
A Demonstration is Worth a Million Exploration Steps?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep reinforcement learning has proven an effective method to solve many intricate
2 tasks, yet it still struggles in data efficiency and generalization to novel scenarios.
3 Recent approaches to deal with this include (1) unsupervised pretraining of the
4 agent in an environment without reward signals, and (2) training the agent using
5 offline data coming from various possible sources. In this paper we propose
6 to consider both of these approaches together and argue that this results in a
7 more realistic setting where different types of data are available, and fast online
8 adaptation to new tasks is required. Towards this goal we extend the Unsupervised
9 RL Benchmark to include access to both unsupervised exploratory data, and offline
10 expert demonstrations, when testing the agents online performance on novel tasks
11 in the environment. Using this setup we solve unaddressed issues in previous work.
12 Specifically, we show how to make unsupervised data more effective by using a
13 reward predictor that is trained from a small amount of supervised offline and online
14 data, and we demonstrate how world models can serve as a way to consolidate
15 agent training from various types of data, leading to faster online adaptation.

16 1 Introduction

17 Deep reinforcement learning (RL) has achieved remarkable success in addressing complex control
18 tasks, yet a significant challenge persists in its limited ability to generalize to novel tasks. While
19 transfer learning practices excel in simpler supervised tasks, RL traditionally treats each task in
20 isolation, hindering the utilization of knowledge from prior experiences.

21 This creates two related big challenges: first, this setup is extremely inefficient because of the vast
22 compute that is necessary to train policies from scratch, when even small changes in the task requires
23 the agent to expose itself to millions of experiences in the environment; and second, this results in
24 brittle policies that fail when facing perturbations to the task or environment.

25 Two recent approaches try to deal with this problem, namely unsupervised pretraining and using offline
26 data. In the first approach, an agent is first pretrained using interactions in a similar environment, but
27 without access to a specific task and its reward function. The assumption is that this can lead to a
28 better initialization of the network weights for the second phase, when the agent starts to receive a
29 reward signal and uses it to finetune its weights towards the desired task.

30 In the second approach, the agent first acquires offline data that contains demonstrations of interactions
31 with the environment and the task reward. The agent can then use this data in different ways to extract
32 the optimal policy for the task at hand.

33 In this paper we propose to unify these two approaches. In our setup an agent has unlimited access to
34 unsupervised data in the environment, and in addition has access to a small amount of expert offline
35 data. The main question we ask is how can the agent use this data to learn a new task faster. We

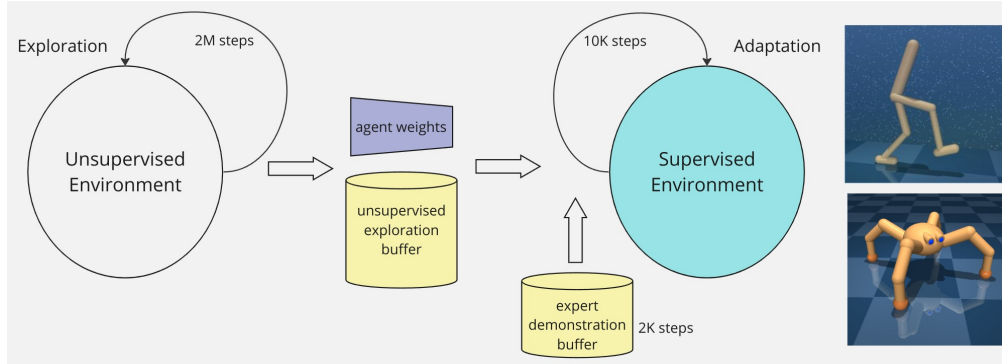


Figure 1: The proposed setup, and the two environments we use.

36 argue that (1) it is better to treat unsupervised data as offline data rather than a method to pretrain the
 37 agent’s weights, and (2) rather than measuring the ability to learn from offline data alone, it is better
 38 to measure its usefulness by the acceleration it provides to subsequent online reinforcement learning.

39 Recent work [18] has already shown a step in the direction of unifying offline and unsupervised
 40 data, as it shows that for offline RL, using data that was acquired through unsupervised exploration
 41 has benefits over using the data from some reward maximizing policy. However the setting is not
 42 completely realistic as it assumes the ground truth reward function is known after the data was
 43 collected. Here we show how to get away from this assumption, by training a reward predictor using
 44 a few supervised offline expert demonstrations, and a small amount of supervised online interactions.

45 To further explore possible methods for marrying unsupervised and offline data, we consider world
 46 models. A recent paper [13] showed that world models are effective in an unsupervised pretraining
 47 setup, achieving excellent results in the Unsupervised RL Benchmark (URLB) [8]. We claim that the
 48 reason for this is that unsupervised pretraining is best used as a source of exploratory data rather than
 49 a way to achieve better initialization of the policy weights. By extending the method in Rajeswar
 50 et al. [13], we show significant improvements in adapting to novel tasks in an order of magnitude
 51 less online steps. We also show that access to a few expert demonstrations further accelerates online
 52 training in the URLB.

53 Following our findings, we propose this setup as an extension to the Unsupervised RL Benchmark that
 54 studies the unified space of using unsupervised exploratory data, supervised expert demonstrations,
 55 and supervised online interactions. We believe this is a natural setup to start benchmarking, as it
 56 represents the realistic setting where different possible data sources are available, and extremely
 57 rapid adaptation is required. Furthermore, our results on this setup suggests that a world model is a
 58 convenient method to fuse different data sources, leading to a method that effectively consolidates
 59 policy training using unsupervised exploratory data, expert demonstrations and online interactions.

60 Our main contributions can be summarized as following:

- 61 1. We propose to study a unified problem dealing with unsupervised pretraining, offline expert
 62 demonstrations, and online interactions in the environment. Towards this goal we publish an
 63 extension to the Unsupervised RL Benchmark, that can serve as an initial test-bed for this setup.
- 64 2. Compared to previous work on unsupervised offline RL [18], we drop the assumption of a known
 65 ground-truth reward function, and propose a way to train a reward predictor using a small amount
 66 of supervised offline and online data. Furthermore, in contrast to previous work on state space
 67 only, we show results in pixel space, making the overall setup more realistic.
- 68 3. Compared to previous work on using world models for unsupervised pretraining [13], we improve
 69 the training method, leading to better performance in fast adaptation, and we show how to
 70 incorporate offline data of a few expert demonstrations, leading to an even larger improvement of
 71 the results on the URLB. Our results suggest that world models are an effective tool to merge
 72 different types of data sources for RL, without losing information in the way.

73 All data and code implementing the benchmark and tested methods will be made available upon
 74 publication.

75 2 Preliminaries and Related Work

76 2.1 Unsupervised Reinforcement Learning Benchmark

77 The Unsupervised Reinforcement Learning Benchmark (URLB) Laskin et al. [8] is a benchmark that
78 compares the adaptation capabilities of different RL unsupervised exploration algorithms. The setup
79 is separated into two stages: pretraining and fine-tuning. During the pretraining stage, a policy is
80 trained on an intrinsic reward, which is specified by the tested exploration method. At the second
81 stage, the resulted policy from the previous stage is used for fine-tuning on a specific task from the
82 same domain as the one used during pre-training. The underlining RL algorithm is DrQv2 [17],
83 which is a variant of DDPG [10].

84 In the pretraining stage, the tested exploration methods are separated into knowledge based, data based
85 and competence methods. Knowledge-based methods focus on maximizing the entropy of visited
86 states. Competence methods focus on maximizing the entropy of the learned policy. Data-based
87 methods learn an explicit skill vector w by maximizing the mutual information between the encoded
88 observation and skill. For each exploration method, the pretraining stage is performed environments
89 from the *dm_control* suit: Walker, Quadraped and Jaco. Here we focus on Walker and Quadraped, as
90 they present a more interesting challenge of exploration vs. fast adaptation.

91 In the finetuning stage, the exploration strategy used during pretraining is evaluated on set of tasks.
92 First, the policy weights and part of the critic network are loaded from the pre-training stage. Then,
93 the policy is trained on the task using the baseline RL algorithm DrQv2. The method is evaluated
94 after 100K steps of finetuning, which is used as a measure of adaptation to a novel tasks.

95 2.2 Intrinsic Reward Models

96 Intrinsic reward method are used to facilitate better exploration in sparse reward environments. The
97 idea is to use an additional reward to promote the visiting of novel states. Intrinsic reward methods
98 are specially important when dealing with unsupervised environments, as they form the only source
99 of reward. In this work we focus on four leading methods for exploration as used in Laskin et al. [7]
100 and Rajeswar et al. [13]. The methods we consider are ICM [12], RND [1], RE3 [14], and Latent
101 Bayesian Surprise (LBS) [11]. The latter is used in DreamerMPC [13], on which we base our method.

102 2.3 ExORL

103 ExORL [18] is a different benchmark that tests ways to utilize unsupervised exploration data from
104 the perspective of offline RL. In URLB, the setup builds on pretraining a policy and later finetuning
105 on a specific task. In contrast, ExORL assumes there is access to the reward function, that can be
106 used to fix the unsupervised exploration buffer, which originally does not contain the reward signal.
107 This is done without collecting any additional data. In this work, they argue that using the whole
108 unsupervised buffer as the finetuning data source is more beneficial then using the policy that was
109 used during the unsupervised data collection stage. All this, assuming the reward function is known
110 and is easy to re-compute for each task.

111 2.4 World Model Pretraining

112 A major development in model based RL is the success of learned world models, which allow
113 simulating environment transitions, and can facilitate sample efficient training. A learned world
114 model is a data driven model that can estimate the distribution over the next state given a history of
115 previous transitions:

$$p_{\theta}(s_{t+1}, r_{t+1} | s_t, a_t, s_{t-1}, \dots). \quad (1)$$

116 In essence, a well trained world model can replace a complex simulation of the environment. This,
117 in turn, allows training the agent with almost no interaction with the actual environment, when the
118 underlining environment dynamics stay the same. Additionally, planning techniques can be used by
119 generating possible futures from the current states.

120 An example of such an architecture for world models is Dreamer [3] which uses an RNN to learn
121 the environment dynamics in a sequential way. In addition, it learns the reward model $R(s, a)$. This
122 allows generating partial trajectories, using some policy, and training an off-policy agent on the

123 generated data. Recently, excellent results were achieved on the tasks tested in URLB. The method,
124 which we call here **DreamerMPC** [13], uses the Dreamer architecture as a world model, where
125 the data collected during the pretraining stage is used to train it, and LBS is used as the intrinsic
126 reward. Additionally, during the fine-tuning stage, a planning strategy called Model Predictive
127 Control (MPC) [4] is used to generate trajectories.

128 In MPC, the planning policy consists of a network that predicts the value function $V_\xi(s)$ and Gaussian
129 action sequence distribution estimator $p(a_1, a_2, \dots, a_t | s_0)$. For each state, several planning iterations
130 are done, to gradually find the optimal trajectory and then the first planned action is used to proceed
131 in the actual environment. For each planning iteration, given the sequence distribution, an expected
132 return is estimated using the world model, and the value function. The return estimation is done in
133 the following way:

- 134 1. the state s_0 is received from the environment
- 135 2. an initialized Gaussian sequence distribution $p(a_1, a_2, \dots, a_t | s_0)$ is used to generate a
136 sequence of actions.
- 137 3. the world model $p_\theta(s_{t+1}, r_{t+1} | s_t, a_t)$ is used to predict the reward at each time step until s_t
- 138 4. the rest of the expected return is calculated by $V_\xi(s_{t+1})$

139 Top sampled sequences are used to update the distribution, and another set of sequences are sampled.
140 This process continues for a set amount of iterations. DreamerMPC used a Dreamer based world
141 model and MPC, and showed superior performance on a variety of tasks compared to standard RL
142 policy optimization via SGD.

143 2.5 Offline RL

144 In offline RL, rather than collecting data by interacting with the environment, data is collected by a
145 different, usually unknown method, and used by the agent as an dataset to learn the policy. Data can
146 come from expert demonstrations, historical data, or simulated environments. The aim is to reduce
147 the need of costly interaction with the environment. For an overview of the approach see Levine et al.
148 [9], and for different datasets see Fu et al. [2].

149 One of the main issues with offline RL relates to the ability to generalize from one setup to another.
150 Several methods were developed to address this problem, such as CQL [6], IQL [5] and CRR [16].

151 3 Motivation: Unsupervised vs. Offline Data

152 The setup we use in this paper combines both unsupervised and offline data. We argue that this is a
153 natural setup to consider, as it deals with related issues that arise in recent work, namely URLB [8],
154 ExORL [18] and DreamerMPC [13]. In this section we discuss these issues based on experiments we
155 performed on URLB, which serve as the motivation for the final setup we propose.

156 The results shown in Fig. 2, and Fig. 3, demonstrate different ways in which unsupervised exploration
157 can be used. The standard method, as proposed in the original benchmark of URLB, is to pretrain an
158 agent in the unsupervised environment, using various possible exploration methods, and then use
159 the weights as initial values before finetuning them using a reward function for some given task.
160 This is denoted by **Finetune** in the figures, where the performance is measured for each task after
161 100K environment steps including a reward signal. While finetuning a pretrained agent leads to some
162 acceleration in training time for some of the tasks, it is not always significantly better than initializing
163 the weights randomly and completely discard the output of the pretraining phase (**baseline**).

164 We argue that this happens because the benefit of the pretrained agent does not necessarily come
165 from the value of its weights, but rather from its behavior in the environment in the initial steps. In
166 other words, not because its weights are already closer to their optimal values, but rather because the
167 pretrained agent can start generating useful exploration trajectories from step one. To test this, we
168 use a different method to exploit unsupervised exploration, where the exploration trajectories in the
169 pretraining phase are stored in a buffer, and then used together with the online data while training the
170 agent on a specific task.

171 Essentially this means the pretraining exploration is treated as offline data, however because it
172 comes from unsupervised exploration, the corresponding reward signal is not right. To correct this,

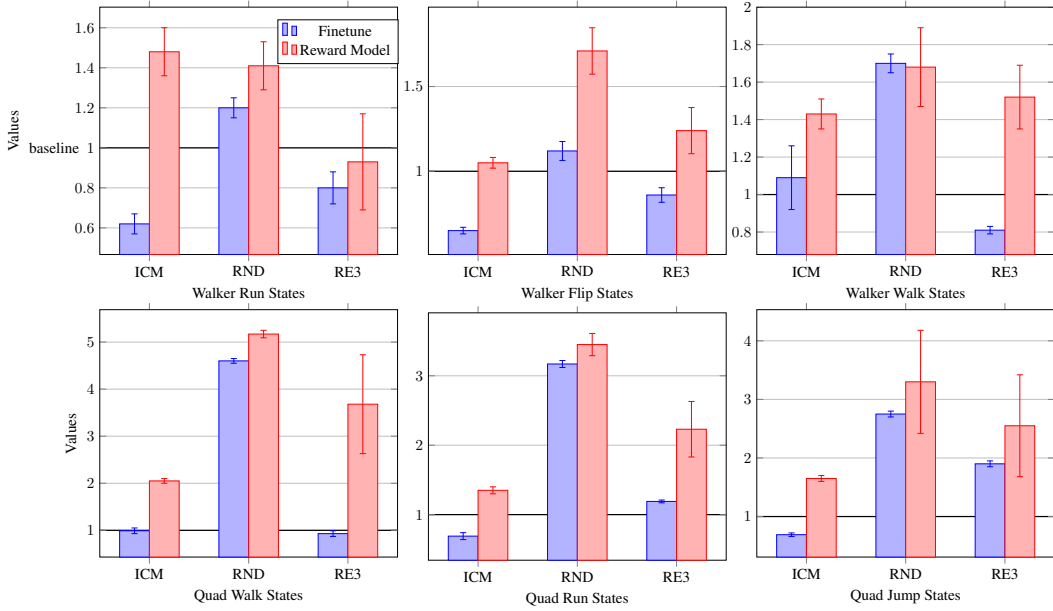


Figure 2: Average episodic returns (over 10 episodes) on the URLB benchmark with **state** observations after 100K training steps. We show two environments and three different tasks for each, and values are normalized to a baseline that does not use any pretraining information. Results show that (1) finetuning a pretrained agent (Finetune) is not always better than training from scratch; and (2) Using the unsupervised pretraining to collect offline data (Reward Model) can lead to better performance. This is implemented by training a reward model to predict the reward of the unsupervised data.

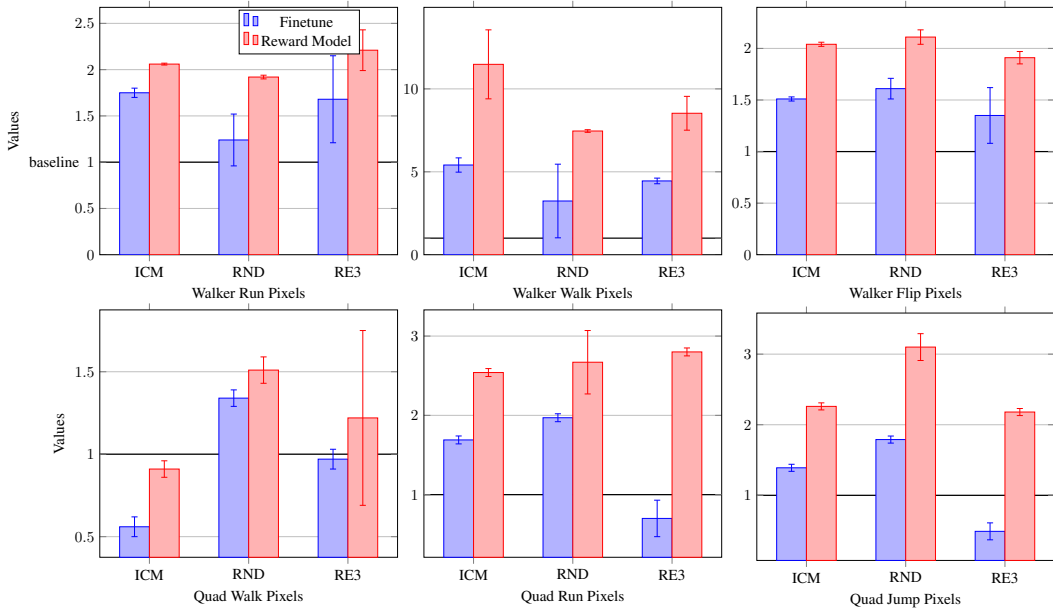


Figure 3: Average episodic returns (over 10 episodes) on the URLB benchmark for **pixel** observations after 100K training steps (normalized to a baseline with no pretraining). The behavior is similar to state observation in Fig. 2, namely (1) finetuning a pretrained agent is not always beneficial (Finetune), and (2) using the unsupervised pretraining as an offline buffer (where the reward signal is predicted by a learned reward model) can lead to better results (Reward Model).

173 we implement a reward-model component, that uses the online supervised data, to learn a reward
 174 predictor. In turn, the reward model can be used to predict the reward in each step of the exploratory
 175 trajectories in the pretrained buffer. The results in Fig. 2 and Fig. 3 show this method (**Reward**
 176 **Model**) outperforms simple finetuning of the agent weights. For more details on the implementation
 177 of this method see the appendix.

178 It is interesting to note the relation between the Reward Model method and two recently proposed
 179 approaches, namely ExORL [18], and DreamerMPC [13] (see Sec. 2 for their description). The
 180 ExORL approach claims that unsupervised exploratory data can serve as better data for offline RL,
 181 provided that the ground truth reward function can be used in the offline training phase. Our proposed
 182 Reward Model method differs from ExORL in two different ways: (1) it does not rely on access
 183 to the true reward function and instead it learns it from the few supervised online steps, and (2) its
 184 performance is measured after a few online steps in the environment rather than using the offline data
 185 only.

186 In DreamerMPC, an agent based on planning with a world model is shown to significantly outperform
 187 all other methods in URLB. While this can be a result of the general efficiency of planning in RL,
 188 we argue that the reasons for the favorable performance of DreamerMPC are related to the results
 189 discussed above, because training a world model on unsupervised exploration can be viewed as a smart
 190 way to store the exploration trajectories. The world model is essentially a buffer of the exploration
 191 trajectories that also allows generalization and planning. We show results with DreamerMPC in later
 192 sections.

193 4 URLB10K: A Benchmark for Fast Adaptation

Table 1: Average episodic return (over 10 episodes) after 10K training steps: we compare learning the reward model from 10K online steps vs. replacing 2K with expert steps (4 episodes). In both cases results are comparable or better than ExORL which uses a GT reward function.

| Task/Method | ExORL | RewardModel | RewardModel+Expert |
|----------------|------------|-------------|--------------------|
| Walker Walk | 260 | 238 | 375 |
| Walker Run | 121 | 96 | 106 |
| Walker Flip | 207 | 236 | 285 |
| Quadruped Walk | 250 | 298 | 402 |
| Quadruped Run | 279 | 299 | 412 |
| Quadruped Jump | 324 | 640 | 630 |

194 Motivated by the results and the points discussed above, we propose a new setup to test fast adaptation
 195 to novel tasks. The idea is to use different sources of information to achieve even faster adaptation
 196 than what is tested in URLB. We base our setup on the following:

- 197 1. The agent has access to a very large number of unsupervised interaction with the environment.
 198 This is similar to the URLB setup, however this can be used either for pretraining an agent,
 199 populating a buffer, or training a world model.
- 200 2. The agent has access to a few expert trajectories, containing the task’s reward. This data
 201 can serve two purposes. First, it can be used as initial information about the task, e.g
 202 to train a reward predictor. Second, it can disentangle the problem of exploration in the
 203 environment, as using this data is somewhat equivalent to a successful exploration that
 204 ensures the important states in the environment were covered.
- 205 3. Given the above data, the agent is assessed on its performance after a small amount of super-
 206 vised online interactions with the environment, measuring its capacity for fast adaptation.

207 Specifically, we use the URLB with the DeepMind Control Suite environments [15], and we propose
 208 to use 2M unsupervised exploration steps and 2K supervised expert demonstration steps (4 episodes),
 209 and then measure the performance of the agent after 10K online supervised steps. This is an order of
 210 magnitude less steps than what is tested in URLB.

211 Compared to URLB, the only addition are 4 episodes of supervised expert demonstration (2K steps)
 212 that we provide for each task. These episodes are achieved by training the DreamerMPC model on

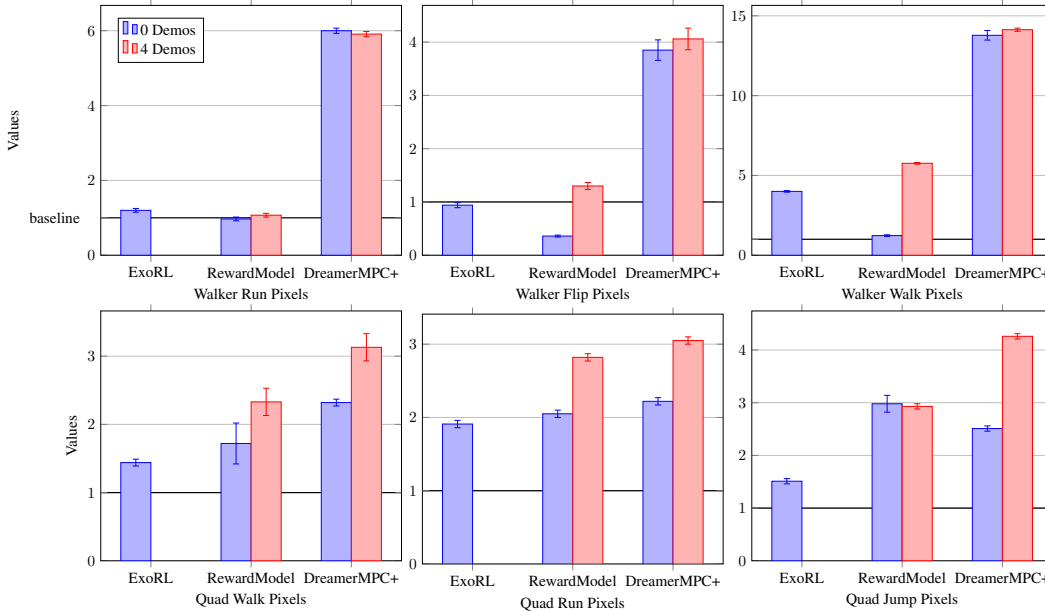


Figure 4: Results for the main methods on URLB10K. Average episodic return (over 10 episodes). ExoRL uses a GT reward function and no online data, RewardModel and DreamerMPC+ are trained on 10K online steps with/without 2K expert demonstrations, and all results are normalized relative to a baseline that was trained on 100K online steps. The results show: (1) Fast adaptation is possible (2) The method based on a world model significantly outperforms other methods. (3) Using expert data is not always beneficial, but can lead to significant improvements.

213 a supervised environment for 2M steps. We use a random 4 episodes from the last 100 episodes of
 214 training.

215 The capacity of fast adaptation can be seen in the results in table 1, presenting episode rewards after
 216 10K supervised steps for two different environments, and three different tasks. For this we
 217 implemented a method that uses a reward model to predict the reward of the unsupervised steps. The
 218 unsupervised steps are collected using LBS exploration [11] (see Sec. 2), and the method then trains
 219 the reward model on data from the online supervised trajectories. We test this method both with and
 220 without extra trajectories coming from the offline expert demonstrations, where the first option uses
 221 8K online steps and 2K offline expert steps, and the second uses 10K online steps only.

222 The method is compared ExoRL [18], which relies on the same unsupervised buffer, but instead
 223 of using offline expert data and online interactions, it fixes the unsupervised reward signal using
 224 privileged knowledge of the ground-truth reward function.

225 The first thing to note in the results in table 1 is the value of replacing 4 episodes from the online
 226 interactions with expert demonstrations, when training the reward model. This generally leads to a
 227 significant improvement, although it is not always the case. The second thing to note is the comparable
 228 and sometimes favorable results of the learned reward model compared to ExoRL. This shows that
 229 in this realistic setup, learning the reward from very few episodes is possible, demonstrating the
 230 potential for developing methods that exhibit extremely fast adaptation. Results can be viewed also in
 231 the bar-plots in figure 4. Details of the implementation of the method can be found in the appendix.

232 5 Using a World Model

233 Following the results of training a reward predictor for URLB10K, a natural next step is to consider a
 234 full world model. While DreamerMPC [13] showed excellent results on the original URLB, we set to
 235 test it on the fast adaptation setup we propose. In addition, we make modifications to the training
 236 method, based on the approach that the model should treat unsupervised exploration as offline data
 237 rather than pretraining data.

238 The modifications we make are: (1) Utilize the exploration buffer directly, sampling batches similarly
 239 to the method of Reward Model; (2) make offline updates after the online steps collection is finished;
 240 (3) increase the environment steps to network update ratio from 10 to 5; and (4) add a few finetuning
 241 updates to the world model on the online episodes and the optional expert episodes. We denote this
 242 modified model as **DreamerMPC+**, and use **DreamerMPC+Expert** when the method is also using
 243 expert demonstrations for training. Additional information on the method implementation can be
 244 found in the appendix.

245 Table 2 shows the results after 10K supervised online steps, or 8K online and 2K offline expert
 246 steps (marked as **Expert**). Comparing the modified method (**DreamerMPC+**) to the vanilla method
 247 (**DreamerMPC**), shows a significant improvement in all tested tasks. Moreover, replacing 4 episodes
 248 with expert demonstrations leads to a further increase in performance for almost all tasks. These
 249 results show that our training method leads a new state of the art on the URLB with or without expert
 250 data, outperforming all previous methods to date.

251 A summary of results on the URLB10K benchmark are presented in table 3 and figure 4, comparing
 252 our methods based on a world model and reward model that use 10K online steps, to ExORL that uses
 253 privileged information of the reward function, and to the baseline of training on 100K online steps.
 254 These results highlight the potential for fast adaptation methods, and demonstrate the superiority of
 255 world models in merging data of various types into an adaptable policy.

Table 2: Fast adaptation in DreamerMPC. Average episodic return (over 10 episodes) after 10K training steps. Our updates to the DreamerMPC model (DreamerMPC+) lead to significant improvements in fast adaptation, and can be even further improved by incorporating 4 expert demonstration episodes (DreamerMPC+Expert).

| Task/Method | DreamerMPC | DreamerMPC+ | DreamerMPC+Expert |
|------------------|------------|-------------|-------------------|
| Walker Walk | 630 | 896 | 918 |
| Walker Run | 328 | 594 | 585 |
| Walker Flip | 607 | 841 | 889 |
| Quadruped Walk | 339 | 401 | 540 |
| Quadruped Run | 200 | 324 | 446 |
| Quadruped Jump | 449 | 539 | 915 |
| Normalized Score | 1.0 | 1.4 | 1.7 |

Table 3: Results for the main methods on URLB10K. Average episodic returns after 10K training steps. DreamerMPC+Expert effectively combines unsupervised data and expert demonstration and significantly outperforms other methods.

| Task/Method | Pix ExORL (GT reward) | RM+Expert | DreamerMPC+Expert |
|----------------|-----------------------|-----------|-------------------|
| Walker Walk | 260 | 375 | 918 |
| Walker Run | 121 | 106 | 585 |
| Walker Flip | 207 | 285 | 889 |
| Quadruped Walk | 250 | 402 | 540 |
| Quadruped Run | 279 | 412 | 446 |
| Quadruped Jump | 324 | 630 | 915 |

256 6 Discussion

257 In this paper we proposed to test fast adaptation to novel tasks, using only 10K online interactions
 258 with the environment. In order to achieve this, we argue that a natural setup is when the agent
 259 has access to a large amount of unsupervised experience, and a small amount of supervised expert
 260 episodes, demonstrating the desired task.

261 Towards this goal we presented an extension to the URLB benchmark, and studied methods that
 262 treat the pretraining data as offline data, specifically using world models. We show this approach
 263 can be related to previous work, and develop training methods that demonstrate the capacity of fast
 264 adaptation in such models.

265 We believe that the proposed setup can serve as an initial benchmark to study and develop methods
 266 that can efficiently process various types of available data, and improve fast adaptation in RL.

267 **References**

- 268 [1] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by Random
269 Network Distillation, October 2018.
- 270 [2] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for
271 deep data-driven reinforcement learning, 2021.
- 272 [3] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering Atari with
273 Discrete World Models, February 2022.
- 274 [4] Nicklas Hansen, Xiaolong Wang, and Hao Su. Temporal Difference Learning for Model
275 Predictive Control, July 2022.
- 276 [5] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit
277 q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
- 278 [6] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning
279 for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:
280 1179–1191, 2020.
- 281 [7] Michael Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas.
282 Reinforcement Learning with Augmented Data, November 2020.
- 283 [8] Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang,
284 Lerrel Pinto, and Pieter Abbeel. URLB: Unsupervised Reinforcement Learning Benchmark,
285 October 2021.
- 286 [9] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:
287 Tutorial, review, and perspectives on open problems, 2020.
- 288 [10] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval
289 Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning,
290 July 2019.
- 291 [11] Pietro Mazzaglia, Ozan Catal, Tim Verbelen, and Bart Dhoedt. Curiosity-Driven Exploration
292 via Latent Bayesian Surprise, February 2022.
- 293 [12] Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven Explora-
294 tion by Self-supervised Prediction, May 2017.
- 295 [13] Sai Rajeswar, Pietro Mazzaglia, Tim Verbelen, Alexandre Piché, Bart Dhoedt, Aaron Courville,
296 and Alexandre Lacoste. Mastering the Unsupervised Reinforcement Learning Benchmark from
297 Pixels, May 2023.
- 298 [14] Younggyo Seo, Lili Chen, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. State
299 Entropy Maximization with Random Encoders for Efficient Exploration, June 2021.
- 300 [15] Saran Tunyasuvunakool, Alistair Muldal, Yotam Doron, Siqi Liu, Steven Bohez, Josh Merel,
301 Tom Erez, Timothy Lillicrap, Nicolas Heess, and Yuval Tassa. dm_control: Software and
302 tasks for continuous control. *Software Impacts*, 6:100022, 2020. ISSN 2665-9638. URL
303 <https://www.sciencedirect.com/science/article/pii/S2665963820300099>.
- 304 [16] Ziyu Wang, Alexander Novikov, Konrad Zolna, Josh S Merel, Jost Tobias Springenberg, Scott E
305 Reed, Bobak Shahriari, Noah Siegel, Caglar Gulcehre, Nicolas Heess, et al. Critic regularized
306 regression. *Advances in Neural Information Processing Systems*, 33:7768–7778, 2020.
- 307 [17] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering Visual Continuous
308 Control: Improved Data-Augmented Reinforcement Learning, July 2021.
- 309 [18] Denis Yarats, David Brandfonbrener, Hao Liu, Michael Laskin, Pieter Abbeel, Alessandro
310 Lazaric, and Lerrel Pinto. Don't Change the Algorithm, Change the Data: Exploratory Data for
311 Offline Reinforcement Learning, April 2022.

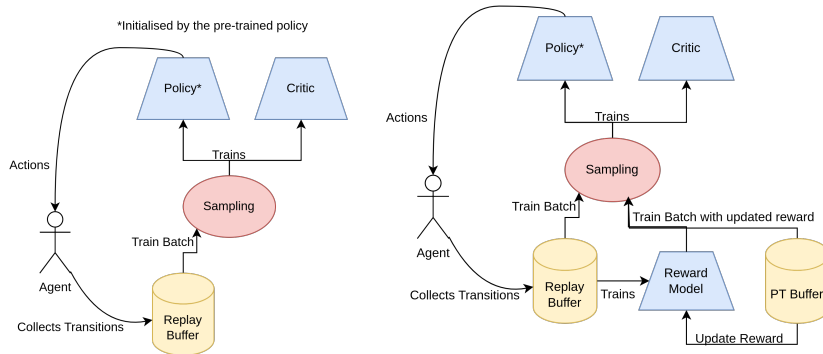


Figure 5: **Left:** Baseline off-policy RL, where we sample only from the replay buffer. **Right:** PT Buffer & Reward Model, where we additionally sample from the PT buffer and predict the task oriented reward using the reward model.

312 A Additional Details of Methods

313 A.1 Reward Model

314 When using the reward model approach, we base the model on the Q network architecture, but instead
 315 of training it on Q values, we train it to output $R(s, a)$. After collecting the seed frames, we start with
 316 200 batch updates to train the model using the optional expert episodes and the seed frames.

317 Next, when starting the online training stage, we intermittently sample online data with the true
 318 reward, and train the world model, or when we sample from the unsupervised buffer we use the
 319 reward model to get the estimated task reward.

320 Training an agent for 10K online steps, including sampling from the exploration buffer takes about 2
 321 hours and about 10 hours for additional 400K offline updates on a single A100 GPU.

322 A.2 DreamerMPC+

323 Here we do the same sampling and seed frames pre-finetuning as done in the reward model. The
 324 difference here is that instead of just learning $R(s, a)$, we learn a world model. The world model
 325 is based on Dreamer and is learned in a supervised fashion from the exploration buffer. During the
 326 task training stage, the world model is finetuned to correct its reward head, as done in [13]. One
 327 difference, is that as we include expert data, we do additional 10 pre-finetuning updates on the world
 328 model just after collecting the seed frames.

329 Training an agent for 10K online steps, including sampling from the exploration buffer takes about 5
 330 hours and about 20 hours for additional 300K offline updates, on a single A100 GPU.

Algorithm 1 Finetuning With a Reward Model on Unsupervised Data

Require: Unsupervised data U , replay buffer D
Require: Reward model R_θ , pretrained policy π_ϕ , initialized critic Q_ξ
Require: number of seed steps N_s , finetune steps N_{ft} , offline updates steps N_{ol}
Require: reward model pre-train steps N_{pt}
Require: offline batches counter $c = 0$

- 1: **for** $t \leftarrow 1$ to N_s **do**
- 2: $D \leftarrow D \cup (s, a, s', r)$
- 3: **end for**
- 4: **for** $t \leftarrow 1$ to N_{pt} **do**
- 5: Update R_θ using mini-batches from D
- 6: **end for**
- 7: **for** $t \leftarrow 1$ to $N_{ft} + N_{ol}$ **do**
- 8: sample batch $b \sim P_{bs}(U, D + c)$
- 9: **if** $b \subset U$ **then** ▷ we sampled from the unsupervised data
- 10: $c \leftarrow c + 1$
- 11: $(s, a, s', r) \leftarrow b$
- 12: $r' \leftarrow R_\theta(s)$
- 13: Update π_ϕ and Q_ξ using the mini-batch (s, a, s', r')
- 14: **continue**
- 15: **else if** $b \subset D$ **then** ▷ we sampled from the replay buffer
- 16: Update π_ϕ and Q_ξ using b
- 17: Update R_θ using b
- 18: **if** $t < N_{ft}$ **then**
- 19: $a' \leftarrow \pi_\phi(s')$ ▷ collect next transition from the environment
- 20: $o, r \sim P(\cdot | s', a')$
- 21: $D \leftarrow D \cup (s', a', o, r)$
- 22: **end if**
- 23: **end if**
- 24: **end for**

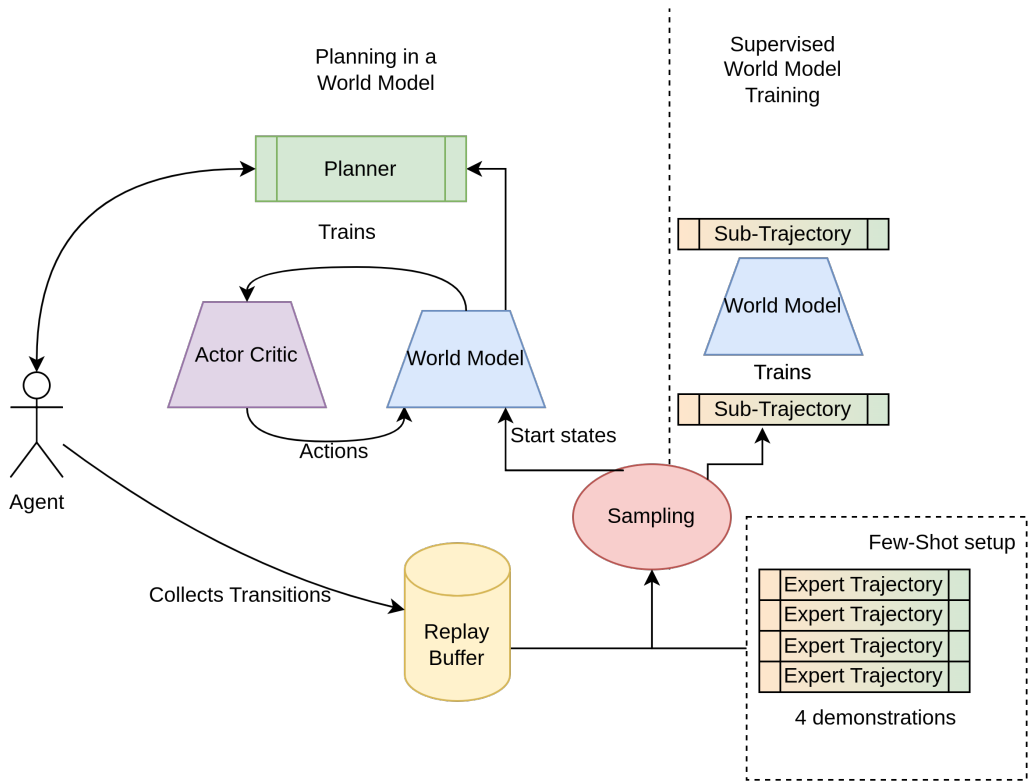


Figure 6: DreamerMPC architecture with additional information from expert demonstration