

# MEMORY OF UNIMAGINABLE OUTCOMES IN EXPERIENCE REPLAY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1 Model-based reinforcement learning (MBRL) applies a single-shot dynamics  
 2 model to imagined actions to select those with best expected outcome. The dy-  
 3 namics model is an unfaithful representation of the environment physics, and its  
 4 capacity to predict the outcome of a future action varies as it is trained iteratively.  
 5 An experience replay buffer collects the outcomes of all actions executed in the  
 6 environment and is used to iteratively train the dynamics model. With growing  
 7 experience, it is expected that the model becomes more accurate at predicting the  
 8 outcome and expected reward of imagined actions. However, training times and  
 9 memory requirements drastically increase with the growing collection of experi-  
 10 ences. Indeed, it would be preferable to retain only those experiences that could  
 11 not be anticipated by the model while interacting with the environment. We argue  
 12 that doing so results in a lean replay buffer with diverse experiences that corre-  
 13 spond directly to the model’s predictive weaknesses at a given point in time.  
 14 We propose strategies for: i) determining reliable predictions of the dynamics  
 15 model with respect to the imagined actions, ii) retaining only the unimaginable  
 16 experiences in the replay buffer, and iii) training further only when sufficient novel  
 17 experience has been acquired. We show that these contributions lead to lower  
 18 training times, drastic reduction of the replay buffer size, fewer updates to the  
 19 dynamics model and reduction of catastrophic forgetting. All of which enable the  
 20 effective implementation of continual-learning agents using MBRL.

## 21 1 INTRODUCTION

22 Model-Based Reinforcement Learning (MBRL) is attractive because it tends to have a lower sample  
 23 complexity compared to model-free algorithms like Soft Actor Critic (SAC) (Haarnoja et al. (2018)).  
 24 MBRL agents function by building a model of the environment in order to predict trajectories of fu-  
 25 ture states based off of imagined actions. An MBRL agent maintains an extensive history of its  
 26 observations, its actions in response to observations, the resulting reward, and new observation in  
 27 an experience replay buffer. The information stored in the replay buffer is used to train a single-shot  
 28 dynamics model that iteratively predicts the outcomes of imagined actions into a trajectory of future  
 29 states. At each time step, the agent executes only the first action in the trajectory, and then the model  
 30 re-imagines a new trajectory given the result of this action (Nagabandi et al. (2018)). Yet, many  
 31 real-world tasks consist in sequences of subtasks of arbitrary length accruing repetitive experiences,  
 32 for example driving over a long straight and then taking a corner. Capturing the complete dynamics  
 33 here requires longer sessions of continual learning. (Xie & Finn (2021))  
 34 Optimization of the experience replay methodology is an open problem. Choice of size and mainte-  
 35 nance strategy for the replay buffer both have considerable impact on asymptotic performance and  
 36 training stability (Zhang & Sutton (2017)). From a resource perspective, the size and maintenance  
 37 strategy of the replay buffer pose major concerns for longer learning sessions.  
 38 The issue of overfitting is also a concern when accumulating similar or repetitive states. The buffer  
 39 can become inundated with redundant information while consequently under-representing other im-  
 40 portant states. Indefinite training on redundant data can result in an inability to generalize to, or  
 41 remember, less common states. Conversely, too small a buffer will be unlikely to retain sufficient  
 42 relevant experience into the future. Ideally, a buffer’s size would be the exact size needed to cap-  
 43 ture sufficient detail for all relevant states (Zhang & Sutton (2017)). Note that knowing a priori all  
 44 relevant states is unfeasible without extensive exploration.

45 We argue that these problems can be subverted by employing a strategy that avoids retaining expe-  
 46 riences that the model already has sufficiently mastered. Humans seem to perform known actions  
 47 almost unconsciously (e.g., walking) but they reflect on actions that lead to unanticipated events  
 48 (e.g. walking over seemingly solid ice and falling through). Such is our inspiration to attempt to  
 49 curate the replay buffer based on whether the experiences are predictable for the model.

50 Through this work, we propose techniques to capture both common and sporadic experiences with  
 51 sufficient detail for prediction in longer learning sessions. The approach comprises strategies for:  
 52 i) determining reliable predictions of the dynamics model with respect to the imagined actions, ii)  
 53 retaining only the unimaginable experiences in the replay buffer, iii) training further only when  
 54 sufficient novel experience has been acquired, and iv) reducing the effects of catastrophic forget-  
 55 ting. These strategies enable a model to self-manage both its buffer size and its decisions to train,  
 56 drastically reducing the wall-time needed to converge. These are critical improvements toward the  
 57 implementation of effective and stable continual-learning agents.

58 Our contributions can be summarized as follows: i) contributions towards the applicability of MBRL  
 59 in continual learning settings, ii) a method to keep the replay buffer size to a minimum without  
 60 sacrificing performance, iii) a method that reduces the training time. These contributions result in  
 61 keeping only useful information in a balanced replay buffer even during longer learning sessions.

## 62 2 RELATED WORK

63 Compared to MFRL, MBRL tends to be more sample-efficient (Deisenroth et al. (2013)) at a cost of  
 64 reduced performance. Recent work by Nagabandi et al. (2018) showed that neural networks effi-  
 65 ciently reduce sample complexity for problems with high-dimensional non-linear dynamics. MBRL  
 66 approaches need to induce potential actions which will be evaluated with a dynamics model to  
 67 choose those with best reward. Random shooting methods artificially generate large number of ac-  
 68 tions (Rao (2010)) and model predictive control (MPC) can be used to select actions (Camacho et al.  
 69 (2004)). Neural networks (NNs) are a suitable alternative to families of equations used to model the  
 70 environment dynamics in MBRL (Williams et al. (2017)). But, overconfident incorrect predictions,  
 71 which are common in DNNs, can be harmful. Thus, quantifying predictive uncertainty, a weak-  
 72 ness in standard NN, becomes crucial. Ensembles of probabilistic NNs proved a good alternative to  
 73 Bayesian NNs in determining predictive uncertainty (Lakshminarayanan et al. (2016)). Further-  
 74 more, an extensive analysis about the types of model that better estimate uncertainty in the MBRL  
 75 setting favored ensembles of probabilistic NNs (Chua et al. (2018)). The authors identified two  
 76 types of uncertainty: aleatoric (inherent to the process) and epistemic (resulting from datasets with  
 77 too few data points). Combining uncertainty aware probabilistic ensembles in the trajectory sam-  
 78 pling of the MPC with a cross entropy controller the authors demonstrated asymptotic performance  
 79 comparable to SAC but with sample efficient convergence. The MPC, however, is still computationally  
 80 expensive (Chua et al. (2018); Zhu et al. (2020)). Quantification of predictive uncertainty serves  
 81 as a notion of confidence in an imagined trajectory. Remonda et al. (2021), used this concept to pre-  
 82 vent unnecessary recalculation, effectively using sequences of actions the model is confident in and  
 83 reducing computations. Our approach also seeks to determine reliable predictions of the dynamics  
 84 model with respect to the imagined actions, but as a basis to manage growth of the experience replay.  
 85 **Use of Experience Replay in MBRL:** While an uncertainty aware dynamics model helps to miti-  
 86 gate the risks of prediction overconfidence, other challenges remain. Another considerable issue  
 87 when training an MBRL agent is the shifting of the state distribution as the model trains. Experi-  
 88 ence replay was introduced by Lin (1992), and has been further improved upon. Typically in RL,  
 89 transitions are sampled uniformly from the replay buffer at each step. Prioritized experience replay  
 90 (PER) (Schaul et al. (2016)) attempts to make learning more efficient by sampling more frequently  
 91 transitions that are more relevant for learning. PER improves how the model samples experiences  
 92 from the already-filled replay buffer, but it does not address how the replay buffer is filled in the  
 93 first place. In addition, neither work addresses the importance of the size of the replay buffer as a  
 94 hyperparameter (Zhang & Sutton (2017)). Our method attempts to balance the replay buffer by only  
 95 adding experiences that should improve the future prediction capacity and keeps the training time  
 96 bounded to a minimum.

97 **Task Agnostic Continual Learning:** The context of our work originates in tasks consisting in com-  
 98 binations of possibly repetitive subtasks of arbitrary length. In the terminology of Normandin et al.  
 99 (2021), we aim for continuous task-agnostic continual reinforcement learning. Meaning that the

100 task boundaries are not observed and transitions may occur gradually (Zeno et al. (2021)). In our  
 101 case, the task latent variable is not observed and the model has no explicit information about task  
 102 transitions. In such context, a continual learner can be seen as an autonomous agent learning over an  
 103 endless stream of tasks, where the agent has to: i) continually adapt in a non-stationary environment,  
 104 ii) retain memories which are useful, iii) manage compute and memory resources over a long period  
 105 of time ( Khetarpal et al. (2020), Thrun (1994)). Our proposed strategies satisfy these requirements.  
 106 Matisen et al. (2020) address the issue of retaining useful memories in a curriculum learning setting  
 107 by training a "teacher" function that mandates a learning and re-learning schedule for the agent  
 108 assuming that the agent will not frequently revisit old experiences/states and will eventually forget  
 109 them. Ammar et al. (2015) focus on agents that acquire knowledge incrementally by learning mul-  
 110 tiple tasks consecutively over their lifetime. Their approach rapidly learns high performance safe  
 111 control policies based on previously learned knowledge and safety constraints on each task, accu-  
 112 mulating knowledge over multiple consecutive tasks to optimize overall performance. Bou Ammar  
 113 & Taylor (2014) developed a lifelong learner for policy gradient RL. Instead of learning a control  
 114 policy for a task from scratch, they leverage on the agent's previously learned knowledge. Knowl-  
 115 edge is shared via a latent basis that captures reusable components of the learned policies. The latent  
 116 basis is then updated with newly acquired knowledge. This resulted in an accelerated learning of  
 117 new task and an improvement in the performance of existing models without retraining on their re-  
 118 spective tasks. With our method, we imbue the RL agent with the ability to self-evaluate and decide  
 119 in real-time if it has sufficiently learned the current state. Unlike the method presented by Matisen  
 120 et al. (2020), our method requires no additional networks to be trained in parallel.

121 Xie & Finn (2021) identified two core challenges in the lifelong learning setting: enabling forward  
 122 transfer, i.e. reusing knowledge from previous tasks to improve learning new tasks, and to improve  
 123 backward transfer which can be seen as avoiding catastrophic forgetting (Kirkpatrick et al. (2017)).  
 124 They developed a method that exploits data collected from previous tasks to cumulatively grow the  
 125 agent's skill-set using importance sampling. Their method requires the agent to know when the  
 126 task changes whereas our method does not have this constrain. Additionally, they focus in forward  
 127 transfer only. Our method addresses both forward and backward transfer.

### 128 3 PRELIMINARIES

129 At each time  $t$ , the agent is at a state  $s_t \in S$ , executes an action  $a_t \in A$  and receives from  
 130 the environment a reward  $r_t = r(s_t, a_t)$  and a state  $s_{t+1}$  according to some environment transi-  
 131 tion function  $f : S \times A \rightarrow S$ . RL consists in training a policy towards maximizing the accu-  
 132 mulated reward obtained from the environment. The goal is to maximize the sum of discounted  
 133 rewards  $\sum_{i=t}^{\infty} \gamma^{(i-t)} r(s_i, a_i)$ , where  $\gamma \in [0, 1]$ . Instead, given a current state  $s_t$ , MBRL arti-  
 134 ficially generates a huge amount of potential future actions, for instance using random shooting ( Rao  
 135 (2010)) or cross entropy( Chua et al. (2018)). Clarification of these methods is beyond the scope  
 136 of this paper; we defer the interested reader to the bibliography. MBRL attempts to learn a dis-  
 137 crete time dynamics model  $\hat{f} = (s_t, a_t)$  to predict the future state  $\hat{s}_{t+\Delta_t}$  of executing action  $a_t$   
 138 at state  $s_t$ . To reach a state into the future, the dynamics model *iteratively* evaluates sequences of  
 139 actions,  $a_{t:t+H} = (a_t, \dots, a_{t+H-1})$  over a longer horizon  $H$ , to maximize their discounted reward  
 140  $\sum_{i=t}^{t+H-1} \gamma^{(i-t)} r(s_i, a_i)$ . These sequences of actions with predicted outcomes are called imagined  
 141 trajectories. The dynamics model  $\hat{f}$  is an inaccurate representation of the transition function  $f$  and  
 142 the future is only partially observable. So, the controller executes only a single action  $a_t$  in the tra-  
 143 jectory before solving the optimization again with the updated state  $s_{t+1}$ . The process is formalized  
 144 in Algorithm 1. The dynamics model  $\hat{f}_\theta$  is learned with data  $\mathcal{D}_{env}$ , collected on the fly. With  $\hat{f}_\theta$ ,  
 145 the simulator starts and the controller is called to plan the best trajectory resulting in  $a_{t:t+H}^*$ . Only  
 146 the first action of the trajectory  $a_t^*$  is executed in the environment and the rest is discarded. This is  
 147 repeated for *TaskHorizon* number of steps. The data collected from the environment is added to  
 148  $\mathcal{D}_{env}$  and  $\hat{f}_\theta$  is trained further. The process repeats for *NIterations*. Note that generating imag-  
 149 ined trajectories requires subsequent calls to the dynamics model to chain predicted future states  
 150  $s_{t+n}$  with future actions up to the task horizon, and so it is only partially parallelizable.

151 **Dynamics model.** We use a probabilistic model to model a probability distribution of next state  
 152 given current state and an action. To be specific, we use a regression model realized using a neural  
 153 network similar to Lakshminarayanan et al. (2016) and Chua et al. (2018). The last layer of the

154 model outputs parameters of a Gaussian distribution that models the aleatoric uncertainty (the un-  
 155 certainty due to the randomness of the environment). Its parameters are learned together with the  
 156 parameters of the neural network. To model the epistemic uncertainty (the uncertainty of the dy-  
 157 namics model due to generalization errors), we use ensembles with bagging where the members of  
 158 the ensemble are identical and only differ in the initial weight values. Each element of the ensemble  
 159 has as input the current state  $s_t$  and action  $a_t$  and is trained to predict the difference between  $s_t$  and  
 160  $s_{t+1}$ , instead of directly predicting the next step. Thus the learning objective for the dynamics model  
 161 becomes,  $\Delta s = s_{t+1} - s_t$ .  $\hat{f}_\theta$  outputs the probability distribution of the future state  $p_{s(t+1)}$  from  
 162 which we can sample the future step and its confidence  $\hat{s}, \hat{s}_\sigma = \hat{f}_\theta(s, [\mathbf{a}])$ . Where the confidence  $s_\sigma$   
 163 captures both, epistemic and aleatoric uncertainty.

164  
 165  
 166 **Algorithm 1** MBRL  
 167 

---

 168 Init  $\mathcal{D}$  with one iteration of a random controller  
 169 **for** Iteration  $i = 1$  **to**  $NIterations$  **do**  
 170   Train  $\hat{f}$  given  $\mathcal{D}$   
 171   **for** Time  $t = 0$  **to**  $TaskHorizon$  **do**  
 172     Get  $a_{t:t+H}^*$  from  
 173     Compute  $OptimalTrajectory(s_t, \hat{f})$   
 174     Execute  $a_t^*$  from optimal actions  $a_{t:t+H}^*$   
 175     Record outcome:  $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t^*, s_{t+1}\}$   
 176 

---

**Trajectory Generation.** Each ensemble element outputs the parameters of a normal distribution. To generate trajectories,  $P$  particles are created from the current state,  $s_t^p = s_t$ , which are then propagated by:  $s_{t+1}^p \sim \hat{f}_b(s_t^p, a_t)$ , using a particular bootstrap element  $b \in \{1, \dots, B\}$ . Chua et al experimented with diverse methods to propagate particles through the ensemble. The  $TS_\infty$  method delivered the best results. It refers to particles never changing the initial bootstrap element. Doing so, results in having both uncertainties separated at the end of the trajectory. Specifically, aleatoric state variance is the average variance of particles of same bootstrap, whilst epistemic state variance is the variance of the average of

177 particles of same bootstrap indexes. We use also  $TS_\infty$ .

179 **Control.** To select the best course of action leading to  $s_H$ , MBRL generates a large number of  
 180 trajectories  $K$  and evaluates them in terms of reward. To find the actions that maximize reward, we  
 181 used the cross entropy method (CEM) Botev et al. (2013), an algorithm for solving optimization  
 182 problems based on cross-entropy minimization. CEM gradually changes the sampling distribution of  
 183 the random search so that the rare-event is more likely to occur and estimates a sequence of sampling  
 184 distributions that converges to a distribution with probability mass concentrated in a region of near-  
 185 optimal solutions. Appendix A details the use of CEM to get the optimal sequence of actions  $a_{t:t+H}^*$

## 186 4 TOWARDS CONTINUAL LEARNING

187 Applying MBRL to a continual learning setting is a promising venue for research. The dynam-  
 188 ics model could be constantly improving and adapting dynamically to changes in the environment.  
 189 Many real-world tasks can be broken in sequences of subtasks of arbitrary length. Capturing the  
 190 complete dynamics then requires exposure to longer sessions of continual learning. Arbitrarily long  
 191 repetitive tasks lead to increasing redundancy in the experience replay constantly increasing of the  
 192 amount of experience collected. These issues hinder the use of MRBL in continual learning settings.  
 193 **What to add to the replay buffer:** We posit that it would be preferable to retain only those experi-  
 194 ences that could not be adequately anticipated by the model during each episode in the environment.  
 195 Essentially, we would only like to add to the replay buffer observations for which the model issued  
 196 a poor prediction. On the contrary, we would like to avoid filling the replay buffer or updating the  
 197 model on observations that the model is good at predicting. We contend that these two elements will  
 198 lead eventually to a balanced replay buffer, which will contain only relevant observations. This will  
 199 contribute to the objective of continual learning.

## 200 5 UARF: UNCERTAINTY AWARE REPLAY FILTERING

201 Continual learning requires the MBRL agent to adapt in a non-stationary environment, retaining  
 202 memories that are useful whilst avoiding catastrophic forgetting, and it can manage compute and  
 203 memory resources over a long period of time ( Khetarpal et al. (2020)). The proposed method,  
 204 UARF, addresses these issues with a variety of strategies. Algorithm 2 is the main algorithm used to  
 205 select which observations to append in the replay buffer. The optimal actions  $a_{t:t+H}^*$  are computed

206 by the *ComputeOptimalTrajectory* function (See Appendix A) given the current state of the  
 207 environment  $s_t$  and  $\hat{f}$ . The future trajectory and its uncertainty,  $p_r^*(t+1:t+1+H)$ , is then obtained by  
 208 using  $a_{t:t+H}^*$  and  $s_t$  with  $\hat{f}$ . The variable *unreliableModel* is set to true when the algorithm believes  
 209 the imagined trajectory not to be trustworthy. Depending on its value, calculation of new trajectories  
 210 and additions to the replay buffer will be avoided and therefore computation time and size of the  
 211 replay buffer will be reduced. If *unreliableModel* is False, the next predicted action is executed in  
 212 the environment. Subsequent actions from  $a_{t:t+H}^*$  are executed until the *unreliableModel* flag is  
 213 set to False or the environment reaches *TaskHorizon* number of steps. The process is repeated for  
 214 the maximum iterations allowed for the task. After the first action, every time an action  $a_{t+1:t+H}^*$   
 215 is executed trajectory computation is avoided and this new observation is *not* added to the replay  
 216 buffer on the basis that the model can already predict its outcome. If *unreliableModel* is True, the  
 217 algorithm calculates a new trajectory and adds the current observation to the replay buffer. Hereby,  
 218 the buffer stores only observations for which the model could not predict (*imagine*) the outcome.

219 **Trustworthy imagination (Algorithm 2 L:18-21).** The algorithm that assigns a value to  
 220 *unreliableModel* is named BICHO. BICHO will essentially return True as long as the reward  
 221 projected in the future does not differ significantly with respect to the imagined future reward  $p_r^*$   
 222 and the confidence of the model remains high. BICHO is built assuming that if parts of the trajec-  
 223 tory do not vary, their projected reward will be as imagined by the model with some confidence.  
 224 After calculating a trajectory, the distribution of rewards  $p_r^*$  is calculated for H steps in the future.  
 225 Whereas, at each step of the environment, independent if the recalculation was skipped or not, a  
 226 new trajectory  $p_r'$  of H steps is projected, starting from state  $s_t$  which is given by the environment  
 227 and using actions  $a_{t+i}^*$  in the imagined trajectory. We use the Wasserstein distance (Bellemare et al.  
 228 (2017)) to find how much these two distributions change after each time step in the environment.  
 229 If the change is  $> \beta$  (which is a hyper parameter to tune) then *unreliableModel* is True. We can  
 230 control how many steps ahead we would like to compare the two distributions. The comparison is  
 231 done for just  $c$  steps ( $< H$ ), which is a hyper parameter to tune. If they differ significantly, then the  
 232 trajectory is unreliable. That is, if the projected reward differs from the imagined one the outcome  
 233 of the actions is uncertain and the trajectory should be recalculated.

234 Even for a model that has converged, predicting trajectories of great length is impossible. Recalculations  
 235 inevitably occur at the end of trajectories. Such recalculations do not necessarily represent  
 236 the appearance of unseen information, but rather a limitation of the successful model in a complex  
 237 environment. Hence, we would not want to add them to the buffer. The *maximum prediction dis-*  
 238 *tance* (MPD) defines a cutoff for a trajectory, and adjusts the strictness of the filtering mechanism.  
 239 Refer to Appendix E for an extensive analysis.

240 **Updates on novel information (Algorithm 2 L:24-25)** over-training the dynamics model leads to  
 241 instabilities due to overfitting. This problem is exacerbated when the replay buffer contains just the  
 242 minimum essential data. If we only filter the replay buffer, continuously updating the parameters of  
 243 the dynamics model will eventually lead to overfitting. Instead, our method updates the parameters  
 244 of the dynamics model only when there is sufficient new information in the replay buffer. We train  
 245 the dynamics model only when new data exceeds the *new\_data\_threshold* hyper parameter. For  
 246 our experiments we set this variable to 0.01 training only when 1% of the experiences in the replay  
 247 buffer are new since the last update of the parameters of the dynamics model.

## 248 6 EXPERIMENTS

249 The primary purpose of the proposed algorithm is for the resulting replay buffer to retain only  
 250 relevant, non-redundant, experiences that will be useful for learning the task. We envision applying  
 251 this method to tasks that require longer training sessions and in continual learning settings.

252 We designed three experimental procedures. The first experiment seeks to establish that our method  
 253 indeed retains a reduced buffer sufficient for achieving expected rewards when learning a single  
 254 task throughout long training sessions. To this end, we evaluate the proposed method in benchmark  
 255 environments for higher number of episodes than in Chua et al. (2018). The second experiment  
 256 seeks to prove that UARF retains a small number of complementary experiences compared to non-  
 257 filtering baseline algorithms when training on a sequence of different but related tasks in a continual  
 258 learning setting. We evaluate our method in a combined task including unseen subtasks. The third  
 259 experiment seeks to show how UARF addresses the effects of catastrophic forgetting.

**Algorithm 2** UARF

```

1: Initialize dynamics model  $\hat{f}$  parameters; Initialize replay buffer  $\mathcal{D}$  with an iteration of a random controller
2:  $unreliableModel = True$  and  $trainModel = False$ 
3: for Iteration  $l = 1$  to  $NIterations$  do
4:   if  $trainModel$  then Train  $\hat{f}$  given  $\mathcal{D}$ 
5:   for Time  $t = 0$  to  $TaskHorizon$  do
6:     if  $unreliableModel$  then
7:       Get  $a_{t:t+H}^*$  from  $ComputeOptimalTrajectory(s_t, \hat{f})$ 
8:       Get  $p_{r(t+1:t+H)}^*$  given  $(s_t, \hat{f}, a_{t:t+H}^*)$  // Use  $\hat{f}$  to predict H rewards ahead
9:        $i = 0$ 
10:    else  $i += 1$ 
11:    Get first action  $a_t$  from available optimal actions  $a_{t:t+H}^*$ 
12:    Execute in the environment  $a_t$  to obtain  $s_{t+1}$  and  $r_{t+1}$ 
13:    Discard first action and keep the rest  $a_t^* = a_{t+1:t+H}^*$ 
14:    Get  $p'_{r(t+i+1:t+H)}$  given  $(s_t, \hat{f}, a_{t+i:t+H}^*)$ 
15:    // Trustworthy imagination
16:     $L = \min(H, c - i)$  // Calculate the number steps ahead to consider
17:     $r\_error = WassersteinDistance(p'_{r(t+i+1:t+i+L)} || p_{r(t+1:t+L)}^*)$ 
18:    if  $r\_error > \beta$  then  $unreliableModel = TRUE$  else  $unreliableModel = False$ 
19:    if  $unreliableModel$  then
20:      Record outcome:  $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, s_{t+1}\}$ 
21:      // Updates on novel information
22:      if  $new\_data\_in\ \mathcal{D} > new\_data\_threshold * length(\mathcal{D})$  then
23:         $trainModel = True$ 

```

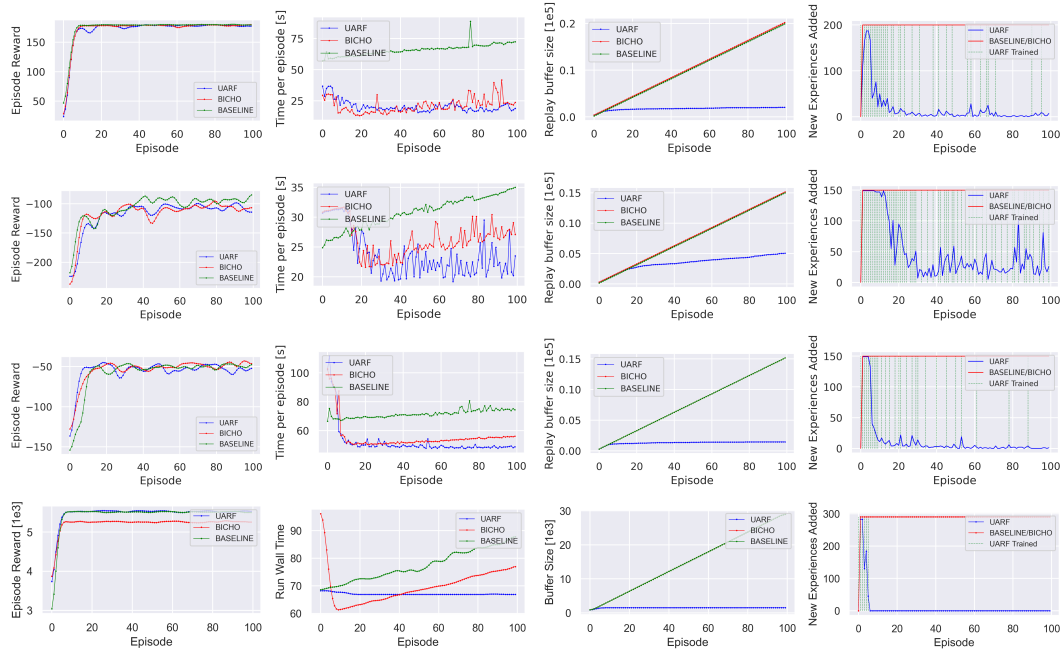


Figure 1: Performance of algorithms (BL: green, BICHO: red, UARF: blue) in (top to bottom) Cartpole, Pusher, Reacher, and Masspoint sector1. From left to right column: episode reward, time per episode (s), cumulative number of observations stored in the replay buffer, new experiences added to the buffer per episode. The rightmost plots illustrate with dashed vertical lines episodes that resulted in UARF updating its model parameters.

260 6.1 E1- CONTINUING TO LEARN A TASK AFTER CONVERGENCE

261 This experiment is intended to show that our method retains sufficient experience to solve the task  
262 while curtailing buffer growth and unnecessary model updates. We intend to prove that this results in

263 a dramatic reduction in the replay buffer size (which is free of any artificially-imposed limits) while  
 264 retaining strong performance (per-episode reward) and reducing per-episode wall clock run-time.

265 We use the MuJoCo (Todorov et al. (2012))  
 266 physics engine and environments Cartpole  
 267 (CP), Pusher (PU) and Reacher (RE) with task  
 268 length ( $TaskH$ ) and trajectory horizon ( $H$ )  
 269 chosen for a valid comparison with Chua et al.  
 270 (2018). With similar training scenarios, Re-  
 271 monda et al. (2021) trained CP for 30 episodes,  
 272 PU and RE for 150. Instead, we trained each  
 273 for 100 episodes. We also included a modified  
 274 version of the Masspoint environment (Thanan-  
 275 jeyan et al. (2020)) (also used in E2). Mass-  
 276 point is a navigation task in which a point mass  
 277 navigates to a given goal. It is a 5-dimensional  
 278  $(x, y, v_x, v_y, \rho)$  state domain. Where  $(x, y)$  is  
 279 the position of the agent,  $(v_x, v_y)$  its speed, and  
 280  $\rho$  is the distance between the agent and the clos-  
 281 est point to a given path. The agent can exert force in cardinal directions and experiences drag coef-  
 282 ficient  $\psi$ . We use  $\psi = 0.6$  and included noise in the starting position. We modified the goal of the  
 283 agent so that it must move as fast as possible without deviating from a given path. Each task and its  
 284 complexity is then determined by the geometry of the path to be followed. The reward is calculated  
 285 as  $r = V(1 - |\rho|)$ . Where  $V$  is the speed of the agent and  $\rho$  the distance to the task’s path. This  
 286 experiment used sector1 (Figure in Appendix B) and Hyperparameters shown in Appendix F.

287 We assess performance in terms of per-episode reward, per-episode wall time, and replay buffer size.  
 288 We evaluate three algorithms: baseline (BL) is a conventional MBRL (PETS Chua et al. (2018)),  
 289 BICHO uses functionality to avoid unnecessary computation, and UARF. BICHO and UARF used  
 290 the same values of  $\beta$  and look-ahead, estimated empirically to produce a reasonable balance in  
 291 terms of per-episode reward and percentage of recalculation. All experiments use random seeds and  
 292 randomized initial conditions for each run, and ran in workstations with Nvidia 3080TI GPUs.

293 **Results:** Fig 1 top shows the results obtained in CP. Fig 1-mid-right shows the size of the replay  
 294 buffer during training. We observe that while the replay buffer keeps grows in the case of BL and  
 295 BICHO, the size of the buffer derived from UARF is comparably flat: the buffer resulting from  
 296 UARF is 10x smaller. The training time per episode (Fig 1 mid-left) remains nearly constant and  
 297 lower for UARF. BL takes substantially longer than both BICHO and UARF to complete an episode.  
 298 The wall time of both the BL and BICHO exhibit linear growth. It takes longer to update the model  
 299 as the replay buffer grows linearly. Fig 1-left shows comparable reward per episode for all methods.  
 300 Results in Fig 1 for PU (row 2), RE (row 3) and Masspoint (bottom) are consistent with those of CP.  
 301 Fig 2 illustrates the management of buffer growth in Masspoint by showing exactly at which steps  
 302 experiences are added to the replay buffer during untrained (E1) and trained (E99) episodes. These  
 303 plots reveal that when the model is untrained, many experiences are added to the buffer throughout  
 304 the episode. After the model is trained (E99), UARF stops adding experiences to the buffer as the  
 305 model is able to predict them. Hence, new experiences are deemed redundant and not useful to  
 306 the model. Results support our claim that UARF obtains a drastically smaller replay buffer that is  
 307 intelligently populated with only relevant information. This is achieved while maintaining strong  
 308 performance in the environment compared to BL. Note that while the curves are plotted per episode,  
 309 it is misleading to assume that all methods converge roughly at the same time. The time per episode  
 310 in the case of UARF is at least half that of BL for approximately in every environment and it remains  
 311 stable, while it increases linearly for BL and BICHO with increasing buffer size. The total wall time  
 312 in average for CP was BL=1.83h, BICHO=0.59h and UARF=0.57h. For PU, it was BL=0.97h,  
 313 BICHO=0.73h and UARF=0.66h. For RE, it was BL=1.99h BICHO=1.55h and UARF=1.45h, and  
 314 for MP it was BL=2.15h, BICHO=1.94h and UARF=1.68h.

## 315 6.2 EX 2. CONTINUAL LEARNING EXPERIMENT

316 This experiment is set in Task Agnostic Continual Reinforcement Learning (the model is not aware  
 317 of tasks or task transitions). In this setting, we sought to prove that UARF maintains a leaner and

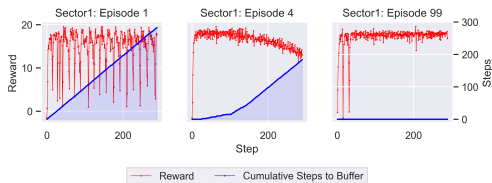


Figure 2: Per-step reward and cumulative steps added to the replay buffer for episodes 1 (left), 4 (middle), and 99 (right) in the Masspoint Sector 1 maneuver. These plots show that UARF adds fewer redundant experiences to the replay buffer as the model converges.

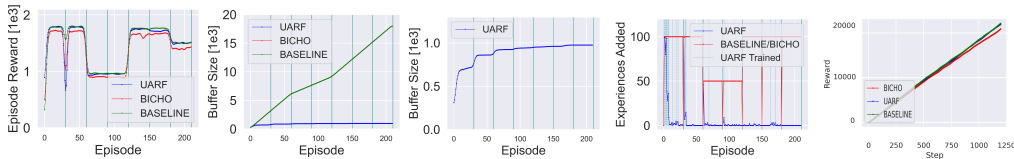
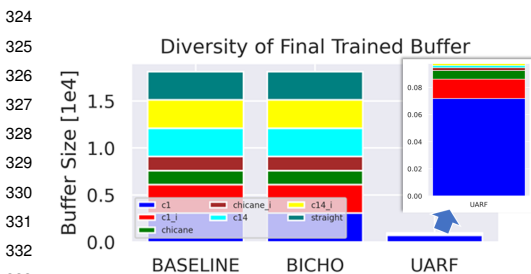


Figure 3: Training performance of Baseline, BICHO and UARF in a Task Agnostic Continual Learning. Models trained on seven maneuvers: corner1, corner1 inverse, chicane, chicane inverse, corner14, corner14 inverse, and straight. Vertical lines indicate a task switch. In the middle-right plot, cyan vertical lines also indicate when UARF triggers a model update. Full-track evaluation in the far-right plot; cumulative reward achieved at each step on the full-track without further training.

318 more relevant collection of experiences in the replay buffer than do baseline algorithms. These char-  
 319 acteristics of the proposed algorithm, we posit, result in strong test performance with less data and  
 320 greater stability. The existence of these characteristics can be verified by observing (after training)  
 321 the size of the buffer, the number of experiences from each maneuver present in the buffer, and the  
 322 performance of the models on the test task. We used the Masspoint racing environment, defining  
 323 different simple tasks that can be composed to solve a complex, unseen one.



324  
 325  
 326  
 327  
 328  
 329  
 330  
 331  
 332  
 333  
 334  
 335  
 336  
 337  
 338  
 339  
 340  
 Figure 4: Distribution of experiences from each sub-task in the replay buffers of each algorithm immediately following training. Detail shows a zoomed-in version for UARF.

Each algorithm trains a model on a sequence of seven separate sub-tasks: two corners and their in-  
 verses, a chicane and its inverse, and a straight (details illustrated in Appendix B). The models retain  
 their parameters and replay buffers between training on each task individually. After training on the last  
 task, the methods are each tested on the full track, which contains some of the sub-tasks seen during  
 training (colored in the full-track image, Appendix B) and tasks unseen during training (shown in black  
 in the full-track image). The model must remember what it learned by training on each sub-task and ap-  
 ply this knowledge to navigate a more complex, unseen task. All of the algorithms had a virtually un-  
 limited replay buffer size. Each model was trained for 30 episodes on each sub-task and then tested on  
 the test task.

341 **Results** Figure 3 shows episode reward, wall-time, buffer size during training, and new experiences  
 342 added to the buffer per episode. Vertical lines illustrate task divisions. High episode reward indicates  
 343 that each model adequately learns each subtask. UARF maintains almost a constant wall-time, while  
 344 BL and BICHO increase as experience accumulates. Buffer growth for BL and BICHO is linear, but  
 345 UARF evidences asymptotic growth (13x smaller) adding no new experiences at the end of training.  
 346 Figure 3-4 shows the buffer growth of UARF. A larger amount of additions to the replay buffer  
 347 occur while training the first tasks. Growth slows to a near halt during the last tasks. This is the  
 348 case for example with the fourth task (chicane inverted). The previous task (chicane) is similar, and  
 349 the information to solve the previous task is enough that the algorithm does not require a significant  
 350 amount of new experience to solve chicane inverted. Figure 4 shows the distribution of experiences  
 351 from each sub-task present in each algorithm’s replay buffer immediately following training. BL  
 352 and BICHO employ a naive approach, resulting in replay buffers with distributions of experience  
 353 determined exclusively by the length of the various maneuvers. The filtering mechanism of UARF  
 354 results in a distribution of experience with some maneuvers having limited representation (e.g., the  
 355 inverse maneuvers) This is because the UARF algorithm intelligently decides to omit redundant  
 356 experiences from the buffer and leaves only the relevant ones. Figure 3 right shows that all three  
 357 algorithms result in a model that adequately solves the test task. UARF continues to manage buffer  
 358 growth while achieving high performance. The results support our initial hypothesis by illustrating  
 359 clearly the proposed algorithm’s propensity to maintain a smaller and more relevant replay buffer  
 360 while achieving the performance of the baseline in a continual learning setting.



## 361 6.3 EX 3. CATASTROPHIC FORGETTING

362 Our approach helps to mitigate catastrophic forgetting. When using a fixed replay buffer size, it is  
 363 important to ensure that the appropriate maximum buffer size is chosen (Zhang & Sutton (2017)).  
 364 If this value is undertuned, important experiences can be jettisoned, and catastrophic forgetting can  
 365 occur. To illustrate how UARF helps to alleviate this risk, we ran the same experiment shown in  
 366 section Ex.2 but with a replay buffer of fixed size (5000 samples; roughly 4x the replay buffer size  
 367 used by UARF in the unlimited size setting). Table 1 compares rewards achieved by each algorithm  
 368 with both unlimited and fixed buffers. The models were validated on the full track and also on a  
 369 maneuver that was trained early on in the training process (c1 inverse). Results reveal that with  
 370 an undertuned fixed buffer size, BL loses about 10% performance both on the full track and on c1  
 371 inverse. This is indicative of the fact that the non-filtering algorithms are hitting the buffer size  
 372 cap, throwing away valuable experiences, and forgetting how to properly solve maneuvers that were  
 trained early on. This impacts performance on the full track as well.

	Unlimited Buffer Full Track	Fixed Buffer Full Track	Fixed Buffer c1 inverse First Pass	Fixed Buffer c1 inverse Post-Training
BASELINE	22172	20235	1787	1561
UARF	21975	22102	1781	1795

Table 1: Fixed Buffer Experiment. Results demonstrate susceptibility to catastrophic forgetting when not using UARF. The BL forgets previous maneuvers after the FIFO mechanism of the fixed-size replay buffer eliminates experiences from them with an impact of about 10% in reward.

373

## 374 7 DISCUSSION AND CONCLUSION

375 The results in E1 reveal that continuing to run our algorithm in a repetitive environment with re-  
 376 dundant or monotonous actions leads to, in some tasks, no increase in buffer and reduced dynamics  
 377 model updates. This has the consequence of reduced running and training times, while reducing the  
 378 effects of catastrophic forgetting and keeping the replay buffer size to a bare minimum. In E2, a con-  
 379 tinual learning setting, we demonstrated that using our approach leads outcomes with 1/25th of the  
 380 experiences without performance degradation. UARF effectively deals with an unbounded growth  
 381 of the replay buffer, which again reduces training time and instabilities. This effect is accentuated  
 382 when training on a continual learning setting. UARF uses a buffer 43x smaller than the baseline.

383 The replay buffer is an instrument that makes the use of deep neural networks in RL more stable  
 384 and it is an essential part in algorithms such as PETs. Such analyses of replay buffer are scarce. But  
 385 recently, research has turned to analyze the contents and strategies to manage the replay buffer of  
 386 RL agents Fedus et al. (2020), and also in supervised learning Aljundi et al. (2019). We contribute  
 387 to such body of work analyzing and offering strategies to manage growth of replay buffer in model  
 388 based RL. Having managed growth, there are several aspects we would like to turn to in the future:  
 389 i) identifying task boundary from the novelty of experiences, ii) managing what to forget for limited  
 390 size buffers, iii) managing what to remember / refresh when a change in task is evident. All this  
 391 would allow to run agents for arbitrary time without having to deal with size of the buffer and would  
 392 offer promising opportunities for deploying MBRL in a continual learning setting.

393 BICHO could be used to prioritize entries in the RB where the model was uncertain. Indeed, prior-  
 394 itized buffer strategies support the usage of experience once it is in the buffer, but as the authors of  
 395 the PER paper state, strategies for what to add and when (our work) are important open avenues for  
 396 research. We did not explore our methods in environments where the tasks have interfering dynam-  
 397 ics. But, if the dynamics change, poor predictions by the model will result in adding experiences to  
 398 the replay buffer. What happens if interfering tasks occur permanently is an interesting follow up.

399 In summary, we proposed strategies that comply with requirements for continual learning. Our ap-  
 400 proach retains only memories which are useful: it obtains lean and diverse replay buffers capturing  
 401 both common and sporadic experiences with sufficient detail for prediction in longer learning ses-  
 402 sions. Our approach manages compute and memory resources over longer periods: it deals with the  
 403 unbounded growth of the replay buffer, its training time and instability due to catastrophic forgetting.  
 404 These results offer promising opportunities for deploying MBRL in a continual learning setting.

## 405 8 REPRODUCIBILITY STATEMENT

406 To make our experiments reproducible, we provide the source code in the supplementary material.  
 407 We include instructions describing how to run all the experiments and to create the images. We in-  
 408 clude the source code of the proposed algorithms, the MassPoint environment and clear instructions  
 409 showing how to install extra packages and dependencies needed to reproduce our experiments.

## 410 REFERENCES

- 411 Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection  
 412 for online continual learning, 2019.
- 413 Haitham Bou Ammar, Rasul Tutunov, and Eric Eaton. Safe policy search for lifelong reinforcement  
 414 learning with sublinear regret, 2015.
- 415 Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement  
 416 learning. In *International Conference on Machine Learning*, pp. 449–458. PMLR, 2017.
- 417 Zdravko I. Botev, Dirk P. Kroese, Reuven Y. Rubinstein, and Pierre L’Ecuyer. Chapter 3 - the  
 418 cross-entropy method for optimization. In C.R. Rao and Venu Govindaraju (eds.), *Handbook of*  
 419 *Statistics*, volume 31 of *Handbook of Statistics*, pp. 35 – 59. Elsevier, 2013. doi: [https://doi.](https://doi.org/10.1016/B978-0-444-53859-8.00003-5)  
 420 [org/10.1016/B978-0-444-53859-8.00003-5](http://www.sciencedirect.com/science/article/pii/B9780444538598000035). URL [http://www.sciencedirect.com/](http://www.sciencedirect.com/science/article/pii/B9780444538598000035)  
 421 [science/article/pii/B9780444538598000035](http://www.sciencedirect.com/science/article/pii/B9780444538598000035).
- 422 Haitham Bou Ammar and Matthew Taylor. Online multi-task learning for policy gradient methods.  
 423 01 2014.
- 424 E.F. Camacho, C. Bordons, and C.B. Alba. *Model Predictive Control*. Advanced Textbooks in  
 425 Control and Signal Processing. Springer London, 2004. ISBN 9781852336943. URL [https://](https://books.google.at/books?id=Sc1H3f3E8CQC)  
 426 [books.google.at/books?id=Sc1H3f3E8CQC](https://books.google.at/books?id=Sc1H3f3E8CQC).
- 427 Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement  
 428 learning in a handful of trials using probabilistic dynamics models, 2018.
- 429 M. P. Deisenroth, G. Neumann, and J. Peters. 2013. doi: 10.1561/23000000021.
- 430 William Fedus, Prajit Ramachandran, Rishabh Agarwal, Yoshua Bengio, Hugo Larochelle, Mark  
 431 Rowland, and Will Dabney. Revisiting fundamentals of experience replay, 2020.
- 432 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
 433 maximum entropy deep reinforcement learning with a stochastic actor, 2018.
- 434 Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards continual reinforce-  
 435 ment learning: A review and perspectives, 2020.
- 436 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A  
 437 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcom-  
 438 ing catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*,  
 439 114(13):3521–3526, 2017.
- 440 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
 441 uncertainty estimation using deep ensembles, 2016.
- 442 Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching.  
 443 *Machine Learning*, pp. 293–321, 1992.
- 444 Tambet Matiisen, Avital Oliver, Taco Cohen, and John Schulman. Teacher–student curriculum learn-  
 445 ing. *IEEE Transactions on Neural Networks and Learning Systems*, 31(9):3732–3740, 2020. doi:  
 446 10.1109/TNNLS.2019.2934906.
- 447 Anusha Nagabandi, G. Kahn, Ronald S. Fearing, and S. Levine. Neural network dynamics for  
 448 model-based deep reinforcement learning with model-free fine-tuning. *2018 IEEE International*  
 449 *Conference on Robotics and Automation (ICRA)*, pp. 7559–7566, 2018.

- 450 Fabrice Normandin, Florian Golemo, Oleksiy Ostapenko, Pau Rodriguez, Matthew D Riemer, Julio  
451 Hurtado, Khimya Khetarpal, Dominic Zhao, Ryan Lindenberg, Timothée Lesort, et al. Sequoia: A  
452 software framework to unify continual learning research. *arXiv e-prints*, pp. arXiv-2108, 2021.
- 453 Anvil V. Rao. A survey of numerical methods for optimal control. *Advances in the Astronautical  
454 Science*, 135:497–528, 2010.
- 455 Adrian Remonda, Eduardo E. Veas, and Granit Luzhnica. Acting upon imagination: when to trust  
456 imagined trajectories in model based reinforcement learning. *CoRR*, abs/2105.05716, 2021. URL  
457 <https://arxiv.org/abs/2105.05716>.
- 458 Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. In  
459 *ICLR (Poster)*, 2016.
- 460 Brijen Thananjeyan, Ashwin Balakrishna, Ugo Rosolia, Felix Li, Rowan McAllister, Joseph E. Gon-  
461 zalez, Sergey Levine, Francesco Borrelli, and Ken Goldberg. Safety augmented value estimation  
462 from demonstrations (saved): Safe deep model-based rl for sparse cost robotic tasks, 2020.
- 463 S. Thrun. A lifelong learning perspective for mobile robot control. In *Proceedings of IEEE/RSJ In-  
464 ternational Conference on Intelligent Robots and Systems (IROS'94)*, volume 1, pp. 23–30 vol.1,  
465 1994. doi: 10.1109/IROS.1994.407413.
- 466 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.  
467 In *IROS*, pp. 5026–5033. IEEE, 2012. ISBN 978-1-4673-1737-5.
- 468 Grady Williams, Paul Drews, Brian Goldfain, James M. Rehg, and Evangelos A. Theodorou. In-  
469 formation theoretic model predictive control: Theory and applications to autonomous driving,  
470 2017.
- 471 Annie Xie and Chelsea Finn. Lifelong robotic reinforcement learning by retaining experiences.  
472 *arXiv preprint arXiv:2109.09180*, 2021.
- 473 Chen Zeno, Itay Golan, Elad Hoffer, and Daniel Soudry. Task-Agnostic Continual Learning Using  
474 Online Variational Bayes With Fixed-Point Updates. *Neural Computation*, 33(11):3139–3177, 10  
475 2021. ISSN 0899-7667. doi: 10.1162/neco\_a\_01430. URL [https://doi.org/10.1162/  
476 neco\\_a\\_01430](https://doi.org/10.1162/neco_a_01430).
- 477 Shangdong Zhang and Richard S. Sutton. A Deeper Look at Experience Replay. *arXiv e-prints*, art.  
478 arXiv:1712.01275, December 2017.
- 479 Guangxiang Zhu, Minghao Zhang, Honglak Lee, and Chongjie Zhang. Bridging imagination and  
480 reality for model-based deep reinforcement learning. *NeurIPS*, 2020.

## 481 A OPTIMAL TRAJECTORY GENERATION

Algorithm 3 shows the use of CEM to compute the optimal sequence of actions  $a_{t:t+H}^*$ .

---

### Algorithm 3 Compute Optimal Trajectory

---

**Input:**  $s_{init}$ : current state of the environment, dynamics model  $\hat{f}$

- 1: Initialize  $P$  particles,  $s_\tau^p$ , with the initial state,  $s_{init}$
- 2: **for** Actions sampled  $a_{t:t+H} \sim CEM(\cdot)$ , 1 **to**  $CEMSamples$  **do**
- 3: Propagate state particles  $s_\tau^p$  using TS and  $\hat{f}|\{\mathcal{D}, a_{t:t+H}\}$
- 4: Evaluate actions as  $\sum_{\tau=t}^{t+H} \frac{1}{P} \sum_{p=1}^P r(s_\tau^p, a_\tau)$
- 5: Update CEM(.) distribution
- 6: **return**  $a_{t:t+H}^*$

---

## 483 B MASS POINT TASKS

484 Each algorithm trains a model on a sequence of seven separate sub-tasks: two corners and their  
 485 inverses, a chicane and its inverse, and a straight (Figure 5). The full track contains some of the sub-  
 486 tasks seen during training (Shown with different colors in the full-track image (Appendix Figure 5))  
 487 in addition to tasks unseen during training (shown in black in the full-track image).

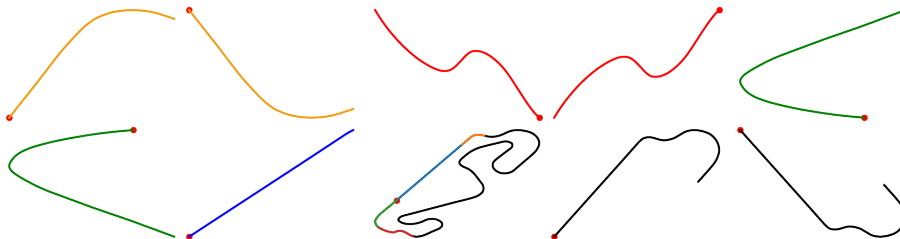


Figure 5: Tasks for the Masspoint environment. The x-axis and the y-axis of each figure represents the x,y coordinates of the path to be followed by the mass point bot. The red dot represents the starting point. Top left-to-right: c1, c1 inverted, chicane, chicane inverted and c14. Bottom left-to-right: c14 inverted, straight, full track (comprising sub-tasks. chicane, c14, straight, c1), sector1 and sector1 inverted

## 488 C ENVIRONMENTS

489 We evaluate the methods on agents in the MuJoCo Todorov et al. (2012) physics engine. To establish  
 490 a valid comparison with Chua et al. (2018) we use four environments with corresponding task length  
 491 ( $TaskH$ ) and trajectory horizon ( $H$ ).

- 492 • Cartpole (CP):  $S \in \mathbb{R}^4, A \in \mathbb{R}^1, TaskH$  200,  $H$  25
- 493 • Reacher (RE):  $S \in \mathbb{R}^{17}, A \in \mathbb{R}^7, TaskH$  150,  $H$  25
- 494 • Pusher (PU):  $S \in \mathbb{R}^{20}, A \in \mathbb{R}^7, TaskH$  150,  $H$  25
- 495 • Masspoint:  $S \in \mathbb{R}^5, A \in \mathbb{R}^2, TaskH$  290,  $H$  25

496 This means that each iteration will run for  $TaskH$ , task horizon, steps, and that imagined trajectories  
 497 include  $H$  trajectory horizon steps.  $S \in \mathbb{R}^i, A \in \mathbb{R}^j$  refers to the dimensions of the environment  
 498 state consisting in a vector of  $i$  components and the action consisting in a vector of  $j$  components.

## 499 D EX 2. CONTINUAL LEARNING EXPERIMENT. ADDITIONAL RESULTS

500 Figure 6 shows additional results with the wall-time during the training process for the continual  
 501 learning experiment.

## 502 E MAXIMUM PREDICTION DISTANCE

503 An additional parameter of interest when using UARF is what we call the "maximum prediction  
 504 distance" or MPD. This parameter operates on the assumption that even for a model that has reached  
 505 convergence, in some environments, predicting trajectories of great length is impossible. As such,  
 506 recalculations must inevitably occur at the end of such long trajectories. These recalculations do  
 507 not necessarily represent the appearance of new, unseen information, but rather a limitation of the  
 508 successful model in a complex environment. Hence, we would not want to add these experiences to  
 509 the buffer.

510 Where we define the cutoff for a trajectory of "great length" can be changed, and it serves to adjust  
 511 the strictness of UARF's filtering mechanism. For Ex.1 and Ex.2, we chose to set the maximum  
 512 prediction distance to 1 to ensure the strictest filtering of the replay buffer.

513 In 7, we evaluate the effect of the MPD on the  
 514 performance of UARF in the cartpole environ-  
 515 ment. We were particularly interested in the  
 516 effect on the rate of recalculation and on the  
 517 size of the replay buffer. In 7 one can see that  
 518 the models converge with no issue, but they do  
 519 differ slightly in the rates of recalculation and  
 520 buffer filtering. The strictest MPD, MPD=1,  
 521 results in the leanest buffer, but its recalcula-  
 522 tion rate is slightly higher than the models with  
 523 MPD=2 and MPD4.

524 These results show that the MPD serves as a  
 525 way to tune the strictness of UARF’s buffer fil-  
 526 tering mechanism. It would be an area of future  
 527 research to find the optimal way to tune this pa-  
 528 rameter automatically throughout training such  
 529 as to best balance recalculation rate and replay  
 530 buffer filtering.

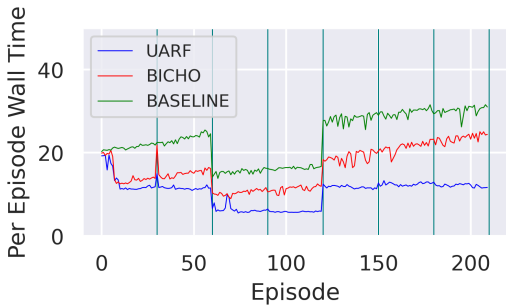


Figure 6: Per episode wall time for the three meth-  
 ods during the training process of Ex.2. Vertical  
 lines indicate task switch points.

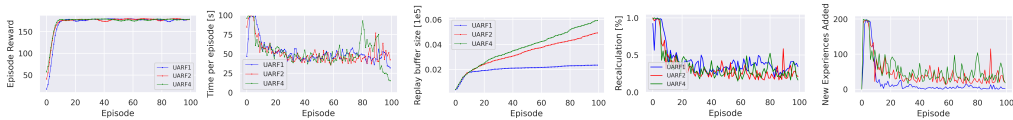


Figure 7: Performance of the examined algorithms in Cartpole using different maximum prediction distances (MPD). The blue line represents UARF with an MPD=1. The red line is UARF with an MPD=2. The green line is UARF with an MPD=4. From left to right column: episode reward, time per episode (s), cumulative number of observations stored in the replay buffer, new experiences added to the buffer per episode.

531 F HYPERPARAMETERS

532 Table 2 shows the hyper parameters used to train UARF. Look-ahead refers to the number of steps  
 533 ahead BICHO and UARF are using to asses the quality of the imagined trajectories.  $\beta$  controls the  
 534 sensitivity of BICHO and UARF to inform whether a trajectory is still valid or not. "New Data  
 535 Train Threshold" refers to the amount of fresh data that must be added to the replay buffer before  
 536 the UARF algorithm triggers the training of the dynamics model.

	Cartpole	Pusher	Reacher	Masspoint
Look-Ahead	10	10	10	10
$\beta$	0.005	0.005	0.005	0.5
New Data Threshold	1%	1%	1%	1%
Training episodes	100	100	10	30/task
CEM population	400	500	400	400
CEM # elites	40	50	40	40
CEM # iterations	5	5	5	5
CEM $\alpha$	0.1	0.1	0.1	0.1
MPD	10	10	10	1

Table 2: Hyperparameters used for UARF implementation.