

Finer Behavioral Foundation Models via Auto-Regressive Features and Advantage Weighting

Edoardo Cetin, Ahmed Touati, Yann Ollivier

December 6, 2024

Abstract

The forward-backward representation (FB) is a recently proposed framework (Touati et al., 2023; Touati & Ollivier, 2021) to train *behavior foundation models* (BFMs) that aim at providing zero-shot efficient policies for any new task specified in a given reinforcement learning (RL) environment, without training for each new task. Here we address two core limitations of FB model training.

First, FB, like all successor-feature-based methods, relies on a *linear* encoding of tasks: at test time, each new reward function is linearly projected onto a fixed set of pre-trained features. This limits expressivity as well as precision of the task representation. We break the linearity limitation by introducing *auto-regressive features* for FB, which let fine-grained task features depend on coarser-grained task information. This can represent arbitrary nonlinear task encodings, thus significantly increasing expressivity of the FB framework.

Second, it is well-known that training RL agents from offline datasets often requires specific techniques. We show that FB works well together with such offline RL techniques, by adapting techniques from (Nair et al., 2020b; Cetin et al., 2024) for FB. This is necessary to get non-flatlining performance in some datasets, such as DMC Humanoid.

As a result, we produce efficient FB BFMs for a number of new environments. Notably, in the D4RL locomotion benchmark, the generic FB agent matches the performance of standard single-task offline agents (IQL, XQL). In many setups, the offline techniques are needed to get any decent performance at all. The auto-regressive features have a positive but moderate impact, concentrated on tasks requiring spatial precision and task generalization beyond the behaviors represented in the trainset.

Together, these results establish that generic, reward-free FB BFMs can be competitive with single-task agents on standard benchmarks, while suggesting that expressivity of the BFM is not a key limiting factor in the environments tested.

1 Introduction

The forward-backward representation (FB) is a recently proposed framework (Touati et al., 2023; Touati & Ollivier, 2021) to train behavior foundation models (BFMs) from offline data. BFMs promise to provide zero-shot efficient policies for any new task specified in a given reinforcement learning (RL) environment, beyond the tasks and behaviors in the training set. This contrasts with traditional offline RL and imitation learning approaches, which are trained to accomplish individual target tasks, with no mechanism to tackle new tasks without repeating the full training procedure.

The FB approach strives to learn an agent that recovers many possible behaviors in a given environment, based on learning *successor measure* representations, without any reward signal. After training, an FB agent can be prompted via several kinds of task description: an explicit reward function, a goal state, or even a single demonstration (Pirodda et al., 2023).

However, in its current formulation, FB has been shown effective only for toy problems and relatively simple locomotion tasks and when trained on undirected datasets collected via unsupervised exploration (Burda et al., 2019).

Here, we tackle two core limitations of the “vanilla” FB framework, namely, the difficulty to learn from complex offline datasets, and the linear correspondence between tasks and features. As a result, we can build high-performing FB BFMs for a series of new environments. Our main contributions are the following:

- We show that the vanilla FB policy optimization leads to poor performance when learning from datasets made of a few near-optimal examples for a few specific tasks. This failure is exactly analogous to naively using online RL algorithms in the offline setting, a well-studied problem (Fujimoto et al., 2019; Wu et al., 2019; Levine et al., 2020). This explains the poor performance of vanilla FB on the *D4RL benchmark*, as reported in recent unsupervised RL works (Park et al., 2024; Frans et al., 2024).
- Accordingly, we introduce a new policy optimization step for FB, to improve learning from offline datasets demonstrating complex behaviors. In particular, we integrate an improved version of *advantage-weighted regression* (Nair et al., 2020a), together with recent advancements from the offline RL literature (Cetin et al., 2024) and additional algorithmic refinements to FB (Section 3.2).

We show these changes are crucial to train FB with common offline datasets beyond pure RND exploration and scale to more challenging environments, often making the difference between near-zero and satisfactory performance (Section 4.3).

- We overcome a core theoretical limitation of FB and, more generally, of all *successor features* frameworks (Barreto et al., 2017; Borsa et al., 2018): their linear correspondence between reward functions and task representation vectors. Indeed, in these frameworks, at test time, the reward function is linearly projected onto a fixed set of pre-trained features.

This results in “reward blurring” and limits spatial precision in the task representation (Touati & Ollivier, 2021).

We introduce a new *auto-regressive* encoding of task features (Section 3.1), that breaks the linearity constraint by letting fine-grained task features depend on coarser-grained task information. This allows for universal approximation of any arbitrary task space (Appendix, Theorem A.3).

We show that auto-regressive features make a moderate but systematic difference when learning *new* test tasks far from ones considered to build the datasets, or for tasks requiring precise goal-reaching (eg, 15% relative increase for goal-reaching in the Jaco arm environment).

- With these improvements, we show that advantage-weighted autoregressive FB (FB-AWARE) extends FB performance to new environments such as Humanoid and the locomotion environments in the canonical *D4RL benchmark* (Fu et al., 2020). On the latter, FB-AWARE matches the performance of standard offline RL agents trained on a single task with full access to rewards (Section 4.3.3), further vindicating the use of behavior foundation models for zero-shot RL.

2 Preliminaries : Notation, Forward-Backward Framework for Behavioral Foundation Models

Markov decision process, notation. We consider a reward-free Markov decision process (MDP) $\mathcal{M} = (S, A, P, \gamma)$ with state space S , action space A , transition probabilities $P(s'|s, a)$ from state s to s' given action a , and discount factor $0 < \gamma < 1$ (Sutton & Barto, 2018). A policy π is a function $\pi: S \rightarrow \text{Prob}(A)$ mapping a state s to the probabilities of actions in A . Given $(s_0, a_0) \in S \times A$ and a policy π , we denote $\Pr(\cdot|s_0, a_0, \pi)$ and $\mathbb{E}[\cdot|s_0, a_0, \pi]$ the probabilities and expectations under state-action sequences $(s_t, a_t)_{t \geq 0}$ starting at (s_0, a_0) and following policy π in the environment, defined by sampling $s_t \sim P(s_t|s_{t-1}, a_{t-1})$ and $a_t \sim \pi(a_t|s_t)$. Given any reward function $r: S \rightarrow \mathbb{R}$, the Q -function of π for r is $Q_r^\pi(s_0, a_0) := \sum_{t \geq 0} \gamma^t \mathbb{E}[r(s_t)|s_0, a_0, \pi]$. The *value function* of π for r is $V_r^\pi(s) := \mathbb{E}_{a \sim \pi(s)} Q_r^\pi(s, a)$, and the *advantage function* is $A_r^\pi(s, a) := Q_r^\pi(s, a) - V_r^\pi(s)$.

We assume access to a dataset consisting of *reward-free* observed transitions (s_t, a_t, s_{t+1}) in the environment. We denote by ρ the distribution of states s_{t+1} in the training set.

Behavioral foundation models, zero-shot RL. A *behavioral foundation model* for a given reward-free MDP, is an agent that can produce approximately optimal policies for any reward function r specified at test time in the environment, without performing additional learning or fine-tuning for each new reward function. An early example of such a model includes universal successor features (SFs) (Borsa et al., 2018), which depend on a set of (sometimes handcrafted)

features: at test time, the reward is linearly projected onto the features, and a pre-trained policy is applied. Forward-backward representations (defined below) are another one, mathematically related to SFs. [Touati et al. \(2023\)](#) compares a number of variants of SFs and FB on a number of empirical problems.

The forward-backward framework. The FB framework ([Touati & Ollivier, 2021](#); [Touati et al., 2023](#)) is a theoretically and empirically well-supported way to train BFMs, based on learning an efficient representation of the *successor measures* M^π for various policies π . For each state-action $(s_0, a_0) \in S \times A$, this is a measure over states, describing the distribution of future states visited by starting at (s_0, a_0) and following policy π . It is defined as

$$M^\pi(s_0, a_0, X) := \sum_{t \geq 1} \gamma^t \Pr(s_t \in X \mid s_0, a_0, \pi) \quad (1)$$

for any subset $X \subset S$. M^π satisfies a measure-valued Bellman equation ([Blier et al., 2021](#)), which can be used to learn approximate parametric models of M .

[Touati & Ollivier \(2021\)](#) propose to learn a finite-rank parametric model of M , as follows:

$$M^{\pi_z}(s_0, a_0, X) \approx \int_{s \in X} F(s_0, a_0, z)^\top B(s) \rho(ds) \quad (2)$$

where ρ is the data distribution, where F and B take values in \mathbb{R}^d , where $z \in \mathbb{R}^d$ is a task encoding vector, and where

$$\pi_z(s) = \arg \max_a F(s, a, z)^\top z \quad (3)$$

is a parametric policy depending on z . F , B and π_z are learned at train time. At test time, given a reward function r , one estimates the task representation vector

$$z = \mathbb{E}_{s \sim \rho}[r(s)B(s)] \quad (4)$$

and then the policy π_z is applied.

The main result of [Touati & Ollivier \(2021\)](#) is that when (2)–(3) hold, then for *any* reward function r , the policy π_z so obtained is optimal. At test time, reward functions for FB may also be specified through an expert demonstration ([Pirodda et al., 2023](#)).

The full algorithm for FB training is provided in [Algo. 1](#) ([Appendix A.3](#)).

3 Breaking Some Key Limitations of the Forward-Backward Framework

3.1 Auto-Regressive Features for Non-Linear Task Encoding

Intuition for auto-regressive FB: nonlinear task encoding. Forward-backward (FB) representations and their predecessor, universal successor features

(SFs), attempt to solve zero-shot RL by linearly projecting new tasks (reward functions r) onto a set of features $B: S \rightarrow \mathbb{R}^d$. At test time, when facing a new reward function r , a task encoding $z \in \mathbb{R}^d$ is computed by $z = \mathbb{E}[r(s)B(s)]$ (FB) or $z = (\mathbb{E}[\phi(s)\phi(s)^\top])^{-1}\mathbb{E}[r(s)\phi(s)]$ (SFs). Then a pretrained policy π_z is applied.

FB aims at learning the features B that minimize the error from this process: B is obtained by a finite-rank approximation of the operator that sends a reward r to its Q -function. Bringing the FB loss to 0 (which requires infinitely many features) guarantees successful zero-shot RL for any reward function r . Theoretically, the features B in FB “most linearize” the computation of Q -functions, and empirically this brings better performance than other feature choices (Touati et al., 2023).

Still, even with the best features B , the task encoding z is *linear*, because $z = \mathbb{E}[r(s)B(s)]$ is linear in r : tasks are identified by the size- d vector of their correlations with a fixed set of d pre-trained features B .

The standard FB framework learns a rank d approximation by focusing on the main eigenvectors of the environment dynamics (Blier et al., 2021). Projecting the reward onto these eigenvectors can remove spatial precision, creating short-term reward blurring (Touati & Ollivier, 2021).

This is clearly suboptimal. Intuitively, if we first acquire information that the rewards are located in the top-left corner of S , we would like to use more precise features located in the top-left corner to better identify the reward function there.

Auto-regressive features make this possible, while still keeping most of the theoretical properties of plain FB. The idea is to compute the task encoding $z = \mathbb{E}[r(s)B(s)]$ progressively, and let the later-computed features B depend on the early components of z . We decompose z and B into K blocks $z = (z_1, z_2, \dots, z_K)$ and $B = (B_1, B_2, \dots, B_K)$. We first compute $z_1 = \mathbb{E}[r(s)B_1(s)]$ as in plain FB. But then we compute $z_2 = \mathbb{E}[r(s)B_2(s, z_1)]$ where the second block of features B_2 is allowed to depend on z_1 , thus conditioning the features on the task information provided by z_1 . This can be iterated: $z_3 = \mathbb{E}[r(s)B_3(s, z_1, z_2)]$, etc. The resulting vector $= (z_1, z_2, \dots, z_K)$ encodes the task in an auto-regressive manner, where the meaning of z_i depends on $z_{1:i-1}$.

Intuitively, z_1 provides a “coarse” task encoding by linear features. Then we compute a further, finer task encoding z_2 by computing the correlation of r with features B_2 that depend on the coarse task encoding z_1 . Hopefully the features B_2 can become more specialized and provide a better task encoding.

In practice, the main change with respect to plain FB training is that B depends on z . We now represent the successor measures M^{π_z} by $F(z)^\top B(z)$, instead of simply $F(z)^\top B$ which shares the same B for all policies. This allows for a better fit of the FB model. This is formalized in the next section and in Appendix A.

Contrary to plain FB, the task encoding $r \mapsto z$ becomes fully nonlinear: the set of tasks r represented exactly becomes a nonlinear submanifold of all possible tasks, instead of a d -dimensional subspace. Even with just two levels of features, this model is able to represent an arbitrary nonlinear mapping between reward functions r and task representations z (Appendix, Theorem A.3), instead of just

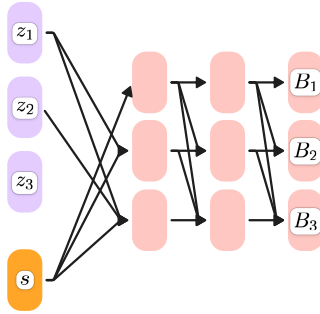


Figure 1: An auto-regressive architecture for $B(s, z)$. The i -th block of the output B only depends on blocks z_1, \dots, z_{i-1} of the input z . In each layer, the weights from each block to the lower-ranking blocks of the next layer have been removed. The state s is still fed to every block on the input layer.

linearly projecting the reward onto a fixed basis of features. This greatly extends the theoretical expressivity of the FB and successor feature frameworks.

This model also encodes a hierarchical prior on tasks, favoring tasks that can be described through a cascade of more and more specialized task features.

FB with auto-regressive task encoding: formal description. Auto-regressive features extend plain FB by letting B depend on z . In ordinary FB, this would be problematic, since the task encoding $z = \mathbb{E}[r(s)B(s)]$ used at test time becomes a fixed point equation if B depends on z . However, this fixed point equation can be handled easily if B has a hierarchical or auto-regressive structure.

Definition 3.1. A feature map $B: S \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ is called *auto-regressive* if, for any $1 \leq i \leq d$ and any $(s, z) \in S \times \mathbb{R}^d$, the i -th component of $B(s, z)$ only depends on (z_1, \dots, z_{i-1}) and not on (z_i, \dots, z_d) .

For such models, we can easily compute fixed-points values of the type $z = B(s, z)$, by first computing the component B_1 of the output (which does not depend on z), which determines z_1 , which allows us to compute the component B_2 of the output, which determines z_2 , etc.

In practice, auto-regressive models $B(s, z)$ can be built by splitting both the representation vector $z \in \mathbb{R}^d$ and the output $B(s, z) \in \mathbb{R}^d$ into k “auto-regressive groups” of dimension d/k . The first group $B_1(s, z)$ in the output of B is actually independent of z , and the i -th group $B_i(s, z)$ of the output of B only takes as inputs the previous groups z_1, \dots, z_{i-1} of z . At test time, this allows us to compute the fixed point $z = \mathbb{E}_{s \sim \rho}[r(s)B(s, z)]$ by first estimating the first group, $z_1 = \mathbb{E}_{s \sim \rho}[r(s)B_1(s)]$ similar to plain FB. Then the other groups are computed iteratively: $z_{i+1} = \mathbb{E}_{s \sim \rho}[r(s)B_{i+1}(s, z_1, \dots, z_i)]$. In the experiments, we focus on $k = 4$ or $k = 8$ auto-regressive groups.

We employ a network architecture (Fig. 1) in which each layer of B_i has access to the previous layers of all previous blocks $B_{1\dots i}$: this ensures good expressivity while preserving the auto-regressive property. This allows for efficient evaluation: this is implemented as masks on fully-connected layers for the full model B .

The following result extends the theorem from [Touati & Ollivier \(2021\)](#) for vanilla FB, to allow B to depend on z .

Theorem 3.2. *Assume we have learned representations $F: S \times A \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $B: S \times \mathbb{R}^d \rightarrow \mathbb{R}^d$, as well as a parametric family of policies π_z depending on $z \in \mathbb{R}^d$, satisfying*

$$\begin{cases} M^{\pi_z}(s_0, a_0, X) = \int_X F(s_0, a_0, z)^\top B(s, z) \rho(ds), & \forall s_0 \in S, a_0 \in A, X \subset S, z \in \mathbb{R}^d \\ \pi_z(s) = \arg \max_a F(s, a, z)^\top z, & \forall (s, a) \in S \times A, z \in \mathbb{R}^d. \end{cases} \quad (5)$$

Then the following holds. For any reward function r , if we can find a value $z_r \in \mathbb{R}^d$ such that

$$z_r = \mathbb{E}_{s \sim \rho}[r(s)B(s, z_r)] \quad (6)$$

then π_{z_r} is an optimal policy for reward r , and the optimal Q -function is $Q_r^(s, a) = F(s, a, z_r)^\top z_r$.*

Moreover, if B is auto-regressive, then the fixed point (6) always exists, and can be computed directly, by iteratively computing each component $z_i = \mathbb{E}[r(s)B_i(s, z_1, \dots, z_{i-1})]$ for $i = 1, \dots, d$.

Further theoretical properties and proofs are given in Appendix A. In particular, Theorem A.3 establishes that autoregressive FB with two blocks is enough to represent *any* task encoding $r \mapsto z_r$, while vanilla FB is contrained to a linear task encoding $r \mapsto z_r$. Thus, autoregressive FB is inherently more expressive.

Training F and B for this setup is similar to [Touati & Ollivier \(2021\)](#), except that B depends on z , which has some consequences for minibatch sampling, and results in higher variance. The details are given in Appendix A.2 and Algorithm 2. Training is based on the measure-valued Bellman equation satisfied by M^{π_z} : we plug in the model $M^{\pi_z} \approx F(z)^\top B(z)\rho$ in this equation and minimize the Bellman gaps.

There is little computational overhead compared to vanilla FB. In practice, we enforce the auto-regressive property via a single neural network B , by dropping a specific subset of the layer connections across neurons (Fig. 1). This implementation allows for efficient training: given access to any specific z and s , the output $B(s, z)$ can be computed in a single forward pass. On the other hand, the computation of the fixed point z_r from (6) requires several forward passes through the network B in Fig. 1, but this occurs only at test time when the reward function is known.

Auto-regressive FB models the successor features M^{π_z} via a model $M^{\pi_z} \approx F(z)^\top B(z)\rho$ with full dependency on z , versus $M^{\pi_z} \approx F(z)^\top B\rho$ for vanilla FB. This is more natural, especially for large γ . Indeed, for $\gamma \rightarrow 1$ we have $M^{\pi_z}(s_0, a_0, ds) = \frac{1}{1-\gamma}\mu_z(ds) + o(1/(1-\gamma))$ with μ_z the stationary distribution

of π_z , namely, approximately rank-one with z dependency on the s part. The vanilla FB model has no z dependency on the s part, only on the (s_0, a_0) part, which means all stationary distributions μ_z must be approximated using the shared features $B(s)$.

3.2 Better Offline Optimization for FB: Advantage Weighting and Other Improvements

Like off-policy algorithms designed for the offline RL setting, vanilla FB training appears prone to distribution shift, hindering its ability to scale and to learn on datasets exhibiting mixtures of behaviors for various tasks. Inspired by the recent analysis (Cetin et al., 2024), we propose a set of modifications to FB training to overcome these limitations.

Improved advantage weighting objective We introduce an alternative policy optimization step for FB, based on recent analysis and advancements in offline RL algorithms with policy constraints. We use an improved version of the *advantage-weighted* (AW) regression loss (Peng et al., 2019), commonly used in popular recent algorithms (Nair et al., 2020a; Kostrikov et al., 2022; Wang et al., 2020; Garg et al., 2023; Cetin et al., 2024). Following Nair et al. (2020a), a first version starts with sampling a batch of n transitions from the data, and updates the parametric policy π_θ to optimize

$$\arg \max_{\theta} \mathbb{E}_{(a_{1:n}, s_{1:n}) \sim \mathcal{B}} \left[\sum_{i=1}^n w(s_i, a_i) \log \pi_\theta(a_i | s_i) \right], \quad (7)$$

$$\text{where } w(s_i, a_i) = \frac{\exp(A_\phi(s_i, a_i)/\beta)}{\sum_{j=0}^n \exp(A_\phi(s_j, a_j)/\beta)},$$

$$\text{and } A_\phi(s, a) = Q_\phi(s, a) - \mathbb{E}_{a' \sim \pi} [Q_\phi(s, a')]$$

is the advantage function as estimated by the critic model Q_ϕ . The weights $w(s, a)$ are a weighted importance sampling (WIS) approximation of $\exp(A_\phi(s, a)/\beta)/Z$, where $Z = \mathbb{E}_{s, a \sim \mathcal{B}} [\exp(A_\phi(s, a)/\beta)]$.

In the FB framework, the policies are conditioned by the latent variable z , and the Q -function estimate is $Q_\phi(s, a, z) = F_\phi(s, a, z)^T z$. Therefore, a direct transposition of (7) to the FB framework leads to the following objective for training the policy:

$$\arg \max_{\theta} \mathbb{E}_{(a_{1:n}, s_{1:n}) \sim \mathcal{B}, z_{1:n} \sim Z} \left[\sum_{i=1}^n w(s_i, a_i, z_i) \log \pi_\theta(a_i | s_i, z_i) \right], \quad (8)$$

$$\text{where } w(s_i, a_i, z_i) = \frac{\exp(A_\phi(s_i, a_i, z_i)/\beta)}{\sum_{j=0}^n \exp(A_\phi(s_j, a_j, z_j)/\beta)},$$

$$A_\phi(s, a, z) = F_\phi(s, a, z)^T z - \mathbb{E}_{a' \sim \pi_z} [F_\phi(s, a', z)^T z],$$

The choice of sampling an independent z_i for each sample (s_i, a_i) , as opposed to using a single z across a minibatch, was made based on empirically faster learning.

However, the variance and bias of this weighted importance sampling approach have a linear inverse relationship with the batch size n . Instead, we propose to use modified weights $w'(s, a, z)$ that implement *improved weighted importance sampling* (IWIS), a simple change to WIS proposed by Skare et al. (2003), shown to reduce the bias of WIS from $O(n^{-1})$ to $O(n^{-2})$. Integrating IWIS yields our final *policy improvement objective*:

$$\arg \max_{\theta} \mathbb{E}_{(a_{1:n}, s_{1:n}) \sim \mathcal{B}, z_{1:n} \sim Z} \left[\sum_{i=1}^n w'(s_i, a_i, z_i) \log \pi_{\theta}(a_i | s_i, z_i) \right], \quad (9)$$

where $w'(s_i, a_i, z_i) \propto \frac{w(s_i, a_i, z_i)}{\sum_{j \neq i} w(s_j, a_j, z_j)}$ and $\sum_{j=0}^n w'(s_j, a_j, z_j) = 1$,

and $w(s_i, a_i, z_i)$ is as in (8).

We validate the effect of IWIS in Table 7 (Appendix D): FB-AW, which uses the IWIS weights (9), performs slightly but consistently better than with the WIS weights (8). Figure 8 (Appendix D) illustrates how AW negates the overoptimism of vanilla FB for predicting future rewards.

Evaluation-based sampling. Furthermore, following Cetin et al. (2024), we integrate *evaluation-based sampling* (ES), an additional component to mitigate the undesirable consequences of learning a Gaussian policy, which is generally insufficient to capture the distribution from the exponentiated advantages. Namely, when deploying π_{θ} at test time, we approximate the argmax in (3) by sampling M actions a_1, \dots, a_M from the trained policy $\pi(s)$, and perform the one with the largest Q -value as predicted by $F_{\phi}(s, a_i, z)^T z$. The specific impact of this change is illustrated in Fig. 7 (Appendix D).

Uncertainty representation. To represent uncertainty in the model, we train two different networks F_1 and F_2 for the forward embedding, inspired by (Fujimoto et al., 2018; Touati et al., 2023). However, we introduce two changes.

In the Bellman equation, we use the *average* of the resulting two estimates of the target successor measures, namely, $\frac{1}{2}(F_1(s_{t+1}, a_{t+1}, z)^{\top} B(s') + F_2(s_{t+1}, a_{t+1}, z)^{\top} B(s'))$. This departs from Touati et al. (2023), which used the min between the target successor measures, namely, $\min\{F_1(s_{t+1}, a_{t+1}, z)^{\top} B(s'), F_2(s_{t+1}, a_{t+1}, z)^{\top} B(s')\}$, in line with (Fujimoto et al., 2018). Indeed, for Q -function estimates, a min might encode some form of conservatism, but for successor measures, the interpretation of a min is less direct.¹

Finally, we use two fully parallel networks for F_1 and F_2 , while Touati et al. (2023) opted for a shared processing network with two separate shallow heads for F_1 and F_2 .

¹For instance, since the Q -function for reward r is $Q = M.r$, taking the min of M might encode a min for a reward r but a max for the reward $-r$.

The specific impact of these changes is reported in Table 7 (Appendix D).

4 Evaluation

4.1 Algorithms and Baselines

We mainly compare the following algorithms:

- “Vanilla” FB: the classical implementation of FB from Touati et al. (2023) that employs TD3 policy improvement loss.
- FB-AW (FB with advantage weighting): The FB method using the advantage weighting components described in Section 3.2.
- FB-AWARE (FB with advantage weighting and auto-regressive encoding): The FB method using both AW and the auto-regressive component from Section 3.1. For the auto-regressive part, we test either 4 or 8 consecutive auto-regressive blocks for B .
- On some environments (those where AW is not necessary to reach good performance) we also include FB-ARE without the AW component.

Very few baselines provide true zero-shot behavioral foundation models that can tackle any new task at test time with no fine-tuning. (More baselines exist for the restricted case of goal-reaching tasks, namely, reaching a given target state, see Section 5.) In addition to the vanilla FB baseline, we include the following non-FB baseline:

- Universal successor features (Borsa et al., 2018) based on *Laplacian eigenfunctions* as the base features (Touati et al., 2023). This version of successor features was found to perform best in Touati et al. (2023). We denote it by LAP-AW, since we use the advantage weighting as in FB-AW.

All the aforementioned variants of FB use the same architecture and consistent hyperparameters.

4.2 Datasets and Benchmarks

We consider a series of environments, tasks, and datasets for these environments, as follows.

- We start with the *Jaco arm* domain (Laskin et al., 2021), a simple robotic arm model. The tasks consist in reaching various target positions (Section 4.3.1). This provides a test of spatial precision.

For this domain, we build a training dataset via the *RND* unsupervised exploration method from Yarats et al. (2022), which provides good data diversity if exploration is not too difficult in an environment.

- Next, we consider four standard domains from the *DeepMind Control (DMC) Suite* (Tassa et al., 2018): Walker, Cheetah, Quadruped, and Humanoid. Since we want to build behavior foundation models and not task-specific agents, on top of the classical tasks for these environments (walk, run...) we introduce a number of additional tasks such as bounce, flip, pullup..., described in Appendix B.2.

For these domains, we consider two training datasets:

- We build a first dataset using RND, as for Jaco. However, RND appears to provide insufficient exploration (particularly on Humanoid); moreover, the RND trainset does not contain any purposeful trajectories.
- Therefore, we also train on the *MOOD* datasets from Cetin et al. (2024). MOOD contains a mixture of behaviors, obtained as follows. For each environment, a small number of “classical” tasks are selected. Then an online TD3 algorithm is used to train a classical agent for each of these tasks. The set of trajectories produced by these agents during training are then pooled and merged into a single dataset for the environment. Thus, the mixed-objective MOOD datasets include high-quality examples for a few tasks in each environment.

Evaluation on the MOOD dataset must distinguish between tasks that contributed to the dataset (*in-dataset* tasks), and tasks that did not (*out-of-dataset* tasks).² A priori, one would expect the former to be easier, as information from the original single-task agents is present in the data. To evaluate the ability of the FB models to generalize beyond in-dataset tasks, we used the new tasks from Appendix B.2 as out-of-dataset tasks.

- Finally, to test the generality of the approach, we also train FB, FB-AW and FB-AWARE agents on the locomotion tasks of the *D4RL benchmark* (Fu et al., 2020). Here we stick to the original tasks in the benchmark, and compare FB performance to the best task-specific offline RL agents in the literature. Since those are single-task while FB is a generalist agent, this is a natural *topline* for FB, so we expect FB to reach a good fraction of the performance of the best task-specific agents, in line with the methodology of Touati et al. (2023).

4.3 Empirical Evaluation

We train FB, FB-AW, FB-AWARE (4 and 8 blocks) and LAP-AW on the four DMC locomotion environments (Walker, Cheetah, Quadruped and Humanoid), as well as the Jaco arm domain. We pretrain each model on both offline datasets (MOOD and RND), and repeat each training 5 times (with different random seeds).

²We avoid “in-distribution” and “out-of-distribution”, since FB is not trained on a distribution of tasks but in an unsupervised way given the data.

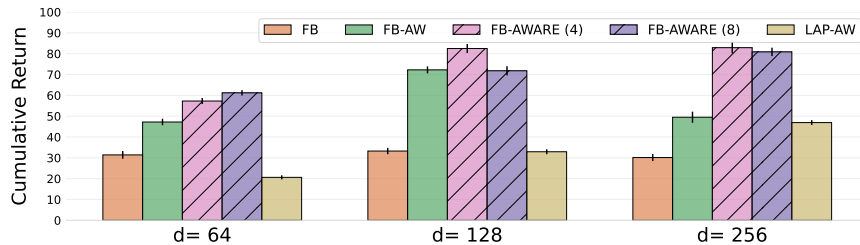


Figure 2: Average cumulative reward achieved by the algorithms, trained on RND dataset for different representation dimensions when aiming to reach goals (four randomly selected goals and four corner goals), in the Jaco arm environment.

We evaluate each model on several downstream tasks per environment. For each model and task, we sample 100,000 states $\{s\}$ from the offline dataset and compute their corresponding task reward $\{r(s)\}$ in order to infer the task encoding vector z_r ((4) or (6)). Then we compute the cumulated reward achieved by the policy π_{z_r} , computed using the task-specific reward, and averaged over 100 episodes. Finally, we report the average and variance of the cumulative reward over the 5 pre-trained models (with different seeds).

4.3.1 Jaco Arm Results

We train each algorithm in the Jaco environment using the RND dataset, for three choices of representation dimension: $d = 64, 128$, and 256 . The tasks involve reaching a goal within an episode length of 250 time steps, with the agent receiving a reward of approximately 1 when the arm’s gripper is close to the target goal specified by its (x, y, z) coordinates. For the goals, we included the four corners of the environment, plus four goals selected at random (once and for all, common to all the algorithms tested).

In Figure 2, we depict the average rewards attained by each algorithm for reaching this mixture of goals. The resulting goal-reaching rewards for dimensions 64, 128, and 256 are presented in Tables 2, 3, and 4, respectively.

FB-AW significantly outperforms FB, more than doubling the score for the best dimension $d = 128$. FB-AWARE with 4 autoregressive blocks further enhances performance by a relative margin of about 15%.

4.3.2 DMC Locomotion Results

For Cheetah, Quadruped, Walker and Humanoid, the MOOD dataset results in substantially better models than the RND dataset, whatever the algorithm (Appendix C, Table 6 vs Table 5). This is especially striking for Humanoid, where RND does not explore enough and even a classical single-task TD3 agent is hard to train. Therefore, we focus the discussion on MOOD, with full RND results in Appendix C.

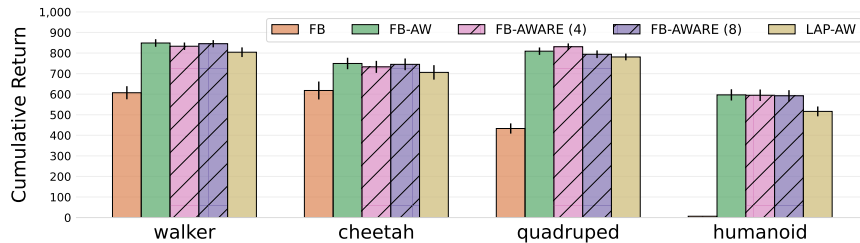


Figure 3: Averaged cumulative reward achieved by the algorithms on *in-dataset tasks*, trained on MOOD dataset for DMC Locomotion.

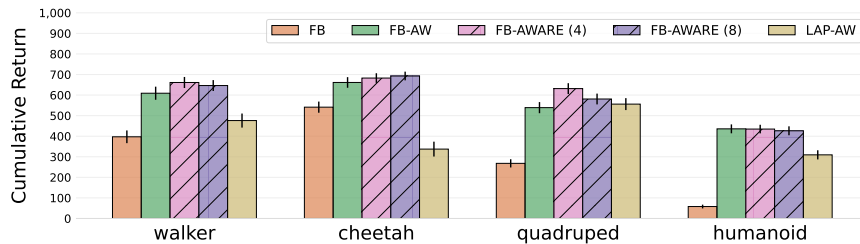


Figure 4: Average cumulative reward achieved by the algorithms on *out-of-dataset tasks*, trained on MOOD dataset for DMC Locomotion.

Figures 3 and 4 report the results for the Cheetah, Quadruped, Walker and Humanoid environments, using MOOD data for training, evaluated on in-dataset tasks and out-of-dataset tasks respectively. We used a fixed representation dimension for all algorithms ($d = 64$ for Walker, Cheetah and Quadruped, $d = 128$ for Humanoid).

In this setup, the advantage weighting component proves crucial for achieving satisfactory performance, particularly in the Humanoid environment, where vanilla FB performance is near-zero. However, on the lower-quality RND dataset, the advantage weighting component appears to hurt performance (Appendix C, Table 5). This is consistent with observations in Yarats et al. (2022) for the single-task setting, where conservative offline-RL methods hurt performance on RND data.

LAP-AW does well at in-dataset tasks, but lags behind FB for out-of-dataset tasks. This may be because LAP-AW’s learned features are closely tied to the in-dataset tasks, derived from the eigenfunctions of the Laplacian of the behavior policy present in the dataset.

FB-AW and FB-AWARE exhibit the most favorable and consistent overall performance. The autoregressive component provides a slight enhancement in out-of-dataset tasks across all environments except from Humanoid. This slight difference is not observed on the RND dataset (Table 5). The best-performing model overall is obtained with the MOOD dataset and FB-AWARE algorithm.

Table 1: Performance on the popular locomotion-v2 and FrankaKitchen datasets from the D4RL benchmark, comparing with recent offline RL algorithms (with performance as reported in the literature). FB-AWARE uses 8 AR groups. Other hyper-parameters are tuned per-environment, consistently with the offline baselines.

Dataset/Algorithm	BC	10%BC	DT	1StepRL	AWAC	TD3+BC	CQL	IQL	XQL	Reward free learning		
										FB	FB-AW	FB-AWARE
halfcheetah-medium-v2	42.6	42.5	42.6	48.4	43.5	48.3	44.0	47.4	48.3	49.0±1.93	60.0±0.9	62.7±0.9
hopper-medium-v2	52.9	56.9	67.6	59.6	57.0	59.3	58.5	66.3	74.2	0.9±0.69	59.1±5.0	59.9±21.4
walker2d-medium-v2	75.3	75.0	74.0	81.8	72.4	83.7	72.5	78.3	84.2	0.5±0.9	80.5±11.7	89.6±0.8
halfcheetah-medium-replay-v2	36.6	40.6	36.6	38.1	40.5	44.6	45.5	44.2	45.2	30.8±23.4	52.7±1.2	50.8±1.2
hopper-medium-replay-v2	18.1	75.9	82.7	97.5	37.2	60.9	95.0	94.7	100.7	16.4±2.9	87.1±3.7	89.6±4.0
walker2d-medium-replay-v2	26.0	62.5	66.6	49.5	27.0	81.8	77.2	73.9	82.2	9.9±4.9	91.7±6.3	98.8±0.5
halfcheetah-medium-expert-v2	55.2	92.9	86.8	93.4	42.8	90.7	91.6	86.7	94.2	91.7±6.7	99.6±0.8	100.1±0.7
hopper-medium-expert-v2	52.5	110.9	107.6	103.3	55.8	98.0	105.4	91.5	111.2	1.8±1.2	55.9±11.6	62.2±9.2
walker2d-medium-expert-v2	107.5	109.0	108.1	113.0	74.5	110.1	108.8	109.6	112.7	0.3±0.8	109.6±1.3	105.8±0.6
Locomotion-v2 total	466.7	666.2	672.6	684.6	450.7	677.4	698.5	692.6	752.9	22.3	696.2	719.5
kitchen-partial-v0	38.0	-	-	-	-	-	49.8	46.3	73.7	4±4	47.0±4.5	52.5±9.4
kitchen-mixed-v0	51.5	-	-	-	-	-	51.0	51.0	62.5	5±5	48.5±7.2	53.5±3.8

4.3.3 Performance of FB-AW and FB-AWARE on D4RL

Finally, to test the generality and robustness of these methods, we test performance on the D4RL benchmark after reward-free training. D4RL is a ubiquitous benchmark, used by many recent offline RL research for evaluation and comparison.

Here, we are comparing the multitask, unsupervised FB-AWARE agent to task-specific agents, so the performance of the latter are a natural *topline* for FB-AW and FB-AWARE. In line with [Touati et al. \(2023\)](#), we expect FB-AWARE to reach a good fraction of the performance of the agents trained specifically on each task.

So we compare FB-AWARE’s performance with the results available for a large pool of offline RL algorithms optimizing for an individual objective with full access to rewards. This setting is quite different from Section 4.3.2 and especially the MOOD datasets, since D4RL datasets mostly comprise trajectories from agents trying to accomplish a single task.

FB-AW and FB-AWARE’s overall performance matches the task-specific recent state-of-the-art from XQL and IQL. There is a slight advantage to FB-AWARE over FB-AW, although this falls within the overall margin of error.

This establishes that zero-shot, task-agnostic behavior foundation models trained via FB can compete with top task-specific agents for offline RL on standard benchmarks.

4.4 Ablations and Discussion

Appendix D contains additional tests and ablations concerning each of the components introduced, such as the impact of dimension d , the number of autoregressive blocks for FB-AWARE, specific design choices for FB training (B normalization, z sampling), and the offline RL methods introduced in Section 3.2 (advantage weighting, improved weighted importance sampling, evaluation-based

sampling, and uncertainty representation).

The impact of training data for learning behavior foundation models.

Perhaps unsurprisingly, the dataset has a large impact on the performance of behavior foundation models, as exemplified by the higher scores of FB methods trained on MOOD vs RND for the Locomotion tasks. On Humanoid, the combination of both MOOD and advantage weighting appears necessary to reach any reasonable performance at all.

On the other hand, the best algorithm to train a behavior foundation model also depends on the data available: with only RND data, advantage weighting actually hurts performance (Table 5), although with these data, performance is relatively poor anyway. This is consistent with existing observations for classical single-task agents: on RND data, TD3 is better than conservative offline-RL approaches (Yarats et al., 2022).

Therefore, to train good behavior foundation models, it may be necessary to train variants via several algorithms, and check the model’s performance on a few validation tasks.

Autoregressive FB and fine-grained behaviors. The advantage provided by autoregressive FB is clearer on the Jaco arm, and on out-of-dataset tasks for DMC Locomotion. We can offer some speculative explanations.

Part of the original motivation behind autoregressive FB was to provide more precise task encoding. Indeed, vanilla FB projects tasks linearly onto a subset of tasks, and this results in *reward blurring* (loss of detail in the actual task being optimized) (Touati & Ollivier, 2021). This may correspond to selecting the largest eigenvalues of the successor measure, since FB fixed points correspond to its eigenvalues (Blier et al., 2021). Intuitively, autoregressive features should provide more detailed features based on a first set of coarse-grained features: for instance, if a target is roughly located in a region via the top-level features B_1 , then B_2 should learn more specific features for the region given by B_1 (as encoded by z_1).

The Jaco arm domain aims at testing this intuition, by providing a task in which a target must be precisely reached. The observed advantage of autoregressive FB for this task is consistent with this intuition, somewhat suggesting autoregressive features are effective at providing finer spatial precision.

One possible, though speculative, explanation for the effect on out-of-dataset tasks is the following. In vanilla FB, the set of perfectly optimized tasks is linearly spanned by B : thus, it is a linear space of dimension d among all possible tasks. In autoregressive FB, the subset of rewards optimized by the features is nonlinear among all possible reward functions, because the features B_2 depend on the continuous variable z_1 ; for instance, with two auto-regressive blocks, one can represent all rewards $r(s) = z_1^T B_1(s) + z_2^T B_2(s, z_1)$. This results in a curved manifold in the set of reward functions when varying z_1 and z_2 . When facing a new arbitrary reward function, it may be easier to find a close match on this manifold than on a linear subspace.

Limitations. The environments considered here are all noise-free, Markovian (history-free) continuous control environments.

The effect of autoregressive FB is relatively modest in these experiments. This is surprising given the huge theoretical change in expressivity compared to vanilla FB. This suggests that the main limiting factor in our suite of experiments may not be the expressivity of the behavior foundation model, perhaps due to limited exploration in the training datasets, or from the relative simplicity of the environments tested.

5 Related Work

The forward-backward framework (Touati & Ollivier, 2021) builds upon the principles of successor features (Barreto et al., 2017; Zhang et al., 2017; Grimm et al., 2019), the continuous extension to the canonical successor representation Dayan (1993) and its continuous state-space extension. However, in contrast to FB, this line of work has mostly focused on constructing a set of features, linear w.r.t. downstream rewards, using apriori knowledge and heuristic measures such as Laplacian eigenfunctions (Borsa et al., 2018). To this end Touati et al. (2023) showed, empirically, the superiority of end-to-end learning with FB as compared to many such heuristics, in line with results in this same work. The proposed autoregressive extension to FB could be also applied to this broader class of methods, a noteworthy direction for future work.

Goal-conditioned RL (GCRL) (Liu et al., 2022) is another area of research closely related to FB, which has seen notable successful applications in real-world robotics (Shah et al., 2021; Ma et al., 2022; Zheng et al., 2023a). As with successor features, GCRL has traditionally relied on a priori knowledge taking the form of explicit demonstrations (Ding et al., 2019; Gupta et al., 2019), handcrafted subgoals (Andrychowicz et al., 2017; Nachum et al., 2018; Chane-Sane et al., 2021), together with other coverage heuristics (Ghosh et al., 2019; Hansen-Estruch et al., 2022) – which have been used both to construct the space of goals and learn its relative multi-task policy. Moreover, in a similar fashion to FB, recent work also strived to model and tackle the GCRL problem with a principled contrastive-like end-to-end objective (Eysenbach et al., 2021, 2022; Zheng et al., 2023b). However, GCRL is, by design, more restrictive than successor features and FB as it cannot capture tasks that go beyond reaching individual points in the space of goals.

Other works avoid the linearity constraint of successor features and vanilla FB by explicitly relying on a prior over tasks. For instance, in an approach akin to meta-RL but without hand-crafted tasks, Frans et al. (2024) use a mixture of random MLPs, random linear functions, and random goal-reaching to pre-train a set of policies together with an encoder that quickly identifies a reward function from a few reward samples. Contrary to FB, the dynamics of the environment plays no role in building the set of features.

6 Conclusions

Specific offline RL training techniques are necessary to build efficient FB behavior foundation models in environments such as DMC Humanoid, and can make the difference between near-zero and good performance.

Employing auto-regressive features greatly enhances the theoretical expressivity of these foundation models, and improves spatial precision and task generalization. The improvement is moderate in our setup, perhaps indicating that BFM expressivity is not a key limiting factor for these tasks.

These improvements bring zero-shot, reward-free FB BFMs on par with single-task, reward-trained offline agents on a number of locomotion environments.

References

- Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Abbeel, O. P., and Zaremba, W. Hindsight experience replay. In *Advances in neural information processing systems*, pp. 5048–5058, 2017.
- Barreto, A., Dabney, W., Munos, R., Hunt, J. J., Schaul, T., van Hasselt, H. P., and Silver, D. Successor features for transfer in reinforcement learning. *Advances in neural information processing systems*, 30, 2017.
- Blier, L., Tallec, C., and Ollivier, Y. Learning successor states and goal-dependent values: A mathematical viewpoint. *arXiv preprint arXiv:2101.07123*, 2021.
- Borsa, D., Barreto, A., Quan, J., Mankowitz, D. J., van Hasselt, H., Munos, R., Silver, D., and Schaul, T. Universal successor features approximators. In *International Conference on Learning Representations*, 2018.
- Burda, Y., Edwards, H., Storkey, A., and Klimov, O. Exploration by random network distillation. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=H1lJJnR5Ym>.
- Cetin, E., Tirinzoni, A., Pirotta, M., Lazaric, A., Ollivier, Y., and Touati, A. Asac: Simple ingredients for offline reinforcement learning with diverse data. *ICML*, 2024.
- Chane-Sane, E., Schmid, C., and Laptev, I. Goal-conditioned reinforcement learning with imagined subgoals. In *International Conference on Machine Learning*, pp. 1430–1440. PMLR, 2021.
- Dayan, P. Improving generalization for temporal difference learning: The successor representation. *Neural computation*, 5(4):613–624, 1993.
- Ding, Y., Florensa, C., Abbeel, P., and Phielipp, M. Goal-conditioned imitation learning. *Advances in neural information processing systems*, 32, 2019.

- Eysenbach, B., Salakhutdinov, R., and Levine, S. C-learning: Learning to achieve goals via recursive classification. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=tc5qisoB-C>.
- Eysenbach, B., Zhang, T., Levine, S., and Salakhutdinov, R. R. Contrastive learning as goal-conditioned reinforcement learning. *Advances in Neural Information Processing Systems*, 35:35603–35620, 2022.
- Frans, K., Park, S., Abbeel, P., and Levine, S. Unsupervised zero-shot reinforcement learning via functional reward encodings. *arXiv preprint arXiv:2402.17135*, 2024.
- Fu, J., Kumar, A., Nachum, O., Tucker, G., and Levine, S. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- Fujimoto, S., van Hoof, H., and Meger, D. Addressing function approximation error in actor-critic methods. In *ICML*, pp. 1582–1591, 2018. URL <http://proceedings.mlr.press/v80/fujimoto18a.html>.
- Fujimoto, S., Meger, D., and Precup, D. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*, pp. 2052–2062. PMLR, 2019.
- Garg, D., Hejna, J., Geist, M., and Ermon, S. Extreme q-learning: Maxent RL without entropy. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=SJOLde3tRL>.
- Ghosh, D., Gupta, A., Reddy, A., Fu, J., Devin, C., Eysenbach, B., and Levine, S. Learning to reach goals via iterated supervised learning. *arXiv preprint arXiv:1912.06088*, 2019.
- Grimm, C., Higgins, I., Barreto, A., Teplyashin, D., Wulfmeier, M., Hertweck, T., Hadsell, R., and Singh, S. Disentangled cumulants help successor representations transfer to new tasks. *arXiv preprint arXiv:1911.10866*, 2019.
- Gupta, A., Kumar, V., Lynch, C., Levine, S., and Hausman, K. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. *arXiv preprint arXiv:1910.11956*, 2019.
- Hansen-Estruch, P., Zhang, A., Nair, A., Yin, P., and Levine, S. Bisimulation makes analogies in goal-conditioned reinforcement learning. In *International Conference on Machine Learning*, pp. 8407–8426. PMLR, 2022.
- Kostrikov, I., Nair, A., and Levine, S. Offline reinforcement learning with implicit q-learning. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=68n2s9ZJWF8>.

- Laskin, M., Yarats, D., Liu, H., Lee, K., Zhan, A., Lu, K., Cang, C., Pinto, L., and Abbeel, P. Urlb: Unsupervised reinforcement learning benchmark. *arXiv preprint arXiv:2110.15191*, 2021.
- Levine, S., Kumar, A., Tucker, G., and Fu, J. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- Liu, M., Zhu, M., and Zhang, W. Goal-conditioned reinforcement learning: Problems and solutions. In Raedt, L. D. (ed.), *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pp. 5502–5511. International Joint Conferences on Artificial Intelligence Organization, 7 2022. doi: 10.24963/ijcai.2022/770. URL <https://doi.org/10.24963/ijcai.2022/770>. Survey Track.
- Ma, Y. J., Sodhani, S., Jayaraman, D., Bastani, O., Kumar, V., and Zhang, A. Vip: Towards universal visual reward and representation via value-implicit pre-training. *arXiv preprint arXiv:2210.00030*, 2022.
- Nachum, O., Gu, S. S., Lee, H., and Levine, S. Data-efficient hierarchical reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- Nair, A., Gupta, A., Dalal, M., and Levine, S. Awac: Accelerating online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020a.
- Nair, A., Gupta, A., Dalal, M., and Levine, S. Awac: Accelerating online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020b.
- Park, S., Kreiman, T., and Levine, S. Foundation policies with Hilbert representations. In *Forty-first International Conference on Machine Learning*, 2024.
- Peng, X. B., Kumar, A., Zhang, G., and Levine, S. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- Pirotta, M., Tirinzoni, A., Touati, A., Lazaric, A., and Ollivier, Y. Fast imitation via behavior foundation models. In *The Twelfth International Conference on Learning Representations*, 2023.
- Shah, D., Eysenbach, B., Kahn, G., Rhinehart, N., and Levine, S. Rapid exploration for open-world navigation with latent goal models. *arXiv preprint arXiv:2104.05859*, 2021.
- Skare, Ø., Bølviken, E., and Holden, L. Improved sampling-importance resampling and reduced bias importance sampling. *Scandinavian Journal of Statistics*, 30(4):719–737, 2003.

- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018. 2nd edition.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D. d. L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- Touati, A. and Ollivier, Y. Learning one representation to optimize all rewards. *Advances in Neural Information Processing Systems*, 34:13–23, 2021.
- Touati, A., Rapin, J., and Ollivier, Y. Does zero-shot reinforcement learning exist? In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=MYEap_0cQI.
- Wang, Z., Novikov, A., Zolna, K., Merel, J. S., Springenberg, J. T., Reed, S. E., Shahriari, B., Siegel, N., Gulcehre, C., Heess, N., and de Freitas, N. Critic regularized regression. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 7768–7778. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/588cb956d6bbe67078f29f8de420a13d-Paper.pdf>.
- Wu, Y., Tucker, G., and Nachum, O. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- Yarats, D., Brandfonbrener, D., Liu, H., Laskin, M., Abbeel, P., Lazaric, A., and Pinto, L. Don’t change the algorithm, change the data: Exploratory data for offline reinforcement learning. *arXiv preprint arXiv:2201.13425*, 2022.
- Zhang, J., Springenberg, J. T., Boedecker, J., and Burgard, W. Deep reinforcement learning with successor features for navigation across similar environments. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2371–2378. IEEE, 2017.
- Zheng, C., Eysenbach, B., Walke, H., Yin, P., Fang, K., Salakhutdinov, R., and Levine, S. Stabilizing contrastive rl: Techniques for offline goal reaching. *arXiv preprint arXiv:2306.03346*, 2023a.
- Zheng, C., Salakhutdinov, R., and Eysenbach, B. Contrastive difference predictive coding. *arXiv preprint arXiv:2310.20141*, 2023b.

A Auto-Regressive FB: Extensions, Proofs, Algorithmic Considerations

A.1 Proof and Extension of Theorem 3.2

Here we prove Theorem 3.2 and extend it in two directions similar to Touati & Ollivier (2021).

The first extension concerns *goal spaces*, and is useful when we know in advance that the rewards functions of interest will not depend on the whole state. For instance, in a multi-agent setting, the reward of an agent may depend only on its own state, but it must still observe the whole state to make decisions. This is formalized by assuming that the reward function only depends on some variable $g = \phi(s)$ rather than the whole state s . Then we can learn with $B(g)$ instead of $B(s)$ (while F and policies still require the full state).

The second extension only uses the finite-rank $F^\top B$ model for the *advantage* functions, namely, the model is $M^{\pi_z}(s_0, a_0, ds) \approx F(s_0, a_0, z)^\top B(s', z) + \bar{m}(s_0, z, s)$ where \bar{m} is any function independent of the actions. This lifts part of the finite-rank restriction, since \bar{m} itself is unconstrained: the finite-rank FB model is only applied to the differential effect of actions on top of the baseline model \bar{m} .

Definition A.1 (Extended forward-backward representation of an MDP). Consider an MDP with state space S and action space A . Let $\phi: S \times A \rightarrow G$ be a function from state-actions to some goal space $G = \mathbb{R}^k$.

Let $Z = \mathbb{R}^d$ be some representation space. Let

$$F: S \times A \times Z \rightarrow Z, \quad B: G \times Z \rightarrow Z, \quad \bar{m}: S \times Z \times G \rightarrow \mathbb{R} \quad (10)$$

be three functions. For each $z \in Z$, define the policy

$$\pi_z(a|s) := \arg \max_a F(s, a, z)^\top z. \quad (11)$$

Let ρ be any measure over the goal space G .

We say that F , B , and \bar{m} are an *extended forward-backward representation* of the MDP with respect to ρ , if the following holds: for any $z \in Z$, any state-action (s, a) the successor measure M^{π_z} of policy π_z is given by

$$M^{\pi_z}(s, a, X) = \int_{g \in X} (F(s, a, z)^\top B(g, z) + \bar{m}(s, z, g)) \rho(dg) \quad (12)$$

for any goal subset $X \subset G$.

Theorem A.2 (Forward-backward representation of an MDP, with features as goals). Consider an MDP with state space S and action space A . Let $\phi: S \times A \rightarrow G$ be a function from state-actions to some goal space $G = \mathbb{R}^k$.

Let F , B , and \bar{m} be an extended forward-backward representation of the MDP with respect to some measure ρ over G .

Then the following holds. Let $R: S \times A \rightarrow \mathbb{R}$ be any bounded reward function, and assume that this reward function depends only on $g = \phi(s, a)$, namely, that there exists a function $r: G \rightarrow \mathbb{R}$ such that $R(s, a) = r(\phi(s, a))$.

Assume that there exists $z_R \in \mathbb{R}^d$ satisfying

$$z_R = \int_{g \in G} r(g)B(g, z_R) \rho(dg). \quad (13)$$

This is always the case if B is auto-regressive.

Then:

1. π_{z_R} is an optimal policy for reward R in the MDP.
2. The optimal Q -function Q_R^* for reward R is

$$Q_R^*(s, a) = F(s, a, z_R)^\top z_R + \int_{g \in G} \bar{m}(s, z_R, g) r(g) \rho(dg). \quad (14)$$

The last term does not depend on the action a , so computing the \bar{m} term is not necessary to obtain the advantages $Q^*(s, a) - Q^*(s, a')$ or the optimal policies.

Theorem A.2 implies Theorem 3.2, by taking $\phi = \text{Id}$ and $\bar{m} = 0$.

There is an important difference between this result and the corresponding statement in non-auto-regressive FB (Theorem 4 in Touati & Ollivier (2021)). For non-autoregressive FB, the FB model provides the Q -functions of all policies π_z for all rewards R , even if $z \neq z_R$. Namely, $Q_R^{\pi_z}(s, a) = F(s, a, z)^\top z_R$ for all pairs (z, R) (in the case $\bar{m} = 0$). Here, this only holds when $z = z_R$. Classical, non-auto-regressive FB provides more Q -functions than strictly needed to obtain the policies: it models the Q -functions of π_z even for rewards unrelated to z . With auto-regressive FB, the additional expressivity of the model comes at the price of getting less information about the other Q -functions.

Proof of Theorem A.2. Let M^π be the successor measure of policy π . Let m^π be the density of M^π with respect to ρ . Let $R(s, a) = r(\phi(s, a))$ be a reward function as in the statement of the theorem.

By Proposition 16 in Touati & Ollivier (2021), The Q -function of π for the reward R is

$$Q_R^\pi(s, a) = \int_g r(g)M^\pi(s, a, dg) \quad (15)$$

Let z_R satisfy the fixed point property (13), and let us take $\pi = \pi_{z_R}$. By definition of an extended FB representation, we have

$$Q_R^{\pi_{z_R}}(s, a) = \int_g r(g)M^{\pi_{z_R}}(s, a, dg) \quad (16)$$

$$= \int_g r(g)(F(s, a, z_R)^\top B(g, z_R) + \bar{m}(s, z_R, g))\rho(dg). \quad (17)$$

$$= F(s, a, z_R)^\top \int_g r(g)B(g, z_R)\rho(dg) + \int_g r(g)\bar{m}(s, z_R, g)\rho(dg). \quad (18)$$

But thanks to the fixed point property (13), we have $\int_g r(g)B(g, z_R)\rho(dg) = z_R$. Therefore, the Q -function of π_{z_R} for reward R is

$$Q_R^{\pi_{z_R}}(s, a) = F(s, a, z_R)^\top z_R + \int_g r(g)\bar{m}(s, z_R, g)\rho(dg). \quad (19)$$

We have to prove that this is the optimal Q -function for R . A pair of a Q -function and policy π are optimal for R if and only if simultaneously $\pi(a|s) = \arg \max_a Q(s, a)$ and $Q = Q_R^\pi$. Here, by definition of the policies π_z , we have

$$\pi_{z_R}(a|s) = \arg \max_a F(s, a, z_R)^\top z_R \quad (20)$$

$$= \arg \max_a Q_R^{\pi_{z_R}}(s, a) \quad (21)$$

since the additional term $\int_g r(g)\bar{m}(s, z_R, g)\rho(dg)$ in $Q_R^{\pi_{z_R}}$ does not depend on a .

Therefore, $Q_R^{\pi_{z_R}}$ and π_{z_R} are optimal for reward R , which ends the proof. \square

Theorem A.3 (Auto-regressive features with two levels are a universal approximator for task encoding). *Assume the state space is finite, so that a reward function is an element of $\mathbb{R}^{\#S}$. Then, for any continuous task encoding function $\zeta: \mathbb{R}^{\#S} \rightarrow \mathbb{R}^d$ mapping rewards r to task encodings $z = \zeta(r)$, such that $z_r = 0$ for $r = 0$, there exist neural networks $B_1(s)$ and $B_2(s, z)$ approximating ζ , namely, for any r ,*

$$\zeta(r) \approx \mathbb{E}_{s \sim \rho}[r(s)B_2(s, z_1)], \quad z_1 = \mathbb{E}_{s \sim \rho}[r(s)B_1(s)] \quad (22)$$

up to an arbitrary precision.

Lemma A.4. *Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^d$ be a C^3 function such that $f(0) = 0$. Then there exists a continuous matrix-valued function $g: \mathbb{R}^n \rightarrow \mathbb{R}^{n \times d}$ such that for any $x \in \mathbb{R}^n$,*

$$f(x) = g(x) \cdot x. \quad (23)$$

Proof of the lemma. By working on each output component of f separately, we can assume that $d = 1$. Thus, we have to prove that for any C^3 function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ with $f(0) = 0$, there exists a vector-valued function $g: \mathbb{R}^n \rightarrow \mathbb{R}$ such that $f(x) = g(x)^\top x$.

Since f is C^3 with $f(0) = 0$, we can write its Taylor expansion

$$f(x) = D^\top x + \frac{1}{2}x^\top Hx + R(x) \quad (24)$$

where $D = \partial_x f(0)$ is the gradient of f at $x = 0$, $H = \partial_x^2 f(0)$ its Hessian, and where the remainder R is $O(\|x\|^3)$.

A priori, the function $R(x)/\|x\|^2$ is defined everywhere except $x = 0$. But since R is $O(\|x\|^3)$, $R(x)/\|x\|^2$ tends to 0 for $x = 0$, so it is a well-defined continuous function on the whole domain.

Thus, let us set

$$g(x) := D + \frac{1}{2}Hx + \frac{R(x)}{\|x\|^2}x. \quad (25)$$

Then, by construction, $g(x)^\top x = f(x)$ as needed. \square

Proof. Proof of Theorem A.3 By the lemma, there exists a matrix-valued function g such that $\zeta(r) = g(r) \cdot r$ for $r \in \mathbb{R}^{\#S}$.

Define $B_1(s)$ to be the one-hot encoding $B_1(s) = \mathbb{1}_s / \rho(s) \in \mathbb{R}^{\#S}$, where $\mathbb{1}_i$ denotes the vector with all zeroes except a 1 at position i . By universal approximation theorems for neural networks, this choice of B_1 can be realized by a neural network with arbitrary good approximation (actually a one-layer neural network with identity weights and no activation function).

Then

$$z_1 = \mathbb{E}_{s \sim \rho}[r(s)B_1(s)] = \sum_s \rho(s)r(s)\mathbb{1}_s / \rho(s) = r. \quad (26)$$

Then in turn, for any B_2 ,

$$z_2 = \mathbb{E}_{s \sim \rho}[r(s)B_2(s, z_1)] = \sum_s \rho(s)r(s)B_2(s, r). \quad (27)$$

For each s and r , $B_2(s, r)$ is an element of \mathbb{R}^d . For each component $1 \leq i \leq d$, define

$$B_2(s, r)_i := g_{is}(r) / \rho(s) \quad (28)$$

where $g(r)$ is the matrix defined above. Then we have, by construction

$$\sum_s \rho(s)r(s)B_2(s, r)_i = \sum_s r(s)g_{is}(r) \quad (29)$$

namely

$$B_2(s, r) = g(r) \cdot r = \zeta(r) \quad (30)$$

as needed.

By universal approximation theorems for neural networks, it is possible to realize this choice of B_2 by a neural network, with arbitrarily good approximation.

Note: the principle of this proof extends to continuous state spaces, by taking a partition of unity for B_1 instead of a one-hot encoding, though this results in further approximation errors. \square

A.2 Training Loss and Algorithmic Considerations for Auto-Regressive FB

Loss for FB-AR; sampling. While plain FB represents the successor measures as $M^{\pi_z}(s, a, ds') \approx F(s, a, z)^\top B(s')\rho(ds')$, auto-regressive FB uses $M^{\pi_z}(s, a, ds') \approx F(s, a, z)^\top B(s', z)\rho(ds')$. The training principle remains the same: learn F and B to minimize the error in this approximation. However, this leads to several changes in practice.

The error of the FB model can be measured as in plain FB, based on the Bellman equation satisfied by M^{π_z} (see Appendix B in [Touati et al. \(2023\)](#)). For each value of z , the Bellman loss on $M^{\pi_z}(s, a, ds') - F(s, a, z)^\top B(s')\rho(ds')$ is

$$\begin{aligned} \mathcal{L}(F, B) = & \mathbb{E}_{\substack{(s_t, a_t, s_{t+1}) \sim \rho \\ s' \sim \rho}} \left[\left(F(s_t, a_t, z)^\top B(s', z) - \gamma \bar{F}(s_{t+1}, \pi_z(s_{t+1}), z)^\top \bar{B}(s', z) \right)^2 \right] \\ & - 2 \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \rho} [F(s_t, a_t, z)^\top B(s_{t+1}, z)] + \text{Const} \end{aligned} \quad (31)$$

where **Const** is a constant term that we can discard since it does not depend on F and B . The only difference with plain FB is that now B depends on z .

For training, the variable z is sampled as in [Touati et al. \(2023\)](#): namely, z is sampled 50% of the time from a standard d -dimensional Gaussian (and later normalized) and 50% of the time from the B representation of a state randomly sampled from the training data (computed every step for normal FB and computed every 32 steps for AR FB for speed). Fig. 13 (Appendix D) tests the effect of only sampling z from a Gaussian.

For reasons discussed below, for FB we follow the original strategy of sampling a different z for each sampled transition, but for FB-AR we sample a unique z for the whole batch.

Algorithm 2 implements this loss, together with sampling of z .

Minibatch handling, and increased variance for auto-regressive FB.

Having B depend on z has practical consequences for variance and minibatch sampling. In vanilla FB, starting from (31), it is possible to sample a minibatch of N transitions (s_t, a_t, s_{t+1}) , a minibatch of N values of s' , choose a different value of z for each s_t in the minibatch, compute the N values of $F(s_t, a_t, z)$ and $B(s')$, and compute the N^2 dot products $F(s_t, a_t, z)^\top B(s')$ involved in the loss, at a cost of N forward passes through F and B (Appendix A in [Touati & Ollivier \(2021\)](#)).

In auto-regressive FB, the value of z must be the same for F and B in the loss (31). Therefore, contrary to plain FB, we only sample a *single* value of z for the minibatch ($I = 1$ in Algorithm 2). Then we can compute the N values $F(s_t, a_t, z)$ and $B(s', z)$, and the N^2 dot products $F(s_t, a_t, z)^\top B(s', z)$ for the loss (31). This appears in Algorithm 2.

Using a single value of z per minibatch results in increased variance of auto-regressive FB with respect to vanilla FB and longer training times, which we observe in practice.

Indeed, if we want to keep using different values of z for each $F(s_t, a_t, z)$, we must either compute many more values $B(s', z)$ (one for each pair s' and z), or use fewer dot products by computing fewer values of $B(s', z)$ and only using them with $F(s_t, a_t, z)$ with the same z .

More precisely, in general we can proceed as follows: Let k be a hyperparameter controlling the number of distinct values of z we will use in the minibatch. We sample N transitions (s_t, a_t, s_{t+1}) , N values of s' , and split these samples into k groups of size N/k . For each group, we sample a value of z , we compute

the values $F(s_t, a_t, z)$ and $B(s', z)$ for the states in that group, and we use all dot products $F(s_t, a_t, z)^\top B(s', z)$ within that group.

Thus, the solution we used just has $k = 1$: the same value of z is used throughout the minibatch, so we can compute N^2 dot products $F(s_t, a_t, z)^\top B(s', z)$ with only N forward passes through F and B . The other extreme would be $k = N$: for each (s_t, a_t, s_{t+1}) , we pick a value of z and a value of s' , compute $F(s_t, a_t, z)$ and $B(s', z)$, and form the dot product $F(s_t, a_t, z)^\top B(s', z)$. This uses more distinct values of z , but has only N dot products. In practice this would mean a reduced variance from sampling z , but an increased variance from sampling s' . Our chosen option ($k = 1$) has the opposite trade-off. This choice was based on preliminary results showing that, while using a higher k appeared to learn slightly faster at the beginning, the fewer dot products led to convergence to slightly lower performance.

To some extent, this effect may represent a hindrance for autoregressive FB compared with vanilla FB. However, Fig. 12 (Appendix D) shows that this effect is limited: indeed, vanilla FB with a single z per minibatch performs only slightly worse than vanilla FB with many z 's.

Normalization of B . As in Touati et al. (2023), to improve numerical conditioning on B , we use an auxiliary orthonormalization loss which ensures that B is approximately an orthogonal matrix. Indeed, one can change F and B without changing the B model, by $F \leftarrow FC$ and $B \leftarrow B(C^\top)^{-1}$ for any invertible matrix B , because the FB model only depends on the values of $F^\top B$ (Touati & Ollivier, 2021). In the case of auto-regressive features, since B depends on z , we can do this normalization separately for each z , without impacting the FB model. Explicitly, the orthonormalization loss is

$$\mathcal{L}_{\text{norm}}(B) := \mathbb{E}_z \left\| \mathbb{E}_{s \sim \rho} [B(s, z)B(s, z)^\top] - \text{Id} \right\|_{\text{Frobenius}}^2 \quad (32)$$

$$= \mathbb{E}_z \mathbb{E}_{s \sim \rho, s' \sim \rho} \left[(B(s, z)^\top B(s', z))^2 - \|B(s, z)\|_2^2 - \|B(s', z)\|_2^2 \right] + \text{Const.} \quad (33)$$

where z is sampled as for the main loss.

Moreover, we further normalize the output of B : for each input (s, z) and each auto-regressive block in the output of B , we set $B(s, z) \leftarrow B(s, z) \frac{\sqrt{d_k}}{\|B(s, z)\|}$ with d_k the size of the block k , so that each output component is of size approximately 1. Fig. 13 (Appendix D) tests the effect of not using this output normalization.

Scale invariance in the reward, and normalization of z . In reinforcement learning, the optimal behavior is the same for reward r or reward $\alpha \cdot r$ with any $\alpha > 0$. This inherent invariance property can be directly enforced within the architecture of FB to help with learning.

In non-autoregressive FB models, since $z = \mathbb{E}_s[r \times B(s)]$, the scaling $r \leftarrow \alpha \cdot r$ directly translates to a scaling in the task representation $z \leftarrow \alpha z$. Thus, one simple way to enforce invariance with respect to the rewards' scale is to always normalize z to a fixed norm (in practice, norm \sqrt{d} , so each of z value value

expectedly has a magnitude close to 1). We denote this operation by the preprocessing function $fz = \bar{z} = \sqrt{d} \times \frac{z}{\|z\|}$, which we apply when feeding any z to π_θ and F_ϕ . Hence, this ensures that $F(s, a, z) = F(s, a, \alpha \times z)$ and $\pi_z = \pi_{\alpha \times z}$ by construction. As a result, the predictions and behaviors made when facing rewards r and $\alpha \times r$ are exactly the same.

However, in our auto-regressive FB models where $z = \mathbb{E}_s[r \times B(s, z)]$, the same strategy cannot be applied when feeding z as input to B . This is because we never have access to the full z until the very end of the inference procedure. Hence, given a reward r , we still want to have scale-invariance in B but have no means of performing a standard input normalization, as we have no access to the magnitude of the resulting z .

A first strategy to counteract this limitation could be to normalize z only within each auto-regressive group, $f_g(z) = \{f_g(z)_1, \dots, f_g(z)_n\} = \{\bar{z}_1, \dots, \bar{z}_n\}$. While this would enforce a fixed scale and not pose test-time issues it would non-trivially affect the information available to B about the actual task z , losing any notion of relative magnitude between different groups. For instance, in case for $n = d$, each normalized component \bar{z}_k would be a binary scalar only preserving the *sign* from the corresponding task representation z_k .

We overcome these limitations by designing a new ‘*residual autoregressive*’ normalization strategy, f_{ar} compatible with the requirements of B ’s inference while still fully and exclusively preserving the information about z ’s direction as with traditional normalization.

Our strategy achieves these properties by using an iterative normalization scheme for each autoregressive component in its output $f_r(z) = \{f_r(z)_1, \dots, f_r(z)_n\}$. As with the aforementioned naive approach, we start with normalizing the first component z_1 within itself: $f_r(z)_1 = \bar{z}_1$. Then, we proceed to *residually normalize* all other z_k , also making use of all previous autoregressive components z_1, \dots, z_{k-1} :

$$f_r(z)_k = \frac{z_k}{\|z_{1:k}\|},$$

where $z_{1:k}$ simply corresponds to the concatenated first k auto-regressive components of z . Finally, we also rescale each k^{th} component by a constant factor $\sqrt{\sum_{i=0}^k d_i}$ to avoid biasing later components to have a smaller magnitude at initialization and incentivize $\|f_r(z)\| \approx \sqrt{d}$. We note there is a bijective map between traditional normalization and this auto-regressive scheme, thus, preserving the full information of the direction component of z without requiring the full vector.

A.3 Algorithms

FB-AWARE training is described in Algorithm 2. For reference, Algorithm 1 describes vanilla FB training.

Algorithm 1 FB training

```
1: Inputs Offline dataset  $\mathcal{D}$ , number of ensemble networks  $M$  for  $F$ , randomly
   initialized network  $\{F_{\theta_m}\}_{m \in [M]}$ ,  $B_\omega$  and  $\pi_\phi$ , transition mini-batch size  $I$  mixing
   probability  $\tau_{\text{mix}}$ , Polyak coefficient  $\zeta$ , orthonormality regularisation coefficient  $\lambda$ .

2: for  $t = 1, \dots$  do
3:   /* Sampling
4:   Sample  $I$  latent vectors
5:
6:    $z \sim \begin{cases} \mathcal{N}(0, I_d) & \text{with prob } 1 - \tau_{\text{mix}} \\ B(s) & \text{where } s \sim \mathcal{D}, \text{ with prob } \tau_{\text{mix}} \end{cases}$ 
7:    $z \leftarrow \sqrt{d} \frac{z}{\|z\|}$ 
8:   Sample a mini-batch of  $I$  transitions  $\{(s_i, a_i, s'_i)\}_{i \in [I]}$  from  $\mathcal{D}$ 
9:   /* Compute FB loss
10:  Sample  $a'_i \sim \pi_\phi(s'_i, z_i)$  for all  $i \in [I]$ 
11:   $\mathcal{L}_{\text{FB}}(\theta_m, \omega) = \frac{1}{2I(I-1)} \sum_{j \neq k} \left( F_{\theta_k}(s_i, a_i, z_i)^\top B_\omega(s'_k) - \gamma \frac{1}{M} \sum_{m \in [M]} F_{\theta_m^-}(s'_i, a'_i, z_i)^\top B_{\omega^-}(s'_k) \right)^2$ 
12:      $- \frac{1}{I} \sum_i F_{\theta_k}(s_i, a_i, z_i)^\top B_\omega(s'_i), \quad \forall m \in [M]$ 
13:  /* Compute orthonormality regularization loss
14:   $\mathcal{L}_{\text{ortho}}(\omega) = \frac{1}{2I(I-1)} \sum_{i \neq k} (B_\omega(s'_i)^\top B_\omega(s'_k))^2 - \frac{1}{I} \sum_i B_\omega(s'_i)^\top B_\omega(s'_i)$ 
15:  /* Compute actor loss
16:  Sample  $a_i^\phi \sim \pi_\phi(s_i, z_i)$  for all  $i \in [I]$ 
17:   $\mathcal{L}_{\text{actor}}(\phi) = -\frac{1}{I} \sum_i \left( \min_{m \in [M]} F_{\theta_m}(s_i, a_i^\phi, z_i)^\top z_i \right)$ 
18:  /* Update all networks
19:   $\theta_m \leftarrow \theta_m - \xi \nabla_{\theta_m} (\mathcal{L}_{\text{FB}}(\theta_k, \omega))$  for all  $m \in [M]$ 
20:   $\omega \leftarrow \omega - \xi \nabla_\omega (\sum_{l \in [m]} \mathcal{L}_{\text{FB}}(\theta_l, \omega) + \lambda \cdot \mathcal{L}_{\text{ortho}}(\omega))$ 
21:   $\phi \leftarrow \phi - \xi \nabla_\phi \mathcal{L}_{\text{actor}}(\phi)$ 
22:  /* Update target networks
23:   $\theta_m^- \leftarrow \zeta \theta_m^- + (1 - \zeta) \theta_m$  for all  $m \in [M]$ 
24:   $\omega^- \leftarrow \zeta \omega^- + (1 - \zeta) \omega$ 
25: end for
```

B Experimental Details

B.1 Network Architecture

- For the backward representation network $B(s, z)$, we first preprocess separately s and (s, z) . For the preprocessing of s , we use a feedforward neural network with one single hidden layer of 256 units. For (s, z) preprocessing, we use a masked network with one single hidden dimension of 256 units. The masked network employs multiplicative binary masks to remove some connections, such that each output layer unit of an autoregressive block is only predicted from the input units of previous blocks. After preprocessing, we concatenate the two outputs and pass them into a two hidden layer masked network that outputs a d -dimensional embedding.
- For the forward network $F(s, a, z)$, we first preprocess separately (s, a) and (s, z) by two feedforward networks with one single hidden layer (with 1024

Algorithm 2 FB-AWARE training

1: **Inputs** Offline dataset \mathcal{D} , number of ensemble networks M , randomly initialized network $\{F_{\theta_m}\}_{m \in [M]}$, B_ω and π_ϕ , transition mini-batch size J , latent vector mini-batch size I , number of autoregressive blocks K , mixing probability τ_{mix} , Polyak coefficient ζ , orthonormality regularisation coefficient λ , temperature β .

2: **for** $t = 1, \dots$ **do**

3: /* Sampling

4: Sample I latent vectors

5:

6: $z \sim \begin{cases} \mathcal{N}(0, I_d) & \text{with prob } 1 - \tau_{\text{mix}} \\ (z_1, \dots, z_K) = (B_1(s), \dots, B_K(s, z_1, \dots, z_{K-1})) & \text{where } s \sim \mathcal{D}, \text{ with prob } \tau_{\text{mix}} \end{cases}$

7: $z \leftarrow \sqrt{d} \frac{z}{\|z\|}$

8: Sample a mini-batch of $I \times J$ transitions $\{(s_{i,j}, a_{i,j}, s'_{i,j})\}_{i \in [I], j \in [J]}$ from \mathcal{D}

9: Sample $a'_{i,j} \sim \pi_\phi(s'_{i,j}, z_i)$ for all $i \in [I], j \in [J]$

10: $\mathcal{L}_{\text{FB}}(\theta_m, \omega) = \frac{1}{2IJ(J-1)} \sum_i \sum_{j \neq k} \left(F_{\theta_k}(s_{i,j}, a_{i,j}, z_i)^\top B_\omega(s'_{i,k}, z_i) - \gamma \frac{1}{M} \sum_{m \in [M]} F_{\theta_m^-}(s'_{i,j}, a'_{i,j}, z_i)^\top B_{\omega^-}(s'_{i,k}, z_i) \right)^2$

11: $- \frac{1}{IJ} \sum_i \sum_j F_{\theta_k}(s_{i,j}, a_{i,j}, z_i)^\top B_\omega(s'_{i,j}, z_i), \quad \forall m \in [M]$

12: /* Compute orthonormality regularization loss

13: $\mathcal{L}_{\text{ortho}}(\omega) = \frac{1}{2IJ(J-1)} \sum_i \sum_{j \neq k} (B_\omega(s'_{i,j}, z_i)^\top B_\omega(s'_{i,k}, z_i))^2 -$

14: $\frac{1}{IJ} \sum_i \sum_j B_\omega(s'_{i,j}, z_i)^\top B_\omega(s'_{i,j}, z_i)$

15: /* Compute actor loss

15: $A(s_{i,j}, a_{i,j}, z_i) \leftarrow \sum_m F_{\theta_m}(s_{i,j}, a_{i,j}, z_i)^\top z_i -$

15: $\mathbb{E}_{a'_{i,j} \sim \pi_\phi(s_{i,j}, z_i)} [\min_m F_{\theta_m}(s_{i,j}, a'_{i,j}, z_i)^\top z_i]$

16: $w(s_{i,j}, a_{i,j}, z_i) \leftarrow \frac{\exp(A(s_{i,j}, a_{i,j}, z_i)/\beta)}{\sum_{i', j'} \exp(A(s_{i', j'}, a_{i', j'}, z_{i'})/\beta)}$

17: $w'(s_{i,j}, a_{i,j}, z_i) \propto \frac{w(s_{i,j}, a_{i,j}, z_i)}{\sum_{(i', j') \neq (i, j)} w(s_{i', j'}, a_{i', j'}, z_{j'})}$

18: $\mathcal{L}_{\text{actor}}(\phi) = -\frac{1}{IJ} \sum_{i,j} w'(s_{i,j}, a_{i,j}, z_i) \log \pi_\phi(a_{i,j} | s_{i,j}, z_i)$

19: /* Update all networks

20: $\theta_m \leftarrow \theta_m - \xi \nabla_{\theta_m} (\mathcal{L}_{\text{FB}}(\theta_k, \omega))$ for all $m \in [M]$

21: $\omega \leftarrow \omega - \xi \nabla_\omega (\sum_{l \in [m]} \mathcal{L}_{\text{FB}}(\theta_l, \omega) + \lambda \cdot \mathcal{L}_{\text{ortho}}(\omega))$

22: $\phi \leftarrow \phi - \xi \nabla_\phi \mathcal{L}_{\text{actor}}(\phi)$

23: /* Update target networks

24: $\theta_m^- \leftarrow \zeta \theta_m^- + (1 - \zeta) \theta_m$ for all $m \in [M]$

25: $\omega^- \leftarrow \zeta \omega^- + (1 - \zeta) \omega$

26: **end for**

units), and a 512-dimensional output space. Then we concatenate these two outputs and pass it into a three-hidden-layer feedforward network (with 1024 units) to output a d -dimensional vector. We use an ensemble of two networks for F .

- For the policy network $\pi(s, z)$, we first preprocess separately s and (s, z) by two feedforward networks with one single hidden layer (with 1024 units) into a 512-dimensional output space. Then we concatenate these two outputs and pass it into another four hidden layer feedforward network (with 1024 units) to output a d_A -dimensional vector. Then we apply a **Tanh** activation as the action space is $[-1, 1]^{d_A}$.

For all the architectures, we apply a layer normalization and **Tanh** activation in the first layer in order to standardize the states and actions. We use **Relu** for the rest of layers.

B.2 Extended DeepMind Control Tasks

To evaluate our new unsupervised RL algorithm across a set of diverse unseen problems, we extend the DeepMind Control suite with 15 new unseen tasks as defined by the following objectives:

- **cheetah bounce** – Simulated cheetah agent is rewarded for advancing while elevating its trunk and maximizing vertical velocity.
- **cheetah march** – Simulated cheetah agent is rewarded for advancing at a constant pace.
- **cheetah stand** – Simulated cheetah agent is rewarded for standing upright on its back leg.
- **cheetah headstand** – Simulated cheetah agent is rewarded for standing on its head while raising its back leg.
- **quadruped bounce** – Simulated quadruped agent is rewarded for advancing while elevating its trunk and maximizing vertical velocity.
- **quadruped skip** – Simulated quadruped agent is rewarded for moving in a diagonal pattern across the environment.
- **quadruped march** – Simulated quadruped agent is rewarded for advancing at a constant pace.
- **quadruped trot** – Simulated quadruped agent is rewarded for advancing while minimizing feet contact with the ground.
- **walker flip** – Simulated walker for performing a cartwheel, flipping its body 360 degrees.
- **walker march** – Simulated walker agent is rewarded for advancing at a constant pace.
- **walker skyreach** – Simulated walker agent is rewarded for pushing either of its legs to maximize vertical reach.
- **walker pullup** – Simulated walker agent is rewarded for pushing its upper trunk vertically while keeping its feet firmly grounded.
- **humanoid dive** – Simulated humanoid agent is rewarded for diving head-first to maximize vertical velocity.
- **humanoid march** – Simulated humanoid agent is rewarded for advancing at a constant pace.

- **humanoid skip** – Simulated humanoid agent is rewarded for moving in a diagonal pattern across the environment.

We hope this new set of problems might facilitate the evaluation of simulated robotics agents for the broader RL field, even beyond the unsupervised setting.

C Full Tables of Results

C.1 Jaco Arm Results

Domain	Task	FB	FB-AW	FB-AWARE (4)	FB-AWARE (8)	LAP-AW
jaco	reach_bottom_left	49.0 \pm 25.5	43.9 \pm 8.6	56.3 \pm 8.6	63.6 \pm 6.4	25.9 \pm 5.9
jaco	reach_bottom_right	30.8 \pm 7.5	71.5 \pm 18.2	57.6 \pm 16.5	58.6 \pm 21.1	34.0 \pm 13.9
jaco	reach_random1	18.0 \pm 8.0	42.9 \pm 15.7	64.4 \pm 17.0	63.9 \pm 10.7	20.4 \pm 12.6
jaco	reach_random2	23.4 \pm 6.4	55.5 \pm 5.6	72.8 \pm 10.7	63.7 \pm 8.1	14.3 \pm 5.1
jaco	reach_random3	43.2 \pm 27.7	39.6 \pm 5.9	53.1 \pm 6.4	59.0 \pm 12.1	14.6 \pm 5.6
jaco	reach_random4	32.6 \pm 23.3	57.4 \pm 11.5	68.4 \pm 11.0	69.9 \pm 10.3	24.1 \pm 2.8
jaco	reach_top_left	32.6 \pm 12.3	41.0 \pm 5.4	41.9 \pm 8.3	62.7 \pm 14.9	10.3 \pm 2.2
jaco	reach_top_right	21.5 \pm 11.6	25.9 \pm 9.2	43.6 \pm 9.7	48.3 \pm 12.1	21.4 \pm 5.2
jaco	Average	31.4 \pm 15.3	47.2 \pm 10.0	57.3 \pm 11.0	61.2 \pm 12.0	20.6 \pm 6.7

Table 2: JACO results on RND dataset, with dimension $d = 64$

Domain	Task	FB	FB-AW	FB-AWARE (4)	FB-AWARE (8)	LAP-AW
jaco	reach_bottom_left	33.8 \pm 17.3	76.0 \pm 12.0	88.1 \pm 18.5	64.6 \pm 13.4	41.3 \pm 10.2
jaco	reach_bottom_right	51.3 \pm 10.2	86.3 \pm 9.4	87.7 \pm 15.0	96.8 \pm 6.9	47.8 \pm 18.1
jaco	reach_random1	32.6 \pm 18.1	75.3 \pm 10.5	85.4 \pm 9.2	87.0 \pm 13.0	30.9 \pm 5.3
jaco	reach_random2	22.9 \pm 10.0	86.3 \pm 9.1	104.1 \pm 7.5	95.5 \pm 12.7	28.8 \pm 6.3
jaco	reach_random3	31.2 \pm 9.0	68.3 \pm 11.3	89.5 \pm 17.7	61.9 \pm 8.2	24.7 \pm 7.2
jaco	reach_random4	21.6 \pm 6.3	82.2 \pm 9.6	101.2 \pm 17.6	82.9 \pm 10.5	34.7 \pm 10.9
jaco	reach_top_left	44.4 \pm 16.6	59.5 \pm 18.3	56.5 \pm 9.6	46.0 \pm 17.7	32.1 \pm 10.2
jaco	reach_top_right	28.3 \pm 13.0	44.2 \pm 12.5	47.5 \pm 5.6	39.7 \pm 8.4	23.1 \pm 7.1
jaco	Average	33.3 \pm 12.6	72.2 \pm 11.6	82.5 \pm 12.6	71.8 \pm 11.4	32.9 \pm 9.4

Table 3: JACO results on RND dataset, with dimension $d = 128$

Domain	Task	FB	FB-AW	FB-AWARE (4)	FB-AWARE (8)	LAP-AW
jaco	reach_bottom_left	33.1 \pm 23.3	60.2 \pm 26.8	101.0 \pm 7.5	91.5 \pm 16.8	42.7 \pm 15.0
jaco	reach_bottom_right	10.5 \pm 2.4	69.2 \pm 32.5	116.9 \pm 9.7	90.9 \pm 12.9	63.9 \pm 13.8
jaco	reach_random1	31.9 \pm 21.8	49.5 \pm 25.5	86.0 \pm 19.2	84.7 \pm 12.4	45.4 \pm 4.8
jaco	reach_random2	32.1 \pm 16.4	50.9 \pm 26.2	99.7 \pm 17.6	94.9 \pm 14.5	47.2 \pm 8.4
jaco	reach_random3	43.6 \pm 18.4	39.6 \pm 22.4	72.4 \pm 13.6	81.5 \pm 9.9	44.9 \pm 9.4
jaco	reach_random4	32.8 \pm 15.4	58.3 \pm 34.5	98.8 \pm 15.7	93.1 \pm 14.2	51.3 \pm 8.2
jaco	reach_top_left	29.4 \pm 11.2	38.4 \pm 18.6	43.6 \pm 10.3	65.1 \pm 21.1	43.0 \pm 6.1
jaco	reach_top_right	27.7 \pm 6.2	29.8 \pm 19.2	44.8 \pm 20.1	45.4 \pm 13.8	37.0 \pm 9.4
jaco	Average	30.1 \pm 14.4	49.5 \pm 25.7	82.9 \pm 14.2	80.9 \pm 14.4	46.9 \pm 9.4

Table 4: JACO results on RND dataset, with dimension $d = 256$

Domain	Task	LAP	LAP-AW	FB-AW	FB	FB-ARE (4)	FB-ARE (8)
cheetah	walk	641.0 \pm 137.7	528.9 \pm 22.9	520.8 \pm 56.6	780.3 \pm 182.7	737.5 \pm 204.7	686.1 \pm 44.8
cheetah	run	156.5 \pm 36.8	116.0 \pm 6.8	141.2 \pm 19.8	306.8 \pm 87.3	261.1 \pm 95.3	241.8 \pm 57.3
cheetah	walk_backward	930.9 \pm 77.9	452.3 \pm 212.6	839.5 \pm 49.3	732.5 \pm 167.3	769.9 \pm 209.9	762.3 \pm 205.6
cheetah	run_backward	230.7 \pm 42.6	109.7 \pm 37.2	196.2 \pm 22.4	136.5 \pm 35.6	174.1 \pm 48.6	186.4 \pm 66.4
cheetah	in_dataset_avg	489.8 \pm 73.8	301.7 \pm 69.9	424.4 \pm 37.0	489.0 \pm 118.2	485.7 \pm 139.6	469.2 \pm 93.5
quadruped	walk	509.1 \pm 38.9	418.3 \pm 42.7	438.2 \pm 192.6	608.4 \pm 72.0	630.1 \pm 96.9	604.0 \pm 116.8
quadruped	run	457.6 \pm 27.7	355.6 \pm 86.5	391.2 \pm 91.9	392.7 \pm 31.4	417.4 \pm 30.6	376.1 \pm 29.2
quadruped	stand	681.6 \pm 221.6	731.3 \pm 166.6	762.2 \pm 152.8	687.9 \pm 29.6	761.7 \pm 75.9	705.4 \pm 58.1
quadruped	jump	464.5 \pm 167.3	493.4 \pm 147.1	563.7 \pm 139.1	567.0 \pm 10.6	609.2 \pm 42.3	580.3 \pm 37.1
quadruped	in_dataset_avg	528.2 \pm 113.9	499.6 \pm 110.7	538.9 \pm 144.1	564.0 \pm 35.9	604.6 \pm 61.4	566.4 \pm 60.3
walker	stand	963.6 \pm 15.3	803.3 \pm 61.9	452.6 \pm 85.8	728.5 \pm 83.0	632.7 \pm 151.7	516.1 \pm 191.0
walker	walk	908.8 \pm 28.1	605.3 \pm 36.4	572.0 \pm 25.3	669.9 \pm 46.6	607.9 \pm 140.2	552.2 \pm 268.4
walker	run	318.7 \pm 15.0	196.6 \pm 13.1	181.2 \pm 16.5	356.2 \pm 20.9	290.4 \pm 22.7	240.0 \pm 122.7
walker	spin	982.9 \pm 3.5	627.9 \pm 135.1	963.7 \pm 5.3	974.9 \pm 10.0	983.4 \pm 1.3	788.2 \pm 391.2
walker	in_dataset_avg	793.5 \pm 15.5	558.3 \pm 61.6	542.4 \pm 33.2	682.4 \pm 40.1	628.6 \pm 79.0	524.1 \pm 243.3
cheetah	bounce	600.4 \pm 23.0	428.2 \pm 210.5	539.9 \pm 25.8	415.7 \pm 119.2	472.8 \pm 44.2	462.2 \pm 23.8
cheetah	march	290.8 \pm 63.2	233.6 \pm 12.9	279.2 \pm 39.0	561.4 \pm 183.7	531.5 \pm 187.6	460.9 \pm 119.2
cheetah	stand	790.1 \pm 107.9	249.5 \pm 131.6	738.7 \pm 66.2	780.9 \pm 105.7	629.3 \pm 49.3	762.9 \pm 179.9
cheetah	headstand	577.9 \pm 145.1	288.9 \pm 236.8	728.4 \pm 83.3	794.7 \pm 9.4	791.1 \pm 52.0	765.1 \pm 70.3
cheetah	out_of_dataset_avg	564.8 \pm 84.8	300.0 \pm 147.9	571.6 \pm 53.6	638.2 \pm 104.5	606.2 \pm 83.3	612.8 \pm 98.3
quadruped	bounce	179.5 \pm 76.1	123.2 \pm 64.4	189.3 \pm 192.0	276.1 \pm 57.1	196.0 \pm 105.6	251.1 \pm 47.9
quadruped	skip	365.3 \pm 108.4	458.0 \pm 114.5	559.5 \pm 220.3	603.3 \pm 30.7	635.2 \pm 34.4	615.3 \pm 35.2
quadruped	march	478.7 \pm 14.1	370.4 \pm 80.7	396.4 \pm 123.7	458.3 \pm 20.4	466.0 \pm 38.7	419.7 \pm 32.0
quadruped	trot	310.6 \pm 7.0	246.7 \pm 54.6	278.3 \pm 122.4	357.6 \pm 12.1	380.9 \pm 47.6	335.3 \pm 35.4
quadruped	out_of_dataset_avg	333.5 \pm 51.4	299.6 \pm 78.5	355.9 \pm 164.6	423.8 \pm 30.1	419.5 \pm 56.6	405.4 \pm 37.6
walker	flip	605.4 \pm 42.4	435.0 \pm 27.7	293.6 \pm 83.8	445.5 \pm 77.4	462.0 \pm 68.7	322.6 \pm 147.7
walker	march	695.8 \pm 47.8	364.3 \pm 29.0	359.9 \pm 15.2	518.4 \pm 95.9	400.5 \pm 178.0	390.3 \pm 190.4
walker	skyreach	653.8 \pm 64.1	423.1 \pm 61.6	406.0 \pm 13.8	417.0 \pm 34.9	331.7 \pm 44.4	261.5 \pm 151.7
walker	pullup	264.8 \pm 80.5	58.4 \pm 39.2	264.5 \pm 60.5	305.9 \pm 109.4	463.1 \pm 107.5	214.6 \pm 155.1
walker	out_of_dataset_avg	554.9 \pm 58.7	320.2 \pm 39.4	331.0 \pm 43.3	421.7 \pm 79.4	414.3 \pm 99.7	297.2 \pm 161.2

Table 5: DMC Locomotion results on RND dataset, with dimension 64, averaged over 100 episodes. Humanoid is not included, as RND produces insufficient exploration for Humanoid: even classical single-task (non-FB) training fails.

C.2 DMC Locomotion Results

Domain	Task	FB	FB-AW	FB-AWARE (4)	FB-AWARE (8)	LAP-AW
cheetah	walk	985.3 \pm 3.1	983.6 \pm 6.6	967.5 \pm 28.4	982.0 \pm 4.2	978.5 \pm 14.0
cheetah	run	213.2 \pm 123.5	560.2 \pm 40.7	525.7 \pm 53.6	547.6 \pm 20.5	448.9 \pm 222.9
cheetah	walk_backward	971.0 \pm 24.3	979.7 \pm 5.1	984.1 \pm 0.7	985.1 \pm 1.2	982.8 \pm 3.5
cheetah	run_backward	302.5 \pm 58.9	473.9 \pm 6.0	454.0 \pm 17.9	465.8 \pm 9.9	413.8 \pm 25.1
cheetah	in_dataset_avg	618.0 \pm 52.4	749.4 \pm 14.6	732.8 \pm 25.1	745.1 \pm 8.9	706.0 \pm 66.4
quadruped	walk	389.5 \pm 238.2	935.5 \pm 6.6	926.8 \pm 4.2	919.3 \pm 8.4	819.1 \pm 132.3
quadruped	run	298.1 \pm 105.4	580.8 \pm 62.5	606.2 \pm 33.0	566.0 \pm 47.8	610.6 \pm 89.3
quadruped	stand	615.4 \pm 191.0	941.1 \pm 6.6	947.9 \pm 4.7	940.3 \pm 9.4	911.5 \pm 27.9
quadruped	jump	429.6 \pm 134.5	779.1 \pm 48.1	841.8 \pm 8.2	751.0 \pm 75.9	782.8 \pm 48.4
quadruped	in_dataset_avg	433.2 \pm 167.3	809.1 \pm 30.9	830.7 \pm 12.5	794.2 \pm 35.4	781.0 \pm 74.5
walker	stand	744.0 \pm 119.3	962.2 \pm 14.7	963.9 \pm 3.7	963.4 \pm 4.3	961.2 \pm 8.2
walker	walk	780.0 \pm 310.8	943.8 \pm 20.8	941.4 \pm 7.4	922.4 \pm 16.4	934.3 \pm 12.7
walker	run	422.5 \pm 167.4	594.9 \pm 12.1	606.5 \pm 6.2	600.8 \pm 37.8	518.9 \pm 39.7
walker	spin	481.6 \pm 226.3	894.7 \pm 84.4	820.9 \pm 114.6	894.8 \pm 63.6	802.0 \pm 180.7
walker	in_dataset_avg	607.0 \pm 205.9	848.9 \pm 33.0	833.2 \pm 33.0	845.3 \pm 30.5	804.1 \pm 60.4
humanoid	walk	9.5 \pm 11.8	793.5 \pm 16.1	789.4 \pm 18.0	791.5 \pm 7.3	715.4 \pm 35.1
humanoid	stand	7.7 \pm 3.8	720.0 \pm 23.5	728.3 \pm 30.1	711.6 \pm 24.4	587.2 \pm 34.0
humanoid	run	2.4 \pm 1.6	276.5 \pm 11.0	266.7 \pm 4.8	273.6 \pm 5.2	246.4 \pm 9.7
humanoid	in_dataset_avg	6.5 \pm 5.7	596.7 \pm 16.9	594.8 \pm 17.6	592.3 \pm 12.3	516.3 \pm 26.3
cheetah	bounce	351.8 \pm 119.9	479.0 \pm 22.3	506.9 \pm 18.9	494.6 \pm 9.4	338.4 \pm 38.3
cheetah	march	521.4 \pm 161.8	897.8 \pm 38.1	903.6 \pm 30.6	921.9 \pm 8.9	819.7 \pm 60.8
cheetah	stand	731.5 \pm 248.4	419.1 \pm 48.7	472.7 \pm 49.4	548.7 \pm 37.7	184.1 \pm 44.9
cheetah	headstand	560.4 \pm 185.9	849.7 \pm 49.9	848.0 \pm 33.8	806.0 \pm 21.4	7.2 \pm 11.4
cheetah	out_of_dataset_avg	541.3 \pm 179.0	661.4 \pm 39.8	682.8 \pm 33.2	692.8 \pm 19.4	337.3 \pm 38.9
quadruped	bounce	114.7 \pm 98.1	181.2 \pm 61.4	284.2 \pm 31.7	223.7 \pm 20.6	202.8 \pm 52.1
quadruped	skip	425.0 \pm 139.2	654.6 \pm 88.0	835.2 \pm 86.9	705.4 \pm 67.8	769.2 \pm 94.6
quadruped	march	304.1 \pm 133.2	747.1 \pm 94.2	791.8 \pm 30.6	800.0 \pm 7.9	742.5 \pm 123.0
quadruped	trot	228.3 \pm 127.1	573.0 \pm 46.9	614.8 \pm 15.9	594.6 \pm 6.1	509.8 \pm 73.1
quadruped	out_of_dataset_avg	268.0 \pm 124.4	539.0 \pm 72.6	631.5 \pm 41.3	580.9 \pm 25.6	556.1 \pm 85.7
walker	flip	404.7 \pm 234.4	909.9 \pm 20.0	913.6 \pm 14.5	914.5 \pm 18.3	780.2 \pm 66.2
walker	march	663.3 \pm 245.6	826.7 \pm 41.3	841.0 \pm 45.1	811.8 \pm 17.7	725.8 \pm 25.7
walker	skyreach	396.5 \pm 35.2	365.2 \pm 42.7	366.7 \pm 25.2	404.3 \pm 51.5	284.6 \pm 52.7
walker	pullup	124.1 \pm 101.7	334.4 \pm 112.7	523.7 \pm 31.4	454.8 \pm 46.9	113.9 \pm 48.8
walker	out_of_dataset_avg	397.2 \pm 154.2	609.0 \pm 54.2	661.3 \pm 29.1	646.4 \pm 33.6	476.1 \pm 48.3
humanoid	dive	165.8 \pm 7.9	404.1 \pm 11.1	409.8 \pm 17.9	396.5 \pm 16.1	242.6 \pm 18.1
humanoid	march	6.3 \pm 8.4	669.6 \pm 23.5	659.5 \pm 23.6	661.2 \pm 18.3	559.6 \pm 53.1
humanoid	skip	2.1 \pm 0.9	233.7 \pm 39.2	234.9 \pm 11.4	221.4 \pm 15.1	126.0 \pm 18.2
humanoid	out_of_dataset_avg	58.1 \pm 5.7	435.8 \pm 24.6	434.7 \pm 17.6	426.4 \pm 16.5	309.4 \pm 29.8

Table 6: DMC locomotion results on MOOD dataset, with dimension = 64 for walker, cheetah, quadruped, and dimension = 128 for humanoid, averaged over 100 episodes

C.3 Additional reward prompts

We demonstrate the adaptability of our FB-AWARE model on the DMC humanoid by showcasing its behavior in response to various reward functions. In Figure 5, we illustrate the agent’s actions when prompted by the following reward functions:

- LEFT_HAND: the task consists in raising the left hand while standing. Specifically, the reward function is defined as having a velocity close to zero

(exponential term), having an upright torso, and maintaining the height of the left wrist above a certain threshold while keeping the height of the right wrist below a different threshold.

$$R_{\text{LEFT_HAND}} = \exp(-(v_x^2 + v_y^2)) * \text{upright} * \mathbb{I}\{\text{left_wrist_z} > 2\} * \mathbb{I}\{\text{right_wrist_z} < 0.9\}$$

- **RIGHT_HAND**: the task consists in raising the right hand while standing. Specifically, the reward function is defined as having a velocity close to zero (exponential term), having an upright torso, and maintaining the height of the right wrist above a certain threshold while keeping the height of the left wrist below a different threshold.

$$R_{\text{RIGHT_HAND}} = \exp(-(v_x^2 + v_y^2)) * \text{upright} * \mathbb{I}\{\text{left_wrist_z} < 0.9\} * \mathbb{I}\{\text{right_wrist_z} > 0.9\}$$

- **WALK_OPEN_HAND**: the task consists in walking while keeping the two hands open. The reward function is defined as having a velocity above some threshold and the absolute distance between the y coordinate of the left and right wrist above some threshold.

$$R_{\text{WALK_OPEN_HAND}} = \mathbb{I}\{v_x^2 + v_y^2 > 5\} * \mathbb{I}\{|\text{left_wrist_y} - \text{right_wrist_y}| > 1.2\}$$

- **SPLIT**: the task consists in doing a split on the ground. The reward can be described as having a velocity close to zero, the height of the pelvis below some threshold and the absolute distance between the y coordinate of the left and right ankle above some threshold.

$$R_{\text{SPLIT}} = \exp(-(v_x^2 + v_y^2)) * \mathbb{I}\{z_{\text{pelvis}} < 0.2\} * \mathbb{I}\{|\text{left_ankle_y} - \text{right_ankle_y}| > 0.5\}$$

D Ablations

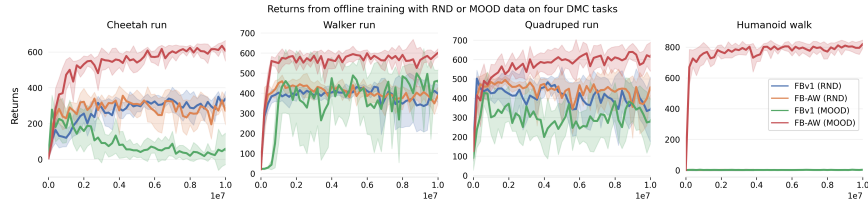


Figure 6: Performance on four representative tasks of the DMC when training FB-AW and vanilla FB (FBv1) either from mixed objective MOOD or pure RND data. The advantage of AW is clear on the mixed-objective MOOD datasets. The RND dataset does not allow FB to reach top performance.

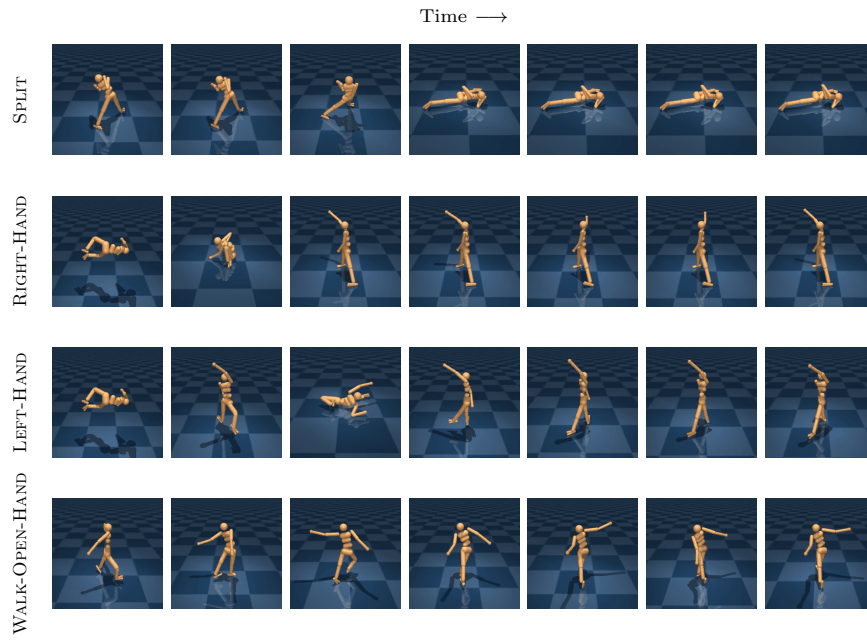


Figure 5: Example of behaviors inferred by from reward equations.

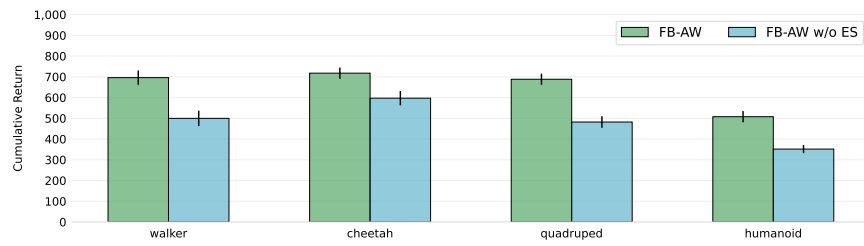


Figure 7: Cumulative reward averaged over all DMC Locomotion tasks, achieved by FB-AW with or without Evaluation-Sampling, trained on MOOD dataset.

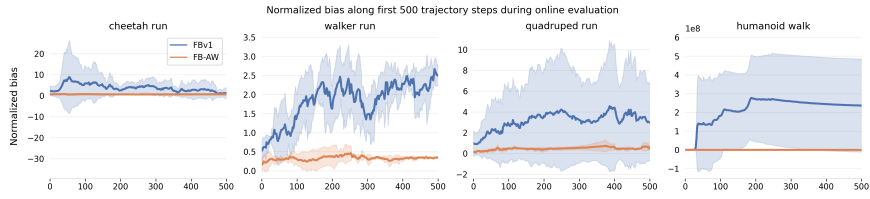


Figure 8: AW improves accuracy of reward prediction by the B model. We plot the bias of estimated rewards when progressing through a trajectory, namely, the difference between the actual trajectory return and the return $\sum_t B(s_t)^T z$ predicted by FB, after offline training on MOOD (averaged across 5 agents, 10 trajectories each). Vanilla FB provides overoptimistic values.

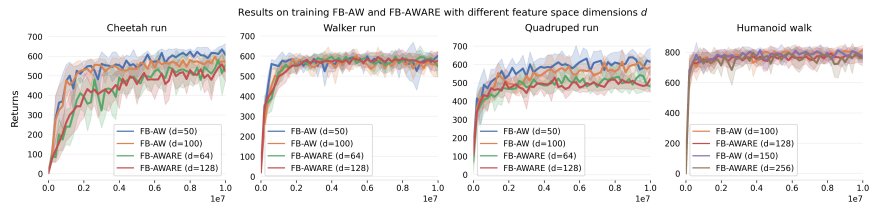


Figure 9: Effect of the z dimension for FB-AW and FB-AWARE on the MOOD mixed objective datasets

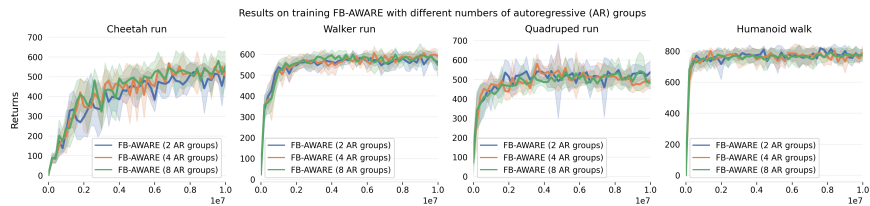


Figure 10: Effect of the number of autoregressive groups for FB-AWARE on the MOOD mixed objective datasets

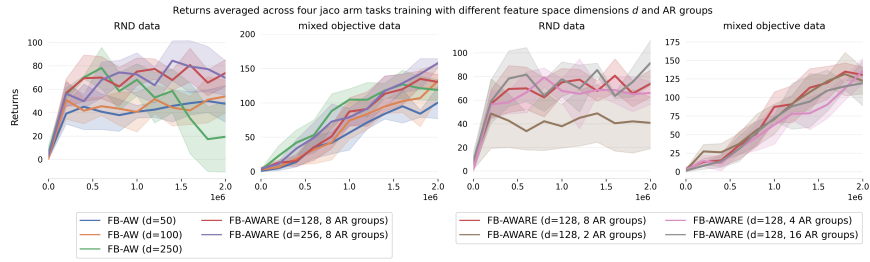


Figure 11: Effect of modifying the z dimension and the number of autoregressive groups for FB-AW and FB-AWARE, for performance in the Jaco arm environment

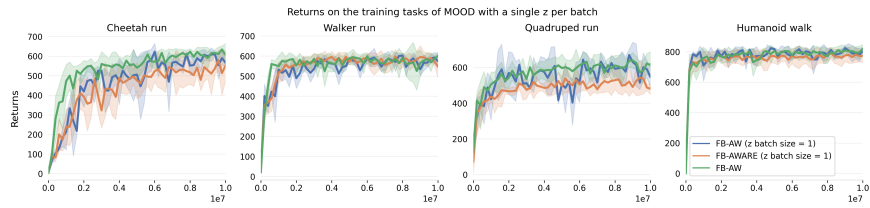


Figure 12: Effect of training FB-AW with a single z -per-batch like FB-AWARE (Section A.2)

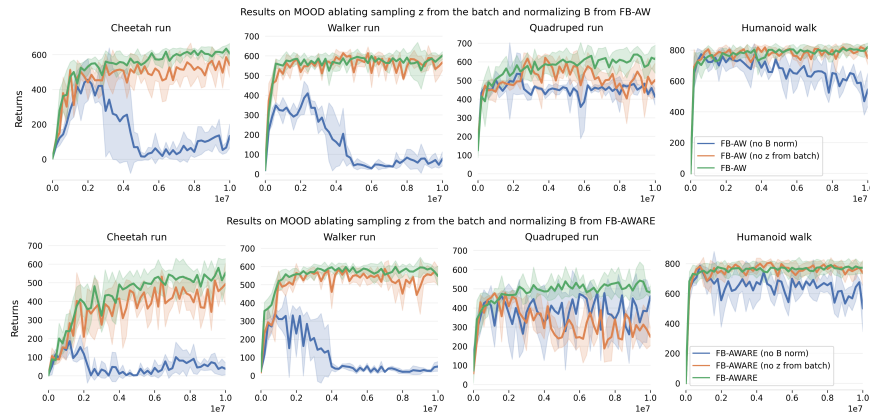


Figure 13: Ablations for Appendix A.2: Effects from training without the B normalization and without sampling 50% of z from other states in the minibatch for FB-AW (Top) and FB-AWARE (Bottom)

Domain	Task	Vanilla FB	FB(fully par. + no min.)	FB-AW(no fully par. + min.)	FB-AW (WIS)	FB-AW
cheetah	walk	275.5 \pm 343.5	500.2 \pm 285.2	987.5 \pm 0.6	975.0 \pm 23.5	975.6 \pm 24.6
	run	77.1 \pm 87.3	106.0 \pm 50.4	548.3 \pm 38.6	548.4 \pm 30.0	610.7 \pm 49.2
	walk_backward	199.3 \pm 209.1	349.7 \pm 394.5	984.8 \pm 0.4	984.5 \pm 1.2	984.7 \pm 0.4
	run_backward	49.1 \pm 58.6	78.6 \pm 118.3	461.0 \pm 6.6	477.4 \pm 7.8	481.4 \pm 4.2
cheetah	in_dataset_avg	0.175 \pm 0.203	0.295 \pm 0.249	0.939 \pm 0.019	0.944 \pm 0.022	0.971 \pm 0.028
quadruped	walk	318.3 \pm 65.7	763.9 \pm 145.2	889.4 \pm 80.5	938.8 \pm 34.7	958.0 \pm 9.1
	run	279.4 \pm 45.8	439.8 \pm 90.2	446.1 \pm 56.6	599.4 \pm 70.4	673.4 \pm 33.1
	stand	619.6 \pm 121.1	831.5 \pm 147.8	901.3 \pm 84.7	944.9 \pm 44.0	975.9 \pm 4.5
	jump	435.0 \pm 79.8	644.1 \pm 127.7	708.2 \pm 20.2	767.4 \pm 74.9	798.4 \pm 55.0
quadruped	in_dataset_avg	0.493 \pm 0.093	0.800 \pm 0.153	0.877 \pm 0.070	0.976 \pm 0.070	1.026 \pm 0.032
walker	stand	906.5 \pm 80.4	976.2 \pm 5.5	964.8 \pm 3.5	978.9 \pm 3.1	956.1 \pm 34.7
	walk	892.8 \pm 102.7	939.9 \pm 44.7	946.7 \pm 10.3	960.1 \pm 7.3	955.4 \pm 16.6
	run	462.2 \pm 37.5	487.0 \pm 44.1	480.5 \pm 52.6	583.3 \pm 20.9	579.8 \pm 45.8
	spin	422.1 \pm 121.4	464.9 \pm 196.4	923.2 \pm 30.1	759.2 \pm 173.6	789.6 \pm 117.7
walker	in_dataset_avg	0.749 \pm 0.093	0.799 \pm 0.081	0.913 \pm 0.032	0.918 \pm 0.055	0.917 \pm 0.061
humanoid	walk	3.6 \pm 5.0	2.3 \pm 1.1	677.6 \pm 92.8	779.9 \pm 37.1	785.0 \pm 20.4
	stand	5.1 \pm 1.9	4.8 \pm 0.8	481.3 \pm 59.5	750.3 \pm 43.3	801.7 \pm 45.8
	run	1.1 \pm 0.8	1.1 \pm 0.6	256.8 \pm 22.9	274.4 \pm 20.4	294.9 \pm 17.5
humanoid	in_dataset_avg	0.005 \pm 0.004	0.004 \pm 0.002	0.793 \pm 0.074	0.965 \pm 0.058	1.014 \pm 0.049
cheetah	bounce	109.6 \pm 83.9	315.4 \pm 102.8	436.3 \pm 47.1	469.3 \pm 39.7	502.5 \pm 19.2
	march	135.8 \pm 182.5	257.8 \pm 110.2	766.9 \pm 92.7	892.0 \pm 35.0	917.4 \pm 13.4
	stand	206.6 \pm 383.7	684.7 \pm 267.4	426.9 \pm 191.6	288.9 \pm 162.6	387.6 \pm 98.4
	headstand	233.8 \pm 369.0	923.7 \pm 57.1	488.0 \pm 345.1	407.0 \pm 369.8	854.5 \pm 42.0
cheetah	out_dataset_avg	0.214 \pm 0.308	0.684 \pm 0.166	0.660 \pm 0.207	0.642 \pm 0.188	0.832 \pm 0.053
quadruped	bounce	133.6 \pm 82.8	248.9 \pm 41.0	204.3 \pm 61.0	246.6 \pm 31.3	231.8 \pm 18.4
	skip	478.9 \pm 113.2	601.0 \pm 72.3	554.0 \pm 1.3	722.7 \pm 112.4	836.1 \pm 104.4
	march	280.8 \pm 73.1	517.2 \pm 108.0	478.8 \pm 41.8	802.1 \pm 55.9	860.2 \pm 26.7
	trot	178.3 \pm 38.0	381.5 \pm 82.6	412.1 \pm 48.6	605.3 \pm 16.9	615.9 \pm 13.6
quadruped	out_dataset_avg	0.524 \pm 0.175	0.897 \pm 0.159	0.834 \pm 0.105	1.187 \pm 0.105	1.245 \pm 0.075
walker	flip	630.3 \pm 111.7	744.4 \pm 105.4	771.8 \pm 104.5	891.2 \pm 24.9	896.5 \pm 41.2
	march	709.4 \pm 182.1	744.6 \pm 183.3	593.9 \pm 12.6	788.3 \pm 63.7	797.2 \pm 36.8
	skyreach	321.7 \pm 130.0	423.0 \pm 70.3	392.7 \pm 15.1	364.3 \pm 48.0	389.8 \pm 6.5
	pullup	114.7 \pm 98.5	101.3 \pm 104.1	33.9 \pm 8.2	309.7 \pm 145.5	376.5 \pm 219.5
walker	out_dataset_avg	0.645 \pm 0.196	0.747 \pm 0.163	0.674 \pm 0.052	0.849 \pm 0.101	0.889 \pm 0.101
humanoid	dive	159.2 \pm 12.2	166.6 \pm 29.5	357.1 \pm 36.6	387.9 \pm 39.8	404.0 \pm 5.7
	march	2.9 \pm 3.5	2.9 \pm 1.8	479.6 \pm 48.5	623.1 \pm 64.6	645.9 \pm 49.7
	skip	2.0 \pm 1.2	1.7 \pm 0.4	81.3 \pm 26.7	254.4 \pm 23.8	250.1 \pm 24.0
humanoid	out_dataset_avg	0.112 \pm 0.011	0.117 \pm 0.021	0.536 \pm 0.072	0.788 \pm 0.079	0.805 \pm 0.049

Table 7: Ablations regarding the various components from Section 3.2: advantage weighting, using the average versus the min of the two target networks for representing uncertainty, using fully parallel architectures for the two target networks, and using improved weighted importance sampling (IWIS) versus ordinary WIS. As described in the text, FB-AW (right column) has advantage weighting, uses the average instead of the min, has fully parallel architectures, and uses IWIS. Vanilla FB (left column) has the opposite settings. We compare other combinations in between. We report performance on the mixed objective datasets from MOOD, on both in-dataset and out-of-dataset tasks. The representation dimension is $d = 50$ for Cheetah, Quadruped, Walker, and $d = 100$ for Humanoid.