]. In many studies, the performance of CL agents is compared against *multi-task* (MTL) agents that are jointly trained on all available tasks. During learning and evaluation, these multi-task agents are typically provided with the identity of the current task (e.g. each datum is coupled with its task ID). The performance of multi-task agents is thought to provide a soft upper bound on the performance of continual learning

*corresponding author: massimo.p.caccia@gmail.com †work done while at Amazon Web Services

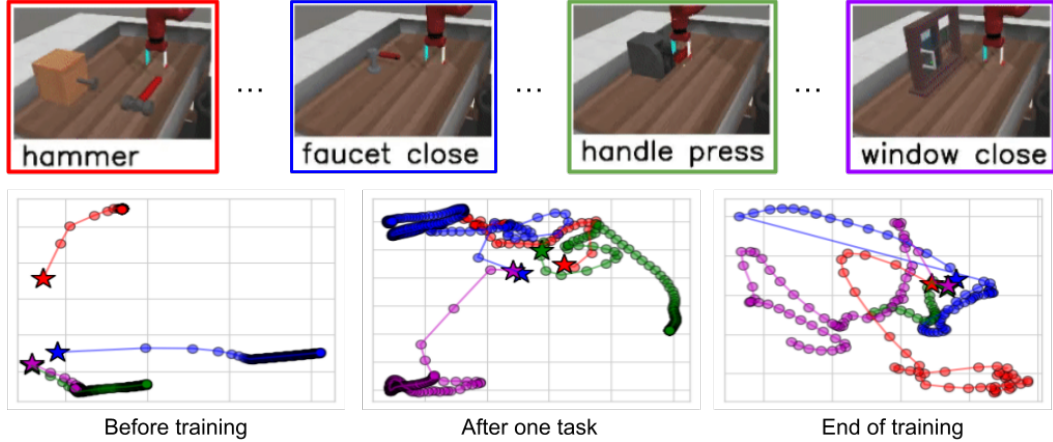


Figure 1: **The Continual-World^c benchmark (top) as well as evolving RNN representations (bottom).** Continual World consists of 10 robotic manipulation environments (four of which are shown above) within the same state space and built on a common reward structure composed of shared components, i.e., reaching, grasping, placing and pushing. We explore the task-agnostic setting in which agents need a full trajectory (rather than a single state) for task identification. We performed a PCA analysis of 3RL’s RNN representations at different stages of training. One episode is shown per task, and the initial task representation (at $t = 0$) is represented by a star. As training progresses, the model learns i) a task-invariant initialization, drawing the initial states closer together, and ii) richer, more diverse representations. Furthermore, the representations constantly evolve throughout the episodes, suggesting the RNN performs more than task inference: it provides useful local information to the policy and critic (see Hypothesis #3 in Sec. 4).

agents since the latter are restricted to learn from tasks in a given sequential order that introduces new challenges, in particular *catastrophic forgetting* [39]. Moreover, continual-learning agents are often trained without knowing the task ID, a challenging setting motivated by practical constraints and known as *task-agnostic CL* [68, 22, 5, 4].

This paper considers continual learning agents that learn using reinforcement learning (RL), i.e. the setting of *continual reinforcement learning* (CRL) or *lifelong reinforcement learning* [25, 47, 48, 2]. We study the task-agnostic continual reinforcement learning (TACRL) setting in challenging robotic manipulation tasks from Meta-World (see top of Fig. 1).

We find two observations challenging common beliefs in CL. First, we surprisingly discover that TACRL agents endowed with a recurrent memory can outperform task-aware agents. We refer to this methodology as *replay-based recurrent reinforcement learning* (3RL) (see bottom of Fig. 1 for a visualization of its representations).

We report a second surprising discovery in which 3RL reaches the performance of its multi-task upper bound (as well as other multi-task RL methods), despite only being exposed to the tasks in a sequential fashion. Conducting a large-scale empirical study^b in which many RL methods are compared in different regimes, we explore several hypotheses to understand the factors underlying these results. Our study indicates that 3RL: 1) quickly infers how new tasks relate to the previous ones, enabling forward transfer as well as 2) learns representations of the underlying MDPs that reduce *task interference*, i.e., when the gradients of different tasks are in conflict [66].

These findings question the need for forgetting-alleviating and task-inference tools when we bring CL closer to some of its real-world applications, i.e., TACRL on a diverse sequence of challenging and inter-related tasks. While it is conventional to assume that task-agnostic continual RL is strictly more difficult than task-aware multi-task RL, this may not actually be the case for representative multi-task RL benchmarks like Meta-World. Despite being far more broadly applicable, TACRL methods may nonetheless be just as performant as their task-aware and multi-task counterparts.

^bThe codebase to reproduce the results will be made available upon publication.

^cThe figures depict the rendering of Meta-World, and not what the agent observes. The agent’s observation space is mainly composed of object, targets and gripper position. Because of the randomness of those positions, the agents need more than one observation to properly infer the hidden state.

2 Background & Task-agnostic Continual Reinforcement Learning

Here we formally define task-agnostic continual reinforcement learning (TACRL), and contrast it against multi-task RL as well as task-aware settings.

MDP. The RL problem is often formulated using a Markov decision process (MDP) [45]. An MDP is defined by the five-tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, r, \gamma \rangle$ with \mathcal{S} the state space, \mathcal{A} the action space, $\mathcal{T}(s'|s, a)$ the transition probabilities, $r(s, a) \in \mathbb{R}$ the reward obtained by taking action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$, and $\gamma \in [0, 1)$ a scalar constant that discounts future rewards. In RL, the transition probabilities and the rewards are typically unknown and the objective is to learn a policy, $\pi(a|s)$ that maximizes the sum of discounted rewards $\mathcal{R}_t^\pi = \sum_{i=t}^{\infty} \gamma^{i-t} r_i = \sum_{i=t}^{\infty} \gamma^{i-t} r(s_i, a_i)$ generated by taking a series of actions $a_t \sim \pi(\cdot|s_t)$. The Q-value $Q^\pi(s, a)$ corresponding to policy π , is defined as the expected return starting at state s , taking a , and acting according to π thereafter:

$$Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] = r(s, a) + \gamma \mathbb{E}_{s', a'} \left[Q^\pi(s', a') \right]$$

POMDP. In most real world applications, if not all, the full information about an environment or a task is not always available to the agent due to various factors such as limited and/or noisy sensors, different states with identical observations, object occlusion, etc. [36, 16]. For this class of problems in which environment states are not fully observable by the agent, partially-observable Markov decision processes (POMDPs) [24] are used to model the problem. A POMDP is defined by a seven-tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{X}, \mathcal{O}, r, \gamma \rangle$ that can be interpreted as an MDP augmented with an observation space \mathcal{X} and a observation-emission function $\mathcal{O}(x'|s)$. In a POMDP, an agent cannot directly infer the current state of the environment s_t from the current observation x_t . We split the state space into two distinct parts: the one that can be directly unveiled from the observation x_t , which we refer to as s_t^o , and the remainder as the hidden state s_t^h , similarly to [42]. To infer the correct hidden state, the agent has to take its history into account: the policy thus becomes $\pi(a_t|s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$. An obvious choice to parameterize such a policy is with a recurrent neural network [35, 58, 3, 15, 42], as described in Sec. 3. Like MDPs, the objective in POMDP is to learn a policy that maximizes the expected return $\mathbb{E}_{s^h} \left[\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] | s^h \right]$.

Task-agnostic Continual Reinforcement Learning (TACRL). TACRL agents operate in a POMDP special case, explained next, designed to study the *catastrophic forgetting* that plagues neural networks [39] when they learn on non-stationary data distributions, as well as *forward transfer* [60], i.e., a method’s ability to leverage previously acquired knowledge to improve the learning of new tasks [37]. First, TACRL’s environments assume that the agent does not have a causal effect on s^h . This assumption increases the tractability of the problem. It is referred to as an hidden-mode MDP (HM-MDP) [9]. Table 1 provides its mathematical description.

The following assumptions helps narrow down on the forgetting problem and knowledge accumulation abilities of neural networks. TACRL’s assumes that s^h follows a non-backtracking chain. Specifically, the hidden states are locally stationary and are never revisited. Finally, TACRL’s canonical evaluation reports the anytime performance of the methods on all tasks, which we will refer to as *global return*. In this manner, we can tell precisely which algorithm has accumulated the most knowledge about all hidden states at the end of its *life*. The hidden state is often referred to as *context*, but more importantly in CRL literature, it represents a *task*. As each context can be reformulated as a specific MDP, we treat tasks and MDP as interchangeable.

Awareness of the Task Being Faced In practical scenarios, deployed agents cannot always assume full observability, i.e. to have access to a *task label* or *ID* indicating which task they are solving or analogously which hidden state they are in. They might not even have the luxury of being “told” when the task changes (task boundary): agents might have to infer it themselves in a data-driven way. We call this characteristic *task agnosticism* [68]. Although impractical, CL research often treats *task-aware* methods, which observe the task label, as a soft upper bound to their task-agnostic counterpart [55, 68]. Augmented with task labels, the POMDP becomes fully observable, collapsing to an MDP problem.

Multi-task Learning (MTL) For neural network agents, catastrophic forgetting can be simply explained by the stationary data distribution assumption of stochastic gradient descent being violated, such that the network parameters become specific to data from the most recent task. Thus it is generally preferable to train on data from all tasks jointly as in MTL [69]. However this may not be possible in many settings, and thus CL is typically viewed as a more broadly applicable methodology that is expected to perform worse than MTL [48, 7].

For RL specifically, multi-task RL (MTRL) often refers to scenarios with families of similar tasks (i.e. MDPs) where the goal is to learn a policy (which can be contextualized on each task’s ID) that maximizes returns across *all* the tasks [65, 6, 28]. While seemingly similar to CRL, the key difference is that MTRL assumes data from all tasks are readily available during training and each task can be visited as often as needed. These are often impractical requirements, which CRL methods are not limited by. Table 1 summarizes the settings we have discussed in this section. See App. A for a more thorough discussion on TACRL and its related settings.

	T	π	Objective	Evaluation
MDP [51]	$p(s_{t+1} s_t, a_t)$	$\pi(a_t s_t)$	$\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t]$	-
POMDP [24]	$p(s_{t+1}^h, s_{t+1}^o s_t^h, s_t^o, a_t)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h} [\mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	-
HM-MDP [9]	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t) p(s_{t+1}^h s_t^h)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h} [\mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	-
Task-agnostic CRL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t) p(s_{t+1}^h s_t^h)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h} [\mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	$\mathbb{E}_{\tilde{s}^h} [\mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$
Task-Aware CRL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t) p(s_{t+1}^h s_t^h)$	$\pi(a_t s_t^h, s_t^o)$	$\mathbb{E}_{s^h} [\mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	$\mathbb{E}_{\tilde{s}^h} [\mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$
Multi-task RL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t) p(s_{t+1}^h)$	$\pi(a_t s_t^h, s_t^o)$	$\mathbb{E}_{\tilde{s}^h} [\mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	-

Table 1: Summarizing table of the settings relevant to TACRL. For readability purposes, \tilde{s}^h denotes the stationary distribution of s^h . The Evaluation column is left blank when it is equivalent to the Objective one.

3 Methods

In this section, we detail the base algorithm and different model architectures used for assembling different continual and multi-task learning baselines.

3.1 Algorithms

We use off-policy RL approaches which have two advantages for (task-agnostic) continual learning. First, they are more sample efficient than online-policy ones [19, 14]. Learning from lower-data regimes is important for CRL as a task is likely to only be seen once as data comes in stream. Second, task-agnostic CRL most likely requires some sort of replay function [54, 33]. This is in contrast to task-aware methods which can, at the expense of computational efficiency, *freeze-and-grow*, e.g. PackNet [38], to incur no forgetting. Off-policy methods, by decoupling the learning policy from the acting policy, support replaying of past data. In short, off-policy learning is the approach of choice in CRL.^d

Base algorithm The Soft Actor-Critic (SAC) [20] is an off-policy actor-critic algorithm for continuous actions. SAC adopts a maximum entropy framework that learns a stochastic policy which maximizes the expected return and also encourages the policy to contain some randomness. To accomplish this, SAC utilizes an actor/policy network π_{ϕ} and critic/Q network Q_{θ} , parameterized by ϕ and θ respectively. Q-values are learnt by minimizing one-step temporal difference (TD) error by sampling previously collected data from the replay buffer [34]. For more details on SAC, please look at App. B.

^dNote that the findings of this paper are not limited to off-policy methods, in fact our 3RL model can be extended to any on-policy method as long as it utilizes a replay buffer [14] Having the capability to support a replay buffer is more important than being on-policy or off-policy.

3.2 Models

We consider various architectures to handle multi-task learning (MTL) as well as continual learning (CL) in both task-aware and task-agnostic setting.

Task ID modeling (TaskID). We assume that a model such as SAC can become task adaptive by providing task information to the networks. Task information such as task ID (e.g. one-hot representation), can be fed into the critics and actor networks as an additional input: $Q_\theta(s, a, k)$ and $\pi_\phi(a|s, k)$ where k is the task ID. We refer to this baseline as Task ID modeling (TaskID). This method is applicable to both multi-task learning and continual learning.

Multi-head modeling (MH). For multi-task learning (which is always task-aware), the standard SAC is typically extended to have multiple heads [65, 67, 61, 66, 19], where each head is responsible for a single distinctive task, i.e. $Q_\Theta = \{Q_{\theta_k}\}_k^K$ and $\pi_\Phi = \{\pi_{\phi_k}\}_k^K$ where K denotes total number of tasks. MH is also applicable to all reinforcement learning algorithms. That way, the networks can be split into 2 parts: (1) a shared state representation network (feature extractor) and (2) multiple prediction networks (heads). This architecture can also be used for task-aware CL, where a new head is newly attached (initialized) when an unseen task is detected during learning.

We also use this architecture in task-agnostic setting for both MTL and CL. Specifically, the number of heads is fixed a priori (we fix it to the number of total tasks) and the most confident actor head, w.r.t. the entropy of the policy, and most optimistic critic head are chosen. Task-agnostic multi-head (TAMH) can help us fraction the potential MH gains over the base algorithm: if MH and TAMH can improve performance, some of MH gains can be explained by its extra capacity instead of the additional task information.

Task-agnostic recurrent modeling (RNN). Recurrent neural networks are able to encode the history of past data [35, 58, 3, 15, 42]. Their encoding can implicitly identify different tasks (or MDP). Thus, we introduce RNNs as a history encoder where history is defined by $\{(s_i, a_i, r_i)\}_i^N$ and we utilize the hidden states z as additional input data for the actor $\pi_\phi(a|s, z)$ and critic $Q_\theta(s, a, z)$ networks. This allows us to train models without any explicit task information, and therefore we use this modeling especially for task-agnostic continual learning. More details about the RNN are provided at the end of the next subsection.

3.3 Baselines

FineTuning is a simple approach to a CL problem. It learns each incoming task without any mechanism to prevent forgetting. Its performance on past tasks indicates how much forgetting is incurred in a specific CL scenario.

Experience Replay (ER) accumulates data from previous tasks in a buffer for retraining purposes, thus slowing down forgetting [48, 1, 8, 30]. Although simple, it is often a worthy adversary for CL methods. One limitation of replay is that, to approximate the data distribution of all tasks, its compute requirements scale linearly with the number of tasks. To alleviate this problem, we use a strategy that caps replay by oversampling the current task from the buffer explained in Alg. 1 L8-9.

Multi-task (MTL) trains on all tasks simultaneously and so it does not suffer from the challenges arising from learning on a non-stationary task distribution. It serves as a soft upper bound for CL methods.

Independent learns a set of separate models for each task whereby eliminating the CL challenges as well as the MTL ones, e.g., learning with conflicting gradients [66].

The aforementioned baselines are mixed-and-matched with the architectural choices to form different baselines, e.g. MTL with TaskID (MTL-TaskID) or FineTuning with MH (FineTuning-MH). At the core of this work lies a particular combination, explained next.

Replay-based Recurrent RL (3RL) A general approach to TACRL is to combine ER—one of CL’s most versatile baseline—with an RNN, one of RL’s most straightforward approach to handling partial observability. We refer to this baseline as *replay-based recurrent RL* (3RL). As an episode unfolds, 3RL’s RNN representations $z_t = \text{RNN}(\{(s_i, a_i, r_i)\}_{i=1}^{t-1})$ should predict the task with increasing

accuracy, thus helping the actor $\pi_\theta(a|s, z)$ and critic $Q_\phi(s, a, z)$ in their respective approximations. We will see in Sec. 4, however, that the RNN delivers more than expected: it enables forward transfer by decomposing new tasks and placing them in the context of previous ones. We provide pseudocode for 3RL in Alg. 1, which we kept agnostic to the base algorithm and not tied to episodic RL. Note that in our implementation, the actor and critics enjoy their own RNNs, as in [15, 42]: they are thus parameterize by θ and ϕ , respectively.

Algorithm 1: 3RL in TACRL

Environment: a set of K MPDs, allowed timesteps T

Input: initial parameters θ , empty FIFO replay buffers $\{\mathcal{D}^i\}_i^K$, replay cap β , batch size b , history length h

```

1 for task  $k$  in  $K$  do
2   set environment to  $n^{th}$  MDP
3   for times-steps  $t$  in  $T$  do
4     /* Sampling stage */
5     compute dynamic task representation  $z_t = \text{RNN}_\theta(\{(s_i, a_i, r_i)\}_{i=t-h-1}^{t-1})$ 
6     observe state  $s_t$  and execute action  $a_t \sim \pi_\theta(\cdot|s_t, z_t)$ 
7     observe reward  $r_t$  and next state  $s_{t+1}$ 
8     store  $(s_t, a_t, r_t, s_{t+1})$  in buffer  $\mathcal{D}^k$ 
9     /* Updating stage */
10    sample a batch  $B$  of  $b \times \min(\frac{1}{n}, 1 - \beta)$  trajectories from the current replay buffer  $\mathcal{D}^k$ 
11    append to  $B$  a batch of  $b \times \min(\frac{n-1}{n}, \beta)$  trajectories from the previous buffers  $\{\mathcal{D}^i\}_i^{k-1}$ 
12    Compute loss on  $B$  and accordingly update parameters  $\theta$  with one step of gradient descent

```

4 Empirical Findings

We now investigate some alluring behaviours we have come upon, namely that replay-based recurrent reinforcement learning (3RL), a task-agnostic continual reinforcement learning (TACRL) baseline, can outperform other task-agnostic but more importantly task-aware baselines, as well as match its MTL soft upper bound.

Benchmarks The benchmark at the center of our empirical study is Meta-World [67], which has become the canonical evaluation protocol for multi-task reinforcement learning (MTRL) [66, 65, 29, 49]. Meta-World offers a suite of 50 distinct robotic manipulation environments. What differentiates Meta-World from previous MTRL and meta-reinforcement learning benchmarks [46] is the broadness of its task distribution. Specifically, the different manipulation tasks are instantiated in the same state and action space^e and share a reward structure, i.e., the reward functions are combinations of reaching, grasping, and pushing different objects with varying shapes, joints and connectivity. Meta-World is thus fertile ground for algorithms to transfer skills across tasks, while representing the types of tasks likely relevant for real-world RL applications (see Fig. 1 for a rendering of some of the environments). Consequently, its adoption in CRL is rapidly increasing [60, 40, 4].

In this work, we study [CW10](#), a benchmark introduced in [60] with a particular focus on *forward transfer*, namely, by comparing a method’s ability to outperform one trained from scratch on new tasks. [CW10](#) is composed of a particular subset of Meta-World conducive for forward transfer and prescribes 1M steps per task, where a step corresponds to a sample collection and an update.

We also study a new benchmark composed of the 20 first alphabetical tasks of Meta-World which we will use to explore more challenging regime: the task sequence is twice as long and data and compute are constrained to half, i.e., 500k steps are allowed per task. We refer to this benchmark as [MW20](#).

In terms of metrics, the reported global and current success are the average success on all tasks and average success on the task that the agent is currently learning, respectively.

Experimental Details We use the hyperparameters prescribed by Meta-World for their Multi-task SAC (MTL-SAC) method. We ensured the performance of our SAC implementation on the [MT10](#), one of Meta-World’s prescribed MTRL benchmark, matches theirs (see App. E). For more details

^ethe fixed action space is an important distinction with traditional incremental supervised learning

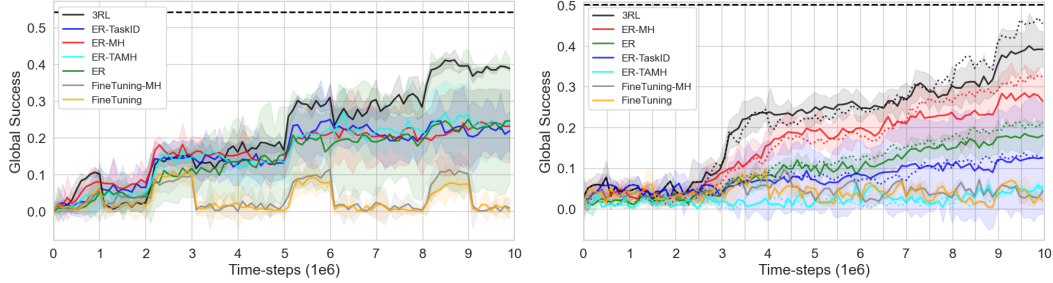


Figure 2: **3RL outperforms all baselines in both CW10 (left) and MW20 (right).** The horizontal line is the reference performance of training independent models on all tasks. The dotted lines (right plot) represent the methods’ performance when they oversample recently collected data. 3RL outperforms all other task-agnostic and more interestingly task-aware baselines. As a side note, ER is high variance as it attempts to solve the POMDP directly without having any explicit or implicit mechanism to do so.

on the hyperparameters and training, see App. D. We test the methods using 8 seeds and report 90% confidence intervals as the shaded area of the figures.

As explained in Alg. 1, we have employed an scheme that oversamples recently gathered data scheme. Whenever we mention oversampling, we mean that the ER algorithm never spends more than 80% of its compute budget replaying old tasks. Without this feature, an ER algorithm would, for example, spend 90% of its compute budget when learning the 10th task.

Task Agnosticism overcomes Task Awareness Our first experiments are conducted on CW10 and MW20 and are reported in Fig. 2. Interestingly, 3RL outperforms all other methods in both benchmarks. This is surprising for two reasons. First, at training time, the task-aware methods learn task-specific parameters to adapt to each individual task. Hence, they suffer less or no forgetting (in general). Second, at test time, the task-agnostic method has to first infer the task through exploration, at the expense of exploitation (see bottom of Fig. 1 for a visualization).

Of course, one could add an RNN to a task-aware method. We intend to compare the effect of learning task-specific parameters compared to learning a common task inference network (the RNN). Nevertheless, a combination of the RNN with ER-MH which was unfruitful (see App. G).

To ensure that these results are not a consequence of the methods having different number of parameters, we learned the CW10 benchmark with bigger networks and found the performance to drop across all methods (see App. H). Further, in our experiments ER-MH is the method with the most parameters and it is outperformed by 3RL.

Continual Learning can Match Multi-Task Learning Multi-task learning is often used as a soft upper bound in evaluating CL methods in both supervised [1, 37, 11] and reinforcement learning [48, 54, 60]. The main reason is that in the absence of additional constraints, multi-task learning (jointly training with data from all tasks in a stationary manner) does not suffer from the catastrophic forgetting that typically plagues neural networks trained on non-stationary data distributions.

3RL can reach this multi-task learning upper bound. In Fig. 3 we report, for the second time, the results of the MW20 experiments. This time, we focus on methods that oversample the current task and more importantly, we report the performance of each method’s multi-task analog, i.e. their soft-upper bound (dashed line of the same color). 3RL, is the only approach that matches the performance of its MTRL equivalent. We believe it is the first time that a specific method achieves the same performance in a non-stationary task regime compared to the stationary one, amidst the introduced challenges like catastrophic forgetting.

For the remainder of the section, we investigate some hypotheses that might explain 3RL’s alluring behavior. We first hypothesise that the RNN boosts performance because it is simply better at learning a single MDP (Hypothesis #1). Next, we investigate the hypothesis that 3RL reduces parameter movement, as it is often a characteristic of successful continual-learning methods (Hypothesis #2). We then explore the hypothesis that the RNN correctly places the new tasks in the context of previous ones (Hypothesis #3). We discuss one last hypothesis in App. K.

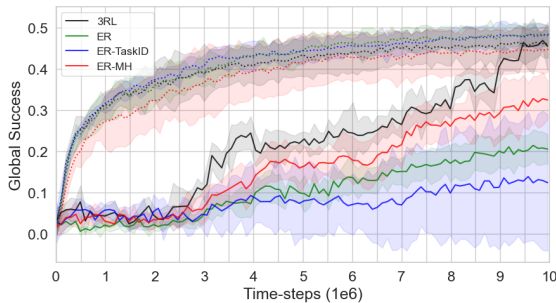


Figure 3: **3RL reaches its MTRL soft-upper bound.** MW20 in solid lines vs MTRL-MW20 in dotted lines. 3RL methods matches its soft-upper bound MTL analog as well as the other MTRL baselines. In contrast, other baselines’ performance are drastically hindered by the non-stationary task distribution.

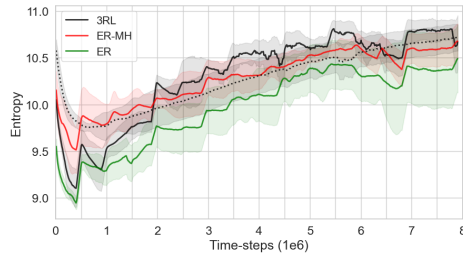


Figure 4: **3RL doesn’t achieve superior continual learning performance through increased parameter stability.** We show the evolution of the methods’ entropy in the parameters updates. We include MTL-RNN (dotted line) as a ref. We do not observe an increase in parameter stability: on the contrary, all methods, increasingly update more weights as new tasks (or data) come in.

Hypothesis #1: RNN individually improves the single-task performance A simple explanation is that the RNN enhances SAC’s ability to learn each task independently, perhaps providing a different inductive bias beneficial to each individual task. To test this hypothesis, we run the CW10 benchmark this time with each task learned separately, which we refer to as the Independent baselines. Note that this hypothesis is unlikely since the agent observes the complete state which is enough to act optimally (i.e. the environments are MDPs not POMDPs). Unsurprisingly, the RNN does not appear to improve performance in STL settings and we discard this hypothesis. Appendix F provides a complete analysis of STL results.

Hypothesis #2: RNN increases parameter stability, thus decreasing forgetting The plasticity-stability tradeoff is at the heart of continual learning: plasticity eases the learning of new tasks. Naive learning methods assume stationary data and so are too plastic in non-stationary regimes leading to catastrophic forgetting. To increase stability, multiple methods enforce [38] or regularize for [27] *parameter stability*, i.e., the tendency of a parameter to stay within its initial value while new knowledge is incorporated. Carefully tuned task-aware methods, e.g. PackNet [38] in [60], have the ability to prevent forgetting.^f

Considering the above, we ask: could 3RL implicitly increase parameter stability? To test this hypothesis we measure the entropy of the weight changes throughout an epoch of learning defined by all updates in between an episode collection. Details about this experiment are found in App. I.

Fig. 4 shows the entropy of 3RL, ER, and ER-MH. We use these baselines since the gap between ER and its upper bound is the largest and the ER-MH gap is in between ER’s and 3RL’s. We find strong evidence to reject our hypothesis. After an initial increase in parameter stability, weight movement increases as training proceeds across all methods and even spikes when a new task is introduced (every 500K steps). MTL-RNN follows the same general pattern as the ER methods.

Hypothesis #3: RNN correctly places the new tasks in the context of previous ones, enabling forward transfer and improving optimization As in real robotic use-cases, MW tasks share a set of low-level reward components like grasping, pushing, placing, and reaching, as well as set of object with varying joints, shapes, and connectivity. As the agent experiences a new task, the RNN could quickly infer how the new data distribution relates with the previous ones and provide to the actor and critics a useful representation. Assume the following toy example: task one’s goal is to grasp a door handle, and task two’s to open a door. The RNN could infer from the state-action-reward trajectory that the second task is composed of two subtasks: the first one as well as a novel pulling one. Doing so would increase the policy learning’s speed, or analogously enable forward transfer.

^fThe observation that PackNet outperforms *an* MTRL baseline in [60] is different from our stronger observation that a *single* method, namely 3RL, achieves the same performance in CRL than in MTRL

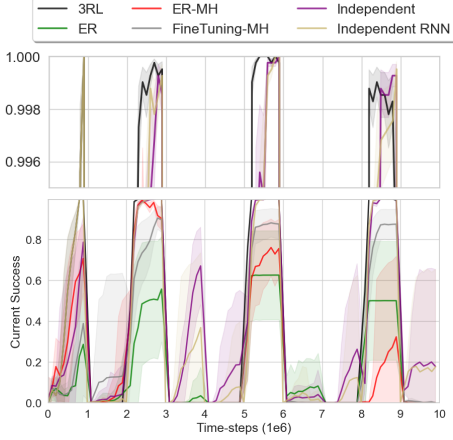


Figure 5: **3RL is the fastest learner.** Current success rate on **CW10**. The Independent method, which trains on each task individually, is still the best approach to maximise single-task performance. However, on the task that the continual-learning methods succeed at, 3RL is the fastest learner. In these cases, its outperformance over Independent and Independent RNN indicates that forward transfer is achieved.

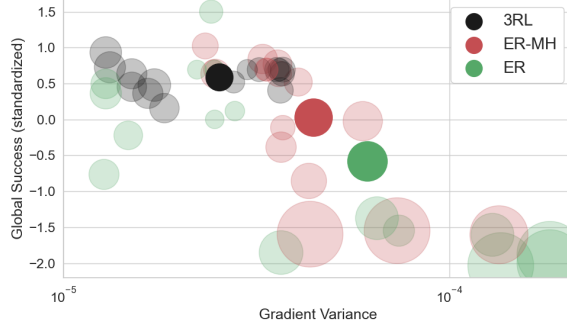


Figure 6: **3RL decreases gradient conflict leading to an increase in training stability and performance.** The global success and gradient as measured by the variance of the gradients are shown are plotted against each other. Training instability as measured by the variance of the Q-values throughout learning is represented by the markers’ size, in a log-scale. Transparent markers depict seeds, whereas the opaque one the means. We observe a negative correlation between performance and gradient conflict (-0.75) as well as performance and training stability (-0.81), both significant under a 5% significance threshold. The hypothesis is that 3RL improves performance by reducing gradient conflict via dynamic task representations.

Now consider a third task in which the agent has to close a door. Again, the first part of the task consists in grasping the door handle. However, now the agent needs to subsequently push and not pull, has was required in task two. In this situation, task interference [66] would occur. Once more, if the RNN could dynamically infer from the context when pushing or pulling is required, it could modulate the actor and critics to have different behaviors in each tasks thus reducing the interference. Note that a similar task interference reduction should be achieved by task-aware methods. E.g., a multi-head component can enable a method to take different actions in similar states depending on the tasks, thus reducing the task interference.

Observing and quantifying that 3RL learns a representation space in which the new tasks are correctly *decomposed* and placed within the previous ones is challenging. Our initial strategy is to look for effects that should arise if this hypothesis was true (so observing the effect would confirm the hypothesis).

First, we take a look at the time required to adapt to new tasks: if 3RL correctly infers how new tasks relate to previous ones, it might be able to learn faster by re-purposing learned behaviors. Fig. 5 depicts the current performance of different methods throughout the learning of **CW10**. For reference, we provide the results of training separate models, which we refer to as Independent and Independent RNN.

The challenges of Meta-World v2 compounded with the ones from learning multiple policies in a shared network, and handling an extra level of non-stationary, i.e. in the task distribution, leaves the continual learners only learning task 0, 2, 5, and 8. On those tasks (except the first one in which no forward transfer can be achieved) 3RL is the fastest continual learner. Interestingly, 3RL showcases some forward transfer by learning faster than the Independent methods on those tasks. This outperformance is more impressive when we remember that 3RL is spending most of its compute replaying old tasks. We thus find some support for Hypothesis #3.

Second, simultaneously optimizing for multiple tasks can lead to conflicting gradients or task interference [66]. To test for this effect, we use the variance of the gradients on the mini-batch throughout the training as a proxy of gradient conflict (see App. J for a discussion on why the gradient’s variance is superior to the gradients’ angle as a proxy for gradient conflict). In Fig. 6 we show the normalized global success metric plotted against the gradient variance. In line with our intuition, we do find that the RNN increases gradient agreement over baselines. As expected, adding a multi-head scheme can also help, to a lesser extent. We find a significant negative correlation of -0.75

between performance and gradient conflict. Fig. 6 also reports training stability, as measured by the standard deviation of Q-values throughout training (not to be confused with the parameter stability, at the center of Hypothesis #2, which measures how much the parameters move around during training). We find 3RL enjoys more stable training as well as a the significant negative correlation of -0.81 between performance and training stability. Note that The plausibility of Hypothesis #3 is thus further increased.

We wrap up the hypothesis with some qualitative support for it. Fig. 1 showcases the RNN representations as training unfolds. If the RNN was merely performing task inference, we would observe the trajectories getting further from each other, not intersecting, and collapsing once the task is inferred. Contrarily, the different task trajectories constantly evolve and seem to intersect in particular ways. Although only qualitative, this observation supports the current hypothesis.

5 Related Work

To study CRL in realistic settings, [61] introduce the Continual World benchmark and discover that many CRL methods that reduce forgetting lose their transfer capabilities in the process, i.e. that policies learned from scratch generally learn new tasks faster than continual learners. Previous works study CL and compare task-agnostic methods to their upper bounds [68, 55] as well as CL methods compare to their multi-task upper bound [47, 48, 2]. Refer to [25] for an in-depth review of continual RL as well as [31, 21] for a CL in general.

RNN were used in the context of continual supervised learning in the context of language modeling [62, 63] as well as in audio [13, 57]. We refer to [10] for a in-depth review of RNN in continual supervised learning.

As in our work, RNN models have been effectively used as policy networks for reinforcement learning, especially in POMDPs where they can effectively aggregate information about states encountered over time [59, 15, 42, 23]. RNN were used in the context of MTRL [41]. Closer to our work, [50] leverages RNN in a task-aware way to tackle a continual RL problem. To the best of our knowledge, RNN have not been employed within TACRL nor combined with Experience Replay in the context of CRL.

6 Conclusion

We have shown that adding a recurrent memory to task-agnostic continual reinforcement learning allows TACRL methods like 3RL to match their multi-task upper bound and even outperform similar task-aware methods. Our large experiments suggest that 3RL manages to decompose the given task distribution into finer-grained subtasks that recur between different tasks, and that 3RL learns representations of the underlying MDP that reduce task interference.

Our findings question the conventional assumption that TACRL is strictly more difficult than task-aware multi-task RL. Despite being far more broadly applicable, TACRL methods like 3RL may nonetheless be just as performant as their task-aware and multi-task counterparts. The need for the forgetting-alleviating and task-inference tools developed by the CL community is now questioned, at least in CRL.

References

- [1] R. Aljundi, L. Caccia, E. Belilovsky, M. Caccia, M. Lin, L. Charlin, and T. Tuytelaars. Online continual learning with maximal interfered retrieval. In *Advances in Neural Information Processing Systems 32*, pages 11849–11860. Curran Associates, Inc., 2019.
- [2] H. B. Ammar, E. Eaton, P. Ruvolo, and M. Taylor. Online multi-task learning for policy gradient methods. In *International conference on machine learning*, pages 1206–1214, 2014.
- [3] B. Bakker. Reinforcement learning with long short-term memory. In *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic*, NIPS’01, page 1475–1482, Cambridge, MA, USA, 2001. MIT Press.

- [4] G. Berseth, Z. Zhang, G. Zhang, C. Finn, and S. Levine. Comps: Continual meta policy search. *ArXiv*, abs/2112.04467, 2021.
- [5] M. Caccia, P. Rodriguez, O. Ostapenko, F. Normandin, M. Lin, L. Page-Caccia, I. H. Laradji, I. Rish, A. Lacoste, D. Vázquez, et al. Online fast adaptation and knowledge accumulation (osaka): a new approach to continual learning. *Advances in Neural Information Processing Systems*, 33, 2020.
- [6] D. Calandriello, A. Lazaric, and M. Restelli. Sparse multi-task reinforcement learning. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in neural information processing systems 27*, pages 819–827. Curran Associates, Inc., 2014.
- [7] A. Chaudhry, P. K. Dokania, T. Ajanthan, and P. H. Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *European Conference on Computer Vision (ECCV)*, 2018.
- [8] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. S. Torr, and M. Ranzato. Continual learning with tiny episodic memories. *ArXiv*, abs/1902.10486, 2019.
- [9] S. P. Choi, D.-Y. Yeung, and N. L. Zhang. Hidden-mode markov decision processes for nonstationary sequential decision making. In *Sequence Learning*, pages 264–287. Springer, 2000.
- [10] A. Cossu, A. Carta, V. Lomonaco, and D. Bacciu. Continual learning for recurrent neural networks: an empirical evaluation. *Neural Networks*, 143:607–627, 2021.
- [11] M. Delange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [12] F. Doshi-Velez and G. D. Konidaris. Hidden parameter markov decision processes: A semiparametric regression approach for discovering latent task parametrizations. *IJCAI : proceedings of the conference*, 2016:1432–1440, 2016.
- [13] B. Ehret, C. Henning, M. R. Cervera, A. Meulemans, J. Von Oswald, and B. F. Grewe. Continual learning in recurrent neural networks. *arXiv preprint arXiv:2006.12109*, 2020.
- [14] R. Fakoor, P. Chaudhari, and A. J. Smola. P3o: Policy-on policy-off policy optimization. In *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, volume 115, pages 1017–1027. PMLR, 2020.
- [15] R. Fakoor, P. Chaudhari, S. Soatto, and A. J. Smola. Meta-q-learning. In *International Conference on Learning Representations*, 2020.
- [16] R. Fakoor and M. Huber. Improving tractability of pomdps by separation of decision and perceptual processes. In *2012 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 593–598, 2012.
- [17] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, 2017.
- [18] R. M. French. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4):128–135, 1999.
- [19] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870. PMLR, 2018.
- [20] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 1861–1870. PMLR, 10–15 Jul 2018.
- [21] R. Hadsell, D. Rao, A. Rusu, and R. Pascanu. Embracing change: Continual learning in deep neural networks. *Trends in Cognitive Sciences*, 24:1028–1040, 12 2020.

- [22] X. He, J. Sygnowski, A. Galashov, A. A. Rusu, Y. W. Teh, and R. Pascanu. Task agnostic continual learning via meta learning. *ArXiv*, abs/1906.05201, 2019.
- [23] N. Heess, J. J. Hunt, T. P. Lillicrap, and D. Silver. Memory-based control with recurrent neural networks. *arXiv preprint arXiv:1512.04455*, 2015.
- [24] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 101(1-2):99–134, 1998.
- [25] K. Khetaarpal, M. Riemer, I. Rish, and D. Precup. Towards continual reinforcement learning: A review and perspectives. *arXiv preprint arXiv:2012.13490*, 2020.
- [26] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. 2017.
- [27] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [28] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017.
- [29] A. Kumar, A. Gupta, and S. Levine. Discor: Corrective feedback in reinforcement learning via distribution correction. *arXiv preprint arXiv:2003.07305*, 2020.
- [30] T. Lesort. Continual learning: Tackling catastrophic forgetting in deep neural networks with replay processes. 2020.
- [31] T. Lesort, M. Caccia, and I. Rish. Understanding continual learning settings with data distribution drift analysis. *arXiv preprint arXiv:2104.01678*, 2021.
- [32] T. Lesort, V. Lomonaco, A. Stoian, D. Maltoni, D. Filliat, and N. Díaz-Rodríguez. Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges. *Information Fusion*, 58:52 – 68, 2020.
- [33] T. Lesort, A. Stoian, and D. Filliat. Regularization shortcomings for continual learning. *ArXiv*, abs/1912.03049, 2019.
- [34] L.-J. Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine learning*, 8(3-4):293–321, 1992.
- [35] L.-J. Lin and T. M. Mitchell. Reinforcement learning with hidden states. In *Proceedings of the Second International Conference on From Animals to Animats 2: Simulation of Adaptive Behavior: Simulation of Adaptive Behavior*, page 271–280, Cambridge, MA, USA, 1993. MIT Press.
- [36] M. L. Littman, A. R. Cassandra, and L. P. Kaelbling. Learning policies for partially observable environments: Scaling up. In *Machine Learning Proceedings 1995*, pages 362–370. 1995.
- [37] D. Lopez-Paz and M. Ranzato. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2017.
- [38] A. Mallya and S. Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. 2018.
- [39] M. McCloskey and N. J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989.
- [40] J. A. Mendez, B. Wang, and E. Eaton. Lifelong policy gradient learning of factored policies for faster training without forgetting. 2020.

- [41] T. Nguyen and T. Obafemi-Ajayi. Recurrent network and multi-arm bandit methods for multi-task learning without task specification (abstract). In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2019.
- [42] T. Ni, B. Eysenbach, and R. Salakhutdinov. Recurrent model-free rl is a strong baseline for many pomdps. 2021.
- [43] F. Normandin, F. Golemo, O. Ostapenko, P. Rodriguez, M. D. Riemer, J. Hurtado, K. Khetarpal, R. Lindeborg, L. Cecchi, T. Lesort, L. Charlin, I. Rish, and M. Caccia. Sequoia: A software framework to unify continual learning research. 2022.
- [44] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [45] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. USA, 1st edition, 1994.
- [46] K. Rakelly, A. Zhou, D. Quillen, C. Finn, and S. Levine. Efficient off-policy meta-reinforcement learning via probabilistic context variables. *arXiv:1903.08254*, 2019.
- [47] J. Ribeiro, F. S. Melo, and J. Dias. Multi-task learning and catastrophic forgetting in continual reinforcement learning. *arXiv preprint arXiv:1909.10008*, 2019.
- [48] D. Rolnick, A. Ahuja, J. Schwarz, T. Lillicrap, and G. Wayne. Experience replay for continual learning. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [49] S. Sodhani, A. Zhang, and J. Pineau. Multi-task reinforcement learning with context-based representations. *arXiv preprint arXiv:2102.06177*, 2021.
- [50] A. Y. Sorokin and M. S. Burtsev. Continual and multi-task reinforcement learning with shared episodic memory. *arXiv preprint arXiv:1905.02662*, 2019.
- [51] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.
- [52] S. Thrun and T. M. Mitchell. Lifelong robot learning. *Robotics and autonomous systems*, 15(1-2):25–46, 1995.
- [53] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012.
- [54] R. Traoré, H. Caselles-Dupré, T. Lesort, T. Sun, G. Cai, N. D. Rodríguez, and D. Filliat. Discorl: Continual reinforcement learning via policy distillation. *CoRR*, abs/1907.05855, 2019.
- [55] G. M. van de Ven and A. S. Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- [56] H. van Hasselt, Y. Doron, F. Strub, M. Hessel, N. Sonnerat, and J. Modayil. Deep reinforcement learning and the deadly triad. 2018.
- [57] Z. Wang, C. Subakan, E. Tzinis, P. Smaragdis, and L. Charlin. Continual learning of new sound classes using generative replay. In *2019 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 308–312. IEEE, 2019.
- [58] S. D. Whitehead and L.-J. Lin. Reinforcement learning of non-markov decision processes. *Artificial Intelligence*, 73(1):271–306, 1995. Computational Research on Interaction and Agency, Part 2.
- [59] D. Wierstra, A. Foerster, J. Peters, and J. Schmidhuber. Solving deep memory pomdps with recurrent policy gradients. In *International conference on artificial neural networks*, pages 697–706. Springer, 2007.
- [60] M. Wolczyk, M. Zajac, R. Pascanu, L. Kucinski, and P. Milos. Continual world: A robotic benchmark for continual reinforcement learning. *CoRR*, abs/2105.10919, 2021.

- [61] M. Wolczyk, M. Zajac, R. Pascanu, L. Kucinski, and P. Milos. Continual world: A robotic benchmark for continual reinforcement learning. In *Advances in Neural Information Processing Systems*, 2021.
- [62] T. Wolf, J. Chaumond, and C. Delange. Continuous learning in a hierarchical multiscale neural network. In *ACL*, 2018.
- [63] T. Wu, M. Caccia, Z. Li, Y.-F. Li, G. Qi, and G. Haffari. Pretrained language model in continual learning: A comparative study. In *International Conference on Learning Representations*, 2021.
- [64] A. Xie, J. Harrison, and C. Finn. Deep reinforcement learning amidst lifelong non-stationarity. *arXiv preprint arXiv:2006.10701*, 2020.
- [65] R. Yang, H. Xu, Y. WU, and X. Wang. Multi-task reinforcement learning with soft modularization. *Advances in Neural Information Processing Systems*, 33:4767–4777, 2020.
- [66] T. Yu, S. Kumar, A. Gupta, S. Levine, K. Hausman, and C. Finn. Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*, 33, 2020.
- [67] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning (CoRL)*, 2019.
- [68] C. Zeno, I. Golan, E. Hoffer, and D. Soudry. Task agnostic continual learning using online variational bayes. 2019.
- [69] Y. Zhang and Q. Yang. A survey on multi-task learning. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1, 2021.

Appendix: Task-Agnostic Continual Reinforcement Learning: In Praise of a Simple Baseline

A Extending TACRL’s related settings

We continue the discussion on TACRL’s related settings. In Meta-RL [17], the training task, or analogously the training hidden-states s^h , are different from the testing ones. We thus separate them into disjoint variable s_{train}^h and s_{test}^h . In this setting, some fast adaptation to s^h is always required. Meta-RL doesn’t deal with a non-stationary training task distribution. Its continual counterpart however, i.e. Continual Meta-RL [4], does. Table 2 summarizes the settings.

Noteworthy, the Hidden Parameter MDP (HiP-MDP) [12] is similar to the HM-MDP but assumes a hidden states, i.e., the hidden states are samples i.i.d. Another setting similar to the HM-MD is the Dynamic Parameter MDP [64] in which the hidden-state are non-stationary but change at every episode.

	T	π	Objective	Evaluation
MDP [51]	$p(s_{t+1} s_t, a_t)$	$\pi(a_t s_t)$	$\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t]$	-
POMDP [24]	$p(s_{t+1}^h, s_{t+1}^o s_t^h, s_t^o, a_t)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	-
HM-MDP [9]	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t)p(s_{t+1}^h s_t^h)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	-
Task-agnostic CRL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t)p(s_{t+1}^h s_t^h)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	$\mathbb{E}_{s^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$
Task-Aware CRL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t)p(s_{t+1}^h s_t^h)$	$\pi(a_t s_t^h, s_t^o)$	$\mathbb{E}_{s^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	$\mathbb{E}_{s^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$
Multi-task RL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t)p(s_{t+1}^h s_t^h)$	$\pi(a_t s_t^h, s_t^o)$	$\mathbb{E}_{s^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s^h]$	-
Meta RL [17]	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t)p(s_{t+1}^h s_t^h)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s_{\text{train}}^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s_{\text{train}}^h]$	$\mathbb{E}_{s_{\text{test}}^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s_{\text{test}}^h]$
Continual Meta-RL	$p(s_{t+1}^o s_{t+1}^h, s_t^o, a_t)p(s_{t+1}^h s_t^h)$	$\pi(a_t s_{1:t}^o, a_{1:t-1}, r_{1:t-1})$	$\mathbb{E}_{s_{\text{train}}^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s_{\text{train}}^h]$	$\mathbb{E}_{s_{\text{test}}^h}[\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t] s_{\text{test}}^h]$

Table 2: Summarizing table of the settings relevant to TACRL. For readability purposes, \tilde{s}^h denotes the stationary distribution of s^h . The Evaluation column is left blank when it is equivalent to the Objective one.

B Soft-Actor Critic

The Soft Actor-Critic (SAC) [20] is an off-policy actor-critic algorithm for continuous actions. SAC adopts a maximum entropy framework that learns a stochastic policy which not only maximizes the expected return but also encourages the policy to contain some randomness. To accomplish this, SAC utilizes an actor/policy network (i.e. π_{ϕ}) and critic/Q network (i.e. Q_{θ}), parameterized by ϕ and θ respectively. Q-values are learnt by minimizing one-step temporal difference (TD) error by sampling previously collected data from the replay buffer [34] denoted by \mathcal{D} .

$$\mathcal{J}_Q(\theta) = \mathbb{E}_{s,a} \left[\left(Q_{\theta}(s, a) - y(s, a) \right)^2 \right], \quad a' \sim \pi_{\phi}(\cdot | s') \quad (1)$$

where $y(s, a)$ is defined as follows:

$$y(s, a) = r(s, a) + \gamma \mathbb{E}_{s', a'} \left[Q_{\hat{\theta}}(s', a') - \alpha \log(a' | s') \right]$$

And then, the policy is updated by maximizing the likelihood of actions with higher Q-values:

$$\mathcal{J}_{\pi}(\phi) = \mathbb{E}_{s, \hat{a}} \left[Q_{\theta}(s, \hat{a}) - \alpha \log \pi_{\phi}(\hat{a} | s) \right], \quad \hat{a} \sim \pi_{\phi}(\cdot | s) \quad (2)$$

where $(s, a, s') \sim \mathcal{D}$ (in both (1) and (2)) and α is entropy coefficient. Note that although SAC is used in this paper, other off-policy methods for continuous control can be equally utilized for CRL. SAC is selected here as it has a straightforward implementation and few hyper-parameters.

C Baselines Definitions

FineTuning-MH is FineTuning with task-specific heads. For each new task, it spawns and attaches an additional output head to the actor and critics. Since each head is trained on a single task, this baseline allows to decompose forgetting happening in the representation of the model (trunk) compared to forgetting in the last prediction layer (head). It is a task-aware method.

ER-TaskID is a variant of ER that is provided with task labels as inputs (i.e. each observation also contains a task label). It is a task-aware method that has the ability to learn a task representation in the first layer(s) of the model.

ER-MH is ER strategy which spawns tasks-specific heads [60], similar to FineTuning-MH. ER-MH is often the hardest to beat task-aware baseline []. ER-TaskID and ER-MH use two different strategies for modelling task labels. Whereas, MH uses $|h| \times |A|$ task-specific parameters (head) taskID only uses $|h|$, with $|h|$ the size of the network’s hidden space (assuming the hidden spaces at each layer have the same size) and $|A|$ the number of actions available to the agent.

ER-TAMH (task-agnostic multi-head) is similar to ER-MH, but the task-specific prediction heads are chosen in a task-agnostic way. Specifically, the number of heads is fixed a priori (in the experiments we fix it to the number of total tasks) and the most confident actor head, w.r.t. the entropy of the policy, and most optimistic critic head are chosen. ER-TAMH has the potential to outperform ER, another task-agnostic baseline, if it can correctly infer the tasks from the observations.

MTL is our backbone algorithm, namely SAC, trained via multi-task learning. It is the analog of ER.

MTL-TaskID is MTL, but the task label is provided to the actor and critic. It is the analog of ER-TaskID and is a standard method, e.g. [19].

MTL-MH is MTL with a task-specific prediction network. It is the analog of ER-MH and is also standard, e.g. [66, 19, 67].

MTL-TAMH is similar to MTL-MH, but the task-specific prediction heads are chosen in the same way as in ER-TAMH.

MTL-RNN is similar to MTL, but the actor and critic are mounted with an RNN. It is the analog of 3RL.

D Experimental Details

We’ve used the SAC hyperparameters prescribed by Meta-World [67]. Specifically, we used a learning rate of 1×10^{-3} ; mini-batch size of 1028; actor and critics are 2-layer MLPs with hidden sizes of [400, 400]; episode length of 500 (except 200 in [CW10](#) ; ReLU activation function, soft target interpolation parameter of 5×10^{-3} .

We used automatic entropy tuning except in the MTRL experiments, where we found it to be detrimental. Because their MT-SAC implementation learns a task-specific entropy term, we think this is the reason why they do not observe the same behavior.

As prescribed, we use a minimum buffer batch size of 1500 when doing 10 tasks and 7500 when doing more. We also use a burn in period of 10,000 time-steps, now prescribed by Continual World [61].

Finally, after struggling with some deadly triad [56] problems in CRL and MTRL, we decided to clip the gradients’ norm at 1, which turned out an effective solution.

D.1 Computing Resources

All experiments were performed on Amazon EC2’s P2 instances which incorporates up to 16 NVIDIA Tesla K80 Accelerators and is equipped with Intel Xeon 2.30GHz cpu family.

All experiments included in the paper can be reproduced by running 43 method/setting configurations with 8 seeds, each running for 4.2 days on average.

D.2 Software and Libraries

In the codebase we’ve used to run the experiments, we have leverage some important libraries and software. We used Mujoco [53] and Meta-World [67] to run the benchmarks. We used Sequoia [43] to assemble the particular CRL benchmarks, including [CW10](#) . We used Pytorch [44] to design the neural networks.

E Validating our SAC implementation on MT10

In Fig. 7 we validate our SAC implementation on Meta-World v2’s MT10.

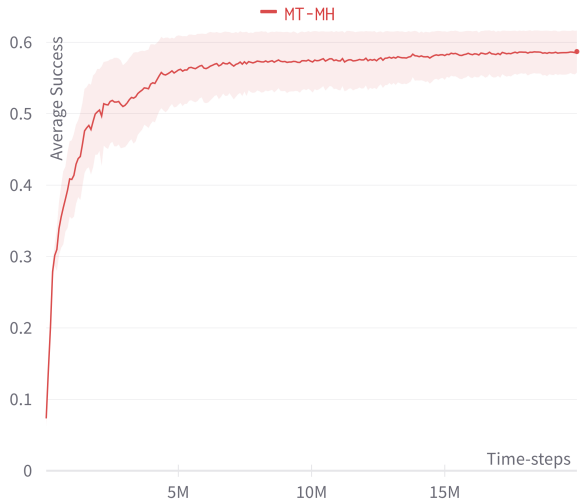


Figure 7: **MT10 experiment.** We repeat the popular MT10 benchmark with our MT-MH implementation. After 20M time-steps, the algorithm reaches a success rate of 58%. This is in line with Meta-World reported results. In Figure 15 of their Appendix, their MT-SAC is trained for 200M time-steps the first 20M time-steps are aligned with our curve. We further note that our CW10 result might seem weak vis à vis the reported ones in [61]. This is explained by [61] using Meta-World-v1 instead of the more recent Meta-World-v2.

F Extended CW10 Single-task results

Results for single-task experiments are shown in Fig. 8 and Fig. 9.

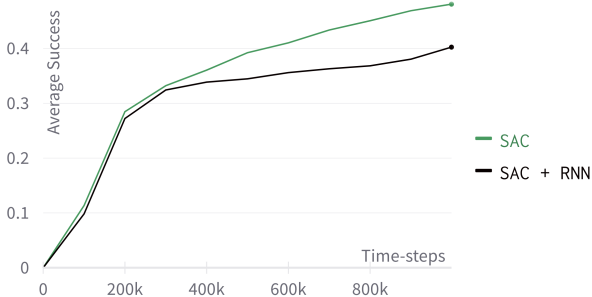


Figure 8: **Single-task learning CW10 experiments.** Average success on all task trained independently. In this regime, the RNN doesn’t help.

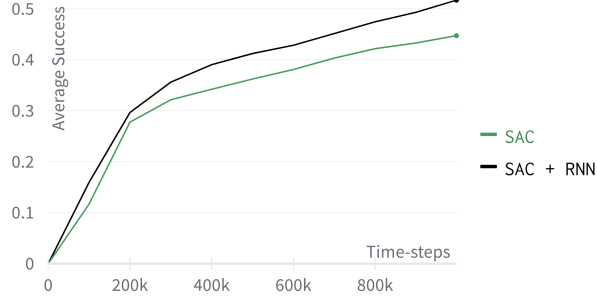


Figure 9: **Single-task learning CW10 experiments without gradient clipping.** We enstored gradient clipping in CRL and MTRL to alleviate the deadly triad problem, a problem we did not find in single-task learning (STL). For completeness, we reran the STL experiments without gradient clipping. We found the performance of the RNN to dramatically increase. Note that this is not the same algorithms used in the CRL and MTRL experiments because of the gradient clipping discrepancy.

G Task-aware meets task-agnostic

In Fig. 10 we show that combining the RNN with MH is not a good proposition in [CW10](#).

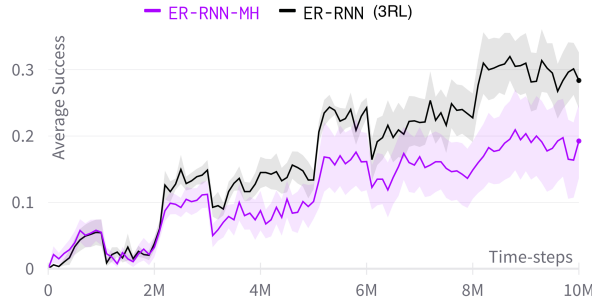


Figure 10: **CW10 experiment combining the RNN and MH.** Combining the best task-agnostic (3RL) and task-aware (ER-MH) CRL methods did not prove useful. Note that this experiment was ran before we enstored gradient clipping, which explains why the performance is lower than previously reported.

H Larger networks don't improve performance

In Fig. 11 we report that increasing the neural net capacity doesn't increase performance.

I Gradient Entropy Experiment

To assess parameter stability, we look at the entropy of the parameters for each epochs. Because the episode are of size 500, the epoch corresponds to 500 updates. To remove the effect of ADAM [26], our optimizer, we approximate the parameters' movement by summing up their absolute gradients throughout the epoch. To approximate the sparsity of the updates, we report the entropy of the absolute gradient sum. For example, a maximum entropy would indicate all parameters are moving equally. If the entropy drops, it means the algorithm is applying sparser updates to the model, similarly to PackNet.

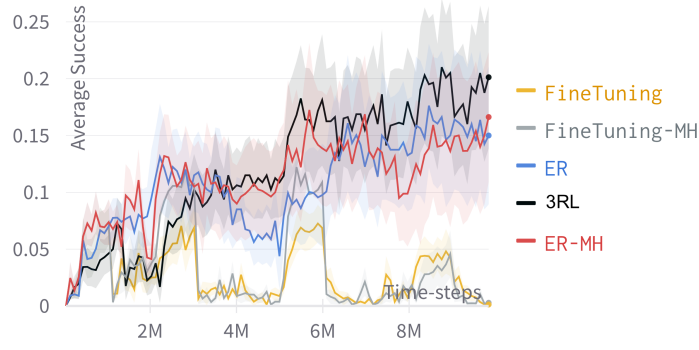


Figure 11: **CW10 experiment with larger neural networks..** We repeated the CW10 experiment, this time with larger neural networks. Specifically, we added a third layer to the actor and critics. Its size is the same as the previous two, i.e. 400. The extra parameters have hindered the performance of all baselines.

J Gradient conflict through time

One might ask why we use the gradients’ variance as a proxy for gradient conflict. Indeed, [66] measures the conflict between two tasks via the angle between their gradients: the tasks conflict if the angle is obtuse. However, we argue that this thinking is inadequate because it does not consider the magnitude of the gradients.

Assume a 1D optimization problem. In case 1, assume that the first task’s gradient is 0.01 and the second’s -0.01. In case 2, assume the gradients are now 0.01 and 0.5. Measuring the conflict via the angle would lead us to think that the tasks are in conflict in case 1 and are not in case 2. This is, however, not the case. The agreement is much higher in case 1: both tasks agree that they should not move too far from the current parameter. In case 2, although the two tasks agree on the direction of the step, they do not agree about the curvature of the loss landscape. The update step will thus be too small and too big for the two tasks. We thus think standard deviation is a better proxy for gradient conflict. We have updated the manuscript to defend our position better.

We’ve included the evolution of gradient conflict in the actor and critics in Fig. 12.

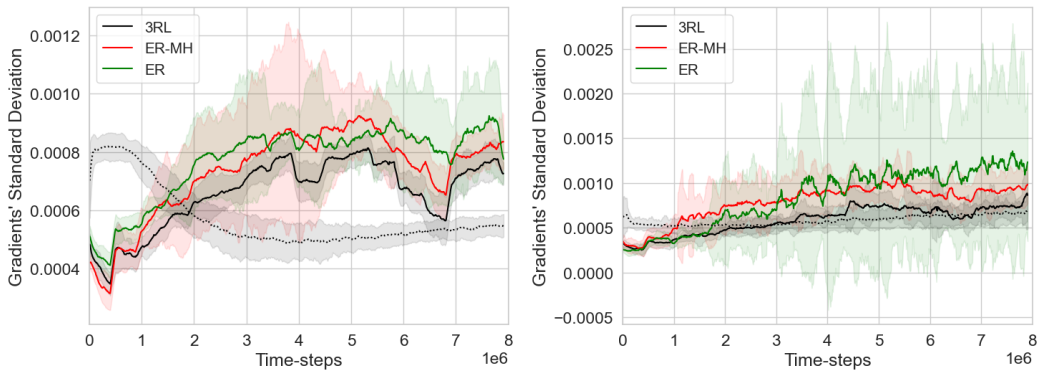


Figure 12: **Gradient variance analysis on CW20.** Comparison of the normalized standard deviation of the gradients for the actor (left) and critics (right) in for different CRL methods. For reference, we included MTL-RNN as the dotted line. The gradient alignment’s rank is perfectly correlated with the performance rank.

K Hypothesis #4: we are operating in regimes where additional task-specific parameter can hurt performance

The task-aware methods learn task-specific parameters which might require more training data. In Figure 2 we compare the relative performance the methods as we increase compute and data from 500k time-steps to 1M. We observe that increasing the data/compute narrows the gap between the 3RL and ER-MH. This increases the plausibility of hypothesis #4.

However, some observations suggest we should look elsewhere. First, in [CW10](#) (Figure Fig. 2) another, another high-data regime, the task-aware methods do not shine. Second, ER-TaskID, which has fewer extra parameters compared to ER-MH, does not work as well in general, and doesn't improve when data/compute is increased two-fold.

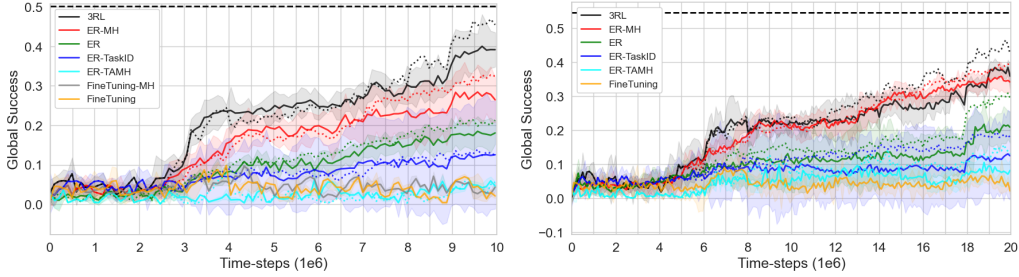


Figure 13: **Increased data/compute experiment.** We show the methods' performance on MW20 with 500k steps (left) and 1M steps (right). Hypothesis #4 is inconclusive.