# Uncovering RL Integration in SSL Loss: Objective-Specific Implications for Data-Efficient RL

**Anonymous authors**
Paper under double-blind review

**Keywords:** Data Efficient RL, Self Predictive RL, Self Supervised Learning

## Summary

This paper presents a systematic analysis of the role of self-supervised learning (SSL) objectives and their modifications in data-efficient reinforcement learning. We investigate previously undocumented modifications in the Self-Predictive Representations (SPR) (Schwarzer et al., 2020) framework that significantly impact agent performance. We demonstrate that feature decorrelation-based SSL objectives can achieve comparable performance without relying on domain-specific modifications, and show that the impact of these modifications persists even in more advanced models.

By conducting extensive experiments on the Atari 100k benchmark and DeepMind Control Suite, we provide insights into how different SSL objectives and their modifications affect learning efficiency across diverse environments. Our findings reveal that the choice and adaptation of SSL objectives play a crucial role in achieving data efficiency in self-predictive reinforcement learning, with implications for the design of future algorithms in this space.

## Contribution(s)

1. We demonstrate that previously undocumented SSL modifications in SPR (Schwarzer et al., 2020) - terminal state masking and prioritized replay weighting - are crucial for performance, with their removal leading to an 18% decrease in IQM score on Atari 100k
   **Context:** These modifications were silently adopted by subsequent work (D'Oro et al., 2023; Nikishin et al., 2022; Schwarzer et al., 2023) and their impact was not previously analyzed

2. We show that the Barlow Twins SSL objective (Zbontar et al., 2021) can come within 5% of SPR's performance without using domain-specific modifications, and VICReg (Bardes et al., 2021) can match PlayVirtual's (Yu et al., 2021) performance in continuous control tasks.
   **Context:** Prior work on SSL in reinforcement learning relied heavily on problem-specific modifications to achieve strong performance (Schwarzer et al., 2020; D'Oro et al., 2023; Schwarzer et al., 2023).

3. We establish that the impact of SSL modifications remains proportionally consistent in more sophisticated models, with unmodified versions of SR-SPR and BBF showing similar relative performance degradation despite having base IQM scores 3x and 2x higher than SPR respectively.
   **Context:** Previous work on SR-SPR (D'Oro et al., 2023; Nikishin et al., 2022) and BBF (Schwarzer et al., 2023) did not investigate the role of these modifications in their improved performance.

# Uncovering RL Integration in SSL Loss: Objective-Specific Implications for Data-Efficient RL

**Anonymous authors**
Paper under double-blind review

## Abstract

In this study, we investigate the effect of SSL objective modifications within the SPR framework, focusing on specific adjustments such as terminal state masking and prioritized replay weighting, which were not explicitly addressed in the original design. While these modifications are specific to RL, they are not universally applicable across all RL algorithms. Therefore, we aim to assess their impact on performance and explore other SSL objectives that do not accommodate these adjustments like Barlow Twins and VICReg. We evaluate six SPR variants on the Atari 100k benchmark, including versions both with and without these modifications. Additionally, we test the performance of these objectives on the DeepMind Control Suite, where such modifications are absent. Our findings reveal that incorporating specific SSL modifications within SPR significantly enhances performance, and this influence extends to subsequent frameworks like SR-SPR and BBF, highlighting the critical importance of SSL objective selection and related adaptations in achieving data efficiency in self-predictive reinforcement learning.

## 1 Introduction

Self-supervised learning (SSL) has become increasingly popular in data-efficient reinforcement learning (RL) due to its benefits in enhancing both efficiency and performance (Schwarzer et al., 2023; Ye et al., 2021; Hafner et al., 2023; Srinivas et al., 2020; Tomar et al., 2021; Li et al., 2023; Cagatan & Akgun, 2023). However, the application of SSL methods is often problem/domain-specific to maximize the performance of the RL agents. Although this approach is rational given the nature of these methods, it raises questions about generalization and transferability.

One of the key challenges in Deep RL is understanding the factors driving performance improvements, whether through hyperparameter tuning or novel algorithmic approaches (Obando-Ceron et al., 2024). The lack of transparency in hyperparameter selection often causes issues while algorithmic innovations are usually well-documented. However, our study of different SSL objectives within the Self-Predictive Representations (SPR) framework (Schwarzer et al., 2020) revealed that the SSL loss used in SPR differs from what is described in the original publication and its following works (Nikishin et al., 2022; D'Oro et al., 2023; Schwarzer et al., 2023) built upon it. This motivated us to investigate the effects of the undocumented modifications and further evaluate additional SSL objectives.

Unlike conventional SSL methods in RL, which often follow vision pretraining approaches (Chen et al., 2020) and directly combine SSL and RL losses (Srinivas et al., 2020), SPR modifies the SSL loss before integrating it with the RL objective. To further clarify, SPR employs the BYOL/SimSiam (Grill et al., 2020; Chen & He, 2020) auxiliary objective and incorporates two algorithm-specific adjustments to the SSL objective: (i) masking SSL loss with a boolean non-terminal state matrix and (ii) applying prioritized replay weighting to the batch loss. Consequently, this poses an essential question: How do these modifications affect the base performance of SSL objectives in the RL agent,

and can they be effectively applied to other SSL techniques in the RL domain? In addition, could this be a recurring phenomenon across the following models (Nikishin et al., 2022; Schwarzer et al., 2023) that adopt SPR as their baseline?

Concurrently, a plethora of novel self-supervised representation learning objectives has emerged (Zbontar et al., 2021; Bardes et al., 2021; Ozsoy et al., 2022; Caron et al., 2021), demonstrating performance improvements beyond image pretraining (Lee et al., 2023b; Goulão & Oliveira, 2023; Zhou et al., 2022; Ömer Veysel Çağatan, 2024). These objectives, based on feature decorrelation, do not inherently support the modifications used in SPR because the loss is computed along the feature dimension instead of the batch dimension, which we detail in Section 4.

This divergence raises another important question: How do these alternative objectives perform relative to the original SPR without SSL modifications? This inquiry is particularly significant because the information required to modify SSL objectives may not always be available in the environment. Understanding the performance of these unmodified objectives could provide valuable insights into the generalizability and robustness of different SSL approaches in RL contexts. Towards this end, we incorporate Barlow Twins and VICReg SSL objectives within SPR.

In essence, we frame our investigation around the following questions:

1. **How do these modifications affect the performance of SPR, and do their impacts extend to SPR-based models such as SR-SPR and BBF? Additionally, how do these alternative objectives compare to the original SPR when no SSL modifications are implemented?**

    Our findings reveal that modifications to SSL significantly affect SPR performance, leading to an 18% decrease in IQM when these modifications are removed. Additionally, SR-SPR and BBF exhibit a similar decline in performance. Among these modifications, prioritized replay weighting stands out as the most influential. Notably, Barlow Twins achieves results comparable to those of the original SPR, while VICReg's performance aligns with that of prioritized replay weighting. This indicates that these problem-specific modifications can be mitigated by employing alternative SSL objectives.Overall, our results underscore the importance of SSL modifications in SPR, which persist in strong models that utilize SPR

2. **How effectively do these SSL objectives perform in an environment in which SPR modifications are not applicable?**

    To address this, we examine VICReg, Barlow Twins, and SPR (BYOL/SimSiam) within the DeepMind Control Suite, where the popular SAC agent does not utilize prioritized replay weighting and the environment lacks a terminal state. Unlike in the Atari 100k benchmark, our results show VICReg as the top performer, even outpacing PlayVirtual, a more sophisticated variant of SPR. Meanwhile, SPR and Barlow Twins exhibit comparable performance levels. These findings highlight that algorithms tailored for specific domains may not consistently excel across different problem sets. Therefore, transferability should be a key factor in the design of new Deep RL algorithms.

## 2 Related Work

Tomar et al. (2021) tackles a more challenging setting for representation learning within RL with background distractors, using a simple baseline approach that avoids metric-based learning, data augmentations, world-model learning, and contrastive learning. They analyze why previous methods may fail or perform similarly to the baseline in this tougher scenario and stress the importance of detailed benchmarks based on reward density, planning horizon, and task-irrelevant components. They propose new metrics for evaluating algorithms and advocate for a data-centric approach to better apply RL to real-world tasks.

Li et al. (2023) explore whether SSL can enhance online RL from pixel data. By extending the contrastive reinforcement learning framework (Srinivas et al., 2020) to jointly optimize SSL and RL losses, and experimenting with various SSL losses, they find that the current SSL approaches
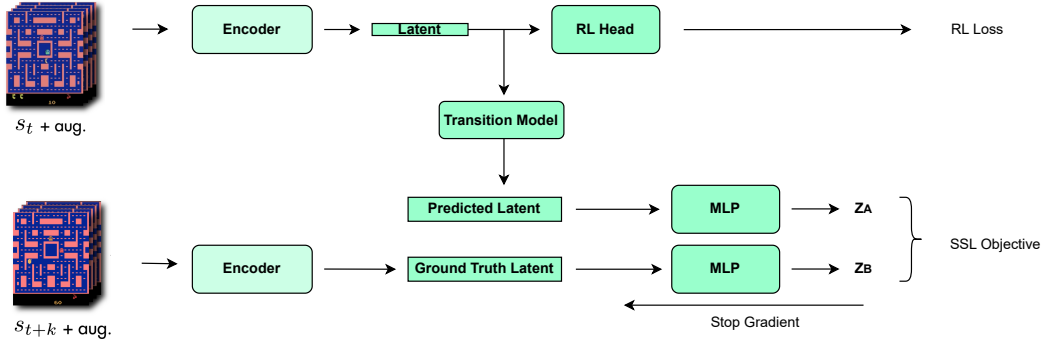
Figure 1: General flow diagram of SPR based methods. An encoder is used to create representations used for reinforcement learning and predicting future representations via a transition model and ground truth representations are created by the same encoder. MLPs differ when the predictor layer is used as in the case of BYOL/SimSiam. While we show the $k^{th}$ step here, the actual loss computation covers steps 1 to $K$. The SSL objective and RL loss changes between specific methods.

offer no significant improvement over baselines that use image augmentation alone, given the same data and augmentation. Even after evolutionary searches for optimal SSL loss combinations, these methods do not outperform carefully designed image augmentations. Their evaluation across various environments, including real-world robots, reveals that no single SSL loss or augmentation method consistently excels.

## 2.1 Data Efficient RL in Atari 100k

The introduction of the Atari 100k benchmark (Kaiser et al., 2019) has expedited the advancement of sample-efficient reinforcement learning algorithms. Model-based approach, SimPLe (Kaiser et al., 2019), outperformed Rainbow DQN (Hessel et al., 2017), showcasing superior performance. Building on Rainbow's framework, Hasselt et al. (2019) enhanced its efficacy through minor hyperparameter adjustments, resulting in Data-Efficient Rainbow (DER), which achieved a higher score compared to SimPLe.

DrQ (Kostrikov et al., 2020) employs a multi-augmentation strategy to regularize the value function during training of both Soft Actor-Critic (Haarnoja et al., 2018) and Deep Q-Network (Mnih et al., 2015). This approach effectively reduces overfitting and enhances training efficiency, leading to performance improvements for both algorithm

Several prevalent methods adopt the Atari 100k dataset, and these can be classified as follows: Model-Based (Hafner et al., 2023; Robine et al., 2023; Micheli et al., 2022; Ayton & Asai, 2021; Robine et al., 2021), Pretraining (Goulão & Oliveira, 2022; Schwarzer et al., 2021b; Lee et al., 2023a; Liu & Abbeel, 2021), Model-Free (Schwarzer et al., 2023; Huang et al., 2022; Nikishin et al., 2022; Cetin et al., 2022a; Lee et al., 2023a; Liang et al., 2022)

## 2.2 Representation Learning in Atari 100k

Cetin et al. (2022b) presents a deep reinforcement learning method using hyperbolic space for latent representations. Their innovative approach tackles optimization challenges in existing hyperbolic deep learning, ensuring stable end-to-end learning through deep hyperbolic representations.

Huang et al. (2022) proposes a Multiview Markov Decision Process (MMDP) with View-Consistent Dynamics (VCD), a method that enhances traditional MDPs by considering multiple state perspectives. VCD trains a latent space dynamics model for consistent state representations, achieved through data augmentation.

115 Srinivas et al. (2020) incorporate the InfoNCE (van den Oord et al., 2019) as an auxiliary component
116 within DER. Cagatan & Akgun (2023) uses Barlow Twins (Zbontar et al., 2021) instead of a con-
117 trastive objective to further improve results. This integration serves to enhance the learning process.
118 SPR (Schwarzer et al., 2020) outperforms all previous model-free approaches by predicting its latent
119 state representations multiple steps into the future with BYOL (Grill et al., 2020).

120 PlayVirtual (Yu et al., 2021) introduces a novel transition model as an alternative to the simplis-
121 tic module in SPR. The methodology enriches actual trajectories by incorporating a multitude of
122 cycle-consistent virtual trajectories. These virtual trajectories, generated using both forward and
123 backward dynamics models, collectively form a closed 'trajectory cycle.' The crucial aspect is en-
124 suring the consistency of this cycle, validating the projected states against real states and actions.
125 This approach significantly improves data efficiency by acquiring robust feature representations with
126 reduced reliance on real-world experiences. This method proves particularly advantageous for tasks
127 where obtaining real-world data is costly or challenging.

## 3   SPR

129 SPR is a performant data-efficient agent and a baseline of many other performant agents (Schwarzer
130 et al., 2023; Nikishin et al., 2022; D'Oro et al., 2023; Yu et al., 2021) and its general architecture
131 is depicted in Figure 1. The approach trains an agent by having it predict the latent state based on
132 the current state. It encodes the present state, forecasts the latent representation of the next state
133 using a transition model, and calculates loss by measuring the mean squared error between normal-
134 ized embeddings. Additionally, SPR adjusts its loss through terminal masking and prioritized replay
135 weighting. These two modifications inject RL-specific information into the auxiliary self-supervised
136 learning task. While the utilization of these ideas is not explicitly mentioned by Schwarzer et al.
137 (2020), it is possible that these techniques were considered self-evident and consequently were in-
138 cluded in their implementation  (Schwarzer et al., 2021a). We mention them here so as to be able to
139 better differentiate between SPR and other SPR variants.

140 SSL loss matrix in SPR denoted as $L$, encompasses negative cosine similarities between predicted
141 latent representations and ground truth latent representations, with dimensions of $B \times (K+1)$, where
142 $B$ is the batch size, and $K$ is the prediction horizon with 1 coming from the current observation. The
143 batch of interactions is drawn from the replay buffer, and their terminal status is known. The terminal
144 mask matrix, $M$, is composed of 0s and 1s denoting terminal and non-terminal states. The process
145 involves updating $L$ through a Hadamard product with $M$, denoted as $L \circ M$, effectively modifying
146 the loss matrix.

147 The loss matrix is divided into two components: SPR loss and Model SPR loss. SPR loss is between
148 the latent representations of the augmented views of the present state. Model SPR loss is between
149 the latent representations of the augmented views of the future states and the predicted future latent
150 representations, generated by the transition model. Model SPR is averaged across the temporal
151 dimension and as a result, both components have $N \times 1$ dimensionality.

152 The loss of each transition is multiplied by the prioritized replay weight, determined by the temporal
153 difference errors. Then the final loss is computed as the weighted sum of the average SPR loss and
154 half the average of the Model SPR loss across a batch as follows:

$$\mathcal{L}_{SPR} = \frac{1}{N} \sum_{i=1}^{N} \omega_i (\lambda \text{SPR}_i + \gamma \text{Model SPR}_i) \tag{1}$$

155 where $N$ is the batch size, $\omega_i$ is the priority weight ($\sum_i \omega_i = 1$), and $i$ indexes individual transitions,
156 where $\lambda, \gamma$ are hyperparameters.

|          | Median | IQM   | Mean   | Opt.Gap |
|----------|--------|-------|--------|---------|
| Barlow   | **0.324** | **0.320** | **0.605** | **0.593** |
| VICReg   | 0.281  | 0.289 | 0.600  | 0.610   |
| VICReg+Non | 0.221 | 0.279 | 0.554  | 0.617   |
| Barlow+Non | -0.009 | -0.011 | -0.171 | 1.171  |
| ZeroJump | 0.270  | 0.262 | 0.528  | 0.636   |

Table 1: Human-normalized aggregate metrics in Atari 100k. Scores were collected from 10 random runs.

|          | Median | IQM   | Mean   | Opt.Gap |
|----------|--------|-------|--------|---------|
| Stop-Grad | **0.271** | **0.303** | **0.615** | **0.577** |
| No Stop-Grad | 0.266 | 0.282 | 0.595 | 0.611 |

Table 2: Human-normalized aggregate metrics in Atari 100k by VICReg-High. Scores, collected from 10 random runs to assess the efficacy of including stop-gradient.

## 4  SPR-*

Despite variations in SSL objectives and RL algorithms across different benchmarks, the architecture remains largely consistent, as depicted in Figure 1. SPR employs a BYOL (Grill et al., 2020) objective with a momentum of 1, essentially adopting the SimSiam (Chen & He, 2020) approach. The primary architectural distinction lies in the inclusion of an extra predictor layer in the online MLP of BYOL or SimSiam to prevent collapse, a feature omitted in the original Barlow Twins and VICReg formulations as their objectives inherently mitigate the risk of collapse.

**SPR-Nakeds**  While SPR demonstrates considerable efficacy, the fundamental question remains unanswered—what is the impact of pure self-supervised learning and potential adaptations leading to SPR? Consequently, we introduce SPR-Naked, representing pure SSL. To assess the effects of prioritized replay weighting and terminal masking, we further establish SPR-Naked+Prio and SPR-Naked+Non, respectively.

In addition to the original SPR and its naked versions, we implement two additional types of agents with different SSL objectives.

**SPR-Barlow**  To extend the Barlow Twins to future predictions, we compute individual cross-correlation matrices for both the current and predicted latent representations at each time step. This results in a total of $K + 1$ matrices, each with dimensions $d \times d$, where $d$ denotes the embedding dimension within a single batch. Subsequently, we calculate the loss for each matrix and average the results. To make it easier to compare, we can define $\overline{\text{SPR}}$ Loss and $\overline{\text{Model SPR}}$ Loss analogously to their SPR counterparts, where the first is about the current state and the latter is about the future states. The final loss is then;

$$\mathcal{L}_{SPR-Barlow} = \overline{\text{SPR}} + \frac{1}{K} \sum_{k=1}^{K} \overline{\text{Model SPR}}_k \tag{2}$$

where $K$ is the number of predicted future observations.

**SPR-VICRegs**  We employ a parallel procedure as in Barlow Twins for VICReg. We introduce two variations of VICReg-High and VICReg-Low, featuring high or low covariance weights in the VICReg loss (Equation 11), while maintaining consistency in other hyperparameters. The primary objective is to observe the impact of feature decorrelation without inducing model collapse.

**Why not employ replay weighting and terminal state masking in Barlow/VICReg?**  The key limitation preventing the use of replay weighting or terminal masking in feature decorrelation-based methods lies in their reliance on covariance regularization. These methods employ either a cross-correlation matrix or a covariance matrix, both with dimensions matching the feature dimension. This structure prohibits applying the weighting of a feature dimension matrix using a batch dimension matrix. Consequently, these methods produce a unified loss for the entire batch, unlike approaches such as BYOL or SimSiam, which generate losses on a per-sample basis.

**Why use stop-gradient in Barlow/VICReg?**  Barlow Twins and VICReg effectively prevent collapse without resorting to symmetry-breaking architectural techniques such as predictor layers or stop-gradient mechanisms. While not strictly necessary in this scenario, we choose to include a

stop-gradient due to its empirically observed performance improvement, as depicted in Table 2. A more grounded reason stems from the architectural asymmetry introduced by the transition model. In the absence of a stop-gradient, gradients from the encoder's upper branch flow through the transition model, whereas gradients from the lower branch directly influence the encoder. This asymmetry can potentially lead to suboptimal encoder updates. Despite collapse avoidance in both cases, the inclusion of a stop-gradient is maintained for its superior performance outcomes.

**Why not other objectives?** Even though there are newly proposed SSL objectives (Silva et al., 2024; Zhang et al., 2024; Weng et al., 2024), it is impractical to include all objectives in experiments due to limited computational resources and the need to prioritize rigorous evaluation to draw precise conclusions however, we attempt to cover the two main families of SSL methods within SPR. The first is self-distillation, represented by BYOL (Grill et al., 2020) or SimSiam (Chen & He, 2020), which are already incorporated into SPR. The second family includes canonical correlation methods, such as VICReg and Barlow. Another category is Deep Metric Learning, which includes contrastive learning variants (Balestriero et al., 2023). However, we do not separately test contrastive objectives, as they have already been shown to be ineffective in SPR (Schwarzer et al., 2020).

**Removing Features with Masking** We discussed why post-loss-calculation modifications cannot be applied to objectives that involve components in the feature dimension rather than the batch dimension. However, non-terminal masking can be employed to exclude samples from the batch before calculating the SSL loss. Thus, we masked features during the training of the SPR-VICReg and SPR-Barlow agents, leading to unexpected results. As shown in Table 1, the SPR-Barlow agent performed even worse than the random agent. A likely explanation is that the Barlow Twins' objective relies on batch normalization to compute the cross-covariance matrix. Since masking causes the batch size to vary dynamically, the batch statistics become inconsistent, adversely affecting the batch normalization process.However, this degradation is not observed to the same extent in the SPR-VICReg agent, as the VICReg objective does not rely on batch normalization.

**Continuous Control Formulation** Although SPR is created specifically for discrete control, delving into the impact of SSL objectives solely within discrete control domains doesn't provide a comprehensive understanding. This is why we adopt a parallel setup to that of PlayVirtual (Yu et al., 2021), where they establish an SPR-like scheme referred to as SPR† as a baseline for continuous control. They utilize the soft actor-critic algorithm (Haarnoja et al., 2018), instead of q-learning due to the continuous nature of the actions. They do not use terminal state masking (since terminal states for control problems are target states) and prioritized replay weighting (since they use a uniform buffer). This shows the importance of generally applicable auxiliary tasks for data-efficient RL.

We evaluate PlayVirtual and SPR† from scratch since we were not able to replicate Yu et al. (2021)'s results, potentially due to different benchmark versions. Furthermore, we assess the performance of VICReg-High and Barlow Twins within the SPR† configuration. We exclude VICReg-Low in this setting due to the minimal performance difference observed in Atari.

Finally, we explore the potential impact of incorporating the predictor network into Barlow Twins and VICReg, even though they inherently do not need it to prevent dimension collapse. Although the addition of a predictor network is novel in Barlow Twins, VICReg becomes similar to the SPR with this addition like SPR with variance-covariance regularization. The decision to refrain from conducting similar experiments in Atari stems from the substantially higher experimental costs, which are at least 10 times greater than those in the control setting.

# 5 Evaluation Setup

## 5.1 Benchmarking: Rliable Framework

Agarwal et al. (2021) discusses the limitations of using mean and median scores as singular estimates in RL benchmarks and highlights the disparities between conventional single-point estimates and

241 the broader interval estimates, emphasizing the potential ramifications for benchmark dependability
242 and interpretation. In alignment with their suggestions, we provide a succinct overview of human-
243 normalized scores, furnished with stratified bootstrap confidence intervals, in Figures 2 and 3.

## 5.2 Atari 100k

245 We assess the SPR framework in a reduced-sample Atari setting, called the Atari 100k bench-
246 mark (Kaiser et al., 2019). In this setting, the training dataset comprises 100,000 environment
247 steps, which is equivalent to about 400,000 frames or slightly under two hours of equivalent hu-
248 man experience. This contrasts with the conventional benchmark of 50,000,000 environment steps,
249 corresponding to approximately 39 days of accumulated experience.

250 The main metric for this setting, widely acknowledged for assessing performance in the Atari 100k
251 context, is the human-normalized score. This measure is mathematically defined as in equation 3,
252 where random score pertains to outcomes achieved through a random policy and the human score is
253 derived from human players (Wang et al., 2015).

$$\frac{score_{\text{agent}} - score_{\text{random}}}{score_{\text{human}} - score_{\text{random}}} \tag{3}$$

## 5.3 Deep Mind Control Suite

255 In the Deep Mind Control Suite (Tassa et al., 2018), the agent is configured to function solely
256 based on pixel inputs. This choice is justified by several reasons: the environments involved offer a
257 reasonably challenging and diverse array of tasks, the sample efficiency of model-free reinforcement
258 learning algorithms is notably low when operating directly from pixels in these benchmarks and the
259 performance on the DM control suite is comparable to the context of robot learning in real-world
260 benchmarks.

261 We use the following six environments (Yarats et al., 2020) for benchmarking: ball-in-cup, finger-
262 spin, reacher-easy, cheetah-run, walker-walk and cartpole-swingup, for 100k steps each.

# 6 Results and Discussion

## 6.1 Atari 100k

265 We mainly investigate the following new SPR models, along with the original SPR: (i) SPR-
266 Naked, featuring no modifications, (ii) SPR-Naked+Non, incorporating terminal masking, (iii) SPR-
267 Naked+Prio, integrating prioritized replay weighting, (iv) SPR-Barlow, (v) SPR-VICReg-High,
268 characterized by a high covariance weight, and (vi) SPR-VICReg-Low, characterized by a low co-
269 variance weight. Moreover, we discuss SR-SPR and BBF with their no modifications versions.

270 Figure 2 shows the performance of the seven agents in the Atari 100k benchmark, calculated us-
271 ing the rliable framework (Agarwal et al., 2021). The individual game performances are given in
272 Appx. 11 and we describe evaluation setup in Section 5.

273 **SPR and SSL Modifications**. The original SPR-agent performs the best (top row of Fig. 2). The
274 modifications to the SPR's SSL objective (see Section 3) have significant impact on the performance
275 but they are not mentioned in the relevant papers (SPR (Schwarzer et al., 2020), SR-SPR (D'Oro
276 et al., 2023; Nikishin et al., 2022), or BBF (Schwarzer et al., 2023)). The no modifications ver-
277 sion, SPR-Naked, performs the worst with a nearly 20% performance drop based on the IQM score
278 (last row of Fig. 2). This is crucial because such modifications may not be suitable for all problem
279 domains, which limits their transferability and generalizability. On the other hand, the role of ter-
280 minal masking and prioritized replay weighting in SPR is especially interesting, as they help boost
281 performance in situations where pure representation learning struggles.

282 Incorporating prioritized replay weights has a positive effect on SPR ($5^{th}$ row of Fig. 2). These
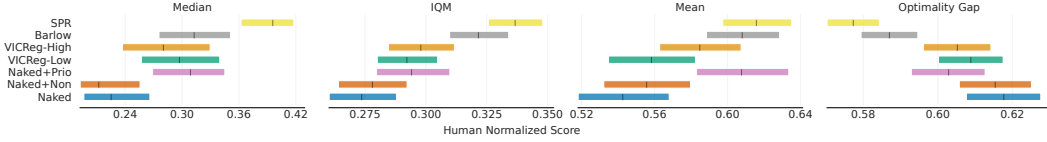283 weights act as markers for Bellman errors that mirror the agent's Q-value approximation perfor-

Figure 2: Mean, median, interquartile mean human normalized scores and optimality gap (lower is better) computed with stratified bootstrap confidence intervals in Atari 100k. 50 runs for SPR-Barlow, SPR-VICReg-High, SPR-VICReg-Low, SPR-Naked+Prio, SPR-Naked+Non,SPR-Naked, 100 runs for SPR from (Agarwal et al., 2021).
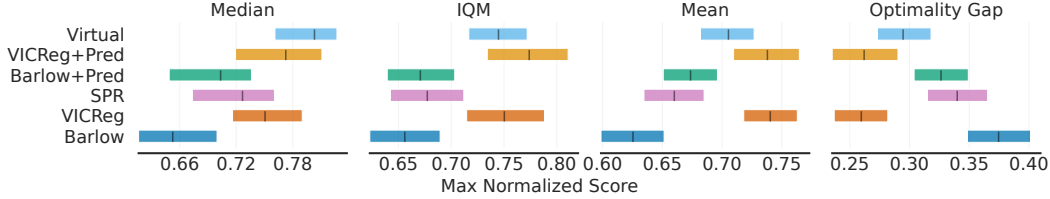


Figure 3: Mean, median, interquartile mean max normalized scores and optimality gap (lower is better) computed with stratified bootstrap confidence intervals in Deep Mind Control Suite 100k, 10 runs for all agents.

mance on particular transitions. Introducing these weights into the representation loss intensifies the emphasis on refining representations that the agent struggles with.

Empirically, terminal state masking shows negligible positive effects, unlike replay weighting, ($6^{th}$ row of Fig. 2). The limited impact of masking might be attributed to the episode lengths, where the agent encounters many regular states but only a single terminal state. The SSL loss may be primarily influenced by intermediate states, which could reduce the effectiveness of masking in these scenarios.

On the other hand, there is a clear synergy between these modifications within SPR. Masking terminal states might be advantageous when agents encounter frequent failures during the initial stages of training or due to the nature of the games. In such cases, terminal states may dominate the replay buffer, which could introduce biased representations that become challenging to correct later on and make it harder for the agent to adapt and improve as it progresses

**SPR-Barlow**. The performance of the Barlow Twins agent is close to the SPR's ($2^{nd}$ row of Fig. 2), with only a 5% difference, where as SPR-Naked has a 20% gap. As described in Section 4, modifications related to SSL do not directly apply to Barlow Twins, VICReg, or any other method regularization in the feature dimension. As such, performing similar to a method with RL specific modifications suggests that Barlow Twins has the potential to serve as a substitute, indicating its promise as a versatile SSL objective for data-efficient RL.

The performance gap between SPR-naked and the feature decorrelation methods (Barlow and VICReg) in this context is somewhat surprising since BYOL or Simsiam outperform them in image classification. In vision pretraining, the goal is to obtain embeddings with well-defined clusters based on the training corpora, enhancing classification performance, where feature decorrelation may be of hindrance. In RL, it is important to differentiate between states (good, bad, or promising if they have not been explored yet) which may not be too different in the image space. As such, methods that emphasize the use of the entire embedding space potentially have a better chance of state separation.

To test this, we evaluate the rank (Kumar et al., 2021) of the advantage and value heads, as well as the output of the convolution head, which is shared by both the RL and SSL objectives. We evaluated multiple methods like Barlow Twins and VICReg, in addition to a variant without SSL
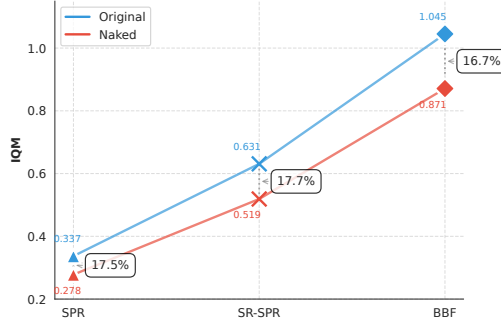
Figure 4: Comparison of IQM performance for the SPR, SR-SPR, and BBF agents alongside their corresponding naked versions. Naked results of SR-SPR and BBF are averaged out across 10 different runs

loss. We found that the rank converges similarly across different games and even if they don't, this does not correlate with performance. We also measured dormant neurons (Sokar et al., 2023) and observed that the results were consistent with the rank findings. These evaluations are detailed in Appx. 9.

**SPR-VICRegs**. Initially, we used the default VICReg hyperparameters given in the original paper (Bardes et al., 2021). Surprisingly, VICReg exhibits a 13% lower performance ($4^{th}$ row of Fig. 2) compared to SPR although it surpasses SPR-Naked. It also falls short of Barlow Twins. This outcome is not immediately evident given that it has a high similarity to the Barlow Twins' objective. One plausible explanation could be the presence of multiple loss components, possibly undermining covariance. To address this, we explore alternative hyperparameters, selecting the set with the highest covariance hyperparameter that avoids collapse and denote it as SPR-VICReg-High, while the previous one is referred to as SPR-VICReg-Low. However, the performance only marginally increases by 2% ($3^{rd}$ row of Fig. 2), lacking behind Barlow Twins once again. The underlying reasons for this performance gap remain subject to further exploration. Nonetheless, it still showcases the effectiveness of feature decorrelation based objectives since both types outperform SPR-Naked.

**BBF and SR-SPR**. It could be argued that modifications to SPR significantly influence performance, particularly due to its relatively low score on Atari 100k, where such changes may have an amplified effect, whereas they might have a more limited impact on stronger models. BBF, the leading value-based agent achieving human-level results on Atari 100k, is built upon SR-SPR, a variant of SPR. Notably, both SR-SPR and BBF exhibit IQM values nearly 3x and 2x higher than SPR, respectively. Thus, their unmodified results will provide insight into whether modifications still play a significant role, even when the model is highly efficient and performing at a human level.

As shown in Figure 4, we observe that modifications result in a fairly consistent performance decline across all models. Due to computational constraints, we did not conduct experiments to determine which modifications have the greatest impact or whether certain SSL objectives could reduce the need for modifications. However, our findings further support and strengthen our earlier conclusions regarding the impact of modifications on SPR.

## 6.2 DMControl

We further evaluate the SSL objectives with the DMControl suite, described in Section 5) since this domain can provide additional insights into the efficacy of SSL objectives in RL. However, since there is no terminal state in this environment and a uniform replay buffer is used, modifications to the SPR loss are not feasible. As such, this evaluation will focus on the generalization of used objectives across domains without targeted optimization for specific problems.

346 Moreover, SPR is not explicitly designed for continuous control. As such, we use a different set of
347 agents modified for continuous control as described in Section 4 but keep the same SSL hyperparam-
348 eters from the Atari benchmark. We pick SPR-VICReg-High due to its better performance over the
349 lower covariance version. We additionally evaluate SPR-Barlow and SPR-Vicreg with an MLP layer
350 as an additional predictor, reflecting Bardes et al. (2021)'s findings on the enhanced performance of
351 BYOL with variance regularization. We build upon the PlayVirtual (Yu et al., 2021) methodology,
352 which is an SPR equipped with an improved transition model, and use it as our baseline.

353 We observe from Fig. 3 that the Barlow Twins objective exhibits the lowest performance, although
354 it closely aligns with SPR, with IQM scores of 0.656, and 0.677 respectively. An interesting obser-
355 vation is that VICReg with an IQM of 0.75 is as good as PlayVirtual (Yu et al., 2021) with 0.744.
356 This underscores the potential of SSL objectives in continuous control. While their impact is vi-
357 tal in discrete control as well, the overall effect, especially when considering the maximum score
358 (representing human performance), is relatively modest. Nevertheless, a substantial improvement
359 is evident in continuous control, even when compared to the highest achievable score. We also see
360 that adding a predictor network has a minimal but positive impact on the IQM performances of both
361 Barlow and VICReg.

## 7  Conclusion

363 Our study demonstrates the significant impact of SSL objective modifications within the SPR frame-
364 work for reinforcement learning, particularly in data-efficient scenarios. We show that specific ad-
365 justments like terminal state masking and prioritized replay weighting substantially improve per-
366 formance on the Atari 100k benchmark, with benefits extending to derivative frameworks such as
367 SR-SPR and BBF. However, our experiments on the DeepMind Control Suite reveal that these en-
368 hancements are not universally applicable across all RL environments. Investigation of alternative
369 SSL objectives (e.g., Barlow Twins, VICReg) further elucidates the nuanced relationship between
370 objective choice and RL task characteristics. These findings emphasize the critical role of carefully
371 tailored SSL objectives in achieving data efficiency in self-predictive reinforcement learning, high-
372 lighting the need for a context-sensitive approach to SSL modification in RL algorithm development.
373 Our work provides valuable insights for researchers and practitioners seeking to optimize RL algo-
374 rithms across diverse applications, potentially leading to more efficient and effective reinforcement
375 learning systems.

## References

377 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C. Courville, and Marc G. Belle-
378 mare. Deep reinforcement learning at the edge of the statistical precipice. In *Neural Information
379 Processing Systems*, 2021.

380 Benjamin J. Ayton and Masataro Asai. Width-based planning and active learning for atari. In
381 *International Conference on Automated Planning and Scheduling*, 2021. URL https://api.
382 semanticscholar.org/CorpusID:238226837.

383 Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Flo-
384 rian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, Avi Schwarzschild, Andrew Gor-
385 don Wilson, Jonas Geiping, Quentin Garrido, Pierre Fernandez, Amir Bar, Hamed Pirsiavash,
386 Yann LeCun, and Micah Goldblum. A cookbook of self-supervised learning, 2023.

387 Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization
388 for self-supervised learning. *ArXiv*, abs/2105.04906, 2021.

389 Omer Veysel Cagatan and Baris Akgun. Barlowrl: Barlow twins for data-efficient reinforcement
390 learning, 2023.

391 Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.
392 Unsupervised learning of visual features by contrasting cluster assignments, 2021.

Edoardo Cetin, Philip J. Ball, Steve Roberts, and Oya Çeliktutan. Stabilizing off-policy deep reinforcement learning from pixels. In *International Conference on Machine Learning*, 2022a. URL https://api.semanticscholar.org/CorpusID:250265109.

Edoardo Cetin, Benjamin Paul Chamberlain, Michael M. Bronstein, and Jonathan J. Hunt. Hyperbolic deep reinforcement learning. *ArXiv*, abs/2210.01542, 2022b. URL https://api.semanticscholar.org/CorpusID:252693361.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. *ArXiv*, abs/2002.05709, 2020. URL https://api.semanticscholar.org/CorpusID:211096730.

Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15745–15753, 2020.

Pierluca D'Oro, Max Schwarzer, Evgenii Nikishin, Pierre-Luc Bacon, Marc G. Bellemare, and Aaron C. Courville. Sample-efficient reinforcement learning by breaking the replay ratio barrier. In *International Conference on Learning Representations*, 2023. URL https://api.semanticscholar.org/CorpusID:259298604.

Manuel Goulão and Arlindo L. Oliveira. Pretraining the vision transformer using self-supervised methods for vision based deep reinforcement learning. *ArXiv*, abs/2209.10901, 2022. URL https://api.semanticscholar.org/CorpusID:252439214.

Manuel Goulão and Arlindo L. Oliveira. Pretraining the vision transformer using self-supervised methods for vision based deep reinforcement learning, 2023.

Jean-Bastien Grill, Florian Strub, Florent Altch'e, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Ávila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning. *ArXiv*, abs/2006.07733, 2020. URL https://api.semanticscholar.org/CorpusID:219687798.

Tuomas Haarnoja, Aurick Zhou, P. Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *ArXiv*, abs/1801.01290, 2018. URL https://api.semanticscholar.org/CorpusID:28202810.

Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.

H. V. Hasselt, Matteo Hessel, and John Aslanides. When to use parametric models in reinforcement learning? *ArXiv*, abs/1906.05243, 2019. URL https://api.semanticscholar.org/CorpusID:186206746.

Matteo Hessel, Joseph Modayil, H. V. Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Gheshlaghi Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2017. URL https://api.semanticscholar.org/CorpusID:19135734.

Tao Huang, Jiacheng Wang, and Xiao Chen. Accelerating representation learning with view-consistent dynamics in data-efficient reinforcement learning. *ArXiv*, abs/2201.07016, 2022. URL https://api.semanticscholar.org/CorpusID:246035501.

Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H. Campbell, K. Czechowski, D. Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz Mohiuddin, Ryan Sepassi, G. Tucker, and Henryk Michalewski. Model-based reinforcement learning for atari. *ArXiv*, abs/1903.00374, 2019. URL https://api.semanticscholar.org/CorpusID:67856232.

Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *ArXiv*, abs/2004.13649, 2020. URL https://api.semanticscholar.org/CorpusID:216562627.

Aviral Kumar, Rishabh Agarwal, Dibya Ghosh, and Sergey Levine. Implicit under-parameterization inhibits data-efficient deep reinforcement learning, 2021. URL https://arxiv.org/abs/2010.14498.

Hojoon Lee, Hanseul Cho, Hyunseung Kim, Daehoon Gwak, Joonkee Kim, Jaegul Choo, Se-Young Yun, and Chulhee Yun. Enhancing generalization and plasticity for sample efficient reinforcement learning. *ArXiv*, abs/2306.10711, 2023a. URL https://api.semanticscholar.org/CorpusID:259203876.

Hojoon Lee, Koanho Lee, Dongyoon Hwang, Hyunho Lee, Byungkun Lee, and Jaegul Choo. On the importance of feature decorrelation for unsupervised representation learning in reinforcement learning, 2023b.

Xiang Li, Jinghuan Shang, Srijan Das, and Michael S. Ryoo. Does self-supervised learning really improve reinforcement learning from pixels?, 2023. URL https://arxiv.org/abs/2206.05266.

Litian Liang, Yaosheng Xu, Stephen McAleer, Dailin Hu, Alexander T. Ihler, P. Abbeel, and Roy Fox. Reducing variance in temporal-difference value estimation via ensemble of deep networks. *ArXiv*, abs/2209.07670, 2022. URL https://api.semanticscholar.org/CorpusID:250341019.

Hao Liu and P. Abbeel. Aps: Active pretraining with successor features. In *International Conference on Machine Learning*, 2021. URL https://api.semanticscholar.org/CorpusID:235825462.

Vincent Micheli, Eloi Alonso, and Franccois Fleuret. Transformers are sample efficient world models. *ArXiv*, abs/2209.00588, 2022. URL https://api.semanticscholar.org/CorpusID:251979354.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas Kirkeby Fidjeland, Georg Ostrovski, Stig Petersen, Charlie Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518:529–533, 2015. URL https://api.semanticscholar.org/CorpusID:205242740.

Evgenii Nikishin, Max Schwarzer, Pierluca D'Oro, Pierre-Luc Bacon, and Aaron C. Courville. The primacy bias in deep reinforcement learning. In *International Conference on Machine Learning*, 2022. URL https://api.semanticscholar.org/CorpusID:248811264.

Johan Obando-Ceron, João G. M. Araújo, Aaron Courville, and Pablo Samuel Castro. On the consistency of hyper-parameter selection in value-based deep reinforcement learning, 2024. URL https://arxiv.org/abs/2406.17523.

Serdar Ozsoy, Shadi S. Hamdan, Sercan Ö. Arik, Deniz Yuret, and Alper Tunga Erdogan. Self-supervised learning with an information maximization criterion. *ArXiv*, abs/2209.07999, 2022.

Jan Robine, Tobias Uelwer, and Stefan Harmeling. Smaller world models for reinforcement learning, 2021.

Jan Robine, Marc Hoftmann, Tobias Uelwer, and Stefan Harmeling. Transformer-based world models are happy with 100k interactions. *ArXiv*, abs/2303.07109, 2023. URL https://api.semanticscholar.org/CorpusID:257496038.

Max Schwarzer, Ankesh Anand, Rishab Goel, R. Devon Hjelm, Aaron C. Courville, and Philip Bachman. Data-efficient reinforcement learning with self-predictive representations. In *International Conference on Learning Representations*, 2020. URL https://api.semanticscholar.org/CorpusID:222163237.

Max Schwarzer, Ankesh Anand, Rishab Goel, R. Devon Hjelm, Aaron C. Courville, and Philip Bachman. Repository published by the spr. https://github.com/mila-iqia/spr, 2021a.

Max Schwarzer, Nitarshan Rajkumar, Michael Noukhovitch, Ankesh Anand, Laurent Charlin, Devon Hjelm, Philip Bachman, and Aaron C. Courville. Pretraining representations for data-efficient reinforcement learning. In *Neural Information Processing Systems*, 2021b. URL https://api.semanticscholar.org/CorpusID:235377401.

Max Schwarzer, Johan S. Obando-Ceron, Aaron C. Courville, Marc G. Bellemare, Rishabh Agarwal, and Pablo Samuel Castro. Bigger, better, faster: Human-level atari with human-level efficiency. *ArXiv*, abs/2305.19452, 2023. URL https://api.semanticscholar.org/CorpusID:258987895.

Thalles Silva, Helio Pedrini, and Adín Ramírez Rivera. Learning from memory: Non-parametric memory augmented self-supervised learning of visual features. 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. 45451–45467. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/silva24c.html.

Ghada Sokar, Rishabh Agarwal, Pablo Samuel Castro, and Utku Evci. The dormant neuron phenomenon in deep reinforcement learning, 2023. URL https://arxiv.org/abs/2302.12902.

A. Srinivas, Michael Laskin, and P. Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. *ArXiv*, abs/2004.04136, 2020. URL https://api.semanticscholar.org/CorpusID:215415964.

Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Riedmiller. Deepmind control suite, 2018.

Manan Tomar, Utkarsh A. Mishra, Amy Zhang, and Matthew E. Taylor. Learning representations for pixel-based control: What matters and why?, 2021. URL https://arxiv.org/abs/2111.07775.

Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding, 2019.

Ziyun Wang, Tom Schaul, Matteo Hessel, H. V. Hasselt, Marc Lanctot, and Nando de Freitas. Dueling network architectures for deep reinforcement learning. In *International Conference on Machine Learning*, 2015. URL https://api.semanticscholar.org/CorpusID:5389801.

Xi Weng, Yunhao Ni, Tengwei Song, Jie Luo, Rao Muhammad Anwer, Salman Khan, Fahad Shahbaz Khan, and Lei Huang. Modulate your spectrum in self-supervised learning, 2024. URL https://arxiv.org/abs/2305.16789.

Denis Yarats, Amy Zhang, Ilya Kostrikov, Brandon Amos, Joelle Pineau, and Rob Fergus. Improving sample efficiency in model-free reinforcement learning from images, 2020.

Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games with limited data, 2021.

Tao Yu, Cuiling Lan, Wenjun Zeng, Mingxiao Feng, Zhizheng Zhang, and Zhibo Chen. Playvirtual: Augmenting cycle-consistent virtual trajectories for reinforcement learning. *Advances in Neural Information Processing Systems*, 34, 2021.

Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. *ArXiv*, abs/2103.03230, 2021. URL https://api.semanticscholar.org/CorpusID:232110471.

Yifan Zhang, Zhiquan Tan, Jingqin Yang, Weiran Huang, and Yang Yuan. Matrix information theory for self-supervised learning. 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. 59897–59918. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/zhang24bi.html.

Jinghao Zhou, Li Dong, Zhe Gan, Lijuan Wang, and Furu Wei. Non-contrastive learning meets language-image pre-training, 2022. URL https://arxiv.org/abs/2210.09304.

Ömer Veysel Çağatan. Unsee: Unsupervised non-contrastive sentence embeddings, 2024. URL https://arxiv.org/abs/2401.15316.

# Supplementary Materials

*The following content was not necessarily subject to peer review.*

## 8 Background

### 8.1 Barlow Twins

The Barlow Twins (Zbontar et al., 2021) employs a symmetric network with twin branches, each processing a different augmented perspective of input data. It aims to minimize off-diagonal components and align diagonal elements of a cross-covariance matrix derived from the representations of these branches. The process involves generating two altered views ($Y^A$ and $Y^B$) using data augmentations, inputting them into a function $f_\theta$ to produce embeddings ($Z^A$ and $Z^B$).

The Barlow Twins loss is defined as:

$$\mathcal{L}_{\mathcal{BT}} \triangleq \underbrace{\sum_i (1 - \mathcal{C}_{ii})^2}_{\text{invariance term}} + \; \lambda \underbrace{\sum_i \sum_{j \neq i} \mathcal{C}_{ij}^2}_{\text{redundancy reduction term}} \tag{4}$$

where $\lambda > 0$ balances the invariance (diagonal elements) and redundancy reduction (off-diagonal) in the loss function. $\mathcal{C}$ is the cross-correlation matrix from embedding outputs of identical networks in the batch. A matrix element is defined as:

$$\mathcal{C}_{ij} \triangleq \frac{\sum_b z_{b,i}^A z_{b,j}^B}{\sqrt{\sum_b \left(z_{b,i}^A\right)^2} \sqrt{\sum_b \left(z_{b,j}^B\right)^2}} \tag{5}$$

where $b$ represents the samples in the batch, and $i$ and $j$ represent dimension indices of the networks' output. Each dimension of the square covariance matrix, $\mathcal{C}$, is the same as the embedding dimension (output dimensionality of the networks). Its values range between -1 (indicating complete anti-correlation) and 1 (representing perfect correlation).

### 8.2 VICReg

VICReg (Bardes et al., 2021) is a method designed to tackle the challenge of collapse directly. It achieves this by introducing a straightforward regularization term that specifically targets the variance of the embeddings along each dimension independently. In addition to addressing the variance, VICReg includes a mechanism to diminish redundancy and ensure decorrelation among the embeddings, accomplished through covariance regularization.

The variance regularization term is a hinge function on the standard deviation of the embeddings along the batch dimension:

$$v(Z) = \frac{1}{d} \sum_{j=1}^{d} \max(0, \gamma - S(z^j, \epsilon)) \tag{6}$$

where $S$ is the regularized standard deviation defined by:

$$S(x; \epsilon) = \sqrt{\text{Var}(x) + \epsilon} \tag{7}$$

Covariance matrix of $Z$ is defined as:

$$C(Z) = \frac{1}{n-1} \sum_{i=1}^{n} (z_i - \bar{z})(z_i - \bar{z})^T \tag{8}$$

15

573 where $\bar{z} = \frac{1}{n}\sum_{i=1}^{n} z_i$. Covariance regularization is defined as:

$$c(Z) = \frac{1}{d}\sum_{i}\sum_{j\neq i}\mathcal{C}_{ij}{}^2 \tag{9}$$

574 where $d$ is the feature dimension. The invariance criterion between $Z$ and $Z'$ is the mean-squared
575 Euclidean distance between each pair of vectors, without any normalization.

$$s(Z, Z') = \frac{1}{n}\sum_{i=1}^{n}||z_i - z_i'||^2 \tag{10}$$

576 The overall loss function is a weighted average of the invariance, variance, and covariance terms:

$$l(Z, Z') = \alpha v(Z) + \beta c(Z) + \gamma s(Z, Z') \tag{11}$$

577 where $\alpha$, $\lambda$, and $\gamma$ hyper-parameters control the importance of each term in the loss.

578 VICReg is quite similar to Barlow Twins in terms of its loss formulation. However, instead of
579 decorrelating the cross-correlation matrix directly, it regularizes the variance along each dimension
580 of the representation, reduces correlation and minimizes the difference of embeddings. This prevents
581 dimension collapse and also forces the two views to be encoded similarly. Additionally, reducing
582 covariance encourages different dimensions of the representation to capture distinct features.

## 9 Rank and Dormant Neuron

584 Kumar et al. (2021) introduced the concept of *effective rank* for representations, represented as
585 $srank_\delta(\phi)$, with $\delta$ being a threshold parameter, set to 0.01 as per their study. They proposed that
586 effective rank is linked to the expressivity of a network, where a decrease in effective rank implies
587 an implicit under-parameterization. The study provides evidence indicating that bootstrapping is the
588 primary factor contributing to the collapse of effective rank, which in turn degrades performance.

589 To investigate how SSL objectives might mitigate rank collapse, we computed the rank of the con-
590 volution output and the outputs of the penultimate layers from the advantage and value heads of
591 three different agents: SPR-VICReg, SPR-Barlow, and ZeroJump (SPR without a transition model),
592 scores in 1. Our observations indicate that, although there are some rank differences among the
593 agents, they often converge to the same rank, and these differences do not correlate with the perfor-
594 mance scores. Figure 5, 7 and 6 include ranks across all games.

595 To explore this further, we examined the proportion of dormant neurons, which are neurons that have
596 near-zero activations. Sokar et al. (2023) showed that deep reinforcement learning agents experience
597 a rise in the number of dormant neurons within their networks. Additionally, a higher prevalence of
598 dormant neurons is associated with poorer performance.

599 We also do not observe a clear pattern in the fractions of dormant neurons, in Figure 8 that could
600 account for the disparities in performance scores, similar to what was seen in the case of neuron
601 ranks. Unlike rank-based observations, where patterns may emerge, the distribution of dormant
602 neurons does not offer an explanation for the differences in the scores across models. This suggests
603 that the relationship between neuron activity and performance metrics might be more complex and
604 not directly attributable to the proportion of inactive neurons.

## 10 Experimental Details

606 We retain all hyperparameters of SPR, SR-SPR, and BBF, except for SPR-Barlow and SPR-VICReg,
607 where we adjust the SPR loss weight and increase the batch size from 32 to 64. The official reposi-
608 tories of the models are used, and all experiments are conducted on a Tesla T4 GPU.

## 11 Full Results on Atari 100k

Table 3: Returns on the 26 games of Atari 100k after 2 hours of real-time experience, and human-normalized aggregate metrics. (VR: VICReg, results with 5 integral digits are rounded to the first integer to fit the table)

| Game | Rand. | Human | Naked | Non | Prio | VR-L | VR-H | Barlow | SPR |
|---|---|---|---|---|---|---|---|---|---|
| Alien | 227.8 | 7127.7 | 868.9 | 881.7 | 872.7 | 902.9 | 922.4 | 891.8 | 841.9 |
| Amidar | 5.8 | 1719.5 | 165.6 | 179.1 | 164.2 | 181.1 | 176.4 | 177.1 | 179.7 |
| Assault | 222.4 | 742.0 | 544.5 | 564.6 | 589.2 | 536.4 | 575.7 | 581.4 | 565.6 |
| Asterix | 210 | 8503.3 | 972.0 | 951.0 | 977.8 | 955.4 | 1021.7 | 981.2 | 962.5 |
| BankHeist | 14.2 | 753.1 | 61.6 | 70.1 | 60.2 | 79.9 | 82.9 | 73.5 | 345.4 |
| BattleZone | 2360 | 37188 | 7552.4 | 9424.2 | 13102 | 12557 | 14892 | 14954 | 14834 |
| Boxing | 0.1 | 12.1 | 27.3 | 30.4 | 36.4 | 31.3 | 33.9 | 35.1 | 35.7 |
| Breakout | 1.7 | 30.5 | 16.7 | 18.0 | 18.2 | 16.9 | 16.3 | 17.0 | 19.6 |
| ChopComm | 811 | 7387.8 | 906.8 | 949.8 | 901.0 | 832.9 | 929.9 | 938.9 | 946.3 |
| CrzyClmbr | 10781 | 35829 | 30056 | 32667 | 35829 | 27035 | 29023 | 29229 | 36701 |
| DemonAtt | 152.1 | 1971.0 | 514.7 | 511.0 | 522.9 | 461.2 | 547.2 | 519.2 | 517.6 |
| Freeway | 0.0 | 29.6 | 17.4 | 13.71 | 16.3 | 28.0 | 27.7 | 29.5 | 19.3 |
| Frostbite | 65.2 | 4334.7 | 1137.2 | 1010.9 | 1014.2 | 1353.0 | 1181.4 | 1191.3 | 1170.7 |
| Gopher | 257.6 | 2412.5 | 585.0 | 660.1 | 548.4 | 737.9 | 713.5 | 691.2 | 660.6 |
| Hero | 1027 | 30826 | 6937.8 | 6497.8 | 5686.6 | 5495.1 | 5559.6 | 5746.8 | 5858.6 |
| Jamesbond | 29 | 302.8 | 327.2 | 359.9 | 349.1 | 357.6 | 384.3 | 404.2 | 366.5 |
| Kangaroo | 52 | 3035.0 | 2970.9 | 2812.1 | 3016.5 | 2290.6 | 1998.3 | 1771.2 | 3617.4 |
| Krull | 1598 | 2665.5 | 3980.4 | 4061.8 | 4213.1 | 4166.6 | 4513.9 | 4363.2 | 3681.6 |
| KFMaster | 258.5 | 22736 | 13126 | 14595 | 15757 | 1488.4 | 15548 | 15998 | 14783 |
| MsPacman | 307.3 | 6951.6 | 1262.1 | 1162.6 | 1324.6 | 1366.8 | 1588.2 | 1388.2 | 1318.4 |
| Pong | -20.7 | 14.6 | -1.8 | -6.0 | -7.2 | -6.3 | -10.1 | -6.7 | -5.4 |
| PrivateEye | 24.9 | 69571 | 85.6 | 77.0 | 88.0 | 100.9 | 96.6 | 99.6 | 86.0 |
| Qbert | 163.9 | 13455 | 847.2 | 758.6 | 759.8 | 796.9 | 687.6 | 765.8 | 866.3 |
| RoadRunner | 11.5 | 7845.0 | 12595 | 12713 | 11211 | 10683 | 9531.5 | 12412 | 12213 |
| Seaquest | 68.4 | 42055 | 524.0 | 524.2 | 523.2 | 576.3 | 651.0 | 669.1 | 558.1 |
| UpNDown | 533.4 | 11693 | 9569.3 | 8130.6 | 10331 | 7952.7 | 9415.3 | 10818 | 10859 |
| #Sprhmn(↑) | 0 | N/A | 4 | 3 | 3 | 4 | 4 | 4 | 6 |
| Mean (↑) | 0.00 | 1.000 | 0.542 | 0.555 | 0.608 | 0.558 | 0.585 | 0.608 | 0.616 |
| Median (↑) | 0.00 | 1.000 | 0.225 | 0.221 | 0.308 | 0.297 | 0.280 | 0.312 | 0.396 |
| IQM (↑) | 0.00 | 1.000 | 0.273 | 0.278 | 0.298 | 0.292 | 0.298 | 0.321 | 0.337 |
| Opt. Gap (↓) | 1.00 | 0.000 | 0.617 | 0.615 | 0.603 | 0.609 | 0.605 | 0.587 | 0.577 |

## 12   Full Results on DMControl 100k

Table 4: Returns on the of DMControl 100k, and Max-normalized aggregate metrics.

| Environment | Virtual | VICReg+Pred | Barlow+Pred | SPR | VICReg | Barlow |
|---|---|---|---|---|---|---|
| FINGER, SPIN | 896.2 | 760.6 | 781.0 | 755.9 | 730.0 | 861.8 |
| CARTPOLE, SWINGUP | 815.1 | 791.6 | 784.0 | 826.0 | 780.1 | 778.6 |
| REACHER, EASY | 827.0 | 790.7 | 589.6 | 671.5 | 736.1 | 526.5 |
| CHEETAH, RUN | 489.6 | 504.3 | 461.6 | 435.2 | 493.5 | 478.6 |
| WALKER, WALK | 404.7 | 622.8 | 521.7 | 404.7 | 765. | 182.2 |
| BALL IN CUP, CATCH | 835.4 | 891.6 | 622.8 | 835.4 | 937.5 | 924.9 |
| Mean (↑) | 0.705 | 0.738 | 0.673 | 0.660 | 0.740 | 0.625 |
| Median (↑) | 0.803 | 0.772 | 0.703 | 0.726 | 0.750 | 0.652 |
| IQM (↑) | 0.744 | 0.773 | 0.670 | 0.677 | 0.750 | 0.656 |
| Optimality Gap (↓) | 0.294 | 0.260 | 0.326 | 0.339 | 0.29 | 0.374 |

## 13 Rank and Dormant Neuron Results

Figure 5: Rank of the output from the penultimate layer of the value head, measured every 10,000 steps and averaged across 10 different runs for every game.
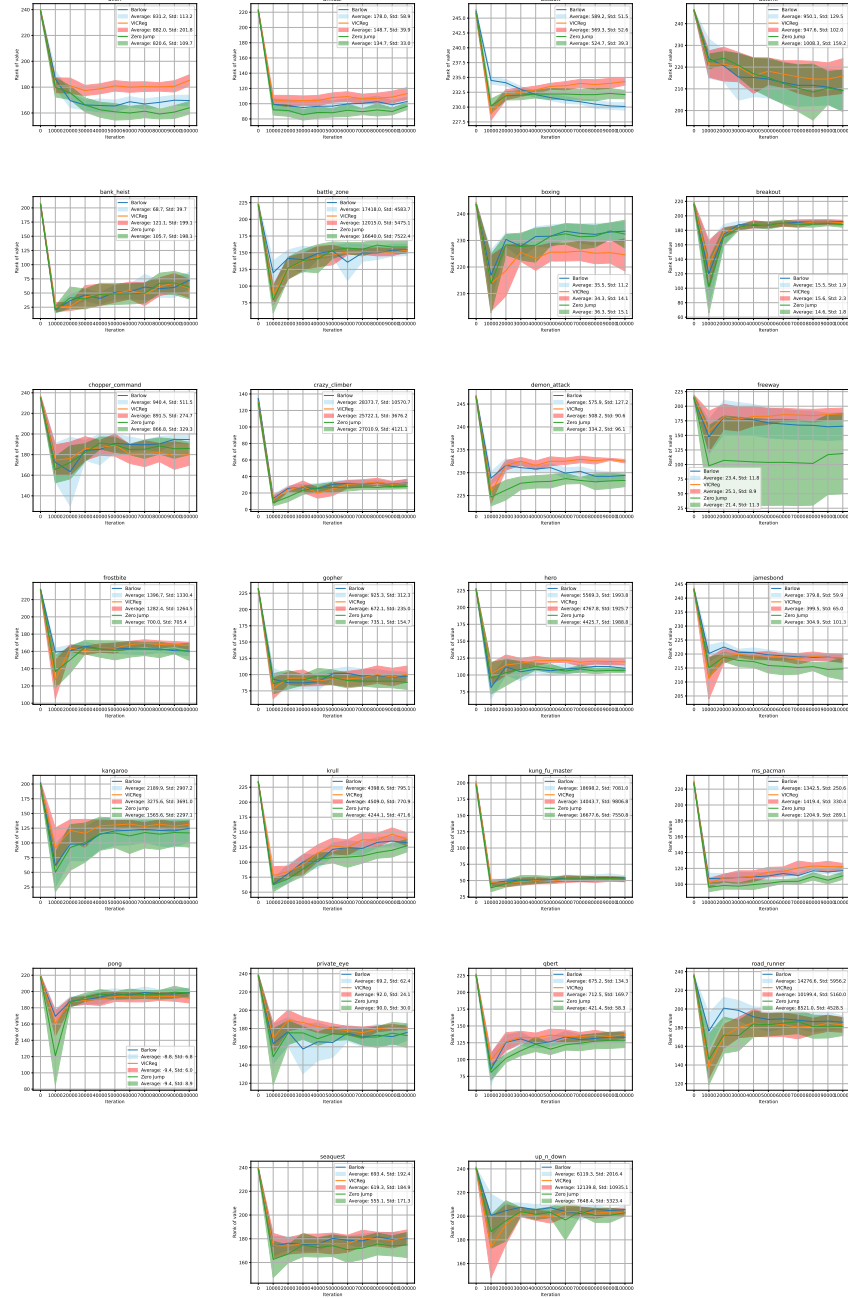


19

Figure 6: Rank of the output from the convolution encoder, measured every 10,000 steps and averaged across 10 different runs for every game.
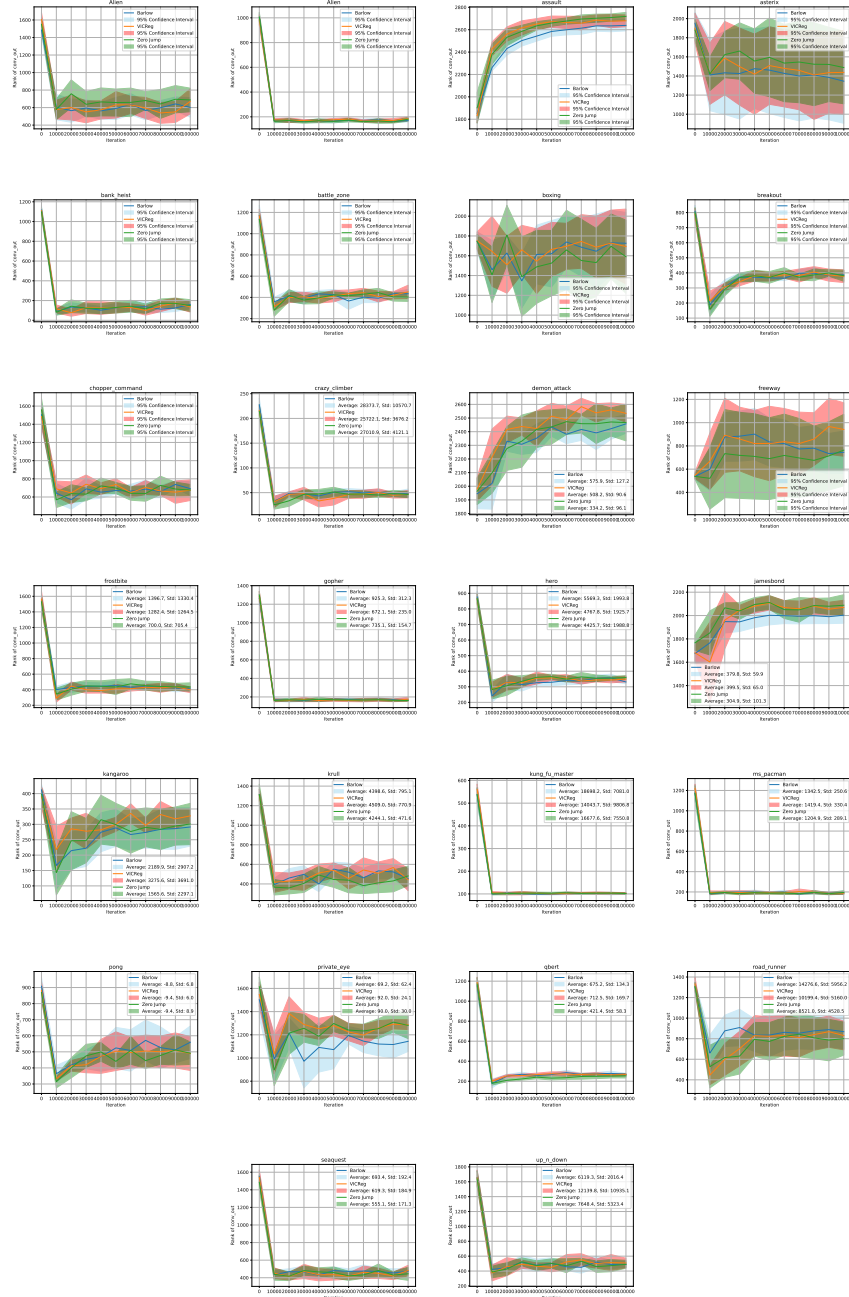
Figure 7: Rank of the output from the penultimate layer of the advantage head, measured every 10,000 steps and averaged across 10 different runs for every game.
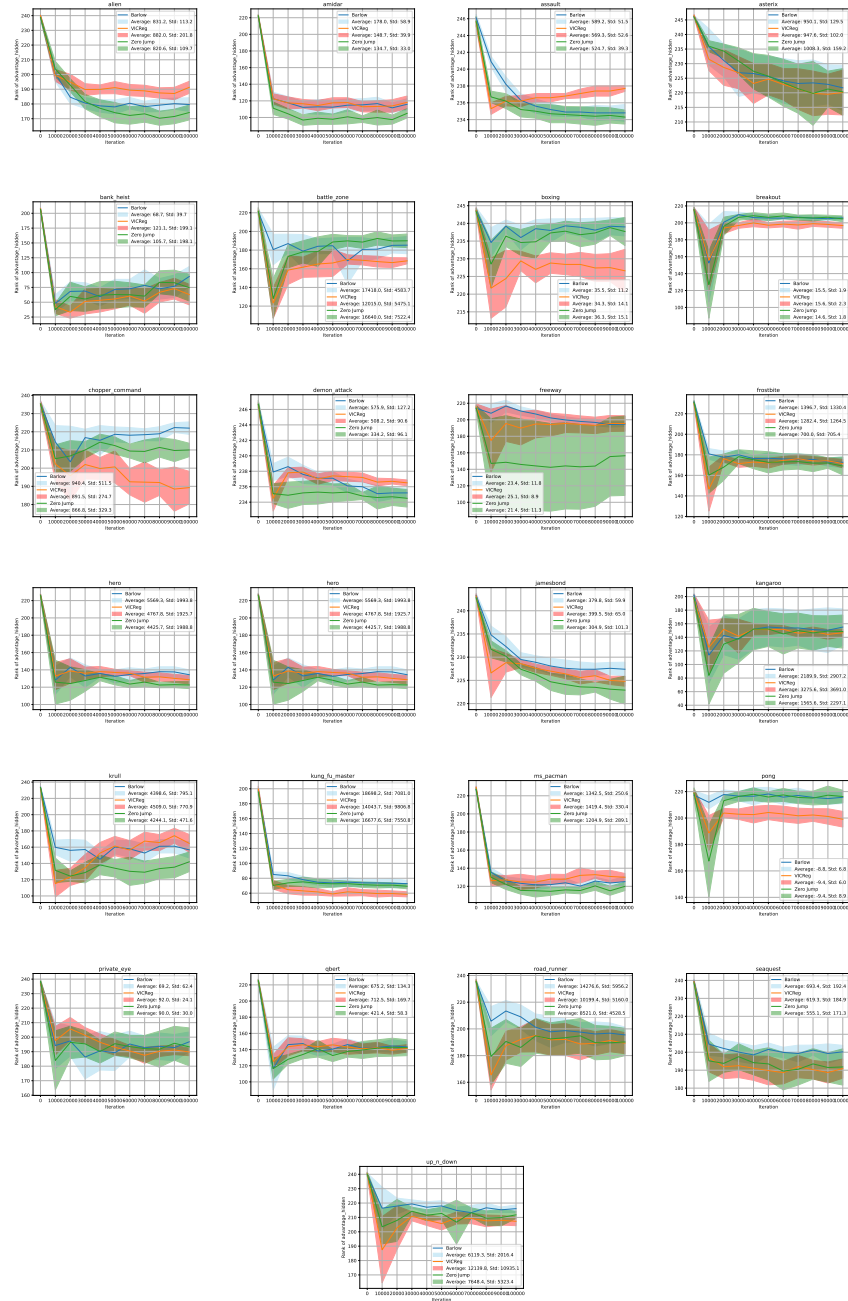
Figure 8: Fraction of dormant neurons averaged across 10 different runs for every game.