

---

# Continual Model-Based Reinforcement Learning with Hypernetworks

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

Effective planning in model-based reinforcement learning (MBRL) and model-predictive control (MPC) relies on the accuracy of the learned dynamics model. In many instances of MBRL and MPC, this model is assumed to be stationary and is periodically re-trained from scratch on state transition experience collected from the beginning of environment interactions. This implies that the time required to train the dynamics model – and the pause required between plan executions – grows linearly with the size of the collected experience. We argue that this is too slow for lifelong robot learning and propose HyperCRL, a method that continually learns the encountered dynamics in a sequence of tasks using task-conditional hypernetworks. Our method has three main attributes: first, it enables constant-time dynamics learning sessions between planning and only needs to store the most recent fixed-size portion of the state transition experience; second, it uses fixed-capacity hypernetworks to represent non-stationary and task-aware dynamics; third, it outperforms existing continual learning alternatives that rely on fixed-capacity networks, and does competitively with baselines that remember an ever increasing coresset of past experience. We show that HyperCRL is effective in continual model-based reinforcement learning in robot locomotion and manipulation scenarios, such as tasks involving pushing and door opening.

19 **1 Introduction**

20 Lifelong model-based robot learning is predicated upon continual adaptation to the dynamics of new  
21 tasks. For example, robots need to learn to manipulate unseen objects with various mass distributions,  
22 walk on new types of terrains with different friction, elasticity, and other physical properties,  
23 or even learn to adapt to different tasks, such as walking, running, or climbing stairs. This presents  
24 at least two challenges for many model-based reinforcement learning (MBRL) and model-predictive  
25 control (MPC) formulations, which typically comprise of a dynamics learning phase followed by a  
26 planning/policy optimization and execution phase.

27 First, these methods are not scalable because the time required to train the dynamics model grows  
28 linearly with the size of the collected experience. Second, as the robot learner encounters and  
29 adapts to new tasks, it has to avoid catastrophic forgetting of the dynamics of old tasks, and should  
30 ideally exhibit both forward transfer (old tasks improve the learning performance on the new task)  
31 and backward transfer (new task improves the performance on old tasks). Many MBRL and MPC  
32 methods lack this type of adaptation and positive transfer.

33 In this work, we propose to extend the task-aware continual learning approach based on hypernet-  
34 works in [1] to adapt to changing environment dynamics and to address the scalability and positive  
35 transfer challenges mentioned above in a reinforcement learning setting. We use task-conditional  
36 hypernetworks, which are neural network models that accept a learned task encoding as an input,

37 and output the weights of another (target) network. In our case, the output is the dynamics model  
38 for that task. No additional information other than task transition boundary is needed. We consider  
39 the setting where task boundaries are known in order to simplify the problem.

40 Our work makes the following contributions: we show that task-aware continual learning with hy-  
41 pernetworks is an effective and practical way to adapt to new tasks and changing dynamics for  
42 model-based reinforcement learning without keeping state transitions from old tasks nor adding ca-  
43 pacity to the dynamics model. We evaluate our method on locomotion and manipulation scenarios,  
44 where we show that our method outperforms related continual learning baselines.

## 45 2 Related Works

46 **Continual Learning of Neural Networks** Continual learning studies the problem of incrementally  
47 learning from a sequential stream of data with only a small portion of the data available at once [2].  
48 A simple yet effective approach is finetuning, which directly tunes the trained source task network  
49 on the target task [3]. The efficacy of this approach for continual learning suffers from the well-  
50 established phenomenon of catastrophic forgetting [4]. Sequential Bayesian posterior updates are a  
51 principled way to perform continual learning and naturally avoids the forgetting problem since the  
52 exact posterior fully incorporates all previous data but in practice approximations have to be made  
53 that may be prone to forgetting. Elastic Weight Consolidation (EWC) uses a Laplace approximation  
54 to the posterior, storing previous tasks’ empirical fisher matrices and regularizing future task weight  
55 deviations under their induced norms [5]. Other works have also employed variational mean-field  
56 approximations [6] or block-diagonal Kronecker factored Laplace approximations [7]. Synaptic  
57 Intelligence (SI) forgoes an obvious approximate Bayesian interpretation but operates similarly to  
58 EWC in that it computes a relative parameter importance measure, but through a linear approxi-  
59 mation to the contribution in loss reduction due to each parameter on previous tasks [8]. Coreset  
60 methods prevent catastrophic forgetting by choosing and storing a significantly smaller subset of  
61 the previous task’s data, which is used to rehearse the model during or after finetuning [9, 10, 11].  
62 Similarly, the inducing points used in sparse Gaussian Process (GP) formulations, which can be seen  
63 as a type of coreset, has been used for continual learning [12, 13]. Another type of approach learns  
64 separate task-specific network components. The most common version of this are multi-head net-  
65 works that learn and switch between separate output layers depending on the task [14]. Progressive  
66 Neural Networks (PNN) [4] are an extreme version of this approach, in which an entirely new copy  
67 of the network is appended for each task, thereby eliminating any forgetting. These methods can  
68 incur significant memory and compute cost especially for larger models and many tasks.

69 **Continual RL** Memory-efficient continual learning methods in the reinforcement learning setting  
70 have also been proposed. PNN was used in an on-policy actor-critic method and was demonstrated  
71 on sequential discrete action Atari games. The authors of [15] build on top of PNN and continual  
72 policy compression methods [16] by compressing the extended model from PNN into a fixed size  
73 network after each task. For the compression stage, they propose a more scalable online EWC algo-  
74 rithm that eschews the linear cost of storing past fisher matrices. The use of coresets has also been  
75 explored in this setting [17]. In [18], a mixture model of separate task-specific neural networks is  
76 used for the environment model that requires adding a new model each time a task is instantiated.  
77 A similar approach was also demonstrated using a mixture of GPs using an online clustering algo-  
78 rithm [19]. Our work is also related to MPC interpretations as a reduction to online learning [20].

## 79 3 Preliminaries

80 **Hypernetworks for Continual Learning** A hypernetwork [21, 22] is a network that generates  
81 the weights of another neural network. The hypernetwork  $H_\Theta(e) = \theta$  with weights  $\Theta$  can  
82 be conditioned on an embedding vector  $e$  to output the weights  $\theta$  of the main (target) network  
83  $f_\theta(x) = f(x; \theta) = f(x; H_\Theta(e))$  by varying the embedding vector  $e$ . The hypernetwork is typically  
84 smaller with respect to the number of trainable parameters in comparison to the main network. Hy-  
85 pernetworks have been shown to be useful in the continual learning setting [1] for classification and  
86 generative models. This has been shown to alleviate some of the issues of *catastrophic forgetting*.  
87 They have also been used to enable gradient-based hyperparameter optimization [23].

88 **Planning with CEM and MPC** The Cross-Entropy Method (CEM) [24] is a widely used online  
89 planning algorithm that samples action sequences from a time-evolving distribution which is usu-

ally considered to be a diagonal Gaussian  $a_{1:h} \sim \mathcal{N}(\mu_{1:h}, \text{diag}(\sigma_{1:h}^2))$ , where  $h$  is the planning horizon. Action sequences are iteratively re-sampled and evaluated under the currently learned dynamics model, and the sampling distribution parameters  $\mu_{1:h}, \sigma_{1:h}$  are re-fitted to the top percentile of trajectories. CEM for planning in MBRL has been successfully used in a number of previous approaches [25, 26], as it alleviates exploitation of model bias compared to purely gradient based optimizations [27] and can better adapt to varying dynamics as compared to fully amortized policies [28].

## 4 The Proposed Approach

### 4.1 Problem Setting and Method Overview

We consider the following **problem setting**: A robot interacts with the environment to solve a sequence of  $T$  goal-directed tasks, each of which brings about different dynamics while having the same state-space  $\mathcal{S}$  and action space  $\mathcal{A}$ . The robot is exposed to the tasks sequentially online without revisiting data collected in a previous task. The robot also has finite memory and is not allowed to maintain a full history of state transitions for the purpose of re-training. Since the distribution of tasks changes over time and the agent does not know about it a priori, it must continually adapt to the streaming data of observations that it encounters, while trying to solve each task. The robot knows when a task switch occurs.

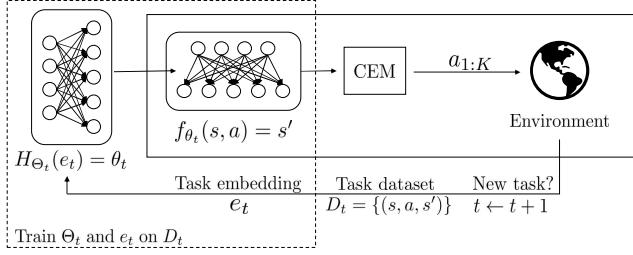


Figure 1: Overview of our proposed solution

We consider the **solution setting** of MBRL with a learned dynamics model, the parameters of which are inferred through a task-conditioned hypernetwork. Given learned task embeddings  $e_t$  and parameters  $\Theta_t$  of the hypernetwork  $H(\cdot)$ , we infer parameters  $\theta_t$  of the dynamics neural network  $f_{\theta_t}(\cdot)$ . Using this dynamics model, we perform CEM optimization to generate action sequences and execute them in the environment for  $K$  time-steps with MPC. We store the observed transitions in the replay dataset and update the parameters of the hypernetwork  $\Theta_t$  and task-embeddings  $e_t$  (off-policy optimization). We repeat this for  $M$  episodes per task, and for each of the  $T$  tasks sequentially.

### 4.2 Training Procedure

**Dynamics Learning** The learned dynamics model is a feed-forward neural network whose parameters vary across tasks. One way to learn a dynamics network  $f_{\theta}(\cdot)$  across tasks is to update it sequentially as training progresses. However, since our problem setting is such that the agent is not allowed to retain state-transition data from previous tasks in the replay buffer, adapting the weights of a single network sequentially across tasks is likely to lead to catastrophic forgetting [1]. In order to alleviate issues of catastrophic forgetting while trying to adapt the weights of the network, we learn a hypernetwork that takes task embeddings as input, and outputs parameters for the dynamics network corresponding to every task, learning different dynamics networks  $f_{\theta_t}(\cdot)$  for each task  $t$ .

We assume that the agent has finite memory and does not have access to state-transition data across tasks. So, the task-specific replay buffer  $\mathcal{D}_t$  is reset at the start of every task  $t$ . For the current episode, the agent generates a dynamics network  $f_{\theta_t}$  using  $\theta_t = H_{\Theta_t}(e_t)$ . Then, for  $k = 1 \dots K$  timesteps and planning horizon  $h$ , the agent optimizes action sequences  $a_{k:k+h}$  using CEM, and executes the first action  $a_k$  (MPC).  $\mathcal{D}_t$  is augmented by a tuple  $(s_k, a_k, s_{k+1})$ , where  $s_k$  is the current state,  $a_k$  is the executed action, and  $s_{k+1}$  is the next observed state under task  $t$ .

The parameters  $\Theta_t$  of the hypernetwork and the task embeddings  $e_t$  are updated by backpropagating gradients with respect to the sum of a dynamics loss  $\mathcal{L}_{\text{dyn}}$  and a regularization term. We define the dynamics loss as  $\mathcal{L}_{\text{dyn}}(\Theta_t, e_t) = \sum_{\mathcal{D}_t} \|\hat{s}_{k+1} - s_{k+1}\|_2$ , where the predicted next states are  $\hat{s}_{k+1} = f_{\theta_t}(s_k, a_k)$  and  $\theta_t = H_{\Theta_t}(e_t)$ . In practice, we infer the difference  $\Delta_{k+1}$  through the dynamics network ( $\Delta_{k+1} = f_{\theta_t}(s_k, a_k)$ ) such that  $\hat{s}_{k+1} = s_k + \Delta_{k+1}$  for stable training. In addition, inputs to the  $f_{\theta_t}$  network are normalized, following the procedure in previous works [25]. Similarly, a new  $e_t$  is initialized at the start of every task and updated every episode during the task by gradient descent. Older task embeddings ( $e_{1:t-1}$ ) are kept fixed.

144 **Regularizing the Hypernetwork** To alleviate catastrophic forgetting, we regularize the output of  
 145 the hypernetwork for all previous task embeddings  $e_{1:t-1}$ . After training for task  $t - 1$ , a snapshot  
 146 of the hypernetwork weights is saved as  $\Theta_{t-1}$ . For each of the past tasks  $i = 1 \dots t - 1$ , we use a reg-  
 147 ularization loss to keep the outputs of the snapshot  $H_{\Theta_{t-1}}(e_i)$  and the current output  $H_{\Theta_t}(e_i)$ . This  
 148 approach sidesteps the need to store all past data across tasks, preserves the predictive performance  
 149 of dynamic networks  $f_{\theta_t}$ , and only requires a single point in the weight space (a copy of the hyper-  
 150 network) to be stored. The task embeddings are differentiable vectors learned along with parameters  
 151 of the hypernetwork. The overall loss function for updating  $\Theta_t$  and  $e_t$  is given by the sum of the  
 152 dynamics loss  $\mathcal{L}_{\text{dyn}}(\cdot)$ , which is evaluated on the data collected from task  $t$  and the regularization  
 153 term  $\mathcal{L}_{\text{reg}}(\cdot)$ :

$$\mathcal{L}_t(\Theta_t, e_t) = \mathcal{L}_{\text{dyn}}(\Theta_t, e_t) + \mathcal{L}_{\text{reg}}(\Theta_{t-1}, \Theta_t, e_{1:t-1}) \quad (1)$$

$$= \mathcal{L}_{\text{dyn}}(\Theta_t, e_t) + \frac{\beta_{\text{reg}}}{t-1} \sum_{i=1}^{t-1} \|H_{\Theta_{t-1}}(e_i) - H_{\Theta_t}(e_i)\|_2^2 \quad (2)$$

154 The planning objective for CEM optimization of action sequences is given by the sum of rewards  
 155 obtained by executing the sequence of actions  $a_{k:k+h}$  under the learned dynamics model  $f_{\theta_t}(\cdot)$  for  
 156 the task. The reward function  $r(s, a)$  is assumed to be known, but nothing precludes learning it from  
 157 data under our current framework.

## 158 5 Experiments

159 We perform multiple robot simulation experiments to answer the following questions: (a) How does  
 160 HyperCRL compare with existing continual learning baselines in terms of overall performance across  
 161 tasks? (b) How effective is HyperCRL in forward transfer and backward transfer across tasks?

### 162 5.1 Pushing a block of non-uniform mass

163 First, we look at an intuitively simple experiment simulated using Surreal Robotics Suite [29], in  
 164 which a Panda robot tries to push a non-uniform-density block across the table. The objective is to  
 165 push the block to the goal position while **maintaining its initial orientation**. We vary the densities  
 166 of the left and right parts of the block across different tasks ( $T = 5$ ), changing the center of mass  
 167 and moment of inertia. The robot needs to learn to maneuver the position of its end-effector on the  
 168 side of the block to correct for orientation deviations while pushing the block forward.

169 **Baselines** We compare the hypernetwork to the following baselines: (i) Multi-task learning with  
 170 access to a buffer of all data from all previous tasks (an oracle) (ii) Coreset that remembers one

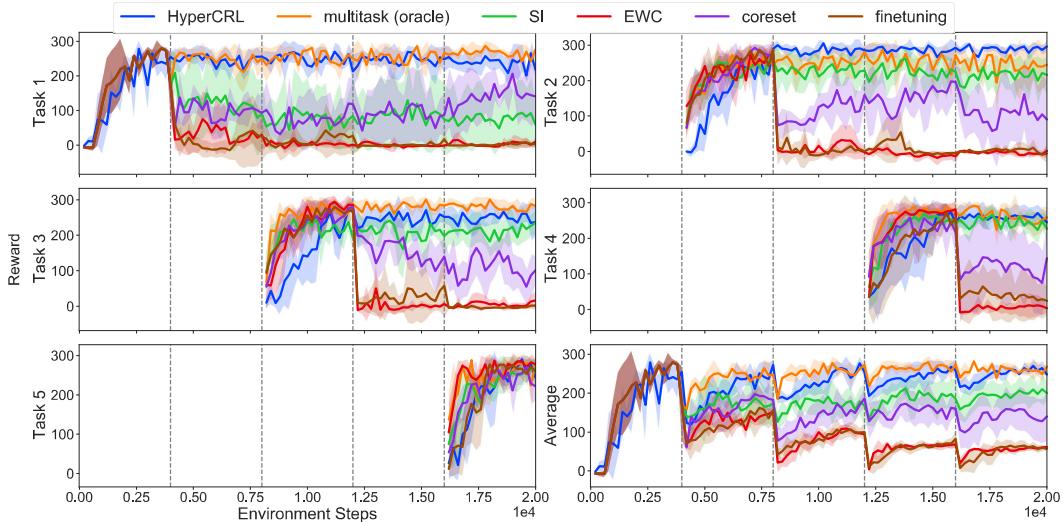


Figure 2: **Reward on Pusher Environment.** Shown are episodic rewards evaluated during training. Results are averaged across four random seeds, and the shaded areas represent one standard deviation. Each task is trained for 4k steps, summing to 20k steps in total. Vertical dotted lines indicate task switches. The bottom-right subplot shows the average reward across all task  $\leq t$  seen so far.

percent of the state-action transition data per task, sampled randomly (iii) Synaptic Intelligence [8] (iv) Elastic Weight Consolidation [5] (v) Fine-tuning, where we optimize  $f_{\theta_t}(\cdot)$  on each task's data without regularization. All the baseline models resemble the target model in our hypernetwork setup, except the multi-head output layer with one head per task. For coresnet or multi-task learning, an additional batch of past data is sampled from the coresnet or oracle every update step, and contributes to the total dynamics loss.

**Results** HyperCRL is able to learn to push all 5 of the blocks across the table with minimal forgetting (Table 1), and even shows signs of positive backward transfer. We provide detailed definition of backward and forward transfer in the Appendix A. The average performance of our method is on par with the multitask learning baseline (Figure 2), which has access to the entire history of data. HyperCRL also outperforms other continual learning baselines, either regularization-based (SI, EWC), or replay-based (Coreset). Simple finetuning, however, catastrophically forgets and is unable to perform pushing for all 5 types of block at the end. To further illustrate how our model benefits from past experience (positive forward transfer), we compare our method against a single-task baseline trained from scratch in the appendix.

## 5.2 Door Opening

Next, we experiment on a more complex task to better demonstrate the flexibility of HyperCRL. We choose a door opening experiment, where the Panda robot has to open doors with different types of handles. Unlike the pushing example, the dynamic switches between tasks are much more discontinuous and the tasks require different motions to solve, due to the rotational joints imposed by some of the handles. A total of  $T = 5$  tasks need to be solved sequentially, each involving a different handle type. Tasks 2 and 4 (similarly tasks 3 and 5) both have a round (lever) handle, but with different turning directions. The environment is modified from the DoorGym environment [30] and simulated in the Surreal Robotics Suite.

**Results** HyperCRL outperforms all continual learning baselines on the door opening task, and the multi-task baseline with oracle. Figure 3 shows the learning curves of all our evaluated methods, compared to a single-task baseline trained from scratch on each task. HyperCRL virtually sees no performance degradation in terms of reward (Table 1). We further evaluate the ability to positively transfer past knowledge in the appendix (Table 3).

## 5.3 Further Analysis

From the two simulation experiments, we demonstrate that our proposed method outperforms state of the art baselines on the following objectives of continual reinforcement learning: 1) High reward across all previously seen tasks. and 2) Minimal forgetting of previous tasks as new tasks need to be learned. From the learning curves (Figure 2 and 3), we observe that HyperCRL consistently outperforms competing baseline throughout the whole learning process. From the retention metric in Table 1, we show that HyperCRL can remember equal or more knowledge than other methods.

Table 1: **Backward Transfer on Pusher and Door Envs.** We measure performance retained at the end of training all 5 tasks compared to the end of training on task  $t$ . Results show the mean and one std. deviation, evaluated across 4 seeds and 10 episodes per seed. Greater than 100 indicates positive backward transfer, otherwise it indicates forgetting.

Task	% Retention in terms of Reward (Pusher)				
	1	2	3	4	Average
Multi-task (Oracle)	<b>142 ± 74</b>	91 ± 25	96 ± 9	91 ± 11	<b>106 ± 20</b>
HyperCRL	99 ± 10	<b>107 ± 9</b>	<b>98 ± 13</b>	<b>103 ± 15</b>	102 ± 6
SI ( $c = 0.1$ )	40 ± 54	88 ± 21	95 ± 21	92 ± 14	79 ± 16
EWC ( $\lambda = 10^5$ )	7 ± 15	13 ± 8	6 ± 8	0 ± 3	4 ± 5
Coreset (1% Data)	87 ± 66	32 ± 46	39 ± 13	61 ± 43	55 ± 23
Finetuning	0 ± 2	-2 ± 6	1 ± 3	13 ± 16	3 ± 4
Task	% Retention in terms of Normalized Reward (Door, see Figure 3)				
	1	2	3	4	Average
Multi-task (Oracle)	37 ± 70	81 ± 88	59 ± 60	<b>93 ± 129</b>	68 ± 45
HyperCRL	<b>113 ± 75</b>	<b>100 ± 76</b>	<b>99 ± 36</b>	86 ± 68	<b>100 ± 26</b>
SI ( $c = 0.1$ )	-17 ± 27	11 ± 61	6 ± 37	16 ± 60	4 ± 24
EWC ( $\lambda = 10^5$ )	-6 ± 33	17 ± 59	4 ± 25	16 ± 50	7 ± 22
Coreset (1% Data)	-5 ± 28	26 ± 53	21 ± 34	45 ± 68	22 ± 24
Finetuning	-14 ± 28	11 ± 58	3 ± 23	16 ± 65	4 ± 23

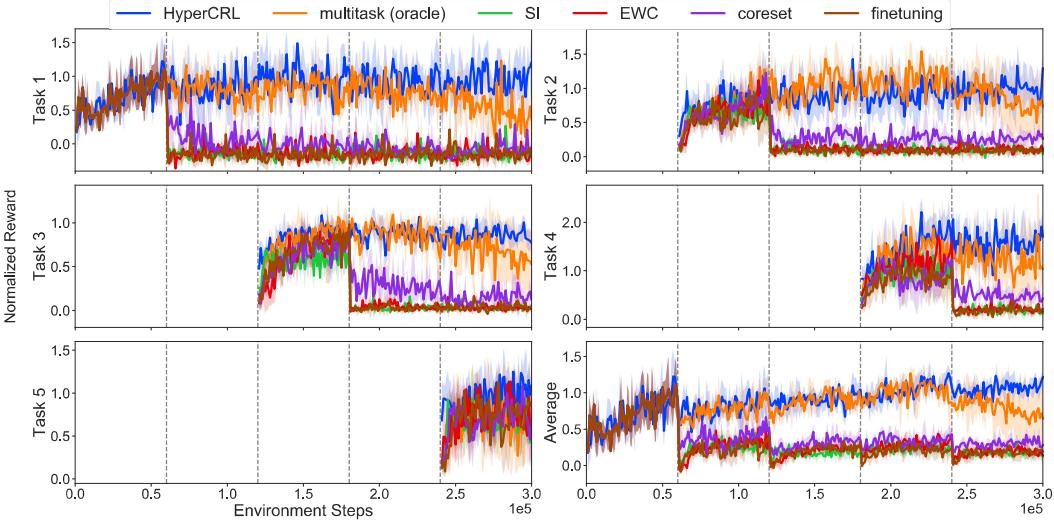


Figure 3: **Normalized Reward on Door Environment** during training. Results are averaged across four random seeds, and the error bars represent one standard deviation. Each task is trained for 60k steps, totaling 300k steps. The reward is normalized with respect to a model trained from scratch separately on each task.

207 One observation here is that the multi-head multitask  
 208 baseline underperforms HyperCRL (Figure 5). We hy-  
 209 pothesize that the hypernetwork architecture might be  
 210 more effective for learning the task since it allows mod-  
 211 ellng multiple distinct dynamics in one model. We  
 212 consider another multi-task baseline, HyperCRL-MT, that  
 213 shares the same architecture of the hypernetwork used in  
 214 HyperCRL, but is trained without regularization and has  
 215 access to the entire replay buffer similar to the multitask  
 216 baseline. We show that this baseline outperforms the mul-  
 217 titask baseline while sharing the same training procedure, in Table 2. Thus, we infer that the hyper-  
 218 network architecture is what makes the difference.

Method	Final Reward	Average
HyperCRL-MT	$1.25 \pm 0.39$	
HyperCRL	$1.08 \pm 0.31$	
Multitask (Oracle)	$0.73 \pm 0.36$	

Table 2: Final performance comparison for the door opening task. Results are normalized episodic reward, averaged over all tasks

## 219 6 Discussion, Limitations, and Conclusion

220 In this paper we present HyperCRL, a task-aware method for continual model-based reinforcement  
 221 learning using hypernetworks. In all of our experiments, we have demonstrated that HyperCRL con-  
 222 stantly outperforms alternative continual learning baselines in terms of overall performance at the  
 223 end of training, as well as in terms of retaining performance on previous tasks. By allowing the en-  
 224 tire network to change between tasks, rather than just the output head layer of the dynamics network,  
 225 HyperCRL is more effective at accurately representing different dynamics, even with significant dis-  
 226 continuity between them, while only requiring a fixed-size hypernetwork and constant time updates  
 227 to generate task-conditional dynamics for planning.

228 While our proposed approach shows good results against the baselines we tested, there are many in-  
 229 teresting prospects for future work. First, better techniques to train hypernetworks will significantly  
 230 aid this line of work. Currently, the size of our hypernetwork is at least an order of magnitude larger  
 231 than the target network, since it directly outputs all the weights. Hypernetworks are often sensi-  
 232 tive to the choice of random seeds and architecture. Second, extending our method to image-based  
 233 RL environments is worth investigating, as it will enable higher-capacity target networks. Finally,  
 234 HyperCRL is not task agnostic, nor can it automatically detect task switching that happens contin-  
 235 uously with no clear task boundaries. A possible direction is to use probabilistic inference models  
 236 (i.e. Bayesian non-parametrics [31]) or changepoint detection methods to perform task identifica-  
 237 tion.

238 **References**

- 239 [1] J. von Oswald, C. Henning, J. Sacramento, and B. F. Grewe. Continual learning with hyper-  
240 networks. In *International Conference on Learning Representations*, 2019.
- 241 [2] M. D. Lange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and  
242 T. Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *arXiv: Computer Vision and Pattern Recognition*, 2020.
- 243 [3] Y. Bengio. Deep learning of representations for unsupervised and transfer learning. In *Proceedings of ICML workshop on unsupervised and transfer learning*, pages 17–36, 2012.
- 244 [4] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu,  
245 R. Pascanu, and R. Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*,  
246 2016.
- 247 [5] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan,  
248 J. Quan, T. Ramalho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in  
249 neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- 250 [6] C. V. Nguyen, Y. Li, T. D. Bui, and R. E. Turner. Variational continual learning. *arXiv preprint*  
251 *arXiv:1710.10628*, 2017.
- 252 [7] H. Ritter, A. Botev, and D. Barber. Online structured laplace approximations for overcoming  
253 catastrophic forgetting, 2018.
- 254 [8] F. Zenke, B. Poole, and S. Ganguli. Continual learning through synaptic intelligence. *Proceedings of machine learning research*, 70:3987, 2017.
- 255 [9] D. Lopez-Paz and M. Ranzato. Gradient episodic memory for continual learning. In *Advances*  
256 *in neural information processing systems*, pages 6467–6476, 2017.
- 257 [10] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny. Efficient lifelong learning with  
258 a-gem, 2018.
- 259 [11] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. Torr, and M. Ran-  
260 zato. Continual learning with tiny episodic memories. 2019.
- 261 [12] M. K. Titsias, J. Schwarz, A. G. de G. Matthews, R. Pascanu, and Y. W. Teh. Functional  
262 regularisation for continual learning using gaussian processes. *ArXiv*, abs/1901.11356, 2020.
- 263 [13] S. Kapoor, T. Karaletsos, and T. D. Bui. Variational auto-regressive gaussian processes for  
264 continual learning, 2020.
- 265 [14] Z. Li and D. Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and*  
266 *machine intelligence*, 40(12):2935–2947, 2017.
- 267 [15] J. Schwarz, J. Luketina, W. M. Czarnecki, A. Grabska-Barwinska, Y. W. Teh, R. Pascanu, and  
268 R. Hadsell. Progress & compress: A scalable framework for continual learning. *arXiv preprint*  
269 *arXiv:1805.06370*, 2018.
- 270 [16] G. Berseth, K. Xie, P. Cernek, and M. van de Panne. Progressive reinforcement learning with  
271 distillation for multi-skilled motion control. *ArXiv*, abs/1802.04765, 2018.
- 272 [17] D. Abel, D. Arumugam, L. Lehnert, and M. Littman. State abstractions for lifelong reinforce-  
273 ment learning. In *International Conference on Machine Learning*, pages 10–19, 2018.
- 274 [18] A. Nagabandi, C. Finn, and S. Levine. Deep online learning via meta-learning: Continual  
275 adaptation for model-based rl. *arXiv preprint arXiv:1812.07671*, 2018.
- 276 [19] A. Abdollahi, R. Allamraju, and G. Chowdhary. Adaptive-optimal control of nonstationary  
277 dynamical systems. In *Proceedings of the 2015 European Guidance, Navigation, and Control*  
278 *Conference*, 2015.
- 279 [20] N. Wagener, C. Cheng, J. Sacks, and B. Boots. An online learning approach to model predictive  
280 control. *CoRR*, abs/1902.08967, 2019. URL <http://arxiv.org/abs/1902.08967>.

- 284 [21] D. Ha, A. Dai, and Q. V. Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- 285 [22] J. U. Schmidhuber. Learning to control fast-weight memories: an alternative to dynamic re-  
286 current networks. 1991.
- 287 [23] M. MacKay, P. Vicol, J. Lorraine, D. Duvenaud, and R. B. Grosse. Self-tuning networks:  
288 Bilevel optimization of hyperparameters using structured best-response functions. *CoRR*,  
289 abs/1903.03088, 2019. URL <http://arxiv.org/abs/1903.03088>.
- 290 [24] R. Y. Rubinstein. Optimization of computer simulation models with rare events. *European  
291 Journal of Operational Research*, 99(1):89–112, 1997.
- 292 [25] K. Chua, R. Calandra, R. McAllister, and S. Levine. Deep reinforcement learning in a handful  
293 of trials using probabilistic dynamics models. In *Advances in Neural Information Processing  
294 Systems*, pages 4754–4765, 2018.
- 295 [26] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent  
296 dynamics for planning from pixels. In *International Conference on Machine Learning*, pages  
297 2555–2565, 2019.
- 298 [27] H. Bharadhwaj, K. Xie, and F. Shkurti. Model-predictive control via cross-entropy and  
299 gradient-based optimization. *Learning for Dynamics and Control*, 2020.
- 300 [28] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent  
301 imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 302 [29] L. Fan, Y. Zhu, J. Zhu, Z. Liu, O. Zeng, A. Gupta, J. Creus-Costa, S. Savarese, and L. Fei-Fei.  
303 Surreal: Open-source reinforcement learning framework and robot manipulation benchmark.  
304 In *Conference on Robot Learning*, 2018.
- 305 [30] Y. Urakami, A. Hodgkinson, C. Carlin, R. Leu, L. Rigazio, and P. Abbeel. Doorgym: A scal-  
306 able door opening environment and baseline agent. *arXiv preprint arXiv:1908.01887*, 2019.
- 307 [31] M. Xu, W. Ding, J. Zhu, Z. Liu, B. Chen, and D. Zhao. Task-agnostic online reinforcement  
308 learning with an infinite mixture of gaussian processes, 2020.
- 309 [32] O. Khatib. A unified approach for motion and force control of robot manipulators: The opera-  
310 tional space formulation. *IEEE Journal on Robotics and Automation*, 3(1):43–53, 1987.
- 311 [33] P. Henderson, W.-D. Chang, F. Shkurti, J. Hansen, D. Meger, and G. Dudek. Benchmark  
312 environments for multitask learning in continuous domains. *arXiv preprint arXiv:1708.04352*,  
313 2017.
- 314 [34] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural  
315 networks. In *Proceedings of the thirteenth international conference on artificial intelligence  
316 and statistics*, pages 249–256, 2010.
- 317 [35] G. M. van de Ven and A. S. Tolias. Three scenarios for continual learning. *arXiv preprint  
318 arXiv:1904.07734*, 2019.