

Bootstrap Your Own Latent A New Approach to Self-Supervised Learning

Jean-Bastien Grill^{*1} Florian Strub^{*1} Florent Altché^{*1} Corentin Tallec^{*1} Pierre H. Richemond^{*1,2}

Elena Buchatskaya¹ Carl Doersch¹ Bernardo Avila Pires¹ Zhaohan Daniel Guo¹

Mohammad Gheshlaghi Azar¹ Bilal Piot¹ Koray Kavukcuoglu¹ Rémi Munos¹ Michal Valko¹

¹DeepMind ²Imperial College

[jbgrill,fstrub,altche,corentint,richemond]@google.com

Abstract

We introduce **Bootstrap Your Own Latent (BYOL)**, a new approach to self-supervised image representation learning. BYOL relies on two neural networks, referred to as *online* and *target* networks, that interact and learn from each other. From an augmented view of an image, we train the online network to predict the target network representation of the same image under a different augmented view. At the same time, we update the target network with a slow-moving average of the online network. While state-of-the-art methods rely on negative pairs, BYOL achieves a new state of the art *without them*. BYOL reaches 74.3% top-1 classification accuracy on ImageNet using a linear evaluation with a ResNet-50 architecture and 79.6% with a larger ResNet. We show that BYOL performs on par or better than the current state of the art on both transfer and semi-supervised benchmarks. Our implementation and pretrained models are given on GitHub.³

1 Introduction

Learning good image representations is a key challenge in computer vision [1, 2, 3] as it allows for efficient training on downstream tasks [4, 5, 6, 7]. Many different training approaches have been proposed to learn such representations, usually relying on visual pretext tasks. Among them, state-of-the-art contrastive methods [8, 9, 10, 11, 12] are trained by reducing the distance between representations of different augmented views of the same image ('positive pairs'), and increasing the distance between representations of augmented views from different images ('negative pairs'). These methods need careful treatment of negative pairs [13] by either relying on large batch sizes [8, 12], memory banks [9] or customized mining strategies [14, 15] to retrieve the negative pairs. In addition, their performance critically depends on the choice of image augmentations [8, 12].

In this paper, we introduce **Bootstrap Your Own Latent (BYOL)**, a new algorithm for self-supervised learning of image representations. BYOL achieves higher performance than state-of-the-art contrastive methods

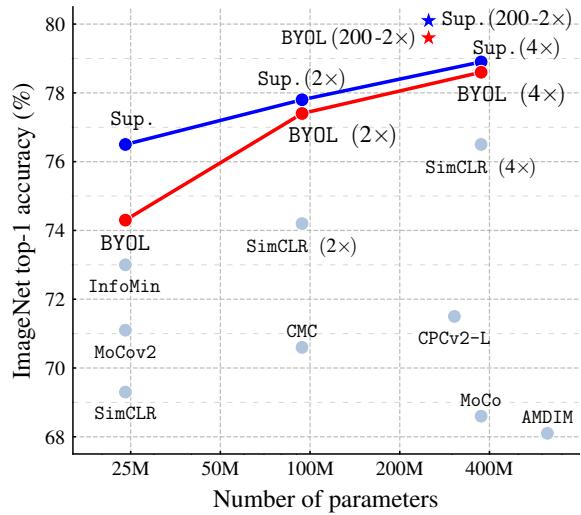


Figure 1: Performance of BYOL on ImageNet (linear evaluation) using ResNet-50 and our best architecture ResNet-200 (2x), compared to other unsupervised and supervised (Sup.) baselines [8].

^{*}Equal contribution; the order of first authors was randomly selected.

³<https://github.com/deepmind/deepmind-research/tree/master/byol>

representation collapse - like trivial solutions like constant representation for both the similar images when no negative were used in contrastive loss, only positive pairs are used to learn the representations

"bootstrapping" can be used metaphorically to describe a process where someone or something achieves success or self-improvement without relying on external help. The metaphor comes from the idea of "pulling oneself up by one's bootstraps," implying a self-sufficient and self-directed effort

without using negative pairs. It iteratively bootstraps⁴ the outputs of a network to serve as targets for an enhanced representation. Moreover, BYOL is more robust to the choice of image augmentations than contrastive methods; we suspect that not relying on negative pairs is one of the leading reasons for its improved robustness. While previous methods based on bootstrapping have used pseudo-labels [16], cluster indices [17] or a handful of labels [18, 19, 20], we propose to directly bootstrap the representations. In particular, BYOL uses two neural networks, referred to as online and target networks, that interact and learn from each other. Starting from an augmented view of an image, BYOL trains its online network to predict the target network's representation of another augmented view of the same image. While this objective admits collapsed solutions, e.g., outputting the same vector for all images, we empirically show that BYOL does not converge to such solutions. We hypothesize (see Section 3.2) that the combination of (i) the addition of a predictor to the online network and (ii) the use of a slow-moving average of the online parameters as the target network encourages encoding more and more information within the online projection and avoids collapsed solutions.

We evaluate the representation learned by BYOL on ImageNet [21] and other vision benchmarks using ResNet architectures [22]. Under the linear evaluation protocol on ImageNet, consisting in training a linear classifier on top of the frozen representation, BYOL reaches 74.3% top-1 accuracy with a standard ResNet-50 and 79.6% top-1 accuracy with a larger ResNet (Figure 1). In the semi-supervised and transfer settings on ImageNet, we obtain results on par or superior to the current state of the art. Our contributions are: (i) We introduce BYOL, a self-supervised representation learning method (Section 3) which achieves state-of-the-art results under the linear evaluation protocol on ImageNet without using negative pairs. (ii) We show that our learned representation outperforms the state of the art on semi-supervised and transfer benchmarks (Section 4). (iii) We show that BYOL is more resilient to changes in the batch size and in the set of image augmentations compared to its contrastive counterparts (Section 5). In particular, BYOL suffers a much smaller performance drop than SimCLR, a strong contrastive baseline, when only using random crops as image augmentations.

2 Related work

Most unsupervised methods for representation learning can be categorized as either generative or discriminative [23, 8]. Generative approaches to representation learning build a distribution over data and latent embedding and use the learned embeddings as image representations. Many of these approaches rely either on auto-encoding of images [24, 25, 26] or on adversarial learning [27], jointly modelling data and representation [28, 29, 30, 31]. Generative methods typically operate directly in pixel space. This however is computationally expensive, and the high level of detail required for image generation may not be necessary for representation learning.

Among discriminative methods, contrastive methods [9, 10, 32, 33, 34, 11, 35, 36] currently achieve state-of-the-art performance in self-supervised learning [37, 8, 38, 12]. Contrastive approaches avoid a costly generation step in pixel space by bringing representation of different views of the same image closer ('positive pairs'), and spreading representations of views from different images ('negative pairs') apart [39, 40]. Contrastive methods often require comparing each example with many other examples to work well [9, 8] prompting the question of whether using negative pairs is necessary.

DeepCluster [17] partially answers this question. It uses bootstrapping on previous versions of its representation to produce targets for the next representation; it clusters data points using the prior representation, and uses the cluster index of each sample as a classification target for the new representation. While avoiding the use of negative pairs, this requires a costly clustering phase and specific precautions to avoid collapsing to trivial solutions.

Some self-supervised methods are not contrastive but rely on using auxiliary handcrafted prediction tasks to learn their representation. In particular, relative patch prediction [23, 40], colorizing gray-scale images [41, 42], image inpainting [43], image jigsaw puzzle [44], image super-resolution [45], and geometric transformations [46, 47] have been shown to be useful. Yet, even with suitable architectures [48], these methods are being outperformed by contrastive methods [37, 8, 12].

Our approach has some similarities with *Predictions of Bootstrapped Latents* (PBL, [49]), a self-supervised representation learning technique for reinforcement learning (RL). PBL jointly trains the agent's history representation and an encoding of future observations. The observation encoding is used as a target to train the agent's representation, and the agent's representation as a target to train the observation encoding. Unlike PBL, BYOL uses a slow-moving average of its representation to provide its targets, and does not require a second network.

The idea of using a slow-moving average target network to produce stable targets for the online network was inspired by deep RL [50, 51, 52, 53]. Target networks stabilize the bootstrapping updates provided by the Bellman equation,

⁴Throughout this paper, the term *bootstrap* is used in its idiomatic sense rather than the statistical sense.

bootstrapping in RL means that you update a value based on some estimates and not on some exact values
<https://datascience.stackexchange.com/questions/26938/what-exactly-is-bootstrapping-in-reinforcement-learning>

making them appealing to stabilize the bootstrap mechanism in BYOL. While most RL methods use fixed target networks, BYOL uses a weighted moving average of previous networks (as in [54]) in order to provide smoother changes in the target representation.

In the semi-supervised setting [55, 56], an unsupervised loss is combined with a classification loss over a handful of labels to ground the training [19, 20, 57, 58, 59, 60, 61, 62]. Among these methods, *mean teacher* (MT) [20] also uses a slow-moving average network, called *teacher*, to produce targets for an online network, called *student*. An ℓ_2 consistency loss between the softmax predictions of the teacher and the student is added to the classification loss. While [20] demonstrates the effectiveness of MT in the semi-supervised learning case, in Section 5 we show that a similar approach collapses when removing the classification loss. In contrast, BYOL introduces an additional predictor on top of the online network, which prevents collapse.

Finally, in self-supervised learning, MoCo [9] uses a slow-moving average network (*momentum encoder*) to maintain consistent representations of negative pairs drawn from a memory bank. Instead, BYOL uses a moving average network to produce prediction targets as a means of stabilizing the bootstrap step. We show in Section 5 that this mere stabilizing effect can also improve existing contrastive methods.

3 Method

We start by motivating our method before explaining its details in Section 3.1. Many successful self-supervised learning approaches build upon the cross-view prediction framework introduced in [63]. Typically, these approaches learn representations by predicting different views (e.g., different random crops) of the same image from one another. Many such approaches cast the prediction problem directly in representation space: the representation of an augmented view of an image should be predictive of the representation of another augmented view of the same image. However, predicting directly in representation space can lead to collapsed representations: for instance, a representation that is constant across views is always fully predictive of itself. Contrastive methods circumvent this problem by reformulating the prediction problem into one of discrimination: from the representation of an augmented view, they learn to discriminate between the representation of another augmented view of the same image, and the representations of augmented views of different images. In the vast majority of cases, this prevents the training from finding collapsed representations. Yet, this discriminative approach typically requires comparing each representation of an augmented view with many negative examples, to find ones sufficiently close to make the discrimination task challenging. In this work, we thus tasked ourselves to find out whether these negative examples are indispensable to prevent collapsing while preserving high performance.

To prevent collapse, a straightforward solution is to use a fixed randomly initialized network to produce the targets for our predictions. While avoiding collapse, it empirically does not result in very good representations. Nonetheless, it is interesting to note that the representation obtained using this procedure can already be much better than the initial fixed representation. In our ablation study (Section 5), we apply this procedure by predicting a fixed randomly initialized network and achieve 18.8% top-1 accuracy (Table 5a) on the linear evaluation protocol on ImageNet, whereas the randomly initialized network only achieves 1.4% by itself. This experimental finding is the core motivation for BYOL: from a given representation, referred to as *target*, we can train a new, potentially enhanced representation, referred to as *online*, by predicting the target representation. From there, we can expect to build a sequence of representations of increasing quality by iterating this procedure, using subsequent online networks as new target networks for further training. In practice, BYOL generalizes this bootstrapping procedure by iteratively refining its representation, but using a slowly moving exponential average of the online network as the target network instead of fixed checkpoints.

3.1 Description of BYOL

BYOL’s goal is to learn a representation y_θ which can then be used for downstream tasks. As described previously, BYOL uses two neural networks to learn: the *online* and *target* networks. The *online* network is defined by a set of weights θ and is comprised of three stages: an *encoder* f_θ , a *projector* g_θ and a *predictor* q_θ , as shown in Figure 2 and Figure 8. The *target* network has the same architecture as the *online* network, but uses a different set of weights ξ . The *target* network provides the regression targets to train the *online* network, and its parameters ξ are an exponential moving average of the *online* parameters θ [54]. More precisely, given a target decay rate $\tau \in [0, 1]$, after each training step we perform the following update,

$$\xi \leftarrow \tau \xi + (1 - \tau) \theta. \quad (1)$$

Given a set of images \mathcal{D} , an image $x \sim \mathcal{D}$ sampled uniformly from \mathcal{D} , and two distributions of image augmentations \mathcal{T} and \mathcal{T}' , BYOL produces two augmented views $v \triangleq t(x)$ and $v' \triangleq t'(x)$ from x by applying respectively image

asymmetric in architecture(predictor in inline), asymmetric learning rule(different weights and moving avg of online weights for target net, gradients weight update for online), asymmetry in optimization(stop gradient in target)

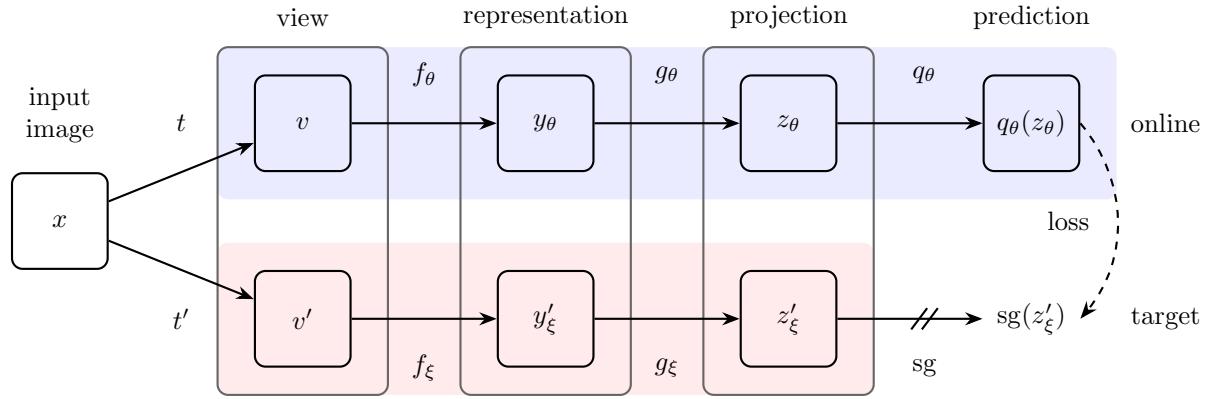


Figure 2: BYOL’s architecture. BYOL minimizes a similarity loss between $q_\theta(z_\theta)$ and $sg(z'_\xi)$, where θ are the trained weights, ξ are an exponential moving average of θ and sg means stop-gradient. At the end of training, everything but f_θ is discarded, and y_θ is used as the image representation.

augmentations $t \sim \mathcal{T}$ and $t' \sim \mathcal{T}'$. From the first augmented view v , the online network outputs a *representation* $y_\theta \triangleq f_\theta(v)$ and a *projection* $z_\theta \triangleq g_\theta(y)$. The target network outputs $y'_\xi \triangleq f_\xi(v')$ and the *target projection* $z'_\xi \triangleq g_\xi(y')$ from the second augmented view v' . We then output a *prediction* $q_\theta(z_\theta)$ of z'_ξ and ℓ_2 -normalize both $q_\theta(z_\theta)$ and z'_ξ to $\bar{q}_\theta(z_\theta) \triangleq q_\theta(z_\theta)/\|q_\theta(z_\theta)\|_2$ and $\bar{z}'_\xi \triangleq z'_\xi/\|z'_\xi\|_2$. Note that this predictor is only applied to the *online branch*, making the *architecture asymmetric* between the *online* and *target* pipeline. Finally we define the following mean squared error between the normalized predictions and target projections,⁵

Minimizing the MSE(the first eq) ==
maximizing the cosine similarity(second eq).
In secnd eq, we are using unnormalized
vectors in the numerator, as they are
normalizing anyway in the cosine sim with l2
norm in denominator

$$\mathcal{L}_{\theta,\xi} \triangleq \|\bar{q}_\theta(z_\theta) - \bar{z}'_\xi\|_2^2 = 2 - 2 \cdot \frac{\langle q_\theta(z_\theta), z'_\xi \rangle}{\|q_\theta(z_\theta)\|_2 \cdot \|z'_\xi\|_2}. \quad (2)$$

<https://stats.stackexchange.com/questions/146221/is-cosine-similarity-identical-to-l2-normalized-euclidean-distance>

We symmetrize the loss $\mathcal{L}_{\theta,\xi}$ in Eq. 2 by separately feeding v' to the online network and v to the target network to compute $\tilde{\mathcal{L}}_{\theta,\xi}$. At each training step, we perform a stochastic optimization step to minimize $\mathcal{L}_{\theta,\xi}^{\text{BYOL}} = \mathcal{L}_{\theta,\xi} + \tilde{\mathcal{L}}_{\theta,\xi}$ with respect to θ only, but *not* ξ , as depicted by the stop-gradient in Figure 2. BYOL’s dynamics are summarized as

$$\begin{aligned} \theta &\leftarrow \text{optimizer}(\theta, \nabla_\theta \mathcal{L}_{\theta,\xi}^{\text{BYOL}}, \eta), \\ \xi &\leftarrow \tau \xi + (1 - \tau) \theta, \end{aligned}$$

Stochastic optimization is an optimization approach that involves using randomness or probability in the optimization process. In traditional optimization, the objective is to find the minimum or maximum of a deterministic function. In contrast, stochastic optimization deals with optimizing objectives that involve randomness or uncertainty.

where *optimizer* is an optimizer and η is a learning rate.

At the end of training, we only keep the encoder f_θ ; as in [9]. When comparing to other methods, we consider the number of inference-time weights only in the final representation f_θ . The full training procedure is summarized in Appendix A, and python pseudo-code based on the libraries JAX [64] and Haiku [65] is provided in in Appendix J.

3.2 Intuitions on BYOL’s behavior

As BYOL does not use an explicit term to prevent collapse (such as negative examples [10]) while minimizing $\mathcal{L}_{\theta,\xi}^{\text{BYOL}}$ with respect to θ , it may seem that BYOL should converge to a minimum of this loss with respect to (θ, ξ) (e.g., a collapsed constant representation). However BYOL’s target parameters ξ updates are **not** in the direction of $\nabla_\xi \mathcal{L}_{\theta,\xi}^{\text{BYOL}}$. More generally, we hypothesize that there is no loss $L_{\theta,\xi}$ such that BYOL’s dynamics is a gradient descent on L jointly over θ, ξ . This is similar to GANs [66], where there is no loss that is jointly minimized w.r.t. both the discriminator and generator parameters. There is therefore no *a priori* reason why BYOL’s parameters would converge to a minimum of $\mathcal{L}_{\theta,\xi}^{\text{BYOL}}$.

While BYOL’s dynamics still admit undesirable equilibria, we did not observe convergence to such equilibria in our experiments. In addition, when assuming BYOL’s predictor to be optimal⁶ i.e., $q_\theta = q^*$ with

$$q^* \triangleq \arg \min_q \mathbb{E} [\|q(z_\theta) - z'_\xi\|_2^2], \quad \text{where } q^*(z_\theta) = \mathbb{E}[z'_\xi | z_\theta], \quad (4)$$

⁵While we could directly predict the representation y and not a projection z , previous work [8] have empirically shown that using this projection improves performance.

⁶For simplicity we also consider BYOL without normalization (which performs reasonably close to BYOL, see Appendix F.6) nor symmetrization

Similarly to GANs, BYOL uses two sets of parameters that are not minimizing the same objective. Thus, there is no *a priori* reason for BYOL’s dynamics to converge to a global minimum of $\|q^*(z_\theta) - z'_\xi\|_2^2$, as they are not following the gradient of this loss (BYOL uses z'). While these dynamics still admit undesirable equilibria where all images are mapped to the same constant projection (e.g., all zeros), BYOL’s empirical performance seems to indicate that such equilibria may be unstable. We hypothesize that maintaining a near optimal predictor at all times is key to avoid collapsed solutions. When using an optimal predictor, BYOL minimizes the(expected) conditional variance of the target projection given the online projection. With a fixed target network, adding more information to the online projection can reduce this conditional variance, but cannot increase it. For example, training dynamics will always tend not to collapse features from the online network, as for any constant C and variables X and Y , $\text{Var}(Y|X) = \text{Var}(Y|C)$. More generally BYOL is encouraged to keep features from the online projection diverse by latching onto any source of variability Z (stemming, e.g., from noise in training dynamics) distinct from existing features, as $\text{Var}(Y|X, Z) > \text{Var}(Y|X, h(X))$ for any variables X, Y, Z and any function h . We will add these additional discussions and experiments to our submission to clarify the role of the predictor.

expectation of the gradient refers to the average gradient of a function with respect to its parameters, taken over some distribution. It is often calculated in the context of stochastic optimization, where the objective function involves random variables

constant representation

we hypothesize that the undesirable equilibria are unstable. Indeed, in this optimal predictor case, BYOL’s updates on θ follow in expectation, the gradient of the expected conditional variance (see Appendix H for details),

$$\nabla_{\theta} \mathbb{E} \left[\|q^*(z_{\theta}) - z'_{\xi}\|_2^2 \right] = \nabla_{\theta} \mathbb{E} \left[\|\mathbb{E}[z'_{\xi}|z_{\theta}] - z'_{\xi}\|_2^2 \right] = \nabla_{\theta} \mathbb{E} \left[\sum_i \text{Var}(z'_{\xi,i}|z_{\theta}) \right], \quad (5)$$

where $z'_{\xi,i}$ is the i -th feature of z'_{ξ} .

Note that for any random variables X , Y , and Z , $\text{Var}(X|Y, Z) \leq \text{Var}(X|Y)$. Let X be the target projection, Y the current online projection, and Z an additional variability on top of the online projection induced by stochasticities in the training dynamics: purely discarding information from the online projection cannot decrease the conditional variance.

In particular, BYOL avoids constant features in z_{θ} as, for any constant c and random variables z_{θ} and z'_{ξ} , $\text{Var}(z'_{\xi}|z_{\theta}) \leq \text{Var}(z'_{\xi}|c)$; hence our hypothesis on these collapsed constant equilibria being unstable. Interestingly, if we were to minimize $\mathbb{E}[\sum_i \text{Var}(z'_{\xi,i}|z_{\theta})]$ with respect to ξ , we would get a collapsed z'_{ξ} as the variance is minimized for a constant z'_{ξ} . Instead, BYOL makes ξ closer to θ , incorporating sources of variability captured by the online projection into the target projection.

Furthermore, notice that performing a hard-copy of the online parameters θ into the target parameters ξ would be enough to propagate new sources of variability. However, sudden changes in the target network might break the assumption of an optimal predictor, in which case BYOL’s loss is not guaranteed to be close to the conditional variance. We hypothesize that the main role of BYOL’s moving-averaged target network is to ensure the near-optimality of the predictor over training; Section 5 and Appendix I provide some empirical support of this interpretation.

3.3 Implementation details

Image augmentations BYOL uses the same set of image augmentations as in SimCLR [8]. First, a random patch of the image is selected and resized to 224×224 with a random horizontal flip, followed by a color distortion, consisting of a random sequence of brightness, contrast, saturation, hue adjustments, and an optional grayscale conversion. Finally Gaussian blur and solarization are applied to the patches. Additional details on the image augmentations are in Appendix B.

Architecture We use a convolutional residual network [22] with 50 layers and post-activation (ResNet-50(1×) v1) as our base parametric encoders f_{θ} and f_{ξ} . We also use deeper (50, 101, 152 and 200 layers) and wider (from 1× to 4×) ResNets, as in [67, 48, 8]. Specifically, the representation y corresponds to the output of the final average pooling layer, which has a feature dimension of 2048 (for a width multiplier of 1×). As in SimCLR [8], the representation y is projected to a smaller space by a multi-layer perceptron (MLP) g_{θ} , and similarly for the target projection g_{ξ} . This MLP consists in a linear layer with output size 4096 followed by batch normalization [68], rectified linear units (ReLU) [69], and a final linear layer with output dimension 256. Contrary to SimCLR, the output of this MLP is not batch normalized. The predictor q_{θ} uses the same architecture as g_{θ} .

Optimization We use the LARS optimizer [70] with a cosine decay learning rate schedule [71], without restarts, over 1000 epochs, with a warm-up period of 10 epochs. We set the base learning rate to 0.2, scaled linearly [72] with the batch size (LearningRate = $0.2 \times \text{BatchSize}/256$). In addition, we use a global weight decay parameter of $1.5 \cdot 10^{-6}$ while excluding the biases and batch normalization parameters from both LARS adaptation and weight decay. For the target network, the exponential moving average parameter τ starts from $\tau_{\text{base}} = 0.996$ and is increased to one during training. Specifically, we set $\tau \triangleq 1 - (1 - \tau_{\text{base}}) \cdot (\cos(\pi k/K) + 1)/2$ with k the current training step and K the maximum number of training steps. We use a batch size of 4096 split over 512 Cloud TPU v3 cores. With this setup, training takes approximately 8 hours for a ResNet-50(1×). All hyperparameters are summarized in Appendix J; an additional set of hyperparameters for a smaller batch size of 512 is provided in Appendix G.

4 Experimental evaluation

We assess the performance of BYOL’s representation after self-supervised pretraining on the training set of the ImageNet ILSVRC-2012 dataset [21]. We first evaluate it on ImageNet (IN) in both linear evaluation and semi-supervised setups. We then measure its transfer capabilities on other datasets and tasks, including classification, segmentation, object detection and depth estimation. For comparison, we also report scores for a representation trained using labels from the train ImageNet subset, referred to as Supervised-IN. In Appendix E, we assess the

generality of BYOL by pretraining a representation on the Places365-Standard dataset [73] before reproducing this evaluation protocol.

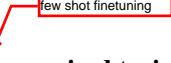
Linear evaluation on ImageNet We first evaluate BYOL’s representation by training a linear classifier on top of the frozen representation, following the procedure described in [48, 74, 41, 10, 8], and appendix C.1; we report top-1 and top-5 accuracies in % on the test set in Table 1. With a standard ResNet-50 ($\times 1$) BYOL obtains 74.3% top-1 accuracy (91.6% top-5 accuracy), which is a 1.3% (resp. 0.5%) improvement over the previous self-supervised state of the art [12]. This tightens the gap with respect to the supervised baseline of [8], 76.5%, but is still significantly below the stronger supervised baseline of [75], 78.9%. With deeper and wider architectures, BYOL consistently outperforms the previous state of the art (Appendix C.2), and obtains a best performance of 79.6% top-1 accuracy, ranking higher than previous self-supervised approaches. On a ResNet-50 ($\times 4$) BYOL achieves 78.6%, similar to the 78.9% of the best supervised baseline in [8] for the same architecture.

Method	Top-1	Top-5	Method	Architecture	Param.	Top-1	Top-5
Local Agg.	60.2	-	SimCLR [8]	ResNet-50 ($\times 2$)	94M	74.2	92.0
PIRL [35]	63.6	-	CMC [11]	ResNet-50 ($\times 2$)	94M	70.6	89.7
CPC v2 [32]	63.8	85.3	BYOL (ours)	ResNet-50 ($\times 2$)	94M	77.4	93.6
CMC [11]	66.2	87.0	CPC v2 [32]	ResNet-161	305M	71.5	90.1
SimCLR [8]	69.3	89.0	MoCo [9]	ResNet-50 ($\times 4$)	375M	68.6	-
MoCo v2 [37]	71.1	-	SimCLR [8]	ResNet-50 ($\times 4$)	375M	76.5	93.2
InfoMin Aug. [12]	73.0	91.1	BYOL (ours)	ResNet-50 ($\times 4$)	375M	78.6	94.2
BYOL (ours)	74.3	91.6	BYOL (ours)	ResNet-200 ($\times 2$)	250M	79.6	94.8

(a) ResNet-50 encoder.

(b) Other ResNet encoder architectures.

Table 1: Top-1 and top-5 accuracies (in %) under linear evaluation on ImageNet.


Semi-supervised training on ImageNet Next, we evaluate the performance obtained when fine-tuning BYOL’s representation on a classification task with a small subset of ImageNet’s train set, this time using label information. We follow the semi-supervised protocol of [74, 76, 8, 32] detailed in Appendix C.1, and use the same fixed splits of respectively 1% and 10% of ImageNet labeled training data as in [8]. We report both top-1 and top-5 accuracies on the test set in Table 2. BYOL consistently outperforms previous approaches across a wide range of architectures. Additionally, as detailed in Appendix C.1, BYOL reaches 77.7% top-1 accuracy with ResNet-50 when fine-tuning over 100% of ImageNet labels.

Method	Top-1		Top-5		Method	Architecture	Param.	Top-1		Top-5	
	1%	10%	1%	10%				1%	10%	1%	10%
Supervised [77]	25.4	56.4	48.4	80.4	CPC v2 [32]	ResNet-161	305M	-	-	77.9	91.2
InstDisc	-	-	39.2	77.4	SimCLR [8]	ResNet-50 ($\times 2$)	94M	58.5	71.7	83.0	91.2
PIRL [35]	-	-	57.2	83.8	BYOL (ours)	ResNet-50 ($\times 2$)	94M	62.2	73.5	84.1	91.7
SimCLR [8]	48.3	65.6	75.5	87.8	SimCLR [8]	ResNet-50 ($\times 4$)	375M	63.0	74.4	85.8	92.6
BYOL (ours)	53.2	68.8	78.4	89.0	BYOL (ours)	ResNet-50 ($\times 4$)	375M	69.1	75.7	87.9	92.5
					BYOL (ours)	ResNet-200 ($\times 2$)	250M	71.2	77.7	89.5	93.7

(a) ResNet-50 encoder.

(b) Other ResNet encoder architectures.

Table 2: Semi-supervised training with a fraction of ImageNet labels.

Transfer to other classification tasks We evaluate our representation on other classification datasets to assess whether the features learned on ImageNet (IN) are generic and thus useful across image domains, or if they are ImageNet-specific. We perform linear evaluation and fine-tuning on the same set of classification tasks used in [8, 74], and carefully follow their evaluation protocol, as detailed in Appendix D. Performance is reported using standard metrics for each benchmark, and results are provided on a held-out test set after hyperparameter selection on a validation set. We report results in Table 3, both for linear evaluation and fine-tuning. BYOL outperforms SimCLR on all benchmarks and the Supervised-IN baseline on 7 of the 12 benchmarks, providing only slightly worse performance on the 5 remaining benchmarks. BYOL’s representation can be transferred over to small images, e.g., CIFAR [78], landscapes, e.g., SUN397 [79] or VOC2007 [80], and textures, e.g., DTD [81].

Method	Food101	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	VOC2007	DTD	Pets	Caltech-101	Flowers
<i>Linear evaluation:</i>												
BYOL (ours)	75.3	91.3	78.4	57.2	62.2	67.8	60.6	82.5	75.5	90.4	94.2	96.1
SimCLR (repro)	72.8	90.5	74.4	42.4	60.6	49.3	49.8	81.4	75.7	84.6	89.3	92.6
SimCLR [8]	68.4	90.6	71.6	37.4	58.8	50.3	50.3	80.5	74.5	83.6	90.3	91.2
Supervised-IN [8]	72.3	93.6	78.3	53.7	61.9	66.7	61.0	82.8	74.9	91.5	94.5	94.7
<i>Fine-tuned:</i>												
BYOL (ours)	88.5	97.8	86.1	76.3	63.7	91.6	88.1	85.4	76.2	91.7	93.8	97.0
SimCLR (repro)	87.5	97.4	85.3	75.0	63.9	91.4	87.6	84.5	75.4	89.4	91.7	96.6
SimCLR [8]	88.2	97.7	85.9	75.9	63.5	91.3	88.1	84.1	73.2	89.2	92.1	97.0
Supervised-IN [8]	88.3	97.5	86.4	75.8	64.3	92.1	86.0	85.0	74.6	92.1	93.3	97.6
Random init [8]	86.9	95.9	80.2	76.1	53.6	91.4	85.9	67.3	64.8	81.5	72.6	92.0

Table 3: Transfer learning results from ImageNet (IN) with the standard ResNet-50 architecture.

Transfer to other vision tasks We evaluate our representation on different tasks relevant to computer vision practitioners, namely semantic segmentation, object detection and depth estimation. With this evaluation, we assess whether BYOL’s representation generalizes beyond classification tasks.

We first evaluate BYOL on the VOC2012 semantic segmentation task as detailed in Appendix D.4, where the goal is to classify each pixel in the image [7]. We report the results in Table 4a. BYOL outperforms both the Supervised-IN baseline (+1.9 mIoU) and SimCLR (+1.1 mIoU).

Similarly, we evaluate on object detection by reproducing the setup in [9] using a Faster R-CNN architecture [82], as detailed in Appendix D.5. We fine-tune on trainval2007 and report results on test2007 using the standard AP₅₀ metric; BYOL is significantly better than the Supervised-IN baseline (+3.1 AP₅₀) and SimCLR (+2.3 AP₅₀).

Finally, we evaluate on depth estimation on the NYU v2 dataset, where the depth map of a scene is estimated given a single RGB image. Depth prediction measures how well a network represents geometry, and how well that information can be localized to pixel accuracy [40]. The setup is based on [83] and detailed in Appendix D.6. We evaluate on the commonly used test subset of 654 images and report results using several common metrics in Table 4b: relative (rel) error, root mean squared (rms) error, and the percent of pixels (pct) where the error, $\max(d_{gt}/d_p, d_p/d_{gt})$, is below 1.25ⁿ thresholds where d_p is the predicted depth and d_{gt} is the ground truth depth [40]. BYOL is better or on par with other methods for each metric. For instance, the challenging pct.<1.25 measure is respectively improved by +3.5 points and +1.3 points compared to supervised and SimCLR baselines.

Method	AP ₅₀	mIoU	Higher better					Lower better	
			Method	pct.<1.25	pct.<1.25 ²	pct.<1.25 ³	rms	rel	
Supervised-IN [9]	74.4	74.4	Supervised-IN [83]	81.1	95.3	98.8	0.573	0.127	
MoCo [9]	74.9	72.5	SimCLR (repro)	83.3	96.5	99.1	0.557	0.134	
SimCLR (repro)	75.2	75.2	BYOL (ours)	84.6	96.7	99.1	0.541	0.129	
BYOL (ours)	77.5	76.3							

(a) Transfer results in semantic segmentation and object detection.

(b) Transfer results on NYU v2 depth estimation.

Table 4: Results on transferring BYOL’s representation to other vision tasks.

5 Building intuitions with ablations

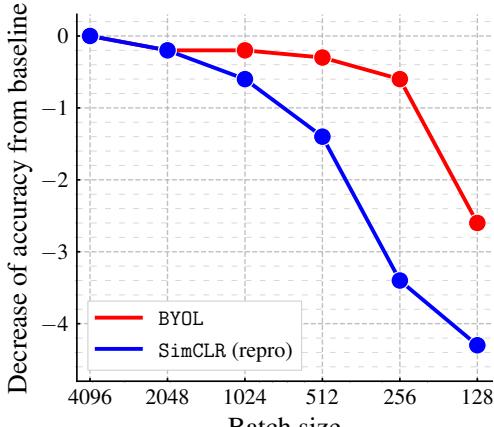
We present ablations on BYOL to give an intuition of its behavior and performance. For reproducibility, we run each configuration of parameters over three seeds, and report the average performance. We also report the half difference between the best and worst runs when it is larger than 0.25. Although previous works perform ablations at 100 epochs [8, 12], we notice that relative improvements at 100 epochs do not always hold over longer training. For this reason, we run ablations over 300 epochs on 64 TPU v3 cores, which yields consistent results compared to our baseline training of 1000 epochs. For all the experiments in this section, we set the initial learning rate to 0.3 with batch size 4096, the weight decay to 10⁻⁶ as in SimCLR [8] and the base target decay rate τ_{base} to 0.99. In this section we report results in top-1 accuracy on ImageNet under the linear evaluation protocol as in Appendix C.1.

Batch size Among contrastive methods, the ones that draw negative examples from the minibatch suffer performance drops when their batch size is reduced. BYOL does not use negative examples and we expect it to be more

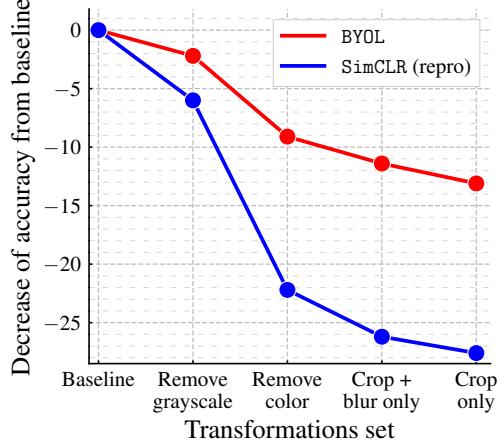
by means of observation or experience rather than theory or pure logic.

robust to smaller batch sizes. To empirically verify this hypothesis, we train both BYOL and SimCLR using different batch sizes from 128 to 4096. To avoid re-tuning other hyperparameters, we average gradients over N consecutive steps before updating the online network when reducing the batch size by a factor N . The target network is updated once every N steps, after the update of the online network; we accumulate the N -steps in parallel in our runs.

As shown in Figure 3a, the performance of SimCLR rapidly deteriorates with batch size, likely due to the decrease in the number of negative examples. In contrast, the performance of BYOL remains stable over a wide range of batch sizes from 256 to 4096, and only drops for smaller values due to batch normalization layers in the encoder.⁷



(a) Impact of batch size



(b) Impact of progressively removing transformations

Figure 3: Decrease in top-1 accuracy (in % points) of BYOL and our own reproduction of SimCLR at 300 epochs, under linear evaluation on ImageNet.

Image augmentations Contrastive methods are sensitive to the choice of image augmentations. For instance, SimCLR does not work well when removing color distortion from its image augmentations. As an explanation, SimCLR shows that crops of the same image mostly share their color histograms. At the same time, color histograms vary across images. Therefore, when a contrastive task only relies on random crops as image augmentations, it can be mostly solved by focusing on color histograms alone. As a result the representation is not incentivized to retain information beyond color histograms. To prevent that, SimCLR adds color distortion to its set of image augmentations. Instead, BYOL is incentivized to keep any information captured by the target representation into its online network, to improve its predictions. Therefore, even if augmented views of a same image share the same color histogram, BYOL is still incentivized to retain additional features in its representation. For that reason, we believe that BYOL is more robust to the choice of image augmentations than contrastive methods.

Results presented in Figure 3b support this hypothesis: the performance of BYOL is much less affected than the performance of SimCLR when removing color distortions from the set of image augmentations (-9.1 accuracy points for BYOL, -22.2 accuracy points for SimCLR). When image augmentations are reduced to mere random crops, BYOL still displays good performance (59.4%, i.e. -13.1 points from 72.5%), while SimCLR loses more than a third of its performance (40.3%, i.e. -27.6 points from 67.9%). We report additional ablations in Appendix F.3.

Bootstrapping BYOL uses the projected representation of a target network, whose weights are an exponential moving average of the weights of the online network, as target for its predictions. This way, the weights of the target network represent a delayed and more stable version of the weights of the online network. When the target decay rate is 1, the target network is never updated, and remains at a constant value corresponding to its initialization. When the target decay rate is 0, the target network is instantaneously updated to the online network at each step. There is a trade-off between updating the targets too often and updating them too slowly, as illustrated in Table 5a. Instantaneously updating the target network ($\tau = 0$) destabilizes training, yielding very poor performance while never updating the target ($\tau = 1$) makes the training stable but prevents iterative improvement, ending with low-quality final representation. All values of the decay rate between 0.9 and 0.999 yield performance above 68.4% top-1 accuracy at 300 epochs.

⁷The only dependency on batch size in our training pipeline sits within the batch normalization layers.

Target	τ_{base}	Top-1	Method	Predictor	Target network	β	Top-1
Constant random network	1	18.8±0.7	BYOL	✓	✓	0	72.5
Moving average of online	0.999	69.8	—	✓	✓	1	70.9
Moving average of online	0.99	72.5	—		✓	1	70.7
Moving average of online	0.9	68.4	SimCLR			1	69.4
Stop gradient of online [†]	0	0.3	—	✓		1	69.1
			—	✓		0	0.3
			—		✓	0	0.2
			—			0	0.1

(a) Results for different target modes. [†]In the *stop gradient of online*, $\tau = \tau_{\text{base}} = 0$ is kept constant throughout training.

(b) Intermediate variants between BYOL and SimCLR.

Table 5: Ablations with top-1 accuracy (in %) at 300 epochs under linear evaluation on ImageNet.

Ablation to contrastive methods In this subsection, we recast SimCLR and BYOL using the same formalism to better understand where the improvement of BYOL over SimCLR comes from. Let us consider the following objective that extends the InfoNCE objective [10, 84] (see Appendix F.4),

$$\text{InfoNCE}_{\theta}^{\alpha, \beta} \triangleq \frac{2}{B} \sum_{i=1}^B S_{\theta}(v_i, v'_i) - \beta \cdot \frac{2\alpha}{B} \sum_{i=1}^B \ln \left(\sum_{j \neq i} \exp \frac{S_{\theta}(v_i, v_j)}{\alpha} + \sum_j \exp \frac{S_{\theta}(v_i, v'_j)}{\alpha} \right), \quad (6)$$

where $\alpha > 0$ is a fixed temperature, $\beta \in [0, 1]$ a weighting coefficient, B the batch size, v and v' are batches of augmented views where for any batch index i , v_i and v'_i are augmented views from the same image; the real-valued function S_{θ} quantifies pairwise similarity between augmented views. For any augmented view u we denote $z_{\theta}(u) \triangleq f_{\theta}(g_{\theta}(u))$ and $z_{\xi}(u) \triangleq f_{\xi}(g_{\xi}(u))$. For given ϕ and ψ , we consider the normalized dot product

$$S_{\theta}(u_1, u_2) \triangleq \frac{\langle \phi(u_1), \psi(u_2) \rangle}{\|\phi(u_1)\|_2 \cdot \|\psi(u_2)\|_2}. \quad (7)$$

Up to minor details (cf. Appendix F.5), we recover the SimCLR loss with $\phi(u_1) = z_{\theta}(u_1)$ (no predictor), $\psi(u_2) = z_{\theta}(u_2)$ (no target network) and $\beta = 1$. We recover the BYOL loss when using a predictor and a target network, i.e., $\phi(u_1) = p_{\theta}(z_{\theta}(u_1))$ and $\psi(u_2) = z_{\xi}(u_2)$ with $\beta = 0$. To evaluate the influence of the target network, the predictor and the coefficient β , we perform an ablation over them. Results are presented in Table 5b and more details are given in Appendix F.4.

The only variant that performs well without negative examples (i.e., with $\beta = 0$) is BYOL, using both a bootstrap target network and a predictor. Adding the negative pairs to BYOL’s loss without re-tuning the temperature parameter hurts its performance. In Appendix F.4, we show that we can add back negative pairs and still match the performance of BYOL with proper tuning of the temperature.

Simply adding a target network to SimCLR already improves performance (+1.6 points). This sheds new light on the use of the target network in MoCo [9], where the target network is used to provide more negative examples. Here, we show that by mere stabilization effect, even when using the same number of negative examples, using a target network is beneficial. Finally, we observe that modifying the architecture of S_{θ} to include a predictor only mildly affects the performance of SimCLR.

Network hyperparameters In Appendix F, we explore how other network parameters may impact BYOL’s performance. We iterate over multiple weight decays, learning rates, and projector/encoder architectures to observe that small hyperparameter changes do not drastically alter the final score. We note that removing the weight decay in either BYOL or SimCLR leads to network divergence, emphasizing the need for weight regularization in the self-supervised setting. Furthermore, we observe that changing the scaling factor in the network initialization [85] did not impact the performance (higher than 72% top-1 accuracy).

Relationship with Mean Teacher Another semi-supervised approach, Mean Teacher (MT) [20], complements a supervised loss on few labels with an additional consistency loss. In [20], this consistency loss is the ℓ_2 distance between the logits from a *student* network, and those of a temporally averaged version of the student network, called *teacher*. Removing the predictor in BYOL results in an unsupervised version of MT with no classification loss that

uses image augmentations instead of the original architectural noise (e.g., dropout). This variant of BYOL collapses (Row 7 of Table 5) which suggests that the additional predictor is critical to prevent collapse in an unsupervised scenario.

Importance of a near-optimal predictor Table 5b already shows the importance of combining a predictor and a target network: the representation does collapse when either is removed. We further found that we can remove the target network without collapse by making the predictor near-optimal, either by (i) using an optimal *linear* predictor (obtained by linear regression on the current batch) before back-propagating the error through the network (52.5% top-1 accuracy), or (ii) increasing the learning rate of the predictor (66.5% top-1). By contrast, increasing the learning rates of both projector *and* predictor (without target network) yields poor results ($\approx 25\%$ top-1). See Appendix I for more details. This seems to indicate that keeping the predictor near-optimal at all times is important to preventing collapse, which may be one of the roles of BYOL’s target network.

6 Conclusion

We introduced BYOL, a new algorithm for self-supervised learning of image representations. BYOL learns its representation by predicting previous versions of its outputs, without using negative pairs. We show that BYOL achieves state-of-the-art results on various benchmarks. In particular, under the linear evaluation protocol on ImageNet with a ResNet-50 (1 \times), BYOL achieves a new state of the art and bridges most of the remaining gap between self-supervised methods and the supervised learning baseline of [8]. Using a ResNet-200 (2 \times), BYOL reaches a top-1 accuracy of 79.6% which improves over the previous state of the art (76.8%) while using 30% fewer parameters.

Nevertheless, BYOL remains dependent on existing sets of augmentations that are specific to vision applications. To generalize BYOL to other modalities (e.g., audio, video, text, ...) it is necessary to obtain similarly suitable augmentations for each of them. Designing such augmentations may require significant effort and expertise. Therefore, automating the search for these augmentations would be an important next future step to generalize BYOL to other modalities.

Broader impact

The presented research should be categorized as research in the field of unsupervised learning. This work may inspire new algorithms, theoretical, and experimental investigation. The algorithm presented here can be used for many different vision applications and a particular use may have both positive or negative impacts, which is known as the dual use problem. Besides, as vision datasets could be biased, the representation learned by BYOL could be susceptible to replicate these biases.

Acknowledgements

The authors would like to thank the following people for their help throughout the process of writing this paper, in alphabetical order: Aaron van den Oord, Andrew Brock, Jason Ramapuram, Jeffrey De Fauw, Karen Simonyan, Katrina McKinney, Nathalie Beauguerlange, Olivier Henaff, Oriol Vinyals, Pauline Luc, Razvan Pascanu, Sander Dieleman, and the DeepMind team. We especially thank Jason Ramapuram and Jeffrey De Fauw, who provided the JAX SimCLR reproduction used throughout the paper.

References

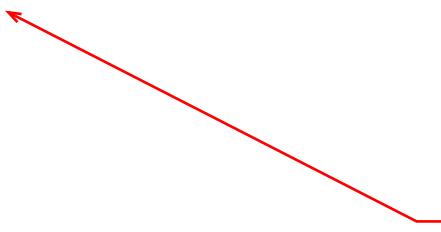
- [1] Kunihiko Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36(4):193–202, 1980.
- [2] Laurenz Wiskott and Terrence J Sejnowski. Slow feature analysis: Unsupervised learning of invariances. *Neural Computation*, 14(4), 2002.
- [3] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18(7):1527–1554, 2006.
- [4] Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *Computer Vision and Pattern Recognition*, 2014.
- [5] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [6] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Computer Vision and Pattern Recognition*, 2014.
- [7] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Computer Vision and Pattern Recognition*, 2015.
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.
- [9] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, 2019.
- [10] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [11] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849v4*, 2019.
- [12] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning. *arXiv preprint arXiv:2005.10243*, 2020.
- [13] 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*, 2019.
- [14] R. Manmatha, Chao-Yuan Wu, Alexander J. Smola, and Philipp Krähenbühl. Sampling matters in deep embedding learning. In *International Conference on Computer Vision*, 2017.
- [15] Ben Harwood, Vijay B. G. Kumar, Gustavo Carneiro, Ian Reid, and Tom Drummond. Smart mining for deep metric learning. In *International Conference on Computer Vision*, 2017.
- [16] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *International Conference on Machine Learning*, 2013.
- [17] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *European Conference on Computer Vision*, 2018.
- [18] Philip Bachman, Ouais Alsharif, and Doina Precup. Learning with pseudo-ensembles. In *Advances in neural information processing systems*, pages 3365–3373, 2014.
- [19] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*, 2016.

- [20] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advances in neural information processing systems*, pages 1195–1204, 2017.
- [21] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition*, 2016.
- [23] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Computer Vision and Pattern Recognition*, 2015.
- [24] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *International Conference on Machine Learning*, 2008.
- [25] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [26] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic back-propagation and variational inference in deep latent gaussian models. *arXiv preprint arXiv:1401.4082*, 2014.
- [27] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Neural Information Processing Systems*, 2014.
- [28] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.
- [29] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Alex Lamb, Martín Arjovsky, Olivier Mastropietro, and Aaron C. Courville. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2017.
- [30] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning. In *Neural Information Processing Systems*, 2019.
- [31] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [32] Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, and Aäron van den Oord. Data-efficient image recognition with contrastive predictive coding. In *International Conference on Machine Learning*, 2019.
- [33] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. *arXiv preprint arXiv:1808.06670*, 2019.
- [34] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. In *Neural Information Processing Systems*, 2019.
- [35] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. *arXiv preprint arXiv:1912.01991*, 2019.
- [36] Junnan Li, Pan Zhou, Caiming Xiong, Richard Socher, and Steven CH Hoi. Prototypical contrastive learning of unsupervised representations. *arXiv preprint arXiv:2005.04966*, 2020.
- [37] Rishabh Jain, Haoqi Fan, Ross B. Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [38] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020.
- [39] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Computer Vision and Pattern Recognition*, 2018.
- [40] Carl Doersch and Andrew Zisserman. Multi-task self-supervised visual learning. In *International Conference on Computer Vision*, 2017.
- [41] Richard Zhang, Phillip Isola, and Alexei A. Efros. Colorful image colorization. In *European Conference on Computer Vision*, 2016.
- [42] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *European Conference on Computer Vision*, 2016.
- [43] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *Computer Vision and Pattern Recognition*, 2016.
- [44] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, 2016.
- [45] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Computer Vision and Pattern Recognition*, 2017.

- [46] Alexey Dosovitskiy, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with convolutional neural networks. In *Neural Information Processing Systems*, 2014.
- [47] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728*, 2018.
- [48] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual representation learning. In *Computer Vision and Pattern Recognition*, 2019.
- [49] Daniel Guo, Bernardo Avila Pires, Bilal Piot, Jean-Bastien Grill, Florent Alché, Rémi Munos, and Mohammad Gheshlaghi Azar. Bootstrap latent-predictive representations for multitask reinforcement learning. In *International Conference on Machine Learning*, 2020.
- [50] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles 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.
- [51] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*, 2016.
- [52] Matteo Hessel, Joseph Modayil, Hado Van 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*, 2018.
- [53] Hado Van Hasselt, Yotam Doron, Florian Strub, Matteo Hessel, Nicolas Sonnerat, and Joseph Modayil. Deep reinforcement learning and the deadly triad. *Deep Reinforcement Learning Workshop NeurIPS*, 2018.
- [54] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [55] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. Semi-supervised learning. *IEEE Transactions on Neural Networks*, 20(3):542–542, 2009.
- [56] Xiaojin Zhu and Andrew B Goldberg. Introduction to semi-supervised learning. *Synthesis lectures on artificial intelligence and machine learning*, 3(1):1–130, 2009.
- [57] Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *Advances in neural information processing systems*, 2014.
- [58] Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semi-supervised learning with ladder networks. In *Advances in neural information processing systems*, 2015.
- [59] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *Advances in Neural Information Processing Systems*, 2019.
- [60] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018.
- [61] David Berthelot, N. Carlini, E. D. Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution matching and augmentation anchoring. In *ICLR*, 2020.
- [62] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *arXiv preprint arXiv:2001.07685*, 2020.
- [63] Suzanna Becker and Geoffrey E. Hinton. Self-organizing neural network that discovers surfaces in random-dot stereograms. *Nature*, 355(6356):161–163, 1992.
- [64] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, and Skye Wanderman-Milne. JAX: composable transformations of Python+NumPy programs, 2018.
- [65] Tom Hennigan, Trevor Cai, Tamara Norman, and Igor Babuschkin. Haiku: Sonnet for JAX, 2020.
- [66] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [67] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [68] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, 2015.
- [69] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *International Conference on Machine Learning*, 2010.
- [70] Yang You, Igor Gitman, and Boris Ginsburg. Scaling SGD batch size to 32k for imagenet training. *arXiv preprint arXiv:1708.03888*, 2017.

- [71] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *International Conference on Learning Representations*, 2017.
- [72] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- [73] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [74] Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better ImageNet models transfer better? In *Computer Cision and Pattern Recognition*, 2019.
- [75] Chengyue Gong, Tongzheng Ren, Mao Ye, and Qiang Liu. Maxup: A simple way to improve generalization of neural network training. *arXiv preprint arXiv:2002.09024*, 2020.
- [76] Xiaohua Zhai, Joan Puigcerver, Alexander I Kolesnikov, Pierre Ruyssen, Carlos Riquelme, Mario Lucic, Josip Djolonga, André Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, Lucas Beyer, Olivier Bachem, Michael Tschannen, Marcin Michalski, Olivier Bousquet, Sylvain Gelly, and Neil Houlsby. A large-scale study of representation learning with the visual task adaptation benchmark. *arXiv: Computer Vision and Pattern Recognition*, 2019.
- [77] Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S4L: Self-supervised semi-supervised learning. In *International Conference on Computer Vision*, 2019.
- [78] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [79] Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *Computer Vision and Pattern Recognition*, 2010.
- [80] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The Pascal visual object classes (VOC) challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010.
- [81] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Computer Vision and Pattern Recognition*, 2014.
- [82] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Neural Information Processing Systems*, 2015.
- [83] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab. Deeper depth prediction with fully convolutional residual networks. In *International Conference on 3D Vision*, 2016.
- [84] Ben Poole, Sherjil Ozair, Aaron van den Oord, Alexander A Alemi, and George Tucker. On variational bounds of mutual information. *arXiv preprint arXiv:1905.06922*, 2019.
- [85] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, 2015.
- [86] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. On the difficulty of training recurrent neural networks. In *International Conference on Machine Learning*, 2013.
- [87] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. *arXiv preprint arXiv:1909.13719*, 2019.
- [88] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Computer Vision and Pattern Recognition*, pages 2818–2826, 2016.
- [89] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *European Conference on Computer Vision*, 2014.
- [90] Thomas Berg, Jiongxin Liu, Seung Woo Lee, Michelle L. Alexander, David W. Jacobs, and Peter N. Belhumeur. Birdsnap: Large-scale fine-grained visual categorization of birds. In *Computer Vision and Pattern Recognition*, 2014.
- [91] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3D object representations for fine-grained categorization. In *Workshop on 3D Representation and Recognition*, Sydney, Australia, 2013.
- [92] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew B. Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- [93] O. M. Parkhi, A. Vedaldi, A. Zisserman, and C. V. Jawahar. Cats and dogs. In *Computer Vision and Pattern Recognition*, 2012.
- [94] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. *Computer Vision and Pattern Recognition Workshop*, 2004.
- [95] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*, 2008.
- [96] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *International Conference on Machine Learning*, 2014.
- [97] Art B Owen. A robust hybrid of lasso and ridge regression. *Contemporary Mathematics*, 443(7):59–72, 2007.

- [98] Chengxu Zhuang, Alex Lin Zhai, and Daniel Yamins. Local aggregation for unsupervised learning of visual embeddings. In *International Conference on Computer Vision*, 2019.
- [99] Deepak Pathak, Ross Girshick, Piotr Dollár, Trevor Darrell, and Bharath Hariharan. Learning features by watching objects move. In *Conference on Computer Vision and Pattern Recognition*, 2017.
- [100] Yuxin Wu and Kaiming He. Group normalization. In *European Conference on Computer Vision*, 2018.



Using an optimal linear predictor in the context of BYOL involves employing linear regression to find a linear mapping that minimizes the prediction error on the current batch of data. Let's break down the process:

Objective:
The goal is to predict representations in an optimal way without relying on a slowly moving target network. Linear regression aims to find the optimal linear predictor, which is a linear function that maps input representations to predicted representations.

Linear Regression:
Linear regression is a supervised learning algorithm used for predicting a continuous target variable (output) based on one or more predictor variables (inputs). In the context of BYOL, the representations learned by the online network can be considered as the predictor variables.

Mathematical Formulation:
The linear predictor is represented as a linear equation: $y=mx+b$, where: y is the predicted representation, x is the input representation, m is the slope (weights) of the linear function, b is the intercept.

Optimization:
The optimization process involves adjusting the parameters (weights and intercept) to minimize the difference between the predicted representations and the true representations on the current batch. The objective is to minimize the mean squared error (MSE) or another suitable loss function.

Results:
After the optimization process, the linear predictor becomes near-optimal for the given batch of data. The resulting linear predictor is then used for making predictions on subsequent batches.

Benefits:
Using an optimal linear predictor is suggested as a strategy to remove the reliance on a slowly moving target network while maintaining or improving the quality of predictions.

In summary, the process involves applying linear regression to find the optimal linear predictor based on the representations learned by the online network. This predictor serves as a replacement for the target network in guiding the learning process. The key idea is to derive an optimal linear mapping that captures the relationship between input and predicted representations.

A Algorithm

Algorithm 1: BYOL: Bootstrap Your Own Latent

Inputs :

\mathcal{D}, \mathcal{T} , and \mathcal{T}'	set of images and distributions of transformations
$\theta, f_\theta, g_\theta$, and q_θ	initial online parameters, encoder, projector, and predictor
ξ, f_ξ, g_ξ	initial target parameters, target encoder, and target projector
optimizer	optimizer, updates online parameters using the loss gradient
K and N	total number of optimization steps and batch size
$\{\tau_k\}_{k=1}^K$ and $\{\eta_k\}_{k=1}^K$	target network update schedule and learning rate schedule

```

1 for  $k = 1$  to  $K$  do
2    $\mathcal{B} \leftarrow \{x_i \sim \mathcal{D}\}_{i=1}^N$                                 // sample a batch of  $N$  images
3   for  $x_i \in \mathcal{B}$  do
4      $t \sim \mathcal{T}$  and  $t' \sim \mathcal{T}'$                                          // sample image transformations
5      $z_1 \leftarrow g_\theta(f_\theta(t(x_i)))$  and  $z_2 \leftarrow g_\theta(f_\theta(t'(x_i)))$  // compute projections
6      $z'_1 \leftarrow g_\xi(f_\xi(t'(x_i)))$  and  $z'_2 \leftarrow g_\xi(f_\xi(t(x_i)))$  // compute target projections
7      $l_i \leftarrow -2 \cdot \left( \frac{\langle q_\theta(z_1), z'_1 \rangle}{\|q_\theta(z_1)\|_2 \cdot \|z'_1\|_2} + \frac{\langle q_\theta(z_2), z'_2 \rangle}{\|q_\theta(z_2)\|_2 \cdot \|z'_2\|_2} \right)$  // compute the loss for  $x_i$ 
8   end
9    $\delta\theta \leftarrow \frac{1}{N} \sum_{i=1}^N \partial_\theta l_i$                                 // compute the total loss gradient w.r.t.  $\theta$ 
10   $\theta \leftarrow \text{optimizer}(\theta, \delta\theta, \eta_k)$                                      // update online parameters
11   $\xi \leftarrow \tau_k \xi + (1 - \tau_k) \theta$                                          // update target parameters
12 end
Output:encoder  $f_\theta$ 

```

B Image augmentations

During self-supervised training, BYOL uses the following image augmentations (which are a subset of the ones presented in [8]):

- **random cropping**: a random patch of the image is selected, with an area uniformly sampled between 8% and 100% of that of the original image, and an aspect ratio logarithmically sampled between 3/4 and 4/3. This patch is then resized to the target size of 224×224 using bicubic interpolation;
- optional **left-right flip**;
- **color jittering**: the brightness, contrast, saturation and hue of the image are shifted by a uniformly random offset applied on all the pixels of the same image. The order in which these shifts are performed is randomly selected for each patch;
- **color dropping**: an optional conversion to grayscale. When applied, output intensity for a pixel (r, g, b) corresponds to its luma component, computed as $0.2989r + 0.5870g + 0.1140b$;
- **Gaussian blurring**: for a 224×224 image, a square Gaussian kernel of size 23×23 is used, with a standard deviation uniformly sampled over $[0.1, 2.0]$;
- **solarization**: an optional color transformation $x \mapsto x \cdot \mathbf{1}_{\{x < 0.5\}} + (1 - x) \cdot \mathbf{1}_{\{x \geq 0.5\}}$ for pixels with values in $[0, 1]$.

Augmentations from the sets \mathcal{T} and \mathcal{T}' (introduced in Section 3) are compositions of the above image augmentations in the listed order, each applied with a predetermined probability. The image augmentations parameters are listed in Table 6.

During evaluation, we use a **center crop** similar to [8]: images are resized to 256 pixels along the shorter side using bicubic resampling, after which a 224×224 center crop is applied. In both training and evaluation, we normalize color channels by subtracting the average color and dividing by the standard deviation, computed on ImageNet, after applying the augmentations.

Parameter	\mathcal{T}	\mathcal{T}'
Random crop probability	1.0	1.0
Flip probability	0.5	0.5
Color jittering probability	0.8	0.8
Brightness adjustment max intensity	0.4	0.4
Contrast adjustment max intensity	0.4	0.4
Saturation adjustment max intensity	0.2	0.2
Hue adjustment max intensity	0.1	0.1
Color dropping probability	0.2	0.2
Gaussian blurring probability	1.0	0.1
Solarization probability	0.0	0.2

Table 6: Parameters used to generate image augmentations.

C Evaluation on ImageNet training

C.1 Self-supervised learning evaluation on ImageNet

Linear evaluation protocol on ImageNet As in [48, 74, 8, 37], we use the standard linear evaluation protocol on ImageNet, which consists in training a linear classifier on top of the frozen representation, *i.e.*, without updating the network parameters nor the batch statistics. At training time, we apply spatial augmentations, *i.e.*, random crops with resize to 224×224 pixels, and random flips. At test time, images are resized to 256 pixels along the shorter side using bicubic resampling, after which a 224×224 center crop is applied. In both cases, we normalize the color channels by subtracting the average color and dividing by the standard deviation (computed on ImageNet), after applying the augmentations. We optimize the cross-entropy loss using SGD with Nesterov momentum over 80 epochs, using a batch size of 1024 and a momentum of 0.9. We do not use any regularization methods such as weight decay, gradient clipping [86], tclip [34], or logits regularization. We finally sweep over 5 learning rates $\{0.4, 0.3, 0.2, 0.1, 0.05\}$ on a local validation set (10009 images from ImageNet train set), and report the accuracy of the best validation hyperparameter on the test set (which is the public validation set of the original ILSVRC2012 ImageNet dataset).

Variant on linear evaluation on ImageNet In this paragraph only, we deviate from the protocol of [8, 37] and propose another way of performing linear evaluation on top of a frozen representation. This method achieves better performance both in top-1 and top-5 accuracy.

- We replace the spatial augmentations (random crops with resize to 224×224 pixels and random flips) with the pre-train augmentations of Appendix B. This method was already used in [32] with a different subset of pre-train augmentations.
- We regularize the linear classifier as in [34]⁸ by clipping the logits using a hyperbolic tangent function

$$\text{tclip}(x) \triangleq \alpha \cdot \tanh(x/\alpha),$$

where α is a positive scalar, and by adding a logit-regularization penalty term in the loss

$$\text{Loss}(x, y) \triangleq \text{cross_entropy}(\text{tclip}(x), y) + \beta \cdot \text{average}(\text{tclip}(x)^2),$$

where x are the logits, y are the target labels, and β is the regularization parameter. We set $\alpha = 20$ and $\beta = 1e-2$.

We report in Table 7 the top-1 and top-5 accuracy on ImageNet using this modified protocol. These modifications in the evaluation protocol increase the BYOL’s top-1 accuracy from 74.3% to 74.8% with a ResNet-50 ($1\times$).

Semi-supervised learning on ImageNet We follow the semi-supervised learning protocol of [8, 77]. We first initialize the network with the parameters of the pretrained representation, and fine-tune it with a subset of ImageNet labels. At training time, we apply spatial augmentations, *i.e.*, random crops with resize to 224×224 pixels and random flips. At test time, images are resized to 256 pixels along the shorter side using bicubic resampling, after which a 224×224 center crop is applied. In both cases, we normalize the color channels by subtracting the average color and dividing by the standard deviation (computed on ImageNet), after applying the augmentations. We optimize the cross-entropy loss using SGD with Nesterov momentum. We used a batch size of 1024, a momentum of

⁸<https://github.com/Philip-Bachman/amdim-public/blob/master/costs.py>

Architecture	Pre-train augmentations	Logits regularization	Top-1	Top-5
ResNet-50 (1×)	✓		74.3	91.6
		✓	74.4	91.8
	✓	✓	74.8	91.8
ResNet-50 (4×)	✓		78.6	94.2
		✓	78.6	94.3
	✓	✓	79.0	94.5
ResNet-200 (2×)	✓		79.6	94.8
		✓	79.6	94.8
	✓	✓	79.8	95.0
			80.0	95.0

Table 7: Different linear evaluation protocols on ResNet architectures by either replacing the spatial augmentations with pre-train augmentations, or regularizing the linear classifier. No pre-train augmentations and no logits regularization correspond to the evaluation protocol of the main paper, which is the same as in [8, 37].

0.9. We do not use any regularization methods such as weight decay, gradient clipping [86], tclip [34], or logits rescaling. We sweep over the learning rate {0.01, 0.02, 0.05, 0.1, 0.005} and the number of epochs {30, 50} and select the hyperparameters achieving the best performance on our local validation set to report test performance.

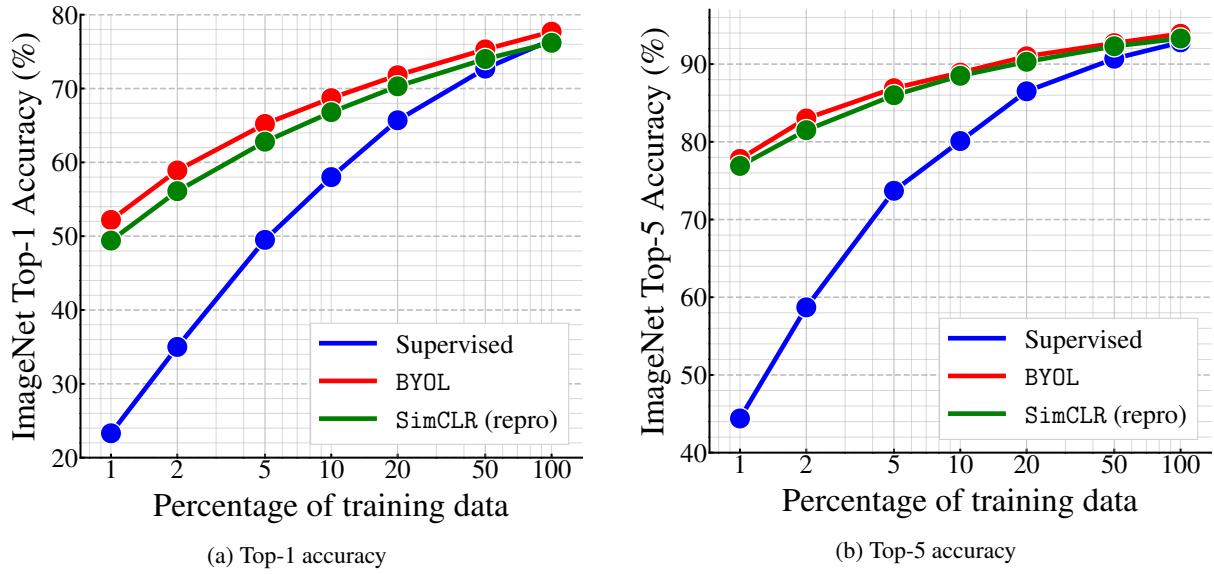


Figure 4: Semi-supervised training with a fraction of ImageNet labels on a ResNet-50 (1×).

In Table 2 presented in the main text, we fine-tune the representation over the 1% and 10% ImageNet splits from [8] with various ResNet architectures.

In Figure 4, we fine-tune the representation over 1%, 2%, 5%, 10%, 20%, 50%, and 100% of the ImageNet dataset as in [32] with a ResNet-50 (1×) architecture, and compare them with a supervised baseline and a fine-tuned SimCLR representation. In this case and contrary to Table 2 we don't reuse the splits from SimCLR but we create our own via a balanced selection. In this setting, we observed that tuning a BYOL representation always outperforms a supervised baseline trained from scratch. In Figure 5, we then fine-tune the representation over multiple ResNet architectures. We observe that the largest networks are prone to overfitting as they are outperformed by ResNets with identical depth but smaller scaling factor. This overfitting is further confirmed when looking at the training and evaluation loss: large networks have lower training losses, but higher validation losses than some of their slimmer counterparts. Regularization methods are thus recommended when tuning on large architectures.

Supervised:			Semi-supervised (100%):		
Method	Top-1	Top-5	Method	Top-1	Top-5
Supervised[8]	76.5	—	SimCLR [8]	76.0	93.1
AutoAugment [87]	77.6	93.8	SimCLR (repro)	76.5	93.5
MaxUp [75]	78.9	94.2	BYOL	77.7	93.9

Table 8: Semi-supervised training with the full ImageNet on a ResNet-50 ($\times 1$). We also report other fully supervised methods for extensive comparisons.

Finally, we fine-tune the representation over the full ImageNet dataset. We report the results in Table 8 along with supervised baselines trained on ImageNet. We observe that fine-tuning the SimCLR checkpoint does not yield better results (in our reproduction, which matches the results reported in the original paper [8]) than using a random initialization (76.5 top-1). Instead, BYOL’s initialization checkpoint leads to a high final score (77.7 top-1), higher than the vanilla supervised baseline of [8], matching the strong supervised baseline of AutoAugment[87] but still 1.2 points below the stronger supervised baseline [75], which uses advanced supervised learning techniques.

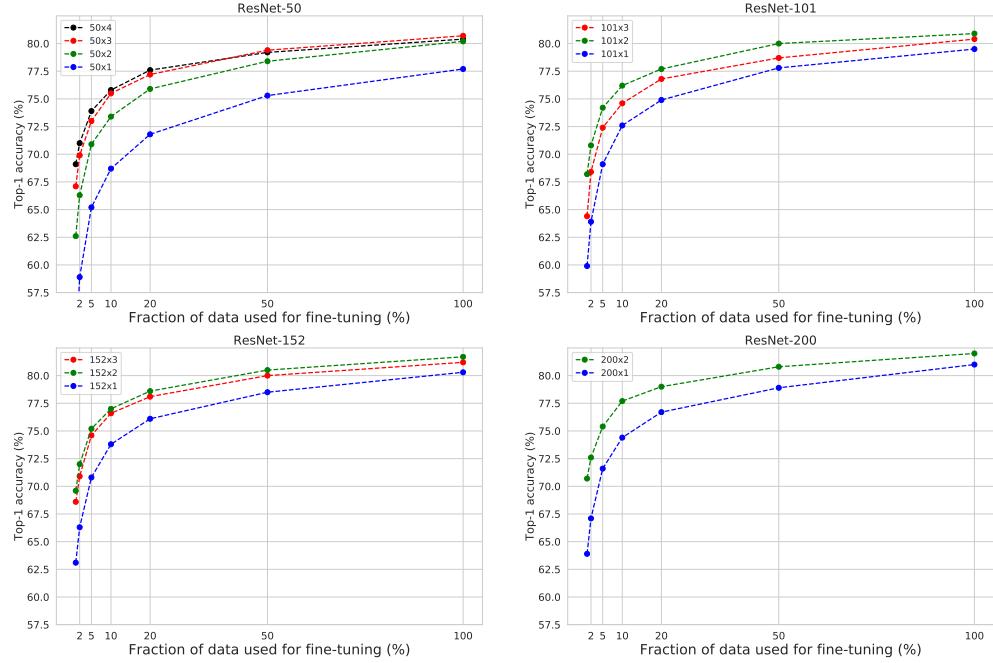


Figure 5: Semi-supervised training with a fraction of ImageNet labels on multiple ResNets architecture pretrained with BYOL. Note that large networks are facing overfitting problems.

C.2 Linear evaluation on larger architectures and supervised baselines

Here we investigate the performance of BYOL with deeper and wider ResNet architectures. We compare ourselves to the best supervised baselines from [8] when available (rightmost column in table 9), which are also presented in Figure 1. Importantly, we close in on those baselines using the ResNet-50 ($\times 2$) and the ResNet-50 ($\times 4$) architectures, where we are within 0.4 accuracy points of the supervised performance. To the best of our knowledge, this is the first time that the gap to supervised has been closed to such an extent using a self-supervised method under the linear evaluation protocol. Therefore, in order to ensure fair comparison, and suspecting that the supervised baselines’ performance in [8] could be even further improved with appropriate data augmentations, we also report on our own reproduction of strong supervised baselines. We use RandAugment [87] data augmentation for all large

Architecture	Multiplier	Weights	BYOL		Supervised (ours)		Supervised [8]
			Top-1	Top-5	Top-1	Top-5	Top-1
ResNet-50	1×	24M	74.3	91.6	76.4	92.9	76.5
ResNet-101	1×	43M	76.4	93.0	78.0	94.0	-
ResNet-152	1×	58M	77.3	93.7	79.1	94.5	-
ResNet-200	1×	63M	77.8	93.9	79.3	94.6	-
ResNet-50	2×	94M	77.4	93.6	79.9	95.0	77.8
ResNet-101	2×	170M	78.7	94.3	80.3	95.0	-
ResNet-50	3×	211M	78.2	93.9	80.2	95.0	-
ResNet-152	2×	232M	79.0	94.6	80.6	95.3	-
ResNet-200	2×	250M	79.6	94.9	80.1	95.2	-
ResNet-50	4×	375M	78.6	94.2	80.7	95.3	78.9
ResNet-101	3×	382M	78.4	94.2	80.7	95.3	-
ResNet-152	3×	522M	79.5	94.6	80.9	95.2	-

Table 9: **Linear evaluation of BYOL on ImageNet using larger encoders.** Top-1 and top-5 accuracies are reported in %.

ResNet architectures (which are all version 1, as per [22]). We train our supervised baselines for up to 200 epochs, using SGD with a Nesterov momentum value of 0.9, a cosine-annealed learning rate after a 5 epochs linear warmup period, weight decay with a value of $1e - 4$, and a label smoothing [88] value of 0.1. Results are presented in Figure 6.

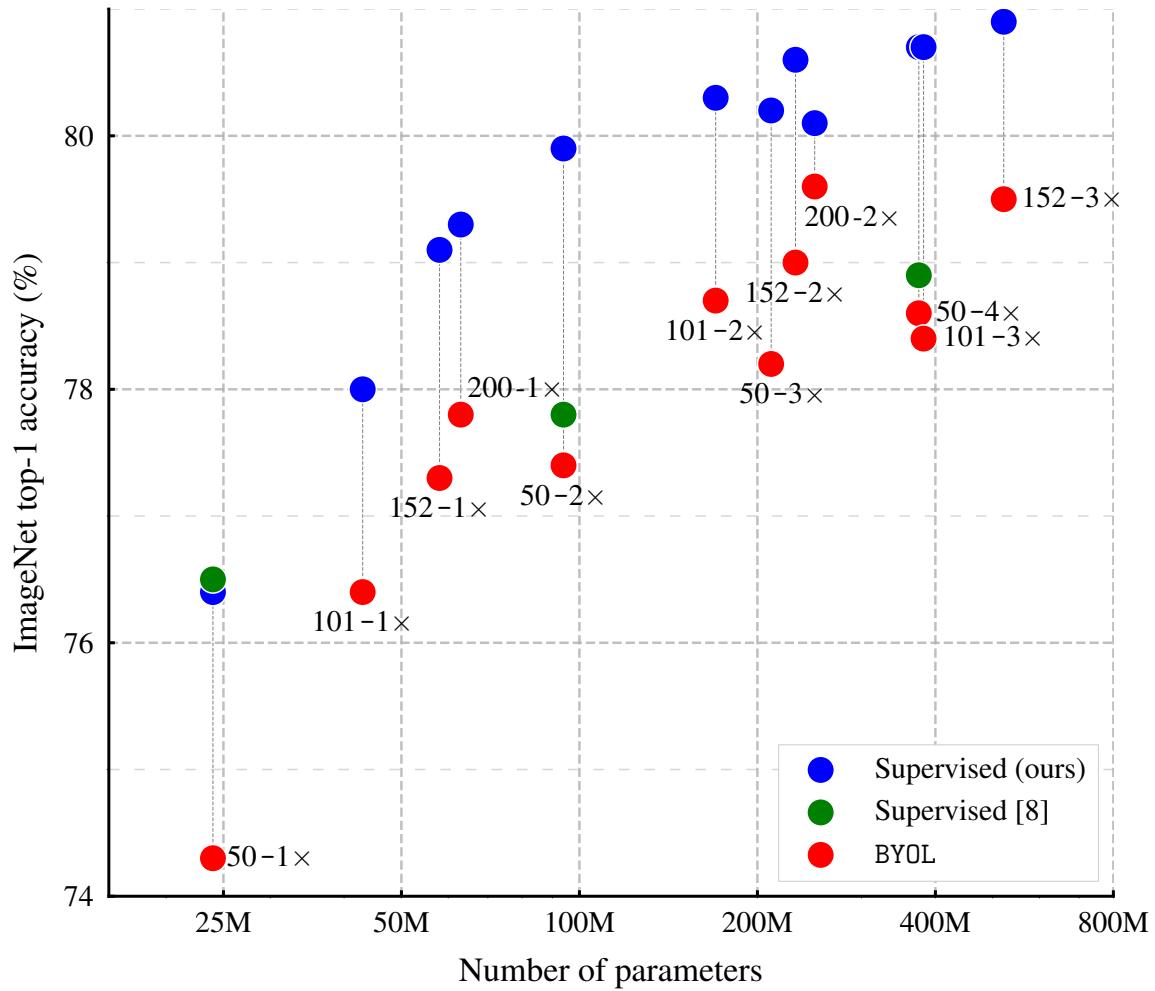


Figure 6: Results for linear evaluation of BYOL compared to fully supervised baselines with various ResNet architectures. Our supervised baselines are ran with RandAugment [87] augmentations.

D Transfer to other datasets

D.1 Datasets

Dataset	Classes	Original train examples	Train examples	Valid. examples	Test examples	Accuracy measure	Test provided
ImageNet [21]	1000	1281167	1271158	10009	50000	Top-1 accuracy	-
Food101 [89]	101	75750	68175	7575	25250	Top-1 accuracy	-
CIFAR-10 [78]	10	50000	45000	5000	10000	Top-1 accuracy	-
CIFAR-100 [78]	100	50000	44933	5067	10000	Top-1 accuracy	-
Birdsnap [90]	500	47386	42405	4981	2443	Top-1 accuracy	-
Sun397 (split 1) [79]	397	19850	15880	3970	19850	Top-1 accuracy	-
Cars [91]	196	8144	6494	1650	8041	Top-1 accuracy	-
Aircraft [92]	100	3334	3334	3333	3333	Mean per-class accuracy	Yes
PASCAL-VOC2007 [80]	20	5011	2501	2510	4952	11-point mAP / AP50	-
PASCAL-VOC2012 [80]	21	10582	-	2119	1449	Mean IoU	-
DTD (split 1) [81]	47	1880	1880	1880	1880	Top-1 accuracy	Yes
Pets [93]	37	3680	2940	740	3669	Mean per-class accuracy	-
Caltech-101 [94]	101	3060	2550	510	6084	Mean per-class accuracy	-
Places365 [73]	365	1803460	1803460	-	36500	Top-1 accuracy	-
Flowers [95]	102	1020	1020	1020	6149	Mean per-class accuracy	Yes

Table 10: Characteristics of image datasets used in transfer learning. When an official test split with labels is not publicly available, we use the official validation split as test set, and create a held-out validation set from the training examples.

We perform transfer via linear classification and fine-tuning on the same set of datasets as in [8], namely Food-101 dataset [89], CIFAR-10 [78] and CIFAR-100 [78], Birdsnap [90], the SUN397 scene dataset [79], Stanford Cars [91], FGVC Aircraft [92], the PASCAL VOC 2007 classification task [80], the Describable Textures Dataset (DTD) [81], Oxford-IIIT Pets [93], Caltech-101 [94], and Oxford 102 Flowers [95]. As in [8], we used the validation sets specified by the dataset creators to select hyperparameters for FGVC Aircraft, PASCAL VOC 2007, DTD, and Oxford 102 Flowers. On other datasets, we use the validation examples as test set, and hold out a subset of the training examples that we use as validation set. We use standard metrics for each datasets:

- **Top-1:** We compute the proportion of correctly classified examples.
- **Mean per class:** We compute the top-1 accuracy for each class separately and then compute the empirical mean over the classes.
- **Point 11-mAP:** We compute the empirical mean *average precision* as defined in [80].
- **Mean IoU:** We compute the empirical mean Intersection-Over-Union as defined in [80].
- **AP50:** We compute the Average Precision as defined in [80].

We detail the validation procedures for some specific datasets:

The Pascal visual object classes (VOC) challenge

- For Sun397 [79], the original dataset specifies 10 train/test splits, all of which contain 50 examples/images of 397 different classes. We use the first train/test split. The original dataset specifies no validation split and therefore, the training images have been further subdivided into 40 images per class for the train split and 10 images per class for the valid split.
- For Birdsnap [90], we use a random selection of valid images with the same number of images per category as the test split.
- For DTD [81], the original dataset specifies 10 train/validation/test splits, we only use the first split
- For Caltech-101 [94], the original dataset does not specifies any train/test splits. We have followed the approach used in [96]: This file defines datasets for 5 random splits of 25 training images per category, with 5 validation images per category and the remaining images used for testing.
- For ImageNet, we took the last 10009 last images of the official tensorflow ImageNet split.
- For Oxford-IIIT Pets, the valid set consists of 20 randomly selected images per class.

Information about the dataset are summarized in Table 10.

D.2 Transfer via linear classification

We follow the linear evaluation protocol of [48, 74, 8] that we detail next for completeness. We train a regularized multinomial logistic regression classifier on top of the frozen representation, i.e., with frozen pretrained parameters

and without re-computing batch-normalization statistics. In training and testing, we do not perform any image augmentations; images are resized to 224 pixels along the shorter side using bicubic resampling and then normalized with ImageNet statistics. Finally, we minimize the cross-entropy objective using LBFGS with ℓ_2 -regularization, where we select the regularization parameters from a range of 45 logarithmically-spaced values between 10^{-6} and 10^5 . After choosing the best-performing hyperparameters on the validation set, the model is retrained on combined training and validation images together, using the chosen parameters. The final accuracy is reported on the test set.

D.3 Transfer via fine-tuning

We follow the same fine-tuning protocol as in [32, 48, 76, 8] that we also detail for completeness. Specifically, we initialize the network with the parameters of the pretrained representation. At training time, we apply spatial transformation, i.e., random crops with resize to 224×224 pixels and random flips. At test time, images are resized to 256 pixels along the shorter side using bicubic resampling, after which a 224×224 center crop is extracted. In both cases, we normalize the color channels by subtracting the average color and dividing by the standard deviation (computed on ImageNet), after applying the augmentations. We optimize the loss using SGD with Nesterov momentum for 20000 steps with a batch size of 256 and with a momentum of 0.9. We set the momentum parameter for the batch normalization statistics to $\max(1 - 10/s, 0.9)$ where s is the number of steps per epoch. The learning rate and weight decay are selected respectively with a grid of seven logarithmically spaced learning rates between 0.0001 and 0.1, and 7 logarithmically-spaced values of weight decay between 10^{-6} and 10^{-3} , as well as no weight decay. These values of weight decay are divided by the learning rate. After choosing the best-performing hyperparameters on the validation set, the model is retrained on combined training and validation images together, using the chosen parameters. The final accuracy is reported on the test set.

D.4 Implementation details for semantic segmentation

We use the same fully-convolutional network (FCN)-based [7] architecture as [9]. The backbone consists of the convolutional layers in ResNet-50. The 3×3 convolutions in the conv5 blocks use dilation 2 and stride 1. This is followed by two extra 3×3 convolutions with 256 channels, each followed by batch normalization and ReLU activations, and a 1×1 convolution for per-pixel classification. The dilation is set to 6 in the two extra 3×3 convolutions. The total stride is 16 (FCN-16s [7]).

We train on the `train_aug2012` set and report results on `val2012`. Hyperparameters are selected on a 2119 images held-out validation set. We use a standard per-pixel softmax cross-entropy loss to train the FCN. Training is done with random scaling (by a ratio in [0.5, 2.0]), cropping, and horizontal flipping. The crop size is 513. Inference is performed on the [513, 513] central crop. For training we use a batch size of 16 and weight decay of 0.0001. We select the base learning rate by sweeping across 5 logarithmically spaced values between 10^{-3} and 10^{-1} . The learning rate is multiplied by 0.1 at the 70-th and 90-th percentile of training. We train for 30000 iterations, and average the results on 5 seeds.

D.5 Implementation details for object detection

For object detection, we follow prior work on Pascal detection transfer [40, 23] wherever possible. We use a Faster R-CNN [82] detector with a R50-C4 backbone with a frozen representation. The R50-C4 backbone ends with the conv4 stage of a ResNet-50, and the box prediction head consists of the conv5 stage (including global pooling). We preprocess the images by applying multi-scale augmentation (rescaling the image so its longest edge is between 480 and 1024 pixels) but no other augmentation. We use an asynchronous SGD optimizer with 9 workers and train for 1.5M steps. We used an initial learning rate of 10^{-3} , which is reduced to 10^{-4} at 1M steps and to 10^{-5} at 1.2M steps.

D.6 Implementation details for depth estimation

For depth estimation, we follow the same protocol as in [83], and report its core components for completeness. We use a standard ResNet-50 backbone and feed the conv5 features into 4 fast up-projection blocks with respective filter sizes 512, 256, 128, and 64. We use a reverse Huber loss function for training [83, 97].

The original NYU Depth v2 frames of size [640, 480] are down-sampled by a factor 0.5 and center-cropped to [304, 228] pixels. Input images are randomly horizontally flipped and the following color transformations are applied:

- *Grayscale* with an application probability of 0.3.

- Brightness with a maximum brightness difference of 0.1255.
- Saturation with a saturation factor randomly picked in the interval [0.5, 1.5].
- Hue with a hue adjustment factor randomly picked in the interval [-0.2, 0.2].

We train for 7500 steps with batch size 256, weight decay 0.0005, and learning rate 0.16 (scaled linearly from the setup of [83] to account for the bigger batch size).

D.7 Further comparisons on PASCAL and NYU v2 Depth

For completeness, Table 11 and 12 extends Table 4 with other published baselines which use comparable networks. We see that in almost all settings, BYOL outperforms these baselines, even when those baselines use more data or deeper models. One notable exception is RMS error for NYU Depth prediction, which is a metric that’s sensitive to outliers. The reason for this is unclear, but one possibility is that the network is producing higher-variance predictions due to being more confident about a test-set scene’s similarities with those in the training set.

Method	AP ₅₀	mIoU
Supervised-IN [9]	74.4	74.4
RelPos [23], by [40] [*]	66.8	-
Multi-task [40] [*]	70.5	-
LocalAgg [98]	69.1	-
MoCo [9]	74.9	72.5
MoCo + IG-1B [9]	75.6	73.6
CPC[32] ^{**}	76.6	-
SimCLR (repro)	75.2	75.2
BYOL (ours)	77.5	76.3

Table 11: Transfer results in semantic segmentation and object detection.

* uses a larger model (ResNet-101). ** uses an even larger model (ResNet-161).

Method	pct.<1.25	Higher better		Lower better	
		pct.<1.25 ²	pct.<1.25 ³	rms	rel
Supervised-IN [83]	81.1	95.3	98.8	0.573	0.127
RelPos [23], by [40] [*]	80.6	94.7	98.3	0.399	0.146
Color [41], by [40] [*]	76.8	93.5	97.7	0.444	0.164
Exemplar [46, 40] [*]	71.3	90.6	96.5	0.513	0.191
Mot. Seg. [99], by [40] [*]	74.2	92.4	97.4	0.473	0.177
Multi-task [40] [*]	79.3	94.2	98.1	0.422	0.152
SimCLR (repro)	83.3	96.5	99.1	0.557	0.134
BYOL (ours)	84.6	96.7	99.1	0.541	0.129

Table 12: Transfer results on NYU v2 depth estimation.

E Pretraining on Places 365

To ascertain that BYOL learns good representations on other datasets, we applied our representation learning protocol on the scene recognition dataset Places365-Standard [73] before performing linear evaluation. This dataset contains 1.80 million training images and 36500 validation images with labels, making it roughly similar to ImageNet in scale. We reuse the *exact* same parameters as in Section 4 and train the representation for 1000 epochs, using BYOL and our SimCLR reproduction. Results for the linear evaluation setup (using the protocol of Appendix C.1 for ImageNet and Places365, and that of Appendix D on other datasets) are reported in Table 13.

Interestingly, the representation trained by using BYOL on Places365 (BYOL-PL) consistently outperforms that of SimCLR on the same dataset, but underperforms the BYOL representation trained on ImageNet (BYOL-IN) on all tasks except Places365 and SUN397 [79], another scene understanding dataset. Interestingly, all three unsupervised representation learning methods achieve a relatively high performance on the Places365 task; for comparison,

confidence intervals are calculated as the half-difference between the maximum and minimum score obtained from multiple runs or experiments using different seeds. The purpose of calculating confidence intervals is to provide a range that reflects the uncertainty or variability in the results due to different random seeds. However, there is a condition mentioned: confidence intervals are omitted for half-differences that are lower than 0.25 accuracy points. This means that if the difference between the maximum and minimum scores is very small (less than 0.25 accuracy points), it is considered negligible, and confidence intervals are not reported or considered in such cases. In statistical terms, confidence intervals give you a range of values within which you can be reasonably confident that the true value lies. In this specific case, the focus is on the half-difference between maximum and minimum scores, and intervals are excluded when this difference is deemed too small to be meaningful (below 0.25 accuracy points).

reference [73] (in its linked repository) reports a top-1 accuracy of 55.2% for a ResNet-50v2 trained from scratch using labels on this dataset.

Method	Places365	ImageNet	Food101	CIFAR10	CIFAR100	Birdsnap	SUN397	Cars	Aircraft	DTD	Pets	Caltech-101	Flowers
BYOL-IN	51.0	74.3	75.3	91.3	78.4	57.3	62.6	67.2	60.6	76.5	90.4	94.3	96.1
BYOL-PL	53.2	58.5	64.7	84.5	66.1	28.8	64.2	55.6	55.9	68.5	66.1	84.3	90.0
SimCLR-PL	53.0	56.5	61.7	80.8	61.1	21.2	62.5	40.1	44.3	64.3	59.4	77.1	85.9

Table 13: Transfer learning results (linear evaluation, ResNet-50) from Places365 (PL). For comparison purposes, we also report the results from BYOL trained on ImageNet (BYOL-IN).

F Additional ablation results

To extend on the above results, we provide additional ablations obtained using the same experimental setup as in Section 5, i.e., 300 epochs, averaged over 3 seeds with the initial learning rate set to 0.3, the batch size to 4096, the weight decay to 10^{-6} and the base target decay rate τ_{base} to 0.99 unless specified otherwise. Confidence intervals correspond to the half-difference between the maximum and minimum score of these seeds; we omit them for half-differences lower than 0.25 accuracy points.

F.1 Architecture settings

Table 14 shows the influence of projector and predictor architecture on BYOL. We examine the effect of different depths for both the projector and predictor, as well as the effect of the projection size. We do not apply a ReLU activation nor a batch normalization on the final linear layer of our MLPs such that a depth of 1 corresponds to a linear layer. Using the default projector and predictor of depth 2 yields the best performance.

Proj. g_θ depth	Pred. q_θ depth	Top-1	Top-5	Projector g_θ output dim	Top-1	Top-5
1	1	61.9	86.0	16	69.9 \pm 0.3	89.9
	2	65.0	86.8			
	3	65.7	86.8			
2	1	71.5	90.7			
	2	72.5	90.8			
	3	71.4	90.4			
3	1	71.4	90.4			
	2	72.1	90.5			
	3	72.1	90.5			

(b) Projection dimension.

(a) Projector and predictor depth (i.e. the number of Linear layers).

Table 14: Effect of architectural settings where top-1 and top-5 accuracies are reported in %.

Table 15a shows the influence of the initial learning rate on BYOL. Note that the optimal value depends on the number of training epochs. Table 15b displays the influence of the weight decay on BYOL.

F.2 Batch size

We run a sweep over the batch size for both BYOL and our reproduction of SimCLR. As explained in Section 5, when reducing the batch size by a factor N , we average gradients over N consecutive steps and update the target network once every N steps. We report in Table 16, the performance of both our reproduction of SimCLR and BYOL for batch sizes between 4096 (BYOL and SimCLR default) down to 64. We observe that the performance of SimCLR deteriorates faster than the one of BYOL which stays mostly constant for batch sizes larger than 256. We believe that the performance at batch size 256 could match the performance of the large 4096 batch size with proper parameter tuning when accumulating the gradient. We think that the drop in performance at batch size 64 in table 16 is mainly related to the ill behaviour of batch normalization at low batch sizes [100].

Learning rate	Top-1	Top-5
0.01	34.8 \pm 3.0	60.8 \pm 3.2
0.1	65.0	87.0
0.2	71.7	90.6
0.3	72.5	90.8
0.4	72.3	90.6
0.5	71.5	90.1
1	69.4	89.2

(a) Base learning rate.

Weight decay coefficient	Top-1	Top-5
$1 \cdot 10^{-7}$	72.1	90.4
$5 \cdot 10^{-7}$	72.6	91.0
$1 \cdot 10^{-6}$	72.5	90.8
$5 \cdot 10^{-6}$	71.0 \pm 0.3	90.0
$1 \cdot 10^{-5}$	69.6 \pm 0.4	89.3

(b) Weight decay.

Table 15: Effect of learning rate and weight decay. We note that BYOL’s performance is quite robust within a range of hyperparameters. We also observe that setting the weight decay to zero may lead to unstable results (as in SimCLR).

Batch size	Top-1		Top-5	
	BYOL (ours)	SimCLR (repro)	BYOL (ours)	SimCLR (repro)
4096	72.5	67.9	90.8	88.5
2048	72.4	67.8	90.7	88.5
1024	72.2	67.4	90.7	88.1
512	72.2	66.5	90.8	87.6
256	71.8	64.3 \pm 2.1	90.7	86.3 \pm 1.0
128	69.6 \pm 0.5	63.6	89.6	85.9
64	59.7 \pm 1.5	59.2 \pm 2.9	83.2 \pm 1.2	83.0 \pm 1.9

Table 16: Influence of the batch size.

F.3 Image augmentations

Table 17 compares the impact of individual image transformations on BYOL and SimCLR. BYOL is more resilient to changes of image augmentations across the board. For completeness, we also include an ablation with symmetric parameters across both views; for this ablation, we use a Gaussian blurring w.p. of 0.5 and a solarization w.p. of 0.2 for both \mathcal{T} and \mathcal{T}' , and recover very similar results compared to our baseline choice of parameters.

Image augmentation	Top-1		Top-5	
	BYOL (ours)	SimCLR (repro)	BYOL (ours)	SimCLR (repro)
Baseline	72.5	67.9	90.8	88.5
Remove flip	71.9	67.3	90.6	88.2
Remove blur	71.2	65.2	90.3	86.6
Remove color (jittering and grayscale)	63.4 \pm 0.7	45.7	85.3 \pm 0.5	70.6
Remove color jittering	71.8	63.7	90.7	85.9
Remove grayscale	70.3	61.9	89.8	84.1
Remove blur in \mathcal{T}'	72.4	67.5	90.8	88.4
Remove solarize in \mathcal{T}'	72.3	67.7	90.8	88.2
Remove blur and solarize in \mathcal{T}'	72.2	67.4	90.8	88.1
Symmetric blurring/solarization	72.5	68.1	90.8	88.4
Crop only	59.4 \pm 0.3	40.3 \pm 0.3	82.4	64.8 \pm 0.4
Crop and flip only	60.1 \pm 0.3	40.2	83.0 \pm 0.3	64.8
Crop and color only	70.7	64.2	90.0	86.2
Crop and blur only	61.1 \pm 0.3	41.7	83.9	66.4

Table 17: Ablation on image transformations.

Loss weight β	Temperature α	Top-1	Top-5
0	0.1	72.5	90.8
0.1	0.01	72.2	90.7
	0.1	72.4	90.9
	0.3	72.7	91.0
	1	72.6	90.9
	3	72.5	90.9
	10	72.5	90.9
0.5	0.01	70.9	90.2
	0.1	72.0	90.8
	0.3	72.7	91.2
	1	72.7	91.1
	3	72.6	91.1
	10	72.5	91.0
1	0.01	53.9 ± 0.5	77.5 ± 0.5
	0.1	70.9	90.3
	0.3	72.7	91.1
	1	72.7	91.1
	3	72.6	91.0
	10	72.6	91.1

Table 18: Top-1 accuracy in % under linear evaluation protocol at 300 epochs of sweep over the temperature α and the dispersion term weight β when using a predictor and a target network.

F.4 Details on the relation to contrastive methods

As mentioned in Section 5, the BYOL loss Eq. 2 can be derived from the InfoNCE loss

$$\text{InfoNCE}_{\theta}^{\alpha, \beta} \triangleq \frac{2}{B} \sum_{i=1}^B S_{\theta}(v_i, v'_i) - \frac{2\alpha \cdot \beta}{B} \sum_{i=1}^B \ln \left(\sum_{j \neq i} \exp \frac{S_{\theta}(v_i, v_j)}{\alpha} + \sum_j \exp \frac{S_{\theta}(v_i, v'_j)}{\alpha} \right), \quad (8)$$

with

$$S_{\theta}(u_1, u_2) \triangleq \frac{\langle \phi(u_1), \psi(u_2) \rangle}{\|\phi(u_1)\|_2 \cdot \|\psi(u_2)\|_2}. \quad (9)$$

The InfoNCE loss, introduced in [10], can be found in factored form in [84] as

$$\text{InfoNCE}_{\theta} \triangleq \frac{1}{B} \sum_{i=1}^B \ln \frac{f(v_i, v'_i)}{\frac{1}{B} \sum_j \exp f(v_i, v'_j)}. \quad (10)$$

As in SimCLR [8] we also use negative examples given by $(v_i, v_j)_{j \neq i}$ to get

$$\frac{1}{B} \sum_{i=1}^B \ln \frac{\exp f(v_i, v'_i)}{\frac{1}{B} \sum_{j \neq i} \exp f(v_i, v_j) + \frac{1}{B} \sum_j \exp f(v_i, v'_j)} \quad (11)$$

$$= \ln B + \frac{1}{B} \sum_{i=1}^B f(v_i, v'_i) - \frac{1}{B} \sum_{i=1}^B \ln \left(\sum_{j \neq i} \exp f(v_i, v_j) + \sum_j \exp f(v_i, v'_j) \right). \quad (12)$$

To obtain Eq. 8 from Eq. 12, we subtract $\ln B$ (which is independent of θ), multiply by 2α , take $f(x, y) = S_{\theta}(x, y)/\alpha$ and finally multiply the second (negative examples) term by β . Using $\beta = 1$ and dividing by 2α gets us back to the usual InfoNCE loss as used by SimCLR.

In our ablation in Table 5b, we set the temperature α to its best value in the SimCLR setting (i.e., $\alpha = 0.1$). With this value, setting β to 1 (which adds negative examples), in the BYOL setting (i.e., with both a predictor and a target network) hurts the performances. In Table 18, we report results of a sweep over both the temperature α and the

Method	Predictor	Target parameters	β	Top-1
BYOL	✓	ξ	0	72.5
	✓	ξ	1	70.9
	✓	ξ	1	70.7
	✓	$\text{sg}(\theta)$	1	70.2
	✓	θ	1	69.4
	✓	$\text{sg}(\theta)$	1	70.1
	✓	$\text{sg}(\theta)$	1	69.2
	✓	θ	1	69.0
	✓	$\text{sg}(\theta)$	0	5.5
	✓	θ	0	0.3
SimCLR		ξ	0	0.2
		$\text{sg}(\theta)$	0	0.1
		θ	0	0.1

Table 19: Top-1 accuracy in %, under linear evaluation protocol at 300 epochs, of intermediate variants between BYOL and SimCLR (with caveats discussed in Appendix F.5). sg means stop gradient.

weight parameter β with a predictor and a target network where BYOL corresponds to $\beta = 0$. No run significantly outperforms BYOL and some values of α and β hurt the performance. While the best temperature for SimCLR (without the target network and a predictor) is 0.1, after adding a predictor and a target network the best temperature α is higher than 0.3. ↗ Joint Loss

Using a target network in the loss has two effects: stopping the gradient through the prediction targets and stabilizing the targets with averaging. Stopping the gradient through the target change the objective while averaging makes the target stable and stale. In Table 5b we only shows results of the ablation when either using the online network as the prediction target (and flowing the gradient through it) or with a target network (both stopping the gradient into the prediction targets and computing the prediction targets with a moving average of the online network). We shown in Table 5b that using a target network is beneficial but it has two distinct effects we would like to understand from which effect the improvement comes from. We report in Table 19 the results already in Table 5b but also when the prediction target is computed with a stop gradient of the online network (the gradient does not flow into the prediction targets). This shows that making the prediction targets stable and stale is the main cause of the improvement rather than the change in the objective due to the stop gradient.

F.5 SimCLR baseline of Section 5

The SimCLR baseline in Section 5 ($\beta = 1$, without predictor nor target network) is slightly different from the original one in [8]. First we multiply the original loss by 2α . For comparison here is the original SimCLR loss,

$$\text{InfoNCE}_\theta \triangleq \frac{1}{B} \sum_{i=1}^B \frac{S_\theta(v_i, v'_i)}{\alpha} - \frac{1}{B} \sum_{i=1}^B \ln \left(\sum_{j \neq i} \exp \frac{S_\theta(v_i, v_j)}{\alpha} + \sum_j \exp \frac{S_\theta(v_i, v'_j)}{\alpha} \right). \quad (13)$$

Note that this multiplication by 2α matters as the LARS optimizer is not completely invariant with respect to the scale of the loss. Indeed, LARS applies a preconditioning to gradient updates on all weights, except for biases and batch normalization parameters. Updates on preconditioned weights are invariant by multiplicative scaling of the loss. However, the bias and batch normalization parameter updates remain sensitive to multiplicative scaling of the loss.

We also increase the original SimCLR hidden and output size of the projector to respectively 4096 and 256. In our reproduction of SimCLR, these three combined changes improves the top-1 accuracy at 300 epochs from 67.9% (without the changes) to 69.2% (with the changes).

F.6 Ablation on the normalization in the loss function

BYOL minimizes a squared error between the ℓ_2 -normalized prediction and target. We report results of BYOL at 300 epochs using different normalization function and no normalization at all. More precisely, given batch of prediction and targets in \mathbb{R}^d , $(p_i, t_i)_{i \leq B}$ with B the batch size, BYOL uses the loss function $\frac{1}{B} \sum_{i=1}^B \|n_{\ell_2}(p_i) - n_{\ell_2}(z_i)\|_2^2$ with $n_{\ell_2} : x \rightarrow x / \|x\|_2$. We run BYOL with other normalization functions: non-trainable batch-normalization

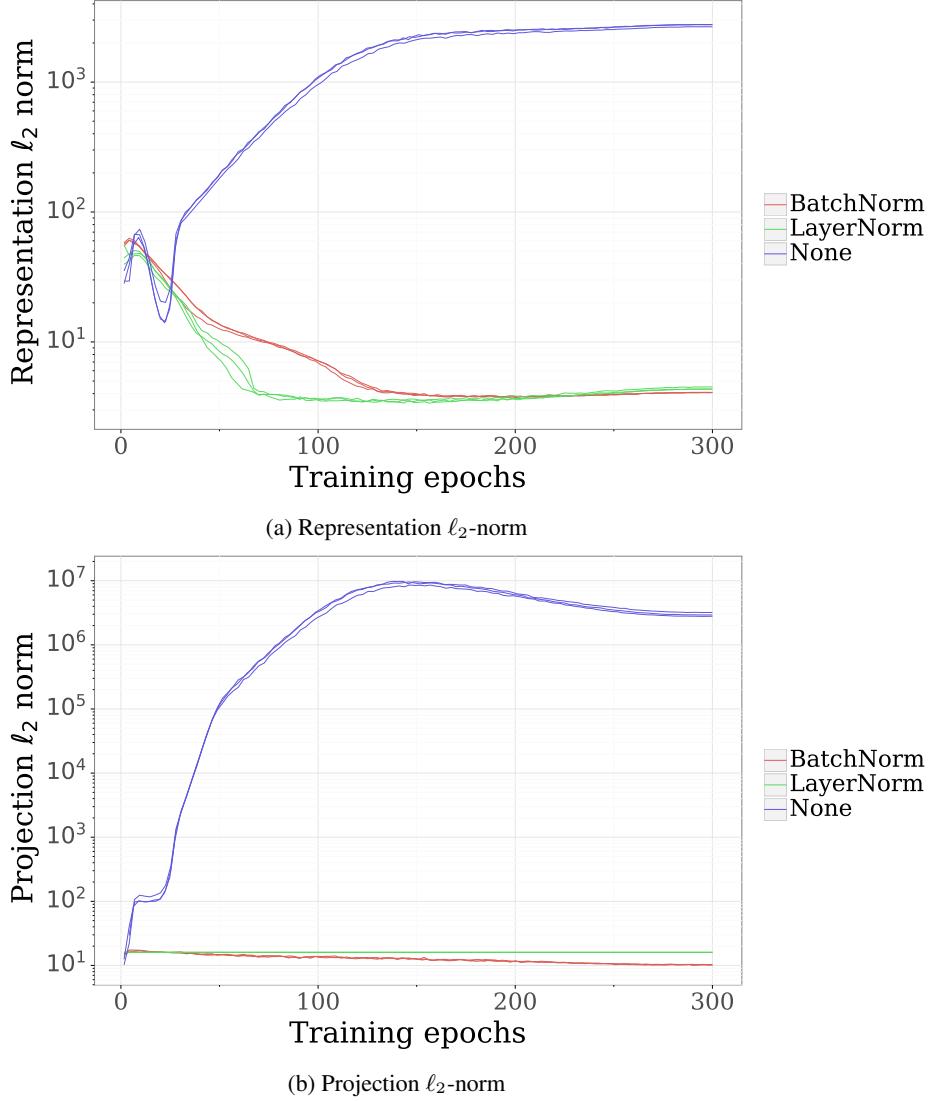


Figure 7: Effect of normalization on the ℓ_2 norm of network outputs.

Normalization	Top-1	Top-5
ℓ_2 -norm	72.5	90.8
LAYERNORM	72.5 ± 0.4	90.1
No normalization	67.4	87.1
BATCHNORM	65.3	85.3

Table 20: Top-1 accuracy in % under linear evaluation protocol at 300 epochs for different normalizations in the loss.

and layer-normalization and no normalization. We divide the batch normalization and layer normalization by \sqrt{d} to have a consistent scale with the ℓ_2 -normalization. We report results in Table 20 where ℓ_2 , LAYERNORM, no normalization and BATCHNORM respectively denote using n_{ℓ_2} , n_{BN} , n_{LN} and n_{IB} with

$$\begin{aligned}
n_{\text{BN}i}^j : x \rightarrow \frac{x_i^j - \mu_{\text{BN}}^j(x)}{\sigma_{\text{BN}}^j(x) \cdot \sqrt{d}}, \quad n_{\text{LN}i}^j : x \rightarrow \frac{x_i^j - \mu_{\text{LN}i}(x)}{\sigma_{\text{LN}i}(x) \cdot \sqrt{d}}, \quad n_{\text{ID}} : x \rightarrow x, \\
\mu_{\text{BN}}^j : x \rightarrow \frac{1}{B} \sum_{i=1}^B x_i^j, \quad \sigma_{\text{BN}}^j : x \rightarrow \sqrt{\frac{1}{B} \sum_{i=1}^B (x_i^j)^2 - \mu_{\text{BN}}^j(x)^2}, \\
\mu_{\text{LN}i} : x \rightarrow \frac{1}{d} \sum_{j=1}^d x_i^j, \quad \sigma_{\text{LN}i} : x \rightarrow \frac{\|x_i - \mu_{\text{LN}i}(x)\|_2}{\sqrt{d}}
\end{aligned}$$

When using no normalization at all, the projection ℓ_2 norm rapidly increases during the first 100 epochs and stabilizes at around $3 \cdot 10^6$ as shown in Figure 7. Despite this behaviour, using no normalization still performs reasonably well (67.4%). The ℓ_2 normalization performs the best.

G Training with smaller batch sizes

The results described in Section 4 were obtained using a batch size of 4096 split over 512 TPU cores. Due to its increased robustness, BYOL can also be trained using smaller batch sizes without significantly decreasing performance. Using the same linear evaluation setup, BYOL achieves 73.7% top-1 accuracy when trained over 1000 epochs with a batch size of 512 split over 64 TPU cores (approximately 4 days of training). For this setup, we reuse the same setting as in Section 3, but use a base learning rate of 0.4 (appropriately scaled by the batch size) and $\tau_{\text{base}} = 0.9995$ with the same weight decay coefficient of $1.5 \cdot 10^{-6}$.

H Details on Equation 5 in Section 3.2

In this section we clarify why BYOL’s update is related to Eq. 5 from Section 3.2,

$$\nabla_{\theta} \mathbb{E} \left[\|q^*(z_{\theta}) - z'_{\xi}\|_2^2 \right] = \nabla_{\theta} \mathbb{E} \left[\|\mathbb{E}[z'_{\xi}|z_{\theta}] - z'_{\xi}\|_2^2 \right] = \nabla_{\theta} \mathbb{E} \left[\sum_i \text{Var}(z'_{\xi,i}|z_{\theta}) \right]. \quad (5)$$

Recall that q^* is defined as

$$q^* \triangleq \arg \min_q \mathbb{E} \left[\|q(z_{\theta}) - z'_{\xi}\|_2^2 \right], \quad \text{where } q^*(z_{\theta}) = \mathbb{E}[z'_{\xi}|z_{\theta}], \quad (4)$$

and implicitly depends on θ and ξ ; therefore, it should be denoted as $q^*(\theta, \xi)$ instead of just q^* . For simplicity we write $q^*(\theta, \xi)(z_{\theta})$ as $q^*(\theta, \xi, z_{\theta})$ the output of the optimal predictor for any parameters θ and ξ and input z_{θ} .

BYOL updates its online parameters following the gradient of Eq. 5, but considering only the gradients of q with respect to its third argument z when applying the chain rule. If we rewrite

$$\mathbb{E} \left[\|q^*(\theta, \xi, z_{\theta}) - z'_{\xi}\|_2^2 \right] = \mathbb{E} [L(q^*(\theta, \xi, z_{\theta}), z'_{\xi})], \quad (14)$$

the gradient of this quantity w.r.t. θ is

$$\frac{\partial}{\partial \theta} \mathbb{E} [L(q^*(\theta, \xi, z_{\theta}), z'_{\xi})] = \mathbb{E} \left[\frac{\partial L}{\partial q} \cdot \frac{\partial q^*}{\partial \theta} + \frac{\partial L}{\partial q} \cdot \frac{\partial q^*}{\partial z} \cdot \frac{\partial z_{\theta}}{\partial \theta} \right], \quad (15)$$

where $\frac{\partial q^*}{\partial \theta}$ and $\frac{\partial q^*}{\partial z}$ are the gradients of q^* with respect to its first and last argument. Using the envelope theorem, and thanks to the optimality condition of the predictor, the term $\mathbb{E} \left[\frac{\partial L}{\partial q} \cdot \frac{\partial q^*}{\partial \theta} \right] = 0$. Therefore, the remaining term $\mathbb{E} \left[\frac{\partial L}{\partial q} \cdot \frac{\partial q^*}{\partial z} \cdot \frac{\partial z_{\theta}}{\partial \theta} \right]$ where gradients are only back-propagated through the predictor’s input is exactly the direction followed by BYOL.

I Importance of a near-optimal predictor

In this part we build upon the intuitions of Section 3.2 on the importance of keeping the predictor near-optimal. Specifically, we show that it is possible to remove the exponential moving average in BYOL’s target network (*i.e.*, simply copy weights of the online network into the target) without causing the representation to collapse, provided the predictor remains sufficiently good.

I.1 Predictor learning rate

In this setup, we remove the exponential moving average (*i.e.*, set $\tau = 0$ over the full training in Eq. 1), and multiply the learning rate of the predictor by a constant λ compared to the learning rate used for the rest of the network; all other hyperparameters are unchanged. As shown in Table 21, using sufficiently large values of λ provides a reasonably good level of performance and the performance sharply decreases with λ to 0.01% top-1 accuracy (no better than random) for $\lambda = 0$.

To show that this effect is directly related to a change of behavior in the predictor, and not only to a change of learning rate in any subpart of the network, we perform a similar experiment by using a multiplier λ on the predictor's learning rate, and a different multiplier μ for the projector. In Table 22, we show that the representation typically collapses or performs poorly when the predictor learning rate is lower or equal to that of the projector. As mentioned in Section 3.2, we further hypothesize that one of the contributions of the target network is to maintain a near optimal predictor at all times.

I.2 Optimal linear predictor in closed form

Similarly, we can get rid of the slowly moving target network if we use a closed form optimal predictor on the batch, instead of a learned, non-optimal one. In this case we restrict ourselves to a linear predictor,

$$q^* = \arg \min_Q \|Z_\theta Q - Z'_\xi\|_2^2 = (Z_\theta^\top Z_\theta)^{-1} Z_\theta^\top Z'_\xi$$

with $\|\cdot\|_2$ being the Frobenius norm, Z_θ and Z'_ξ of shape (B, F) respectively the online and target projections, where B is the batch size and F the number of features; and q^* of shape (F, F) , the optimal linear predictor for the batch.

At 300 epochs, when using the closed form optimal predictor, and directly hard copying the weights of the online network into the target, we obtain a top-1 accuracy of *fill*.

λ	Top-1
0	0.01
1	5.5
2	62.8 \pm 1.5
10	66.6
20	66.3 \pm 0.3
Baseline	72.5

Table 21: Top-1 accuracy at 300 epochs when removing the slowly moving target network, directly hard copying the weights of the online network into the target network, and applying a multiplier to the predictor learning rate.

		λ_{pred}				
		1	1.5	2	5	10
μ_{proj}	1	3.2 \pm 2.9	25.7 \pm 6.6	60.8 \pm 2.9	66.7 \pm 0.4	66.9
	1.5	1.4 \pm 2.0	9.2 \pm 7.0	55.2 \pm 5.8	61.5 \pm 0.6	66.0 \pm 0.3
	2	2.0 \pm 2.8	5.3 \pm 1.9	15.8 \pm 13.4	60.9 \pm 0.8	66.3
	5	1.5 \pm 0.9	2.5 \pm 1.5	2.5 \pm 1.4	20.5 \pm 2.0	60.5 \pm 0.6
	10	0.1	2.1 \pm 0.3	1.9 \pm 0.8	2.8 \pm 0.4	8.3 \pm 6.8

Table 22: Top-1 accuracy at 300 epochs when removing the slowly moving target network, directly hard copying the weights of the online network in the target network, and applying a multiplier μ to the projector and λ to the predictor learning rate. The predictor learning rate needs to be higher than the projector learning rate in order to successfully remove the target network. This further suggests that the learning dynamic of predictor is central to BYOL's stability.

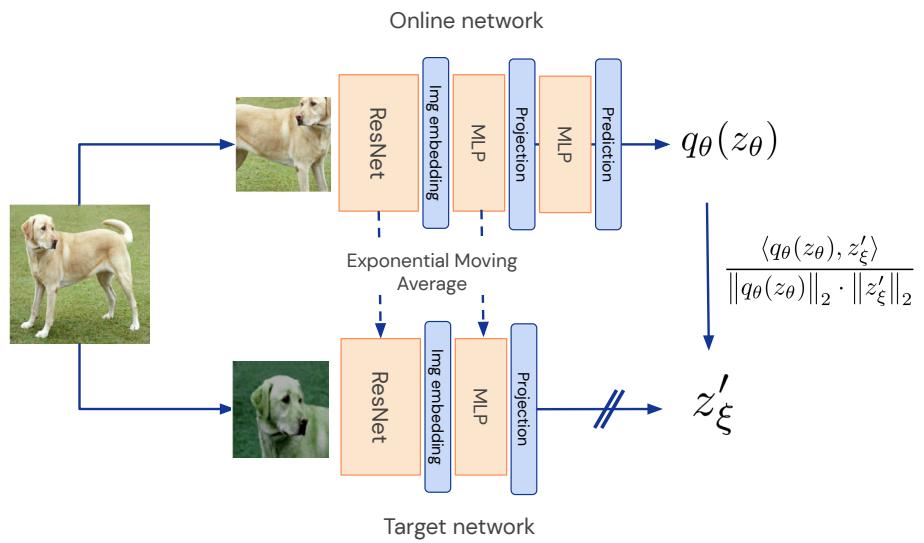


Figure 8: BYOL sketch summarizing the method by emphasizing the neural architecture.

J BYOL pseudo-code in JAX

J.1 Hyper-parameters

```
HPS = dict(
    max_steps=int(1000. * 1281167 / 4096), # 1000 epochs
    batch_size=4096,
    mlp_hidden_size=4096,
    projection_size=256,
    base_target_ema=4e-3,
    optimizer_config=dict(
        optimizer_name='lars',
        beta=0.9,
        trust_coeff=1e-3,
        weight_decay=1.5e-6,
        # As in SimCLR and official implementation of LARS, we exclude bias
        # and batchnorm weight from the Lars adaptation and weightdecay.
        exclude_bias_from_adaption=True),
    learning_rate_schedule=dict(
        # The learning rate is linearly increase up to
        # its base value * batchsize / 256 after warmup_steps
        # global steps and then anneal with a cosine schedule.
        base_learning_rate=0.2,
        warmup_steps=int(10. * 1281167 / 4096),
        anneal_schedule='cosine'),
    batchnorm_kwarg=dict(
        decay_rate=0.9,
        eps=1e-5),
    seed=1337,
)
```

J.2 Network definition

```
def network(inputs):
    """Build the encoder, projector and predictor."""
    embedding = ResNet(name='encoder', configuration='ResNetV1_50x1')(inputs)
    proj_out = MLP(name='projector')(embedding)
    pred_out = MLP(name='predictor')(proj_out)
    return dict(projection=proj_out, prediction=pred_out)

class MLP(hk.Module):
    """Multi Layer Perceptron, with normalization."""

    def __init__(self, name):
        super().__init__(name=name)

    def __call__(self, inputs):
        out = hk.Linear(output_size=HPS['mlp_hidden_size'])(inputs)
        out = hk.BatchNorm(**HPS['batchnorm_kwarg'])(out)
        out = jax.nn.relu(out)
        out = hk.Linear(output_size=HPS['projection_size'])(out)
        return out

# For simplicity, we omit BatchNorm related states.
# In the actual code, we use hk.transform_with_state. The corresponding
# net_init function outputs both a params and a state variable,
# with state containing the moving averages computed by BatchNorm
net_init, net_apply = hk.transform(network)
```

J.3 Loss function

```
def loss_fn(online_params, target_params, image_1, image_2):
    """Compute BYOL's loss function.

Args:
    online_params: parameters of the online network (the loss is later
        differentiated with respect to the online parameters).
    target_params: parameters of the target network.
    image_1: first transformation of the input image.
    image_2: second transformation of the input image.

Returns:
    BYOL's loss function.
"""

online_network_out_1 = net_apply(params=online_params, inputs=image_1)
online_network_out_2 = net_apply(params=online_params, inputs=image_2)
target_network_out_1 = net_apply(params=target_params, inputs=image_1)
target_network_out_2 = net_apply(params=target_params, inputs=image_2)

def regression_loss(x, y):
    norm_x, norm_y = jnp.linalg.norm(x, axis=-1), jnp.linalg.norm(y, axis=-1)
    return -2. * jnp.mean(jnp.sum(x * y, axis=-1) / (norm_x * norm_y))

# The stop_gradient is not necessary as we explicitly take the gradient with
# respect to online parameters only. We leave it to indicate that gradients
# are not backpropagated through the target network.
loss = regression_loss(online_network_out_1['prediction'],
                       jax.lax.stop_gradient(target_network_out_2['projection']))
loss += regression_loss(online_network_out_2['prediction'],
                       jax.lax.stop_gradient(target_network_out_1['projection']))
return loss
```

J.4 Training loop

```
def main(dataset):
    """Main training loop."""

    rng = jax.random.PRNGKey(HPS['seed'])
    rng, rng_init = jax.random.split(rng, num=2)
    dataset = dataset.batch(HPS['batch_size'])
    dummy_input = dataset.next()
    byol_state = init(rng_init, dummy_input)

    for global_step in range(HPS['max_steps']):
        inputs = dataset.next()

        rng, rng1, rng2 = jax.random.split(rng, num=3)
        image_1 = simclr_augmentations(inputs, rng1, image_number=1)
        image_2 = simclr_augmentations(inputs, rng2, image_number=2)
        byol_state = update_fn(
            **byol_state,
            global_step=global_step,
            image_1=image_1,
            image_2=image_2)

    return byol_state['online_params']
```

J.5 Update function

```
optimizer = Optimizer(**HPS['optimizer_config'])

def update_fn(online_params, target_params, opt_state, global_step, image_1,
              image_2):
    """Update online and target parameters.

Args:
    online_params: parameters of the online network (the loss is differentiated
                   with respect to the online parameters only).
    target_params: parameters of the target network.
    opt_state: state of the optimizer.
    global_step: current training step.
    image_1: first transformation of the input image.
    image_2: second transformation of the input image.

Returns:
    Dict containing updated online parameters, target parameters and
    optimization state.
"""

# update online network
grad_fn = jax.grad(loss_fn, argnums=0)
grads = grad_fn(online_params, target_params, image_1, image_2)
lr = learning_rate(global_step, **HPS['learning_rate_schedule'])
updates, opt_state = optimizer(lr).apply(grads, opt_state, online_params)
online_params = optix.apply_updates(online_params, updates)

# update target network
tau = target_ema(global_step, base_ema=HPS['base_target_ema'])
target_params = jax.tree_multimap(lambda x, y: x + (1 - tau) * (y - x),
                                  target_params, online_params)
return dict(
    online_params=online_params,
    target_params=target_params,
    opt_state=opt_state)

def init(rng, dummy_input):
    """BYOL's state initialization.

Args:
    rng: random number generator used to initialize parameters.
    dummy_input: a dummy image, used to compute intermediate outputs shapes.

Returns:
    Dict containing initial online parameters, target parameters and
    optimization state.
"""

online_params = net_init(rng, dummy_input)
target_params = net_init(rng, dummy_input)
opt_state = optimizer(0).init(online_params)
return dict(
    online_params=online_params,
    target_params=target_params,
    opt_state=opt_state)
```