Self-Predictive Representations for Combinatorial Generalization in Behavioral Cloning

Anonymous authors Paper under double-blind review

Keywords: Sequential Decision Making, Combinatorial Generalization, Representation Learning.

Summary

Behavioral cloning (BC) methods trained with supervised learning (SL) are an effective way to learn policies from human demonstrations in domains like robotics. Goal-conditioning these policies enables a single generalist policy to capture diverse behaviors contained within an offline dataset. While goal-conditioned behavior cloning (GCBC) methods can perform well on in-distribution training tasks, they do not necessarily generalize zero-shot to tasks that require conditioning on novel state-goal pairs, i.e. *combinatorial generalization*. In part, this limitation can be attributed to a lack of temporal consistency in the state representation learned by BC; if temporally related states are encoded to similar latent representations, then the out-of-distribution gap for novel state-goal pairs would be reduced. Hence, encouraging this temporal consistency in the representation space should facilitate combinatorial generalization. Successor representations, which encode the distribution of future states visited from the current state, nicely encapsulate this property. However, previous methods for learning successor representations have relied on contrastive samples, temporal-difference (TD) learning, or both. In this work, we propose a simple yet effective representation learning objective, BYOL- γ augmented GCBC.

Contribution(s)

1. We propose BYOL- γ , a novel representation learning objective that relates to the successor measure, which we prove in the finite MDP setting.

Context: While prior theory supports this objective (Tang et al., 2022; Khetarpal et al., 2025), to our knowledge, we are the first to directly relate a BYOL-based objective to the successor measure and to use it in practice for representation learning.

2. Empirically, we demonstrate that such a BYOL objective can be used as auxiliary objective augmenting the capabilities of Behavior Cloning, and we find that BYOL- γ obtains competitive results on navigation tasks in OGBench compared to existing methods. Qualitatively, we show that the BYOL- γ objective learns similar representation structures to contrastive learning, encoding temporal distance between states.

Context: We build on the setting of using auxiliary representation learning objectives for BC from Myers et al. (2025b), however we further illustrate the relationship between generalization and approximations to the successor measure, and demonstrate that alternative representation learning can obtain competitive or better performance for achieving combinatorial generalization.

Self-Predictive Representations for Combinatorial Generalization in Behavioral Cloning

Anonymous authors

Paper under double-blind review

Abstract

1	Behavioral cloning (BC) methods trained with supervised learning (SL) are an effective
2	way to learn policies from human demonstrations in domains like robotics. Goal-
3	conditioning these policies enables a single generalist policy to capture diverse behaviors
4	contained within an offline dataset. While goal-conditioned behavior cloning (GCBC)
5	methods can perform well on in-distribution training tasks, they do not necessarily
6	generalize zero-shot to tasks that require conditioning on novel state-goal pairs, i.e.
7	combinatorial generalization. In part, this limitation can be attributed to a lack of
8	temporal consistency in the state representation learned by BC; if temporally related
9	states are encoded to similar latent representations, then the out-of-distribution gap for
10	novel state-goal pairs would be reduced. Hence, encouraging this temporal consistency
11	in the representation space should facilitate combinatorial generalization. Successor
12	representations, which encode the distribution of future states visited from the current
13	state, nicely encapsulate this property. However, previous methods for learning successor
14	representations have relied on contrastive samples, temporal-difference (TD) learning, or
15	both. In this work, we propose a simple yet effective representation learning objective,
16	BYOL- γ augmented GCBC, which is not only able to theoretically approximate the suc-
17	cessor representation in the finite MDP case without contrastive samples or TD learning,
18	but also, results in competitive empirical performance across a suite of challenging tasks
19	requiring combinatorial generalization.

20 1 Introduction

21 Generalization has been a long-standing goal in machine learning and robotics. Recently, large-scale 22 supervised models for language and vision have demonstrated impressive generalization when trained over vast amounts of data. In the robotics domain, this has motivated the development of large-scale 23 24 supervised behavior cloning (BC) models trained on offline datasets of diverse demonstrations (Ghosh 25 et al., 2024; Kim et al., 2024). However, these models still suffer from a lack of generalization. In 26 particular, while BC methods can perform well on tasks directly observed in the dataset, they often fail 27 to perform zero-shot transfer to tasks requiring novel combinations of in-distribution behavior, known 28 as combinatorial generalization. In the robotics domain, where demonstration data is time-intensive 29 and costly to produce, simply scaling the dataset is often not possible. Hence, achieving this type of 30 generalization algorithmically will be critical to unlocking the potential for large-scale supervised 31 policy training.

The property of combinatorial generalization has been previously formalized as the ability to "stitch" (Ghugare et al., 2024). Specifically, stitching refers to the ability of a policy to reach a goal state from a start state when trained on a dataset of trajectories which, provides sufficient coverage of the path to the goal, but which does not contain a single complete trajectory of the path. The lack of stitching observed in goal-conditioned behavioral cloning (GCBC) and, more generally, supervised learning, can be understood through the inductive biases of the model. By construction, BC methods do not encode the inductive bias that the observed data are generated from a Markov decision process



Figure 1: (a) Self-predictive Representations. Example training trajectories, $s_0 \rightarrow s_h$ and $s_b \rightarrow s_f$, which intersect at w. After training on these trajectories, we evaluate on a task like $s_0 \rightarrow s_f$, requiring combinatorial generalization. To learn better representations for generalization, a self-predictive representation predicts a future state $\phi(w)$ from an earlier state $\phi(e)$ via $\psi(\phi(e))$. (b) Representation learning with BYOL- γ . We predict future state representations $\phi(s_{t+k})$ via $\psi_f(\phi(s_t), a)$, and also predict backwards with $\psi_b(\phi(s_{t+k}))$. The target offset is sampled geometrically: $k \sim \text{geom}(1 - \gamma)$. Stop-gradients are denoted by //. We provide more details on the loss \mathcal{L} in Section 4.2.

39 (MDP). In contrast, reinforcement learning (RL) policies that are trained via temporal difference (TD) 40 learning directly utilize the structure of the MDP, and pass information through time using dynamic 41 programming. Offline RL (Levine et al., 2020) has been proposed as a method for achieving stitching 42 in policies trained on offline datasets. However, these methods are challenging to scale due to the 43 instability of bootstrapping in TD learning when combined with fully offline training. Scaling has 44 been more successful with supervised methods, such as in robotics, where training robot foundation 45 models with BC (Ghosh et al., 2024; Kim et al., 2024) on large-scale datasets (O'Neill et al., 2024; 46 Khazatsky et al., 2024) can lead to more general-purpose policies.

Various methods have attempted to imbue GCBC models with the ability to stitch by augmenting 47 48 the training data using the Markovian assumption (Yamagata et al., 2023; Ghugare et al., 2024). 49 However, these methods ultimately rely on similar training procedures as offline RL or otherwise 50 require a distance metric already aligned with temporal distance in the MDP. Another approach 51 frames stitching as a representation learning problem (Myers et al., 2025b). In this context, the 52 objective is to learn a latent representation of states that reflects temporal proximity in the underlying 53 MDP to facilitate generalization to unseen state-goal pairs. In particular, Myers et al. (2025b) train a 54 representation that approximates the successor measure (SM) (Blier et al., 2021) through contrastive 55 learning (CL) (van den Oord et al., 2019) as an auxiliary loss in GCBC, demonstrating enhanced 56 combinatorial generalization. This is a promising step towards achieving combinatorial generalization 57 in GCBC methods. In this work, we expand on this connection between the successor measure and 58 combinatorial generalization, proposing an alternative objective that captures similar properties with 59 a simpler learning procedure and achieves as good or better generalization performance.

In particular, in other domains, like vision, it has been found that contrastive learning can often be substituted with self-predictive representations (Grill et al., 2020), which have also been found to be useful as auxiliary losses for model-free RL (Schwarzer et al., 2020). Adapting the *Bootstrap Your Own Latent* (BYOL) framework (Grill et al., 2020) to the RL setting, representations can be trained by predicting the latent representation of the next state from the current latent state representation. These BYOL objectives have appealing properties, such as neither relying on negative examples nor reconstruction.

As motivation, we foremost evaluate the BYOL objective as an auxiliary loss for GCBC; however, we find that there exists an empirical and theoretical gap compared to contrastive methods. This

leads us to propose a novel objective, **BYOL**- γ , which removes the gap between self-predictive and

- 70 contrastive objectives. Concretely, BYOL- γ predicts latent representations of states sampled from a
- 71 γ -discounted future state distribution. While the standard BYOL objective has been shown to learn
- representations capturing spectral information about the one-step transition dynamics (Khetarpal et al.,
- 73 2025), we show that the representations learned by BYOL- γ capture spectral information related to
- The successor measure. In the finite MDP case, we show that, in fact, BYOL- γ approximates the
- successor representation. Empirically, we demonstrate on the challenging OGBench dataset (Park et al., 2025) that BYOL- γ augmented GCBC is competitive with contrastive methods in improving
- combinatorial generalization. Key contributions of our work are as follows:
- We propose BYOL-γ, a novel representation learning objective that relates to the successor measure,
 which we prove in the finite MDP setting.
- Empirically, we demonstrate that such a BYOL objective can be used as auxiliary objective augmenting the capabilities of Behavior Cloning, and we find that BYOL- γ obtains competitive results on navigation tasks in OGBench compared to existing methods.
- Qualitatively, we show that the BYOL-γ objective learns similar representation structures to CL,
 encoding temporal distance between states.

85 2 Related Work

86 Stitching in Supervised Methods. Outcome (goals or return)-conditioned behavioral cloning 87 (OCBC) methods (Schmidhuber, 2020; Chen et al., 2021; Emmons et al., 2022) provide a simple and scalable alternative to traditional offline RL (Levine et al., 2020) methods. However, these methods 88 89 do not properly "stitch" and generalize to unseen outcomes Brandfonbrener et al. (2022); Ghugare 90 et al. (2024). To reduce this problem, various works have proposed augmenting training data used 91 by BC methods. Some work incorporates methodlogy from offline RL to label returns or goals for 92 downstream SL (Char et al., 2022; Yamagata et al., 2023). Other work has considered relabeling 93 goals through clustering states (Ghugare et al., 2024), which relies on a good distance metric in 94 state-space and is limited to short trajectory stitches. Other work has utilized planing Zhou et al. 95 (2024) for goal relabeling, or generative models to synthesize new trajectories (Lu et al., 2023; Lee 96 et al., 2024). Rather than using models to generate data for BC, other work directly evaluates the 97 combinatorial generalization achieved by planning with generative models (Luo et al., 2025). In 98 this work, we neither require combining SL with explicit Q-learning, utilize generative models, or 99 perform explicit planning.

100 Visual Representation Learning. Instead of auxiliary representation learning, prior work has 101 considered using pretrained visual representations for BC. Utilizing pretrained BC has been a reliable 102 method to efficiently learn policies and improve generalization (Radosavovic et al., 2022; Majumdar 103 et al., 2023; Nair et al., 2022). Instead of pretraining representations out-of-domain, such as on 104 large image datasets, DynaMo (Cui et al., 2024) evaluates pretraining representations in-domain on 105 specific robotics datasets, and then performs a separate BC finetuning stage. However, we focus on 106 representation learning as an auxiliary objective. For enhancing combinatorial generalization, this 107 can be advantageous, as training as an auxiliary task maintains the structure of the representation 108 space, and prevents overfitting the BC objective.

109 **Representation learning in RL.** Our objective is most closely related to approaches using auxiliary 110 BYOL objectives in online RL (Gelada et al., 2019; Schwarzer et al., 2020; Ni et al., 2024; Voelcker 111 et al., 2024). These objectives can help with sample-efficiency, such as in challenging, partially observed environments with sparse rewards, or with noisy states. Additionally, self-predictive 112 113 dynamics models are used in planning and model-based RL (François-Lavet et al., 2019; Ye et al., 114 2021; Hansen et al., 2022). Various works have also characterized the dynamics of BYOL objectives 115 in the RL setting, showing that BYOL objectives capture spectral information about the policy's 116 transitions (Tang et al., 2022; Khetarpal et al., 2025). In the offline setting, how well Joint Embedding 117 Predictive Architecture (JEPA) world models generalize when used for explicit planning has been 118 studied Sobal et al. (2025), however not for combinatorial generalization. Additionally, certain 119 representation structures for value functions, namely quasimetrics (Liu et al., 2023; Wang et al., 120 2023; Wang and Isola, 2022; Myers et al., 2024) can also lead to policies that better generalize to

longer horizons (Myers et al., 2025a). Quantities related to the Successor Representation (SR) 121

122 (Dayan, 1993), such as successor features (SF) (Barreto et al., 2017), and the successor measure (SM)

123 (Blier et al., 2021) have been widely used for generalization and transfer in reinforcement learning 124 (Carvalho et al., 2024). Similarly to BYOL, these objectives have been used for representation

learning in RL (Lan et al., 2022; Farebrother et al., 2023). While prior BYOL methods either perform 125

126 1-step, or relatively short fixed n-step prediction, neither of these choices directly approximate the

127 successor measure. Our setup is most related to temporal representation alignment (TRA) (Myers

128 et al., 2025b), which recently proposed using contrastive learning as an auxiliary objective for BC to

129 improve combinatorial generalization. In this work, we further build on the relationship between the

130 SM and combinatorial generalization, and propose new representation learning objectives which can

131 lead to better performance.

Background 132 3

133 Controlled Markov Process. We consider goal-conditioned decision-making problems, with state 134 space S, goals $g \in S$, action space A, initial state distribution $p_0(s)$, dynamics $p(s_{t+1} | s_t, a)$, and

135 with policies $\pi(a|s, q)$.

Successor Representation (SR) and Successor Measure (SM). In a finite MDP, the successor 136

137 representation (SR) (Dayan, 1993) of a policy is:

$$M^{\pi}(s,s') := \mathbb{E}\left[\sum_{t \ge 0} \gamma^{t} \mathbb{1}_{(s_{t+1}=s')} \mid s_{0} = s, \pi\right]$$
(1)

We use the convention of counting from s_{t+1} , writing in matrix form $M^{\pi} = \sum_{t\geq 0} \gamma^t (P^{\pi})^{t+1}$. The transition matrix transition for policy π is P^{π} , with $P_{i,j}^{\pi} = \sum_a \pi(a|s=i)P_{i,a,j}$, where $P_{i,a,j} = p(s_{t+1} = j \mid s_t = i, a)$. The successor representation also satisfies the bellman equation, $M^{\pi} = P^{\pi} + \gamma P^{\pi} M^{\pi} = P^{\pi} (I - \gamma P^{\pi})^{-1}$. For a fixed policy, the successor representation 138 139 140 141 describes a type of temporal distance between states. The successor measure (SM) (Blier et al., 142 2021) extends SR to continuous spaces S: $M^{\pi}(s, X) := \sum_{t>0} \gamma^t P(s_{t+1} \in X \mid s) \; \forall X \subset S.$ We 143 also define the *normalized* successor representation, or measure $\tilde{M}^{\pi} = (1 - \gamma)M^{\pi}$. In the finite case, 144 the normalized successor representation \tilde{M}^{π} has rows that sum to one like transitions P^{π} . Another 145 quantity successor features (SF) (Barreto et al., 2017) are the expected discounted sum of future 146 features $\phi(s) \in \mathbb{R}^d$: $\psi^{\pi}(s) = \mathbb{E}\left[\sum_{t \ge 0} \gamma^t \phi(s_{t+1}) \mid s_0 = s, \pi\right]$. We can relate SFs to the SM with $\psi^{\pi}(s) = \int_{s'} M^{\pi}(s, s') \phi(s')$. Each of these quantities can also be defined with conditioning on the 147 148 149 first action and then following the policy, e.g. $M^{\pi}(s, a, s')$.

150 3.1 Representation Learning

- 151 We begin with two representation learning methods that approximate the density of the SM.
- 152 Forward-Backward. We consider a simplified version of the Forward-Backward loss that approxi-153 mates the successor measure for a fixed policy π , discussed by Touati et al. (2023).

$$\min_{\phi,\psi} \mathbb{E}_{\substack{s_t \sim p(s), s' \sim p(s)\\ s_{t+1} \sim p^{\pi}(s_{t+1}|s_t)}} \left[(\psi(s_t)^T \phi(s') - \gamma \bar{\psi}(s_{t+1})^T \bar{\phi}(s'))^2 \right] - 2 \mathbb{E}_{\substack{s_t \sim p(s)\\ s_{t+1} \sim p^{\pi}(s_{t+1}|s_t)}} \left[\psi(s_t)^T \phi(s_{t+1})^T \bar{\phi}(s') \right]$$
(2)

154 FB learns an approximation of the successor measure with factorization $M^{\pi}(s,s_{+}) \approx$ 155 $\psi(s_t)\phi(s_+)p(s_+)$ using TD learning. Given transitions (s_t, s_{t+1}) sampled by a policy π , the second 156 term relates to fitting $M^{\pi}(s_t, s_{t+1})$. Given an independently sampled state s', the first term bootstraps an estimate of $M^{\pi}(s_t, s')$ from $\overline{M}^{\pi}(s_{t+1}, s')$, where $\overline{\phi}, \overline{\psi}$ denote stop-gradient operations. 157

158 Contrastive Learning. Temporal contrastive learning used in MDPs (Eysenbach et al., 2022) is

159 related to a Monte Carlo (MC) approximation of the (discounted) successor measure. This can be

- 160 implemented with a InfoNCE (van den Oord et al., 2019) loss that maximizes the similarity of a
- 161 positive pair between a state s_t and a future state from the same trajectory s_+ , and minimizing the
- 162 similarity of s_t and random states s_- :

$$\max_{\substack{\phi,\psi\\k\sim\text{geom}(1-\gamma)\\s_{+}=s_{t+k},s^{2:N}\sim p(s)}} \mathbb{E}\left[\log\frac{e^{f(\psi(s_{t}),\phi(s_{+}))}}{\sum_{i=2}^{N}e^{f(\psi(s_{t}),\phi(s_{-}^{i})}}\right]$$
(3)

A common choice for the energy function f is the inner product $f(\psi(s)\phi(s_+)) = \psi(s)^T \phi(s_+)$. A 163 key aspect to note is that the positive sample s_+ comes from an MC sample from $s_+ \sim M^{\pi}(s_t, s_+)$. 164 The optimal solution to (3) gives $\tilde{M}^{\pi}(s, s_{+}) \approx C \exp(\psi(s_{t})^{T} \phi(s_{+})) \cdot p(s)$. However in-practice, 165 166 we only have dataset of MC samples from π , i.e. fixed-length trajectories, which means we do not actually estimate the SM of π . In Appendix C, we further elaborate on the relationship between the 167 168 FB loss, and CL. Particularly, in the limit, an n-step version of FB is related to CL. 169 **BYOL.** We now look at an objective that captures information about single-step transition instead of

170 the successor measure. In the context of RL, self-predictive models jointly learn a latent space and 171 a dynamics model through predicting future latent representations. Self-predictive models rely on latent bootstrapped targets (BYOL) (Grill et al., 2020), avoiding reconstruction (generative models), 172 173 or negative samples (contrastive learning). Self-predictive models are also a type of joint-embedding 174 predictive architectures (JEPAs) (LeCun, 2022; Garrido et al., 2024).

175 Given an encoder which produces a representation $z_t = \phi(s_t)$, and dynamics $\psi(z_{t+1}|z_t)$ for a fixed 176 policy π , we minimize the difference between our prediction and target representation in latent-space:

$$\min_{\phi,\psi} \mathbb{E}_{s_t \sim p(s), s_{t+1} \sim p^{\pi}(s_{t+1}|s_t)}, \left[f(\psi(\phi(s_t)), \bar{\phi}(s_{t+1})) \right]$$
(4)

is a convex function such as the squared
$$l_2$$
 norm, and $\bar{\phi}$ refers to an EMA target, or stop-

178 gradient. Variants of this BYOL objective have been widely used to learn state abstractions, and

work as an auxiliary loss when approximating the value function in deep RL (Gelada et al., 2019; 179

180 Schwarzer et al., 2020; Ni et al., 2024). In the finite MDP, this objective captures spectral information

181 about the policy's transition matrix P^{π} (Tang et al., 2022; Khetarpal et al., 2025) which we discuss in

182 Appendix D.1.

Where *f*

177

183 3.2 Combinatorial Generalization from Offline Data

We now shift focus on how we can learn policies from offline data using behavioral cloning, and then 184 185 introduce a combinatorial generalization gap that arises in this setting.

We consider a **dataset** $\mathcal{D} = \{(s_0^i, a_0^i, \cdots, s_T^i, a_T^i)\}_{i=1}^N$, composed of trajectories generated by a set 186 of unknown policies $\{\beta_i(a|s)\}$. Goal Conditioned Behavioral Cloning (GCBC) trains a policy 187 188 with maximum likelihood to reproduce the behaviors from the dataset. After sampling a current state, 189 a goal is sampled as a future state from the same trajectory:

$$\max \mathcal{L}_{BC}(\pi) = \max \mathbb{E} \qquad \qquad \beta_{i} \sim p(\beta_{i}) \qquad s \sim p(s|\beta_{i}) \qquad [\log \pi(a|s, g = s_{+})]$$

 $\underset{a \sim \beta_j(a \mid s), s \leftarrow \mathcal{M}^{\beta_j}(s, s_+)}{ \underset{a \sim \beta_j(a \mid s), s_+ \sim \mathcal{M}^{\beta_j}(s, s_+)}{ \underset{a \sim \beta_j(a \mid s), s_+ \sim \mathcal{M}^{\beta_j}(s, s_+)}} } [\log \pi(a \mid s, g \mid s_j)]$ π Generalization gap. While this policy can perform well in-distribution, the behavior cloning policy

190 191 struggles to generalize to reach goals from states that are not in matching training trajectories. We 192 now review a more formal definition of this type of generalization gap.

We consider Lemma 3.1 from Ghugare et al. (2024), which says there exists a single Markovian 193 194 policy $\beta(a|s)$ that has the same occupancy as the mixture of j policies:

$$M^{\beta}(s) = \mathbb{E}_{p(\beta_{j})} \left[M^{\beta_{j}}(s) \right]$$
(6)

(5)

This policy also has construction: $\beta(a \mid s) := \sum_{j} \beta_j(a \mid s) p(\beta_j \mid s)$, where $p(\beta_j \mid s)$ is the distribution 195 196 over policies in s as reflected by the dataset.

197 Using the successor measure of the individual policies, and the mixture policy, we can quantify a gap

between accomplishing out-of-distribution tasks versus in-distribution training tasks (Ghugare et al.,
 2024):

$$\underbrace{\mathbb{E}_{\substack{s_0 \sim M^{\beta}(s_0) \\ s_g \sim M^{\beta}(s_0, s_g)}}}_{\text{tasks requiring combinatorial generalization}} \begin{bmatrix} u^{\pi}(s_0, s_g) \end{bmatrix} - \underbrace{\mathbb{E}_{\beta_j \sim p(\beta_j), s_0 \sim M^{\beta_j}(s_0)} \left[u^{\pi}(s_0, s_g) \right]}_{\substack{s_g \sim M^{\beta_j}(s_0, s_g) \\ \text{in-distribution training tasks}}$$
(7)

Here, u is a performance metric of the policy π such as the success rate to reach s_g from s_0 . As we perform well on in-distribution tasks due to a correspondence to Equation (5), the BC policy has no guarantees for the first term. This is because after sampling a state, the goal is sampled from the

203 successor measure of the mixture policy.

4 Closing the Generalization Gap with Representations

г

In this section, we aim to close the aforementioned generalization gap. We consider a policy trained with the BC objective π to be made more robust to the tasks requiring combinatorial generalization through representation learning. We begin with a setup similar to Equation (7), but with a shared initial state s_0 for both the in-distribution and out-of-distribution task. For the in-distribution task, we sample a goal as before, labeled as s_w . However, for the out-of-distribution task, we sample a goal s_f to be a state that can be reached by the mixture policy β after s_w . (8):

$$\mathbb{E}_{\substack{\beta_{j} \sim p(\beta_{j}), s_{0} \sim M^{\beta_{j}}(s)\\s_{w} \sim M^{\beta_{j}}(s_{0}, s_{w})}} \left[\underbrace{\mathbb{E}_{s_{f} \sim M^{\beta}(s_{w}, s_{f})} \left[u^{\pi}(s_{0}, s_{f}) \right) \right]}_{\text{extended task requiring generalization}} - \underbrace{u^{\pi}(s_{0}, s_{w})}_{\text{in-distribution task}} \right]$$
(8)

٦

$$= \mathbb{E}_{\substack{\beta_j \sim p(\beta_j), s_0 \sim M^{\beta_j}(s)\\ s_w \sim M^{\beta_j}(s_0, s_w)}} \left[\mathbb{E}_{s_f \sim M^{\beta}(s_w, s_f)} \left[u^{\pi}(s_0, \phi(s_f)) \right] - u^{\pi}(s_0, \phi(s_w)) \right]$$
(9)
want invariance with respect to future goals through ϕ

Then, in Equation (9) we add a goal representation ϕ that processes the goal before going to policy π . Intuitively, a policy could achieve the out-of-distribution task by first going from s_0 to s_w (indistribution), and then completing the remaining task s_w to s_f . In essence, we want that when conditioning on $\phi(s_f)$, the policy should first go to s_w , which can be achieved by learning ϕ , where $\phi(s_w)$ is similar to $\phi(s_f)$ (Myers et al., 2025b). More formally, for $s_f \sim M^\beta(s_w, s_f)$ we want an invariance $\phi(s_f) \approx \phi(s_w)$.

217 From this observation, we can understand that approximating the successor measure of the mixture 218 policy β , when parameterized by ϕ , will build the desired representation. One choice for the 219 representation learning objective can be the FB algorithm, which obtains a factorization $M^{\beta}(s, s_{+}) \approx$ 220 $\psi(s_t)\phi(s_+)p(s_+)$ (Touati et al., 2023). FB utilizes TD-learning to learn the same representations given transitions (s_t, s_{t+1}) , regardless of how these transitions are divided between trajectories. 221 222 However, FB may not scale as well to large datasets and high-dimensional state spaces. This may 223 lead us to prefer an MC approximation of the SM, using CL. A key point here is that we do not 224 have MC samples from β , only the individual policies β_j , so CL does not directly approximate M^{β} . However, in practice, CL can still build a representation space with a compositional structure (Myers 225 226 et al., 2025b), and may be a more scalable option than FB.

227 4.1 BYOL- γ : Connecting self-predictive objectives to the successor representation

To build representations that lead to generalization, we propose a predictive objective, relying on neither TD learning nor negative samples. Specifically, we propose BYOL- γ which allows us to use the BYOL framework to capture information related to *successor representations* and its generalizations. Given a state s_t , a BYOL objective would normally sample a prediction target from a one-step transition s_{t+1} as in Equation (4). However, we make a modification to predict empirical samples from the normalized successor measure:

$$\mathcal{L}_{\text{BYOL-}\gamma}(\phi,\psi) = \mathbb{E}_{\substack{s_t \sim p(s), k \sim \text{geom}(1-\gamma)\\s_{t+k} \sim p^{\pi}(s_{t+k}|s_t)}} \left[f(\psi(\phi(s_t)), \bar{\phi}(s_{t+k})) \right]$$
(10)

Where *f* refers to an energy function, ϕ refers to the encoder, and ψ the predictor. With $\gamma = 0$, we have $s_{t+k} = s_{t+1}$ corresponding to an approximation of the one-step transitions, recovering the base BYOL objective. Figure 1b depicts our overall representation learning objective. We can view this objective as iteratively minimizing an upper-bound on the error between $\psi(\phi((s)))$ and the true successor features of the policy ψ^{π} with changing basis features $\overline{\phi}$. With convex *f*, by Jensen's inequality we have:

$$\mathcal{L}_{\text{BYOL-}\gamma}(\phi,\psi) = \mathbb{E}_{s_t \sim p(s), s_+ \sim \tilde{M}^{\pi}(s_t, s_+)} \left[f(\psi(\phi(s_t)), \bar{\phi}(s_+)) \right]$$
(11)

$$\geq \mathbb{E}_{s_t \sim p(s)}, \left[f(\psi(\phi(s_t)) \mathbb{E}_{s_+ \sim \tilde{M}^{\pi}(s_t, s_+)} \bar{\phi}(s_+)) \right]$$
(12)

$$= \mathbb{E}_{s_t \sim p(s)} \left[f(\psi(\phi(s_t)), (1 - \gamma)\psi_{\phi}^{\pi}(s_t) \right]$$
(13)

Specifically, in the finite MDP, we can precisely describe the relationship of our objective to the SR with the following result:

Theorem 4.1. Given a finite MDP with linear representations $\Phi \in \mathbb{R}^{|S| \times d}$, and predictor $\Psi \in \mathbb{R}^{d \times d}$, under assumptions of orthogonal initialization for Φ (Ass. D.1), a uniform initial state distribution $p_0(s)$ (Ass. D.2), and symmetric transition dynamics (Ass. D.3), minimizing the selfpredictive learning objective $\mathcal{L}_{BYOL-\gamma}(\phi, \psi)$ approximates a matrix decomposition of the successor representation $\tilde{M}^{\pi} \approx \Phi \Psi \Phi^T$, corresponding to successor features $(1 - \gamma)\Psi^{\pi} \approx \Psi \Phi$.

Proof is in Appendix D.2, where we show that existing theory (Khetarpal et al., 2025) also translates to the proposed BYOL- γ objective. Finally, we can see the relation between this objective and CL (3), with the most striking difference being the removal of the denominator involving negative samples. Surprisingly, we reveal that this simplified system still captures similar information and also can lead to empirical generalization in Section 5.1.

BYOL- γ **Variants.** We discuss a few variants on our base objective, namely, we find it beneficial to consider **bidirectional prediction** (Guo et al., 2020; Tang et al., 2022) where we add an additional backwards predictor ψ_b which predicts a past representation from the future:

$$\mathcal{L}_{\text{BYOL-}\gamma}(\phi,\psi) = \mathbb{E}_{s_t \sim p(s), s_t \sim \tilde{M}^{\pi}(s_t, s_+)} \left[f(\psi_f(\phi(s_t)), \bar{\phi}(s_+)) + f(\bar{\phi}(s_t), \psi_b(\phi(s_+))) \right]$$
(14)

We utilize an **action-conditioned** variant of the forward predictor $\psi_f(\phi(s_t), a_t)$, which can be interpreted as a temporally extended latent dynamics model, or capturing information about $\tilde{M}^{\pi}(s, a, s_{\perp})$.

257 4.2 Training a policy with auxiliary BYOL- γ

We consider BYOL- γ as an auxiliary loss for a BC policy $\pi_{\Theta}(a|s,g), \Theta = (\theta, \phi, \psi)$, to better address generalization of a policy. The parameters of the encoder and predictor correspond to ϕ, ψ , and θ includes additional parameters such as a policy head which transforms representations to actions. With this policy, we train with the objective:

$$\mathcal{L}_{BC}(\Theta) + \alpha \mathcal{L}_{BYOL\gamma}(\phi, \psi)$$

$$= \mathbb{E}_{(s_t, a_t, s_+) \sim \mathcal{D}} \left[-\log \pi_{\Theta}(a | \phi(s), g = \phi(s_+)) \right]$$

$$+ \alpha \mathbb{E}_{(s_t, a_t, s_+) \sim \mathcal{D}} \left[f(\psi_f(\phi(s_t), a_t), \bar{\phi}(s_+)) + f(\bar{\phi}(s_t), \psi_b(\phi(s_+))) \right]$$
(15)

For *f*, we choose a cross-entropy loss between the (softmax) normalized representations of the prediction and the target, similar to DINO (Caron et al., 2021): $f_{CE}(a, b) = \text{softmax}(b) \cdot \log \text{softmax}(a)$. We also find a normalized l_2 loss, $f_{l_2} = \|\frac{a}{\|a\|} - \frac{b}{\|b\|}\|_2^2$, commonly used in BYOL setups (Grill et al., 2020; Schwarzer et al., 2020) also works, which we ablate in Section 5.4. We describe additional training details in Appendix A.1.

267 **5** Experiments

268 We have shown the theoretical basis for BYOL- γ as an appropriate choice of representation learn-269 ing objective for combinatorial generalization. Next, we demonstrate its performance empirically. Namely, we compare our proposed method BYOL- γ to alternative representation learning methods. 270 271 We compare representation learning algorithms across three axes: (1) First, we compare qualitatively 272 whether the representations appear to capture temporal relationships (2) Second, we assess represen-273 tations quantitatively by measuring zero-shot generalization performance on unseen tasks that require 274 combinatorial generalization (3) Third, we assess generalization performance over an increasing 275 generalization horizon. Finally, we perform ablations on the various components of our proposed 276 method to demonstrate the relative importance of each algorithmic choice.

Environments We empirically evaluate how well our approach can help with combinatorial generalization on offline goal-reaching tasks on OGBench (Park et al., 2025), which contains both navigation and manipulation tasks, across both low-dimensional and visual observations. We focus on navigation environments, where OGBench provides stitch datasets, that assess combinatorial generalization by training on trajectories that span at most 4 maze cells, while evaluating on tasks that are longer,

282 requiring combining information from multiple smaller trajectories.

283 **Baselines** We benchmark against non-hierarchical methods, that perform end-to-end control from 284 state to low-level actions (e.g. joint-control). In addition to BYOL- γ used as an auxiliary loss for BC, 285 we evaluate several baselines: GCBC is the standard behavioral cloning baseline, which we aim to 286 improve upon with representation learning objectives. Offline RL from OGBench, including implicit {V,Q}-learning (IVL, IQL) (Kostrikov et al., 2022), Quasimetric RL (QRL)(Wang et al., 2023), 287 288 and Contrastive RL (CRL) (Eysenbach et al., 2022). BYOL is a minimal version of our BYOL- γ 289 setup with 1-step prediction ($\gamma = 0$), only forwards prediction (ψ_f) without action-conditioning 290 $(\psi_f(\phi(s_t)))$, and loss f_{l_2} . **TRA** (Myers et al., 2025b) is an auxiliary representation objective using 291 contrastive learning related to an MC approximation of the successor measure FB, an on-policy 292 version of forward-backward representation described in Equation (2) used as an auxiliary objective, 293 and is related to TD approximation of the successor measure.

294 **Experimental Setup** We match the training details of OGBench, and consider a similar represen-295 tation learning setup to TRA. For TRA and FB, we utilize a similar setup described in Section 4.2. 296 We found it was beneficial to add action conditioning to FB, but did not see an overall improvement 297 for TRA, so we use the original setup without action-conditioning. We provide a full comparison 298 for action-conditioning in Appendix B.1. We note that when we perform action-conditioning, we 299 change the representation of the policy from $\pi(\phi(s), \psi(g))$ to $\pi(\phi(s), \phi(g))$. In Table 1, we utilize 300 superscript a to denote methods with action-conditioning. Notably, we find that the weight of the 301 auxiliary representation learning objectives (α) can be sensitive to both the embodiment, and size of 302 environment (medium vs large). For each method, we perform a hyperparameter sweep over 4 α 303 values, and report the best result for each environment in Table 1. We hold other hyperparameters 304 constant across experiments, except with variation between non-visual and visual environments noted 305 in Appendix A.

306 5.1 Qualitative analysis of representations

In Figure 2, we display a qualitative analysis of the representations. We visualize the similarity between the future prediction ψ for each state to $\phi(g)$ for a fixed goal g. We can see that **BYOL**- γ seems to learn a representation that encodes reachability between states, and has a similar structure to **FB**, which is known to approximate the successor measure. **TRA** and base **BYOL** seem to both capture similar structure and learn a less well-defined latent space. However, BYOL- γ and FB have more distinct similarity, and have visible "paths" of similar states. **BYOL**- γ also appears to capture the most similarity among more distant pairs of states. Compared to **TRA**, our hypothesis here is that 314 we have more optimistic similarity between distant states due to the lack of a negative term in the loss, which pushes representations apart.



Figure 2: Visualization of the Learned Representation: depicts the similarity between the prediction of the current state representation to the goal representation. For BYOL- γ and FB, we visualize the cosine similarity between $\psi(\phi(s, a)), \phi(g) \forall s \in D$ for a fixed goal g which is indicated by the star marked in red. For TRA, we compare $\psi(s), \phi(g)$. BYOL- γ captures similar temporal relationships as the baseline methods.

316 5.2 Zero-shot performance on combinatorial generalization tasks

315

In Table 1, we provide the performance results across all methods. Overall, our proposed method

318 BYOL- γ , shows improved performance vs. GCBC across most environments, and is either com-

petitive with or outperforms FB and TRA. Importantly, we find that a minimal BYOL setup does

320 not confer significant benefit over the base GCBC except in non-visual antmaze environments.

321 Generally, auxiliary representation learning with GCBC outperforms existing offline RL methods.

Within the auxillary loss methods, we find that **FB** and **BYOL**- γ tend to outperform **TRA** on most environments. While we find that **FB** outperforms **BYOL**- γ on environments with smaller state spaces (antmaze-{medium, large}), we find that **BYOL**- γ 's simpler training procedure is beneficial in environments with larger state spaces (humanoidmaze-{medium, large}, visual-antmaze-medium and visual-scene-play).

Dataset	BYOL- γ^a	BYOL	TRA	FB ^a	GCBC	GCIVL	GCIQL	QRL	CRL
antmaze-medium-stitch antmaze-large-stitch humanoidmaze-medium-stitch humanoidmaze-large-stitch antsoccer-arena-stitch	$58 \pm 519 \pm 751 \pm 613 \pm 325 \pm 5$	$\begin{array}{c} 59\pm 4 \\ 17\pm 6 \\ 23\pm 3 \\ 3\pm 1 \\ 12\pm 7 \end{array}$	$\begin{array}{c} 54\pm 6\\ 11\pm 8\\ 45\pm 8\\ 5\pm 4\\ 14\pm 4\end{array}$	$\begin{array}{c} {\bf 64} \pm 6 \\ {\bf 23} \pm 4 \\ {\bf 42} \pm 4 \\ {\bf 11} \pm 3 \\ {\bf 22} \pm 10 \end{array}$	$\begin{array}{c} 45 \pm 11 \\ 3 \pm 3 \\ 29 \pm 5 \\ 6 \pm 3 \\ 24 \pm 8 \end{array}$	$\begin{array}{c} 44 \pm 6 \\ 18 \pm 2 \\ 12 \pm 2 \\ 1 \pm 1 \\ 21 \pm 3 \end{array}$	$\begin{array}{c} 29 \pm 6 \\ 7 \pm 2 \\ 12 \pm 3 \\ 0 \pm 0 \\ 2 \pm 0 \end{array}$	$\begin{array}{c} 59\pm7\\ 18\pm2\\ 18\pm2\\ 3\pm1\\ 1\pm1 \end{array}$	$\begin{array}{c} 53 \pm 6 \\ 11 \pm 2 \\ 36 \pm 2 \\ 4 \pm 1 \\ 1 \pm 0 \end{array}$
visual-antmaze-medium-stitch visual-antmaze-large-stitch visual-scene-play	$ \begin{array}{r} {\bf 68} \pm 4 \\ {\bf 26} \pm 5 \\ {\bf 17} \pm 1 \end{array} $	$\begin{array}{c} 57 \pm 8 \\ 26 \pm 5 \\ 13 \pm 3 \end{array}$	$\begin{array}{c} 52\pm3\\17\pm1\\\textbf{16}\pm3\end{array}$	$\begin{array}{c} 49 \pm 2 \\ \textbf{29} \pm 2 \\ 14 \pm 1 \end{array}$	$\begin{array}{c} {\bf 67} \pm 4 \\ 24 \pm 3 \\ 12 \pm 2 \end{array}$		$\begin{array}{c} 2 \pm 0 \\ 0 \pm 0 \\ 12 \pm 2 \end{array}$	$\begin{array}{c} 0 \pm 0 \\ 1 \pm 1 \\ 10 \pm 1 \end{array}$	$\begin{array}{c} 69 \pm 2 \\ 11 \pm 3 \\ 11 \pm 2 \end{array}$
average	35	26	27	32	26	16	8	14	25

Table 1: **OGBench: We find that BYOL-** γ **performs better overall compared to prior methods**. We report mean and standard deviation over 10 training seeds in non-visual environments, and 4 seeds in visual environments. We match the OGBench evaluation setup of 5 evaluation (state,goal) tasks, and 50 episodes per task. The success rate is then averaged over the last 3 checkpoints. We color the best non-RL method, and **bold** values within 95% of its value in the same row. We use superscript *a* to denote methods utilizing action-conditioning.

- 327 Interestingly, in visual-antmaze TRA and FB actually seem to hurt performance in comparison
- to base GCBC. On the other hand, with BYOL- γ we see no performance degradation over GCBC
- 329 on the visual environments, a considerable improvement over other methods.

330 **5.3** Evaluating generalization with increasing horizon

331 We conduct experiments to understand how success rate changes as an agent has to reach more

332 challenging goals further away from its starting position. For each maze environment, we consider

Dataset	BYOL- γ^a	-a	f_{l_2}	$-\psi_b$	$\gamma = 0$
antmaze-medium-stitch antmaze-large-stitch humanoidmaze-medium-stitch humanoidmaze-large-stitch antsoccer-arena-stitch	$61 \pm 621 \pm 554 \pm 514 \pm 221 \pm 4$	$\begin{array}{c} 63 \pm 9 \\ \textbf{27} \pm 7 \\ 48 \pm 5 \\ 12 \pm 6 \\ 20 \pm 5 \end{array}$	$\begin{array}{c} 56 \pm 4 \\ 24 \pm 6 \\ 49 \pm 6 \\ \textbf{15} \pm 7 \\ 11 \pm 5 \end{array}$	$\begin{array}{c} {\bf 67}\pm 2\\ 19\pm 7\\ {\bf 52}\pm 5\\ 13\pm 2\\ {\bf 27}\pm 7 \end{array}$	$59 \pm 5 \\ 8 \pm 4 \\ 18 \pm 2 \\ 3 \pm 1 \\ 25 \pm 7$
visual-antmaze-medium-stitch visual-antmaze-large-stitch	$\frac{68\pm 4}{26\pm 5}$	$\begin{array}{c} 65 \pm 3 \\ 25 \pm 8 \end{array}$	$\begin{array}{c} 63\pm5\\ \textbf{27}\pm7 \end{array}$	$\begin{array}{c} 61\pm 4 \\ \textbf{28}\pm 2 \end{array}$	$\begin{array}{c} 54\pm9\\ 28\pm1 \end{array}$
average	33	33	31	33	24

Table 2: **BYOL**- γ **ablations.** We ablate components of our representation learning objective. For each ablation, we perform a hyperparameter sweep over α , and report the best result per-environment. For all environments, we report results over 4 seeds (for **BYOL**- γ , we use the first 4 of the 10 reported in Table 1). We color the best method, and **bold** values within 95% of its value in the same row.

the same base 5 evaluation tasks used in Table 1, but construct intermediate waypoints along the

shortest path to the final goal determined by Breadth First Search. We also include an additional

335 maze environment, giant on which all methods have zero success rates to reach distant goals. This

336 gives a more holistic view on an agent's performance.

337 We display results in Figure 3 and Appendix E, where we 338 can see how performance drops off for all methods after a 339 generalization threshold denoted by the red bar. While all 340 methods cannot fully reach distant goals on giant, we 341 see that **BYOL-** γ has the slowest drop-off in performance. 342 We note that this is a challenging task, that requires stitch-

343 ing up to approximately 8 different trajectories.

344 5.4 Components

345 affecting generalization for representation learning



346 We ablate key components of the **BYOL**- γ objective in 347 Table 2. This includes removing action conditioning for 348 forward predictor ψ_f (-*a*), swapping the loss from cross-349 entropy to normalized squared l_2 norm (f_{l_2}), removing 350 backwards predictor ψ_b , and predicting the representation 351 of the adjacent state ($\gamma = 0$). Both removing action-



352 conditioning, and backwards prediction overall lead to similar results, but variability per-environment.

For f_{l2} , we obtain slightly worse average performance, and for $\gamma = 0$, we see the largest drop-off, especially on humanoidmaze environments.

355 6 Discussion

Limitations. While we demonstrate that BYOL- γ and other representation learning objectives offer a promising recipe for obtaining combinatorial generalization, we find that there still exists a generalization gap, especially on challenging navigation environments e.g. giant. We find a less significant improvement over BC on visual environments, which may motivate additional investigation. Additionally, we may anticipate more benefit from representation learning when applied to larger visual datasets, which has been fruitful in other domains.

362 **Conclusion.** In this work, we provide a stronger understanding of the relationship between quantities 363 related to successor representations and the generalization of policies trained with behavioral cloning. 364 We propose a new self-predictive representation learning objective, BYOL- γ , and show that it 365 captures information related to the successor measure, resulting in a competitive choice of an 366 auxiliary loss for better generalization. We demonstrate that augmenting behavior cloning with 367 meaningful representations results in new capabilities such as improved combinatorial generalization, 368 especially in larger and more complex environments.

369 **References**

André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt,
 and David Silver. Successor features for transfer in reinforcement learning. *Advances in neural*

- *information processing systems*, 30, 2017. URL https://arxiv.org/abs/1606.05312.
- Léonard Blier, Corentin Tallec, and Yann Ollivier. Learning successor states and goal-dependent
 values: A mathematical viewpoint, 2021. URL https://arxiv.org/abs/2101.07123.
- David Brandfonbrener, Alberto Bietti, Jacob Buckman, Romain Laroche, and Joan Bruna. When does
 return-conditioned supervised learning work for offline reinforcement learning? In Alice H. Oh,
- 377 Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information*
- 378 *Processing Systems*, 2022. URL https://openreview.net/forum?id=XByg4kotW5.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jegou, Julien Mairal, Piotr Bojanowski, and
 Armand Joulin. Emerging properties in self-supervised vision transformers. In 2021 IEEE/CVF
 International Conference on Computer Vision (ICCV), pages 9630–9640, 2021. doi: 10.1109/
 ICCV48922.2021.00951.
- Wilka Carvalho, Momchil S. Tomov, William de Cothi, Caswell Barry, and Samuel J. Gershman.
 Predictive representations: Building blocks of intelligence. *Neural Computation*, 36(11):2225–2298, 10 2024. ISSN 0899-7667. doi: 10.1162/neco_a_01705. URL https://doi.org/10.
 1162/neco_a_01705.
- Yash Chandak, Shantanu Thakoor, Zhaohan Daniel Guo, Yunhao Tang, Remi Munos, Will Dabney,
 and Diana L Borsa. Representations and exploration for deep reinforcement learning using singular
 value decomposition. In *International Conference on Machine Learning*, pages 4009–4034. PMLR,
 2023. URL https://arxiv.org/abs/2305.00654.
- Ian Char, Viraj Mehta, Adam Villaflor, John M. Dolan, and Jeff Schneider. Bats: Best action trajectory
 stitching, 2022. URL https://arxiv.org/abs/2204.12026.
- Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel,
 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
 modeling. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021. URL https://openreview.net/forum?
 id=a7APmM4B9d.
- Zichen Jeff Cui, Hengkai Pan, Aadhithya Iyer, Siddhant Haldar, and Lerrel Pinto. Dynamo: In-domain
 dynamics pretraining for visuo-motor control. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://arxiv.org/abs/2409.12192.
- Peter Dayan. Improving generalization for temporal difference learning: The successor representation.
 Neural Computation, 5(4):613–624, 1993. doi: 10.1162/neco.1993.5.4.613.
- Scott Emmons, Benjamin Eysenbach, Ilya Kostrikov, and Sergey Levine. Rvs: What is essential for
 offline RL via supervised learning? In *International Conference on Learning Representations*,
 2022. URL https://openreview.net/forum?id=S874XAIpkR-.
- Benjamin Eysenbach, Tianjun Zhang, Sergey Levine, and Russ R Salakhutdinov. Contrastive learning
 as goal-conditioned reinforcement learning. *Advances in Neural Information Processing Systems*,
 35:35603–35620, 2022.
- Jesse Farebrother, Joshua Greaves, Rishabh Agarwal, Charline Le Lan, Ross Goroshin, Pablo Samuel
 Castro, and Marc G Bellemare. Proto-value networks: Scaling representation learning with
 auxiliary tasks. In *The Eleventh International Conference on Learning Representations*, 2023.
 URL https://openreview.net/forum?id=oGDKSt9JrZi.

413 Vincent François-Lavet, Yoshua Bengio, Doina Precup, and Joelle Pineau. Combined reinforcement

414 learning via abstract representations. In *Proceedings of the AAAI Conference on Artificial Intelli-*

415 gence, volume 33, pages 3582-3589, 2019. URL https://arxiv.org/abs/1809.04506.

416 Quentin Garrido, Mahmoud Assran, Nicolas Ballas, Adrien Bardes, Laurent Najman, and Yann
 417 LeCun. Learning and leveraging world models in visual representation learning, 2024. URL
 418 https://arxiv.org/abs/2403.00504.

Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. Deepmdp:
Learning continuous latent space models for representation learning. In *International conference on machine learning*, pages 2170–2179. PMLR, 2019. URL https://arxiv.org/abs/
1906.02736.

Dibya Ghosh, Homer Rich Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey
Hejna, Tobias Kreiman, Charles Xu, Jianlan Luo, You Liang Tan, Lawrence Yunliang Chen,
Quan Vuong, Ted Xiao, Pannag R. Sanketi, Dorsa Sadigh, Chelsea Finn, and Sergey Levine.
Octo: An open-source generalist robot policy. In *Robotics: Science and Systems*, 2024. URL
https://doi.org/10.15607/RSS.2024.XX.090.

Raj Ghugare, Matthieu Geist, Glen Berseth, and Benjamin Eysenbach. Closing the gap between TD
learning and supervised learning - a generalisation point of view. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://arxiv.org/abs/2401.
11237.

Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena
Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,
et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. URL https://arxiv.org/abs/
2006.07733.

Zhaohan Daniel Guo, Bernardo Avila Pires, Bilal Piot, Jean-Bastien Grill, Florent Altché, Remi
Munos, and Mohammad Gheshlaghi Azar. Bootstrap latent-predictive representations for multitask
reinforcement learning. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3875–3886. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.
press/v119/guo20g.html.

Nicklas Hansen, Xiaolong Wang, and Hao Su. Temporal difference learning for model predictive
control. In *International Conference on Machine Learning (ICML)*, 2022. URL https://
arxiv.org/abs/2203.04955.

Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth 446 Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, 447 448 Peter David Fagan, Joey Hejna, Masha Itkina, Marion Lepert, Yecheng Jason Ma, Patrick Tree Miller, Jimmy Wu, Suneel Belkhale, Shivin Dass, Huy Ha, Arhan Jain, Abraham Lee, Youngwoon 449 450 Lee, Marius Memmel, Sungjae Park, Ilija Radosavovic, Kaiyuan Wang, Albert Zhan, Kevin Black, 451 Cheng Chi, Kyle Beltran Hatch, Shan Lin, Jingpei Lu, Jean Mercat, Abdul Rehman, Pannag R Sanketi, Archit Sharma, Cody Simpson, Quan Vuong, Homer Rich Walke, Blake Wulfe, Ted 452 Xiao, Jonathan Heewon Yang, Arefeh Yavary, Tony Z. Zhao, Christopher Agia, Rohan Baijal, 453 Mateo Guaman Castro, Daphne Chen, Qiuyu Chen, Trinity Chung, Jaimyn Drake, Ethan Paul 454 455 Foster, Jensen Gao, David Antonio Herrera, Minho Heo, Kyle Hsu, Jiaheng Hu, Donovon Jackson, 456 Charlotte Le, Yunshuang Li, Xinyu Lin, Zehan Ma, Abhiram Maddukuri, Suvir Mirchandani, 457 Daniel Morton, Tony Khuong Nguyen, Abigail O'Neill, Rosario Scalise, Derick Seale, Victor 458 Son, Stephen Tian, Emi Tran, Andrew E. Wang, Yilin Wu, Annie Xie, Jingyun Yang, Patrick 459 Yin, Yunchu Zhang, Osbert Bastani, Glen Berseth, Jeannette Bohg, Ken Goldberg, Abhinav Gupta, Abhishek Gupta, Dinesh Jayaraman, Joseph J Lim, Jitendra Malik, Roberto Martín-460 461 Martín, Subramanian Ramamoorthy, Dorsa Sadigh, Shuran Song, Jiajun Wu, Michael C. Yip, 462 Yuke Zhu, Thomas Kollar, Sergey Levine, and Chelsea Finn. DROID: A large-scale in-the-wild

- robot manipulation dataset. In RSS 2024 Workshop: Data Generation for Robotics, 2024. URL
 https://openreview.net/forum?id=Ml2pTYLNLi.
- Khimya Khetarpal, Zhaohan Daniel Guo, Bernardo Avila Pires, Yunhao Tang, Clare Lyle, Mark
 Rowland, Nicolas Heess, Diana L Borsa, Arthur Guez, and Will Dabney. A unifying framework for
 action-conditional self-predictive reinforcement learning. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025. URL https://arxiv.org/abs/2406.02035.
- Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael
 Rafailov, Ethan P Foster, Pannag R Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel,
- Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. OpenVLA: An
- open-source vision-language-action model. In 8th Annual Conference on Robot Learning, 2024.
- 473 URL https://arxiv.org/abs/2406.09246.
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit
 q-learning. In *International Conference on Learning Representations*, 2022. URL https:
 //arxiv.org/abs/2110.06169.
- Charline Le Lan, Stephen Tu, Adam Oberman, Rishabh Agarwal, and Marc G. Bellemare. On the
 generalization of representations in reinforcement learning, 2022. URL https://arxiv.org/
 abs/2203.00543.
- 480 Yann LeCun. A path towards autonomous machine intelligence version, 2022. URL https: 481 //openreview.net/forum?id=BZ5alr-kVsf.

Jaewoo Lee, Sujin Yun, Taeyoung Yun, and Jinkyoo Park. GTA: Generative trajectory augmentation
with guidance for offline reinforcement learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?
id=kZpNDbZrzy.

486 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
487 review, and perspectives on open problems, 2020. URL https://arxiv.org/abs/2005.
488 01643.

Bo Liu, Yihao Feng, Qiang Liu, and Peter Stone. Metric residual network for sample efficient
goal-conditioned reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 8799–8806, 2023. URL https://arxiv.org/abs/2208.
08133.

- 493 Cong Lu, Philip J. Ball, Yee Whye Teh, and Jack Parker-Holder. Synthetic experience replay.
 494 In Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL https:
 495 //openreview.net/forum?id=6jNQ1AY1Uf.
- Yunhao Luo, Utkarsh A. Mishra, Yilun Du, and Danfei Xu. Generative trajectory stitching through
 diffusion composition, 2025. URL https://arxiv.org/abs/2503.05153.

Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Yecheng Jason Ma, Claire Chen, Sneha Silwal,
Aryan Jain, Vincent-Pierre Berges, Tingfan Wu, Jay Vakil, Pieter Abbeel, Jitendra Malik, Dhruv
Batra, Yixin Lin, Oleksandr Maksymets, Aravind Rajeswaran, and Franziska Meier. Where are we
in the search for an artificial visual cortex for embodied intelligence? In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://arxiv.org/abs/2303.
18240.

Vivek Myers, Chongyi Zheng, Anca Dragan, Sergey Levine, and Benjamin Eysenbach. Learning
 temporal distances: Contrastive successor features can provide a metric structure for decision making. In *Forty-first International Conference on Machine Learning*, 2024. URL https:
 //openreview.net/forum?id=xQiYCmDrjp.

508 Vivek Myers, Catherine Ji, and Benjamin Eysenbach. Horizon Generalization in Reinforcement

Learning. In International Conference on Learning Representations, January 2025a. URL
 https://arxiv.org/pdf/2501.02709.

Vivek Myers, Bill Chunyuan Zheng, Anca Dragan, Kuan Fang, and Sergey Levine. Temporal
 representation alignment: Successor features enable emergent compositionality in robot instruction
 following, 2025b. URL https://arxiv.org/abs/2502.05454.

Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhi Gupta. R3m: A universal
 visual representation for robot manipulation. In *Conference on Robot Learning*, 2022. URL
 https://arxiv.org/abs/2203.12601.

Tianwei Ni, Benjamin Eysenbach, Erfan Seyedsalehi, Michel Ma, Clement Gehring, Aditya Mahajan,
 and Pierre-Luc Bacon. Bridging state and history representations: Understanding self-predictive
 rl. In *The Twelfth International Conference on Learning Representations*, 2024. URL https:
 //arxiv.org/abs/2401.08898.

521 Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham 522 Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex 523 Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Anikait Singh, 524 Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, 525 Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan Wahid, Ben Burgess-Limerick, Beomjoon 526 Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea 527 Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher 528 Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne 529 Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa 530 Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Freek 531 Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, 532 Glen Berseth, Gregory Kahn, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui 533 Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung 534 535 Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan 536 Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, 537 538 Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim, 539 Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana 540 Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin 541 Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan 542 Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, 543 Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi Jim Fan, Lionel Ott, Lisa Lee, 544 Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, 545 Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong 546 Zhang, Mingyu Ding, Minho Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki 547 Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman 548 Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R 549 Sanketi, Patrick Tree Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaresan, Qiuyu Chen, Quan Vuong, Rafael 550 Rafailov, Ran Tian, Ria Doshi, Roberto Martín-Martín, Rohan Baijal, Rosario Scalise, Rose 551 552 Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, 553 Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar 554 Bahl, Shivin Dass, Shubham Sonawani, Shuran Song, Sichun Xu, Siddhant Haldar, Siddharth 555 Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, 556 Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj 557 Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted 558 Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, 559 Trevor Darrell, Trinity Chung, Vidhi Jain, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram

Burgard, Xi Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei,
Xuanlin Li, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu,
Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yoonyoung Cho, Youngwoon Lee, Yuchen
Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang
Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen
Zhang, and Zipeng Lin. Open x-embodiment: Robotic learning datasets and rt-x models : Open
x-embodiment collaboration0. In 2024 IEEE International Conference on Robotics and Automation
(ICRA), pages 6802–6003–2024. doi: 10.1100/ICRA57147.2024.10611477

- 567 (*ICRA*), pages 6892–6903, 2024. doi: 10.1109/ICRA57147.2024.10611477.
- Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking
 offline goal-conditioned rl. In *International Conference on Learning Representations (ICLR)*, 2025.
 URL https://arxiv.org/abs/2007.05929.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
 Learning transferable visual models from natural language supervision. In *International Conference*
- on Machine Learning, 2021. URL https://arxiv.org/abs/2103.00020.
- Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell.
 Real-world robot learning with masked visual pre-training. In *6th Annual Conference on Robot Learning*, 2022. URL https://arxiv.org/abs/2203.06173.
- Juergen Schmidhuber. Reinforcement learning upside down: Don't predict rewards just map them
 to actions, 2020. URL https://arxiv.org/abs/1912.02875.
- Max Schwarzer, Ankesh Anand, Rishab Goel, R. Devon Hjelm, Aaron C. Courville, and Philip
 Bachman. Data-efficient reinforcement learning with self-predictive representations. In Inter *national Conference on Learning Representations*, 2020. URL https://arxiv.org/abs/
 2007.05929.
- Vlad Sobal, Wancong Zhang, Kynghyun Cho, Randall Balestriero, Tim G. J. Rudner, and Yann
 LeCun. Learning from reward-free offline data: A case for planning with latent dynamics models,
 2025. URL https://arxiv.org/abs/2502.14819.
- Yunhao Tang, Zhaohan Daniel Guo, Pierre H. Richemond, Bernardo Ávila Pires, Yash Chandak,
 Rémi Munos, Mark Rowland, Mohammad Gheshlaghi Azar, Charline Le Lan, Clare Lyle, Andr'as
 Gyorgy, Shantanu Thakoor, Will Dabney, Bilal Piot, Daniele Calandriello, and M. Vaíko. Understanding self-predictive learning for reinforcement learning. In *International Conference on Machine Learning*, 2022. URL https://arxiv.org/abs/2212.03319.
- Ahmed Touati, Jérémy Rapin, and Yann Ollivier. Does zero-shot reinforcement learning exist?
 In *The Eleventh International Conference on Learning Representations*, 2023. URL https:
 //openreview.net/forum?id=MYEap_OcQI.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive
 coding, 2019. URL https://arxiv.org/abs/1807.03748.
- 597 Claas A. Voelcker, Tyler Kastner, Igor Gilitschenski, and Amir-massoud Farahmand. When does
 598 self-prediction help? understanding auxiliary tasks in reinforcement learning. *Reinforcement* 599 *Learning Conference*, August 2024. URL https://arxiv.org/abs/2406.17718.
- Tongzhou Wang and Phillip Isola. Improved representation of asymmetrical distances with interval
 quasimetric embeddings. In *NeurIPS 2022 Workshop on Symmetry and Geometry in Neural Representations*, 2022. URL https://arxiv.org/abs/2211.15120.
- Tongzhou Wang, Antonio Torralba, Phillip Isola, and Amy Zhang. Optimal goal-reaching reinforce ment learning via quasimetric learning. In *International Conference on Machine Learning*. PMLR,
 2023. URL https://arxiv.org/abs/2304.01203.

- 606 Taku Yamagata, Ahmed Khalil, and Raúl Santos-Rodríguez. Q-learning decision transformer:
- 607 leveraging dynamic programming for conditional sequence modelling in offline rl. In *Proceedings*
- 608 of the 40th International Conference on Machine Learning, ICML'23. JMLR.org, 2023.
- 609 Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games
- 610 with limited data. Advances in neural information processing systems, 34:25476–25488, 2021.
- 611 URL https://arxiv.org/abs/2111.00210.
- 612 Zhaoyi Zhou, Chuning Zhu, Runlong Zhou, Qiwen Cui, Abhishek Gupta, and Simon Shaolei Du. Free
- from bellman completeness: Trajectory stitching via model-based return-conditioned supervised
- 614 learning. In The Twelfth International Conference on Learning Representations, 2024. URL
- 615 https://arxiv.org/abs/2310.19308.

616 A Experimental Setup

Hyperparameter		Shared			
actor head	MLF	MLP (512,512,512)			
representation encoder (ϕ)	MLP (64,64,64)				
predictor (ψ)	M	MLP (64,64,64)			
encoder ensemble		2			
learning rate		3×10^{-4}			
optimizer		Adam			
	Non-visual	Visual			
Gradient steps	1000k	500k			
Batch size	1024	256			
au (EMA)	1.0	0.99			
γ	0.99	{0.66, 0.99}			
α (alignment)	{1,6,40,100}	{1,6,10,20}			
additional encoder	n/a	impala_small			
encoder output dimension	s	64			

Table 3: Hyperparameters for BYOL- γ

617 A.1 Implementation Details

618 In this section we provide more training details for BYOL- γ , and representation learning baselines.

619 We match the training details of OGBench, including gradient steps, batch size, learning rate.

620 **Network Architecture.** We follow the same network architecture setup as TRA, where we utilize 621 MLP-based encoders, and action head. For the output dimension of the encoder, we use the state 622 dimension for non-visual experiments, and 64 for visual experiments. For the predictor ψ , we utilize 623 an MLP of the same architecture as the encoder. For image-based tasks, there is an additional CNN, 624 which then passes output to the MLP encoder.

625 **Representation Ensemble.** We follow the setup of TRA which utilizes representation ensembling, 626 such that two copies of the encoder ϕ_1 , ϕ_2 are in parallel. We also have two distinct predictors ψ_1 , ψ_2 627 for each ensemble. As input to the policy head, we average the representations, $\bar{z} = \frac{\phi_1(s_t) + \phi_2(s_2)}{2}$. 628 Each representation is trained independently for the BYOL loss, but the BC loss differentiates through 629 both ϕ_8 .

630 **Alignment.** We find that the choice of weight of the auxiliary loss for the representation learning 631 objective is sensitive to both the robot embodiment and the environment size. For comparison, we 632 perform a hyperparameter search over four alignment values for BYOL- γ , TRA, and FB, and then 633 report the best value for each environment in Table 1.

634 **Discount.** For sampling the next-state, we utilize a discount factor of $\gamma = 0.99$ for all non-visual 635 environments. For visual environments, we perform a hyperparameter search over {0.66, 0.99}, 636 however all representation learning methods performed better at $\gamma = 0.66$.

637 **A.2 BYOL-**γ

Target network. For BYOL, we find that exponential moving average (EMA) target networks for the encoder ϕ are not necessary for non-visual environments ($\tau = 1$), but for visual environments,

we find that a fast target stabilizes training ($\tau = 0.99$):

$$\phi_{\text{target}} = \tau \phi_{\text{online}} + (1 - \tau) \phi_{\text{target}}$$

638 A.3 TRA

639 In practice, TRA uses a symmetric version (Radford et al., 2021) of the InfoNCE objective discussed 640 in Equation 3. We write this in batch form, $\mathcal{B} = \{(s_i, s_{+,i})\}_{i=1}^{|\mathcal{B}|}$ rather than in expectation:

$$\mathcal{L}_{\text{TRA}} = \mathbb{E}_{\mathcal{B}} \left[-\frac{1}{B} \sum_{i=1}^{|\mathcal{B}|} \log \frac{e^{f(\psi(s_i),\phi(s_{+,i}))}}{\sum_{j=1}^{|\mathcal{B}|} e^{f(\psi(s_i),\phi(s_{+,j}))}} - \frac{1}{B} \sum_{i=1}^{|\mathcal{B}|} \log \frac{e^{f(\psi(s_i),\phi(s_{+,i}))}}{\sum_{j=1}^{|\mathcal{B}|} e^{f(\psi(s_j),\phi(s_{+,i}))}} \right]$$
(16)

641 Additionally, TRA minimizes the squared norm of representations $\min_{\phi,\psi} \lambda \mathbb{E}_s \left[\frac{\|\phi(s)\|^2}{d} + \frac{\|\psi(s)\|^2}{d} \right]$ 642 with $\lambda = 10^{-6}$. For TRA, we search over $\alpha = \{10, 40, 60, 100\}$.

643 A.4 FB

Prior work which trains FB for zero-shot policy optimization (Touati et al., 2023) typically normalizes ϕ with an additional loss term so that $\mathbb{E} \left[\phi \phi^T \right] \approx I_d$. However, we found that adding this loss term was not beneficial to performance in our setting and hence do not include it.

FB uses an EMA target network as described in A.2 with $\tau = 0.005$. For FB, we search over $\alpha = \{0.01, 0.05, 0.001, 0.005\}.$

649 A.5 Code.

We utilize the OGBench (Park et al., 2025) codebase and benchmark, and its extensions in the TRA codebase (Myers et al., 2025b) for equal comparison.

652 A.6 Compute Requirements

653 We perform all experiments utilizing single GPUs, predominately NVIDIA RTXA8000 and L40S.

654 We utilize 6 CPU cores, 24G of RAM for non-visual environments, and 64G for visual experiments.

Experiments take 2 to 4 hours for non-visual and 6 to 12 hours for visual environments.

656 **B** Ablations.

657 B.1 Action-conditioning

In this section, we ablate the component of performing action-conditioning for the predictor $\psi(s_t)$ vs $\psi(s_t, a_t)$ for TRA and FB. We consider a similar comparison for BYOL- γ in Table 2. For this comparison, when we perform action-conditioning, we utilize a policy representation $\pi(s = \phi(s), g = \phi(g))$, and otherwise $\pi(s = \psi(s), g = \phi(g))$. We find that results can be environment specific. On average, results are not improved for TRA, but we find an improvement for FB, hence in our main Table 1 we include the action-conditioned results for FB and the action-free results for TRA.

Dataset	TRA	TRA ^a	FB	FB ^a
antmaze-medium-stitch antmaze-large-stitch humanoidmaze-medium-stitch humanoidmaze-large-stitch antsoccer-arena-stitch	$\begin{array}{c} 54\pm 6 \\ 11\pm 8 \\ \textbf{45}\pm 8 \\ 5\pm 4 \\ 14\pm 4 \end{array}$	$\begin{array}{c} 57 \pm 12 \\ 7 \pm 7 \\ \textbf{45} \pm 5 \\ 9 \pm 4 \\ \textbf{25} \pm 8 \end{array}$	$\begin{array}{c} {\bf 64} \pm 10 \\ {\bf 17} \pm 6 \\ {\bf 36} \pm 3 \\ {\bf 6} \pm 2 \\ {\bf 17} \pm 5 \end{array}$	$\begin{array}{c} {\bf 64} \pm 6 \\ {\bf 23} \pm 4 \\ {42} \pm 4 \\ {\bf 11} \pm 3 \\ {22} \pm 10 \end{array}$
visual-antmaze-medium-stitch visual-antmaze-large-stitch visual-scene-play	$\begin{array}{c} {\bf 52} \pm 3 \\ {\bf 17} \pm 1 \\ {\bf 16} \pm 3 \end{array}$	$\begin{array}{c} 33\pm4\\ 22\pm5\\ \textbf{18}\pm2 \end{array}$	$\begin{array}{c} 47\pm5\\ \textbf{28}\pm3\\ 12\pm2 \end{array}$	$\begin{array}{c} {\bf 49} \pm 2 \\ {\bf 29} \pm 2 \\ {\bf 14} \pm 1 \end{array}$
average	27	27	28	32

Table 4: Action-conditioning ablations. We ablate the choice to condition on the first action for predictor ψ for TRA and FB over 10 seeds for non-visual and 4 seeds for visual environments.

665 C CL to FB

- Here, illustrate that connection between CL and FB, showing that in the limit an n-step version of FBbecomes similar to CL.
- 668 We can rewrite Equation (3) to see the connection between FB (TD) and CL (MC). Under assumptions
- 669 that ϕ, ψ are centered ($\mathbb{E}[\phi] = \mathbb{E}[\psi] = 0$), and unit normalized $\|\phi\|_2, \|\psi\|_2 = 1$, if we apply a second-
- order Taylor expansion to the denominator of the CL loss (Touati et al., 2023) we have:

$$CL_{InfoNCE} \approx \frac{1}{2} \mathbb{E}_{s \sim p_0(s), s' \sim p_0(s')} \left[(\psi(s)^T \phi(s'))^2 \right] - 2\mathbb{E}_{\substack{k \sim \text{geom}(1-\gamma)\\s_t \sim p_0, s_{t+k} \sim p^\pi(s_{t+k}|s_t)}} \left[\psi(s_t)^T \phi(s_{t+k}) \right]$$
(17)

Next, we can consider an n-step variant of the FB loss (Blier et al., 2021) which we refer to as FB(n):

$$\min_{\phi,\psi} \mathbb{E}_{\substack{s_t \sim p_0 \\ s' \sim p_0}} \left[(\psi(s_t)^T \phi(s') - \gamma^n \bar{\psi}(s_{t+n})^T \bar{\phi}(s'))^2 \right] - 2 \sum_{i=1}^n \mathbb{E}_{s_t \sim p_0, s_{t+i} \sim p^\pi} \left[\gamma^i \psi(s_t)^T \phi(s_{t+i}) \right]$$
(18)

672 We can make the full connection to CL with infinite horizon n:

$$FB_{n\to\infty}(n) = \mathbb{E}_{\substack{s_t \sim p_0 \\ s' \sim p_0}} \left[(\psi(s_t)^T \phi(s'))^2 \right] - 2 \sum_{i=1}^n \mathbb{E}_{\substack{s_t \sim p_0 \\ s_{t+i} \sim p^\pi(s_{t+i}s_0)}} \left[\gamma^i \psi(s_t)^T \phi(s_{t+i}) \right]$$
(19)

$$= \mathbb{E}_{\substack{s_{t} \sim p_{0} \\ s_{t}^{\prime} \sim p_{0}}} \left[(\psi(s_{t})^{T} \phi(s^{\prime}))^{2} \right] - \frac{2\gamma}{(1-\gamma)} \sum_{i=1}^{n} \mathbb{E}_{\substack{s_{t} \sim p^{\pi}(s_{t+i}s_{0})}} \left[(1-\gamma)\gamma^{i-1} \psi(s_{t})^{T} \phi(s_{t+i}) \right]$$
(20)

$$= \mathbb{E}_{\substack{s_t \sim p_0 \\ s' \sim p_0}} \left[(\psi(s_t)^T \phi(s'))^2 \right] - \frac{2\gamma}{(1-\gamma)} \mathbb{E}_{\substack{k \sim \text{geom}(1-\gamma) \\ s_t \sim p_0, s_{t+k} \sim p^\pi(s_{t+k}|s_t)}} \left[\psi(s_t)^T \phi(s_{t+k}) \right]$$
(21)

673 Thus, we can see that in the infinite horizon form of FB(n), it is related to the form of CL_{InfoNCE} in 674 (3), but with the positive contrastive term weighted by factor $\frac{\gamma}{1-\gamma}$.

675 **D** Finite MDP

676 **D.1 BYOL**

677 **BYOL as an Ordinary Differential Equation (ODE)** In finite MDPs, we can characterize the 678 BYOL objective which gives intuition about what information is captured in ϕ , ψ , and conditions 679 that may be useful for stability (Tang et al., 2022; Khetarpal et al., 2025). Consider a finite MDP

with transition P^{π} , linear d-dimensional encoder $\Phi \in \mathbb{R}^{|S| \times d}$, and linear action-free latent-dynamics 680 $\Psi \in \mathbb{R}^{d \times d}$. In a finite MDP, Equation (4) becomes: 681

$$\min_{\Phi,\Psi} \operatorname{BYOL}(\Phi,\Psi) := \min_{\Phi,\Psi} \mathbb{E}_{s_t \sim p(s), s_{t+1} \sim P^{\pi}}, \left[\|\psi^T \Phi^T s_t - \bar{\Phi}^T s_{t+1}\|_2^2 \right]$$
(22)

A property to prevent this objective from collapsing is that Ψ is updated more quickly than Φ . In 682

683 practice, this is commonly realized as the dynamics are generally a smaller network than the encoder.

This system can be analyzed in an ideal setup, where we first find the optimal Ψ , each time before 684 685 taking a gradient step for Φ , which leads to the ODE for representations Φ (Tang et al., 2022):

$$\Psi^* \in \arg\min_{\Psi} \mathsf{BYOL}(\Phi, \Psi), \quad \dot{\Phi} = -\nabla_{\Phi} \mathsf{BYOL}(\Phi, \Psi)|_{\Psi = \Psi^*}$$
(23)

- We are able to analyze this ODE with the following assumptions (Tang et al., 2022): 686
- Assumption D.1 (Orthogonal initialization). $\Phi^{\top} \Phi = I$ 687
- Assumption D.2 (Uniform state distribution). $p_0(s) = \frac{1}{|S|}$ Assumption D.3 (Symmetric dynamics). $P^{\pi} = (P^{\pi})^{\top}$ 688
- 689

Under these three assumptions, Khetarpal et al. (2025) prove that the BYOL ODE is equivalent to 690 691 monotonically minimizing the surrogate objective:

$$\min_{\Psi} \|P^{\pi} - \Phi \Psi \Phi^T\|_F + C \tag{24}$$

Where $\|\cdot\|_F$ is the Frobenius matrix norm. Thus, we can understand that the BYOL objective as 692 learning a d-rank decomposition of the underlying dynamics P^{π} . Additionally, the top d eigenvectors 693 694 of P^{π} match those of $(I - \gamma P^{\pi})^{-1} = M^{\pi}$ Chandak et al. (2023). However, we will highlight that 695 there are key differences when *learning* a low-rank decomposition between P^{π} and M^{π} . This is described by Touati et al. (2023), where we can consider that in a real-world problem with underlying 696 continuous-time dynamics, actions have little effect, and P^{π} is close to the identity, e.g. close to 697 full-rank. However, M^{π} , which takes powers of P_{π}^{t} , has a "sharpening effect" on the difference 698 699 between eigenvalues, which gives a clearer learning signal. This is intuitive on a real-world problem 700 like robotics, even with discrete-time dynamics, where $s_{t+1} \approx s_t$, but we have larger differences 701 between s_t and s_{t+k} .

702 **D.2** BYOL- γ

703 In the finite MDP, we now verify theorem 4.1, where BYOL- γ approximates the successor repre-704 sentation with matrix decomposition $M^{\pi} \approx \Phi \Psi \Phi^{T}$.

We consider the same objective (22), where we need to update the expectation of the sampling 705 706 distribution:

$$\min_{\Phi,\Psi} \text{BYOL-}\gamma(\Phi,\Psi) := \min_{\Phi,\Psi} \mathbb{E}_{s_t \sim p(s), s_+ \sim \tilde{M}^{\pi}}, \left[\|\psi^T \Phi^T s_t - \bar{\Phi}^T s_+\|_2^2 \right]$$
(25)

Assuming that this objective is optimized under the ODE (23). We have that our objective monotoni-707 cally minimizes: 708

$$\min_{\Psi} \|\tilde{M}^{\pi} - \Phi \Psi \Phi^T\|_F + C \tag{26}$$

This directly translates as we can consider $\tilde{M}^{\pi} = P^{\pi}$ as simply a valid transition matrix for a new, 709 temporally abstract, version of the original MDP. We maintain the original assumptions D.1, D.2, and 710 D.3. We do not need an additional assumption for \tilde{M}^{π} , as assumption D.3 for symmetric P^{π} implies 711 a symmetric $M^{\pi} = (1 - \gamma) \sum_{t>0} \gamma^t P_{\pi}^t$, 712

Under this setup, we also have that $\Psi \Phi \in \mathbb{R}^{n \times d}$ relates to the successor feature matrix, where each 713 row $(\Psi\Phi)_i$ contains the vector $(1-\gamma)\psi^{\pi}(s_i)$: 714

$$(1 - \gamma)\psi^{\pi}(s_i) = \sum_{j} \tilde{M}^{\pi}(s_i, s_j)\phi(s_j)$$
(27)

$$= (\tilde{M}^{\pi} \Phi)_i \tag{28}$$

$$\approx (\Phi \Psi \Phi^T \Phi)_i$$

$$= (\Phi \Psi)_i$$
(29)
(30)

$$= (\Phi \Psi)_i \tag{30}$$

- 715 In other words, in the restricted finite MDP, where we minimize (26), we are simultaneously learning
- 716 successor features $\psi^{\pi} \approx \Psi \Phi$ and basis features Φ .

717 E Additional Results for Horizon Generalization







Figure 4: Evaluating Generalization with Increasing Horizons: The distances to the right of the red dotted line require combinatorial generalization. The maze maps show examples of how intermediate goals are selected along the optimal path.

- 718 We include additional results matching the setup in Section 5.3, for antmaze-medium, and
- 719 {humanoidmaze}-{medium, large, giant} in Figure 4. We can observe that BYOL- γ leads
- in performance as the distance between the start and goal grows when compared to other methods.