FOUNDATION POLICIES WITH MEMORY

Anonymous authors

Paper under double-blind review

ABSTRACT

A generalist agent should perform well on novel tasks in unfamiliar environments. While Foundation Policies (FPs) enable generalization across new tasks, they lack mechanisms for handling novel dynamics. Conversely, agents equipped with memory models can adapt to new dynamics, but struggle with unseen tasks. In this work, we bridge this gap by integrating memory models into the FP architecture, allowing policies to condition on both task and environment dynamics. We evaluate FPs enhanced with attention, state-space, and RNN-based memory models on POPGym, a memory benchmark, and ExORL, an unsupervised RL benchmark. Our results show that GRUs achieve the best generalization to unseen tasks and dynamics for a given recurrent state size, approaching the performance of a supervised baseline that has access to task information during training and significantly outperforming memory-free FPs. Additionally, our approach improves FP performance on entirely new environments not encountered during training. Our anonymized code is available at https: //anonymous.4open.science/r/zero-shot-96A1, and our datasets are open-sourced at REDACTED.

1 INTRODUCTION

029 030 031

032

033

034

035

037

038

000

001 002 003

004

005 006 007

008 009 010

011

012

013

014

015

016

017

018

019

021

Reinforcement Learning (RL) agents [92] exhibit superhuman decision-making skill when tasked with a *single* objective in a *single* environment [89, 67, 90, 91]. A new line of work focuses on producing *generalist agents* that replicate such results across *many* tasks and environments [83, 53, 105]. Foundation policies (FPs) [96, 97, 76, 45] are a promising approach for building generalist agents, providing a principled mechanism for generalising to *any* downstream task in an environment after an offline reward-free pre-training phase. However, as yet, FPs are not equipped to deal with a change in dynamics between pre-training and deployment.

A concurrent line of work on *in-context* RL attempts to build generalist agents by using *memory models* to condition policies on reward-labelled trajectories [14, 43, 56, 54, 62, 23] or to reach arbitrary goal states [31]. In principal, these models can perform dynamics generalisation by inferring changes between training and testing from the trajectory used to condition the policy. However, they lack the task generalisation ability of FPs for two reasons. They are either 1) trained with reward supervision and so cannot reliably generalise to new tasks with different reward functions, or are 2) trained without reward supervision to reach any goal-state in an environment and so cannot reliably generalise to new tasks not codified by a goal state.

Here, we reconcile these lines of work and propose *foundation policies with memory*, an architecture that, like in-context RL agents, infers the current dynamics context using powerful memory models and passes it to an FP for solving unseen tasks. We evaluate FPs with attention [98], state-space [32, 33], and RNN-based [24, 17] memory models across a range of experiments testing their ability to infer the dynamics context, and generalise to unseen tasks in unseen dynamics. We find that GRUs achieve the best generalisation to unseen tasks and dynamics for a given recurrent state size, approaching the performance of a supervised baseline that has access to task information during training and significantly outperforming memory-free FPs (Figure 1). Finally, we find that FPs with memory improve FP performance on entirely new environments not seen during training.



Figure 1: **Zero-shot task and dynamics generalisation.** FPs with memory models generalise to test tasks and dynamics not seen during training on the ExORL benchmark. FB-GRU approaches the performance of a supervised baseline, TD3-GRU, despite being trained without rewards. A full discussion is provided in Section 4.3.

2 PRELIMINARIES

054

055 056

058 059

060

061 062

063

064 065

066

067 068

069

071

073 074

075 076

077 Contextual markov decision processes. A Contextual Markov Decision Process (CMDP) is defined by $(C, S, O, \phi, A, R, \rho, \gamma, M(c))$. C is the set of contexts, S is the underlying state space, 078 O and A are sets of observations and actions, $\phi: S \to O$ is the observation function, $R: S \to \mathbb{R}$ 079 is a reward-function specifying a *task*, γ is a discount factor, and ρ is the initial state distribution [35]. M is a function that maps a context $c \in C$ to a Partially Observable Markov Decision Process 081 (POMDP) [2] $M(c) = (S, O, A, R, \rho, \gamma, P^c)$ with a context-dependent transition function $P^c : S \times A \times C \to \Delta(S)$. A Markov policy $\pi : S \to \Delta(A)$ is optimal in context c for reward function R if 083 it maximises the expected discounted future reward *i.e.* $\pi_{c,R}^* = \arg \max_{\pi} \mathbb{E}[\gamma^t R(s_{t+1})|s_0, a_0, \pi, c],$ 084 where $\mathbb{E}[\cdot|s_0, a_0, \pi, c]$ is the expectation under state-action sequence $(s_t, a_t)_{t>0}$ starting at (s_0, a_0) 085 with $s_t \sim P^c(\cdot | s_{t-1}, a_{t-1})$ and $a_t \sim \pi(\cdot | s_t)$. Note that the context $c \in \overline{C}$ cannot be observed directly. 087

Problem setting. We split the CMDP into a set of training contexts C_{train} and testing contexts C_{test} . We assume access to a dataset $\mathcal{D}_{\text{train}}$ of *unlabelled* observation-action trajectories $\tau = (o_0, a_0, o_1, \dots, o_T)$ collected from the training contexts by a highly exploratory behaviour policy. Our goal is to pre-train an adaptive policy $\pi(a|h, z)$, where $h \in \mathbb{R}^m$ is a hidden state summarising both the context c and inferred Markov state s, and $z \in \mathbb{R}^d$ denotes a compact representation of the task. We will pre-train this policy solely from offline data $\mathcal{D}_{\text{train}}$, without online interactions.

We will evaluate the policy on an unseen test task R_{test} in an unseen test context $c_{\text{test}} \in C_{\text{test}}$. The test task is revealed either via $\mathcal{D}_{\text{test}}$, a small dataset of *labelled* observation trajectories $((o_{t-L}, \ldots, o_t), R_{\text{test}}(s_t))$ of length L, or as an explicit function $o \mapsto R_{\text{test}}(s)$ (like 1 at a goal state and 0 elsewhere)¹. Unless, the agent can infer c_{test} from $\mathcal{D}_{\text{test}}$, it will need to infer it from the observation-action history it observes during evaluation. This problem setting is directly equivalent to [97]'s zero-shot RL setting with a change in the environment dynamics between training and testing. As a result, we call it **zero-shot RL under changed dynamics**.

Foundation policies. Foundation policies (FPs) approximate the (universal) successor features [6, 11] of near-optimal policies for any task in an environment. They require access to a feature map $\varphi: S \mapsto \mathbb{R}^d$ that embeds states into a representation space in which the reward is assumed to be linear *i.e.* $R(s) = \varphi(s)^\top z$ with weights $z \in \mathbb{R}^d$ representing a task. The USFs $\psi: S \times A \times \mathbb{R}^d \to \mathbb{R}^d$

¹Note that the agent only sees the observation, but the reward is a function of the underlying state.

are defined as the discounted sum of future features subject to a task-conditioned policy $\pi(s, z)$:

$$\psi(s_0, a_0, z) = \mathbb{E}\left[\sum_{t \ge 0} \gamma^t \varphi(s_{t+1}) | s_0, a_0, \pi(s, z)\right] \quad \forall \ s_0 \in S, a_0 \in A, z \in \mathbb{R}^d.$$
(1)

where the policy is trained in an actor-critic formulation [51] such that

$$\pi(s, z) \approx \arg\max_{a} \psi(s, a, z)^{\top} z, \ \forall \ s \in S, a \in A, z \in \mathbb{R}^d,$$
(2)

where $\psi(s, a, z)^{\top} z$ is the Q function (critic) formed by ψ . During training, candidate task weights are sampled from \mathcal{Z} , a prior over the task space². During evaluation, the test task weights are found by regressing labelled states onto the features: $z_{\text{test}} := \arg \min_z \mathbb{E}_{s \sim d}[(R_{\text{test}}(s) - \varphi(s)^{\top}z)^2]$, before being passed to the policy. The features can be learned with Hilbert representations [76], laplacian eigenfunctions [97], contrastive methods [97], or in service of *successor measure* prediction [10], as is the case for the forward Backward (FB) foundation policy [96] used in this work.

METHOD

Recall that our goal is to pre-train an adaptive policy $\pi(a|h, z)$ that is conditioned on h, a hidden state summarising both the context c and inferred Markov state s, and task z. As we outlined in Section 2, the FP framework provides a principled way of pre-training $\pi(a|s, z)$ *i.e.* a policy conditioned solely on the task and Markov state. In this section, we will discuss amendments to the FP framework that allow policies to be conditioned on h rather than s.

Following past work on RL in CMDPs, we assume that we can produce an estimate of the dynamics context c and Markov state s from a trajectory of observation-action pairs $\tau = x_0, \ldots, x_L$, where $x_n = \epsilon(o_n, a_n)$ is some encoding of an observation-action pair and L is the *context length* [31, 65]. We seek a model of the form

$$y_j, h_j = f(x_j, h_{j-1}), \ j \in [1, \dots, L],$$
(3)

where x_j, y_j are the inputs and outputs at time j, and f updates a hidden state $h \in \mathbb{R}^m$ summarising the current Markov state and dynamics context prediction. This is the standard setup of a memory model in RL [4, 68, 69, 70, 82, 66, 73, 40, 86, 100, 8, 109], because the asymptotic inference time complexity is $\mathcal{O}(1)$ which is helpful for fast data collection, or high-frequency motor control [62]. Until recently, only Recurrent Neural Networks (RNNs) [24, 41, 17] have had this property, but newly proposed structured state-space models (S4) [32, 33, 34] and fast Transformers [98, 19, 48] have runtime complexity approaching that of RNNs, and model histories with large L more accurately. We explore all of these memory models in Section 4.

3.2 FOUNDATION POLICIES WITH MEMORY

Equipped with memory model f, we now condition the FP's actor and critic on the hidden state it produces. We define *contextual* USFs as the discounted sum of future features extracted from the hidden state, subject to a policy conditioned on the inferred Markov state and dynamics context $\pi(h,z)$

$$\psi(h,z) = \mathbb{E}\left[\sum_{t\geq 0} \gamma^t \varphi(h_{t+1}) | h_0, \pi(h,z)\right] \quad h_0 = \mathbf{0}^m, \forall \ z \in \mathbb{R}^d, \tag{4}$$

where $h_t = f(x_t, h_{t-1})$ from Equation 3, x_t is zero-padded for all t < L, and $h_0 = \mathbf{0}^m$ is an initial hidden-state of zeroes. The policy is trained such that

$$\pi(h, z) \approx \arg \max \psi(h, z)^{\top} z, \ \forall \ h \in \mathbb{R}^m, \ z \in \mathbb{R}^d,$$
(5)

²See Appendix B.1.1 for more detail on \mathcal{Z} .



Figure 2: Foundation policies with memory. FPs are optimised in a standard actor critic setup [51]. The policy π selects an action a_t conditioned on a history of observations and actions $o_{t-L}, a_{t-L-1}, \ldots, o_t, a_{t-1}$ of length L encoded by the actor's memory model, and the task vector z. The Q function formed by the USF ψ evaluates the sequence of observations and actions $o_{t-L}, a_{t-L}, \ldots, o_t, a_t$ encoded by the critic's memory model for task z. The architecture of an FP *without* memory is illustrated in Figure 6 in Appendix B for comparison.

180

181

182

183

where $\psi(h, z)^{\top} z$ is the Q function (critic) formed by ψ . Training proceeds exactly as with conven-188 tional USFs, and the test-time task weights are found by regressing labelled states onto the hidden-189 state features: $z_{\text{test}} := \arg \min_{z} \mathbb{E}_{s_t, (o_{t-L:t}, a_{t-L:t}) \sim d} [(R_{\text{test}}(s_t) - \varphi(f((o_{t-L:t}, a_{t-L:t}), h_{0:L})^\top z)^2],$ 190 before being passed to the policy. The full architecture and optimisation procedure is summarised in Figure 2. We found that using a shared memory model for the actor and critic led to model collapse, so use separate memory models for each. This corroborates the findings of [73]. Full im-192 plementation details are provided in Appendix B. In the experiments discussed in Section 4 we use FB representations as our FP which follow a slightly different training procedure. We discuss these 194 details in Appendix B.

195 196 197

198 199

200

201

202

203

204

191

193

EXPERIMENTS 4

In this section we perform an empirical study to evaluate our proposed method. We seek answers to three questions: (Q1) Can our method encode trajectories into a Markov state for use in solving one task in an environment? (Q2) Can our method generalise to unseen tasks in an environment with different dynamics to those seen in training? I.e. can our method perform zero-shot RL under changed dynamics? And (Q3) Can our method generalise to unseen tasks in a completely different environment to those seen in training? I.e. can our method perform zero-shot environment generalisation?

4.1 Setup

209 **Environments.** We respond to Q1 using the POPGym benchmark [68], a set of tests that evaluate 210 an agent's ability to infer Markov states from trajectories of observations and actions. We only eval-211 uate on the "Hard" versions of CartPole, Pendulum, Noisy CartPole, Noisy Pendulum and Repeat 212 Previous environments following [62]. For these experiments, we allow the agent to recondition 213 its policy on the previous L observation-action pairs every step so we can disentangle the memory model's ability to accurately model the Markov state from its ability to carry forward an accurate 214 hidden state. For all other experiments we do not allow such re-conditioning and require the policy 215 to condition on only the previous hidden state, current observation-action pair, and task. Note for

these experiments $C_{\text{train}} = C_{\text{test}}$, so we are not yet testing whether our method can generalise across contexts.

We respond to Q2 using the ExORL benchmark [108], a set of tests that evaluate an agent's ability 219 to generalise to unseen tasks on the DeepMind Control Suite [93]. We evaluate on the same envi-220 ronments as [97]: Walker, Maze, Cheetah and Quadruped, removing the velocities from each of the 221 state spaces to ensure the observations are not Markov, and call these variants occluded. To evaluate 222 dynamics generalisation we train on datasets collected from environment instances where the robot's 223 mass and damping coefficient are scaled to $\{0.5x, 1.5x\}$ of their usual values. We then evaluate on 224 environment instances where the robot's mass and damping coefficient are scaled by $\{1.0x, 2.0x\}$ 225 of their usual values, where 1.0x requires the agent to generalise via interpolation, and 2.0x requires 226 the agent to generalise via extrapolation [74]. We evaluate on all tasks provided by the DeepMind control suite, and increase the number of goals in Maze from 4 to 20 for a total of 32 tasks across 4 227 environments. 228

We also respond to Q3 using the ExORL benchmark [108]. This time we train on Walker-Occluded and Quadruped-Occluded and test on Cheetah-Occluded. The dynamics are unscaled (i.e. 1.0x) and, as before, we evaluate on all tasks provided by the DeepMind control suite.

Aggregation across tasks or environments is always summarised by the Interquartile Mean (IQM) and standard deviation following the recommendations of [1]. On POPGym we report the meanmax epoch reward (MMER) metric used in the original paper. On ExORL, we report scores from the learning step for which the all-task IQM is maximised across seeds. Full experimental details are provided Appendix A, and a full description of our evaluation protocol is provided in Appendix A.3.

238 We use FB [96] and HILP [76] as our FP baselines. FB is the most performant FP **Baselines**. 239 utilising successor measures, and HILP is the most performant FP utilising successor features. Both 240 methods assume access to the Markov state for training as discussed in Section 2. So, instead of 241 conditioning their predictions on a single observation, we provide them a stack of the 4 most recent 242 observations i.e. $s_t = (o_{t-3}, o_{t-2}, o_{t-1}, o_t)$, also known as *frame-stacking* [67]. Frame-stacking 243 is a naive method for inferring a Markov state from a short trajectory of observations, and is the 244 first solution one would use when faced with our problem. We also baseline against a single-task, 245 memory-based, supervised RL method. For this we use Offline TD3 [28] with a GRU memory model [17], which we refer to as **TD3-GRU**. Offline TD3 is the most performant single-task method 246 on the ExORL benchmark [108]; TD3-GRU is the most performant method in [73], and TD3 with an 247 LSTM memory model was shown to be particularly performant in [66]. TD3-GRU should indicate 248 how well an agent optimising for one task performs if provided reward supervision. 249

250 Though FPs can be deployed online they require an exploration policy for data col-Datasets. lection. To disentangle test-time performance from an agent's data collection ability, we collect 251 datasets on their behalf in advance using RND [11], an unsupervised RL algorithm. Agents trained 252 on datasets collected with RND exhibit better performance than comparable methods like APS [60], 253 APT [61], Proto [107] and DIAYN [25] in [108, 97, 76]. RND is run for 5 million learning steps 254 in each of our environments and every transition is cached. For the supervised baseline TD3-GRU, 255 transitions are relabelled with the appropriate rewards for a given task following [108]. All other 256 methods are trained on these datasets reward-free. 257

Memory models. We test the performance of FPs equipped with three memory models. We use 258 the most performant versions from each of the categories discussed in Section 3: attention-based, 259 state-space based, and RNN-based. For our attention-based memory model we use a Transformer 260 [98] with *FlashAttention* [19] for faster inference than a conventional Transformer. For our state-261 space-based memory model we use Diagonalized S4 [33], which uses a diagonal update matrix to 262 perform faster training and inference than the popular S4 model [32]. For our RNN-based memory 263 model we use a GRU [17] as it is the most performant RNN on the POPGym benchmark. Hereafter 264 we refer to the FB models we augment with these as FB-TF, FB-S4 and FB-GRU respectively. 265 To ensure a fair comparison across memory models, we follow [68] and restrict each model to a 266 fixed hidden state size, rather than a fixed parameter count. Concretely, we allow each model a 267 hidden state size of $32^2 = 1024$ dimensions. In Section 4.3 we condition the models on trajectories of length 32, so a hidden state size of 32^2 allows the attention-based, and state-space models to 268 perform their tensor products across the full input trajectory, and gives the RNN two 512-dimension 269

270 layers in which to summarise the trajectory. Full implementation details are provided in Appendix 271 Β.

272 273

274

292

293

304 305

306

307

308

310

311

312

313 314

315

316

Popgym 4.2

275 We report the aggregate performance of all FB-based 276 algorithms on the POPGym environments in Figure 3. 277 Our supervised baseline, TD3-GRU, performs similarly 278 to the PPO-GRU approach in the original POPGym paper. FB with frame-stacking performs poorly, reaching 279 only 30% of TD3-GRU's aggregate score. Our three 280 memory-based methods perform comparitively better, 281 with FB-TF reaching 80% of TD3-GRU's performance, 282 and FB-S4 and FB-GRU matching TD3-GRU's perfor-283 mance. We find that all methods fail on the RepeatPrevi-284 ousHard environment, where other in-context RL agents 285 have shown strong performance [62, 31]. This task requires the agent to remember the suit of a card dealt 64 287 timesteps ago (Appendix A), and our models are trained 288 with context length L = 64. This suggests that the mem-289 ory models are not accurately recalling information from the start of their context. The implications of our choice 290 of length L are discussed in Section 5. 291



Figure 3: **POPGym results.** Aggregate mean-maximum epoch reward (MMER) across all POPGym environments, normalised w.r.t. TD3-GRU performance.

ZERO-SHOT RL UNDER CHANGED DYNAMICS 4.3

We report the aggregate performance of all algorithms

295 on our zero-shot RL under changed dynamics experiments in Figure 4 (left), and the ratios of inter-296 polation/extrapolation performance to train performance in Figure 4 (right). As with our POPGym 297 experiments, FB performs poorly, reaching $\sim 25\%$ of the performance of our supervised baseline 298 on the training environments. HILP performs slightly better, as we would expect given its results on 299 EXORL in [76], but still much poorer than TD3-GRU. Of the three FPs with memory, FB-GRU per-300 forms best on train, interpolation and extrapolation evaluations, with results relative to the supervised 301 baseline similar to FB trained on Markov states in [97]. Aggregate test performance approximately 302 matches TD3-GRU aggregate test performance despite not seeing rewards during training. FB-TF exhibits the best interpolation-to-train ratio, and FB-GRU the best extrapolation-to-train ratio. 303



320 Figure 4: Zero-shot dynamics generalisation on ExORL. (Left) Aggregate performance across all 321 ExORL tasks and domains normalised w.r.t. TD3-GRU performance, averaged over 5 seeds. We train on dynamics $\{0.5x, 1.5x\}$ their typical values and evaluate on 1.0x (interpolation) and 2.0x 322 (extrapolation). (*Right*) The ratios of interpolation and extrapolation performance to train perfor-323 mance.



Figure 5: **Zero-shot environment generalisation on ExORL.** Aggregate performance across all tasks in the train environments (Walker, Quadruped), and the test environment (Cheetah), averaged across 5 seeds.

4.4 ZERO-SHOT ENVIRONMENT GENERALISATION

345 We report the aggregate performance of all FB-based algorithms on our zero-shot environment gen-346 eralisation experiments in Figure 5. Here, FB-GRU performs best on the training environments, but 347 poorest on the testing environment, with FB performing similarly poorly on the testing environment. FB-S4 improves performance on the testing environment by \sim 4x over FB, but at the cost of reduc-348 ing training environment performance by half. FB-TF improves both training performance by 5% 349 over FB and triples test performance. We emphasise that, although FPs with memory do improve 350 zero-shot environment generalisation performance in some cases, the absolute returns remain low (a 351 max of 33 for FB-S4 out of a possible 1000) suggesting their is significant room for improvement. 352

353 354 355

377

324

325 326

327

328

330

331

332 333

334

335

340

341

342 343

344

5 DISCUSSION AND LIMITATIONS

Context length. In Section 4 we train agents with a context length L = 64 timesteps, which is the 356 maximum context length we could afford with our computational budget³. We see two limitations 357 with this. First, we have assumed the dynamics context and task for all of our experiments can be 358 inferred from this context, but it is not clear that this is the case. Indeed, TD3-GRU with reward 359 supervision and a maximally exploratory dataset does not match its performance with Markov states 360 from [97]. Second, successful episodes run for a minimum of 200 timesteps (as in PendulumHard) 361 and a maximum of 1000 timesteps (as in ExORL), meaning we never train the memory model 362 over full episodes, nor can we reliably maintain an episode's full trajectory in context at test-time. This introduces situations where the hidden state will be erroneously initialised mid-episode during 364 training, creating well-established theoretical issues for memory models in RL [68], though these 365 are yet to prove critical empirically [73, 66].

The obvious solution to these problems, were it available to us, would be to increase L until it is the maximum episode length, and train for longer as in [31, 62, 68]. However, even if we were to do this, any real-world deployment may induce episodes longer than this assumed max length, or indeed we may wish to operate in the non-episodic, continual setting. The existing literature implicitly assumes that if L is very large such issues will resolve themselves, but this is not clear to us. Exploring how to deal with such situations is an important future research direction.

Datasets. As outlined in Section 4.1, we train all methods on datasets pre-collected with RND [11] which is a highly exploratory algorithm designed for maximising data heterogeneity. However, deploying such an algorithm in any real setting may be costly, time-consuming or dangerous. As a result, our proposals are more likely to be trained on real-world datasets that are smaller and more

³Our shared resource limits us to a maximum run length of 24 hours per GPU, and the ExORL runs took approximately 20 hours on one A100. See Appendix A.4 for more detail.

homogeneous. It is not clear how our specific proposals will interact with such datasets. If, for
example, the dataset only exhibits parts of the state space from which the dynamics cannot be wellinferred, like a robot stuck stationary, then we would expect our proposals to struggle. Indeed, with
poor coverage of the state-action space, we would expect to see the OOD pathologies seen in the
single-task Offline RL setting [52, 59]. That said, the proposals of [45] for conducting zero-shot RL
from real-world datasets could be integrated into our proposals easily, and may help.

384

386

6 RELATED WORK

387 **Generalist policy pre-training** FPs build upon successor representations [20], universal value 388 function approximators [85], successor features [6] and successor measures [10]. The state-389 of-the-art methods instantiate these ideas as either universal successor features (USFs) [11] or 390 forward-backward (FB) representations [96, 97], with recent work showing they can be trained 391 on low quality datasets [45], or used to perform a range of imitation learning techniques efficiently [80]. A representation learning method is required to learn the features for USFs, with 392 past works using inverse curiosity modules [79], diversity methods [60, 38], Laplacian eigenfunc-393 tions [101], or contrastive learning [16]. No works have yet explored the generalisation capac-394 ity of FPs to unseen dynamics. Two concurrent lines of work on goal-conditioned RL and in-395 context RL also seek to build generalist policies. Goal-conditioned RL methods train policies to 396 reach any goal state from any other goal state [75, 63, 104, 26, 99], but lack the ability to gen-397 eralise to tasks with "dense" reward functions, like those on the locomotion tasks in ExORL. In-398 context RL methods train policies using sequence models conditioned on reward-labelled histories 399 [14, 43, 58, 83, 110, 13, 30, 88, 103, 102, 31, 62, 94, 23], but, unlike FPs, do not have a robust 400 mechanism for training without access to rewards.

Dynamics Generalisation Dynamics generalisation in RL is a well-established problem [50, 65, 74]. Common remedies include: data augmentation [81, 5, 106, 37, 36, 55], domain randomisation [95, 21, 46, 47, 77], learning context-aware policies [87, 57, 9, 44], and meta-learning [15, 83, 27, 71, 72]. Our work is most similar to those that tackle dynamics generalisation by conditioning policies on dynamics inferred with a memory model [73, 18]. Where these methods are concerned with generalising to one unseen task in unseen dynamics contexts, our method can generalise to more than one unseen tasks in unseen dynamics contexts.

408 409

410

7 CONCLUSION

411 In this paper, we explored augmenting Foundation Policies (FPs) with memory models to allow 412 them to condition policies on a dynamics context inferred from a history of observations and ac-413 tions. We evaluated our proposals with attention, state-space, and RNN-based memory models on 414 POPGym, a memory benchmark, and ExORL, an unsupervised RL benchmark. Our results show 415 that GRUs achieve the best generalisation to unseen tasks and dynamics for a given recurrent state 416 size, approaching the performance of a supervised baseline that has access to task information during training and significantly outperforming memory-free FPs. We believe our proposals represent 417 a further step toward the development of generalist, adaptive agents. 418

- 419
- 420
- 421
- 422 423
- 424
- 425
- 426
- 427
- 428 429
- 429
- 431

432 REFERENCES

434

435

436 437

438

439

440

441 442

443

444 445

446

447 448

449

450

451

452

453

454 455

456

457

458

459

460

461 462

463

464

465

466

467 468

469

470

471

472

473 474

475

476

477

478

479

480 481

482

483

484

- [1] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in neural information processing systems*, 34:29304–29320, 2021.
- [2] Karl Johan Åström. Optimal control of markov processes with incomplete state information. *Journal of mathematical analysis and applications*, 10(1):174–205, 1965.
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv* preprint arXiv:1607.06450, 2016.
 - [4] Bram Bakker. Reinforcement learning with long short-term memory. *Advances in neural information processing systems*, 14, 2001.
 - [5] Philip J Ball, Cong Lu, Jack Parker-Holder, and Stephen Roberts. Augmented world models facilitate zero-shot dynamics generalization from a single offline environment. In *International Conference on Machine Learning*, pp. 619–629. PMLR, 2021.
 - [6] André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt, and David Silver. Successor features for transfer in reinforcement learning. Advances in neural information processing systems, 30, 2017.
 - [7] Andrew G Barto, Richard S Sutton, and Charles W Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE transactions on systems, man, and cybernetics*, (5):834–846, 1983.
 - [8] Jacob Beck, Kamil Ciosek, Sam Devlin, Sebastian Tschiatschek, Cheng Zhang, and Katja Hofmann. Amrl: Aggregated memory for reinforcement learning. In *International Conference on Learning Representations*, 2020.
 - [9] Michael Beukman, Devon Jarvis, Richard Klein, Steven James, and Benjamin Rosman. Dynamics generalisation in reinforcement learning via adaptive context-aware policies. Advances in Neural Information Processing Systems, 36, 2024.
 - [10] Léonard Blier, Corentin Tallec, and Yann Ollivier. Learning successor states and goaldependent values: A mathematical viewpoint. arXiv preprint arXiv:2101.07123, 2021.
 - [11] Diana Borsa, André Barreto, John Quan, Daniel Mankowitz, Rémi Munos, Hado Van Hasselt, David Silver, and Tom Schaul. Universal successor features approximators. *arXiv preprint arXiv:1812.07626*, 2018.
 - [12] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
 - [13] Yevgen Chebotar, Quan Vuong, Alex Irpan, Karol Hausman, Fei Xia, Yao Lu, Aviral Kumar, Tianhe Yu, Alexander Herzog, Karl Pertsch, et al. Q-transformer: Scalable offline reinforcement learning via autoregressive q-functions. arXiv preprint arXiv:2309.10150, 2023.
 - [14] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. Advances in neural information processing systems, 34:15084– 15097, 2021.
 - [15] Tao Chen, Adithyavairavan Murali, and Abhinav Gupta. Hardware conditioned policies for multi-robot transfer learning. Advances in Neural Information Processing Systems, 31, 2018.
 - [16] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- [17] Kyunghyun Cho. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.

490

491

492 493

494

495

496

497

498

499 500

501

502

504

505

507

509

510

511

512

513 514

515

516

517

518

519

521

522

523

524

527

528 529

530

531

532

534

536

- 486 [18] Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying 487 generalization in reinforcement learning. In International conference on machine learning, 488 pp. 1282-1289. PMLR, 2019.
 - [19] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. Advances in Neural Information Processing Systems, 35:16344–16359, 2022.
 - [20] Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural computation*, 5(4):613–624, 1993.
 - [21] Michael Dennis, Natasha Jaques, Eugene Vinitsky, Alexandre Bayen, Stuart Russell, Andrew Critch, and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised environment design. Advances in neural information processing systems, 33:13049–13061, 2020.
 - [22] Kenji Doya. Temporal difference learning in continuous time and space. Advances in neural information processing systems, 8, 1995.
 - [23] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. Rl²: Fast reinforcement learning via slow reinforcement learning. arXiv preprint arXiv:1611.02779, 2016.
 - [24] Jeffrey L Elman. Finding structure in time. Cognitive science, 14(2):179–211, 1990.
 - [25] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. arXiv preprint arXiv:1802.06070, 2018.
 - [26] Benjamin Eysenbach, Tianjun Zhang, Sergey Levine, and Russ R Salakhutdinov. Contrastive learning as goal-conditioned reinforcement learning. Advances in Neural Information Processing Systems, 35:35603–35620, 2022.
 - [27] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In International conference on machine learning, pp. 1126–1135. PMLR, 2017.
 - [28] Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. Advances in neural information processing systems, 34:20132–20145, 2021.
 - [29] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In International conference on machine learning, pp. 2052–2062. PMLR, 2019.
 - [30] Hiroki Furuta, Yutaka Matsuo, and Shixiang Shane Gu. Generalized decision transformer for offline hindsight information matching. arXiv preprint arXiv:2111.10364, 2021.
 - [31] Jake Grigsby, Linxi Fan, and Yuke Zhu. Amago: Scalable in-context reinforcement learning for adaptive agents. International Conference on Learning Representations, 2023.
 - [32] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. arXiv preprint arXiv:2111.00396, 2021.
 - [33] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. Advances in Neural Information Processing Systems, 35:35971–35983, 2022.
- 535 [34] Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured state spaces. Advances in Neural Information Processing Systems, 35:22982–22994, 2022.
 - [35] Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes, 2015.

540 [36] Nicklas Hansen and Xiaolong Wang. Generalization in reinforcement learning by soft data 541 augmentation. In 2021 IEEE International Conference on Robotics and Automation (ICRA), 542 pp. 13611–13617. IEEE, 2021. 543 [37] Nicklas Hansen, Hao Su, and Xiaolong Wang. Stabilizing deep q-learning with convnets 544 and vision transformers under data augmentation. Advances in neural information processing systems, 34:3680-3693, 2021. 546 547 [38] Steven Hansen, Will Dabney, Andre Barreto, Tom Van de Wiele, David Warde-Farley, and 548 Volodymyr Mnih. Fast task inference with variational intrinsic successor features. arXiv 549 preprint arXiv:1906.05030, 2019. 550 [39] Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, 551 David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. Array 552 programming with numpy. Nature, 585(7825):357-362, 2020. 553 554 [40] Natalia Hernandez-Gardiol and Sridhar Mahadevan. Hierarchical memory-based reinforce-555 ment learning. Advances in Neural Information Processing Systems, 13, 2000. 556 [41] Sepp Hochreiter and Jurgen Schmidhuber. Long short-term memory. Neural computation, 9 (8):1735–1780, 1997. 558 559 [42] John D Hunter. Matplotlib: A 2d graphics environment. Computing in science & engineering, 9(03):90-95, 2007. 561 [43] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big 562 sequence modeling problem. Advances in neural information processing systems, 34:1273-563 1286, 2021. 565 [44] Scott Jeen and Jonathan M. Cullen. Dynamics generalisation with behaviour foundation 566 models. RL Conference Workshop on Training Agents with Foundation Models, 2024. 567 [45] Scott Jeen, Tom Bewley, and Jonathan M. Cullen. Zero-shot reinforcement learning from low 568 quality data. Advances in Neural Information Processing Systems 38, 2024. 569 570 [46] Minqi Jiang, Michael Dennis, Jack Parker-Holder, Jakob Foerster, Edward Grefenstette, and 571 Tim Rocktäschel. Replay-guided adversarial environment design. Advances in Neural Infor-572 mation Processing Systems, 34:1884–1897, 2021. 573 [47] Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay. In Interna-574 tional Conference on Machine Learning, pp. 4940-4950. PMLR, 2021. 575 576 [48] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers 577 are rnns: Fast autoregressive transformers with linear attention. In International conference 578 on machine learning, pp. 5156–5165. PMLR, 2020. 579 [49] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv 580 preprint arXiv:1412.6980, 2014. 581 582 [50] Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of zero-shot 583 generalisation in deep reinforcement learning. Journal of Artificial Intelligence Research, 76: 584 201-264, 2023. 585 [51] Vijay Konda and John Tsitsiklis. Actor-critic algorithms. Advances in neural information 586 processing systems, 12, 1999. 588 [52] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-589 policy q-learning via bootstrapping error reduction. In Advances in Neural Information Pro-590 cessing Systems, volume 32. Curran Associates, Inc., 2019. [53] Aviral Kumar, Rishabh Agarwal, Xinyang Geng, George Tucker, and Sergey Levine. Of-592 fline q-learning on diverse multi-task data both scales and generalizes. arXiv preprint arXiv:2211.15144, 2022.

594 [54] Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steiger-595 wald, DJ Strouse, Steven Hansen, Angelos Filos, Ethan Brooks, et al. In-context reinforce-596 ment learning with algorithm distillation. arXiv preprint arXiv:2210.14215, 2022. 597 [55] Misha Laskin, Kimin Lee, Adam Stooke, Lerrel Pinto, Pieter Abbeel, and Aravind Srinivas. 598 Reinforcement learning with augmented data. Advances in neural information processing systems, 33:19884-19895, 2020. 600 601 [56] Jonathan Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and 602 Emma Brunskill. Supervised pretraining can learn in-context reinforcement learning. Ad-603 vances in Neural Information Processing Systems, 36, 2023. 604 [57] Kimin Lee, Younggyo Seo, Seunghyun Lee, Honglak Lee, and Jinwoo Shin. Context-aware 605 dynamics model for generalization in model-based reinforcement learning. In International 606 Conference on Machine Learning, pp. 5757–5766. PMLR, 2020. 607 608 [58] Kuang-Huei Lee, Ofir Nachum, Mengjiao Yang, Lisa Lee, Daniel Freeman, Winnie Xu, Sergio Guadarrama, Ian Fischer, Eric Jang, Henryk Michalewski, et al. Multi-game decision 609 transformers. Advances in neural information processing systems, 35, 2022. 610 611 [59] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: 612 Tutorial, review, and perspectives on open problems. arXiv preprint arXiv:2005.01643, 2020. 613 [60] Hao Liu and Pieter Abbeel. Aps: Active pretraining with successor features. In *International* 614 Conference on Machine Learning, pp. 6736–6747. PMLR, 2021. 615 616 [61] Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. Ad-617 vances in Neural Information Processing Systems, 34:18459–18473, 2021. 618 [62] Chris Lu, Yannick Schroecker, Albert Gu, Emilio Parisotto, Jakob Foerster, Satinder Singh, 619 and Feryal Behbahani. Structured state space models for in-context reinforcement learning. 620 Advances in Neural Information Processing Systems, 36, 2024. 621 622 [63] Yecheng Jason Ma, Jason Yan, Dinesh Jayaraman, and Osbert Bastani. How far i'll go: 623 Offline goal-conditioned reinforcement learning via f-advantage regression. arXiv preprint 624 arXiv:2206.03023, 2022. 625 [64] Wes McKinney et al. pandas: a foundational python library for data analysis and statistics. 626 Python for high performance and scientific computing, 14(9):1–9, 2011. 627 628 [65] Ishita Mediratta, Qingfei You, Minqi Jiang, and Roberta Raileanu. A study of generalization in offline reinforcement learning. In The Twelfth International Conference on Learning 629 Representations, 2024. 630 631 [66] Lingheng Meng, Rob Gorbet, and Dana Kulić. Memory-based deep reinforcement learning 632 for pomdps. In 2021 IEEE/RSJ international conference on intelligent robots and systems 633 (IROS), pp. 5619–5626. IEEE, 2021. 634 [67] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G 635 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 636 Human-level control through deep reinforcement learning. Nature, 518(7540):529-533, 637 2015. 638 639 [68] Steven Morad, Ryan Kortvelesy, Matteo Bettini, Stephan Liwicki, and Amanda Pro-640 rok. Popgym: Benchmarking partially observable reinforcement learning. arXiv preprint arXiv:2303.01859, 2023. 641 642 [69] Steven Morad, Ryan Kortvelesy, Stephan Liwicki, and Amanda Prorok. Reinforcement learn-643 ing with fast and forgetful memory. Advances in Neural Information Processing Systems, 36, 644 2024.645 [70] Steven Morad, Chris Lu, Ryan Kortvelesy, Stephan Liwicki, Jakob Foerster, and Amanda 646 Prorok. Revisiting recurrent reinforcement learning with memory monoids. arXiv preprint 647 arXiv:2402.09900, 2024.

652

653 654

655

656

657

658

659

661

662

663

664

665 666

667

668 669

670

671

672

673

674

675 676

677

678

679

680

681

682 683

684

685

686

687

688 689

690 691

692

- [71] Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Learning to adapt in dynamic, real-world environments through meta-reinforcement learning. *arXiv preprint arXiv:1803.11347*, 2018.
 - [72] Anusha Nagabandi, Chelsea Finn, and Sergey Levine. Deep online learning via metalearning: Continual adaptation for model-based rl. *arXiv preprint arXiv:1812.07671*, 2018.
 - [73] Tianwei Ni, Benjamin Eysenbach, and Ruslan Salakhutdinov. Recurrent model-free rl can be a strong baseline for many pomdps. *arXiv preprint arXiv:2110.05038*, 2021.
 - [74] Charles Packer, Katelyn Gao, Jernej Kos, Philipp Krähenbühl, Vladlen Koltun, and Dawn Song. Assessing generalization in deep reinforcement learning. arXiv preprint arXiv:1810.12282, 2018.
 - [75] Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. Hiql: Offline goalconditioned rl with latent states as actions. *Advances in Neural Information Processing Systems*, 36, 2024.
 - [76] Seohong Park, Tobias Kreiman, and Sergey Levine. Foundation policies with hilbert representations. *International Conference on Machine Learning*, 2024.
 - [77] Jack Parker-Holder, Minqi Jiang, Michael Dennis, Mikayel Samvelyan, Jakob Foerster, Edward Grefenstette, and Tim Rocktäschel. Evolving curricula with regret-based environment design. In *International Conference on Machine Learning*, pp. 17473–17498. PMLR, 2022.
 - [78] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary De-Vito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
 - [79] Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *International conference on machine learning*, pp. 2778–2787. PMLR, 2017.
 - [80] Matteo Pirotta, Andrea Tirinzoni, Ahmed Touati, Alessandro Lazaric, and Yann Ollivier. Fast imitation via behavior foundation models. In *International Conference on Learning Repre*sentations, 2024.
 - [81] Roberta Raileanu, Max Goldstein, Denis Yarats, Ilya Kostrikov, and Rob Fergus. Automatic data augmentation for generalization in deep reinforcement learning. *arXiv preprint arXiv:2006.12862*, 2020.
 - [82] Dhruv Ramani. A short survey on memory based reinforcement learning. *arXiv preprint arXiv:1904.06736*, 2019.
 - [83] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *Transactions of Machine Learning Research*, 2022.
 - [84] Michel F Sanner et al. Python: a programming language for software integration and development. J Mol Graph Model, 17(1):57–61, 1999.
 - [85] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approximators. In *International conference on machine learning*, pp. 1312–1320. PMLR, 2015.
 - [86] Juergen Schmidhuber. Reinforcement learning upside down: Don't predict rewards–just map them to actions. *arXiv preprint arXiv:1912.02875*, 2019.
- [87] Younggyo Seo, Kimin Lee, Ignasi Clavera Gilaberte, Thanard Kurutach, Jinwoo Shin, and Pieter Abbeel. Trajectory-wise multiple choice learning for dynamics generalization in reinforcement learning. *Advances in Neural Information Processing Systems*, 33:12968–12979, 2020.
- [88] Max Siebenborn, Boris Belousov, Junning Huang, and Jan Peters. How crucial is transformer in decision transformer? arXiv preprint arXiv:2211.14655, 2022.

702 [89] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van 703 Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanc-704 tot, et al. Mastering the game of go with deep neural networks and tree search. Nature, 529 705 (7587):484-489, 2016. 706 [90] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game 708 of go without human knowledge. Nature, 550(7676):354-359, 2017. 709 710 [91] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. A general 711 reinforcement learning algorithm that masters chess, shogi, and go through self-play. Science, 712 362(6419):1140-1144, 2018. 713 714 [92] Richard Sutton and Andrew Barto. Reinforcement Learning: An Introduction. The MIT 715 Press, second edition, 2018. 716 [93] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David 717 Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. 718 arXiv preprint arXiv:1801.00690, 2018. 719 720 [94] Adaptive Agent Team, Jakob Bauer, Kate Baumli, Satinder Baveja, Feryal Behbahani, 721 Avishkar Bhoopchand, Nathalie Bradley-Schmieg, Michael Chang, Natalie Clay, Adrian Collister, et al. Human-timescale adaptation in an open-ended task space. arXiv preprint 722 arXiv:2301.07608, 2023. 723 724 [95] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. 725 Domain randomization for transferring deep neural networks from simulation to the real 726 world. In 2017 IEEE/RSJ international conference on intelligent robots and systems (IROS), 727 pp. 23-30. IEEE, 2017. 728 [96] Ahmed Touati and Yann Ollivier. Learning one representation to optimize all rewards. Ad-729 vances in Neural Information Processing Systems, 34:13–23, 2021. 730 731 [97] Ahmed Touati, Jérémy Rapin, and Yann Ollivier. Does zero-shot reinforcement learning 732 exist? In The Eleventh International Conference on Learning Representations, 2023. 733 [98] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, 734 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural informa-735 tion processing systems, 30, 2017. 736 [99] Tongzhou Wang, Antonio Torralba, Phillip Isola, and Amy Zhang. Optimal goal-reaching 737 reinforcement learning via quasimetric learning. In International Conference on Machine 738 Learning, pp. 36411–36430. PMLR, 2023. 739 740 [100] Daan Wierstra and Jürgen Schmidhuber. Policy gradient critics. In European Conference on 741 Machine Learning, pp. 466-477. Springer, 2007. 742 [101] Yifan Wu, George Tucker, and Ofir Nachum. The laplacian in rl: Learning representations 743 with efficient approximations. arXiv preprint arXiv:1810.04586, 2018. 744 745 [102] Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu, Ding Zhao, Joshua Tenenbaum, and 746 Chuang Gan. Prompting decision transformer for few-shot policy generalization. In Proceed-747 ings of the 39th International Conference on Machine Learning, pp. 24631–24645, 17–23 Jul 748 2022. 749 [103] Taku Yamagata, Ahmed Khalil, and Raul Santos-Rodriguez. Q-learning decision transformer: 750 Leveraging dynamic programming for conditional sequence modelling in offline RL. In Pro-751 ceedings of the 40th International Conference on Machine Learning, volume 202, pp. 38989– 752 39007, 23-29 Jul 2023. 753 [104] Rui Yang, Lin Yong, Xiaoteng Ma, Hao Hu, Chongjie Zhang, and Tong Zhang. What is essen-754 tial for unseen goal generalization of offline goal-conditioned rl? In International Conference 755 on Machine Learning, pp. 39543-39571. PMLR, 2023.

- [105] Sherry Yang, Ofir Nachum, Yilun Du, Jason Wei, Pieter Abbeel, and Dale Schuurmans. Foun-dation models for decision making: Problems, methods, and opportunities. arXiv preprint arXiv:2303.04129, 2023.
- [106] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous control: Improved data-augmented reinforcement learning. arXiv preprint arXiv:2107.09645, 2021.
- [107] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement learning with prototypical representations. In International Conference on Machine Learning, pp. 11920-11931. PMLR, 2021.
- [108] Denis Yarats, David Brandfonbrener, Hao Liu, Michael Laskin, Pieter Abbeel, Alessandro Lazaric, and Lerrel Pinto. Don't change the algorithm, change the data: Exploratory data for offline reinforcement learning. arXiv preprint arXiv:2201.13425, 2022.
- [109] Marvin Zhang, Zoe McCarthy, Chelsea Finn, Sergey Levine, and Pieter Abbeel. Learning deep neural network policies with continuous memory states. In 2016 IEEE international conference on robotics and automation (ICRA), pp. 520-527. IEEE, 2016.
- [110] Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In international conference on machine learning, pp. 27042–27059. PMLR, 2022.

810	A
811	-
812	
813	A
814	
815	
816	
817	
818	
819	
820	
821	B
822	
823	
824	
825	
826	
827	-
828	C
829	
830	
831	
832	
833	
834	
835	
836	
837	
838	
839	
840	
841	
842	
843	
844	
845	
846	
847	
848	
849	
850	
851	
852	
853	
854	
855	
856	

APPENDICES

A	Experimental Details	17
	A.1 POPGym	17
	A.2 ExORL	17
	A.3 Evaluation Protocol	18
	A.4 Computational Resources	18
В	Implementation Details	18
	B.1 Foundation Policies	18
	B.2 TD3-GRU	19
	B.3 Code References	20
С	Extended Results	21

C Extended Results

864 A EXPERIMENTAL DETAILS

А.1 РОРБУМ

867 868

866

870

871

872

We consider 5 environments from the POPGym benchmark [68] which is built atop the OpenAI gym [12]. Each tests the agents ability to summarise a trajectory of observations and actions into a Markov state for use in solving one downstream task. Following [62] we only consider the *hard* variants, because the other variants are considered too straightforward.

873 Stateless CartPole Hard. The cartpole environment from [7], but with the angular and linear
874 positions removed from the observation. The agent must integrate to compute positions from ve875 locity and balance the pole atop the cart to receive reward. The *hard* variant requires the pole to be
876 balanced for 600 timesteps (the *easy* and *medium* variants require the pole to be balanced for fewer
877 timesteps).

878 Noisy Stateless CartPole Hard. The same as Stateless CartPole Hard but with Gaussian noise 879 added to observations. The *hard* variant sets the standard deviation of the noise $\sigma = 0.3$ (the *easy* 880 and *medium* set $\sigma = 0.1$ and $\sigma = 0.2$ respectively).

881 Stateless Pendulum Hard. The swing-up pendulum [22] with the angular position information
 removed. The agent must integrate to compute positions from velocity and swing the pendulum up
 to receive reward. The *hard* variant requires the pendulum to be balanced for 200 timesteps (the
 easy and *medium* variants require the pole to be balanced for fewer timesteps).

- Noisy Stateless Pendulum Hard. The same as Stateless Pendulum Hard but with Gaussian noise added to observations. The *hard* variant sets the standard deviation of the noise $\sigma = 0.3$ (the *easy* and *medium* set $\sigma = 0.1$ and $\sigma = 0.2$ respectively).
- **Repeat Previous Hard.** Observations contain one of four values. The agent is rewarded for outputting the observation from some constant k timesteps ago, i.e. observation o_{t-k} at time t. The hard variant sets k = 64 (the easy and medium variants set k < 64).

892 893

894

A.2 EXORL

We consider 4 environments (three locomotion and one goal-directed) from the ExORL benchmark [108] which is built atop the DeepMind Control Suite [93]. We occlude their states by removing all velocity components, similar to [73, 66]. Environments are visualised here: https://www. youtube.com/watch?v=rAai4QzcYbs. The domains are summarised in Table 1.

Walker-Occluded. A two-legged robot required to perform locomotion starting from bent-kneed position. The observation and action spaces are 17 and 6-dimensional respectively (after occlusion), consisting of joint torques and positions. ExORL provides 4 tasks stand, walk, run and flip. The reward function for stand motivates straightened legs and an upright torso; walk and run are supersets of stand including reward for small and large degrees of forward velocity; and flip motivates angular velocity of the torso after standing. Rewards are dense.

Quadruped-Occluded. A four-legged robot required to perform locomotion inside a 3D maze. The observation and action spaces are 67 and 12-dimensional respectively (after occlusion), consisting of joint torques and positions. We evaluate on 4 tasks stand, run, walk and jump. The reward function for stand motivates a minimum torso height and straightened legs; walk and run are supersets of stand including reward for small and large degrees of forward velocity; and jump adds a term motivating vertical displacement to stand. Rewards are dense.

911 **Maze-Occluded.** A 2D maze with four rooms where the task is to move a point-mass to one of 912 the rooms. The observation and action spaces are both 2-dimensional (after occlusion); the obser-913 vation space consists of x, y positions of the mass, the action space is the x, y tilt angle. ExORL 914 provides four reaching tasks top left, top right, bottom left and bottom right 915 corresponding to each room. We add four other goals in each room following [97] to provide a total 916 of 20 goal reaching tasks. The mass is always initialised in the top left and the reward is proportional 917 to the distance from the goal, though is sparse i.e. it only registers once the agent is reasonably close 918 to the goal.

918 Cheetah-Occluded. A running two-legged robot. The observation and action spaces are 10 and 919 6-dimensional respectively (after occlusion), consisting of positions of robot joints. We evaluate on 920 4 tasks: walk, walk backward, run and run backward. Rewards are linearly propor-921 tional either a forward or backward velocity-2 m/s for walk and 10 m/s for run.

922 923

924

A.3 EVALUATION PROTOCOL

925 We evaluate the cumulative reward (hereafter called score) achieved by all methods across three 926 seeds in POPGym and 5 seeds in ExORL. We report task scores as per the best practice recom-927 mendations of [1]. Concretely, we run each algorithm for 500k learning steps (1m for ExORL), 928 evaluating task scores at checkpoints of 20,000 steps. At each checkpoint, we perform 10 rollouts, 929 record the score of each, and find the interquartile mean (IQM). We average across seeds at each checkpoint. We extract task scores from the learning step for which the all-task IQM is maximised 930 across seeds. Results are reported with their associated standard deviation. Aggregation across 931 tasks, domains and datasets is always performed by evaluating the IQM. 932

933 934

935

COMPUTATIONAL RESOURCES A.4

936 We train our models on NVIDIA A100 GPUs. Training TD3-GRU to solve one task on one GPU takes approximately 6 hours for POPGym and 8 hours for ExORL. One run of FB-stack and SF-937 stack on one domain (for all tasks) takes approximately 3 hours for POPGym and 5 hours for ExORL 938 on one GPU. One run of the memory-based FPs on one domain (for all tasks) on one GPU in 939 approximately 20 hours. Note the POPGym experiments run for 500k learning steps, whereas the 940 EXORL experiments run for 1m learning steps. As a result, our core experiments on the EXORL 941 benchmark used approximately 65 GPU days of compute. 942

943 944

945 946

947

B IMPLEMENTATION DETAILS

B 1 FOUNDATION POLICIES

948 FB and HILP follow the implementations by [76] which follow [97], other than the batch size which 949 we reduce from 1024 to 512 to reduce the computational expense of each run without limiting 950 performance as per [45]. Hyperparameters are reported in Table 2. An illustration of a standard FP architecture is provided in Figure 6, for comparison with the FP with memory architecture in Figure 951 2. 952

953 Forward Representation $F(o, a, z) / \text{USF} \psi(o, a, z)$. Inputs $(o_{t-L;t}, a)$ and state-task pairs (o, z)954 are preprocessed by feedforward MLPs that embed their inputs into a 512-dimensional space. These 955 embeddings are concatenated and passed through a third feedforward MLP F / ψ which outputs a d-dimensional embedding vector. The Transformer memory model with Flash Attention follows the 956 exact implementation in [31]; the S4d memory model follows the exact implementation in [68], and 957 the GRU memory model follows the exact implementation provided by Torch. 958

959 960

961 Table 1: ExORL domain summary. *Dimensionality* refers to the relative size of state and action 962 spaces. Type is the task categorisation, either locomotion (satisfy a prescribed behaviour until the 963 episode ends) or goal-reaching (achieve a specific task to terminate the episode). Reward is the 964 frequency with which non-zero rewards are provided, where dense refers to non-zero rewards at every timestep and sparse refers to non-zero rewards only at positions close to the goal. Green and 965 966 red colours reflect the relative difficulty of these settings.

967	Environment	Dimensionality	Туре	Reward
968	Walker-Occluded	Low	Locomotion	Dense
969	Ouadruped-Occluded	High	Locomotion	Dense
970	Maze-Occluded	Low	Goal-reaching	Sparse
971	Cheetah-Occluded	Low	Locomotion	Dense

Table 2: FP Hyperparameters.					
Hyperparameter	Value				
Latent dimension d	50				
F / ψ dimensions	(1024, 1024)				
<i>B</i> dimensions	(512, 512)				
Preprocessor dimensions	(512, 512)				
Transformer heads	4				
Transformer / S4d model dimension	32				
GRU dimensions	(512, 512)				
Context length L	32 (Section 4.3), 64 (Sections 4.2 and 4.4)				
Frame stacking (FB & HILP)	4				
Std. deviation for policy smoothing σ	0.2				
Truncation level for policy smoothing	0.3				
Learning steps	1,000,000 (ExORL), 500,000 (POPGym)				
Batch size	512				
Optimiser	Adam [49]				
Learning rate	0.0001				
Discount γ	0.98				
Activations (unless otherwise stated)	ReLU				
Target network Polyak smoothing coefficient	0.01				
z-inference labels	10,000				
z mixing ratio	0.5				
HILP representation discount factor	0.98				
HILP representation expectile	0.5				
HILP representation target smoothing coefficient	0.005				

Backward Representation $B(o_{t-L:t})$ (for FB). Inputs are preprocessed by feedforward MLPs that embed their inputs into a 512-dimensional space then passed to the backward representation Bwhich is a feedforward MLP that outputs a d-dimensional embedding vector.

Actor $\pi(o_{t-L:t}, z)$. Inputs $(o_{t-L:t}, a)$ and state-task pairs (o, z) are preprocessed by feedforward MLPs that embed their inputs into a 512-dimensional space. These embeddings are concatenated and passed through a third feedforward MLP which outputs a *a*-dimensional vector, where *a* is the action-space dimensionality. A Tanh activation is used on the last layer to normalise their scale. As per [29]'s recommendations, the policy is smoothed by adding Gaussian noise σ to the actions during training.

1008 Misc. Layer normalisation [3] and Tanh activations are used in the first layer of all MLPs to 1009 standardise the inputs.

1010 1011 B.1.1 *z* Sampling

FPs require a method for sampling the task vector z at each learning step. [97] employ a mix of two methods, which we replicate:

- 1. Uniform sampling of z on the hypersphere surface of radius \sqrt{d} around the origin of \mathbb{R}^d ,
- 2. Biased sampling of z by passing states $s \sim D$ through the backward representation z = B(s). This also yields vectors on the hypersphere surface due to the L2 normalisation described above, but the distribution is non-uniform.
- 1020 We sample z 50:50 from these methods at each learning step.
- 1021

1023

1015

1016

1017

1018

1019

972

1022 B.2 TD3-GRU

We adopt the same implementation and hyperparameters as is used on the ExORL benchmark.
Hyperparameters are reported in Table 3. The memory module follows the implementation from [73] and uses a seperate encoder for the actor and critic.



[51]. The policy π selects an action a_t conditioned on a the current observation o_t , and the task vector z. The Q function formed by the USF ψ evaluates the action a_t given the current observation o_t and task z.

Critic(s). TD3 employs double Q networks, where the target network is updated with Polyak aver-aging via a momentum coefficient. The critics are feedforward MLPs that take a state-action pair (s, a) as input and output a value $\in \mathbb{R}^1$.

Actor. The actor is a standard feedforward MLP taking the state s as input and outputting an a-dimensional vector, where a is the action-space dimensionality. The policy is smoothed by adding Gaussian noise σ to the actions during training.

Misc. As is usual with TD3, layer normalisation [3] is applied to the inputs of all networks.

1	0	5	7
1	0	5	8

1059	Table 3: TD3-GRU hyperparameters.					
1060	Hyperparameter	Value				
1062	Critic dimensions	(1024, 1024)				
1063	Actor dimensions	(1024, 1024)				
1064	GRU dimensions	(512, 512)				
1004	Preprocessor dimensions	(512, 512)				
1065	Learning steps	1,000,000 (ExORL), 500,000 (POPGym)				
1066	Batch size	512				
1067	Optimiser	Adam				
1068	Learning rate	0.0001				
1069	Discount γ	0.98				
1070	Activations	ReLU				
1071	Target network Polyak smoothing coefficient	0.01				
1072	Std. deviation for policy smoothing σ	0.2				
1073	Truncation level for policy smoothing	0.3				

B.3 CODE REFERENCES

This work was enabled by: Python [84], NumPy [39], PyTorch [78], Pandas [64] and Matplotlib [42].



Figure 7: **Per-environment POPGym results.** The results are aggregated over 3 seeds, visualised by environment, and report the normalised MMER as with Table 4.



Figure 8: **Per-environment zero-shot dynamics generalisation results.** The results are aggregated over 5 seeds and all tasks in each environment, and show the normalised IQM w.r.t. TD3-GRU.

C EXTENDED RESULTS

We report a full breakdown of our results summarised in Sections 4.2, 4.3, 4.4. Table 4 reports results on POPGym from Section 4.2, Table 5 reports results on the zero-shot dynamics generalisation experiments from Section 4.3, and Table 6. Additionally, Figures 7 and 8 show plots where the results are aggregated by environment, and Figure 9 show plots where the zero-shot dynamics generalisation results are aggregated by task.

Table 4: **Full results on the POPGym environments (3 seeds).** We report the *unnormalised* meanmax epoch return (MMER) return \pm the standard deviation averaged over 3 seeds.

Environment	TD3-gru	FB-stack	FB-TF (ours)	FB-S4 (ours)	FB-GRU (ou
NoisyStatelessCartPoleHard	$0.156 \pm {\scriptstyle 0.011}$	$0.05\pm$ 0.0	$0.132 \pm {\scriptstyle 0.012}$	$0.16\pm {\scriptstyle 0.022}$	$0.196 \pm \scriptstyle 0.021$
NoisyStatelessPendulumHard	$0.543 \pm {\scriptstyle 0.004}$	$0.381 \pm$ 0.033	0.572 ± 0.005	0.572 ± 0.007	0.572 ± 0.009
RepeatPreviousHard	$-0.418\pm$ 0.012	$-0.455\pm ext{0.002}$	-0.441 ± 0.002	$-0.436\pm$ 0.016	-0.431 ± 0.0
StatelessCartPoleHard	1.0 ± 0.0	$0.017\pm$ 0.0	$0.515\pm$ 0.259	$1.0\pm$ 0.0	0.983 ± 0.025
StatelessPendulumHard	$0.8\pm$ 0.032	$0.436 \pm {\scriptstyle 0.033}$	$0.601\pm$ 0.015	$0.77\pm {\scriptstyle 0.032}$	0.742 ± 0.01



Figure 9: **Per-task zero-shot dynamics generalisation results.** The results are aggregated over 5 seeds, and show the normalised IQM w.r.t. TD3-GRU.

1189

1190

1191 1192 1193

1194 dataset-domain pair, we report the score at the step for which the all-task IQM is maximised when averaging across 5 seeds \pm the standard deviation. 1195 1196 Dynamics Environment Task TD3-gru HILP-stack FB-stack FB-TF (ours) FB-S4 (ours) FB-GRU (ours) All tasks 111 ± 98 21 ± 6 32 ± 7 25 ± 11 28 ± 3 $43 \pm {}_{18}$ 1197 Run 26 ± 6 7 ± 3 $13 \pm {}_2$ 5 ± 3 10 ± 2 $2\pm {\scriptstyle 14}$ Cheetah 1198 Run Backward 16 ± 9 5 ± 9 5 ± 6 $10 \pm s$ 5 ± 4 $12 \pm {}_{16}$ 233 ± 76 26 ± 33 77 ± 13 16 ± 7 $53 \pm {}_{13}$ 14 ± 67 Walk 1199 $36 \pm {}_{18}$ $89 \pm {}_{63}$ Walk Backward $196 \pm {\scriptstyle 122}$ $27 \pm {}_{18}$ 64 ± 48 $41 \pm {}_{18}$ 1200 26 ± 20 80 ± 45 18 ± 25 20 ± 18 Maze Multi goal $413 \pm {}_{346}$ 153 ± 40 $203 \pm _{31}$ 0.5x All tasks $279 \pm {}_{32}$ 123 ± 25 $327 \pm {}_{21}$ 541 ± 71 $330 \pm {}_{179}$ 1201 315 ± 86 $290 \pm {\scriptstyle 115}$ Jump 104 ± 76 $85 \pm {}_{25}$ 698 ± 60 $311 \pm {}_{273}$ Quadruped 1202 268 ± 88 104 ± 46 96 ± 73 $210 \pm _{74}$ 319 ± 100 178 ± 112 Run $322 \pm s_3$ Stand 311 ± 90 151 ± 55 469 ± 140 824 ± 108 415 ± 337 1203 289 ± 40 Walk 234 ± 113 123 ± 69 367 ± 110 362 ± 84 270 ± 58 446 ± 114 $89 \pm {\scriptstyle 13}$ 451 ± 34 $321 \pm {}_{24}$ All tasks 646 ± 224 $533 \pm {}_{46}$ 1204 Flip 570 ± 16 425 ± 46 $330 \pm {}_{28}$ $409 \pm {}_{182}$ $74 \pm {}_{25}$ 489 ± 123 Walker 1205 $169 \pm {}_{26}$ 167 ± 10 193 ± 10 Run 249 ± 11 33 ± 3 117 ± 8 Stand 847 ± 30 827 ± 98 $182 \pm {}_{29}$ 778 ± 15 $594 \pm {}_{49}$ $934 \pm {}_{17}$ 1206 723 ± 30 $43 \pm {}_{26}$ Walk 372 ± 196 425 ± 74 $243 \pm {\scriptstyle 16}$ $504 \pm {}_{66}$ 1207 All tasks 146 ± 206 36 ± 36 $72 \pm {}_{24}$ 40 ± 11 $47 \pm {}_{14}$ 54 ± 38 Run $37 \pm {}_{21}$ 4 ± 11 $25 \pm {}_{21}$ 5 ± 3 16 ± 7 $3\pm{}_{24}$ 1208 Cheetah Run Backward 11 ± 6 0 ± 17 6 ± 8 16 ± 14 11 ± 6 20 ± 19 1209 Walk 524 ± 254 102 ± 116 163 ± 66 31 ± 30 92 ± 45 49 ± 152 Walk Backward 255 ± 192 18 ± 16 59 ± 59 $92 \pm {}_{42}$ 53 ± 19 69 ± 69 1210 Maze Multi goal $354 \pm {}_{354}$ $11 \pm {}_{18}$ $58 \pm _{32}$ 35 ± 33 21 ± 17 154 ± 47 133 ± 29 1211 $1 \mathbf{x}$ All tasks 232 ± 77 179 ± 19 360 ± 23 445 ± 56 $403 \pm {}_{160}$ Jump $218 \pm {}_{95}$ 108 ± 34 135 ± 69 359 ± 58 557 ± 101 409 ± 210 Ouadruped 1212 106 ± 77 $203 \pm s_3$ Run $246 \pm s_7$ $76 \pm s_2$ 277 ± 44 296 ± 93 Stand $390 \pm {\scriptstyle 130}$ $313 \pm {}_{48}$ $135 \pm {}_{28}$ $565 \pm {}_{140}$ $677 \pm {\scriptstyle 125}$ $532 \pm _{299}$ 1213 168 ± 49 88 ± 59 248 ± 47 $286 \pm {}^{54}$ $362 \pm {}_{135}$ Walk 192 ± 103 1214 All tasks 519 ± 192 385 ± 109 74 ± 5 459 ± 42 301 ± 45 $565 \pm {}_{54}$ Flip 432 ± 55 315 ± 155 64 ± 7 409 ± 36 299 ± 51 480 ± 49 1215 Walker Run 267 ± 29 168 ± 60 26 ± 2 181 ± 21 $113 \pm {}_{16}$ $218 \pm {}_{24}$ 1216 $781 \pm s_0$ 871 ± 38 Stand 671 ± 113 168 ± 20 732 ± 41 554 ± 66 606 ± 46 234 ± 57 Walk 348 ± 186 47 ± 9 475 ± 96 654 ± 127 1217 52 ± 17 52 ± 55 All tasks 240 ± 286 13 ± 42 $56 \pm {}_{25}$ 23 ± 4 1218 Run $52 \pm {}_{29}$ $2 \pm {}_{41}$ $17 \pm {}_{22}$ 5 ± 3 18 ± 11 7 ± 41 Cheetah Run Backward 24 ± 36 6 ± 9 11 ± 9 8 ± 5 14 ± 5 $15 \pm {}_{13}$ 1219 127 ± 92 29 ± 11 $100 \pm {}_{31}$ 45 ± 204 Walk 715 ± 84 $14 \pm {}_{134}$ 1220 Walk Backward 428 ± 290 18 ± 10 25 ± 37 $48 \pm {}_{18}$ $74 \pm {\scriptstyle 27}$ $84 \pm s_2$ 250 ± 367 0±17 $67 \pm {}_{41}$ 39 ± 37 $16 \pm {}_{14}$ $141\pm {\scriptstyle 50}$ Maze Multi goal 1221 1.5x All tasks 217 ± 79 177 ± 27 108 ± 39 320 ± 26 264 ± 70 371 ± 135 1222 Jump 168 ± 69 120 ± 78 74 ± 75 255 ± 76 285 ± 96 297 ± 165 Quadruped 310 ± 101 200 ± 107 204 ± 94 Run 245 ± 139 119 ± 55 81 ± 94 1223 Stand 371 ± 175 291 ± 46 155 ± 64 489 ± 128 $348 \pm {}_{140}$ 546 ± 288 $190 \pm s_{6}$ 147 ± 103 60 ± 37 252 ± 45 265 ± 61 329 ± 103 Walk 1224 All tasks $364 \pm {}_{166}$ 222 ± 27 65 ± 4 $336 \pm {}_{26}$ 232 ± 50 $514 \pm {}_{17}$ 1225 Flip $273 \pm {\scriptstyle 27}$ 130 ± 56 49 ± 7 272 ± 27 $208 \pm {}_{44}$ 384 ± 19 Walker 204 ± 40 82 ± 19 23 ± 3 $136 \pm {}_{24}$ 96 ± 15 $232 \pm {}_{24}$ Run 1226 $630 \pm ss$ $461 \pm {}_{27}$ $148 \pm {}_{14}$ 419 ± 102 790 ± 31 Stand 547 ± 17 1227 Walk $198 \pm {}_{63}$ 38 ± 11 387 ± 55 $191 \pm {}^{53}$ $641 \pm {}_{46}$ 454 ± 82 19 ± 14 58 ± 35 All tasks 312 ± 320 23 ± 7 56 ± 19 24 ± 22 1228 Run 71 ± 55 6 ± 30 9 ± 3 5 ± 1 20 ± 14 8 ± 27 Cheetah Run Backward $\begin{array}{c}5\pm {}_3\\147\pm {}_{107}\end{array}$ 1229 19 ± 37 $6 \pm {}_{14}$ 9 ± 5 20 ± 7 13 ± 11 32 ± 11 91 ± 30 Walk 775 ± 83 $21 \pm {}_{31}$ $43 \pm {}_{21}$ 1230 Walk Backward 20 ± 79 17 ± 19 552 ± 394 35 ± 64 48 ± 27 93 ± 33 218 ± 355 65 ± 50 14 ± 12 Maze 0 ± 11 44 ± 37 131 ± 47 Multi goal 1231 132 ± 22 112 ± 27 341 ± 116 2x All tasks $215 \pm {}_{53}$ 270 ± 36 166 ± 85 1232 $268 \pm _{79}$ Jump $169 \pm {}_{140}$ 62 ± 38 $74 \pm ss$ $170 \pm {}_{118}$ $275 \pm {}_{159}$ Quadruped 221 ± 0 70 ± 66 86 ± 100 304 ± 102 140 ± 96 179 ± 117 Run 1233 353 ± 68 315 ± 193 $294 \pm {\scriptstyle 90}$ 142 ± 30 $223 \pm {\scriptstyle 197}$ $421 \pm {}_{210}$ Stand 1234 $70 \pm {}_{61}$ $57 \pm {}_{37}$ $111 \pm _{97}$ Walk $209 \pm {\scriptstyle 122}$ 222 ± 78 $351 \pm {}_{139}$ 60 ± 5 186 ± 39 $151 \pm {}_{32}$ $432 \pm {}_{24}$ $157 \pm {}_{28}$ All tasks 243 ± 118 1235 Flip 168 ± 40 105 ± 35 44 ± 6 $149 \pm {}_{41}$ $140 \pm {}_{28}$ 268 ± 32 Walker 1236 Run $142 \pm {}_{31}$ 65 ± 13 22 ± 4 78 ± 20 $69 \pm {}_{16}$ 229 ± 27 $143 \pm {}_{16}$ 341 ± 62 237 ± 51 656 ± 12 Stand 434 ± 66 364 ± 53 1237 319 ± 79 $79 \pm {}_{43}$ 139 ± 49 567 ± 48 Walk 29 ± 6 183 ± 57 1238

Table 5: Full results on the ExORL dynamics generalisation experiments (5 seeds). For each

1239

1240

Table 6: **Full results on the ExORL environment generalisation experiments (5 seeds).** For each dataset-domain pair, we report the score at the step for which the all-task IQM is maximised when averaging across 5 seeds \pm the standard deviation. The Baseline FB-GRU represents the scores of an FB-GRU model trained solely on Cheetah-Occluded.

Environment	Task	FB-stack	FB-TF (ours)	FB-S4d (ours)	FB-GRU (ours)	Baseline FB-GRU
	Walk	25 ± 7	$34 \pm s$	21 ± 2	22 ± 9	-
Wallron	Stand	$126 \pm {}_{28}$	$144 \pm {}_{23}$	$95 \pm {}_{32}$	$82 \pm {}_{18}$	-
walker	Run	17 ± 7	26 ± 5	18 ± 2	$15 \pm {}_{12}$	-
	Flip	$23\pm\mathrm{s}$	27 ± 6	21 ± 2	$23 \pm$ 10	-
	Walk	$60\pm {}_{64}$	$99 \pm {}_{39}$	28 ± 30	$110 \pm {}_{116}$	-
Ouedmand	Stand	$150\pm$ 79	$148 \pm {}_{84}$	$80 \pm {}_{29}$	257 ± 128	-
Quadruped	Run	$74 \pm {}_{48}$	$71 \pm {}_{26}$	$25\pm$ 75	$63 \pm {}_{56}$	-
	Jump	71_{35}	$58 \pm _{73}$	$88 \pm _{72}$	$190 \pm {}_{92}$	-
	Walk	2 ± 3	$24 \pm {}_{12}$	$89 \pm {}_{24}$	10 ± 4	$38 \pm {}_{216}$
Cheetah	Walk Backward	12 ± 6	$38 \pm {}_{21}$	22 ± 7	$10 \pm {}_{12}$	3 ± 7
	Run	-	3 ± 2	16 ± 4	2 ± 1	8_{46}
	Run Backward	2 ± 1	8 ± 4	4 ± 1	1 ± 2	0 ± 1