

HOW DOES SIMSIAM AVOID COLLAPSE WITHOUT NEGATIVE SAMPLES? TOWARDS A UNIFIED UNDERSTANDING OF PROGRESS IN SSL

Anonymous authors

Paper under double-blind review

ABSTRACT

Towards avoiding collapse in self-supervised learning (SSL), contrastive loss is widely used but often requires a large number of negative samples. Without negative samples yet achieving competitive performance, a recent work (Chen & He, 2021) has attracted significant attention for providing a minimalist simple Siamese (SimSiam) method to avoid collapse. However, the reason for its success remains not fully clear and our investigation starts by revisiting the explanatory claims in the SimSiam. After refuting their claims, we introduce vector decomposition for analyzing the collapse based on the gradient analysis of l_2 normalized vector. This yields a unified perspective on how negative samples and SimSiam predictor alleviate collapse. Such a unified perspective comes timely for understanding the recent progress in SSL.

1 INTRODUCTION

Beyond the success of NLP (Lan et al., 2020; Radford et al., 2019; Devlin et al., 2019; Su et al., 2020; Nie et al., 2020), self-supervised learning (SSL) has also shown its potential in the field of vision tasks (Li et al., 2021; Chen et al., 2021; El-Nouby et al., 2021). Without the ground-truth label, the core of most SSL methods lies in learning an encoder with augmentation invariant representation (Bachman et al., 2019; He et al., 2020; Chen et al., 2020a; Caron et al., 2020; Grill et al., 2020). Specifically, they often minimize the representation distance between two positive samples, *i.e.* two augmented views of the same image, based on a Siamese network architecture (Bromley et al., 1993). It is widely known that for such Siamese networks there exists a degenerate solution, *i.e.* all outputs “collapsing” to an undesired constant (Chen et al., 2020a; Chen & He, 2021). Early works have attributed the collapse to lacking a repulsive component in the optimization goal and adopted negative samples, *i.e.* views of different samples, to alleviate this problem. Introducing momentum into the target encoder, BYOL shows that Siamese architectures can be trained with only positive pairs. A technical report (Fetterman & Albrecht, 2020) has suggested the importance of batch normalization in BYOL for its success, however, a recent work Richemond et al. (2020) refutes their claim by showing BYOL works without BN.

More recently, SimSiam (Chen & He, 2021) has caught great attention by further simplifying BYOL in terms of removing momentum encoder, which can be seen as a major milestone achievement in SSL through providing a *minimalist* method for achieving competitive performance. However, the reason for the success of SimSiam remains not fully clear with a core question summarized as:

How does SimSiam avoid collapse without negative samples?

Our work attempts to answer the above question and our investigation starts with revisiting the explanatory claims in the original SimSiam paper. Notably, two components, *i.e.* stop gradient and predictor, are observed to be essential for the success of SimSiam (Chen & He, 2021). The reason has been mainly attributed to the stop gradient (Chen & He, 2021) by hypothesizing that it implicitly involves two sets of variables and SimSiam behaves like alternating between optimizing each set. Chen & He argue that the predictor h is helpful in SimSiam because h fills the gap to approximate expectation over augmentations (EOA).

Unfortunately, the above explanatory claims are found to be flawed due to reversing the two paths with and without gradient. This motivates us to find an alternative explanation, for which we in-

roduce a simple yet intuitive framework for facilitating analyzing collapse in SSL. Specifically, we propose to decompose a representation vector into center and residual components. This decomposition facilitates understanding which gradient component is beneficial for avoiding collapse. Under this framework, we attempt to understand why a Naive Siamese architecture causes collapse, and provide a unified perspective on how negative samples help prevent such collapse as well as how predictor in SimSiam achieves it *without* negative samples. For both InfoNCE with negative samples (He et al., 2019; Chen et al., 2020b;a; Tian et al., 2019; Khosla et al., 2020) and SimSiam with asymmetric predictor, their center gradient component behaves like de-centering for preventing collapse. Interestingly, we also find that for both, the residual vector component achieves de-correlation which somewhat surprisingly also prevents collapse. Overall, our work demystifies the success of InfoNCE and SimSiam, which bridges the gap between them and recent frameworks with explicit de-centering and de-correlation. Our work contributes to a unified understanding of the recent progress of different frameworks in SSL.

Finally, our work points out an important detail on whether BN, which removes the mean center, is essential for preventing a collapse in BYOL or SimSiam. Specifically, we reveal that BN alone prevents collapse only when BN is adopted immediately after the MSE loss. Towards more explainable SimSiam our work investigates multiple new variants of predictors for preventing collapse. We follow prior works to set the experiment, see [Appendix A.1](#).

2 REVISITING SIMSIAM AND ITS EXPLANATORY CLAIMS

l_2 normalized vector and optimization goal. SSL trains an encoder f for learning discriminative representation and we denote such representation as a vector z , *i.e.* $f(x) = z$ where x is a certain input. For the augmentation-invariant representation, a straightforward goal is to minimize the distance between the representations of two positive samples, *i.e.* augmented views of the same image, for which mean squared error (MSE) is a default choice. To avoid scale ambiguity, the vectors are often l_2 normalized, *i.e.* $Z = z/\|z\|_2$ (Chen & He, 2021), before calculating the MSE:

$$\mathcal{L}_{MSE} = (Z_a - Z_b)^2/2 - 1 = -Z_a \cdot Z_b = L_{\cosine}, \quad (1)$$

which shows the equivalence of normalized MSE loss to cosine loss (Grill et al., 2020).

Collapse in SSL and solution of SimSiam. Based on a Siamese architecture, the loss in Eq. 1 causes the collapse, *i.e.* f always outputs a constant regardless of the input. We refer this Siamese architecture with loss Eq. 1 as *Naive Siamese* in the remainder of paper. Towards alleviating this phenomenon, multiple works have investigated solutions, among them contrastive loss with negative samples is widely used (Chen et al., 2020a). **SimSiam solves the collapse problem via predictor and stop gradient, based on which the encoder is optimized with a symmetric loss:**

$$L_{SimSiam} = -(P_a \cdot \text{sg}(Z_b) + P_b \cdot \text{sg}(Z_a)), \quad (2)$$

where $\text{sg}(\cdot)$ is *stop gradient* and P is the output of predictor h , *i.e.* $p = h(z)$ and $P = p/\|p\|_2$.

2.1 REVISING EXPLANATORY CLAIMS IN SIMSIAM

Interpreting stop gradient as AO. Chen & He hypothesize that the stop gradient in Eq. 2 is an implementation of Alternating between the Optimization of two sub-problems, which is denoted as AO. Specifically, with the loss considered as $\mathcal{L}(\theta, \eta) = \mathbb{E}_{x, \mathcal{T}} \left[\|\mathcal{F}_\theta(\mathcal{T}(x)) - \eta_x\|_2^2 \right]$, the optimization objective $\min_{\theta, \eta} \mathcal{L}(\theta, \eta)$ can be solved by $\eta^t \leftarrow \arg \min_{\eta} \mathcal{L}(\theta^t, \eta)$ and $\theta^{t+1} \leftarrow \arg \min_{\theta} \mathcal{L}(\theta, \eta^t)$ alternatively. It is acknowledged in (Chen & He, 2021) that their hypothesis does not fully explain why the collapse is prevented. Nonetheless, they mainly attribute SimSiam success to the stop gradient with the interpretation that AO might make it difficult to approach a constant for $\forall x$.

Interpreting predictor as EOA. Their AO problem (Chen & He, 2021) is formulated independent of predictor h , for which they believe that the usage of predictor h is related to approximating EOA for filling the gap of ignoring $\mathbb{E}_{\mathcal{T}}[\cdot]$ in a sub-problem of AO. The approximation of $\mathbb{E}_{\mathcal{T}}[\cdot]$ is summarized in [Appendix A.2](#). Chen & He support their interpretation by proof-of-concept experiments. Specifically, they show that updating η_x with a moving-average $\eta_x^t \leftarrow m * \eta_x^t + (1 - m) * \mathcal{F}_{\theta^t}(\mathcal{T}^t(x))$ can help prevent collapse without predictor (see Fig 1 (b)). Given that the training completely fails when the predictor and moving average are both removed, at first sight, their reasoning seems valid.

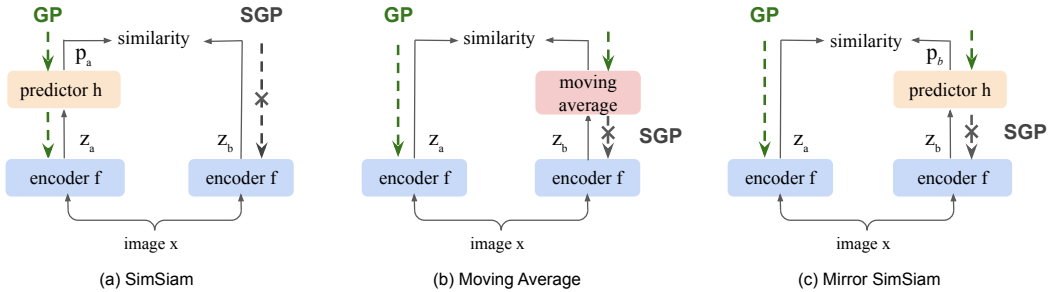


Figure 1: **Reasoning Flaw in SimSiam.** (a) The standard SimSiam architecture. (b) Moving-Average Model proposed in the proof-of-concept experiment (Chen & He, 2021). (c) Mirror SimSiam, which has the same model architecture as SimSiam but with the reverse of GP and SGP.

2.2 DOES THE PREDICTOR FILL THE GAP TO APPROXIMATE EOA?

Reasoning flaw. Taking the stop gradient into account, we divide SimSiam model into two sub-model with different paths and term them Gradient Path (GP) and Stop Gradient Path (SGP). For SimSiam, only the sub-model with GP includes the predictor (see Fig. 1 (a)). We point out that their reasoning flaw of predictor analysis lies in the *reverse of GP and SGP*. By default, the SGP is on the same side of moving-average sub-model, as shown in Fig. 1 (b). Conceptually, Fig. 1 (b) is similar to Fig. 1 (c) instead of Fig. 1 (a). It is worth mentioning that the Mirror SimSiam in Fig. 1 (c) is what stop gradient in the original SimSiam avoids. Therefore, it is problematic to perceive h as EOA.

Method	# aug	Collapse	Std	Accuracy (%)
Moving average	2	×	0.0108	46.57
Same batch	10	✓	0	1
Same batch	25	✓	0	1

Table 1: Influence of Explicit EOA. Detailed setup is reported in [Appendix A.3](#)

Explicit EOA does not prevent collapse. (Chen & He, 2021) points out that “in practice, it would be unrealistic to actually compute the expectation $\mathbb{E}_{\mathcal{T}}[\cdot]$. But it may be possible for a neural network (e.g., the predictor h) to learn to predict the expectation, while the sampling of \mathcal{T} is implicitly distributed across multiple epochs.” If implicitly sampling across multiple epochs is beneficial, explicitly sampling sufficient large N augmentations in a batch with the latest model would be more beneficial to approximate $\mathbb{E}_{\mathcal{T}}[\cdot]$. However, Table 1 showed that the collapse still occurs and suggests the success of moving average should not be interpreted from the EOA perspective.

2.3 ASYMMETRIC INTERPRETATION OF PREDICTOR WITH STOP GRADIENT IN SIMSIAM

Symmetric Predictor does not prevent collapse. The difference between Naive Siamese and SimSiam lies in whether the gradient in backward propagation flows through a predictor, however, we show that this propagation helps avoid collapse only when the predictor is not included in the SGP path. With h being trained the same as Eq. 2, we optimize the encoder f through replacing the \mathbf{Z} in Eq. 2 with \mathbf{P} . The results in [Table. 2](#) show that it still leads to collapse. Actually, this is well expected by perceiving h to be part of the new encoder F , *i.e.* $\mathbf{p} = F(x) = h(f(x))$. In other words, the symmetric architectures with and without predictor h both lead to collapse.

Predictor with stop gradient is asymmetric. Clearly, the success of SimSiam lies in the asymmetric architecture, *i.e.* one path with predictor and the other without predictor. Under this asymmetric architecture, the role of stop gradient is to only allow the path with predictor to be optimized with the encoder output as target, not vice versa. In other words, the SimSiam avoids collapse by excluding Mirror SimSiam (Fig. 1 (c)) which has a loss (mirror-like Eq 2) as $\mathcal{L}_{\text{Mirror}} = -(\text{sg}(\mathbf{P}_a) \cdot \mathbf{Z}_b + \text{sg}(\mathbf{P}_b) \cdot \mathbf{Z}_a)$.

Predictor vs. inverse predictor. We interpret h as a function mapping from z to p , and introduce a conceptual inverse h^{-1} for mapping from p to z , *i.e.* $z = h^{-1}(p)$. Here, as shown in Table 2, SimSiam with symmetric predictor (Fig. 2 (b)) leads to collapse, while SimSiam (Fig. 1 (a)) avoids collapse. With the conceptual h^{-1} , we interpret Fig. 1 (a) the same as Fig. 2 (c) which differs from Fig. 2 (b) via changing the optimized target from p_b to z_b , *i.e.* $z_b = h^{-1}(p_b)$. This interpretation suggests that the collapse can be avoided by processing the optimized target with h^{-1} . By contrast, Fig. 1 (c) and Fig. 2 (a) both lead to **collapse**, suggesting that processing the optimized target with h is *not* beneficial for preventing collapse. Overall, asymmetry alone does not guarantee collapse avoidance, and the optimized target needs to be processed by h^{-1} not h .

Method	Collapse	Top-1
Simsiam	×	66.62
Mirror SimSiam	✓	1
Naive Siamese	✓	1
Symmetric Predictor	✓	1

Table 2: Various siamese architectures. Detailed trend and setup are in Appendix A.4

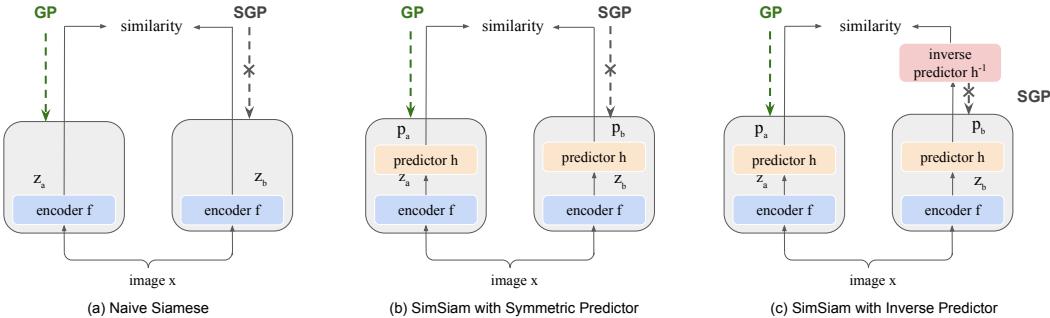


Figure 2: Different architectures of Siamese model (a), (b), and (c). Here in (c), an additional inverse predictor h^{-1} is designed that has same architecture as predictor h .

Trainable h^{-1} and its implication on EOA. In the above, we propose a *conceptual* inverse predictor h^{-1} in Fig. 2 (c), however, it remains yet unknown whether the model with h^{-1} (Fig. 2 (c)) is equivalent to the original SimSiam (Fig. 1 (a)) *experimentally*? In other words, is the model with h^{-1} trainable to achieve performance comparable to SimSiam? To answer this question, we attempt to train h^{-1} , for which a detailed setup is reported in Appendix A.5. The results in Fig. 3 show that a trainable h^{-1} leads to slightly inferior performance, which is expected because h^{-1} cannot make z_b^* perfectly same as z_b , where z_b^* is the output of h^{-1} . Note that they would be perfectly equivalent if $z_b^* = z_b$. Despite a slight performance drop, the results confirm that h^{-1} is trainable. The fact h^{-1} is trainable provides additional evidence that the role h plays in SimSiam is not EOA because theoretically the h^{-1} cannot restore a random augmentation \mathcal{T}' from an expectation p , where $p = h(z) = \mathbb{E}_{\mathcal{T}}[\mathcal{F}_{\theta^t}(\mathcal{T}(x))]$.

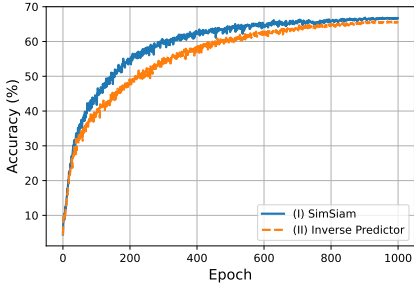


Figure 3: Comparison of original SimSiam and SimSiam with Inverse Predictor.

3 VECTOR DECOMPOSITION FRAMEWORK FOR UNDERSTANDING COLLAPSE

We introduce a gradient analysis framework based on vector decomposition to facilitate their investigation. By default, SimCLR (Chen et al., 2020a) and SimSiam (Chen & He, 2021) both adopt l_2 normalization in their loss for avoiding scale ambiguity. We treat the l_2 normalized vector, *i.e.* Z as the encoder output, which significantly simplifies the gradient derivation and relevant analysis.

3.1 ANALYSIS FRAMEWORK FROM THE VECTOR DECOMPOSITION PERSPECTIVE

Vector decomposition. For the purpose of analysis, we propose to decompose \mathbf{Z} into two parts, $\mathbf{Z} = \mathbf{o} + \mathbf{r}$, where \mathbf{o} , \mathbf{r} denote *center vector* and *residual vector* respectively. Specifically, the center vector \mathbf{o} is defined as an average of \mathbf{Z} over the whole representation space $\mathbf{o}_z = \mathbb{E}[\mathbf{Z}]$, however, we approximate it with all vectors in current mini-batch, i.e. $\mathbf{o}_z = \frac{1}{M} \sum_{m=1}^M \mathbf{Z}_m$, where M is the mini-batch size. We define the residual vector \mathbf{r} as the residual part of \mathbf{Z} , i.e. $\mathbf{r}_a = \mathbf{Z}_a - \mathbf{o}_z$.

Collapse from the vector perspective. A Naive Siamese is well expected to collapse since the loss is designed to minimize the distance between positive samples, for which a constant constitutes an optimal solution to minimize such loss. When such collapse occurs, $\forall i, \mathbf{Z}_i = \frac{1}{M} \sum_{m=1}^M \mathbf{Z}_m = \mathbf{o}_z$, where i denotes a random sample index, which shows the constant vector is \mathbf{o}_z when collapse occurs. This interpretation only suggests a possibility that a dominant \mathbf{o} can be one of the viable solutions, while the optimization, such as SimSiam, might still lead to a non-collapse solution. It merely describes \mathbf{o} as the *consequence* of the collapse, and our work investigates the *cause* of such collapse through analyzing the influence of individual gradient components, i.e. \mathbf{o} and \mathbf{z} .

Competition between \mathbf{o} and \mathbf{r} . Complementary to the Standard Deviation (Std) (Chen & He, 2021), we introduce the ratio of \mathbf{o} in \mathbf{z} , i.e. $m_o = \|\mathbf{o}\|/\|\mathbf{z}\|$, where $\|\cdot\|$ is the L_2 norm, for indicating collapse. Similarly, the ratio of \mathbf{r} in \mathbf{z} is defined as $m_r = \|\mathbf{r}\|/\|\mathbf{z}\|$. When collapse happens, i.e. all vectors \mathbf{Z} are close to each other or \mathbf{o} , m_o approaches 1 and m_r approaches 0, which is not desirable for SSL. A desirable case would be a relatively small m_o and a relatively large m_r , suggesting a relatively small (large) contribution of \mathbf{o} (\mathbf{r}) in each \mathbf{Z} . We interpret the cause of collapse as a competition between \mathbf{o} and \mathbf{r} where \mathbf{o} dominates over \mathbf{r} , i.e. $m_o \gg m_r$. For Eq 1, the derived negative gradient on \mathbf{Z}_a (ignoring \mathbf{Z}_b for simplicity due to symmetry) is shown as:

$$\mathcal{G}_{\text{cosine}} = -\frac{\partial \mathcal{L}_{\text{MSE}}}{\partial \mathbf{Z}_a} = \mathbf{Z}_b - \mathbf{Z}_a \iff -\frac{\partial \mathcal{L}_{\text{cosine}}}{\partial \mathbf{Z}_a} = \mathbf{Z}_b, \quad (3)$$

where the gradient component \mathbf{Z}_a is a *dummy* term because the loss $-\mathbf{Z}_a \cdot \mathbf{Z}_a = -1$ is a constant having zero gradient on the encoder f .

Conjecture1. With $\mathbf{Z}_a = \mathbf{o} + \mathbf{r}_a$, we conjecture that the gradient component of \mathbf{o} is expected to update the encoder to boost the center vector thus increase m_o , while the gradient component of \mathbf{r}_a is expected to behave in the opposite direction to increase m_r thus decrease m_o . On the other hand, a random gradient component is expected to have relatively small influence on the trend.

To verify the conjecture, we revisit the *dummy* gradient term \mathbf{Z}_a . We design loss $-\mathbf{Z}_a \cdot \text{sg}(\mathbf{o}_z)$ and $-\mathbf{Z}_a \cdot \text{sg}(\mathbf{Z}_a - \mathbf{o}_z)$ to show the influence of gradient component \mathbf{o} and \mathbf{r}_a , respectively. The results in Fig. 4 (a) show that gradient component \mathbf{o}_z has the effect of increasing m_o while decreasing m_r . On the contrary, \mathbf{r}_a helps increase m_r while decreasing m_o . Overall, the results verify Conjecture1.

Revisit collapse in a symmetric architecture.

Based on Conjecture1, here, we provide an intuitive interpretation on why a symmetric Siamese architecture, such as Fig. 2 (a) and (b), cannot be trained without collapse. The gradient in Eq 3 can be interpreted as two equivalent forms, from which we choose $\mathbf{Z}_b - \mathbf{Z}_a = (\mathbf{o}_z + \mathbf{r}_b) - (\mathbf{o}_z + \mathbf{r}_a) = \mathbf{r}_b - \mathbf{r}_a$. Since \mathbf{r}_b comes from the same positive sample as \mathbf{r}_a , it is expected that \mathbf{r}_b also increases m_r , however, this effect is expected to be smaller than that of \mathbf{r}_a , thus causing collapse.

Basic gradient and Extra gradient components. Extra Gradient component, denoted as \mathcal{G}_e , needs to be introduced to break the **symmetry** for preventing collapse. As the term suggests, \mathcal{G}_e is defined as a gradient term relative to the basic gradient in a naive (symmetric) Siamese architecture. For example, in Fig. 2 (a), the negative gradient on \mathbf{Z}_a is \mathbf{Z}_b (see Eq 3) and \mathbf{Z}_b can be seen as the basic gradient because \mathbf{Z}_a and \mathbf{Z}_b are represented by the same encoder. The extra gradient introduced by negative samples can thus be perceived as \mathcal{G}_e . Similarly, in Fig. 2 (b), the negative gradient on

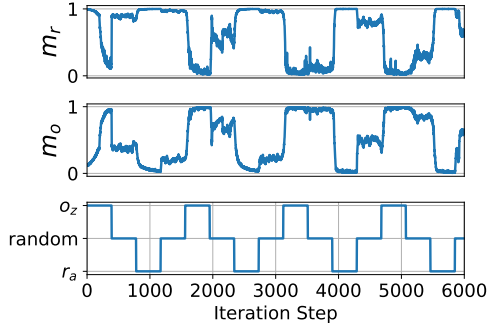


Figure 4: Influence of various gradient components on m_r and m_o .

P_a is P_b (see Eq 3) and P_b can be seen as the basic gradient. With the basic gradient P_b , the extra gradient G_e caused by the asymmetric predictor is then $Z_b - P_b$ for Fig. 1 (a). With $G_e = o_e + r_e$, we analyze how its center vector o_e and residual vector r_e in various setups influence the collapse in the following.

3.2 A TOY EXAMPLE EXPERIMENT WITH NEGATIVE SAMPLE

Which repulsive component helps avoid collapse? Existing works often attribute the collapse in Naive Siamese to lacking a *repulsive part* during the optimization. This explanation has motivated previous works to adopt contrastive learning, *i.e.* attracting the positive samples while *repulsing* the negative samples. We experiment with a simple triplet loss¹, $\mathcal{L}_{tri} = -Z_a \cdot \text{sg}(Z_b - Z_n)$, where Z_n indicates the representation of a Negative sample. The derived negative gradient is $Z_b - Z_n$, where Z_b is the basic gradient component in the symmetric architecture and thus $G_e = -Z_n$ in this setup. For a sample representation Z , what determines it as a positive sample for attracting or a negative sample for repulsing is the component r , *thus it might be tempting to guess that r_e is the key component of repulsive part that avoids the collapse*. However, our results show that the component beneficial for preventing collapse inside G_e is o_e instead of r_e .

Method	$\mathcal{L}_{triplet}$	Std	m_o	m_r	Collapse	Accuracy (%)
Baseline	$-Z_a \cdot \text{sg}(Z_b + G_e)$	0.020	0.026	0.99	×	36.21
Removing r_e	$-Z_a \cdot \text{sg}(Z_b + o_e)$	0.02005	0.026	0.99	×	47.41
Removing o_e	$-Z_a \cdot \text{sg}(Z_b + r_e)$	0	1	0	✓	1

Table 3: Gradient component analysis with a random negative sample.

To explore the individual influence of o_e and r_e in the G_e , we design two experiments by removing one component while keeping the other one. In the first experiment, we remove the r_e in G_e while keeping the o_e . In contrast, the o_e is removed while keeping the r_e in the second experiment. The results are shown in Table 3. In contrast to what existing explanations may expect, we find that the gradient component o_e prevents collapses. Actually, $o_e = -o_z$ in this setup, thus o_e has the de-centering role and prevents collapse. On the contrary, r_e does not prevent collapse and keeping r_e even decreases the performance when the collapse does not occur. Since the negative sample is randomly chosen, r_e just behaves like a random noise on the optimization.

3.3 DECOMPOSED GRADIENT ANALYSIS FOR SIMSIAM

It is worth mentioning that we cannot derive the gradient on z_a in SimSiam because the predictor is a nonlinear MLP module. The negative gradient on P_a for $\mathcal{L}_{SimSiam}$ in Eq. 2 is derived as:

$$\mathcal{G}_{SimSiam} = -\frac{\partial \mathcal{L}_{SimSiam}}{\partial P_a} = Z_b = P_b + (Z_b - P_b) = P_b + G_e \quad (4)$$

where G_e indicates the aforementioned extra gradient component. To investigate the influence of o_e and r_e on the collapse, similar to the analysis with the toy example experiment in Sec. 3.2, we design the experiment by removing one component while keeping the other one. The results are reported in Table 4. As expected, the model collapses when both components in G_e are removed and the best performance is achieved when both components are kept. Interestingly, the model does not collapse when either o_e or r_e is kept.

o_e	r_e	Collapse	Top-1 (%)
✓	✓	×	66.62
✓	×	×	48.08
×	✓	×	66.15
×	×	✓	1

Table 4: Gradient component analysis for SimSiam.

How o_e prevents collapse in SimSiam. Here, o_z and o_p are used to denote the center vector of Z and P , respectively, for differentiation. *With Conjecture 1, it is well expected that o_e helps prevent collapse if o_e contains negative o_p since the analyzed vector is P_a .* To determine the amount of component of o_p existing in o_e , we measure the cosine similarity between $o_e - \eta_p o_p$ and o_p for a wide range of η_p , where $o_e = o_z - o_p$. The results in Fig 5 (a) show that when η_p is around -0.5 for their cosine similarity being close to zero, suggesting o_e has $-0.5 o_p$. Overall, the o_e

¹Note that the triplet loss here does not have clipping form as in Schroff et al. (2015) for simplicity.

indeed contains negative \mathbf{o}_p , thus naturally explaining why the gradient component \mathbf{o}_e helps prevent collapse. Similarly, we also evaluate the amount of component \mathbf{o}_z existing in the negative \mathbf{o}_e , *i.e.* $\mathbf{o}_p - \mathbf{o}_z$, via reporting the similarity between $-\mathbf{o}_e - \eta_z \mathbf{o}_z$ and \mathbf{o}_z and the results show that their cosine similarity is zero when η_z is set to around 0.2. This positive η_z explains why Fig. 2(c) causes collapse. In other words, processing the target vector \mathbf{Z}_b with h^{-1} , as in Fig. 2 (c), alleviates collapse ($\eta_p = -0.5$), while processing it with h , as in Fig. 1(c), actually strengthens the collapse ($\eta_z = 0.2$). Note that a negative η has the de-centering role for preventing collapse.

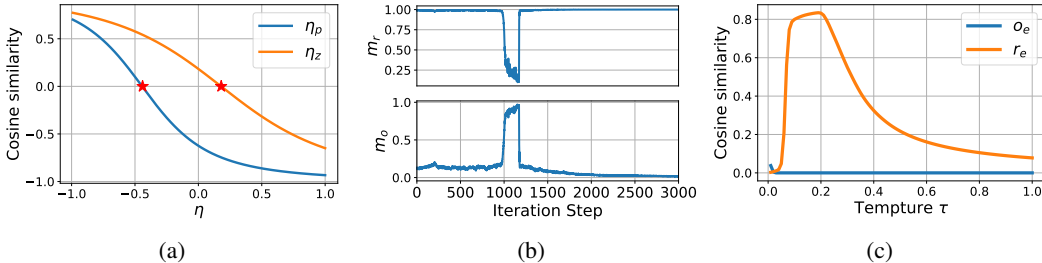


Figure 5: (a) Investigating the amount of \mathbf{o}_p existing in $\mathbf{o}_z - \mathbf{o}_p$ and the amount of \mathbf{o}_z existing in $\mathbf{o}_p - \mathbf{o}_z$. (b) Normally train the model as SimSiam for 5 epochs, then using collapsing loss for 1 epoch to reduce m_r , followed by a correlation regularization loss. (c) Cosine similarity between \mathbf{r}_e (\mathbf{o}_e) and gradient on \mathbf{Z}_a induced by a correlation regularization loss.

Interesting phenomenon for \mathbf{r}_e . In the toy example experiment in Section 3.2, \mathbf{r}_e is found to be *not* beneficial for preventing collapse and keeping \mathbf{r}_e even decreases the performance. Interestingly, as shown in Table 4, we find that \mathbf{r}_e alone is sufficient for preventing collapse and even achieves comparable performance as \mathbf{G}_e . Motivated from preventing the collapse to a constant, multiple prior works, such as W-MSE (Ermolov et al., 2021), Barlow-twins (Zbontar et al., 2021), DINO (Caron et al., 2021), explicitly adopt de-centering to prevent collapse. Despite various motivations, we find that they all implicitly introduce an \mathbf{o}_e that contains negative center vector. The success of their approaches aligns well with our conjecture 1 as well as our above empirical results. They all attribute their reason for not collapsing to de-centering, *i.e.* introducing an \mathbf{o}_e that contains negative center vector. To our knowledge, we are the first to show that \mathbf{r}_e can also prevent collapse. We leave relevant explanation to Section 3.4.

3.4 DIMENSIONAL DE-CORRELATION HELPS PREVENT COLLAPSE

Conjecture 2 and motivation. We conjecture that dimensional de-correlation increases m_r for preventing collapse. The motivation is straightforward as follows. The dimensional correlation would be minimum if only a single dimension has a very high value for every individual class and the dimension changes for different classes. In another extreme case, when all the dimensions have the same values, equivalent to have a single dimension, which already collapses by itself in the sense of losing representation capacity. Conceptually, \mathbf{r}_e has no direct influence on the center vector, thus we interpret that \mathbf{r}_e prevents collapse through increasing m_r .

To verify the above conjecture, we train SimSiam normally with the loss in Eq 2 and train with several epochs with the loss in Eq 1 for intentionally decreasing the m_r to close to zero. Then, we train the loss with only a correlation regularization term, which is detailed in Appendix A.6. The results in Fig. 5 (b) show that this regularization term increases m_r at a very fast rate.

Dimensional de-correlation in InfoNCE. Contrastive learning often requires a large number of negative samples and InfoNCE loss is a default choice in multiple seminal works (Sohn, 2016; Wu et al., 2018; Oord et al., 2018). The derived negative gradient of InfoNCE on \mathbf{Z}_a is proportional to $\mathbf{Z}_b + \sum_{i=0}^N -\lambda_i \mathbf{Z}_i$, where $\lambda_i = \frac{\exp(\mathbf{Z}_a \cdot \mathbf{Z}_i / \tau)}{\sum_{i=0}^N \exp(\mathbf{Z}_a \cdot \mathbf{Z}_i / \tau)}$, and $\mathbf{Z}_0 = \mathbf{Z}_b$ for notation simplicity. The detailed derivation is given in Appendix A.7. The extra gradient component $\mathbf{G}_e = \sum_{i=0}^N -\lambda_i \mathbf{Z}_i = -\mathbf{o}_z - \sum_{i=0}^N \lambda_i \mathbf{r}_i$, for which $\mathbf{o}_e = -\mathbf{o}_z$ and $\mathbf{r}_e = -\sum_{i=0}^N \lambda_i \mathbf{r}_i$. Clearly, \mathbf{o}_e contains negative \mathbf{o}_z as de-centering for avoiding collapse, which is equivalent to the toy example in Section 3.2 when the \mathbf{r}_e is removed. In the toy example, \mathbf{r}_e is found to be not helpful for avoiding collapse, however \mathbf{r}_e for

InfoNCE loss alone is found to prevent the collapse. The results are shown in Figure 6. Regarding r_e , the main difference between \mathcal{L}_{tr_i} in the toy example and InfoNCE is that the latter exploits a batch of negative samples instead of a random one. Since r_e has no center vector, we believe that the effect of r_e for preventing collapse is caused by dimensional de-correlation.

λ_i is proportional to $\exp(\mathbf{Z}_a \cdot \mathbf{Z}_i)$, indicating a large weight is put on the negative sample that is more similar to the anchor \mathbf{Z}_a . When a certain negative sample is more similar to the anchor, its dimensional values tend to have a high correlation with \mathbf{Z}_a . Intuitively, a high weight λ_i being put on such negative representation tends to decrease dimensional correlation. To verify this intuition, we measure the cosine similarity between r_e and the gradient on \mathbf{Z}_a induced by a correlation regularization loss. The results are shown in Fig. 5 (c). The gradient similarity is positive for a wide range of temperature value τ . When τ is very high, the λ_i tends to be equal to each other and thus the decorrelation effect vanishes. Specifically, when τ is sufficiently large, $\sum_{i=1}^N \lambda_i r_i$ is equivalent to $\sum_{i=1}^N \frac{1}{N} r_i$, *i.e.* no de-correlation. Empirically, we find that when τ is around 0.2, r_e achieves high similarity with the gradient induced by the correlation regularization loss. More analysis about the influence of the temperature are shown in Appendix A.7.

Dimensional de-correlation in SimSiam.

Assuming h only has a single FC layer to exclude the influence of o_e , the weights in FC are expected to learn the correlation between different dimensions for the encoder output. This interpretation aligns well with the finding that the eigenspace of h weight align well with that of correlation matrix (Tian et al., 2021). In essence, the h is trained to minimize the cosine similarity between $h(z_a)$ and $I(z_b)$, where I is identity mapping. Thus, h is optimized close to I , which is conceptually equivalent to optimizing the de-correlation for \mathbf{Z} . As shown in Table 4, r_e alone for SimSiam also prevents collapse. Since r_e has no de-centering effect, the reason that prevents is collapse is attributed to the de-correlation effect. From Figure 6, we observe that both SimSiam and InfoNCE decrease the covariance value during training. For completeness, we also report the results of removing r_e in InfoNCE. Due to removing r_e , the covariance value does not decrease during the training, which also leads to a significant performance drop compared with those with such de-correlation.

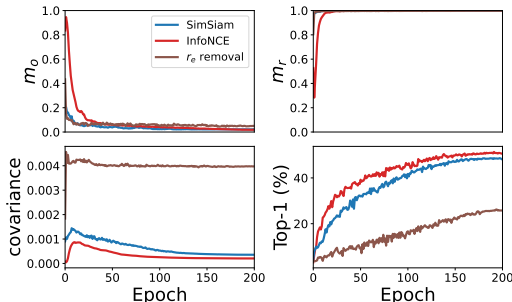


Figure 6: Influence of various gradient components on m_r and m_o .

3.5 TOWARDS A UNIFIED UNDERSTANDING OF RECENT PROGRESS IN SSL

Unifying InfoNCE and SimSiam. At first sight, there is no conceptual similarity between InfoNCE and SimSiam, and this is why the community is intrigued by the success of SimSiam without negative samples. Through decomposing the G_e into o_e and r_e , we find that for both, their o_e plays the role of de-centering and their r_e behaves like de-correlation. In this sense, we bring the two seemingly irrelevant frameworks into a unified perspective in the sense of de-centering and de-correlation.

Beyond InfoNCE and SimSiam. A literature survey on the recent progress, in SSL shows that there is a trend of performing *explicit* manipulation of de-centering and de-correlation. Three representative works are W-MSE (Ermolov et al., 2021), Barlow-twins (Zbontar et al., 2021), DINO (Caron et al., 2021) and that they often achieve performance comparable to those with InfoNCE or SimSiam. Towards a unified understanding of recent progress in SSL, our work is most similar to a concurrent work Bardes et al. (2021). Their work is mainly inspired by Barlow-twins (Zbontar et al., 2021) but decomposes them into three components in a more explicit manner. On the contrary, our work is motivated to answer the question how SimSiam prevents collapse without negative samples. The take-way is also different in the sense that we find that de-correlation itself is sufficient for preventing collapse, while their work claims that variance component (equivalent to de-centering) is an indispensable component for preventing collapse. Overall, through demystifying InfoNCE and SimSiam, our work bridges the understanding gap between those frameworks that implicitly and explicitly exploit de-centering and de-correlation for achieving competitive performance. We believe

that their marginal performance gap can be often ascribed to hyper-parameter tuning effect because they in essence achieve the same de-centering and de-correlation goals.

4 TOWARDS SIMPLE AND EXPLAINABLE PREDICTORS

Two important details: BN and l_2 -normalization. Regarding the reason for SimSiam’s success without negative samples, to our knowledge, our work is the first to revisit and refute the explanatory claims in (Chen & He, 2021). Several works, however, have attempted to demystify the success of BYOL Grill et al. (2020), a close variant of SimSiam. The success has been ascribed to BN in (Fetterman & Albrecht, 2020), however, Richemond et al. (2020) refutes the claim and (Chen & He, 2021) also claims that there is no evidence that BN itself is beneficial for preventing collapse. Here, we point out a critical implementation detail that determines whether BN helps prevent collapse. Since the the role of intermediate BNs is ascribed to stabilize training Richemond et al. (2020); Chen & He (2021), we only investigate the final BN in the encoder. Specifically, BN alone, *i.e.* using the Naive Siamese as shown in Fig 2, prevents collapse only if it is applied immediately before the MSE loss. For example, the cosine loss adopts l_2 -normalization before the MSE loss, for which the BN needs to be applied after the l_2 normalization. The rationale is that, as our conjecture 1 suggests, de-centering is the key for preventing collapse. In other words, BN that removes the mean is beneficial for preventing collapse. However after the l_2 normalization, a new mean is generated. Moreover, we empirically find that unnormalized MSE loss works poorly for SimSiam, thus our further investigation still adopts the default cosine loss.

SimSiam with Explainable predictor. Based on our understanding of how SimSiam prevents collapse, we attempt to demonstrate that simple components (**instead of a non-linear MLP**) in the predictor are sufficient enough for preventing collapse. For example, to achieve de-centering, ideally a single bias layer would be sufficient because a bias vector can represent the center vector. On the other hand, to achieve dimensional de-correlation, a single FC layer might be sufficient because a single FC layer can realize the interaction among various dimensions. With the current SimSiam encoder, we find that the above goals cannot be achieved. However, attaching a l_2 normalization layer at the end of the encoder, *i.e.* the input of the predictor facilitates the investigation.

Bias	(1) single bias	(2) bias in MLP
Similarity	0.99	0.89

Table 5: Similarity between *center vector* and (1) *single learnable bias layer*, (2) *the last bias vector* of MLP predictor.

A predictor with a single bias layer can be trained, where the bias vector is found to a very high cosine similarity with the center vector (see Table 5). A theoretical derivation justifying the 0.99 cosine similarity is reported in Appendix A.8. A bias in the MLP predictor also has high a cosine similarity of 0.89. A single bias layer is found to be sufficient for preventing collapse (see Table 6). To learn the dimensional correlation, theoretically, a FC layer is sufficient. In practice, we find that a single FC layer can be difficult to train. Inspired by Bell-Kligler et al. (2019), we adopt two consecutive FC layers even though two FC layers are mathematically equivalent to a single FC layer. Multiple FC layers have the property to make the training more stable. This implementation is equivalent to removing the BN and ReLU in the original predictor. Moreover, to train a single FC layer, we apply a Tanh layer to constrain the FC weight values after every iteration to make the training more stable. They achieve comparable performance as original SimSiam (see Table 6).

Method	Predictor	Top-1 (%)
SimSiam	MLP	66.9
One bias	Bias	49.82
One FC	Tanh(FC)	64.82
Two FC	FC+FC+Bias	66.7

Table 6: Linear evaluation on CIFAR100.

5 CONCLUSION

Our work has investigated how SimSiam prevents collapse without negative samples. Our investigation starts with revisiting the explanatory claims in SimSiam and point out a hidden flaw in their reasoning. After refuting their claims, our work proposes to decompose the representation vector and analyze the decomposed gradients. We find that for both InfoNCE and SimSiam, the center vector gradient has the de-centering effect and the residual gradient vector has the de-correlation effect. Both de-centering and de-correlation effects are found to be beneficial for preventing collapse. We also investigate SimSiam with more explainable predictors.

REFERENCES

- Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. *arXiv preprint arXiv:1906.00910*, 2019.
- Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- Sefi Bell-Kligler, Assaf Shocher, and Michal Irani. Blind super-resolution kernel estimation using an internal-gan. *NeurIPS*, 2019.
- Jane Bromley, James W Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. Signature verification using a “siamese” time delay neural network. *International Journal of Pattern Recognition and Artificial Intelligence*, 1993.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. *arXiv preprint arXiv:2104.14294*, 2021.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020a.
- Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, 2021.
- Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.
- Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. *ICCV*, 2021.
- Victor G. Turrissi da Costa, Enrico Fini, Moin Nabi, Nicu Sebe, and Elisa Ricci. Solo-learn: A library of self-supervised methods for visual representation learning, 2021.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423.
- Alaaeldin El-Nouby, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. Xcit: Cross-covariance image transformers. *arXiv preprint arXiv:2106.09681*, 2021.
- Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-supervised representation learning. In *ICML*. PMLR, 2021.
- Abe Fetterman and Josh Albrecht. Understanding self-supervised and contrastive learning with “bootstrap your own latent” (byol), 2020.
<https://untitled-ai.github.io/understanding-self-supervised-contrastive-learning.html>.
- Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, 2019.

- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *arXiv preprint arXiv:2004.11362*, 2020.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=H1eA7AEtvS>.
- Chunyu Li, Jianwei Yang, Pengchuan Zhang, Mei Gao, Bin Xiao, Xiyang Dai, Lu Yuan, and Jianfeng Gao. Efficient self-supervised vision transformers for representation learning. *arXiv preprint arXiv:2106.09785*, 2021.
- Ping Nie, Yuyu Zhang, Xiubo Geng, Arun Ramamurthy, Le Song, and Daxin Jiang. Dc-bert: Decoupling question and document for efficient contextual encoding. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 1829–1832, 2020.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
- Pierre H Richemond, Jean-Bastien Grill, Florent Alché, Corentin Tallec, Florian Strub, Andrew Brock, Samuel Smith, Soham De, Razvan Pascanu, Bilal Piot, et al. Byol works even without batch statistics. *arXiv preprint arXiv:2010.10241*, 2020.
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 815–823, 2015.
- Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NeurIPS*, 2016.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. {VL}-{bert}: Pre-training of generic visual-linguistic representations. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SygXPaeYvH>.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.
- Yuandong Tian, Xinlei Chen, and Surya Ganguli. Understanding self-supervised learning dynamics without contrastive pairs. *arXiv preprint arXiv:2102.06810*, 2021.
- Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, 2018.
- Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. *ICML*, 2021.

A APPENDIX

A.1 EXPERIMENTAL SETTINGS

Unsupervised pre-training: Below are settings for self-supervised pre-training stage. All experiments setting are used as default from the solo-learn library (da Costa et al., 2021).

Data augmentation: We use the series of transformations as follows, notations are in Pytorch style: *RandomResizedCrop* with scale [0.2, 1.0], interpolation is bicubic. Randomly apply with probability of 0.8 for *ColorJitter* (brightness (0.4), contrast (0.4), saturation (0.4), hue (0.1)) with strength accordingly. Random gray scale *RandomGrayscale* apply with $p = 0.2$ Random apply horizontal flip with $p = 0.5$ Normalize image to mean and std as (0.4914, 0.4822, 0.4465) and (0.247, 0.243, 0.261), respectively.

Initialization: All layers of the network are initialized by default of Pytorch.

Network architecture: Network backbone **B** is Resnet-18. Projection MLP head contains three fully-connected (FC) layers follow by Batch Norm (BN) and ReLU. We remove ReLU in the final FC layer: $FC_1 + BN + ReLU + FC_2 + BN + ReLU + FC_3 + BN$. All FC layers have 2048 neurons for both input, output and hidden dimensions. Predictor MLP head includes two FC layers as follows: $FC_1 + BN + ReLU + FC_2$. Input and output of predictor both have dimension of 2048, while hidden dimension is 512.

Optimizer: SGD optimizer is used for unsupervised pre-training. Batch size is 256 and learning rate is linearly scaling by formula $lr \times Bsz/256$ from the base learning rate of 0.5. The schedule for learning rate is cosine decay as SimSiam. Momentum 0.9 and weight decay 1.0×10^{-5} are used for SGD. We use one GPU for each pre-training experiment. The learning rate of predictor is fixed (not cosine decay) same as SimSiam trick. We use a warmup training for first 10-epochs. If not specified, by default we train the model for 1000 epochs.

Online linear evaluation: The frozen features (2048 dimensions) from the training set are extracted (from the self-supervised pre-trained model) to feed into train a supervised linear classifier (1 FC layer with input 2048 and output is 10 classes). The test is done on the validation set. Learning rate for linear classifier is 0.1. All online evaluation reports are done using the solo-learn library for self-supervised learning (da Costa et al., 2021).

A.2 TWO SUB-PROBLEMS IN AO OF SIMSIAM

In the sub-problem $\eta^t \leftarrow \arg \min_{\eta} \mathcal{L}(\theta^t, \eta)$, η^t indicating latent representation of images at step t is actually obtained through $\eta_x^t \leftarrow \mathbb{E}_{\mathcal{T}} \left[\mathcal{F}_{\theta^t}(\mathcal{T}(x)) \right]$, where they in practice ignore $\mathbb{E}_{\mathcal{T}}[\cdot]$ and sample only one augmentation \mathcal{T}' , i.e. $\eta_x^t \leftarrow \mathcal{F}_{\theta^t}(\mathcal{T}'(x))$. Conceptually, Chen & He equate the role of predictor to EOA.

A.3 EXPERIMENTAL DETAILS FOR EXPLICIT EOA IN TABLE 1

In the *Moving average* experiment, we follow the setting in SimSiam Chen & He (2021) without predictor. In contrast to Chen & He (2021) that uses the representations of the same image from prior epochs, we perform EOA in the same mini-batch. In the *Same batch* experiment, multiple augmentations, 10 augmentations for instance, are applied on the same image. With multi augmentations, we get the corresponding encoded representation, i.e. $z_i, i \in [1, 10]$. We minimize the cosine distance between the first representation z_1 and the average of the remaining vectors, i.e. $\bar{z} = \frac{1}{9} \sum_{i=2}^{10} z_i$. Gradient detach is applied on the average vector.

A.4 EXPERIMENTAL SETUP AND RESULT TREND FOR TABLE 2.

Here we provide the pseudocode for mirror SimSiam, Symmetric Predictor. In the Mirror SimSiam experiment which relates to Fig. 1 (c), we stop the gradient in the input of predictor. Then we attract the predictor with non-stop gradient projector output. Without taking symmetric loss into account,

the pseudocode is shown in Algorithm 1. Taking symmetric loss into account, the pseudocode is shown in Algorithm 2.

Algorithm 1 Pytorch-like Pseudocode: Mirror SimSiam

```
# f: encoder (backbone + projector)
# h: predictor

for x in loader: # load a minibatch x with n samples
    x_a, x_b = aug(x), aug(x) # augmentation
    z_a, z_b = f(x_a), f(x_b) # projections

    p_b = h(z_b.detach()) # detach z_b but still allowing gradient p_b

    L = D_cosine(z_a, p_b) # loss

    L.backward() # back-propagate
    update(f, h) # SGD update

def D_cosine(z, p): # negative cosine similarity
    z = normalize(z, dim=1) # l2-normalize
    p = normalize(p, dim=1) # l2-normalize
    return -(z*p).sum(dim=1).mean()
```

Algorithm 2 Pytorch-like Pseudocode: Mirror SimSiam

```
# f: encoder (backbone + projector)
# h: predictor

for x in loader: # load a minibatch x with n samples
    x_a, x_b = aug(x), aug(x) # augmentation
    z_a, z_b = f(x_a), f(x_b) # projections

    p_b = h(z_b.detach()) # detach z_b but still allowing gradient p_b
    p_a = h(z_a.detach()) # detach z_a but still allowing gradient p_a

    L = D_cosine(z_a, p_b)/2 + D_cosine(z_b, p_a)/2 # loss

    L.backward() # back-propagate
    update(f, h) # SGD update

def D_cosine(z, p): # negative cosine similarity
    z = normalize(z, dim=1) # l2-normalize
    p = normalize(p, dim=1) # l2-normalize
    return -(z*p).sum(dim=1).mean()
```

To implement the training dynamic of Fig. 2 (b), the pseudocode in Algorithm 3.

Algorithm 3 Pytorch-like Pseudocode: Symmetric Predictor

```
# f: encoder (backbone + projector)
# h: predictor

for x in loader: # load a minibatch x with n samples
    x_a, x_b = aug(x), aug(x) # augmentation
    z_a, z_b = f(x_a), f(x_b) # projections
    p_a, p_b = h(z_a), h(z_b) # predictions

    L = D(p_a, p_b)/2 + D(p_b, p_a)/2 # loss

    L.backward() # back-propagate
    update(f, h) # SGD update

def D(p, z): # negative cosine similarity
    z = z.detach() # stop gradient
    p = normalize(p, dim=1) # l2-normalize
    z = normalize(z, dim=1) # l2-normalize
    return -(p*z).sum(dim=1).mean()
```

The results of SimSiam, Naive Siamese, Mirror SimSiam, Symmetric Predictor are shown in Fig. 7. We observe that all models collapse except for SimSiam.

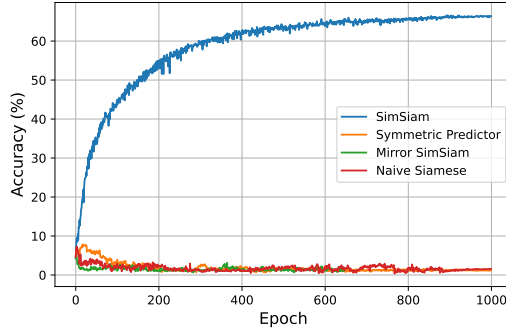


Figure 7: Results of Naive Siamese, Mirror SimSiam, Symmetric Predictor

A.5 EXPERIMENTAL DETAILS FOR INVERSE PREDICTOR.

In the inverse predictor experiment which relates to Fig. 2 (c), we introduce a new predictor which has the same structure as that of original predictor. The training loss consists of 3 parts: predictor training loss, inverse predictor training and new encoder (old encoder+predictor) training. The new encoder consists of old encoder + predictor. Proper detach needs to be put in the implementation. We provide the pseudocode in Algorithm 4.

Algorithm 4 Pytorch-like Pseudocode: Train Inverse Predictor

```

# f: encoder (backbone + projector)
# h: predictor
# h_inv: inverse predictor

for x in loader: # load a minibatch x with n samples
    x_a, x_b = aug(x), aug(x) # augmentation
    z_a, z_b = f(x_a), f(x_b) # projections
    p_a, p_b = h(z_a), h(z_b) # predictions

    d_p_a, d_p_b = h(z_a.detach()), h(z_b.detach()) # detached predictor output
    # predictor training loss
    L_pred = D(d_p_a, z_b)/2 + D(d_p_b, z_a)/2 # to train h

    inv_p_a, inv_p_b = h_inv(p_a.detach()), h_inv(p_b.detach()) # to train h_inv
    # inverse predictor training loss
    L_inv_pred = D(inv_p_a, z_a)/2 + D(inv_p_b, z_b)/2

    # encoder training loss
    L_enc = D(p_a, h_inv(p_b))/2 + D(p_b, h_inv(p_a))

    L = L_pred + L_inv_pred + L_enc

    L.backward() # back-propagate
    update(f, h, h_inv) # SGD update

def D(p, z): # negative cosine similarity with detach on z
    z = z.detach() # stop gradient
    p = normalize(p, dim=1) # l2-normalize
    z = normalize(z, dim=1) # l2-normalize
    return -(p*z).sum(dim=1).mean()

```

A.6 REGULARIZATION LOSS IN CONJECTURE 2

We compute auto-correlation matrix of encoder output along the mini-batch. Recently, Dimensional Correlation in SSL has attached great attention to prevent collapse (Zbontar et al., 2021; Ermolov et al., 2021). Following the (Zbontar et al., 2021), we define \mathcal{C} as the dimensional correlation matrix computed among the representations in a batch. Suppose the dimension of representation vector of a single sample Z is d , \mathcal{C} is matrix with the dimension of $d \times d$:

$$\mathcal{C}_{ij} = \frac{\sum_b Z_{b,i}^A Z_{b,j}^B}{\sqrt{\sum_b (Z_{b,i}^A)^2} \sqrt{\sum_b (Z_{b,j}^B)^2}} \quad (5)$$

where b indicates the sample index in a batch and i, j indicates dimension indexes in Z . Here we use $Z_{b,i}^A = Z_{b,j}^B = Z_1$. The correlation loss calculation in pytorch style as follows:

```
corr = torch.einsum("bi, bj -> ij", z_a, z_a) / N
diag = torch.eye(D, device=corr.device)
cdf = (corr - diag).pow(2)
loss = scale_loss * cdf.sum()
```

A.7 DERIVATIVE OF INFO NCE

Let S be the cosine similarity between Z_a and Z_b . The InfoNCE loss function can be expressed as follows, when cosine similarity are used,

$$\begin{aligned} \mathcal{L}_{InfoNCE} &= -\log \frac{\exp(\mathbf{Z}_a \cdot \mathbf{Z}_b / \tau)}{\exp(\mathbf{Z}_a \cdot \mathbf{Z}_b / \tau) + \sum_{i=1}^N \exp(\mathbf{Z}_a \cdot \mathbf{Z}_i / \tau)} \\ &= -\log \frac{\exp(\mathbf{Z}_a \cdot \mathbf{Z}_b / \tau)}{\sum_{i=0}^N \exp(\mathbf{Z}_a \cdot \mathbf{Z}_i / \tau)}, \end{aligned} \quad (6)$$

where N indicates the number of negative samples and $Z_0 = Z_b$ for simplifying the notation. By treating $Z_a \cdot Z_i$ as the logit in the normal CE loss, we have the corresponding probability for each logit as $\lambda_i = \frac{\exp(\mathbf{Z}_a \cdot \mathbf{Z}_i / \tau)}{\sum_{i=0}^N \exp(\mathbf{Z}_a \cdot \mathbf{Z}_i / \tau)}$, where $i = 0, 1, 2, \dots, N$ and we have $\sum_{i=0}^N \lambda_i = 1$.

The negative gradient of representation Z_a with regard to InfoNCE loss is shown as follows:

$$\begin{aligned} -\frac{\partial \mathcal{L}_{InfoNCE}}{\partial \mathbf{Z}_a} &= \frac{1}{\tau} (1 - \lambda_0) \mathbf{Z}_b - \frac{1}{\tau} \sum_{i=1}^N \lambda_i \mathbf{Z}_i \\ &= \frac{1}{\tau} (\mathbf{Z}_b - \sum_{i=0}^N \lambda_i \mathbf{Z}_i) \\ &= \frac{1}{\tau} (\mathbf{Z}_b - \sum_{i=0}^N \lambda_i (\mathbf{o}_z + \mathbf{r}_i)) \\ &= \frac{1}{\tau} (\mathbf{Z}_b + (-\mathbf{o}_z - \sum_{i=0}^N \lambda_i \mathbf{r}_i)) \\ &\propto \mathbf{Z}_b + (-\mathbf{o}_z - \sum_{i=0}^N \lambda_i \mathbf{r}_i) \end{aligned} \quad (7)$$

where $\frac{1}{\tau}$ can be adjusted through learning rate and is omitted for simple discussion. With Z_b as the basic symmetric gradient, $\mathbf{G}_e = -\mathbf{o}_z - \sum_{i=0}^N \lambda_i \mathbf{r}_i$, for which $\mathbf{o}_e = -\mathbf{o}_z$ and $\mathbf{r}_e = -\sum_{i=0}^N \lambda_i \mathbf{r}_i$.

When the temperature is set to a large value, $\lambda_i = \frac{\exp(\mathbf{Z}_a \cdot \mathbf{Z}_i / \tau)}{\sum_{i=0}^N \exp(\mathbf{Z}_a \cdot \mathbf{Z}_i / \tau)}$, tends to be similar to each other for various i and InfoNCE will degenerate to a simple contrastive loss, *i.e.* $\mathcal{L}_{simple} = -\mathbf{Z}_a \cdot \mathbf{Z}_b + \frac{1}{N} \sum \mathbf{Z}_a \cdot \mathbf{Z}_n$, which repulse every negative sample with an equal force. In contrast, a relative smaller temperature will give more weight, *i.e.* larger λ , to samples that are more similar to the anchor.

Hypothesis: The biased repulsive weight on each negative sample in the gradient leads to a de-correlation effect.

The influence of the temperature on the covariance and accuracy is shown in Fig 8. We observe that a higher temperature tends to decrease effect of de-correlation, indicated by higher covariance value, which also leads to a performance drop. This verifies the above hypothesis because a large temperature causes more balanced weights *et al.*, which is found to alleviate the effect of de-correlation. The model is trained with 200 epochs with the default setting in Solo-learn for SimCLR framework.

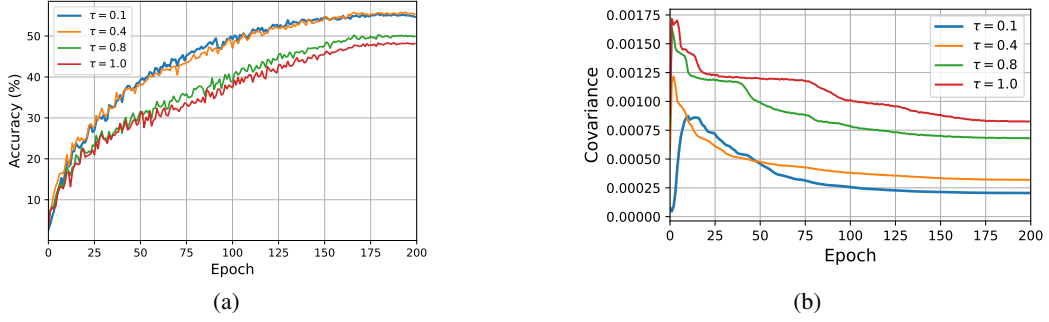


Figure 8: Influence of temperature. (a) Top-accuracy; (b) Covariance value.

A.8 THEORETICAL DERIVATION FOR A SINGLE BIAS LAYER

The cosine similarity loss is defined as Eq 8, and the derivative is presented as Eq 9:

$$\text{cossim}(a, b) = \frac{a \cdot b}{\sqrt{a^2 \cdot b^2}} \quad (8)$$

$$\frac{\partial}{\partial a} \text{cossim}(a, b) = \frac{b_1}{|a| \cdot |b|} - \text{cossim}(a, b) \cdot \frac{a_1}{|a|^2} \quad (9)$$

The above equation is used a prior for our following derivations. The single bias layer in the predictor is denoted as \mathbf{b} . $\mathbf{p}_a = \mathbf{Z}_a + \mathbf{b}$ (note that \mathbf{Z}_a is l_2 -normalized as mentioned in the main manuscript) and $\mathbf{P}_a = \mathbf{p}_a / \|\mathbf{p}_a\|$. The cosine similarity loss is shown as

$$\begin{aligned} \mathcal{L}_{\text{cosine}} &= -\mathbf{P}_a \cdot \mathbf{Z}_b \\ &= -\frac{\mathbf{p}_a}{\|\mathbf{p}_a\|} \cdot \frac{\mathbf{z}_b}{\|\mathbf{z}_b\|} \end{aligned} \quad (10)$$

The gradient on the bias vector is derived as

$$\begin{aligned} -\frac{\partial \mathcal{L}_{\text{cosine}}}{\partial \mathbf{b}} &= -\frac{\partial \mathcal{L}_{\text{cosine}}}{\partial \mathbf{p}_a} \\ &= \frac{\mathbf{z}_b}{\|\mathbf{z}_b\| \cdot \|\mathbf{p}_a\|} - \text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b) \cdot \frac{\mathbf{p}_a}{\|\mathbf{p}_a\|^2} \\ &= \frac{1}{\|\mathbf{p}_a\|} \left(\frac{\mathbf{z}_b}{\|\mathbf{z}_b\|} - \text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b) \cdot \mathbf{P}_a \right) \\ &= \frac{1}{\|\mathbf{p}_a\|} \left(\mathbf{Z}_b - \text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b) \cdot \frac{\mathbf{Z}_a + \mathbf{b}}{\|\mathbf{p}_a\|} \right) \\ &= \frac{1}{\|\mathbf{p}_a\|} \left((\mathbf{o}_z + \mathbf{r}_b) - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|} \cdot (\mathbf{o}_z + \mathbf{r}_a + \mathbf{b}) \right) \\ &= \frac{1}{\|\mathbf{p}_a\|} ((\mathbf{o}_z + \mathbf{r}_b) - m \cdot (\mathbf{o}_z + \mathbf{r}_a + \mathbf{b})) \\ &= \frac{1}{\|\mathbf{p}_a\|} ((1 - m)\mathbf{o}_z - m\mathbf{b} + \mathbf{r}_b - m \cdot \mathbf{r}_a), \end{aligned} \quad (11)$$

where $m = \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|}$. Assuming that the training is stable and the bias vector converges to a certain value when $-\frac{\partial \text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\partial \mathbf{b}} = 0$. Thus, the converged \mathbf{b} satisfies the following constraint:

$$\begin{aligned} \frac{1}{\|\mathbf{p}_a\|} ((1 - m)\mathbf{o}_z - m\mathbf{b} + \mathbf{r}_b - m\mathbf{r}_a) &= 0 \\ \mathbf{b} &= \frac{1 - m}{m} \mathbf{o}_z + \frac{1}{m} \mathbf{r}_b - \mathbf{r}_a. \end{aligned} \quad (12)$$

With a batch of samples, the average of $(\mathbf{r}_b - \mathbf{r}_a)$ is expected to be close to 0. Thus, the bias layer vector is expected to converge to:

$$\mathbf{b} = \frac{1-m}{m} \mathbf{o}_z. \quad (13)$$

In practice, with the training dynamics, \mathbf{b} is expected to fluctuate around $\frac{1-m}{m} \mathbf{o}_z$. Nonetheless, we empirically measure the cosine similarity between \mathbf{b} and \mathbf{o}_z and find they are around 0.99.

We further derive the gradient on \mathbf{Z}_a as follows:

$$\begin{aligned} -\frac{\partial \mathcal{L}_{\text{cosine}}}{\partial \mathbf{Z}_a} &= -\frac{\partial \mathcal{L}_{\text{cosine}}}{\partial \mathbf{p}_a} \\ &= \frac{1}{\|\mathbf{p}_a\|} \left(\mathbf{Z}_b - \text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b) \cdot \frac{\mathbf{Z}_a + \mathbf{b}}{\|\mathbf{p}_a\|} \right) \\ &= \frac{1}{\|\mathbf{p}_a\|} \mathbf{Z}_b - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|^2} \mathbf{Z}_a - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|^2} \mathbf{b}. \end{aligned} \quad (14)$$

Here, for the second part of the gradient, since the loss $-\mathbf{Z}_a \cdot \mathbf{Z}_a = -1$ is a constant, this second part is a *dummy* term. Considering Eq 13 and the value of m , we have $\mathbf{b} = (\frac{\|\mathbf{p}_a\|}{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)} - 1) \mathbf{o}_z$. Then the above equation is equivalent to

$$\begin{aligned} -\frac{\partial \mathcal{L}_{\text{cosine}}}{\partial \mathbf{Z}_a} &= \frac{1}{\|\mathbf{p}_a\|} \mathbf{Z}_b - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|^2} \mathbf{b} \\ &= \frac{1}{\|\mathbf{p}_a\|} \mathbf{Z}_b - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|^2} \left(\frac{\|\mathbf{p}_a\|}{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)} - 1 \right) \mathbf{o}_z \\ &= \frac{1}{\|\mathbf{p}_a\|} \mathbf{Z}_b - \frac{1}{\|\mathbf{p}_a\|} \left(1 - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|} \right) \mathbf{o}_z \\ &\propto \mathbf{Z}_b - \left(1 - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|} \right) \mathbf{o}_z. \end{aligned} \quad (15)$$

With \mathbf{Z}_b as the basic symmetric gradient, the extra gradient component $\mathbf{G}_e = -(1 - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|}) \mathbf{o}_z$, $\mathbf{p}_a = \mathbf{Z}_a + \mathbf{b}$ and $\|\mathbf{Z}_a\| = 1$, thus $\|\mathbf{p}_a\| < 1$ only when \mathbf{Z}_a is negatively correlated with \mathbf{b} . In practice, however, \mathbf{Z}_a and \mathbf{b} are often positively correlated to some extent due to their shared center vector component. In other words, $\|\mathbf{p}_a\| > 1$. Moreover, $\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)$ is smaller than 1, thus $-(1 - \frac{\text{cossim}(\mathbf{Z}_a, \mathbf{Z}_b)}{\|\mathbf{p}_a\|}) < 0$, suggesting \mathbf{G}_e consists of negative \mathbf{o}_z with the effect of de-centerization. This derivation shows why a single bias layer can help avoid collapse.