
[Re] Understanding Self-Supervised Learning Dynamics without Contrastive Pairs

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 **Scope of Reproducibility**

3 The authors in [1] claim that with the underlying learning dynamics of *BYOL* [2] and *SimSiam* [3], a new method
4 *DirectPred* can be derived. We investigate the assumptions made for this derivation and also compare the quality of the
5 produced encoder representations through linear probing of these networks.

6 **Methodology**

7 We reimplemented *BYOL*, *SimSiam* and *DirectPred* from scratch as well as their ablations in TensorFlow. We checked
8 original repository in written PyTorch for some implementation details. In all experiments we used the CIFAR-10 train
9 set for training and the test set for evaluation. We were running our experiments for more than 100 hours on GCP's
10 V100 GPU.

11 **Results**

12 We show that the theoretical assumption regarding eigenspace alignment and symmetry hold also for a different dataset
13 other than the one used in the original paper. In addition, we reproduce ablations regarding learning rate, weight decay
14 and Exponential Moving Average.

15 Since we used CIFAR-10 in all experiments we can not directly compare accuracies. However, we show the same
16 relative behaviour of different networks given hyperparameter changes. We can directly compare performance for
17 one of the experiments (Table 8. in [1] bottom left part). Our models, namely *SGD Baseline*, *DirectPred* (with
18 and without frequency=5), achieve comparable accuracy which differ by at most 1%. We also confirm the claim
19 that *DirectPred* outperforms its one-layer SGD alternative. Our code can be accessed under the following link:
20 <https://anonymous.4open.science/r/SelfSupervisedLearning-FD0F>.

21 **What was easy**

22 The architecture of the Siamese network and training schemes were both straightforward to implement and easy to
23 understand.

24 **What was difficult**

25 We could not run our code on STL-10 dataset due to time and resource constraints. Due to differences between PyTorch
26 and TensorFlow libraries, we had to implement some parts by hand to keep our code as close to the original work as
27 possible. Also, original repository is not easy to read and does not cover all the experiments (e.g. eigenspace alignment
28 experiment). Correctly applying data-augmentation was also a hard task due to assumptions of how the individual data
29 augmentations functions actually work.

30 **Communication with original authors**

31 We did not contact authors of the paper since we did not encounter any major issues during the reproducibility study.

32 1 Introduction

33 Self-Supervised learning has become an important task in many domains, since labeled data is often rare and expensive
34 to get. Many modern methods of Self-Supervised learning are based on Siamese-networks [4] which are weight sharing
35 Neural networks for two or more inputs which representations then will be compared in latent space. The representation
36 created by this approach can then be used for classification by fine-tuning on fewer labelled data-points. Traditionally,
37 during pre-training positive pairs (same image, or two images from the same class) and negative pairs (different images
38 or two images from a different class) are used. The distance of the representation of positive pairs is minimized while
39 the distance of the representation of negative pairs is maximized, which prevents the networks from collapse (i.e
40 mapping all inputs to the same representation). These methods have shown quite some success in the past [5], [6], [7],
41 [8]. However, these methods rely on negative pairs, and large batch sizes which makes the training less feasible.

42 Recently, new methods have been proposed which rely only on positive pairs and yet don't collapse [2], [3]. In the
43 paper "Understanding Self-Supervised Learning Dynamics without Contrastive Pairs" by Tian et.al. [1] the underlying
44 dynamics are explored and based on the theoretical results, a new method, *DirectPred*, was proposed which does not
45 need an update of the predictor via gradient descent but instead is set directly each iteration.

46 The focus of this work is to test several assumptions made in [1] for the theoretical analysis and see if they hold. For
47 this, we will concentrate especially on the eigenvalues of the predictor network and the eigenspace alignment with its
48 input. Also, we will reproduce the results from [1], [2] and [3] on CIFAR-10 to compare their learned representation
49 using linear probing.

50 2 Related work

51 A common approach to representation learning without Siamese networks is generative modelling. Typically these
52 methods model a distribution over the data and a latent space, from which then embeddings can be drawn as data
53 representations. Usually these approaches rely on Auto-encoding [9, 10] or Adversarial networks [11, 12]. However,
54 generative models are often computationally heavy and hard to train.

55 Discriminative methods using Siamese networks like SimCLR [5, 6] and Moco [7] outperform generative models and
56 have lower computational cost. However, these methods rely on very large batch sizes since they use contrastive pairs.
57 Most recent methods, replicated in this work, like *BYOL* [2] and *SimSiam* [3], only rely on positive pairs and therefore
58 can make use of smaller batch sizes. To understand why these methods do not collapse, the dynamics of these networks
59 are analysed with linear models in [1, 13]. From this analysis, the authors could derive ablations of *BYOL* where part of
60 the network is directly set to its optimal solution instead of being trained by gradient descent.

61 3 Method

62 In this section we will describe the methods of *BYOL* and *SimSiam* as well as their successor *DirectPred*.

63 3.1 BYOL & SimSiam

64 The network architecture of the models is shown in Figure 1. First, two augmented views X'_1 and X'_2 of an image X
65 are created and fed into the online network W and target network W_a respectively. Both of these networks have the
66 same architecture, a ResNet-18 (W_{enc}^x) as encoder [14], which is supposed to create hidden features and a projector
67 head W_{pro}^x , which is a two layer MLP, with purpose to map the feature space into a lower dimensional hidden space.
68 The online network also has an additional predictor head, again consisting of a two layer MLP. The target network has
69 a *StopGrad* function instead of a predictor head. Therefore during back propagation, only the weights of the online
70 network are updated via gradient descent. The loss between the output of the online and target network is equal to the
71 cosine-similarity loss function.

$$\mathcal{L}(\hat{Z}_1^{(O)}, \hat{Z}_2^{(T)}) = -\frac{\langle Z'_1, Z'_2 \rangle}{\|Z'_1\|_2 \|Z'_2\|_2} \quad (1)$$

72 Note, that the final loss of one image is the symmetric loss $\mathcal{L}(\hat{Z}_1^{(O)}, \hat{Z}_2^{(T)}) + \mathcal{L}(\hat{Z}_2^{(O)}, \hat{Z}_1^{(T)})$, since each augmentation is
73 given to both networks. As mentioned, the target network is not updated with gradient descent, but with an exponential

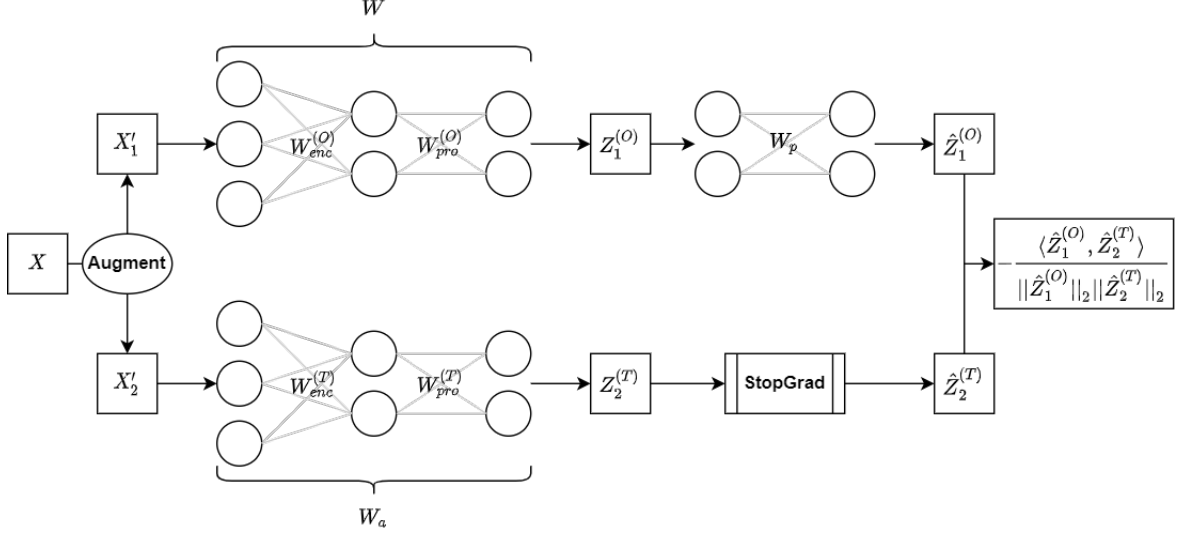


Figure 1: Network architecture for all presented methods

74 moving average (EMA). After each batch the target network will be set to $W_a = W_a + (1 - \tau)(W - W_a)$. In *SimSiam*
 75 the target network is set directly to the online network after each update, i.e $\tau = 0$.

76 3.2 DirectPred

77 [1] derives a one layer predictor head analytically with the analysis of the underlying learning dynamics of the models
 78 presented in Section 3.1 with an approximation of the actual network as a purely linear model. The learning dynamics
 79 of the networks are

$$\dot{W}_p = \alpha_p(-W_p W(X + X') + W_a X)W^\top - \eta W_p \quad (2)$$

$$\dot{W} = W_p^\top(-W_p W(X + X') + W_a X) - \eta W \quad (3)$$

$$\dot{W}_a = \beta(-W_a + W) \quad (4)$$

80 With $X = \mathbb{E}[\hat{x}\hat{x}^\top]$, where \hat{x} is the average augmented view of a datapoint and X' is the covariance matrix of the
 81 augmented views. α_p and β are multiplicative learning rate ratios, i.e $\alpha_p = \frac{\alpha_{\text{pred}}}{\alpha}$ and $\beta = \frac{1-\tau}{\alpha}$ (here α and α_{pred} are the
 82 learning rates for W and W_p respectively). In addition to the linearity of the network, three simplifying assumptions
 83 where made:

- 84 • The target network is always in a linear relationship with the online network (e.g. $W_a(t) = \tau(t)W(t)$)
- 85 • The original data distribution $p(X)$ is isotropic and its augmentation $\hat{p}(X'|X)$ has mean X and covariance $\sigma\mathbf{I}$
- 86 • The predictor W_p is symmetric

87 Based on these assumptions, one can show, that the eigenspaces of the output of the online network and the predictor
 88 W_p align. Let $F = WXW^\top$ (i.e. the output of the online network when it is approximated as a linear model), then it
 89 follows with the three assumptions, that the eigenspaces of these two matrices align over time (e.g. for all non-zero
 90 eigenvalues λ_{W_p}, λ_F of W_p and F , the corresponding normalized eigenvectors v_{W_p}, v_F are parallel, $v_{W_p}^\top v_F = 1$).
 91 With this alignment one can derive decoupled dynamics for the eigenvalues of W and W_p . By analysing this system, it
 92 can be shown that it has, depending on the weight decay parameter, several fixpoints, from which some are stable and
 93 some not. The trivial solution (the collapse) is one of them and the basin of attraction of these fixpoints varies with the
 94 relative learning rate of the predictor $\frac{\alpha_{\text{pred}}}{\alpha}$. With this analysis, [1] derives conditions under which the trivial fixpoint can
 95 be avoided. For a thorough mathematical analysis, we refer to [1]. In Section 5.1 we will present empirical evidence,
 96 that the symmetry assumption holds, and that the eigenspaces align. Furthermore, in Section 5.3 we will investigate
 97 the role of weight decay and the learning rate.

98 From the decoupled dynamics of the eigenvalues, we can also derive an analytical expression for the predictor W_p . Let
 99 $F = U\Omega U^\top$ be the eigen-decomposition of F with $\Omega = \text{diag}(\lambda_F^{(1)}, \dots, \lambda_F^{(d)})$ the diagonal matrix with the eigenvalues
 100 of F , then we can approximate the eigenvalues of W_p with

$$\lambda_{W_p}^{(j)} = \sqrt{\lambda_F^{(j)}} + \epsilon \max_j \lambda_F^{(j)} \quad (5)$$

101 and therefore set W_p to

$$W_p = U \text{diag}(\lambda_{W_p}^{(1)}, \dots, \lambda_{W_p}^{(d)}) U^\top \quad (6)$$

102 Note, that we cannot compute F directly, which is why we use a running average \hat{F} as approximation in practice

$$\hat{F} = \rho \hat{F} + (1 - \rho) \hat{Z} \quad (7)$$

103 where $\hat{Z} = \hat{Z}_1^{(O)} \hat{Z}_2^{(O)\top}$.

104 We denote this method *DirectPred* and in Section 5.2 we show, that *DirectPred* can perform similar to *BYOL* and
 105 *SimSiam*

106 4 Data & Configurations

107 We ran our experiments on Google Cloud Platform using Virtual Machine with a V100 GPU.

108 All experiments are conducted on CIFAR-10 [15], which contains 60 000 RGB images uniformly distributed over 10
 109 classes. The pre-training and the linear evaluation are done on the entire training set, which consists of 50 000 images.
 110 For the linear evaluation, only a linear layer is used on top of the encoder, where the weights of the encoder are frozen
 111 (i.e. we test how linearly separable the encoders output is). The reported accuracy results are produced from a test set
 112 containing 10 000 images. Also, to account for the small dimension of the CIFAR-10 images ($32 \times 32 \times 3$) we use $3 \times$
 113 3 convolutions and stride 1 without maximum pooling in the first block of the encoder.

114 To augment each image, we first do a random flip, take a random crop (up to 8% of the original size) of the image. Then
 115 we randomly adjust brightness, saturation, contrast and hue of the RGB image by a random factor ¹. Finally with a 20%
 116 chance we convert the image to grey scale.

117 **Self-supervised pretraining** In the basic setting, the online network use ResNet-18 as encoder, two layer projector
 118 MLP, two layer predictor MLP, where the first layer consists of 512 nodes, followed by BatchNorm and ReLU, and then
 119 a linear output layer with 128 nodes. For *BYOL* we use EMA to update target network and for *SimSiam* we directly set
 120 encoder and projector of target network to the weights of the online one ($\tau = 0$). We use SGD optimizer with learning
 121 rate 0.03, momentum 0.9 and weight decay (L2 penalty) of 0.0004. The predictor of *DirectPred* is set directly and are
 122 not trained with gradient descent and consist of one linear layer with 128 nodes. By SGD baseline for those methods
 123 we mean a network pre-trained with a one linear layer predictor with or without EMA. In all experiments, we use batch
 124 size of 128. For updating the target network we used the EMA parameter $\tau = 0.996$. For *DirectPred* we use $\epsilon = 0.1$
 125 and $\rho = 0.3$.

126 **Linear evaluation** In order to test the performance of the different models, we use linear evaluation, i.e. we train a
 127 linear layer on top of the ResNet-18 encoder with frozen weights for 100 epochs. This measures how linearly separable
 128 the learned representations of the encoder are. We use Adam optimizer [16] with polynomial decay of learning rate
 129 from $5e-2$ to $5e-4$. Images are normalized but we do not use augmentation for this part of training just as in the original
 130 repository for *DirectPred*.

131 5 Experiments and findings

132 In this section, we will first show that the assumptions and theoretical findings from Section 3.2 hold in practice.
 133 Finally, we will pre-train and use linear evaluation on the different models presented in Section 3 in order to test their
 134 performances.

¹for brightness, saturation and contrast we chose a value uniformly at random between 0.6 and 1.4. For adjusting the hue, we set the maximal value to 0.1

135 **5.1 Eigenspace alignment**

136 First, we pre-train *BYOL* and *SimSiam* keep track of the predictor heads symmetry and eigenspace alignment. In Figure
 137 2 we can see, that the assumption of a symmetric predictor W_p holds. Even without symmetry regularisation, W_p
 138 approaches symmetry during training. Also, we can see that for all non-zero eigenvalues of W_p the eigenspaces between
 139 F and W_p align as the training progresses.

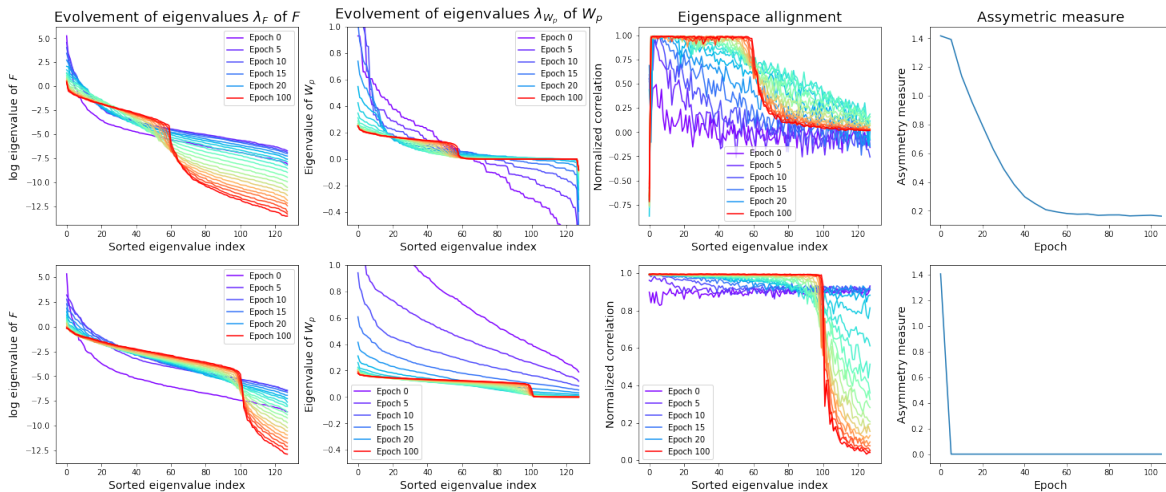


Figure 2: Pre-training *BYOL* for 100 epochs of CIFAR-10. **Top row:** *BYOL* without symmetry regularisation on W_p . **Bottom row:** *BYOL* with symmetry regularisation on W_p . The eigenvalues of F are plotted on the log scale, since the eigenvalues vary a lot. The asymmetry is measured by $\frac{\|W_p - W_p^T\|}{\|W_p\|}$

140 We ran the same Experiment for *SimSiam*, and can also see the same effect on the predictor and the alignment (Figure
 141 3). If we don't use a symmetric predictor, we also see that the eigenspaces for the non-zero eigenvalues align. However,
 142 once we use symmetry regularisation on W_p , all eigenvalues become zero, which shows that the network collapses. We
 143 will see later in Section 5.3 that we can prevent this collapse by using different learning rates α , α_{pred} and weight decay
 144 η , η_{pred} for W and W_p respectively.

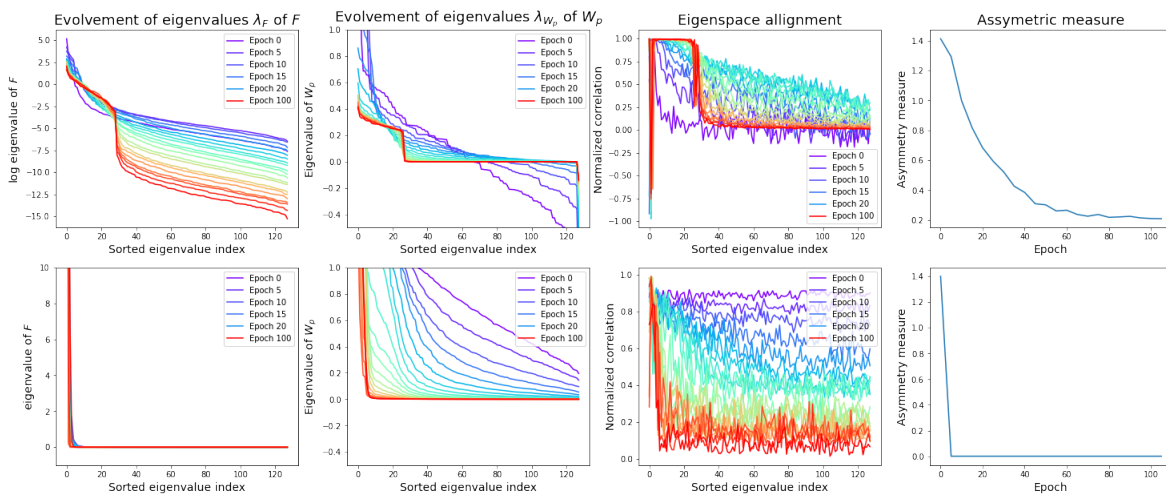


Figure 3: Pre-training *SimSiam* for 100 epochs of CIFAR-10. **Top row:** *SimSiam* without symmetry regularisation on W_p . **Bottom row:** *SimSiam* with symmetry regularisation on W_p . Note that the eigenvalues of F are not plotted on the log scale here, since we get 0 values.

145 **5.2 Performance**

146 **Byol & SimSiam** In table 1 we can see that the performance of *BYOL* increases slightly when using symmetry
 147 regularisation on the predictor. However, as already seen in Figure 3, when using no EMA, we observe that the network
 148 collapses. We observe in general better performance for models trained with EMA, given the same hyperparameters.
 149 However, we did not use extensive hyperparameter tuning, as performance is not the focus of our work.

	symmetric W_p	non symmetric W_p
EMA	85.7	84.2
No EMA	20.3	79.4

Table 1: Comparison of a two layer predictor with and without symmetry regularisation as well as with and without EMA (i.e first row is *BYOL* and second row is *SimSiam*).

150 **DirectPred** As we can see in Figure 2 & 3, the eigenspaces for both models align and therefore the theoretical
 151 assumptions of [1] hold. As we can see in Table 2, all models perform reasonably well, and can achieve almost the
 152 same performance as *BYOL* or *SimSiam*. However, as already mentioned earlier, we can see that models with EMA
 153 outperform models without EMA. In addition, we run an experiments where the predictor is only updated every 5th
 154 step according to Equation 6 and otherwise is updated with gradient decent, we call this method *DirectPred*₅. We
 155 can see that the hybrid method *DirectPred*₅ does not increase performance, however, according to [1] when training
 156 for 500 epochs, *DirectPred*₅ can outperform *DirectPred*. Due to computational constraints we cannot reproduce this
 157 experiment.

	SGD Baseline	DirectPred	DirectPred ₅
EMA	83.3%	84.7%	84.1%
No EMA	77.8%	78.6%	-

Table 2: The accuracies of SGD baselines, DirectPred and DirectPred with Frequency 5 with and without EMA

158 **5.3 Influence of weight decay and learning rate**

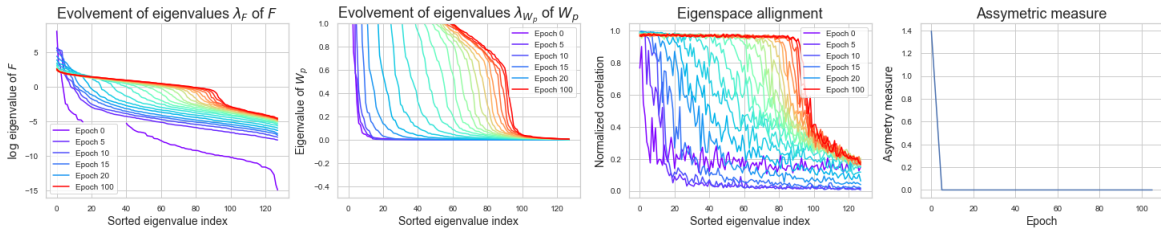


Figure 4: *SimSiam* with symmetric predictor but learning rates $\alpha = 0.2$, $\alpha_{\text{pred}} = 2$ and weight decay $\eta = 0$, $\eta_{\text{pred}} = 4e - 4$

159 As we can see in Figure 3, *SimSiam* with symmetric predictor does collapse. However, we can prevent this by adjusting
 160 the weight decay and learning rate. To make sure the network does converge to a stable non-collapsing fix-point, the
 161 weight decay of the predictor should be set higher than the rest of the network ($\eta_{\text{pred}} > \eta$, for mathematical analysis see
 162 [1]). By omitting weight decay, we are not able to stabilize the training of *SimSiam* with symmetric predictor and we
 163 can also see, that methods without weight decay perform worse, than with weight decay (Table 3). Also, to decrease
 164 the basin of attraction, of the trivial fixpoints, the learning rate of the predictor should be rather large compared to the
 165 learning rate of the rest of the network, i.e $\frac{\alpha_{\text{pred}}}{\alpha} \gg 1$ (see Section 3.2 in [1]).

166 **6 Challenges**

167 The original paper describes the methods and mathematical derivations well. Authors also share which hyperparameters
 168 they used in most of the experiments. Since the authors provided the open-source repository for the paper, we could

	symmetric W_p	regular W_p
$\eta = 0 \ \& \ \eta_{\text{pred}} = 0.0004$		
Byol	81.34 %	81.69 %
SimSiam	79.1 %	81.39 %
$\eta = \eta_{\text{pred}} = 0$		
Byol	80.78 %	80.42 %
SimSiam	20.27 %	79.22 %

Table 3: Byol and SimSiam trained with different values for the weight decay parameters. For all experiments in this Table, we set the learning rates $\alpha_{\text{pred}} = 2$ and $\alpha = 0.2$. Note, that the important condition is $\frac{\alpha_{\text{pred}}}{\alpha} \gg 1$, i.e. we got only slightly worse results with $\alpha_{\text{pred}} = 0.2$ and $\alpha = 0.02$

169 check some of the details of the experiments there. However, as the code is not well-structured it was at times
170 challenging to analyse. Furthermore, not all of the experiments are shared in the repository, for example there is no
171 code which produces eigenspace experiments results or config for weight decay experiment.

172 The reproduced paper did not outlined self-contained description on the methods it used as it built upon previous works.
173 Thanks to the detailed description of *BYOL* by Grill et. al. [2] we were able to reproduce the paper achieving similar
174 results as the authors.

175 Due to time constraints we decided to use CIFAR-10 instead of STL-10 which was used in most of the experiments in
176 the reproduced paper. However, claims tested by us in this work are not restricted to one dataset and we shown that they
177 indeed hold in a different setting. One of the main challenges was the large amount of computations required for all the
178 experiments, it took around 4 hours and 30 minutes to pre-train and fine tune a single model, and in total we trained for
179 around 100+ hours.

180 Our work is implemented in TensorFlow and one of the challenges was differences between TensorFlow and PyTorch
181 libraries. For instance, in PyTorch one of the parameters of the SGD optimizer is weight decay (L2 penalty), in
182 TensorFlow we had to implement it by hand as TensorFlow’s SGDW implements only Decoupled Weight Decay
183 Regularization [17]. Furthermore, image augmentation methods such as ColorJitter from PyTorch do not have exact
184 corresponding methods in Tensroflow. We used a custom way to do it so that augmentations are as close as possible to
185 the original version.

186 7 Conclusion

187 In this work we study and reimplement three architectures used to give insight into self-supervised representation
188 learning without contrastive pairs namely *BYOL*, *SimSiam*, *DirectPred* and their ablations. Our experimental results
189 aligned well with both the theoretical analysis about the eigenspaces and the symmetric assumptions made in [1] and
190 translate to other dataset than used in the paper. Lastly, we confirmed that SimSiam can be prevented from collapsing
191 with the use of weight decay and adjusting a learning rate of predictor.

192 Furthermore, we confirm the claim that *DirectPred* outperforms its one-layer SGD alternative. However, we cannot
193 report that *DirectPred* could outperform *Byol*. This may be due to the fact that we used CIFAR-10 as opposed to STL-10
194 in the original paper. This leaves us with the conclusion, that *DirectPred* gives valuable insights into the dynamics of
195 unsupervised representation learning without contrastive pairs, but do not necessarily build new state of the art models
196 themselves.

197 8 Ethical considerations

198 Self-supervised learning circumvents label scarcity which is one of the most common problems when applying ML
199 to new scenarios. This can have both positive and negative consequences. On one hand, it can accelerate important
200 developments for example in medical diagnosis. However, it can also be used in unethical ways such as in surveillance
201 or military equipment. Furthermore, there will be less need for people labelling datasets which will result in reduction
202 of job positions in this area.

203 **References**

- 204 [1] Yuandong Tian, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning dynamics without
205 contrastive pairs, 2021.
- 206 [2] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya,
207 Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray
208 Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised
209 learning, 2020.
- 210 [3] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning, 2020.
- 211 [4] Jane Bromley, James Bentz, Leon Bottou, Isabelle Guyon, Yann Lecun, Cliff Moore, Eduard Sackinger, and
212 Rookpak Shah. Signature verification using a "siamese" time delay neural network. *International Journal of*
213 *Pattern Recognition and Artificial Intelligence*, 7:25, 08 1993.
- 214 [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive
215 learning of visual representations, 2020.
- 216 [6] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised
217 models are strong semi-supervised learners, 2020.
- 218 [7] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual
219 representation learning, 2020.
- 220 [8] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding,
221 2019.
- 222 [9] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2014.
- 223 [10] Pascal Vincent, Hugo Larochelle, Y. Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust
224 features with denoising autoencoders. pages 1096–1103, 01 2008.
- 225 [11] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning, 2017.
- 226 [12] Jeff Donahue and Karen Simonyan. Large scale adversarial representation learning, 2019.
- 227 [13] Xiang Wang, Xinlei Chen, Simon S. Du, and Yuandong Tian. Towards demystifying representation learning with
228 non-contrastive self-supervision, 2021.
- 229 [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- 230 [15] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10.
- 231 [16] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
232 2014.
- 233 [17] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*,
234 2017.