# Tuned Contrastive Learning

**Chaitanya Animesh** [1,2*]    **Manmohan Chandraker**[1]

[1]UC San Diego    [2]Otter.ai, Inc.

canimesh@ucsd.edu    mkchandraker@ucsd.edu

## Abstract

*In recent times, contrastive learning based loss functions have become increasingly popular for visual self-supervised representation learning owing to their state-of-the-art (SOTA) performance. Most of the modern contrastive learning methods generalize only to one positive and multiple negatives per anchor in a batch. A recent state-of-the-art contrastive loss called supervised contrastive (SupCon) loss, extends self-supervised contrastive learning to supervised setting by generalizing to multiple positives and negatives in a batch and improves upon the cross-entropy loss. In this paper, we propose a novel contrastive loss function — Tuned Contrastive Learning (TCL) loss, that generalizes to multiple positives and negatives in a batch and offers parameters to tune and improve the gradient responses from hard positives and hard negatives. We provide theoretical analysis of our loss function's gradient response and show mathematically how it is better than that of SupCon loss. We empirically compare our loss function with SupCon loss and cross-entropy loss in supervised setting on multiple classification-task datasets to show its effectiveness. We also show the stability of our loss function to a range of hyper-parameter settings. Unlike SupCon loss which is only applied to supervised setting, we show how to extend TCL to self-supervised setting and empirically compare it with various SOTA self-supervised learning methods. Hence, we show that TCL loss achieves performance on par with SOTA methods in both supervised and self-supervised settings.*

## 1. Introduction

Paucity of labeled data limits the application of supervised learning to various visual learning tasks [35]. As a result, unsupervised [17, 19, 26] and self-supervised based learning methods [4, 6, 18, 32] have garnered a lot of attention and popularity for their ability to learn from vast unlabeled data. Such methods can be broadly classified into two categories: generative methods and discriminative methods.
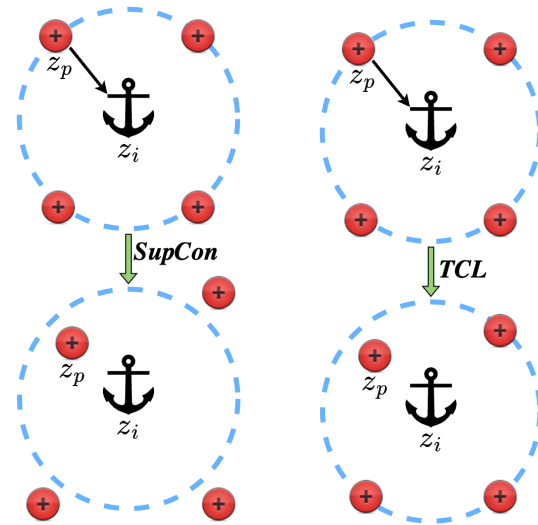


Figure 1. Figure illustrates intuitively how TCL loss differs from SupCon loss [23]. For the terms inside the summation of $L_i^{sup}$ (from Eq. (2)) in the SupCon loss to decrease, the anchor $z_i$ will pull the positive $z_p$ but push away the other positives to some extent in the embedding space. TCL loss introduces parameters that reduce this effect and increase the gradient response from positives. This leads to consistently better performance.

Generative methods [17, 26] train deep neural networks to generate in the input space i.e. the pixel space and hence, are computationally expensive and not necessary for representation learning. On the other hand, discriminative approaches [1, 6, 13, 16, 30, 33] train deep neural networks to learn representations for pretext tasks using unlabeled data and an objective function. Out of these discriminative based approaches, contrastive learning based methods [1, 6, 30] have performed significantly well and are an active area of research.

The common principle of contrastive learning based methods in an unsupervised setting is to create semantic preserving transformations of the input which are called *positives* and treat transformations of other samples in a

---

*Work done as part of Master's Thesis at UC San Diego.

[0]Pytorch implementation of our work released and available at: https://github.com/chaitanyaanimesh/Tuned-Contrastive-Learning.

batch as *negatives* [2, 23]. The contrastive loss objective considers every transformed sample as a reference sample, called an *anchor*, and is then used to train the network architecture to pull the positives (for that anchor) closer to the anchor and push the negatives away from the anchor in latent space [2, 23]. The positives are often created using various data augmentation strategies. Supervised Contrastive Learning [23] extended contrastive learning to supervised setting by using the label information and treating the other samples in the batch having the same label as that of the anchor also as positives in addition to the ones produced through data augmentation strategies. It presents a new loss called supervised contrastive loss (abbreviated as SupCon loss) that can be viewed as a loss generalizing to multiple positives available in a batch.

In this work, we propose a novel contrastive learning loss objective, which we call **Tuned Contrastive Learning (TCL) Loss** that can use multiple positives and multiple negatives present in a batch. We show how it can be used in supervised as well as self-supervised settings. TCL loss improves upon the limitations of the SupCon loss: 1. Implicit consideration of positives as negatives and, 2. No provision of regulating hard negative gradient response. TCL loss thus gives better gradient response to hard positives and hard negatives. This leads to small ($\leq 1\%$ in terms of classification accuracy) but consistent improvements in performance over SupCon loss and comprehensive outperformance over cross-entropy loss. Since TCL generalizes to multiple positives, we then present a novel idea of having and using positive triplets (and possibly more) instead of being limited to positive pairs for self-supervised learning. We evaluate our loss function in self-supervised settings without making use of any label information and show how TCL outperforms SimCLR [6] and performs on par with various SOTA self-supervised learning methods [3, 4, 8, 9, 15, 18, 20, 32, 35]. Our key contributions in the paper are as follows:

1. We identify and analyse in detail two limitations of the supervised contrastive (SupCon) loss.

2. We present a novel contrastive loss function called Tuned Contrastive Learning (TCL) loss that generalizes to multiple positives and multiple negatives in a batch, overcomes the described limitations of the SupCon loss and is applicable in both supervised and self-supervised settings. We mathematically show with clear proofs how our loss's gradient response is better than that of SupCon loss.

3. We compare TCL loss with SupCon loss (as well as cross-entropy loss) in supervised settings on various classification datasets and show that TCL loss gives consistent improvements in top-1 accuracy over SupCon loss. We empirically show the stability of TCL

loss to a range of hyperparameters: network architecture, batch size, projector size and augmentation strategy.

4. At last, we present a novel idea of having positive triplets (and possibly more) instead of positive pairs and show how TCL can be extended to self-supervised settings. We empirically show that TCL outperforms SimCLR, and performs on par with various SOTA self-supervised learning (SSL) methods.

## 2. Related Work

In this section, we cover various popular and recent works in brief involving contrastive learning.

Deep Metric learning methods originated with the idea of contrastive losses and were introduced with the goal of learning a distance metric between samples in a high-dimensional space [2]. The goal in such methods is to learn a function that maps similar samples to nearby points in this space, and dissimilar samples to distant points. There is often a margin parameter, $m$, imposing the distance between examples from different classes to be larger than this value of $m$ [2]. The triplet loss [22] and the proposed improvements [7, 27] on it used this principle. These methods rely heavily on sophisticated sampling techniques for choosing samples in every batch for better training.

SimCLR [6], an Info-NCE [30] loss based framework, learns visual representations by increasing the similarity between the embeddings of two augmented views of the input image. Augmented views generally come from a series of transformations like random resizing, cropping, color jittering, and random blurring. Although they make use of multiple negatives, only one positive is available per anchor. They require large batch sizes in order to have more hard negatives in the batch to learn from and boost the performance. SupCon loss [23] applies contrastive learning in supervised setting by basically extending the SimCLR loss to generalize to multiple positives available in a batch and improves upon the cross-entropy loss which lacks robustness to noisy labels [29, 34] and has the possibility of poor margins [14, 25].

Unlike SimCLR or SupCon, many SOTA SSL approaches only work with positives (don't require negatives) or use different approach altogether. BYOL [18] uses asymmetric networks with one network using an additional predictor module while the other using exponential moving average (EMA) to update its weights, in order to learn using positive pairs only and prevent collapse. SimSiam [9] uses stop-gradient operation instead of EMA and asymmetric networks to achieve the same goal. Barlow Twins [32] objective function on the other hand computes the cross-correlation matrix between the embeddings of two identical networks fed with augmentations of a batch of samples,

and tries to make this matrix close to identity. SwAV uses a clustering approach and enforces consistency between the cluster assignments of multiple positives produced through multi-crop strategy [4].

## 3. Methodology

### 3.1. Supervised Contrastive Learning & Its Issues

The framework for Supervised Contrastive Learning consists of three components: a data augmentation module that produces two augmentations for each sample in the batch, an encoder network that maps the augmentations to their corresponding representation vectors and a projection network that produces normalized embeddings for the representation vectors to be fed to the loss function. The projection network is later discarded and the encoder network is used at inference time by training a linear classifier (attached to the frozen encoder) with cross-entropy loss. Section 3.1 of [23] contains more details on this. The SupCon loss is given by the following two equations (refers to $L_{out}^{sup}$ in [23]):

$$L^{sup} = \sum_{i \in I} L_i^{sup} \tag{1}$$

where

$$L_i^{sup} = \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log\left(\frac{\exp(z_i.z_p/\tau)}{D_{z_i}^{sup}}\right) \tag{2}$$

and

$$D_{z_i}^{sup} = \sum_{p' \in P(i)} \exp(z_i.z_{p'}/\tau) + \sum_{n \in N(i)} \exp(z_i.z_n/\tau) \tag{3}$$

Here $I$ denotes the batch of samples obtained after augmentation and so, will be twice the size of the original input batch. $i \in I$ denotes a sample (anchor) within it. $z_i$ denotes the normalized projection network embedding for the sample $i$ as given by the projector network. $P(i)$ is the set of all positives for the anchor $i$ (except the anchor $i$ itself) i.e. positive from the augmentation module and positives with the same label as anchor $i$ in the batch $I$. $N(i)$ denotes the set of negatives in the batch such that $N(i) \equiv I \setminus (P(i) \cup \{i\})$. As shown in Section 2 of the supplementary material of [23], we have the following lemma:

**Lemma 1** *The gradient of the SupCon loss per sample — $L_i^{sup}$ with respect to the normalized projection network embedding $z_i$ is given by:*

$$\frac{\partial L_i^{sup}}{\partial z_i} = \frac{1}{\tau}\bigg(\underbrace{\sum_{p \in P(i)} z_p(P_{ip}^s - X_{ip})}_{\textit{Gradient from positives}} + \underbrace{\sum_{n \in N(i)} z_n P_{in}^s}_{\textit{Gradient from negatives}}\bigg) \tag{4}$$

*where*

$$X_{ip} = \frac{1}{|P(i)|} \tag{5}$$

$$P_{ip}^s = \frac{exp(z_i.z_p/\tau)}{\sum_{a \in A(i)} exp(z_i.z_a/\tau)} \tag{6}$$

$$P_{in}^s = \frac{exp(z_i.z_n/\tau)}{\sum_{a \in A(i)} exp(z_i.z_a/\tau)} \tag{7}$$

Note that $A(i) \equiv P(i) \cup N(i)$ here. The authors further show in Section 3 of the supplementary [23] that the gradient from a positive while flowing back through the projector into the encoder reduces to almost zero for easy positives and $|P_{ip}^s - X_{ip}|$ for a hard positive because of the normalization consideration in the projection network. Similarly, the gradient from a negative reduces to almost zero for easy negatives and $|P_{in}^s|$ for a hard negative. We now present and analyse the following two limitations of the SupCon loss:

1. **Implicit consideration of positives as negatives**: Having a closer look at the $L_i^{sup}$ (Eq. (2)) loss term reveals that each individual term of $L_i^{sup}$ inside the summation, consists of similarity terms of the anchor $i$ with all the positives in the batch — the set $P(i)$ — in the denominator, thereby implicitly considering all the positives except the positive $p$ in the numerator as negatives. Although $L_i^{sup}$ as a whole will consider all positives as positives indeed, its individual terms at their respective local levels implicitly consider the positives as negatives that leads to reduced gradient response from positives. A glance at the derivation of Lemma 1 in [23] clearly shows that this leads to the magnitude of the gradient response from a hard positive getting reduced to $|X_{ip} - P_{ip}^s|$ instead of simply $|X_{ip}|$. The term $P_{ip}^s$ consists of an exponential term in the numerator and thus can reduce the magnitude of $|X_{ip} - P_{ip}^s|$ considerably, especially because the temperature $\tau$ is generally chosen to be small. Note that the authors of [23] approximate the numerator of $P_{ip}^s$ to 1 while considering the magnitude of $|X_{ip} - P_{ip}^s|$ in their supplementary by assuming $z_i.z_p \approx 0$ for a hard positive which might not always be true. Another way to look at this limitation analytically is to observe the log part in the individual terms of $L_i^{sup}$. For them to decrease and ideally converge to close to zero, the numerator term inside the log function will encourage the anchor $z_i$ to pull the positive $z_p$ towards it while the denominator term will encourage it to push away the other positives present in $P(i)$ by some extent, thereby treating the other positives as negatives implicitly.

2. **No possibility of regulating $P_{in}^s$**: [6, 23] mention that performance in contrastive learning benefits from hard

negatives and gradient contribution from hard negatives should be higher. It is easy to observe from Eq. (7) that the magnitude of the gradient signal from a hard negative — $|P_{in}^s|$ in the SupCon loss decreases with batch size and the number of positives in the batch, and can become considerably small, especially since the denominator consists of similarity terms between the anchor and all the positives in the batch which are temperature scaled and exponentiated. This can limit the gradient contribution from hard negatives.

## 3.2. Tuned Contrastive Learning

In this section, we present our novel contrastive loss function — **Tuned Contrastive Learning (TCL) Loss**. Note that our representation learning framework remains the same as that of Supervised Contrastive Learning discussed above. The TCL loss is given by the following equations:

$$L^{tcl} = \sum_{i \in I} L_i^{tcl} \qquad (8)$$

$$L_i^{tcl} = \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log(\frac{\exp(z_i.z_p/\tau)}{D(z_i)}) \qquad (9)$$

where

$$D(z_i) = \sum_{p' \in P(i)} \exp(z_i.z_{p'}/\tau) + $$
$$k_1(\sum_{p' \in P(i)} \exp(-z_i.z_{p'})) + k_2(\sum_{n \in N(i)} \exp(z_i.z_n/\tau)) \qquad (10)$$

$$k_1, k_2 \geq 1 \qquad (11)$$

$k_1$ and $k_2$ are scalar parameters that are fixed before training. All other symbols have the same meaning as discussed in the previous section. We now present the following lemma:

**Lemma 2** *The gradient of the TCL loss per sample — $L_i^{tcl}$ with respect to the normalized projection network embedding $z_i$ is given by:*

$$\frac{\partial L_i^{tcl}}{\partial z_i} = \frac{1}{\tau}(\underbrace{\sum_{p \in P(i)} z_p(P_{ip}^t - X_{ip} - Y_{ip}^t)}_{\substack{Gradient \\ from\ positives}} + \underbrace{\sum_{n \in N(i)} z_n P_{in}^t}_{\substack{Gradient \\ from\ negatives}})$$

$$(12)$$

*where*

$$X_{ip} = \frac{1}{|P(i)|} \qquad (13)$$

$$P_{ip}^t = \frac{exp(z_i.z_p/\tau)}{D(z_i)} \qquad (14)$$

$$Y_{ip}^t = \frac{\tau k_1 exp(-z_i.z_p)}{D(z_i)} \qquad (15)$$

$$P_{in}^t = \frac{k_2 exp(z_i.z_n/\tau)}{D(z_i)} \qquad (16)$$

From Lemma 2, Theorem 1 and Theorem 2 follow in a straightforward fashion. The proofs for Lemma 2 and the two theorems are provided in our supplementary.

**Theorem 1** *For $k_1, k_2 \geq 1$, the magnitude of the gradient from a hard positive for TCL loss is strictly greater than the magnitude of the gradient from a hard positive for SupCon and hence, the following result follows:*

$$\underbrace{|X_{ip} - P_{ip}^t + Y_{ip}^t|}_{\textit{(TCL's hard positive gradient)}} > \underbrace{|X_{ip} - P_{ip}^s|}_{\textit{(Supcon's hard positive gradient)}} \qquad (17)$$

**Theorem 2** *For fixed $k_1$, the magnitude of the gradient response from a hard negative for TCL loss — $P_{in}^t$ strictly increases with $k_2$.*

**Effects of $k_1$ and $k_2$** The authors of SupCon show (in equation 18 in the supplementary of [23]) that the magnitude of gradient response from a hard positive $|X_{ip} - P_{ip}^s|$ increases with the number of positives and negatives in the batch. This is basically a result of reducing the value of $P_{ip}^s$, a term that results from having positive similarity terms in the denominator of $L_i^{sup}$. But they approximate the numerator of $P_{ip}^s$ to 1 by assuming $z_i.z_p \approx 0$ for a hard positive which might not always be true (especially since $\tau$ is typically chosen to be small like 0.1). As evident from the proof of Theorem 1 in our supplementary, we further push this idea and reduce the value of $P_{ip}^s$ in SupCon loss to $P_{ip}^t$ in TCL loss by having an extra term in the denominator involving $k_1$ — $k_1(\sum_{p' \in P(i)} \exp(-z_i.z_{p'}))$ and choosing a large enough value for $k_1$. Hence, it reduces the effect of implicit consideration of positives as negatives, the first limitation of SupCon loss discussed in the previous section. Note that having the extra term to increase the gradient response from hard positive is not the same as increasing the gradient response by amplifying the learning rate. This is because for the same and fixed learning rate, TCL loss increases the magnitude of the gradient signal over SupCon loss by changing the coefficient of $z_p$ in Eq. (12) which in turn means changing the gradient direction as well. This leads to consistently better performance as shown in the numerous experiments that we perform. Also, it directly follows from Theorem 2 that $k_2$ allows to regulate (increase) the gradient signal from a hard negative and thus, overcomes the second limitation of the SupCon loss.

**Augmentation Strategy for Self-Supervised Setting**
Since TCL loss can use multiple positives, we consider working with positive triplets instead of positive pairs in self-supervised settings. Given a batch $B$ with $N$ samples, we produce augmented batch $I$ of size $3N$ by producing three augmented views (positives) for each sample in $B$. This idea can further be extended in different ways to have more positives per anchor. For example, one can think of combining different augmentation strategies to produce multiple views per sample although we limit ourselves to positive triplets in this work.

## 4. Experiments

We evaluate TCL in three stages: 1. Supervised setting, 2. Hyper-parameter stability and 3. Self-supervised setting. We then present empirical analysis of TCL loss's parameters — $k_1$ and $k_2$. All the relevant training details are mentioned in our supplementary.

### 4.1. Supervised Setting

We start by evaluating TCL in supervised setting first. Since the authors of [23] mention that SupCon loss performs significantly better than triplet loss [22] and N-pair loss [28], we directly compare TCL loss with SupCon and cross-entropy losses on various classification benchmarks that include CIFAR-10, CIFAR-100 [24], Fashion MNIST (FMNIST) [31] and ImageNet-100 [12]. The encoder network chosen is ResNet-50 [21] for CIFAR and FMNIST datasets while Resnet-18 [21] for ImageNet dataset (because of memory constraints). The representation vector is the activation of the final pooling layer of the encoder. ResNet-18 and ResNet-34 encoders give 512 dimensional representation vectors while ResNet-50 and above produce 2048 dimensional vectors. The projector network is a MLP with one hidden layer with sizes being 512 for ResNet-18 and Resnet-34, and 2048 for ResNet-50 and higher networks. The output layer of the projector MLP is 128 dimensional for all the networks. We use the same cross-entropy implementation as used by Supervised Contrastive Learning [23].

**Note that for fair comparison of TCL with Supervised Contrastive Learning, we keep the architecture and all other possible hyper-parameters except the learning rate exactly the same. We also do hyper-parameter tuning significantly more for Supervised Contrastive Learning than for TCL.** As a result, we found that our re-implementation of Supervised Contrastive Learning gave better results than what is reported in the paper [23]. For example, on CIFAR-100 our significantly tuned version of SupCon achieves 79.1% top-1 classification accuracy, 2.6% more than what is reported in SupCon paper. As the authors of SupCon [23] mention that 200 epochs of contrastive training are sufficient for training

a ResNet-50 on complete ImageNet dataset, our observations for the supervised setting case on relatively smaller datasets like CIFAR, FMNIST and ImageNet-100 are consistent with this finding. We train Resnet-50 (and ResNet-18) for a total of 150 epochs – 100 epochs of contrastive training for the encoder and the projector followed by 50 epochs of cross-entropy training for the linear layer. Note that 150 epochs of total training was sufficient for our re-implementation of SupCon loss to achieve better results than reported in the paper (2.6% more on CIFAR-100 and 0.3% more on CIFAR-10). We anyways still provide results for 250 epochs of training in the supplementary. We have also provided 95% confidence intervals calculated over different seeds for this setting in the supplementary. As Tab. 1 shows, TCL loss consistently performs better than SupCon loss and outperforms cross-entropy loss on all the datasets.

| Dataset | Cross-Entropy | SupCon | TCL |
|---|---|---|---|
| CIFAR-10 | 95.0 | 96.3 (96.0) | 96.4 |
| CIFAR-100 | 75.3 | 79.1 (76.5) | 79.8 |
| FashionMNIST | 94.5 | 95.5 | 95.7 |
| ImageNet-100 | 84.2 | 85.9 | 86.7 |

Table 1. Comparison of top-1 accuracies of TCL loss with SupCon loss and cross-entropy loss in supervised settings. The values in parenthesis for SupCon loss denote the values presented in their paper.

### 4.2. Hyper-parameter Stability

We now show the stability of TCL loss to a range of hyper-parameters. We compare TCL loss with SupCon loss on various hyper-parameters — encoder architectures, batch sizes, projection embedding sizes and different augmentations. For all the hyper-parameter experiments we choose CIFAR-100 as the common dataset (unless stated otherwise), set total training epochs to 150 (same as earlier section), temperature $\tau$ to 0.1 and use SGD optimizer with momentum=0.9 and weight decay=$1e-4$.

#### 4.2.1 Encoder Architecture

We choose 4 encoder architectures of varying sizes- ResNet-18, ResNet-34, ResNet-50 and ResNet-101. For both TCL loss and SupCon loss, we choose batch size as 128 and AutoAugment [10] data augmentation method. As evident from Fig. 2b, TCL loss achieves consistent improvements in top-1 test classification accuracy over SupCon loss on all the architectures. We also tested TCL loss and SupCon loss on ImageNet-100 with ResNet-18 (batch size of 256) and ResNet-34 (batch size of 128). Using ResNet-18, TCL loss achieved 86.7% top-1 accuracy while SupCon loss achieved 85.9% top-1 accuracy. When switching to ResNet-
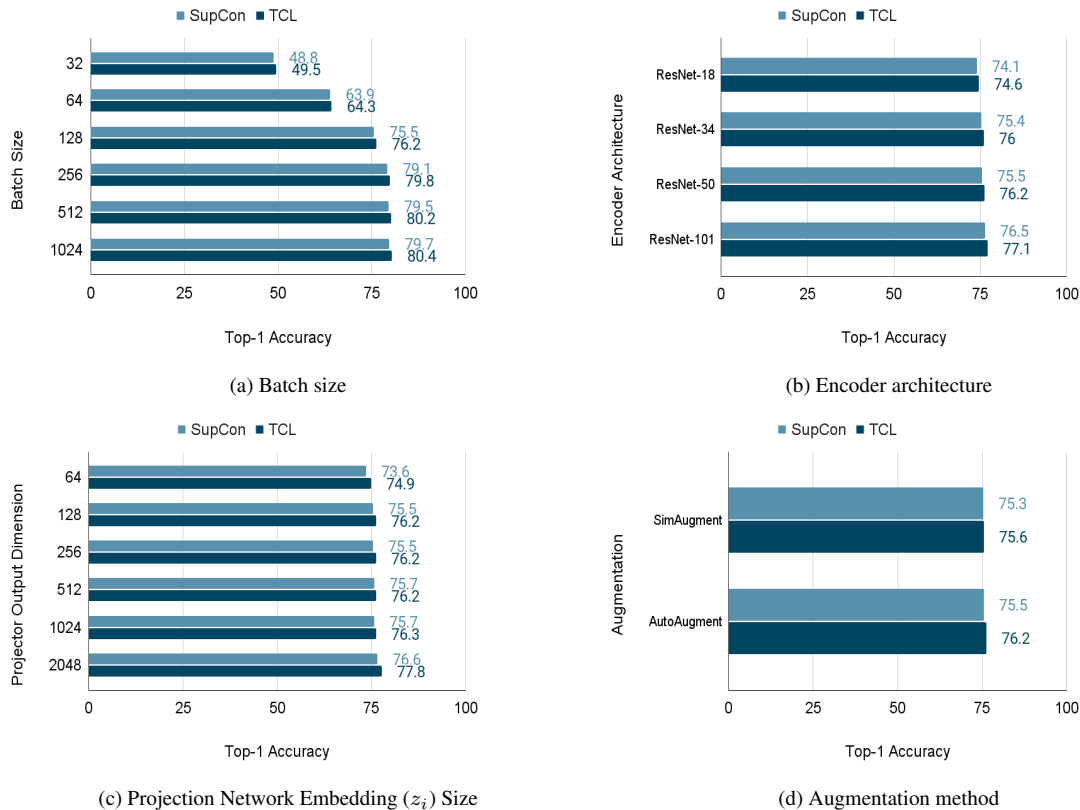
(a) Batch size



(b) Encoder architecture



(c) Projection Network Embedding ($z_i$) Size



(d) Augmentation method

Figure 2. SupCon vs TCL losses on a range of hyper-parameters.

34, TCL loss got 87.2% top-1 accuracy while SupCon loss got 86.5% top-1 accuracy.

### 4.2.2 Batch Size

For comparing TCL loss with SupCon loss on different batch sizes, we choose ResNet-50 as the encoder architecture and AutoAugment [10] data augmentation. As evident from Fig. 2a, we observe that TCL loss consistently performs better than SupCon loss on all batch sizes. All the batch sizes mentioned are after performing augmentation. Note that the authors of SupCon loss use an effective batch size of 256 (after augmentation) for CIFAR datasets in their released code[1]. We select batch sizes equal to, smaller and greater than this value for comparison to demonstrate the effectiveness of Tuned Contrastive Learning.

### 4.2.3 Projection Network Embedding ($z_i$) Size

In this section we analyse empirically how SupCon and TCL losses perform on various projection network output embedding sizes. This particular experiment was not explored as stated by the authors of Supervised Contrastive

Learning [23]. ResNet-50 is the common encoder used with Auto-Augment [10] data augmentation. As evident from Fig. 2c, we observe that TCL loss achieves consistent improvements in top-1 test classification accuracy over SupCon loss for various projector output sizes. We observe that 64 performs the worst while 128, 256, 512 and 1024 give similar results. 2048 performs the best for both with TCL loss achieving 1.2% higher accuracy than SupCon loss for this size.

### 4.2.4 Augmentations

We choose two augmentation strategies — AutoAugment and SimAugment for comparisons. AutoAugment [10] is a two-stage augmentation policy trained with reinforcement learning and gives stronger (aggressive and diverse) augmentations. SimAugment [6] is relatively a weaker augmentation strategy used in SimCLR that applies simple transformations like random flips, rotations, color jitters and gaussian blurring. We don't use gaussian blur in our implementation of SimAugment and train for 100 extra epochs i.e. 250 epochs while using it. Fig. 2d shows that TCL loss performs better than SupCon loss with both augmentations although, the gain is more with AutoAugment – the stronger

---

| Method | Projector Size | CIFAR-10 | CIFAR-100 | ImageNet-100 |
|---|---|---|---|---|
| BYOL [18] | 4096 | 92.6 | **70.2** | **80.1** |
| DINO [5] | 256 | 89.2 | 66.4 | 74.8 |
| SimSiam [9] | 2048 | 90.5 | 65.9 | 77.0 |
| MOCO V2 [8, 20] | 256 | **92.9** | 69.5 | 78.2 |
| ReSSL [36] | 256 | 90.6 | 65.8 | 76.6 |
| VICReg [3] | 2048 | 90.1 | 68.5 | 79.2 |
| SwAV [4] | 256 | 89.2 | 64.7 | 74.3 |
| W-MSE [15] | 256 | 88.2 | 61.3 | 69.1 |
| ARB [35] | 256 | 91.8 | 68.2 | 74.9 |
| ARB [35] | 2048 | 92.2 | 69.6 | 79.5 |
| Barlow-Twins [32] | 256 | 87.4 | 57.9 | 67.2 |
| Barlow-Twins [32] | 2048 | 89.6 | 69.2 | 78.6 |
| SimCLR [6] | 256 | 90.7 | 65.5 | 77.5 |
| **TCL (Self-Supervised)** | 256 | 91.8 | 67.2 | 78.4 |
| **TCL (Supervised)** | 128 | 95.8 | 77.5 | 86.7 |

Table 2. Comparison of top-1 accuracy of TCL with various SSL methods. Values in bold show the best performing method.

augmentation strategy.

## 4.3. Self-Supervised Setting

In this section we evaluate TCL without any labels in self-supervised setting by making use of positive triplets as described earlier. We compare TCL with various SOTA SSL methods as shown in Tab. 2. The results for these methods are taken from the works of [35], [11]. The datasets used for comparison are CIFAR 10, CIFAR-100 and ImageNet-100. ResNet-18 is the common encoder used for every method. For CIFAR-10 and CIFAR-100 every method uses 1000 epochs of contrastive pre-training including TCL. For ImageNet-100, every method does 400 epochs of contrastive pre-training.

Tab. 2 shows the top-1 accuracy achieved by various methods on the three datasets. TCL performs consistently better than SimCLR [6] and performs on par with various other methods. **Note that methods like BYOL [18], VI-CReg [3], ARB [35] and Barlow-Twins [32] use much larger projector size for output embedding and extra hidden layers in the projector MLP to get better performance while MOCO V2 [8] uses a queue size of 32,768 to get better results. Few of the methods like BYOL [18], SimSiam [9], MOCO V2 [8, 20] also maintain two networks and hence, effectively use double the number of parameters and are memory intensive.** TCL will also benefit from larger projector sizes, extra hidden layers or using a large momentum queue [8] but our aim here is to do a fair comparison with its counterpart SimCLR [6] and show that TCL loss is a better InfoNCE-like loss in this setting.

We also add the results of supervised TCL that can make use of labels as it is generalizable to any number of posi-
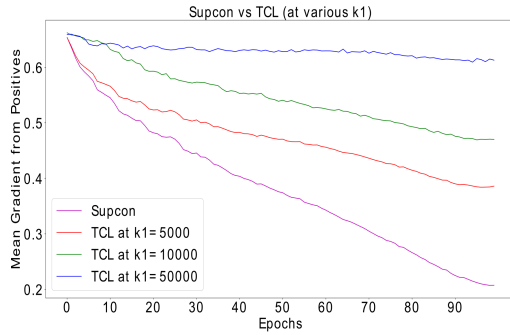
tives. Supervised TCL achieves significantly better results than all other SSL methods. SwAV does use a multi-crop strategy to create multiple augmentations but is not extended to supervised setting to use the labels [4].

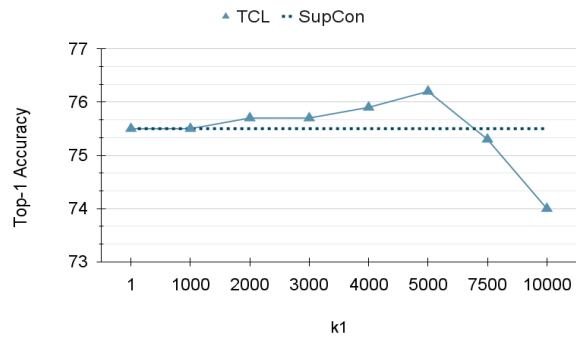## 4.4. Analyzing effects of $k_1$ and $k_2$ on TCL

As we discussed earlier in Sec. 3.2, $k_1$ helps in increasing the magnitude of gradient from positives while $k_2$ helps in regulating (increasing) the gradient from negatives. We verify our claims empirically. We provide values for $k_1$ and $k_2$ for all our experiments and insights on how to choose them in the supplementary.

**Analyzing effects of $k_1$** We calculate the mean gradient from all positives (expressions from Eq. (17)) per anchor averaged across the batch and plot the values for SupCon loss and TCL loss over the course of training of ResNet-50 on CIFAR-100 for 100 epochs. As evident from Fig. 3a, increasing the value of $k_1$ increases the magnitude of gradient response from positives. We also analyze how this correlates with the top-1 accuracy in Fig. 3b. As we see for small values of $k_1$, the top-1 accuracy remains more or less the same as that of SupCon loss. As we increase it further, the gradient from positives increase leading to gains in top-1 accuracy. The top-1 accuracy reaches a peak and then starts to drop with further increase in $k_1$. We hypothesize that this drop is because very large values of $k_1$ start affecting the gradient response from negatives (Eq. (16) and Eq. (10)). We verify this hypothesis while analyzing $k_2$.
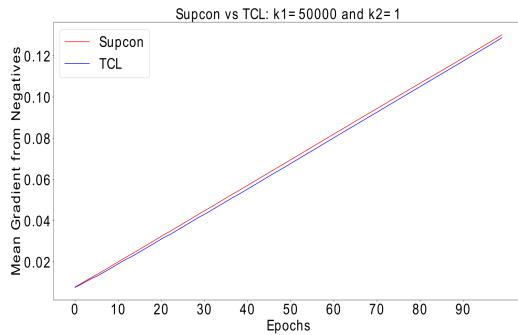
**Analyzing effects of $k_2$** We calculate the mean gradient from all negatives (expressions from Eq. (7) and Eq. (16)) per anchor averaged across the batch for the same setting as above and plot the values for SupCon loss and TCL loss. As we see in Fig. 3c, TCL loss's gradient lags behind Sup-
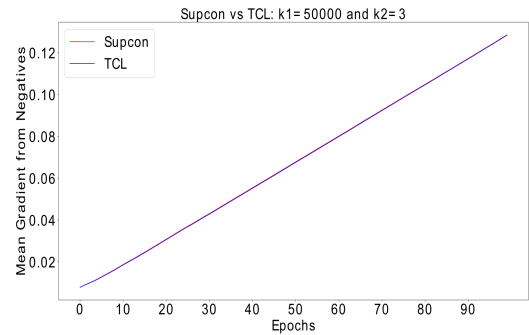
(a) Plot of mean gradient from positives for SupCon loss and TCL loss (at various values of $k_1$)



(b) Top-1 accuracy vs $k_1$ on CIFAR-100



(c) Plot of mean gradient from negatives for SupCon loss and TCL loss ($k_1 = 50000$ and $k_2 = 1$)



(d) Plot of mean gradient from negatives for SupCon loss and TCL loss ($k_1 = 50000$ and $k_2 = 3$)

Figure 3. Analysis of effects of the parameters $k_1$ and $k_2$ on TCL loss

Con loss's gradient by some margin for $k_1 = 50000$ and $k_2 = 1$. This value of $k_1$ actually leads to a top-1 accuracy of 71.8%, a drop in performance. When we start increasing the value of $k_2$, the gradient response from negatives increase for TCL loss. Fig. 3d shows that by increasing $k_2$ to 3 while $k_1 = 50000$, the gap between gradient (from negatives) curves of TCL loss and SupCon loss vanishes. We also observe that the top-1 accuracy for TCL loss increases back to 76.2%, the best possible accuracy that we got for this setting.

## 5. Conclusion & Limitations

In this work, we have presented a novel contrastive loss function called Tuned Contrastive Learning (TCL) loss that generalizes to multiple positives and multiple negatives present in a batch and is applicable to both supervised and self-supervised settings. We showed mathematically how its gradient response to hard positives and hard-negatives is better than that of SupCon loss. We evaluated TCL loss in supervised and self-supervised settings and showed that it performs on par with existing state-of-the-art supervised and self-supervised learning methods. We also showed empirically the stability of TCL loss to a range of hyper-

parameter settings.

A limitation of our work is that the proposed loss objective introduces two extra parameters $k_1$ and $k_2$, for which the values are chosen heuristically. Future direction can include works that try making these parameters learnable as part of the training process so that they are chosen automatically during test time or come up with loss objectives that provide the properties of TCL loss out of the box without introducing any extra parameters.

## 6. Potential Societal Impact

This paper proposes a novel contrastive loss that can help make discriminative models such as classification models more accurate and also prove to be helpful in labeling large amounts of unlabeled data. Hence, this work can find applications in AI systems used in various different industries.

At the same time, it is important to note that contrastive learning, in comparison to cross-entropy based learning, requires longer duration of training meaning higher energy consumption and more carbon emissions. This calls for developing new learning paradigms that offer the benefits of contrastive learning but take fewer epochs and less time to train.

# References

[1] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. 1

[2] Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, Avi Schwarzschild, Andrew Gordon Wilson, Jonas Geiping, Quentin Garrido, Pierre Fernandez, Amir Bar, Hamed Pirsiavash, Yann LeCun, and Micah Goldblum. A cookbook of self-supervised learning, 2023. 2

[3] Adrien Bardes, Jean Ponce, and Yann LeCun. VI-CReg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations*, 2022. 2, 7

[4] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020. 1, 2, 3, 7

[5] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 7

[6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 13–18 Jul 2020. 1, 2, 3, 6, 7

[7] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: a deep quadruplet network for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 403–412, 2017. 2

[8] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning, 2020. 2, 7

[9] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15750–15758, June 2021. 2, 7

[10] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018. 5, 6

[11] Victor Guilherme Turrisi da Costa, Enrico Fini, Moin Nabi, Nicu Sebe, and Elisa Ricci. solo-learn: A library of self-supervised methods for visual representation learning. *Journal of Machine Learning Research*, 23(56):1–6, 2022. 7

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 5

[13] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 1

[14] Gamaleldin Elsayed, Dilip Krishnan, Hossein Mobahi, Kevin Regan, and Samy Bengio. Large margin deep networks for classification. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. 2

[15] Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-supervised representation learning. *CoRR*, abs/2007.06346, 2020. 2, 7

[16] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728*, 2018. 1

[17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. 1

[18] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, koray kavukcuoglu, Remi Munos, and Michal Valko. Bootstrap your own latent - a new approach to self-supervised learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 21271–21284. Curran Associates, Inc., 2020. 1, 2, 7

[19] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742, 2006. 1

[20] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 2, 7

[21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5

[22] Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *Similarity-Based Pattern Recognition: Third International Workshop, SIMBAD 2015, Copenhagen, Denmark, October 12-14, 2015. Proceedings 3*, pages 84–92. Springer, 2015. 2, 5

[23] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*,

volume 33, pages 18661–18673. Curran Associates, Inc., 2020. 1, 2, 3, 4, 5, 6

[24] A. Krizhevsky and G Hinton. Learning multiple layers of features from tiny images, 2009. 5

[25] Weiyang Liu, Yandong Wen, Zhiding Yu, and Meng Yang. Large-margin softmax loss for convolutional neural networks, 2017. 2

[26] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft, 2013. 1

[27] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4004–4012, 2016. 2

[28] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. 5

[29] Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. Training convolutional networks with noisy labels, 2015. 2

[30] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *ArXiv*, abs/1807.03748, 2018. 1, 2

[31] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017. 5

[32] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. PMLR, 2021. 1, 2, 7

[33] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14*, pages 649–666. Springer, 2016. 1

[34] Richard Zhang, Phillip Isola, and Alexei A Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1058–1067, 2017. 2

[35] Shaofeng Zhang, Lyn Qiu, Feng Zhu, Junchi Yan, Hengrui Zhang, Rui Zhao, Hongyang Li, and Xiaokang Yang. Align representations with base: A new approach to self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16600–16609, June 2022. 1, 2, 7

[36] Mingkai Zheng, Shan You, Fei Wang, Chen Qian, Changshui Zhang, Xiaogang Wang, and Chang Xu. Ressl: Relational self-supervised learning with weak augmentation. *Advances in Neural Information Processing Systems*, 34:2543–2555, 2021. 7