ZERO-SHOT ADAPTATION OF BEHAVIORAL FOUNDA-TION MODELS TO UNSEEN DYNAMICS

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

016

017

018

019

021

024

025

026

027

028

031

032

033

034

037

038

040

041

043

044

046

047

048

051

052

ABSTRACT

Behavioral Foundation Models (BFMs) proved successful in producing nearoptimal policies for arbitrary tasks in a zero-shot manner, requiring no test-time retraining or task-specific fine-tuning. Among the most promising BFMs are the ones that estimate the successor measure learned in an unsupervised way from task-agnostic offline data. However, these methods fail to react to changes in the dynamics, making them inefficient under partial observability or when the transition function changes. This hinders the applicability of BFMs in a real-world setting, e.g., in robotics, where the dynamics can unexpectedly change at test time. In this work, we demonstrate that Forward–Backward (FB) representation, one of the methods from the BFM family, cannot produce reasonable policies under distinct dynamics, leading to an interference among the latent policy representations. To address this, we propose an FB model with a transformer-based belief estimator, which greatly facilitates zero-shot adaptation. Additionally, we show that partitioning the policy encoding space into dynamics-specific clusters, aligned with the context-embedding directions, yields additional gain in performance. Those traits allow our method to respond to the dynamics mismatches observed during training and to generalize to unseen ones. Empirically, in the changing dynamics setting, our approach achieves up to a 2x higher zero-shot returns compared to the baselines for both discrete and continuous tasks.

1 Introduction

One very desirable property of reinforcement learning (RL) agents is their ability to adapt during test-time to new tasks or to environment changes, without requiring any fine-tuning or planning. Achieving this in as few trials as possible would be even better: the ideal being the zero-shot adaptation (Touati et al., 2022), where the agent never interacts with the environment at test-time and relies solely on the task-agnostic data. Behavioral Foundational Models (BFMs) (Sikchi et al., 2024; Tirinzoni et al.) may be considered as a step in this direction, because they can learn a variety of policies from offline data without knowing the rewards. During inference, it is possible to extract a task-specific policy that is theoretically optimal in terms of performance. Recent work (Tirinzoni et al.) demonstrates that methods based on *successor measure* estimation through Forward-Backward (FB) decomposition (Touati & Ollivier, 2021), is especially versatile and can successfully imitate diverse behaviors from provided data.

At the same time, FB possesses a fundamental drawback that limits its adaptation ability. In our paper, we show that FB is unable to generalize across different environment configurations (dynamics), such as changes in a transition function (*e.g.*, new obstacles) or some latent factor variation (*e.g.*, wind direction). This limitation stems from the way the *successor measure* (Dayan, 1993) is estimated: FB averages the discounted future-occupancy state distribution over all observed dynamics, which inevitably causes *interference* in a policy representation space. This fact alone may severely constrain the applicability of FB in the real-world scenarios. For example, one of the largest robotics dataset, Open X-Embodiment Collaboration (2023), consists of 22 different robot embodiments, and training FB on each of them independently is infeasible. In Section 3.1, we discuss this limitation and support our claims both theoretically and empirically.

To remedy this, we introduce Belief-FB (BFB), a conditioning method for FB through a *belief* estimation, a popular technique of uncertainty quantification in Meta-RL (Zintgraf et al., 2020; Dorfman et al., 2021). To implement this, we employ a permutation-invariant transformer encoder, denoted

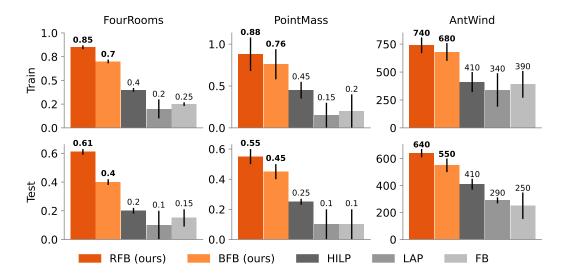


Figure 1: **Summary of results**. Aggregate mean performance over *seen* (train) and *unseen* (test) dynamics for zero-shot RL. The error bars indicate standard deviation over three seeds. Notably, both BFB and RFB adapt not only to the dynamics seen during training but are also able to generalize to unseen dynamics. There are 30 (20) training (test) dynamics for FourRooms and PointMass and 16 (4) for AntWind environments.

as $f_{\rm dyn}$, which processes a given trajectory from the dataset to produce a dynamics-specific vector h. This vector is subsequently utilized as a conditioning input to the future outcomes representation function, expressed as $F(\cdot,\cdot,h,\cdot)$. We pre-train $f_{\rm dyn}$ in a self-supervised fashion, thus posing no additional requirements on the data structure or the trajectory re-labeling, while maintaining theoretical guarantees. We discuss the implementation of Belief-FB in Section 3.2.

Remarkably, Belief-FB enables the generalization capabilities of FB not only through the dynamics seen **in the training dataset**, but also on the **unseen test configurations** never present in the offline data. We also find that in order to align *belief* estimation better with FB, one also needs to partition the policy space encodings prior into dynamics-specific clusters, so we propose Rotation-FB (RFB) that accomplishes this partitioning. We present the theoretical support and the implementation details of Rotation-FB in Section 3.3. Empirically, both BFB and RFB outperform baselines for seen and unseen dynamics, as gathered in Figure 1 and discussed in Section 4.3.

We believe that our work sufficiently broadens the possible applicability of BFMs, yet keeping all of the zero-shot properties unchanged. Our contributions are as follows:

- We demonstrate the limitation of Forward-Backward (FB) representations (Touati & Ollivier, 2021), which lies in its inability to generalize *per se* across different dynamics both from train and test, where dynamics shift constitute of new layout grids or changes in the transition function that are hidden from an agent. Refer to Section 3.1 for more discussion.
- We propose Belief–FB (BFB), which employs a transformer encoder to infer a belief over the current dynamics (Zintgraf et al., 2020; Dorfman et al., 2021). Analyzing BFB's policy encoding space reveals that additional disentanglement is beneficial, motivating our Rotation–FB (RFB) extension. Section 3.2 examines Belief-FB, and Section 3.3 details Rotation-FB's theoretical motivation and implementation.
- We empirically demonstrate that both BFB and RFB can adapt to different dynamics, unlike its counterparts in the zero-shot setup. Refer to Section 4.3 for the discussion and Figure 1 for results.

2 Behavioral Foundation Models

A Behavioral Foundation Model (BFM) (Pirotta et al., 2023; Tirinzoni et al.; Frans et al., 2024; Sikchi et al., 2025) is an RL agent trained in an unsupervised manner on a task-agnostic dataset to approximate optimal policies for various reward functions (tasks) specified at inference (test-time).

Forward-Backward Representation (FB) (Touati & Ollivier, 2021) approximates a discounted successor measure (Blier et al., 2021; Janner et al., 2020) for various behaviors across diverse tasks. The successor measure $M^{\pi}(s_0, a_0, X)$ for subset $X \subset \mathcal{S}$ is defined as cumulative discounted time spend at X starting at (s_0, a_0) and following π thereafter. More formally, for tabular example:

$$M^{\pi}(s_0, a_0, X) = \sum_{t \ge 0} \gamma^t P((s_{t+1}, a_{t+1}) \in X | s_0, a_0, \pi)$$
 (1)

with the corresponding Q-function for a specific task r:

$$Q_r^{\pi}(s_0, a_0) = \sum_{s^+ \in X} r(s^+) M^{\pi}(s_0, a_0, s^+).$$
 (2)

In continuous case, the FB representation aims to approximate successor measure through finite-rank approximation under diverse policies through forward $F: \mathcal{S} \times \mathcal{A} \times \mathcal{Z} \to \mathbb{R}^d$ and backward $B: \mathcal{S} \to \mathbb{R}^d$ functions. Given a set of policies π_z parametrized by task variable drawn uniformly from sphere $z_{\text{FB}} \in \text{Unif}(\mathcal{Z} = \mathbb{S}^d)$. Assuming ρ is a probability distribution over states within the offline dataset, the objective for FB is written as $M^{\pi_z}(s_0, a_0, X) \approx \int_{s^+ \in X} F(s_0, a_0, z)^T B(s^+) \rho(ds)$. Then, policy can be extracted as:

$$\pi_z(s) \approx \arg\max_a F(s, a, z)^T z.$$
 (3)

For continuous case, the greedy policy is approximated via DDPG (Lillicrap et al., 2015). During test time the task policy parametrization is approximated as $z_{test} \approx \mathbb{E}_{(s,a) \sim \rho}[r_{test}(s,a)B(s,a)]$. If the inferred task vector z_{test} lies within the task sampling distribution (in a linear span) of \mathcal{Z} used during training, then the optimal policy for task r_{test} is obtained from Equation 2 as $\pi_z(s) \approx \arg\max_a Q^{\pi_z}_{r_{test}}(s,a)$. For more details on training and inference procedures of FB, we refer reader to Appendix B.3. Extended discussion on other related works is included in the Appendix A

3 Method

Problem Statement. Our goal is to pre-train an agent in unsupervised regime on C_{train} contexts, which define particular MDP (CMDP), so that it is able to generalize to unseen ones during test time, *i.e.*, zero-shot¹. We collect diverse dataset, consisting of mix of highly exploratory unknown policies from varying environment layouts, differing either in dynamics (e.g., wind, friction, etc.) or environment specifications (e.g., positions of obstacles and doors). At test time, the agent is provided with small (up to episode termination steps) reward-free transitions from test context from explorative policy. Provided information must be used by an agent to recalibrate successor measure estimation corresponding to encountered environment. In an ideal scenario, the agent minimizes regret with an optimal policy for this environment dynamics.

maximizes the expected discounted return across both train and test contexts. We refer to Appendix A for details.

To formally study optimality guarantees of the problem above, we employ the following assumption commonly used for dynamics generalization (Eysenbach et al., 2021; Jeen & Cullen, 2024):

Assumption 1 (Coverage). Let $\mathcal{P}^c(s_{t+1}|s_t, a_t)$ be a transition probability for current context c given small dataset of reward-free random interactions either from test or train context. Then, $\mathcal{P}^{c_{\text{test}}}(s_{t+1}|s_t, a_t) > 0 \Rightarrow \mathcal{P}^{c_{\text{train}}}(s_{t+1}|s_t, a_t) > 0 \ \forall s_t, s_{t+1} \in \mathcal{S}, a_t \in \mathcal{A}.$

3.1 INVESTIGATING LATENT DIRECTIONS SPACE UNDER MULTIPLE DYNAMICS

We begin by addressing the following question: Why does FB representations fail to generalize effectively (both for train and test) to different situations under dynamics variations, *i.e.*, if learned on data sampled from diverse CMDPs? While the answer may appear intuitive, a closer look into the geometric structure of learned latent directions $z_{\text{FB}} \in \mathcal{Z}$, which encode possible policies π_z reveals critical insights which will be helpful later. We approach this question both theoretically and empirically on custom didactic discrete partially-observable Randomized Doors (see Appendix C.1) environment. Partial observability adds additional challenges and showcases the need to estimate belief state, which we discuss in the following sections.

¹We use the term "zero-shot RL" following Touati & Ollivier (2021).

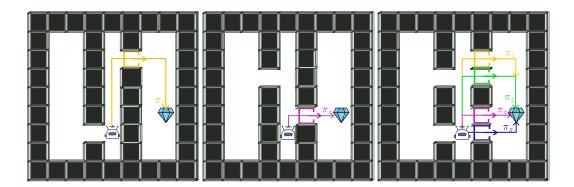


Figure 2: Randomized-Doors environment for three different layouts, each produced through varying the grid structure (exact randomization procedure is a hidden variable) (left-middle) From state s, the goal of an agent is to capture a diamond at target location by picking up the most suitable policy π_z (yellow for the first type and purple for the second) to move to the closest open door based on internal representation. (middle) When there are multiple possible future outcomes in the training data from the same state, the π_z 's (different colors) interfere with each other, leading to picking up an averaged policy.

In this experiment the only source of dynamics variation is the grid layout type. Namely, the positions of doors and walls are changed each new episode, depending on hidden configuration variable c. We collect a dataset of random trajectories drawn from multiple layouts, yielding near-uniform coverage of the entire (x,y) states. Now, consider a particular state s that an agent finds itself in three different layouts (see Figure 2). During FB training, we evaluate the forward representation $F(s,\cdot,z_{\rm FB})$ for latent directions (policy representations) $z_{\rm FB} \sim {\rm Uniform}(\mathbb{S}^{d-1})$, where each $z_{\rm FB}$ indexes a distinct policy starting at s.

In this setting a single grid state can require different optimal policies, depending on the layout an agent is instantiated in. Because $z_{\rm FB}$ does not enforce a separation of layout-specific futures, the FB model suffers from *interference*: latent directions encoding conflicting future outcomes overlap and become entangled in the policy representation space \mathcal{Z} . For each of the layout configuration and fixed state s from above, Figure 3 depicts latent directions $z_{\rm FB}$, colored by optimal policy as $a_{\rm color} = \arg\max_a F(s,a,z_{\rm FB})^T z_{\rm FB}$. When FB is trained on first two layouts in isolation, a unique dominant behavior (colored) emerges in \mathcal{Z} , recovering the optimal goal-reaching policy π_z^* . In contrast, training on data which mixes transitions from various environment instances results in $z_{\rm FB}$ to blend dynamics-specific information and instead to average over the possible futures, yielding a policy that is sub-optimal for every layout even from train set. Those observations are supported theoretically below.

Let $\{M^{\pi_i}\}_{i=1}^k$ be a collection of successor measures corresponding to optimal policies $\{\pi_i^*\}_{i=1}^k$ for distinct CMDPs defined by hidden context configurations $c_i \in C$. Assume that ρ is the state-action distribution supported on the offline dataset used for FB training and $M^{\pi_i^*}(s,a,\cdot) \approx F(s,a,z_i)^T B(\cdot)$ is approximated via rank d factors. Define the worst-case approximation error ϵ_k over context-dependent k successor measures as follows:

$$\epsilon_k := \inf_{F,B} \max_{1 \le i \le k} ||M^{\pi_i^*} - F(\cdot, \cdot, z_i)^T B(\cdot)||_{L^2(\rho)}.$$
(4)

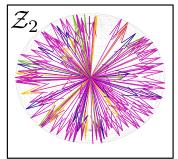
Then, the extracted policy π_{z_i} for (s, a) satisfies:

Theorem 1 (Regret-bound for Multiple Dynamics). For any bounded reward $||r||_{\infty} \leq R$ and particular test-time CMDP,

$$\mathbb{E}_{(s,a)\sim\rho_{test}}[Q_r^{\pi^*}(s,a) - Q_r^{\pi_{z_i}}(s,a)] \le \frac{2\gamma\epsilon_k||r||_{\infty}}{(1-\gamma)^2}.$$
 (5)

Because $\epsilon_{k+1} \ge \epsilon_k$ (monotonicity), the worst case regret per any CMDP at test time increases as more environments are included during training.

We provide a proof in Appendix B. Intuitively, Theorem 1 tells that adding transitions from more CMDPs only increases the worst-case optimality gap: as number of environments k grows, **FB** is forced to average over incompatible future dynamics. The proof relies on monotonicty property of



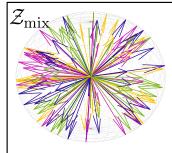


Figure 3: Different learned policy encoding π_z projections for three environment configurations from Figure 2 are visualized (yellow, purple and mixed trajectories). For a fixed state s and same goal across configurations, arrows depict latent directions $z_{\rm FB} \in \mathcal{Z}$ and colored by optimal behavior as $a_{color} = \arg\max_a F(s, a, z_{\rm FB})^T z_{\rm FB}$. (left-middle) When FB is trained on the two distinct configurations in separation, most of the latent directions agree on the optimal policy π_z . (right) When FB is trained on mix of CMDPs and at test time tasked with any particular configuration from train, obtained policy is ambiguous, since most policy-encoding directions do not agree on the action.

error term in Equation 4 and Theorem 9 from Touati & Ollivier (2021). In Section 3.3 we will refine this result and show that it is possible to remove explicit dependence of k, lowering the upper bound.

This interference highlights a fundamental trade-off. FB is expressive enough to model any task, yet when it is trained across environments with distinct unobserved parameters, the lack of contextual conditioning forces it to average different dynamics rather than separate them. The resulting successor measure merges transitions from distinct layouts and entangles directions in the latent space \mathcal{Z} . To disentangle these directions, we must represent uncertainty about the hidden context explicitly. The next section introduces a belief-conditioned objective that infers the latent context and allows FB to maintain environment-specific successor measures.

Takeaway 1

Because FB training inherently averages over all possible future states, it cannot learn a disentangled policy space and, therefore, fails to adapt to changes in dynamics.

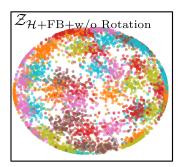
3.2 Belief State Modeling

To resolve the interference issue described in Section 3.1, we infer the latent context of an environment and augment FB input on that belief. We train a transformer encoder $f_{\rm dyn}$, by passing to a set of transitions $\{(s_t, a_t, s'_{t+1})\}_{t=1}^N$ and outputting an $h \in \mathbb{R}^d$. We denote the space of all possible inferred contexts as \mathcal{H} , where each element h encodes dynamics for particular environment. Because the ordering is discarded and no rewards in transitions are provided, the encoder must focus on dynamics specific mismatches (e.g., layout geometry, friction or wind direction), rather than policy specifics. Such context encoder should be permutation invariant, since unobservable factors describing environment are independent of the order of transitions in an episode. This setting provides theoretical ground for zero-shot and few-shot learning Snell et al. (2017).

Concretely, dataset consists of episodes $(\{(s_t, a_t, s'_{t+1})_{c_i}\}_{t=1}^N$ coming from CMDPs with randomly instantiated hidden specification variable c_i (different dynamics). We train a transformer encoder on random episodes (without episodic labels c_i) of context length n to infer contextual (hidden) variable h which fully specifies the dynamics across given episode. The transformer encoder loss involves two main components: 1) h is encouraged to follow a Gaussian prior and is shared across trajectory, and 2) projection head, which combines h with (s_t, a_t) to predict s_{t+1} . Those stages can be either trained end-to-end or separately. We observed that separating FB training from f_{dyn} gives better results.

For each trajectory we concatenate the inferred context vector h with the task vector z_{FB} to obtain augmented input $[h; z_{FB}]$ and condition only forward network as:

$$\hat{M}_{\pi_z}(s_t, a_t, s_{t+1}) = F(s_t, a_t, [h; z_{\text{FB}}])^T B(s_{t+1}).$$
(6)



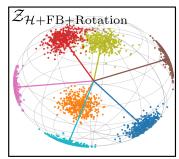


Figure 4: Visualization of inferred contexts h from space of all possible contexts \mathcal{H} (depicted as arrows) and task vectors z_{FB} (depicted as points on sphere boundary). Transitions from same CMDP colored the same. Concentration parameter κ defines spread of clusters. (left) Untrained transformer f_{dyn} output for different transitions is unstructured and same transitions coming from same CMDP (identical colors) are not collinear. (middle) New sampling procedure aligns policy specific vectors z_{FB} with context specific h, but clusters overlap before training. (right) After training, h for transitions from the same context are aligned and policies z_{FB} do not interfere between different environment configurations.

We empirically found that conditioning the backward network B degraded performance, producing smoothed out Q function, so B remains shared across contexts. Algorithm is summarized in Algorithm 1.

At test time, the agent is provided with a short (context length), reward-free trajectory and it is passed to f_{dvn} to obtain h. By plugging the result into Equation 3, the greedy policy is obtained.

Takeaway 2

We train a transformer in a self-supervised regime to estimate a belief over possible contexts, augmenting FB inputs and enabling effective disentanglement of contextual representations.

3.3 Structuring directions in the latent space

Insights from Section 3.1 showed that sampling task-vectors $z_{\rm FB}$ uniformly on the hypersphere encodes averaged policies, while Section 3.3 provided a solution through explicit context identification. We now combine these observations together through enhanced sampling $z_{\rm FB}$ around the inferred context h.

In Vanilla-FB, each state s draws $z_{\rm FB} \sim {\rm Unif}(\mathbb{S}^{d-1})$ with no inductive bias, so resulting policies π_z conflict with each other in CMDP setting, **even if additional explicit conditioning is introduced as before**. We replace uniform prior with a *von Mises-Fisher*(vMP) distribution centered at the context direction for episode $h = f_{\rm dyn}(\{(s_i, a_i, s_{i+1})\})$ as

$$z_{h+\text{FB}} \sim \text{vMF}(\mu = h, \kappa).$$
 (7)

with κ controlling the spread or *diversity* of policies (left and middle figures from Figure 4). In practice, to draw $z_{h+\text{FB}}$ we first pick a simple vector (e.g., the first basis vector), perturb with vMF noise, and finally rotate the result onto h with Householder reflection.

This enhancement has several benefits: 1) because directions h that differ in dynamics now occupy disjoint cones on the hypersphere, FB can fit the successor measure locally inside each cone, avoiding the destructive averaging effect quantified in Section 3.1 and 2) alignment procedure encourages the agent to explore policies that are plausible under its current belief while still injecting controlled diversity through κ .

Importantly, such a procedure not only has empirical benefits as we will show in Section 4, but also lowers bound from above in Theorem 1, making it non dependent on number of environments k.

Theorem 2 (Regret bound under latent space partitioning). Define ϵ_k as worst-case approximation error as in Equation 4. The Gram matrix of the task (policy) representations $\{z_{FB}\}_{i=1}^k$ is block diagonal w.r.t. partition $\{S_j\}$, with each S_j being the set of task-vector indices which satisfy

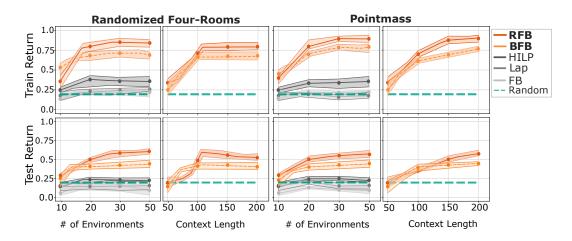


Figure 5: **Ablations on data diversity and context length of transformer encoder.** We show the influence of number of environments (data diversity) and context length on train and test performance in Four-Rooms and Pointmass environments. For data-diversity ablation, we see a clear performance boost up until some point, after which it platoes, as the Theorem 1 predicts. In our context-length ablation, we observe similar behaviour: performance improves as the context grows up to the length of a single episode, and then levels off. The results are averaged across three seeds, the opaque fill indicates standard deviation.

 $\langle z_{FB}, h^j \rangle \geq \cos \theta_{max}$ with θ_{max} being angle between any two latent vectors. Then,

$$\epsilon_k = \max_{j \le L} \epsilon_j, \quad \epsilon_k \le \epsilon_{k_{max}},$$
(8)

with $k_{max} := \max_{j} |S_j|$ being the size of the largest cone block.

Intuitively, Theorem 2 states that after the partitioning procedure of the latent space into non-overlapping clusters based on context representations h, the global worst-case FB approximation error $\epsilon_k = \max_{j \leq L} \epsilon_j$ is determined only by the cluster whose error ϵ_j is largest. Importantly, this bound does not depend on number of training environments k.

We provide a more formal treatment of this statement and a full proof in Appendix B.

Takeaway 3

Adjusting the prior over task vectors z_{FB} further mitigates the averaging effect and disentangles policy representations better based on the inferred dynamics.

4 EXPERIMENTS

In this section, we compare proposed methods, namely: **Belief-FB** (**BFB**) (Section 3.2) and its extension **Rotation-FB** (**RFB**) (Section 3.3), against the baselines in discrete and continuous settings. We outline experiments design below; all other necessary details are provided in Appendix D. Every environment is framed as a contextual MDP (CMDP), where the context differs by the underlying hidden variation (*e.g.*, grid layout, transition dynamics). During test time, we provide a single trajectory from random exploration policy, which enables context inference.

4.1 Environments and Setup

To support claims and theoretical insights made in previous sections, we consider the following experimental setups: (i) discrete, partially observable Randomized Four-Rooms (Appendix C.2), (ii) continuous AntWind (Appendix C.3), and lastly (iii) continuous partially observable Randomized-Pointmass (Appendix C.4). We vary the number of train layouts for each experiment, while fixing the number of held-out *unseen* context settings to 20 for Randomized Four-Rooms and Randomized-Pointmass, and 4 for Ant-Wind. We perform comparisons against following baselines:

HILP (Park et al., 2024) is a method that learns state representations from offline data so that the distance in the learned representation space is proportional to the number of steps between two

states in original space. **FB** (Touati & Ollivier, 2021) is an original version of the FB, described in Section 2. **Laplacian RL** (**LAP**) (Wu et al., 2019) constructs a graph Laplacian over state transitions from experience replay, then computes its eigenvectors to form low-dimensional representations that capture the environment's intrinsic structure. **Random** agent, which randomly explores the environment in a task-independent manner.

Randomized Four-Rooms is a discrete, deterministic, partially observable environment, where the task is to optimally move to the goal location. Training data is collected by executing random policies in N distinct grid layouts, that differ in doorway and wall locations.

Ant-Wind is a continuous environment, where the goal is to make an ant to walk forward as fast as possible. The environment dynamics are determined by the direction (angle) of a wind d.

Randomized-Pointmass is a partially observable continuous environment, where the task is to move to the goal locations. Maze grid structure is generated randomly, where each cell either contains wall or empty, while ensuring there is a path between start and goal locations.

4.2 COULD

BELIEF ESTIMATION ENABLE ADAPTATION IN FB?

Previously, we provided the theoretical foundations and speculated on the matter why FB is unable to differentiate between distinct dynamics and how we can use the belief estimation to overcome this. We refer to Table 1 and Figure 1 that show our empirical findings to support our claims.

We would like to point out that neither FB nor LAP are able to outperform a simple random baseline in PointMass and Four-Room, indicating that the policy they learn is most likely stuck in some obstacle due to averaging (see Section 3.1. Only HILP, which uses a different way to learn policy representations, is able to perform better than random policy.

Belief-FB and Rotation-FB outperform every baseline method, indicating that belief estimation is indeed a missing piece for adaptation. Notably, our methods also demonstrate generalization capabilities beyond train data on unseen test tasks.

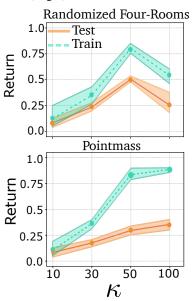


Figure 6: Influence of κ in RFB on performance. The results are averaged across three seed, the opaque fill represents standard deviation.

4.3 DO BFB AND RFB CAPTURE HIDDEN PROPERTIES OF THE ENVIRONMENT?

For an agent to refine its policy, it needs to keep track and update the uncertainty over possible environment configurations. Both Belief-FB and Rotation-FB accomplish this. Figure 7 illustrates this phenomenon visually. In Randomized-Door (left), the episodic trajectories from five layouts form non-overlapping clusters in the first two principal components of h, effectively disentangling different dynamics.

In Ant-Wind, the embeddings lie almost perfectly on a circle whose azimuth matches the underlying wind direction, generalizing smoothly to the 4 held-out wind angles. The quantitative results for evaluation in Table 1 (averaged across all environments) reveal that the baseline methods fail to recover those environment-specific properties and therefore produce sub-optimal policies even for train cases. In particular, HILP tends to predict an average direction in Randomized Four-rooms and ignores obstacles, while FB outputs same policy and Q function for almost all environments. Figure 12 shows that Q function is properly estimated only for BFB and RFB, respecting wall positions.

4.4 Does change in context length input to the f_{DYN} impacts performance?

In this experiment, we examine whether increasing the input trajectory length of improves performance. We vary the context length of $f_{\rm dyn}$ from 50 to 200 and present the results in Figure 5 for both Randomized Four-Rooms and Randomized Pointmass environments, across train and test configurations. The results show that performance is poor when the context length is shorter than a single trajectory episode (100 steps), as short trajectories only capture local, near-term goals. Conversely,

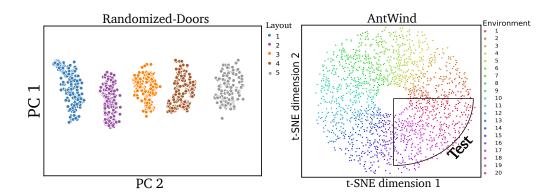


Figure 7: 2D projections of z_{dyn} inferred from different trajectories across number of different contexts (colors), showing effective disentangling environments based on transition function or other mismatches. (*left*) First two principal components are visualized for estimated z_{dyn} from five trajectories, each representing different layout type in Randomized-Doors. (*right*) Inferred context variables z_{dyn} recover hidden wind direction parameter in AntWind environment both for train and test, proving successful extrapolation properties.

excessively long sequences provide no additional benefit due to redundancy, since $f_{\rm dyn}$ already contains all neccessary information. Evaluations on both train and test environments demonstrate that $f_{\rm dyn}$ produces representations h capable of distinguishing between different context instances while maintaining robustness.

4.5 Does increase in dataset diversity make policies more robust?

We investigate if diversifying CMDP training configurations improves performance. Intuitively, broader state-action space coverage enhances successor measure estimation. Experiments confirm this: Figure 5 shows rapid improvement for BFB up to 25 configurations, while baselines match random policy performance. Once learned representations h from $f_{\rm dyn}$ cover all variation modes (contexts), additional data yields minimal gain (< 3%). These results align with Theorem 1.

4.6 How κ in RFB influences performance?

As described in Section 3.3, RFB concentration κ regularizes the diversity of policies for each environment. One the one hand, concentration should be high to ensure non-overlapping policy parametrized clusters π_z for different h, while at the same time it should not exceed certain value to control the diversity of policies in the environment, preventing collapsed solutions. Figure 6 shows that lower values of κ , meaning task-vectors z_{FB} are sampled with high deviation around h, likely producing overlapping clusters. As κ grows, task-vectors become more specialized, lowering variance which results in higher performance.

5 CONCLUSION & LIMITATIONS

We introduce **Belief-FB** (**BFB**) and **Rotation-FB** (**RFB**), two methods that extend the Forward-Backward representation to handle dynamics mismatches. We first identify a critical limitation in existing approaches: interference arises when naively sampling policy-encoding latent directions during training on transitions from conflicting dynamics. To address this, we learn hidden context variables (belief states) via a transformer encoder and use them for additional conditioning (Belief-FB). We improve latent-direction sampling by aligning task-relevant abstractions with environment-specific features, ensuring distinct regions in latent space of policies. Both BFB and RFB demonstrate theoretical and empirical improvements over prior methods. However, limitations include evaluations on a narrow set of dynamics mismatches and the introduction of the additional hyperparameter κ that controls policy diversity across environments. Also, random exploration at test time could fail at more complex environments and combining BFB and RFB together with more clever exploration methods at test time (Grillotti et al., 2024; Urpí et al., 2025) would make methods more scalable.

As future research directions, it would be valuable to investigate whether other zero-shot RL methods, those not based on successor-measure estimation, exhibit similar interference issues, and to scale our approach to more complex benchmarks such as XLand-MiniGrid (Nikulin et al., 2024; 2025) or Kinetix (Matthews et al., 2025).

REFERENCES

- Siddhant Agarwal, Harshit Sikchi, Peter Stone, and Amy Zhang. Proto successor measure: Representing the space of all possible solutions of reinforcement learning, 2025. URL https://openreview.net/forum?id=s9SVlWOcLt.
- Andre Barreto, Will Dabney, Remi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt, and David Silver. Successor features for transfer in reinforcement learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/350db081a661525235354dd3e19b8c05-Paper.pdf.
- Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, and Shimon Whiteson. A survey of meta-reinforcement learning, 2024. URL https://arxiv.org/abs/2301.08028.
- Steffen Bickel, Michael Brückner, and Tobias Scheffer. Discriminative learning for differing training and test distributions. In *Proceedings of the 24th international conference on Machine learning*, pp. 81–88, 2007.
- Léonard Blier, Corentin Tallec, and Yann Ollivier. Learning successor states and goal-dependent values: A mathematical viewpoint. *arXiv preprint arXiv:2101.07123*, 2021.
- Diana Borsa, Andre Barreto, John Quan, Daniel J. Mankowitz, Hado van Hasselt, Remi Munos, David Silver, and Tom Schaul. Universal successor features approximators. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=S1VWjiRcKX.
- Open X-Embodiment Collaboration. Open X-Embodiment: Robotic learning datasets and RT-X models. https://arxiv.org/abs/2310.08864, 2023.
- Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural computation*, 5(4):613–624, 1993.
- Ron Dorfman, Idan Shenfeld, and Aviv Tamar. Offline meta learning of exploration, 2021. URL https://arxiv.org/abs/2008.02598.
- Benjamin Eysenbach, Shreyas Chaudhari, Swapnil Asawa, Sergey Levine, and Ruslan Salakhutdinov. Off-dynamics reinforcement learning: Training for transfer with domain classifiers. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=eqBwg3AcIAK.
- Basura Fernando, Amaury Habrard, Marc Sebban, and Tinne Tuytelaars. Unsupervised visual domain adaptation using subspace alignment. In *Proceedings of the IEEE international conference on computer vision*, pp. 2960–2967, 2013.
- Kevin Frans, Seohong Park, Pieter Abbeel, and Sergey Levine. Unsupervised zero-shot reinforcement learning via functional reward encodings. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 13927–13942. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/frans24a.html.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv* preprint arXiv:2004.07219, 2020.
- Karol Gregor, Danilo Jimenez Rezende, Frederic Besse, Yan Wu, Hamza Merzic, and Aaron van den Oord. Shaping belief states with generative environment models for rl. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/2c048d74b3410237704eb7f93a10c9d7-Paper.pdf.

- Luca Grillotti, Maxence Faldor, Borja G León, and Antoine Cully. Quality-diversity actor-critic: learning high-performing and diverse behaviors via value and successor features critics. *arXiv* preprint arXiv:2403.09930, 2024.
 - Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. Pmlr, 2018.
 - Michael Janner, Igor Mordatch, and Sergey Levine. γ -models: Generative temporal difference learning for infinite-horizon prediction. In *Advances in Neural Information Processing Systems*, 2020.
 - Scott Jeen and Jonathan Cullen. Dynamics generalisation with behaviour foundation models. In Workshop on Training Agents with Foundation Models at RLC 2024, 2024. URL https://openreview.net/forum?id=Alu8YM7vuP.
 - Scott Jeen, Tom Bewley, and Jonathan Cullen. Zero-shot reinforcement learning from low quality data. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=79eWvkLjib.
 - Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of zero-shot generalisation in deep reinforcement learning. *Journal of Artificial Intelligence Research*, 76: 201–264, 2023.
 - Wouter M Kouw and Marco Loog. A review of domain adaptation without target labels. *IEEE transactions on pattern analysis and machine intelligence*, 43(3):766–785, 2019.
 - Michael Laskin, Luyu Wang, Junhyuk Oh, Emilio Parisotto, Stephen Spencer, Richie Steigerwald, DJ Strouse, Steven Hansen, Angelos Filos, Ethan Brooks, Maxime Gazeau, Himanshu Sahni, Satinder Singh, and Volodymyr Mnih. In-context reinforcement learning with algorithm distillation, 2022. URL https://arxiv.org/abs/2210.14215.
 - Jonathan N. Lee, Annie Xie, Aldo Pacchiano, Yash Chandak, Chelsea Finn, Ofir Nachum, and Emma Brunskill. Supervised pretraining can learn in-context reinforcement learning, 2023. URL https://arxiv.org/abs/2306.14892.
 - Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv* preprint arXiv:1509.02971, 2015.
 - Michael Matthews, Michael Beukman, Chris Lu, and Jakob Nicolaus Foerster. Kinetix: Investigating the training of general agents through open-ended physics-based control tasks. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=zCxGCdzreM.
 - Aditya Modi, Nan Jiang, Satinder Singh, and Ambuj Tewari. Markov decision processes with continuous side information. In *Algorithmic learning theory*, pp. 597–618. PMLR, 2018.
 - Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, Artem Agarkov, Viacheslav Sinii, and Sergey Kolesnikov. Xland-minigrid: Scalable meta-reinforcement learning environments in jax, 2024. URL https://arxiv.org/abs/2312.12044.
 - Alexander Nikulin, Ilya Zisman, Alexey Zemtsov, and Vladislav Kurenkov. Xland-100b: A large-scale multi-task dataset for in-context reinforcement learning, 2025. URL https://arxiv.org/abs/2406.08973.
 - Seohong Park, Tobias Kreiman, and Sergey Levine. Foundation policies with hilbert representations. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=LhNsSaAKub.
 - Matteo Pirotta, Andrea Tirinzoni, Ahmed Touati, Alessandro Lazaric, and Yann Ollivier. Fast imitation via behavior foundation models. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023. URL https://openreview.net/forum?id=SHNjk4h0jn.

- Andrey Polubarov, Nikita Lyubaykin, Alexander Derevyagin, Ilya Zisman, Denis Tarasov, Alexander Nikulin, and Vladislav Kurenkov. Vintix: Action model via in-context reinforcement learning, 2025. URL https://arxiv.org/abs/2501.19400.
 - Kate Rakelly, Aurick Zhou, Deirdre Quillen, Chelsea Finn, and Sergey Levine. Efficient off-policy meta-reinforcement learning via probabilistic context variables, 2019. URL https://arxiv.org/abs/1903.08254.
 - Harshit Sikchi, Siddhant Agarwal, Pranaya Jajoo, Samyak Parajuli, Caleb Chuck, Max Rudolph, Peter Stone, Amy Zhang, and Scott Niekum. Rl zero: Zero-shot language to behaviors without any supervision. *arXiv preprint arXiv:2412.05718*, 2024.
 - Harshit Sikchi, Andrea Tirinzoni, Ahmed Touati, Yingchen Xu, Anssi Kanervisto, Scott Niekum, Amy Zhang, Alessandro Lazaric, and Matteo Pirotta. Fast adaptation with behavioral foundation models. *arXiv preprint arXiv:2504.07896*, 2025.
 - Viacheslav Sinii, Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, and Sergey Kolesnikov. In-context reinforcement learning for variable action spaces, 2024. URL https://arxiv.org/abs/2312.13327.
 - Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning, 2017. URL https://openreview.net/forum?id=B1-Hhnslg.
 - Casper Kaae Sønderby, Jose Caballero, Lucas Theis, Wenzhe Shi, and Ferenc Huszár. Amortised map inference for image super-resolution. *arXiv preprint arXiv:1610.04490*, 2016.
 - Denis Tarasov, Alexander Nikulin, Ilya Zisman, Albina Klepach, Andrei Polubarov, Lyubaykin Nikita, Alexander Derevyagin, Igor Kiselev, and Vladislav Kurenkov. Yes, q-learning helps offline in-context RL. In *Scaling Self-Improving Foundation Models without Human Supervision*, 2025. URL https://openreview.net/forum?id=B86JMHZUnc.
 - Jayden Teoh, Pradeep Varakantham, and Peter Vamplew. On generalization across environments in multi-objective reinforcement learning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=tuEP424UQ5.
 - Andrea Tirinzoni, Ahmed Touati, Jesse Farebrother, Mateusz Guzek, Anssi Kanervisto, Yingchen Xu, Alessandro Lazaric, and Matteo Pirotta. Zero-shot whole-body humanoid control via behavioral foundation models.
 - Ahmed Touati and Yann Ollivier. Learning one representation to optimize all rewards. *Advances in Neural Information Processing Systems*, 34:13–23, 2021.
 - Ahmed Touati, Jérémy Rapin, and Yann Ollivier. Does zero-shot reinforcement learning exist? *arXiv* preprint arXiv:2209.14935, 2022.
 - Masatoshi Uehara, Issei Sato, Masahiro Suzuki, Kotaro Nakayama, and Yutaka Matsuo. Generative adversarial nets from a density ratio estimation perspective. *arXiv preprint arXiv:1610.02920*, 2016.
 - Núria Armengol Urpí, Marin Vlastelica, Georg Martius, and Stelian Coros. Epistemically-guided forward-backward exploration. *arXiv preprint arXiv:2507.05477*, 2025.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
 - Yifan Wu, George Tucker, and Ofir Nachum. The laplacian in RL: Learning representations with efficient approximations. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=HJlNpoA5YQ.

- Jinwei Xing, Takashi Nagata, Kexin Chen, Xinyun Zou, Emre Neftci, and Jeffrey L Krichmar. Domain adaptation in reinforcement learning via latent unified state representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 10452–10459, 2021.
- Amy Zhang, Rowan McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning invariant representations for reinforcement learning without reconstruction. *arXiv* preprint *arXiv*:2006.10742, 2020.
- Chuning Zhu, Xinqi Wang, Tyler Han, Simon Shaolei Du, and Abhishek Gupta. Distributional successor features enable zero-shot policy optimization. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=8IysmgZte4.
- Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning, 2020. URL https://arxiv.org/abs/1910.08348.
- Ilya Zisman, Vladislav Kurenkov, Alexander Nikulin, Viacheslav Sinii, and Sergey Kolesnikov. Emergence of in-context reinforcement learning from noise distillation. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=Y8KsHT1kTV.
- Ilya Zisman, Alexander Nikulin, Viacheslav Sinii, Denis Tarasov, Nikita Lyubaykin, Andrei Polubarov, Igor Kiselev, and Vladislav Kurenkov. N-gram induction heads for in-context rl: Improving stability and reducing data needs, 2025. URL https://arxiv.org/abs/2411.01958.

A EXTENDED RELATED WORKS AND BACKGROUND

A.1 BACKGROUND

Contexual Markov Decision Process. Throughout paper we will be dealing with a Contextual Markov Decision Process (CMDP), defined by a tuple $\langle \mathcal{C}, \mathcal{S}, \mathcal{A}, \gamma, \mathcal{M} \rangle$, where \mathcal{C} is a context space and \mathcal{S}, \mathcal{A} are shared state and action spaces across environments. Function \mathcal{M} maps particular context $c \in \mathcal{C}$ to respective MDP, i.e., $\mathcal{M}(c) = \langle \mathcal{S}, \mathcal{A}, \mathcal{T}^c, R^c, \mu^c, \gamma \rangle$ with context-dependent transition function $\mathcal{T}^c : \mathcal{S} \times \mathcal{A} \times \mathcal{C} \to \mathcal{S}$, μ^c being an initial distribution over states and $\gamma \in (0,1)$ a discount factor. Intuitively, the context $c \in \mathcal{C}$ represents a fixed environmental configuration, such as obstacle positions, layout geometry, dynamics vector parameters or seed. Throughout this work, the context remains static within each episode, consistent with prior literature (Modi et al., 2018; Kirk et al., 2023; Teoh et al., 2025). A policy $\pi : \mathcal{S} \to \Delta \mathcal{A}$ is optimal for context c for the reward function R if it maximizes expected discounted future reward, i.e., $\pi^*_{c,R}(s_0, a_0) = \arg \max_{\pi} \mathbb{E}[\sum_{i=1}^{n} \gamma^t R(s_t, a_t) | s_0, a_0, \pi, c]$.

When the context is fully observable, augmenting the state space with the given context reduces the CMDP to a standard MDP, eliminating the need to model distinct dynamics \mathcal{T}^c , rewards R^c or initial states μ^c . However, if the context is partially observable, the learned model must infer and track the uncertainty over true hidden configuration to maintain theoretical optimality guarantees. Such task can be framed as posterior estimation $p(c|\mathcal{H})$ or *belief* over possible contexts c given accumulated history H.

Most successful methods for deriving an optimal policy across arbitrary tasks from a task-agnostic dataset leverage successor features (Dayan, 1993; Barreto et al., 2017; Borsa et al., 2019; Park et al., 2024; Zhu et al., 2024) or their continuous counterpart, successor measures (Blier et al., 2021; Touati & Ollivier, 2021; Touati et al., 2022; Agarwal et al., 2025; Jeen et al., 2024). In this work, we focus on the latter framework, specifically its instantiation via forward-backward representations (Touati & Ollivier, 2021). Below, we briefly outline its key properties.

Zero-Shot RL. Given an offline dataset of transitions $\mathcal{D} = \{(s_i, a_i, s_{i+1})\}_{i=1}^{|\mathcal{D}|}$ generated by an unknown behavior policies, the agent's objective is to learn a compact abstraction of the environment from which it is possible. At test time, this abstraction helps to obtain optimal policy for *any* reward function r_{test} which defines a particular *task*. Reward function can be specified either as a small dataset of reward-labeled states $\mathcal{D}_{test} = \{(s_i, r_{test}(s_i))\}_{i=1}^k$ or as a direct mapping $s \to r_{test}(s)$. While some prior works assume access to the context labels (Gregor et al., 2019), we focus on the setting where the context is unknown and must be inferred from the data. Alternative formulations of zero-shot RL exist under other formalisms, and we refer to (Kirk et al., 2023) for comprehensive overview.

A.2 RELATED LITERATURE

Domain Adaptation and Transfer Learning in RL. While our work will focus on domain adaptation applied to estimating successor measure for various dynamics mismatches, we start by briefly reviewing more general ideas in classic domain adaptation and refer to (Kouw & Loog, 2019) for detailed overview. Most methods for domain adaptation can be categorized into *importance-weighting* (Bickel et al., 2007; Uehara et al., 2016; Sønderby et al., 2016) and *domain-invariant feature learning* (Fernando et al., 2013; Eysenbach et al., 2021; Xing et al., 2021; Zhang et al., 2020) approaches. Former methods estimate the likelihood ratio of examples under samples from target domain versus samples from source, which is then used to recalibrate examples from the source domain. The latter approaches learn a unified representation of the environment, targeting to extract only task-relevant abstraction, negating distracting information.

The most relevant approach which enables FB representations to generalize across dynamics is Contexual FB (Jeen & Cullen, 2024). This approach uses importance-weighting formalism and introduces two classifiers, which estimate the likelihood of transitions (s_t, a_t) and (s_t, a_t, s_{t+1}) being from train or test context and augment the reward function to account for those discrepancies in the dynamics. If augmented reward function lies in the linear span of the \mathbb{Z} space during FB training, then the policy can be extracted as described in Equation 3. However, such an approach requires training classifiers from scratch for each novel layout of the environment, limiting its applicability.

Meta-RL. Another major line of related works, Meta-Reinforcement Learning (Meta-RL), focuses on few-shot domain adaptation to unseen tasks or dynamics (Beck et al., 2024). The significant part of research in Meta-RL is dedicated to explicitly learning the *belief* by collecting a history of interactions with the environment on inference during test-time (Zintgraf et al., 2020; Dorfman et al., 2021; Rakelly et al., 2019). However, recent works show that it is possible to quantify the *belief* without learning the posterior implicitly (Laskin et al., 2022; Lee et al., 2023; Zisman et al., 2024; Sinii et al., 2024; Zisman et al., 2025; Tarasov et al., 2025; Polubarov et al., 2025). Leveraging in-context ability of transformers Vaswani et al. (2017), one can learn an end-to-end supervised model, while the transformer's context will absorb into robust representation the adaptation-relevant information thus enabling fast adaptation. We also leverage this in-context ability to construct the belief representation of the dynamics the agent currently in, but instead operating in a zero-shot manner.

B Proofs

B.1 THEOREM 1

Preserving notation from Section 3.1, we provide a full proof of the Theorem 1. Let $\{M_{\pi_i}\}$ be a collection of successor measure of the optimal policies $\{\pi_i\}_{i=1}^k$ for k distinct CMDPs. Given a reference measure ρ on $\mathcal{S} \times \mathcal{A}$ let worst case regret be defined as

$$\epsilon_k := \inf_{F,B} \max_{i < i < k} ||M_{\pi_i} - F(\cdot, \cdot, z_i)^T B(\cdot)||_{L^2_{\rho}}$$

$$\tag{9}$$

Theorem (Regret-bound for Multiple Dynamics). Then, for any bounded $||r_{\infty}|| \le R$ and any CMDP whose state-action distribution ρ_{test} (assuming absolute continuity, i.e., $d\rho_{test}/\rho$ is bounded), the policy extracted from F, B for that CMDP satisfies:

$$\mathbb{E}_{(s,a) \sim \rho_{\textit{lest}}}[Q^{\pi^*}(s,a) - Q^{\pi_{z_i}}(s,a)] \le \frac{2\gamma \epsilon_k ||r||_{\infty}}{(1-\gamma)^2}$$

Since $\epsilon_{k+1} \ge \epsilon_k$ (monotonicity) the worst case regret per any CMDP at test time increases as more environments are included during training.

Lemma 1. Theorems 8-9 from Touati & Ollivier (2021) prove this inequality for single instance of MDP, showing that if FB approximation error in $L^2(\rho)$ is at most ϵ then pointwise value gap is bounded by:

$$(Q_r^* - Q_r^{\pi_{z_i}}) \le \frac{\gamma}{1 - \gamma} (P_{\pi^*} - P_{\pi_z}) (I - \gamma P_{\pi^*})^{-1} E(z) r) \tag{10}$$

with E(z) being a point-wise error matrix over state-actions as $E(z) = M^{\pi_z}(s,a,s') - F(s,a,z)^T B(s,a)$. Since

$$||(I - \gamma P)^{-1}||_{\infty} \le \frac{1}{1 - \gamma}$$
 (11)

results in coefficient $2\gamma/(1-\gamma)^2$ in Equation 1.

Proof. Define a transition kernel P_i of CMDP at index i and M_{π_i} its successor measure. Let $E_i = M_{\pi_i} - F(s, a, z_i)^T B(\cdot) = M_{\pi_i} - \hat{M}_i$. Then, using $Q^* = (I - \gamma P_{\pi^*})^{-1} r$ value gap decomposes as

$$Q^* - Q^{\pi_{z_i}} = \gamma (I - \gamma P_{\pi^*})^{-1} (P_{\pi_*} - P_{\pi_{z_i}}) (I - \gamma P_{\pi_{z_i}})^{-1} r$$
(12)

Since each of the resolvent factors (denote them as E_i) are at most $1/(1-\gamma)$ in L^{∞} , then from triangle inequality:

$$||Q^* - Q^{\pi_{z_i}}||_{\infty} \le \frac{2\gamma}{(1-\gamma)^2} ||E_i||_{L^2_{\rho}} ||r||_{\infty}$$
(13)

From Assumption 1 on absolute continuity,

$$\mathbb{E}_{(s,a)\sim\rho_{\text{test}}}\{Q^* - Q^{\pi_{z_i}}\} \le ||Q^* - Q^{\pi_{z_i}}||_{\infty}$$
(14)

Substituting this into Equation 13, gives desired inequality bound in Theorem 1. \Box

B.2 THEOREM 2

810

811

812

813

814 815 816

817

818

819

820

821

822

823

824

825

826 827

828

829

830

831 832

833

834 835

836

837

838

839

840 841

842

843

844

845 846

847

848

849

850 851

852

853

854

855 856

857

858

859

860

861

862

863

Section 3.3 introduced a new sampling procedure of $z_{\rm FB}$, which improves upon usual uniform sampling. This procedure can also be studied more formally.

Given an L possible contexual representations h of the environments coming from f_{dyn} , define a *cone* around each of the context axes $\{h^1, h^2 \dots h^L\} \in \mathbb{S}^{d-1}$, with the angle between any two latent vectors $\theta_{\rm max}$ set

$$C_j = \{ z_{FB} \in \mathbb{S}^{d-1} | \langle z_{FB}, h^j \rangle \ge \cos \theta_{\text{max}} \}$$
 (15)

Corresponding policy task vectors are defined for each cone $z_{\rm FB}^i \in {\cal C}_{c(i)}$, with $c(i) \in \{1, \dots L\}$ being a classification function, mapping index i to one of the predifined context axes. For functions F, B define per environment error as:

$$\mathcal{E}_i(F,B) := ||M^{\pi_i} - F(\cdot, \cdot, z_{FB}^i)^T B(\cdot)||_{L^2(\rho)}$$
(16)

With following optimization tasks:

$$\epsilon_k := \inf_{F,B} \max_{1 \le i \le k} \mathcal{E}_i(F,B), \quad \epsilon_j := \inf_{F,B} \max_{i \in \mathcal{S}_j} \mathcal{E}_i(F,B)$$
(17)

with $S_i = \{i | c(i) = j\}$ being a set of task vectors $(z_{\rm FB})$ indices that fall into the j-th cone of the latent space partition.

Theorem (Regret-bound under latent space partitioning). *Under assumptions above, the Gram matrix* of the directions $\{z_{FB}\}_{i=1}^k$ is block diagonal w.r.t. partition $\{S_j\}$ and

$$\epsilon_k = \max_{j \le L} \epsilon_j, \quad \epsilon_k \le \epsilon_{k_{max}}$$
 (18)

with $k_{max} := \max_{i} |S_i|$ being the size of a largest cone block.

In order to prove this theorem, assume that collection of contexual embeddings $\{h_i\}_{i=1}^L$ obtained from L environments are almost orthogonal.

Proof. Define a $k \times k$ Gram matrix as $G = \langle z_{FB}^i, z_{FB}^j \rangle$ with i, j corresponding to cone partition. Because cones, corresponding to different contexual embeddings h, are disjoint and lie in a span $\{h_i\}$, the resulting Gram matrix is block diagonal $G = \operatorname{diag}(G^{(1)}, G^{(2)}, ..., G^{(L)})$. For a fixed rank d of F, B, the worst case approximation error is

$$\epsilon_k(F, B) = \max_{1 \le i \le k} ||M_{\pi_i} - \hat{M}_{\pi_i}||_{L^2(\rho)} = \max_{j \le L} \max_{i \in S_j} ||M_{\pi_i} - \hat{M}_{\pi_i}||_{L^2(\rho)}$$
(19)

Since matrix G is block-diagonal, optimization of F, B decouples over blocks of G. Namely, minimizer on the full set is obtained by minimizing each block separately, hence:

$$\epsilon_k = \inf_{F,B} \epsilon_k(F,B) = \max_{j \le L} \epsilon_j \tag{20}$$
 By taking $k_{\max} = \max_j |S_j|$ and $\epsilon_k \le \epsilon_{k_{\max}}$ for each block, we obtain desired inequality. \square

Notably, such orthogonal cone partitioning eliminates interference. Once each cone has its own slice of the latent space, adding more cones does not enlarge the worst-case error bound, and with representation capacity of F and B being $d \geq k_{\text{max}}$ the FB model can reach zero approximation error in principle.

B.3 FB TRAINING

In this section we describe the training procedure of FB in more details. Everything follows the notation from Touati & Ollivier (2021).

Assume that ρ is supported over all provided data, i.e., it is non-zero everywhere.

$$\mathcal{L}_{FB} = \mathbb{E}_{(s_t, a_t, s_{t+1}, s_+) \sim \mathcal{D}, z \sim \mathcal{Z}} [(F(s_t, a_t, z)^T B(s_+) - \gamma \hat{F}(s_{t+1}, \pi_z(s_{t+1}, z)^T \hat{B}(s_+))^2 - 2F(s_t, a_t, z)^T B(s_{t+1})]$$
(21)

Here, s_+ is a future outcome either from the same trajectory or randomly sampled from data. \hat{F}, \hat{B} are target networks with Z being a task space, encoding all possible policies. The policy π_z is trained in an actor-critic formulation and parametrized as Boltzmann policy $\pi_{z_i}(\cdot|s_i) =$ softmax $(F(s_i, \cdot, z_i)^T z_i/\tau)$ for continuous environments. Additionally, B is forced to be orthogonal for different s, which is enforced by contrastive loss $\mathbb{E}_{(s,s+)\sim\rho}[B(s)^TB(s_+)]$, with ρ being a measure over dataset.

C ENVIRONMENT DESCRIPTIONS

C.1 RANDOMIZED-DOORS

The Randomized-Doors MiniGrid environment (Figure 8) is a discrete-state, discrete-action finite horizon deterministic environment in which agent has an objective to go to goal location with maximum return of 1. Each episode terminates after 100 steps or after reaching goal location. The randomization determines possible open doors locations, fully specifying particular layout. In our experiments, the observation state of an agent consists of (x,y) coordinates tuple, making it partially observable. Such setting requires to properly update beliefs over unobservable layout configuration type. The action space consists of four actions, namely $\{up, down, right, left\}$, while (x,y) coordinates across both axes are bounded by grid size, which we take to be 9×9 .

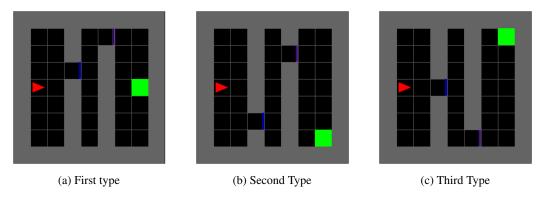


Figure 8: Several possible layouts are visualized, each corresponding to unique possible doors configurations. The agent is denoted as a red triangle. The task specification (goal position) with reward of 1 is denoted by green square and is also randomized. It is a custom implementation based on Empty MiniGrid (https://minigrid.farama.org).

C.2 RANDOMIZED FOUR-ROOMS

The Randomized Four-Rooms MiniGrid environment Figure 9 is a modification of classic Four-Rooms and is a discrete-state, discrete-action, deterministic partially observable environment. For each episode, the maze layout (grid type) is generated randomly, ensuring all of the four rooms are connected with exactly single door. Observation state consists of (x,y) coordinates, making this environment hard and checks whether agent could successfully estimate uncertainty over hidden configurations solely based on number of occurrence of each transition, recovering dynamics. In our experiments, we consider 11×11 bounds for height and width.

Observation space consists of raw discrete (x,y) coordinates on the grid, while actions correspond to a set of possible moves $\{up, down, left, right\}$. For every layout we record 500 episodes of length 100, yielding a dataset that covers almost all possible (s,a) transitions. For testing on unseen configurations, we fix agent starting position to coordinates of the first empty cell and evaluate performance across 3 static goal positions, farhest away from starting position.

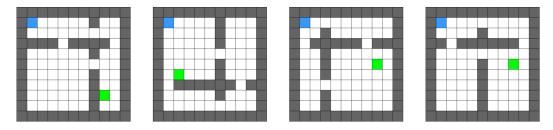


Figure 9: **Different layout configurations from randomized Four-Rooms environment.** During inference, the goal for the agent (depicted in blue) is to achieve green location. In our experiments we fix starting agent position and fix 3 goals, one for each room.

C.3 ANT-WIND

The AntWind environment is a modified version of the Ant locomotion task from the MuJoCo simulator, commonly used to test an agent's adaptability to changing dynamics. In this environment, an ant-like robot must learn to move forward while being subjected to external wind forces varying in magnitude and direction. In our experiments we consider 17 environments for training, covering three quadrants of possible wind directions on the circle, while leaving others for test, checking extrapolation on the fourth quadrant.

For our experiment, we collect dataset by training SAC (Haarnoja et al., 2018) on 3/4 of all possible directions, which results in 16 environments and hold out the other 1/4 for evaluation. Resulting dataset consists of 3400 transition tuples, where each environment configuration is represented as trajectory of length 256.

C.4 RANDOMIZED POINTMASS

Randomized Pointmass is a modification of pointmass environment from D4RL Fu et al. (2020). Each episode the environment grid structure is randomized, ensuring all cells are interconnected. The observation space consists of (x,y) transitions. Start position is determined as a first empty cell, while goal location is chosen to be the fartherst away from start (based on Manhattan distance) and ensuring existence of at least one valid trajectory (e.g., through BFS).

Observation space consists of (global x, global y) position, similar to Four-Rooms. We fix dataset size to be $1e^6$, only varying number of layouts and episodes per layout, while fixing episode length to 250. We use explore policy, which is a random policy with a portion of actions repeated ("sticky-actions").

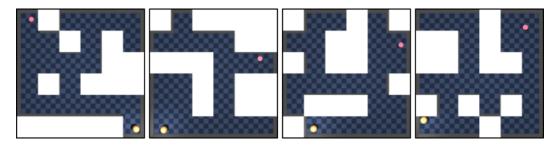


Figure 10: Examples of pointmass grid variations.

D EXPERIMENTS DETAILS

Randomized-Doors. For didactic example from Section 3.1 we collect diverse dataset from different layout configurations (open door locations) such that visitation distribution over all states is non-zero. Black color denotes obstacles. The episode length is set to be 100, which is equal to the context length of the transformer encoder for this experiment. Overall, we collect 500 episodes per layout and coverage heatmap is visualized in Figure 11.

Table 1: Comparison of proposed approaches against baselines on **test** (unseen) environments. Results for Fourrooms and Pointmass are averaged across 20 mazes configurations.

Environment (Test)			Method			
Environment (165t)	Random	Vanilla-FB	HILP	Lap	Belief-FB	Rotation-FB
Randomized-Fourrooms Randomized-Pointmass Ant-Wind	$\begin{array}{c} 0.05 \pm 0.01 \\ 0.03 \pm 0.01 \\ 250 \pm 200.0 \end{array}$	$\begin{array}{c} 0.15 \pm 0.06 \\ 0.1 \pm 0.1 \\ 250 \pm 98.5 \end{array}$	$\begin{array}{c} 0.2 \pm 0.02 \\ 0.25 \pm 0.02 \\ 410 \pm 40.5 \end{array}$	$\begin{array}{c} 0.1 \pm 0.1 \\ 0.1 \pm 0.1 \\ 290 \pm 22.5 \end{array}$	$\begin{array}{c} 0.4 \pm \! 0.02 \\ 0.45 \pm \! 0.05 \\ 550 \pm \! 50.5 \end{array}$	$\begin{array}{c} 0.61 \pm 0.02 \\ 0.55 \pm 0.05 \\ 640 \pm 30.7 \end{array}$

Note on relience on random exploration during test time. Random exploration relience of BFB and RFB in highly complex environments may fail to discover crucial states needed to disambiguate dynamics identification. However, we emphasize that our work addresses a distinct bottleneck: existing behavioral foundation models (BFMs), particularly FB, tend to collapse when trained

Table 2: Comparison of proposed approaches against baselines on **train** environments. Results for Fourrooms and Pointmass are averaged across 20 mazes configurations.

Environment (Train)	Method					
	Random	Vanilla-FB	HILP	Lap	Belief-FB	Rotation-FB
Randomized-Fourrooms Randomized-Pointmass Ant-Wind	$\begin{array}{c} 0.18 \pm 0.02 \\ 0.0 \pm 0.05 \\ -190 \pm 250 \end{array}$	$\begin{array}{c} 0.25 \pm 0.02 \\ 0.2 \pm 0.2 \\ 390 \pm 120 \end{array}$	$\begin{array}{c} 0.4 \pm 0.02 \\ 0.45 \pm 0.1 \\ 410 \pm 90 \end{array}$	$\begin{array}{c} 0.2 \pm 0.1 \\ 0.15 \pm 0.15 \\ 340 \pm 150 \end{array}$	$\begin{array}{c} 0.7 \pm \! 0.02 \\ 0.76 \pm \! 0.18 \\ 680 \pm \! 80 \end{array}$	$\begin{array}{c} 0.85 \pm \! 0.02 \\ 0.88 \pm \! 0.2 \\ 740 \pm \! 70 \end{array}$

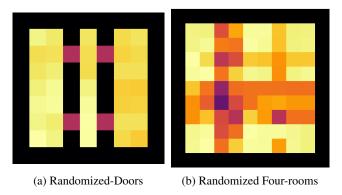


Figure 11: Empirical state occupancy measures (ρ) visualizations of collected datasets for discrete-based environments.

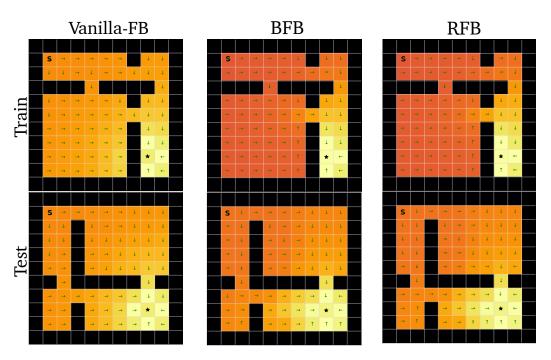


Figure 12: **Q-function and deterministic policy visualizations** (Equation 3) on Randomized Four-Rooms environment. Vanilla-FB ignores environment structure and resulting policy moves through obstacles. BFB and RFB do not have such issue.

on offline data composed of mixed CMDPs. Consequently, training BFMs on large scale mixed multi-modal (in terms of dynamics) data would yield an averaged policy, thus limiting their current applicability to unimodal datasets (in terms of dynamics mismatch). Both BFB and RFB overcome this collapse. Developing smarter test-time exploration strategies to streamline dynamics identification remains an important direction for future research.

D.1 DATASET GENERATION

For Randomized Four-Rooms, we produce four training datasets with the following parameters:

# Transitions	# layouts	# episodes per layout	episode length
1000000	10	1000	100
1000000	20	500	100
1000000	30	250	100
1000000	50	150	100

Table 3: Details for Randomized Four-Rooms datasets

Randomized Four-Rooms. For experiments on Randomized Four-Rooms during dataset collection we generate randomly grid layout, ensuring that each room is interconnected by exactly one door. For evalution we fix agent start position to (1,1) with the goal of reaching 3 other goals, specified at other rooms. Each episode terminates after 100 steps. The evaluation protocol is averaged success rate across 3 across 20 environments.

AntWind. For AntWind we first collect trajectories by varying wind direction d and training an expert-like SAC agent. After training, we collected evaluation trajectories from trained agent. This ensures that all directions are covered and explicitly sets dynamics context. As said in Experiments section, we train on 16 environments with wind directions corresponding to first 3 quadrants of circle, leaving other 4 (last quadrant) for hold out.

E IMPLEMENTATION DETAILS

E.1 FORWARD-BACKWARD REPRESENTATIONS

E.1.1 GPUs

We run each experiment on 1 Nvidia RTX 4090. The overall training time (for both dynamics encoder and FB training) is approximately 1 hour.

E.1.2 ARCHITECTURE

The forward-backward architecture described below mostly follows the implementation by Touati et al. (2022). All other additional hyperparameters for BFB and RFB are reported in Table 4. Moreover, we should emphasize that our choice of transformer architecture for $f_{\rm dyn}$ is mainly based on its abilities to encode large sequences, and other architectural designs (e.g State-Space Models, RNNs) can also be used. This choice does not change our observations from Section 3.2, Section 3.3.

Forward Representation F(s,a,z). The input to the forward representation F is always preprocessed. State-action pairs (s,a) and state-task pairs (s,z) have their own preprocessors which are feedforward MLPs that embed their inputs into a 512-dimensional space. These embeddings are concatenated and passed through a third feedforward MLP F which outputs a d-dimensional embedding vector. Note: the forward representation F is identical to ψ used by USF so their implementations are identical (see Table 4). Also, for stability reasons of TD learning, we make ensemble of F and take their mean as aggregation function.

Backward Representation B(s). The backward representation B is a feedforward MLP that takes a state as input and outputs a d-dimensional embedding vector.

Actor $\pi(s,z)$. Like the forward representation, the inputs to the policy network are similarly preprocessed. State-action pairs (s,a) and state-task pairs (s,z) have their own preprocessors which feedforward MLPs that embed their inputs into a 512-dimensional space. These embeddings are concatenated and passed through a third feedforward MLP which outputs a a-dimensional vector, where a is the action-space dimensionality. A Tanh activation is used on the last layer to normalise their scale. Note the actors used by FB and USFs are identical (see Table 4). For discrete environments, optimal policy is greedy, while for continuous DDPG-style is used for approximating argmax.

Table 4: **Hyperparameters for FB.** Hyperparameters for Belief-FB and Rotation-FB are highlighted in

Hyperparameter	Value
Latent dimension d	150 (100 for discrete)
F / ψ dimensions	(1024, 1024)
B / φ dimensions	(256, 256, 256)
Preprocessor dimensions	(1024, 1024)
Std. deviation for policy smoothing σ	0.2
Truncation level for policy smoothing	0.3
Learning steps	1,000,000
Batch size	1024
Optimiser	Adam
Learning rate	0.0001
Learning rate of $f_{\rm dyn}$	0.0001
Discount γ	0.99, 0.98 (Maze)
Activations (unless otherwise stated)	GeLU
Target network Polyak smoothing coefficient	0.05
z-inference labels	10,000
z mixing ratio	0.5
κ	50, 100 for Pointmass
Contexual representation h dimension	150 (100 for discrete)
Next state predictor g_{pred}	(256, 256, 256)

Misc. Layer normalisation and Tanh activations are used in the first layer of all MLPs to standardise the inputs as recommended in original paper for both discrete and continuous beenhmarks. Baseline is taken from official repository contrallable agent.

E.2 HILP

We take official implementation in JAX from Park et al. (2024) together with all of the hyperparameters.

E.3 TASK SAMPLING DISTRIBUTION ${\cal Z}$

Vanilla-FB. FB representations require a method for sampling the task vector z at each learning step. Touati et al. (2022) employ a mix of two methods, which we replicate:

- 1. Uniform sampling of z on the hypersphere surface of radius \sqrt{d} around the origin of \mathbb{R}^d ,
- 2. Biased sampling of z by passing states $s \sim \mathcal{D}$ through the backward representation z = B(s). This also yields vectors on the hypersphere surface due to the L2 normalization described above, but the distribution is non-uniform.

We sample $z \sim 50:50$ (either randomly or from B) from these methods at each learning step as in original work by Touati & Ollivier (2021).

Rotation-FB. After transformer $f_{\rm dyn}$ pretraining stage, RFB at each gradient step chooses task-conditioning vector $z_{\rm FB}$ based on i) context representation h acting as axes coming from $f_{\rm dyn}$ and ii) drawing task encoding vectors $z_{\rm FB}$ around this axes. We also perform normalization as in Vanilla-FB by projecting resulting vector on a surface of hypersphere of radius \sqrt{d} .

Stage ii) is implemented as drawing samples as $z_{\rm FB} \sim {\rm vMF}(\mu=h,\kappa)$. In order to remove high computational costs, we implement this sampling procedure through Householder reflection around context axes, by first drawing z from one of the basis vectors (e.g., north pole) and then performing rotation.

E.4 PSEUDOCODE **Algorithm 1** Belief-FB Training 1: **Input**: offline diverse dataset \mathcal{D} consisting of transitions based on hidden configuration variable c_i 2: Initialize transformer encoder $f_{\text{dyn}_{\theta}}$, F_{η} , B_{ω} , number of gradient steps for transformer pre-training K, context length T, Polyak coefficient, β , batch size B learning rates $\lambda_f, \lambda_F, \lambda_B$ 3: while update steps < K dosample batch of B trajectories of length T $\{(s_{i,t}, a_{i,t}, s_{i,t+1})\}_{i=1,\dots B, t=1,\dots,T} \sim \mathcal{D}$ $(\boldsymbol{\mu}_i; \log \boldsymbol{\sigma}_i), = f_{\text{dyn}_{\theta}}(\{s_{i,t}, a_{i,t}, s_{i,t+1}\}_{t=1}^M), i = 1, \dots, B,$ $z_i = \mu_i + \epsilon_i \odot \exp(\log \sigma_i),$ $\mathbf{Z}_{i,t} = \mathbf{z}_{ ext{dyn}_i}, \ t = 1, \dots, T$ # Representation $z_{ ext{dyn}}$ is shared across each sequence 7: $\hat{s}_{i,t+1} = g_{\text{pred}}(s_{i,t}, a_{i,t}, \mathbf{Z}_{i,t})$ $t=1,\ldots,T,\ i=1,\ldots,B$ $\mathcal{L}_{\text{context}} = \frac{1}{BT} \sum_{i=1}^{B} \sum_{t=1}^{T} \left\| \hat{s}_{i,t+1} - s_{i,t+1} \right\|_{2}^{2}$ $\theta_{f_{\text{dyn}}} \leftarrow \theta_{f_{\text{dyn}}} - \lambda_{f} \nabla_{\theta} \mathcal{L}_{\text{context}}(\theta)$ 10: 11: end while 12: while not converged do $\eta_F \leftarrow \eta_F - \lambda_F \nabla_{\eta_F} J_{(F,B)}(\eta_F)$ # FB training, Equation 21 $\omega_B \leftarrow \omega_B - \lambda_B \nabla_{\omega_B} J_{(F,B)}(\omega_B)$ 15: end while **Algorithm 2** Sampling z_{FB} for RFB **Input:** B (batch size), d (latent dimension), anchor matrix $\mathbf{H} \in \mathbb{R}^{B \times d}$, κ (concentration) Output: $\mathbf{Z} \in \mathbb{R}^{B \times d}$ 1: Normalize anchors: $\mathbf{u}_i \leftarrow \mathbf{H}_i/(\|\mathbf{H}_i\|_2 + \varepsilon)$ \triangleright for $i = 1, \ldots, B$ 2: $\mathbf{S} \leftarrow \text{VMF_Sample_NorthPole}(B, d, \kappa)$ \triangleright draw B VMF samples 3: **for** $i \leftarrow 1$ **to** B **do** $\mathbf{R}_i \leftarrow \text{Householder}_{\mathbf{R}} \text{Otation}(\mathbf{u}_i)$ $\mathbf{z}_i \leftarrow \mathbf{R}_i \, \mathbf{S}_i$ 6: end for 7: $\mathbf{Z} \leftarrow \text{Project_To_Sphere}(\{\mathbf{z}_i\}_{i=1}^B)$ 8: return Z