# REPRESENTATION LEARNING VIA INVARIANT CAUSAL MECHANISMS

**Jovana Mitrovic**    **Brian McWilliams**    **Jacob Walker**    **Lars Buesing**    **Charles Blundell**

DeepMind, London, UK
{mitrovic, bmcw, jcwalker, lbuesing, cblundell}@google.com

## ABSTRACT

Self-supervised learning has emerged as a strategy to reduce the reliance on costly supervised signal by pretraining representations only using unlabeled data. These methods combine heuristic proxy classification tasks with data augmentations and have achieved significant success, but our theoretical understanding of this success remains limited. In this paper we analyze self-supervised representation learning using a causal framework. We show how data augmentations can be more effectively utilized through explicit *invariance constraints* on the proxy classifiers employed during pretraining. Based on this, we propose a novel self-supervised objective, Representation Learning via Invariant Causal Mechanisms (RELIC), that enforces *invariant prediction* of proxy targets across augmentations through an invariance regularizer which yields improved generalization guarantees. Further, using causality we generalize contrastive learning, a particular kind of self-supervised method, and provide an alternative theoretical explanation for the success of these methods. Empirically, RELIC significantly outperforms competing methods in terms of robustness and out-of-distribution generalization on ImageNet, while also significantly outperforming these methods on Atari achieving above human-level performance on 51 out of 57 games.

## 1 INTRODUCTION

Training deep networks often relies heavily on large amounts of useful supervisory signal, such as labels for supervised learning or rewards for reinforcement learning. These training signals can be costly or otherwise impractical to acquire. On the other hand, unsupervised data is often abundantly available. Therefore, pretraining representations for unknown downstream tasks without the need for labels or extrinsic reward holds great promise for reducing the cost of applying machine learning models. To pretrain representations, self-supervised learning makes use of proxy tasks defined on unsupervised data. Recently, self-supervised methods using contrastive objectives have emerged as one of the most successful strategies for unsupervised representation learning (Oord et al., 2018; Hjelm et al., 2018; Chen et al., 2020a). These methods learn a representation by classifying every datapoint against all others datapoints (negative examples). Under assumptions on how the negative examples are sampled, minimizing the resulting contrastive loss has been justified as maximizing a lower bound on the mutual information (MI) between representations (Poole et al., 2019). However, (Tschannen et al., 2019) has shown that performance on downstream tasks may be more tightly correlated with the choice of encoder architecture than the achieved MI bound, highlighting issues with the MI theory of contrastive learning. Further, contrastive approaches compare different views of the data (usually under different data augmentations) to calculate similarity scores. This approach to computing scores has been empirically observed as a key success factor of contrastive methods, but has yet to be theoretically justified. This lack of a solid theoretical explanation for the effectiveness of contrastive methods hinders their further development.

To remedy the theoretical shortcomings, we analyze the problem of self-supervised representation learning through a causal lens. We formalize intuitions about the data generating process using a causal graph and leverage causal tools to derive properties of the optimal representation. We show that a representation should be an *invariant predictor* of proxy targets under interventions on features that are only correlated, but not causally related to the downstream targets of interest.

Since neither causally nor purely correlationally related features are observed and thus performing actual interventions on them is not feasible, for learning representation with this property we use data augmentations to simulate a subset of possible interventions. Based on our causal interpretation, we propose a regularizer which enforces that the prediction of the proxy targets is invariant across data augmentations. We propose a novel objective for self-supervised representation learning called REpresentation Learning with Invariant Causal mechanisms (RELIC). We show how this explicit invariance regularization leverages augmentations more effectively than previous self-supervised methods and that representations learned using RELIC are guaranteed to generalize well to downstream tasks under weaker assumptions than those required by previous work (Saunshi et al., 2019).

Next we generalize contrastive learning and provide an alternative theoretical explanation to MI for the success of these methods. We generalize the proxy task of instance discrimination commonly used in contrastive learning using the causal concept of *refinements* (Chalupka et al., 2014). Intuitively, a refinement of a task can be understood as a more fine-grained variant of the original problem. For example, a refinement for classifying cats against dogs would be the task of classifying individual cat and dog breeds. The instance discrimination task results from the most fine-grained refinement, e.g. discriminating individual cats and dogs from one another. We show that using refinements as proxy tasks enables us to learn useful representations for downstream tasks. Specifically, using causal tools, we show that learning a representation on refinements such that it is an invariant predictor of proxy targets across augmentations is a *sufficient condition* for these representations to generalize to downstream tasks (cf. Theorem 1). In summary, we provide theoretical support both for the general form of the contrastive objective as well as for the use of data augmentations. Thus, we provide an alternative explanation to mutual information for the success of recent contrastive approaches namely that of causal refinements of downstream tasks.

We test RELIC on a variety of prediction and reinforcement learning problems. First, we evaluate the quality of representations pretrained on ImageNet with a special focus on robustness and out-of-distribution generalization. RELIC performs competitively with current state-of-the-art methods on ImageNet, while significantly outperforming competing methods on robustness and out-of-distribution generalization of the learned representations when tested on corrupted ImageNet (ImageNet-C (Hendrycks & Dietterich, 2019)) and a version of ImageNet that consist of different renditions of the same classes (ImageNet-R (Hendrycks et al., 2020)). In terms of robustness, RELIC also significantly outperforms the supervised baseline with an absolute reduction of $4.9\%$ in error. Unlike much prior work that specifically focuses on computer vision tasks, we test RELIC for representation learning in the context of reinforcement learning on the Atari suite (Bellemare et al., 2013). There we find that RELIC significantly outperforms competing methods and achieves above human-level performance on $51$ out of 57 games.

**Contributions.**

- We formalize problem of self-supervised representation learning using causality and propose to more effectively leverage data augmentations through invariant prediction.

- We propose a new self-supervised objective, REpresentation Learning with Invariance Causal mechanisms (RELIC), that enforces invariant prediction through an explicit regularizer and show improved generalization guarantees.

- We generalize contrastive learning using refinements and show that learning on refinements is a sufficient condition for learning useful representations; this provides an alternative explanation to MI for the success of contrastive methods.

## 2 REPRESENTATION LEARNING VIA INVARIANT CAUSAL MECHANISMS

**Problem setting.** Let $X$ denote the unlabelled observed data and $\mathcal{Y} = \{Y_t\}_{t=1}^T$ be a set of unknown tasks with $Y_t$ denoting the targets for task $t$. The tasks $\{Y_t\}_{t=1}^T$ can represent both a multi-environment as well as a multi-task setup. Our goal is to pretrain with unsupervised data a representation $f(X)$ that will be useful for solving the downstream tasks $\mathcal{Y}$.

**Causal interpretation.** To effectively leverage common assumptions and intuitions about data generation of the unknown downstream tasks for the learning algorithm, we propose to formalize
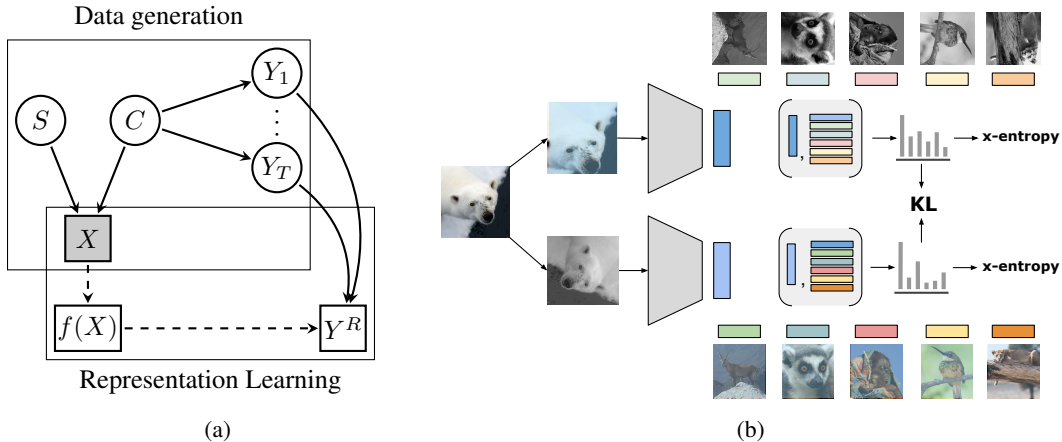
Figure 1: **(a)** Causal graph formalizing assumptions about content and style of the data and the relationship between targets and proxy tasks. **(b)** RELIC objective. KL refers to the Kullback-Leibler divergence, while x-entropy denotes cross entropy.

them using a causal graph. We start from the following assumptions: a) the data is generated from *content* and *style* variables, with b) only content (and not style) being relevant for the unknown downstream tasks and c) content and style are independent, i.e. style changes are content-preserving. For example, when classifying dogs against giraffes from images, different parts of the animals constitute content, while style could be, for example, background, lighting conditions and camera lens characteristics. By assumption, content is a good representation of the data for downstream tasks and we therefore cast the goal of representation learning as estimating content. In the following, we compactly formalize these assumptions with a causal graph[1], see Figure 1a.

Let $C$ and $S$ be the latent variables describing content and style. In Figure 1a, the directed arrows from $C$ and $S$ to the observed data $X$ (e.g. images) indicate that $X$ is generated based on content and style. The directed arrow from $C$ to the target $Y_t$ (e.g. class labels) encodes the assumption that content directly influences the target tasks, while the absence of any directed arrow from $S$ to $Y_t$ indicates that style does not. Thus, content $C$ has all the necessary information to predict $Y_t$. The absence of any directed path between $C$ and $S$ in Figure 1a encodes the intuition that these variables are independent, i.e. $C \perp\!\!\!\perp S$.

Using the independence of mechanisms (Peters et al., 2017), we can conclude that under this causal model performing interventions on $S$ does not change the conditional distribution $P(Y_t|C)$, i.e. manipulating the value of $S$ does not influence this conditional distribution. Thus, $P(Y_t|C)$ is invariant under changes in style $S$. We call $C$ an *invariant representation* for $Y_t$ under $S$, i.e.

$$p^{do(S=s_i)}(Y_t \,|\, C) = p^{do(S=s_j)}(Y_t \,|\, C) \quad \forall\, s_i, s_j \in \mathcal{S}, \tag{1}$$

where $p^{do(S=s)}$ denotes the distribution arising from assigning $S$ the value $s$ with $\mathcal{S}$ the domain of $S$ (Pearl, 2009). Specifically, using $C$ as a representation allows for us to predict targets stably across perturbations, i.e. content $C$ is both a useful and robust representation for tasks $\mathcal{Y}$.

Since the targets $Y_t$ are unknown, we will construct a proxy task $Y^R$ in order to learn representations from unlabeled data $X$ only. In order to learn useful representations for $Y_t$, we will construct proxy tasks that represents more fine-grained problems that $Y_t$; for a more formal treatment of proxy tasks please refer to Section 3. Further, to learn invariant representations, such as $C$, we enforce Equation 1 which requires us to observe data under different style interventions, i.e. we need data that describes the same content under varying style. Since we do not have access to $S$, to simulate style variability we use content-preserving data augmentations (e.g. rotation, grayscaling, translation, cropping for images). Specifically, we utilize *data augmentations as interventions on the style variable $S$*, i.e. applying data augmentation $a_i$ corresponds to intervening on $S$ and setting it to $s_{a_i}$. [2] Although

---

[1]See (Peters et al., 2017) for a review of causal graphs and causality.

[2]Since neither content nor style are a priori known, choosing a set of augmentations implicitly defines which aspects of the data are considered style and which are content.

we are not able to generate all possible styles using a fixed set of data augmentations, we will use augmentations that generate large sets of diverse styles as this allows us to learn better representations. Note that the heuristic of estimating similarity based on different views from contrastive learning can be interpreted as an implicit invariance constraint.

**RELIC objective.** Equation 1 provides a general scheme to estimate content (c.f. Figure 1a). We operationalize this by proposing to learn representations such that prediction of proxy targets from the representation is invariant under data augmentations. The representation $f(X)$ must fulfill the following *invariant prediction* criteria

$$(\textit{Invariant prediction}) \qquad p^{do(a_i)}(Y^R|f(X)) = p^{do(a_j)}(Y^R|f(X)) \quad \forall a_i, a_j \in \mathcal{A}. \qquad (2)$$

$\mathcal{A} = \{a_1, \ldots, a_m\}$ is the set of data augmentations which *simulate* interventions on the style variables and $p^{do(a)}$ denotes $p^{do(S=s_a)}$.

To achieve invariant prediction, we propose to explicitly enforce invariance under augmentations through a regularizer. This gives rise to an objective for self-supervised learning we call Representation Learning via Invariant Causal Mechanisms (RELIC). We write this objective as

$$\mathbb{E}_X \mathbb{E}_{\substack{a_{lk}, a_{qt} \\ \sim \mathcal{A} \times \mathcal{A}}} \sum_{b \in \{a_{lk}, a_{qt}\}} \mathcal{L}_b(Y^R, f(X)) \quad s.t. \quad KL\left(p^{do(a_{lk})}(Y^R | f(X)), p^{do(a_{qt})}(Y^R | f(X))\right) \leq \rho$$

where $\mathcal{L}$ is the proxy task loss and $KL$ is the Kullback-Leibler (KL) divergence. Note that any distance measure on distributions can be used in place of the KL divergence. We explain the remaining terms in detail below.

Concretely, as proxy task we associate to every datapoint $x_i$ the label $y_i^R = i$. This corresponds to the instance discrimination task, commonly used in contrastive learning (Hadsell et al., 2006). We take pairs of points $(x_i, x_j)$ to compute similarity scores and use pairs of augmentations $a_{lk} = (a_l, a_k) \in \mathcal{A} \times \mathcal{A}$ to perform a style intervention. Given a batch of samples $\{x_i\}_{i=1}^N \sim \mathcal{D}$, we use

$$p^{do(a_{lk})}(Y^R = j \mid f(x_i)) \propto \exp\left(\phi(f(x_i^{a_l}), h(x_j^{a_k}))/\tau\right).$$

with $x^a$ data augmented with $a$ and $\tau$ a softmax temperature parameter. We encode $f$ using a neural network and choose $h$ to be related to $f$, e.g. $h = f$ or as a network with an exponential moving average of the weights of $f$ (e.g. target networks similar to (Grill et al., 2020)). To compare representations we use the function $\phi(f(x_i), h(x_j)) = \langle g(f(x_i)), g(h(x_j))\rangle$ where $g$ is a fully-connected neural network often called the critic.

Combining these pieces, we learn representations by minimizing the following objective over the full set of data $x_i \in \mathcal{D}$ and augmentations $a_{lk} \in \mathcal{A} \times \mathcal{A}$

$$-\sum_{i=1}^N \sum_{a_{lk}} \log \frac{\exp\left(\phi(f(x_i^{a_l}), h(x_i^{a_k}))/\tau\right)}{\sum_{m=1}^M \exp\left(\phi(f(x_i^{a_l}), h(x_m^{a_k}))/\tau\right)} + \alpha \sum_{a_{lk}, a_{qt}} KL(p^{do(a_{lk})}, p^{do(a_{qt})}) \qquad (3)$$

with $M$ the number of points we use to construct the contrast set and $\alpha$ the weighting of the invariance penalty. We used the shorthand $p^{do(a)}$ for $p^{do(a)}(Y^R = j \mid f(x_i))$. With appropriate choices for $\phi$, $g$, $f$ and $h$ above, Equation 3 recovers many recent state-of-the-art methods (c.f. Table 5 in Section A). Figure 1b presents a schematic of the RELIC objective.

The explicit invariance penalty encourages the within-class distances (for a downstream task of interest) of the representations learned by RELIC to be tightly concentrated. We show this empirically in Figure 2 and theoretically in Appendix B. In the following section we provide theoretical justification for using an instance discrimination-based contrastive loss using a causal perspective. We also show (cf. Theorem 1 below) that minimizing the contrastive loss alone (i.e. $\alpha = 0$) does not guarantee generalization. Instead, invariance across augmentations must be explicitly enforced.

## 3 GENERALIZING CONTRASTIVE LEARNING

**Learning with refinements.** In contrastive learning, the task of instance discrimination, i.e. classifying the dataset $\{(x_i, y_i^R = i)|x_i \in \mathcal{D}\}$, is used as the proxy task. To better understand contrastive

learning and motivate this proxy task, we generalize instance discrimination using the causal concept of *refinements* (Chalupka et al., 2014). Intuitively, a refinement of one problem is another more fine-grained problem. If task $Y_t$ is to classify cats against dogs, then a refinement of $Y_t$ is the task of classifying cats and dogs into their individual breeds. See Figure 4 for a further visual example. For any set of tasks, there exist many different refinements. However, the most fine-grained refinement corresponds exactly to classifying the dataset $\{(x_i, y_i^R = i) | x_i \in \mathcal{D}\}$. Thus, the instance discrimination task used in contrastive learning is a specific type of refinement. For a definition and formal treatment of refinements please refer to Appendix D.

Let $Y^R$ be targets of a proxy task that is a refinement for all tasks in $\mathcal{Y}$. Leveraging causal tools, we connect learning on refinements to learning on downstream tasks. Specifically, we provide a theoretical justification for exchanging unknown downstream tasks with these specially constructed proxy tasks. We show that if $f(X)$ is an invariant representation for $Y^R$ under changes in style $S$, then $f(X)$ is also an invariant representation for tasks in $\mathcal{Y}$ under changes in style $S$. Thus by enforcing invariance under style interventions on a refinement, we learn representations that generalize to downstream tasks.[3] This is summarized in the following theorem.

**Theorem 1.** *Let $\mathcal{Y} = \{Y_t\}_{t=1}^T$ be a family of downstream tasks. Let $Y^R$ be a refinement for all tasks in $\mathcal{Y}$. If $f(X)$ is an invariant representation for $Y^R$ under style interventions $S$, then $f(X)$ is an invariant representation for all tasks in $\mathcal{Y}$ under style interventions $S$, i.e.*

$$p^{do(s_i)}(Y^R \,|\, f(X)) = p^{do(s_j)}(Y^R \,|\, f(X)) \quad \Rightarrow \quad p^{do(s_i)}(Y_t \,|\, f(X)) = p^{do(s_j)}(Y_t \,|\, f(X)) \quad (4)$$

*for all $s_i, s_j \in \mathcal{S}$ with $p^{do(s_i)} = p^{do(S=s_i)}$. Thus, $f(X)$ is a representation that generalizes to $\mathcal{Y}$.*

Theorem 1 states that if $Y^R$ is a refinement of $\mathcal{Y}$ then learning a representation on $Y^R$ is a *sufficient* condition for this representation to be useful on $\mathcal{Y}$. For a formal exposition of these points and accompanying proofs, please refer to Appendix D. Recall that the instance discrimination proxy task is the most fine-grained refinement, and so the left hand side of 4 is satisfied for any downstream task satisfying the stated assumptions of the theorem.

We generalize contrastive learning through refinements and connect representations learned on refinements and downstream tasks in Theorem 1. Thus, using causality we provide an alternative explanation to mutual information for the success of contrastive learning. Note that our methodology of refinements is not limited to instance discrimination tasks and is thus more general than currently used contrastive losses. Real world data often includes rich sources of metadata which can be used to guide the construction of refinements by grouping the data according to any available meta-data. Note that the coarser we can create a refinement, the more data efficient we can expect to be when learning representations for downstream tasks. Further, we can also expect to require less supervised data to finetune the representation.

## 4   RELATED WORK

**Contrastive objectives and mutual information maximization.** Many recent approaches to self-supervised learning are rooted in the well-established idea of maximizing mutual information (MI), e.g. Contrastive Predictive Coding (CPC) (Oord et al., 2018; Hénaff et al., 2019), Deep InfoMax (DIM) (Hjelm et al., 2018) and Augmented Multiscale DIM (AMDIM) (Bachman et al., 2019). These methods are based on noise contrastive estimation (NCE) (Gutmann & Hyvärinen, 2010) which, under specific conditions, can be viewed as a bound on MI (Poole et al., 2019). The resulting objective functions are commonly referred to as InfoNCE.

The precise role played by mutual information maximization in self-supervised learning is subject to some debate. (Tschannen et al., 2019) argue that the performance on downstream tasks is not correlated with the achieved bound on MI, but may be more tightly correlated with encoder architecture and capacity. Importantly, InfoNCE objectives require custom architectures to ensure the network does not converge to non-informative solutions thus precluding the use of standard architectures. Recently, several works (He et al., 2019; Chen et al., 2020a) successfully combined

---

[3]Note that since refinements are more fine-grained that the original task, if a representation captures a refinement then it also captures the downstream tasks as strictly more information is needed to solve the refinement.

contrastive estimation with a standard ResNet-50 architecture. In particular, SimCLR (Chen et al., 2020a) relies on a set of *strong* augmentations[4], while (He et al., 2019) uses a memory bank. Inspired by target networks in reinforcement learning, (Grill et al., 2020) proposed BYOL: an algorithm for self-supervised learning which remarkably does not use a contrastive objective. Although theoretical explanation for the good performance of BYOL is presently missing, interestingly the objective, an $\ell_2$ distance between two different embeddings of the input data resembles the $\ell_2$ form of our regularizer proposed in Equation 5 in Appendix B.

Recently, (Saunshi et al., 2019) proposed a learning theoretic framework to analyze the performance of contrastive objectives. However, without strong assumptions on intra-class concentration they note that contrastive objectives are fundamentally limited in the representations they are able to learn. RELIC explicitly enforces intra-class concentration via the invariance regularizer, ensuring that it generalizes under weaker assumptions. Unlike (Saunshi et al., 2019) which do not discuss augmentations, we incorporate augmentations into our theoretical explanation of contrastive methods.



Figure 2: Distribution of the linear discriminant ratio ($F_{\mathrm{LDA}}$, see text) of $f$ for RELIC, SimCLR and AMDIM ($y$-axis clipped to aid visualization).

The reasons for the improvement in performance from AMDIM through to SimCLR and BYOL are not easily explained by either the MI maximization or the learning theoretic viewpoint. Further, it is not clear why relatively minor architectural differences between the methods result in significant differences in performance nor is it obvious how current state-of-the-art can be improved. In contrast to prior art, the performance of RELIC is explained by connections to causal theory. As such it gives a clear path for improving results by devising problem appropriate refinements, interventions and invariance penalties. Furthermore, the use of invariance penalties in RELIC as dictated by causal theory yields significantly more robust representations that generalize better than those learned with SimCLR or BYOL.

**Causality and invariance.** Recently, the notion of invariant prediction has emerged as an important operational concept in causal inference (Peters et al., 2016). This idea has been used to learn classifiers which are robust against domain shifts (Gong et al., 2016). Notably, (Heinze-Deml & Meinshausen, 2017) propose to use group structure to delineate between different environments where the aim is to minimize the classification loss while also ensuring that the conditional variance of the prediction function within each group remains small. Unlike (Heinze-Deml & Meinshausen, 2017) who use supervised data and rely on having a grouping in the training data, our approach does not rely on ground-truth targets and can flexibly create groupings of the training data if none are present. Further, we enforce invariant prediction within the group by constraining the distance between distributions resulting from contrasting data across groups.

## 5 EXPERIMENTS

We first visualize the influence of the explicit invariance constraint in RELIC on the linear separability of the learned representations. We then evaluate RELIC on a number of prediction and reinforcement learning tasks for usefulness and robustness. For the prediction tasks, we test RELIC after pretraining the representation in a self-supervised way on the training set of the ImageNet ILSVRC-2012 dataset (Russakovsky et al., 2015). We evaluate RELIC in the linear evaluation setup on ImageNet and test its robustness and out-of-distribution generalization on datasets related to ImageNet. Unlike much prior work in contrastive learning which focuses specifically on computer vision tasks, we test RELIC also in the context of learning representations for reinforcement learning. Specifically, we test RELIC on the suite of Atari games (Bellemare et al., 2013) which consists of 57 diverse games of varying difficulty.

**Linear evaluation.** In order to understand how representations learned by RELIC differ from other methods, we compare it against those learned by AMDIM and SimCLR in terms of Fischer's *linear discriminant ratio* (Friedman et al., 2009): $F_{\mathrm{LDA}} = \|\mu_k - \mu_{k'}\|^2 / \sum_{i,j \in \mathcal{C}_k} \|f(x_i) - f(x_j)\|^2$ where
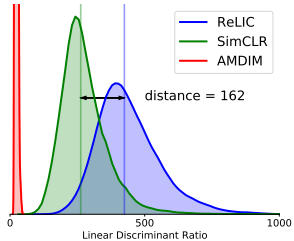
---

[4]The set of augmentations includes Gaussian blurring, various colour distortions, flips and random cropping.

$\mu_k = \frac{1}{|\mathcal{C}_k|} \sum_{i \in \mathcal{C}_k} f(x_i)$ is the mean of the representations of class $k$ and $\mathcal{C}_k$ is the index set of that class. A larger $F_{\text{LDA}}$ implies that classes are more easily separated with a linear classifier. This can be achieved by either increasing distances between classes (numerator) or shrinking within-class variance (denominator).

Figure 2 shows the distribution of $F_{\text{LDA}}$ for RELIC, SimCLR and AMDIM after training as measured on the (downsampled) ImageNet validation set. The distance between medians of RELIC and SimCLR is 162. AMDIM is tightly concentrated close to 20. The invariance penalty ensures that—even though labels are *a-priori* unknown—for RELIC within-class variability of $f$ is concentrated leading to better linear separability between classes in the downstream task of interest. This is reflected in the rightward shift of the distribution of $F_{\text{LDA}}$ in Figure 2 for RELIC compared with SimCLR and AMDIM which do not impose such a constraint.

Next we evaluate RELIC's representation by training a linear classifier of top of the fixed encoder following the procedure in (Kolesnikov et al., 2019; Chen et al., 2020a) and Appendix E.4. In Table 1, we report top-1 and top-5 accuracy on the ImageNet test set. Methods denoted with * use SimCLR augmentations (Chen et al., 2020a), while methods denoted † use custom, stronger augmentations. Comparing methods which use SimCLR augmentations, RELIC outperforms competing approaches on both ResNet-50 and ResNet-50 with target network. For completeness, we report results for SwAV (Caron et al., 2020) and InfoMin (Tian et al., 2020), but note that these methods use stronger augmentations which alone have been shown to boost performance by over 5%. A fair comparison between different objectives can only be achieved under the same architecture and the same set of augmentations.

Table 1: Accuracy (in %) under linear evaluation on ImageNet for different self-supervised representation learning methods. Methods with * use SimCLR augmentations. Methods with † use custom, stronger augmentations.

| Method | | Top-1 | Top-5 |
|---|---|---|---|
| *ResNet-50 architecture* | | | |
| PIRL (Misra & Maaten, 2020) | | 63.6 | - |
| CPC v2 (Hénaff et al., 2019) | | 63.8 | 85.3 |
| CMC (Tian et al., 2019) | | 66.2 | 87.0 |
| SimCLR (Chen et al., 2020a) | * | 69.3 | 89.0 |
| SwAV (Caron et al., 2020) | * | 70.1 | - |
| RELIC (ours) | * | 70.3 | 89.5 |
| InfoMin Aug. (Tian et al., 2020) | † | 73.0 | 91.1 |
| SwAV (Caron et al., 2020) | † | 75.3 | - |
| *ResNet-50 with target network* | | | |
| MoCo v2 (Chen et al., 2020b) | | 71.1 | - |
| BYOL (Grill et al., 2020) | * | 74.3 | 91.6 |
| RELIC (ours) | * | 74.8 | 92.2 |

**Robustness and generalization.** We evaluate robustness and out-of-distribution generalization of RELIC's representation on datasets Imagenet-C (Hendrycks & Dietterich, 2019) and ImageNet-R (Hendrycks et al., 2020), respectively. To evaluate RELIC's representation, we train a linear classifier on top of the frozen representation following the procedure described in (Chen et al., 2020a) and appendix E.5.2. For Imagenet-C we report the mean Corruption Error (mCE) and Corruption Errors for Noise corruptions in Table 3. RELIC has significantly lower mCE than both the supervised ResNet-50 baseline and the unsupervised methods SimCLR and BYOL. Also, it has the lowest Corruption Error on 14 out of 15 corruptions when compared to SimCLR and BYOL. Thus, we see that RELIC learns the most robust representation. RELIC also outperforms SimCLR and BYOL on ImageNet-R showing its superior out-of-distribution generalization ability; see Table 2. For further details and results please consult E.5.

Table 2: Top-1 error rates for different self-supervised representation learning methods on ImageNet-R. All models are trained only on clean ImageNet images and RELIC$_T$ refers to RELIC using a ResNet-50 with target network as in BYOL (Grill et al., 2020).

| Method | Supervised | SimCLR | RELIC (ours) | BYOL | RELIC$_T$ (ours) |
|---|---|---|---|---|---|
| Top-1 Error (%) | 63.9 | 81.7 | 77.4 | 77.0 | 76.2 |

Table 3: Mean Corruption Error (mCE), mean relative Corruption Error (mrCE) and Corruption Errors for the "Noise" class of corruptions (Gaussian, Shot, Impulse) on ImageNet-C. The mCE value is the average across 75 different corruptions. Methods are trained only on clean ImageNet images.

| Method | mCE | mrCE | Gaussian | Shot | Impulse |
|---|---|---|---|---|---|
| Supervised | 76.7 | 105.0 | 80.0 | 82.0 | 83.0 |
| *ResNet-50 architecture:* | | | | | |
| SimCLR | 87.5 | 111.9 | 79.4 | 81.9 | 89.6 |
| ReLIC (ours) | 76.4 | **87.7** | 67.8 | 70.7 | 77.0 |
| *ResNet-50 with target network:* | | | | | |
| BYOL | 72.3 | 90.0 | 65.9 | 68.4 | 73.7 |
| ReLIC (ours) | **70.8** | 88.4 | **63.6** | **65.7** | **69.2** |

**Reinforcement Learning.** Much prior work in contrastive learning has focused specifically on computer vision tasks. In order to compare these approaches in a different domain, we investigate representation learning in the context of reinforcement learning. We compare RELIC as an auxiliary loss against other state of the art self-supervised losses on an agent trained on 57 Atari games. Using human normalized scores as a metric, we use the original architecture and hyperparameters of the R2D2 agent (Kapturowski et al., 2019) and supplement it with a second encoder trained with a given representation learning loss. When auxiliary losses are present, the Q-Network takes the output of the second encoder as an input. The Q-Network and the encoder are trained with separate optimizers. For the augmentation baseline, the Q-Network takes two identical encoders trained end-to-end. Table 4 shows a comparison between RELIC, SimCLR, BYOL, CURL (Srinivas et al., 2020), and feeding augmented observations directly to the agent (Kostrikov et al., 2020). We find that RELIC has a significant advantage over competing self-supervised methods, performing best in 25 out of 57 games. The next best performing method, CURL performs best in 11 games. Full details are presented in Section E.6.

Table 4: Human Normalized Scores over 57 Atari Games.

| Atari Performance | RELIC | SimCLR | CURL | BYOL | Augmentation |
|---|---|---|---|---|---|
| Capped mean | **91.46** | 88.76 | 90.72 | 89.43 | 80.60 |
| Number of superhuman games | **51** | 49 | 49 | 49 | 34 |
| Mean | **3003.73** | 2086.16 | 2413.12 | 1769.43 | 503.15 |
| Median | **832.50** | 592.83 | 819.56 | 483.39 | 132.17 |
| 40% Percentile | 356.27 | 266.07 | **409.46** | 224.80 | 94.35 |
| 30% Percentile | **202.49** | 174.19 | 190.96 | 150.21 | 80.04 |
| 20% Percentile | **133.93** | 120.84 | 126.10 | 118.36 | 57.95 |
| 10% Percentile | **83.79** | 37.19 | 59.09 | 44.14 | 32.74 |
| 5% Percentile | **20.87** | 12.74 | 20.56 | 7.75 | 2.85 |

## 6   CONCLUSION

In this work we have analyzed self-supervised learning using a causal framework. Using a causal graph, we have formalized the problem of self-supervised representation learning and derived properties of the optimal representation. We have shown that representations need to be invariant predictors of proxy targets under interventions on features that are only correlated, but not causally related to the downstream tasks. We have leveraged data augmentations to simulate these interventions and have proposed to explicitly enforce this invariance constraint. Based on this, we have proposed a new self-supervised objective, Representation Learning via Invariant Causal Mechanisms (RELIC), that enforces invariant prediction of proxy targets across augmentations using an invariance regularizer. Further, we have generalized contrastive methods using the concept of refinements and have shown that learning a representation on refinements using the principle of invariant prediction is a sufficient condition for these representations to generalize to downstream tasks. With this, we have provided an alternative explanation to mutual information for the success of contrastive methods. Empirically

we have compared RELIC against recent self-supervised methods on a variety of prediction and reinforcement learning tasks. Specifically, RELIC significantly outperforms competing methods in terms of robustness and out-of-distribution generalization of the representations it learns on ImageNet. RELIC also significantly outperforms related self-supervised methods on the Atari suite achieving superhuman performance on 51 out of 57 games. We aim to investigate the construction of more coarse-grained refinements and the empirical evaluation of different kinds of refinements in future work.

## REFERENCES

Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. In *Advances in Neural Information Processing Systems*, pp. 15509–15519, 2019.

Marc G. Bellemare, Yavar Naddaf, J. Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents (extended abstract). *J. Artif. Intell. Res.*, 47:253–279, 2013.

M. Caron, I. Misra, J. Mairal, Priya Goyal, P. Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *ArXiv*, abs/2006.09882, 2020.

Krzysztof Chalupka, Pietro Perona, and Frederick Eberhardt. Visual causal feature learning. *arXiv preprint arXiv:1412.2309*, 2014.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020a.

Xinlei Chen, Haoqi Fan, Ross B. Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *ArXiv*, abs/2003.04297, 2020b.

Paul Erdős and Alfréd Rényi. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci*, 5(1):17–60, 1960.

Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 2. Springer series in statistics New York, 2009.

Alan Frieze and Michał Karoński. *Introduction to random graphs*. Cambridge University Press, 2016.

Mingming Gong, Kun Zhang, Tongliang Liu, Dacheng Tao, Clark Glymour, and Bernhard Schölkopf. Domain adaptation with conditional transferable components. In *International conference on machine learning*, pp. 2839–2848, 2016.

Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.

Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 297–304, 2010.

Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pp. 1735–1742. IEEE, 2006.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, 2019.

Christina Heinze-Deml and Nicolai Meinshausen. Conditional variance penalties and domain shift robustness. *arXiv preprint arXiv:1710.11469*, 2017.

Olivier J Hénaff, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019.

Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.

Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. *arXiv preprint arXiv:2006.16241*, 2020.

R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2018.

S. Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ArXiv*, abs/1502.03167, 2015.

Steven Kapturowski, Georg Ostrovski, Will Dabney, John Quan, and Remi Munos. Recurrent experience replay in distributed reinforcement learning. *Iclr*, 2019.

A. Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual representation learning. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1920–1929, 2019.

Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing deep reinforcement learning from pixels. *arXiv preprint arXiv:2004.13649*, 2020.

I. Loshchilov and F. Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *ICLR*, 2017.

Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6707–6717, 2020.

V. Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.

Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

Judea Pearl. *Causality*. Cambridge university press, 2009.

Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):947–1012, 2016.

Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. MIT press, 2017.

Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. On variational bounds of mutual information. In *International Conference on Machine Learning*, pp. 5171–5180, 2019.

Olga Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Zhiheng Huang, A. Karpathy, A. Khosla, M. Bernstein, A. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115:211–252, 2015.

Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khandeparkar. A theoretical analysis of contrastive unsupervised representation learning. In *International Conference on Machine Learning*, pp. 5628–5637, 2019.

Aravind Srinivas, Michael Laskin, and Pieter Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning. *arXiv preprint arXiv:2004.04136*, 2020.

Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.

Yonglong Tian, C. Sun, Ben Poole, Dilip Krishnan, C. Schmid, and Phillip Isola. What makes for good views for contrastive learning. *ArXiv*, abs/2005.10243, 2020.

Michael Tschannen, Josip Djolonga, Paul K Rubenstein, Sylvain Gelly, and Mario Lucic. On mutual information maximization for representation learning. *arXiv preprint arXiv:1907.13625*, 2019.

Haohan Wang, Songwei Ge, E. Xing, and Zachary Chase Lipton. Learning robust global representations by penalizing local predictive power. *ArXiv*, abs/1905.13549, 2019.

Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.

## A  RELATIONSHIP BETWEEN ReLIC AND OTHER METHODS

Table 5: The objective in eq. (3) recovers state of the art methods depending on design choices ("-" denotes the identity function and "norml." means $g$ is constrained to have unit norm).

| Method | $\phi$ | $g$ | Regl. |
|---|---|---|---|
| CPC (Hénaff et al., 2019) | $\langle g, Wg \rangle$ | PixelCNN | - |
| AMDIM (Bachman et al., 2019) | $\langle \cdot, \cdot \rangle$ | - | - |
| SimCLR (Chen et al., 2020a) | $\langle g, g \rangle$ | MLP, norml. | - |
| BYOL (Grill et al., 2020) | - | $g_1, g_2$ 1 layer MLP, norml. | $\|g_1(g_2) - g_2\|^2$ |
| ReLIC (ours) | $\langle g, g \rangle$ | MLP, norml. | Eq. (3) |

## B  DISTANCE CONCENTRATION AND GENERALIZATION

Quantifying the generalization performance of representations learned on unlabelled data is a difficult task without imposing assumptions on the underlying structure of the data and the downstream tasks of interest. The results in (Saunshi et al., 2019) assume a latent class structure underlying the data. The similarity of images under each (potentially overlapping) latent class $c$ is measured by a probability distribution $\mathcal{D}_c$. In the contrastive setting a positive pair of points $\{x, x^+\}$ is said to be sampled from a distribution $\mathbb{E}_c \mathcal{D}_c(x) \mathcal{D}_c(x^+)$ and a negative example $x^-$ is sampled from the marginal distribution. The task of interest is multi-class classification using the learned representation. In our setting the augmented data points $\{x_i^{a_l}, x_i^{a_k}\}$ and $\{x_i^{a_l}, x_m^{a_k}\}_{m=1}^M$ take the roles of the pairs of positive and negative points, respectively.

In this section, under the same structural assumptions on the data as (Saunshi et al., 2019) we will show that a similar result holds but under weaker assumptions on the function, $f$.

To intuit the following results, we can view our explicit invariance constraint through the lens of distance concentration. Its effect can be seen intuitively in Figure 3. The shaded region represents the set of augmentations, $\mathcal{A}$ around an image. Depicted are two images $x_i$ and $x_j$ from the ImageNet class Stingray. The points $x_i^{a_l}$ and $x_j^{a_k}$ are augmentations which correspond to a region of overlap between the augmentation sets of $x_i$ and $x_j$. If the augmentations $f(x_i^{a_l})$ and $f(x_j^{a_k})$ are similar enough, encouraging $f(x_i)$ to be close to $f(x_i^{a_l})$ and similarly for $f(x_j)$ and $f(x_j^{a_k})$ indirectly encourages $f(x_i)$ to be close to $f(x_j)$. This has the effect of concentrating distances between similar images. We will make this intuition more formal in the following discussion.



Figure 3: Visual representation of invariance penalty. Shaded region denotes set of augmentations around an image.

Consider a modified, Euclidean distance regularized version of our objective

$$\hat{f} \in \arg\min_{f \in \mathcal{F}} \sum_{i=1}^{N} \sum_{a_{lk}} \ell(\{f(x_i^{a_l})^\top (f(x_i^{a_k}) - f(x_m^{a_k}))\}_{m=1}^M) \quad (5)$$
$$s.t. \quad \|f(x_i) - f(x_i^{a_k})\|^2 \leq \rho.$$

where $f \in \mathcal{F} = \{f : \mathcal{X} \mapsto \mathbb{R}^d \ s.t. \ \|f\|_2 \leq T\}$ with $T \geq 0$. Here $\ell(v) = \log(1 + \sum_m \exp(v_m))$ is the logistic loss. For a single negative, this is equivalent to the standard ReLIC objective with an identity critic.

**Assumptions.**  We require that the following assumptions hold: **(A1)** $\hat{f}$ is $L$-Lipschitz and minimizes eq. (5) such that the constraint is active and **(A2)** x is a bounded variable.

**Lemma 1** (Concentration). *If assumption (A1) holds for $\rho \leq \frac{B}{6L\kappa}$, and (A2) holds for $x$, $\hat{f}(x)$ is a sub-Gaussian random variable with parameter $\sigma_{\hat{f}}^2 \leq \frac{1}{\kappa}\sigma_x^2$.*
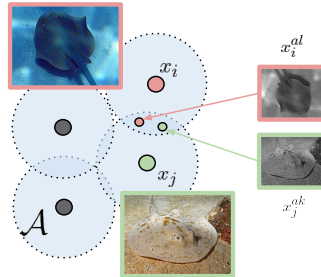
See Appendix C for proof. This result states that the Euclidean version of our invariance regularizer has the effect of contracting the within-class variance of the data. Figure 2 shows that this holds in practise for the original version of our objective in eq. (3). This guarantees that the following generalization result from (Saunshi et al., 2019) holds. For brevity we state an informal version of the Theorem with details deferred to the original publication.

**Theorem 2** (Generalization. Adapted from Lemma B.2. from (Saunshi et al., 2019)). *Let $L^{\mu}_{sup}(f)$ be the standard $(K+1)$-wise hinge loss of the linear classification function $W^{\mu}f$ whose $c^{th}$ column is $\mu_x = \frac{1}{|\mathcal{C}_c|}\sum_{i\in\mathcal{C}_c} f(x_i)$ the mean of representations corresponding to class c. Further, let $L^{\mu}_{\gamma(f),sup}(f)$ use the the hinge loss with margin $\gamma(f) = 1 + c'M\sigma_f(\sqrt{k} + \sqrt{\log\frac{M}{\epsilon}})$ with $c'$ constant and $M = \max_x \|f(x)\|$. If $\hat{f}$ is the minimizer of eq. (5) and if Assumptions **(A1)** and **(A2)** hold then with high probability*

$$L^{\mu}_{sup}(\hat{f}) \leq \gamma(f)L^{\mu}_{\gamma(f),sup}(f) + Gen_N + \epsilon \tag{6}$$

Here, $Gen_N$ is a standard generalization bound which depends on the Rademacher complexity of the function class $\mathcal{F}$ and the sample size, $N$.

For all practical purposes, the final generalization result is identical to (Saunshi et al., 2019) stating that $\hat{f}$—which is learned by minimizing a contrastive objective on unlabelled data—performs well on labelled data. However, this crucially depends on the intraclass concentration of the representation, that $f(x)$ is sub-Gaussian with parameter $\sigma_f^2$. Whereas in (Saunshi et al., 2019) this was assumed to hold, our Lemma 1 shows that the necessary concentration is ensured by our invariance penalty. Experimentally we see this property holds in practise (figure 2).

## C   ADDITIONAL RESULTS

*Proof of Lemma 1.* Assume the data $x$ is $\sigma_x^2$-sub-Gaussian. In practise this holds since $x$ is bounded. It immediately follows that $L$-Lipschitz function $f(x)$ sub-Gaussian with parameter at most $L$. Now we will characterize the reduction in variance from $x$ to $f$. Assume there is a ball of radius $B$ around each point such that for any augmentation $x_i^s$ of $x_i$ $\|x_i - x_i^s\|_2^2 \leq B$. By assumption (A1) we have that $\|f(x_i) - f(x_i^s)\|_2^2 \leq \rho$. This implies that for points $x_i$ and $x_j$ such that $\|x_i - x_j\|_2^2 \leq 2B$, there exists a region of overlap so that $\|f(x_i) - f(x_j)\|_2^2 \leq \|f(x_i) - f(x_i^s)\|_2^2 + \|f(x_i^s) - f(x_j)\|_2^2 \leq 2\rho$.

In practise this says that there are augmentations of $x_i$ which are sufficiently similar to augmentations of $x_j$ so that their representations should be similar, thereby driving $f(x_i)$ and $f(x_j)$ to be closer.

The variance of points in $f$ space is

$$\sigma_f^2 = \frac{1}{2N^2}\sum_i\sum_j \|f(x_i) - f(x_j)\|_2^2$$

The overlap $B < \|x_i - x_j\|_2^2 \leq 2B$ induces a graph where we say $j \in \mathcal{N}(i) \; \forall \; j$ s.t. $\|x_i - x_j\|_2^2 \leq 2B$. For $N$ samples we can decompose the variance as

$$\sigma_f^2 = \frac{1}{2N^2}\sum_i\sum_j \|f(x_i) - f(x_j)\|_2^2$$

$$= \frac{1}{2N^2}\sum_i\sum_{j\in\mathcal{N}(i)} \|f(x_i) - f(x_j)\|_2^2 + \sum_{j'\notin\mathcal{N}(i)} \|f(x_i) - f(x_{j'})\|_2^2$$

By smoothness of $f$ we always have that have $\|f(x_i) - f(x_{j'})\|_2^2 \leq L\|x_i - x_{j'}\|_2^2$. By the constraint we have that $\|f(x_i) - f(x_j)\|_2^2 \leq \frac{2\rho L}{B}\|x_i - x_j\|_2^2 \; \forall j \in \mathcal{N}(i)$ and for $\delta = \frac{2\rho L}{B} < 1$.

**Constant proportion overlap.**   Now, assuming that for each point $i$ there is a constant proportion of the points, $0 \leq \alpha \leq 1$ in the set $\mathcal{N}(i) \; \forall i$ we can obtain the following inequality

$$\sigma_f^2 = \frac{1}{2N^2}\sum_i\sum_j \|f(x_i) - f(x_j)\|_2^2$$

$$\leq \alpha\delta\sigma_x^2 + (1-\alpha)L\sigma_x^2$$

$$= (\alpha\delta + (1-\alpha)L)\sigma_x^2 \tag{7}$$

13

For $\sigma_f^2 \leq \sigma_x^2$ we require $(\alpha\delta + (1-\alpha)L) \leq 1$. Since both terms are positive we separately require $(1-\alpha)L \leq 1$:

$$(1-\alpha)L < 1$$
$$(1-\alpha) < \frac{1}{L}$$
$$\alpha > (1 - \frac{1}{L})$$

This condition makes sense since the larger $\alpha$, the fewer unconnected components in the graph. If the above holds, we also require $\alpha\frac{2\rho L}{B} < 1 - (1-\alpha)L$ to ensure the sum is bounded above by 1. This implies $\rho < \frac{(1-(1-\alpha)L)B}{2L\alpha}$.

However, $\alpha$ is a property of the augmentation set and not directly a user-controllable parameter so if $\alpha$ is too small or the function is not smooth enough, it might not be possible to set $\rho$ in such a way to induce contraction in $\sigma_f^2$.

In the next section we derive a tighter concentration based on the structure of random graphs which are induced by the connectivity between data points and their augmentations.

**Random graphs.** Consider the graph $G(V, E)$ induced by the constraints $(i, j) \in E \; \forall \; \|x_i - x_j\|_2^2 \leq 2B$. Call $\mathcal{N}(i)$ the set of neighbours of point $i$. For $N$ points, if there is a constant probability $\alpha$ that $j \in \mathcal{N}(i)$ then $G_{N,\alpha}$ is an Erdös-Renyi graph.

From Theorem 3, if $\alpha \geq \frac{c \log N}{N}$ for $c > 1$ then with high probability, there are *no* unconnected components in $G$. That is, every vertex in V is reachable from any other vertex in a finite number of steps. We can then decompose the contribution to the variance in terms of components in the graph that are adjacent and those which are reachable within a certain number of steps.

Let the degree—the shortest path—between any two points be at most $D$ we obtain the following refinement of eq. (7)

$$\sigma_f^2 = \frac{1}{2N^2} \sum_i \sum_j \|f(x_i) - f(x_j)\|_2^2$$
$$\leq \alpha\delta\sigma_x^2 + (1-\alpha)D\delta\sigma_x^2$$

From Theorem 4 we have with high probability that $3 \leq D \leq 4$. So for $\sigma_f^2 \leq \frac{1}{\kappa}\sigma_x^2$ with $\kappa \geq 1$ we require $\rho \leq \frac{B}{2L\kappa(\alpha+3(1-\alpha))} \leq \frac{B}{6L\kappa}$. $\qquad\square$

**Theorem 3** (Connectedness (Erdős & Rényi, 1960))**.** *If $p = \frac{c \log n}{n}$ where $c > 1$ with high probability then the graph $G(n, p)$ has no unconnected components.*

**Definition 1** (Diameter)**.** *For a connected graph, $G(V, E)$ the diameter $diam(G) = \max dist(v_i, v_j)$ where $dist(v_i, v_j)$ is the minimum number of edges in the path between $v_i$ and $v_j$.*

**Theorem 4** (Diameter of random graphs (Frieze & Karoński, 2016))**.** *Let $d \geq 2$ be a fixed positive integer. For $c > 0$ and*

$$p^d n^{d-1} = \log(n^2/c)$$

*Then $diam(G_{n,p}) \geq d$ with probability $\exp(-c/2)$ and $diam(G_{n,p}) \leq d + 1$ with probability $1 - \exp(-c/2)$.*

## D   GENERALIZING CONTRASTIVE LEARNING

### D.1   REFINEMENTS

On the unsupervised observed data $\mathcal{D}$, any task as defined by targets $Y_t$ induces an equivalence relation, i.e. $Y_t$ partitions $\mathcal{D}$ into equivalence classes. It divides $\mathcal{D}$ based on values of the target, $\mathcal{D} = \{\{x_a | y_a = y_i\}_{i=1}^M\}$ where $\{y_1, \ldots, y_M\}$ for some $M$ is the set of target values. Here the equivalence relation associates datapoints based on the value of the target they predict. For example,
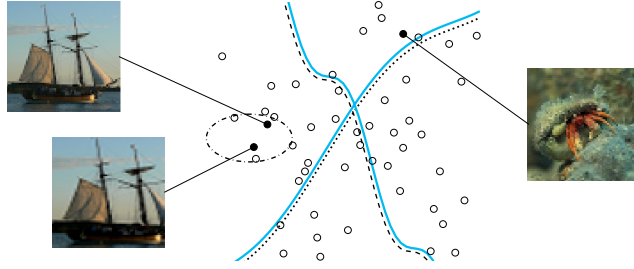
Figure 4: Visualization of a refinement of a set of tasks. The tasks are to classify aquatic vs non-aquatic life and animal vs non-animal with the individual class boundaries denoted by the dashed and dotted black lines. A refinement for these tasks is a to classify aquatic animal vs aquatic non-animal vs non-aquatic animal vs non-aquatic non-animal and the class boundaries are given in teal. The ellipse indicates the set of points induced by augmenting the image of the ship.

if $\mathcal{D}$ is a set of images of cats and dogs and $Y_t$ denotes labels cat and dog, then $\mathcal{D}$ is partitioned into two equivalence classes corresponding to cat and dog images by $Y_t$.

Intuitively, a refinement is a subdivision of an existing partition. For a visualization of a refinement of a set of tasks see Figure 4. To mathematically define refinements, we first need to introduce what it means for an equivalence relation to be finer than another equivalence relation.

**Definition 2. (Fineness).** *Let $\sim$ and $\approx$ be two equivalence relations on the set $\mathcal{D}$. If every equivalence class of $\sim$ is a subset of an equivalence class of $\approx$, we say that $\sim$ is finer than $\approx$.*

Now we define what refinements.

**Definition 3. (Refinement).** *Let $A$, $B$ be sets of equivalence classes induced by equivalence relations $\sim$ and $\approx$ over the set $\mathcal{D}$. If $\sim$ is finer than $\approx$, then we call $A$ a refinement of $B$.*

Furthermore, we can relate the corresponding sets of equivalence classes.

**Lemma 2.** *Let $\sim$ and $\approx$ be two equivalence relationships on the set $\mathcal{D}$ and denote the corresponding induced partitions by $A$ and $B$. If $\sim$ is finer than $\approx$, then every equivalence class of $\approx$ is a union of equivalence classes of $\sim$.*

Coming back to the example of cats and dogs, let $\approx$ be the relation that associates cats with cats and dogs with dogs. Now the relation $\sim$ which associated both cats and dogs with their specific breed (e.g. poodles with other poodles) is finer than $\approx$. Note that $\sim$ partitions $\mathcal{D}$ into breeds and so we can easily generate the sets of cats and dogs (i.e. equivalence classes of $\approx$) by taking a union over all the corresponding breeds.

### D.2 PROOF OF THEOREM 1

**Definition 4. (Invariant Representation).** *Let $X$ and $Y$ be the covariates and target, respectively. We call $f(X)$ an invariant representation for $Y$ under style $S$ if*

$$p^{do(S=s_i)}(Y \mid f(X)) = p^{do(S=s_j)}(Y \mid f(X)) \quad \forall s_i, s_j \in \mathcal{S}, \tag{8}$$

*where $do(S = s)$ denotes assigning $S$ the value $s$ and $\mathcal{S}$ is the domain of $S$.*

**Theorem 1. Theorem 1.** *Let $\mathcal{Y} = \{Y_t\}_{t=1}^T$ be a family of downstream tasks. Let $Y^R$ be a refinement for all tasks in $\mathcal{Y}$. If $f(X)$ is an invariant representation for $Y^R$ under changes in style $S$, then $f(X)$ is an invariant representation for all tasks in $\mathcal{Y}$ under changes in style $S$, i.e.*

$$p^{do(s_i)}(Y^R \mid f(X)) = p^{do(s_j)}(Y^R \mid f(X)) \quad \Rightarrow \quad p^{do(s_i)}(Y_t \mid f(X)) = p^{do(s_j)}(Y_t \mid f(X)) \tag{9}$$

*for all $t \in \{1, \ldots, T\}$ and for all $s_i, s_j \in \mathcal{S}$ with $p^{do(s_i)} = p^{do(S=s_i)}$. Thus, $f(X)$ is a representation that generalizes to $\mathcal{Y}$.*

15

**Proof.** Let $t \in \{1, \ldots, T\}$. We have

$$p^{do(s_i)}(Y_t|f(X)) = \int p^{do(s_i)}(Y_t|Y^R) p^{do(s_i)}(Y^R|f(X)) dY^R = \int p(Y_t|Y^R) p^{do(s_i)}(Y^R|f(X)) dY^R$$

$$= \int p(Y_t|Y^R) p^{do(s_j)}(Y^R|f(X)) dY^R = p^{do(s_j)}(Y_t|f(X)).$$

For the second and last equality, we used that the mechanism of $Y_t|Y^R$ is independent of $S$, i.e. $p^{do(s_i)}(Y_t|Y^R) = p^{do(s_j)}(Y_t|Y^R)$. The third equality follows from the assumption that $f(X)$ is an invariant representation for $Y^R$ under changes in $S$. Thus, we get that $f(X)$ is an invariant representation for $Y_t$ under changes in $S$. Specifically, for a representation to be an invariant representation for $Y_t$ it is a *sufficient condition* for it to be an invariant representation for $Y^R$. □

## E EXPERIMENTAL DETAILS

### E.1 IMAGE AUGMENTATIONS

For pretraining the representations in RELIC, we apply the augmentation scheme proposed in SimCLR (Chen et al., 2020a) and used in (Grill et al., 2020). This consists of the following augmentations applied in the order they are listed

- random crop – we randomly crop the image using an area randomly selected between $8\%$ and $100\%$ of the image with an logarithmically sampled aspect ration between $3/4$ and $4/3$. After this, we resize the patch to $224 \times 224$;
- random horizontal flip;
- color jittering – we apply in random order perturbations to brightness, contrast, saturation and hue of the image by shifting them by a random uniform offset;
- grayscale – we randomly apply grayscaling;
- Gaussian blurring – we blur the image using a $23 \times 23$ square Gaussian kernel with standard deviation uniformly sampled in $[0.1, 0.2]$;
- solarization – we transform all the pixels with $x \to x * 1_{\{x<0.5\}} + (1-x) * 1_{\{x\geq0.5\}}$.

We use the same parameters for the augmentations and probabilities of applying individual augmentations as SimCLR (Chen et al., 2020a). After applying augmentations, we normalize the images with the mean and standard deviation computed on ImageNet across the color channels.

### E.2 ARCHITECTURE

We test RELIC on two different architectures – ResNet-50 (He et al., 2016) and ResNet-50 with target network as in (Grill et al., 2020). For ResNet-50, we use version 1 with post-activation. We take the representation to be the output of the final average pooling layer, which is of dimension 2048. As in SimCLR (Chen et al., 2020a), we use a critic network to project the representation to a lower dimensional space with a multi-layer perceptron (MLP). When using ResNet-50 as encoder, we treat the parameters of the MLP (e.g. depth and width) as hyperparameters and sweep over them. This MLP has batch normalization (Ioffe & Szegedy, 2015) after every layer, rectified linear activations (ReLU) (Nair & Hinton, 2010). We used a 4 layer MLP with widths $[4096, 2048, 1024, 512]$ and output size 128 with ResNet-50. When using a ResNet-50 with target networks as in (Grill et al., 2020), we exactly follow their architecture settings.

### E.3 OPTIMIZATION

We use a batch size of 4096 and the LARS optimizer (You et al., 2017) with a cosine decay learning rate schedule (Loshchilov & Hutter, 2017) for 1000 epochs with 10 epochs for warm-up. We exclude the biases and batch normalization parameters from LARS adaptation. We use as the base learning rate 0.3 for ResNet-50 and 0.2 for ResNet-50 with target network. We scale this learning rate by batch size/256 and use a global weight decay parameter of $1.5 * 10^{-6}$ and exclude the biases and batch normalization parameters. For the target network, we follow the approach of BYOL (Grill

et al., 2020) and start the exponential moving average parameter $\tau$ at $\tau_{base} = 0.996$ and increase it to one during training via $\tau = 1 - (1 - \tau_{base})(\cos(\pi k/K) + 1)/2$ with k the current training step and K the maximum number of training steps.

### E.4 Evaluation on ImageNet

We follow the standard linear evaluation protocol on ImageNet as in (Kolesnikov et al., 2019; Chen et al., 2020a; Grill et al., 2020). We train a linear classifier on top of the fixed representation, i.e. we do not update the network parameters or the batch statistics. For training, we randomly crop and resize images to $224 \times 224$, and randomly horizontally flip the images after that. For testing, the images are resized to 256 pixels along the shorter dimension with bicubic resampling after which we take a center crop of size $224 \times 224$. Both for training and testing, the images are normalized by substracting the mean and standard deviations across the color channels computed on ImageNet after the augmentations. We use Stochastic Gradient Descent with a Nestorov momentum of $0.9$ and train for 80 epochs with a batch size of $1024$. We do not use any regularization techniques, e.g. weight decay.

### E.5 Robustness and Generalization

#### E.5.1 Dataset Details

**ImageNet-C.** The ImageNet-C dataset (Hendrycks & Dietterich, 2019) consists of 15 different types of corruptions from the noise, blur, weather, and digital categories applied to the validation images of ImageNet. This dataset is used for measuring semantic robustness. Figure 5 visualizes the corruption types. Each type of corruption has 5 levels of severity, i.e. there are 75 distinct corruptions in the dataset. In Figure 6, we display the Impulse noise corruption for 5 different severity levels. As can be seen, with increasing severity level the image becomes increasingly corrupted and difficult to parse. In addition to these 75 corruption types, there are an additional 4 corruption types (speckle noise, gaussian blur, spatter and saturate) that are provided as a validation set. We use these additional corruption types for selecting the best hyperparameters. For further details on this dataset, please refer to (Hendrycks & Dietterich, 2019).

**ImageNet-R.** The ImageNet-R dataset (Hendrycks et al., 2020) consists of $30,000$ images depicting various artistic renditions (e.g., paintings, sculpture, origami, cartoon) of 200 ImageNet object classes. This dataset is used to measure out-of-distribution generalization to various abstract visual renditions as it emphasizes shape over texture. The data was collected primarily from Flickr and also includes line drawings from (Wang et al., 2019). The images represent naturally occurring objects and have different textures and local image statistic to those of ImageNet. Figure 7 visualizes different images from the dataset. For further details on this dataset, please refer to (Hendrycks et al., 2020).

#### E.5.2 Evaluation

To evaluate robustness and generalization of the learned representation, we follow the standard linear evaluation protocol on ImageNet as in (Chen et al., 2020b;a; Kolesnikov et al., 2019). We train a linear classifier on top of the frozen representation, i.e. we do not update either the network parameters nor the batch statistics. During training, we augment the data by randomly cropping, resizing to $224 \times 224$ and randomly flipping the image. At test time, images are resized to 256 pixels along the shorter side via bicubic resampling and we take a $224 \times 224$ center crop. Both during training and testing, after applying augmentations we normalize the color channels by subtracting the average color and dividing by the standard deviation that is computed on ImageNet. We optimize the cross-entropy loss using Stochastic Gradient Descent with Nestorov momentum of $0.9$. We sweep over number for epochs $\{30, 50, 60, 90\}$, learning rates $\{0.4, 0.3, 0.2, 0.1, 0.05, 0.01\}$ and batch sizes $\{1024, 2048, 4096\}$. We select hyperparameters on the validation set provided in ImageNet-C and report the performance on ImageNet-R and on the test set of ImageNet-C under the best validation hyperparameters. We do not use any regularization techniques such as weight decay, gradient clipping, $tanh$ clipping or logits regularization.

Figure 5: The ImageNet-C dataset consists of 15 types of corruptions from noise, blur, weather, and digital categories. Each type of corruption has five levels of severity, resulting in 75 distinct corruptions. See different severity levels in Figure 6.



Figure 6: The 5 different levels of severity of Impulse noise corruption available in the ImageNet-C dataset. With increasing severity the dog image is markedly corrupted.

### E.5.3 Robustness metrics and further results

Let $f$ be a classifier that has not been trained on ImageNet-C. For each corruption type $c$ and level of severity $1 \leq s \leq 5$, denote the top-1 error of this classifier as $E_{s,c}^f$. Different corruption types pose different levels of difficulty. To make error rates across corruption types more comparable, the



Figure 7: Example images from the dataset ImageNet-R which contains $30,000$ images of $200$ ImageNet classes. This dataset emphasizes shape over texture and has different textures and local image statistic to those of ImageNet.

error rates are divided by AlexNet's errors. This standardized measure is the Corruption Error and is computed as

$$CE_c^f = \left( \sum_{s=1}^{5} E_{s,c}^f \right) / \left( \sum_{s=1}^{5} E_{s,c}^{AlexNet} \right)$$

The average error across all 15 corruption types is called the mean Corruption Error (mCE). Corruption Errors and mCE measure absolute robustness.

To better assess robustness, we also report the relative Corruption Error which measures relative robustness, i.e. loss in performance under corruptions. Denote by $E_{\text{clean}}^f$ the top-1 error rate for $f$ on the clean test set of ImageNet. The relative Corruption Error is given as

$$rCE_c^f = \sum_{s=1}^{5} \left( E_{s,c}^f - E_{clean}^f \right) / \sum_{s=1}^{5} \left( E_{s,c}^{AlexNet} - E_{clean}^{AlexNet} \right)$$

The mean relative Corruption Error (mrCE) is the mean of the relative Corruption Errors across all the corruption types. For more details and intuitions about there measures please refer to (Hendrycks & Dietterich, 2019).

In Table 6, we report Corruption Errors for Blur, Weather, and Digital corruption types. In Table 7, we report the relative robustness. As per (Hendrycks & Dietterich, 2019), we used the following values as the average AlexNet errors across severities, i.e. $\frac{1}{5} \sum_{s=1}^{5} E_{s,c}^{AlexNet}$, to normalize the Corruption Error values – Gaussian Noise 88.6%, Shot Noise 89.4%, Impulse Noise 92.3%, Defocus Blur 82.0%, Glass Blur 82.6%, Motion Blur 78.6%, Zoom Blur 79.8%, Snow 86.7%, Frost 82.7%, Fog 81.9%, Brightness 56.5%, Contrast 85.3%, Elastic Transformation 64.6%, Pixelate 71.8%, JPEG 60.7%, Speckle Noise 84.5%, Gaussian Blur 78.7%, Spatter 71.8%, Saturate 65.8%.

Table 6: Mean Corruption Error (mCE) and Corruption Error values for Blur, Weather, and Digital corruption types on ImageNet-C. All models are trained only using clean ImageNet images.

| Method | mCE | Blur | | | | Weather | | | | Digital | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Defocus | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
| Supervised | 76.7 | 75 | **89** | **78** | **80** | 78 | 75 | 66 | 57 | 71 | **85** | 77 | 77 |
| *Using ResNet-50:* | | | | | | | | | | | | | |
| SimCLR | 87.5 | 94.8 | 103.3 | 101.8 | 101.9 | 83.7 | 80.6 | 65.6 | 71.5 | 54 | 106.8 | 105.2 | 93 |
| ReLIC (ours) | 76.4 | 81.4 | 96.9 | 92.7 | 93.2 | 73.7 | 71.2 | 54.5 | 60.2 | **46.9** | 97.4 | 85.5 | 77.2 |
| *ResNet-50 with target network:* | | | | | | | | | | | | | |
| BYOL | 72.3 | 75 | 93.6 | 86.3 | 87.9 | 74.3 | 69.1 | 48.5 | 55 | 48.6 | 90.4 | **74.3** | 73 |
| ReLIC (ours) | **70.8** | **73.2** | 94 | 81.9 | 87 | **73.2** | **68** | **47.5** | **54.2** | 48.4 | 89.5 | 75.6 | **71.8** |

Table 7: Mean relative Corruption Error (mrCE) and relative Corruption Error values for different corruptions and methods on ImageNet-C. The mrCE value is the mean relative Corruption Error of the corruptions in Noise, Blur, Weather, and Digital columns. All models are trained only using clean ImageNet images. RELIC-t denotes using RELIC with a ResNet-50+target network architecture as in BYOL (Grill et al., 2020).

| Method | mrCE | Noise | | | Blur | | | | Weather | | | | Digital | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Gauss. | Shot | Impulse | Defocus | Glass | Motion | Zoom | Snow | Frost | Fog | Bright | Contrast | Elastic | Pixel | JPEG |
| Supervised | 105 | 104 | 107 | 107 | 97 | **126** | **107** | **110** | 101 | 97 | 79 | 62 | 89 | **146** | 111 | 132 |
| *Using ResNet-50:* | | | | | | | | | | | | | | | | |
| SimCLR | 111.9 | 88 | 92.7 | 106.6 | 122.2 | 139.6 | 140.5 | 139.4 | 96.9 | 91.6 | 59.9 | 74.8 | 36.7 | 181.5 | 158.3 | 149.8 |
| ReLIC (ours) | 87.7 | 67.3 | 73.1 | 84.8 | 96.1 | 128.7 | 123 | 123.1 | 79.3 | 74.4 | 38.9 | 33.4 | 24.6 | 157.5 | 112 | 99.8 |
| *ResNet-50 with target network:* | | | | | | | | | | | | | | | | |
| BYOL | 90 | 72.5 | 77.2 | 86.7 | 93 | 132 | 120 | 122.5 | 89.7 | 80.2 | 36.6 | 41.5 | 37.8 | 155 | **97.6** | 108.4 |
| ReLIC (ours) | 88.4 | 69.1 | **73** | **79.2** | **90.4** | 134.2 | 111.6 | 121.9 | 88.5 | 79.2 | **35.6** | 41.7 | 38.5 | 154.5 | 102.7 | 106.8 |

## E.6  EVALUATION ON ATARI

For our experiments on Atari, we use the agent from R2D2 (Kapturowski et al., 2019) with standard hyperparameters noted below. We train each agent on approximately 15 billion frames and add a second encoder with the same architecture used in the Q-Network of the original agent. This second encoder is trained with a separate optimizer with only a representation learning objective. The agent then takes the output of this encoder as a given input. We use standard augmentations used in prior work (Kostrikov et al., 2020) where we pad the frames on all sides with 4 pixels copied from the

borders and then randomly cropping 84 windows. We randomly shift pixel intensity according to the distribution $s = 1.0 + 0.1 * \mathcal{N}'$ where $\mathcal{N}'$ is the standard Normal distribution with values clipped between -2 and 2. $s$ is then multiplied by the original image to return the augmented image.

**RELIC and SimCLR** For our implementation of RELIC and SimCLR, we do not use a critic embedding at all and utilize the last layer of the encoder for the objective. As in CURL (Srinivas et al., 2020) we utilize a target encoder for the second augmentation where we update the weights with a momentum of .99. We also clipped the gradients of our optimizer using a global norm ratio of 40. We report the hyperparameters in Table 8.

Table 8: RELIC and SimCLR Details

| Parameter | Value |
| --- | --- |
| Normalize Inputs | True |
| Temperature | 1.0 Constant |
| Scaling of Embeddings | False |
| Optimizer | Adam |
| Learning Rate | 5e-4 |
| Epsilon | 0.01 |
| Beta 1 | 0.9 |
| Beta 2 | 0.999 |

**CURL** For CURL, we use a second encoder as noted before. With the exception of the encoder architecture and the optimizer parameters, all hyperparameters are the same as in (Srinivas et al., 2020) including the momentum value for the target network weight updates. We utilize the same architecture in the paper with a linear layer as a critic embedding for the target encoder.

Table 9: CURL Details

| Parameter | Value |
| --- | --- |
| Optimizer | Adam |
| Learning Rate | 1e-3 |
| Epsilon | 0.01 |
| Beta 1 | 0.9 |
| Beta 2 | 0.999 |

**BYOL** In BYOL, we utilize two-layer perceptron networks as our predictor and projection layers. For both networks, the number of hidden units in the two layers was 1024 and 512. We use a target network update momentum of .99. The optimizer parameters are the same as in Table 8.

**Direct Augmentation** We also compared against direct augmentation of the observations in the replay buffer as in DrQ (Kostrikov et al., 2020). We keep the architecture the same in this instance and use two duplicate encoders as input to the agent. In this case, the optimizer can jointly update both encoders and train them end-to-end.

Table 10: Individual Mean Episode Return on Atari.

| Games | Average Human | Random | RELIC (ours) | SimCLR | CURL | BYOL | Augmentation |
|---|---|---|---|---|---|---|---|
| alien | 7127.70 | 227.80 | 8766.57 | **10082.54** | 8506.48 | 9671.89 | 5201.93 |
| amidar | 1719.50 | 5.80 | **28449.26** | 28141.18 | 27213.75 | 25965.05 | 867.66 |
| assault | 742.00 | 222.40 | **92963.07** | 36109.84 | 7139.67 | 13565.20 | 1539.71 |
| asterix | 8503.30 | 210.00 | **998426.72** | 997305.51 | 661431.39 | 986307.92 | 26239.64 |
| asteroids | 47388.70 | 719.10 | 83669.38 | 7299.90 | 76612.17 | 55936.02 | **101340.17** |
| atlantis | 29028.10 | 12850.00 | 1575940.94 | 1584392.76 | **1584698.01** | 1530122.45 | 794011.79 |
| bank heist | 753.10 | 14.20 | 1521.38 | 2467.62 | **4095.29** | 1659.94 | 771.60 |
| battle zone | 37187.50 | 2360.00 | **452831.48** | 278903.14 | 287792.06 | 338695.47 | 31511.75 |
| beam rider | 16926.50 | 363.90 | **136695.24** | 98551.42 | 116794.58 | 87454.20 | 46894.14 |
| berzerk | 2630.40 | 123.70 | **146213.60** | 1301.36 | 73754.38 | 1265.21 | 73645.52 |
| bowling | 160.70 | 23.10 | 205.09 | 193.50 | **230.31** | 172.21 | 164.68 |
| boxing | 12.10 | 0.10 | 100.00 | 100.00 | 100.00 | 100.00 | **100.00** |
| breakout | 30.50 | 1.70 | 405.05 | 404.06 | 407.14 | **409.48** | 150.67 |
| centipede | 12017.00 | 2090.90 | **220886.86** | 99544.92 | 167779.11 | 146735.67 | 20152.01 |
| chopper command | 7387.80 | 811.00 | 999900.00 | 999900.00 | **999900.00** | 962003.61 | 5399.56 |
| crazy climber | 35829.40 | 10780.50 | 272179.68 | 266870.81 | **301689.62** | 210477.39 | 96538.00 |
| defender | 18688.90 | 2874.50 | **576405.57** | 522617.05 | 560816.84 | 493410.36 | 78750.19 |
| demon attack | 1971.00 | 152.10 | 143774.79 | **143786.19** | 143737.36 | 143574.86 | 821.98 |
| double dunk | -16.40 | -18.60 | **24.00** | 24.00 | 24.00 | 24.00 | 14.82 |
| enduro | 860.50 | 0.00 | 2371.27 | 2366.19 | **2373.12** | 2368.00 | 1361.66 |
| fishing derby | -38.70 | -91.70 | 68.17 | **83.00** | 72.21 | 70.11 | 19.93 |
| freeway | 29.60 | 0.00 | 33.00 | 32.93 | **33.04** | 33.00 | 32.00 |
| frostbite | 4334.70 | 65.20 | 10156.41 | **11171.49** | 3693.20 | 5793.80 | 5708.35 |
| gopher | 2412.50 | 257.60 | **123170.74** | 122368.21 | 122371.64 | 120317.04 | 43711.82 |
| gravitar | 3351.40 | 173.00 | 4186.09 | 3601.14 | **4997.87** | 4048.25 | 2014.59 |
| hero | **30826.40** | 1027.00 | 13615.35 | 13523.98 | 13620.78 | 13558.04 | 8957.00 |
| ice hockey | 0.90 | -11.20 | 56.39 | 48.27 | 45.06 | **59.70** | -2.43 |
| jamesbond | 302.80 | 29.00 | **15632.87** | 5714.62 | 10052.04 | 10099.81 | 1441.95 |
| kangaroo | 3035.00 | 52.00 | 14342.59 | 14215.11 | 11674.19 | **14471.65** | 7249.73 |
| krull | 2665.50 | 1598.00 | **137099.65** | 100426.69 | 86049.99 | 80414.04 | 16626.09 |
| kung fu master | 22736.30 | 258.50 | **230241.57** | 220076.57 | 228943.94 | 208064.38 | 64632.42 |
| montezuma revenge | **4753.30** | 0.00 | 1066.67 | 733.33 | 1072.30 | 419.54 | 26.67 |
| ms pacman | 6951.60 | 307.30 | 13367.55 | 12053.76 | **13465.80** | 12726.79 | 3238.90 |
| name this game | 8049.00 | 2292.30 | **48669.30** | 46657.55 | 47417.82 | 44848.29 | 13416.57 |
| phoenix | 7242.60 | 761.40 | 803108.37 | 253542.40 | 580969.56 | 20317.80 | 6264.39 |
| pitfall | **6463.70** | -229.40 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| pong | 14.60 | -20.70 | **21.00** | 21.00 | 21.00 | 21.00 | 21.00 |
| private eye | **69571.30** | 24.90 | 10154.93 | 5115.34 | 5190.28 | 470.68 | 111.77 |
| qbert | 13455.00 | 163.90 | **353197.13** | 24340.75 | 208207.97 | 57261.24 | 11051.97 |
| riverraid | 17118.00 | 1338.50 | **23525.44** | 20400.83 | 20230.02 | 22206.57 | 10487.59 |
| road runner | 7845.00 | 11.50 | 213173.15 | 236235.30 | 241917.98 | 238880.54 | **440430.17** |
| robotank | 11.90 | 2.20 | 97.65 | 82.60 | **98.13** | 62.54 | 49.98 |
| seaquest | 42054.70 | 68.40 | **999999.00** | 999999.00 | 666700.67 | 29160.93 | 37397.26 |
| skiing | **-4336.90** | -17098.10 | -24761.06 | -23076.73 | -15497.66 | -26028.08 | -22162.91 |
| solaris | **12326.70** | 1236.30 | 4594.37 | 4571.27 | 4276.39 | 4331.03 | 4142.69 |
| space invaders | 1668.70 | 148.00 | **3625.52** | 3619.94 | 3542.48 | 3613.93 | 835.37 |
| star gunner | 10250.00 | 664.00 | 283499.72 | **289099.89** | 129720.84 | 175486.67 | 43167.07 |
| surround | 6.50 | -10.00 | **10.00** | 9.96 | 1.60 | 9.56 | -0.64 |
| tennis | -8.30 | -23.80 | 0.00 | 0.00 | 0.00 | 0.00 | **0.12** |
| time pilot | 5229.20 | 3568.00 | 309297.74 | 92888.66 | **400326.69** | 48011.44 | 14198.37 |
| tutankham | 167.60 | 11.40 | **371.17** | 306.45 | 337.61 | 285.36 | 144.30 |
| up n down | 11693.20 | 533.40 | **577256.03** | 520666.59 | 566912.89 | 552110.67 | 143512.38 |
| venture | 1187.50 | 0.00 | 1929.53 | **1945.20** | 1906.84 | 1881.76 | 733.29 |
| video pinball | 17667.90 | 0.00 | 978292.52 | **993332.08** | 932523.58 | 623223.24 | 37584.71 |
| wizard of wor | 4756.50 | 563.50 | **123513.74** | 89462.62 | 106801.20 | 68256.44 | 5940.82 |
| yars revenge | 54576.90 | 3092.90 | 228704.52 | 99636.25 | **229221.52** | 86847.75 | 48041.63 |
| zaxxon | 9173.30 | 32.50 | **120830.77** | 57379.66 | 85906.74 | 48067.61 | 23688.22 |