

# An Introduction to Contrastive Learning

Zhiming Xu

zx2rw@virginia.edu



UNIVERSITY  
of VIRGINIA

ENGINEERING

School of Engineering and Applied Science  
University of Virginia

June 4, 2021

- ▶ Learning representations  $f_{\theta} : \mathbf{x} \rightarrow \mathbf{v}$  via separating positives and negatives (contrast) measured by score.
- ▶ Optimize parameters  $\theta$   
st.  $\text{score}(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}_+)) > \text{score}(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}_-))$ .

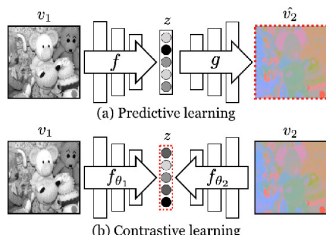


Figure 1: Predictive (autoencoder) v.s. Contrastive [1].



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

## Summary

- ▶ Problem formation: non-parametric instance-level classification.
- ▶ Learning objective: embedding function  $f_{\theta}(\cdot)$  that induces a metric over image space  $d_{\theta}(x, y) = \|f_{\theta}(x) - f_{\theta}(y)\|$ .
- ▶ Novelty: train a non-parametric classifier that distinguishes each image instance as its own class.

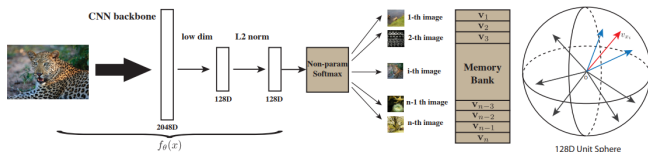


Figure 2: Non-parametric instance discrimination framework.



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

## Summary



- ▶ Suppose we have  $n$  images  $\{x_1, x_2, \dots, x_n\}$  and their features  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ ,  $\mathbf{v}_i = f_{\theta}(x_i)$ .
- ▶ *Parametric Classifier*

$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{w}_i^{\top} \mathbf{v})}{\sum_{j=1}^n \exp(\mathbf{w}_j^{\top} \mathbf{v})},$$

$\mathbf{w}_j$  is a weight vector (parameter) for each class  $j$ ,  $\mathbf{w}_j^{\top} \mathbf{v}$  measures how similar  $\mathbf{v}$  matches class (instance)  $j$ .

- ▶ *Non-Parametric Classifier*

$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{v}_i^{\top} \mathbf{v} / \tau)}{\sum_{j=1}^n \exp(\mathbf{v}_j^{\top} \mathbf{v} / \tau)}, \quad (1)$$

$\tau$  is a temperature parameter controlling the concentration level of distribution, and no training parameters are involved.



- ▶ Weight vectors (parameters)  $\{\mathbf{w}_j\}$  are only valid for training classes and do not generalize.
- ▶ Non-parametric features  $\{\mathbf{v}_i\}$  can be added from new instances.
- ▶ Non-parametric formulation eliminates the need for weight vectors and thus reduces computing and storing costs.

- ▶ Maximum the joint probability that all instances are classified as themselves:  $\prod_{i=1}^n P_{\theta}(i|f_{\theta}(x_i))$ .
- ▶ Equivalently, minimize the negative log-likelihood:

$$J(\theta) = - \sum_{i=1}^n \log P(i|f_{\theta}(x_i)). \quad (2)$$



To compute Eq. 1, all images' features,  $\{\mathbf{v}_i\}$  are needed and used multiple times. A *memory bank*  $V$  is used to reduce the computing and storing cost.

## Memory Bank

Suppose  $\mathbf{f}_i = f_{\theta}(x_i)$  is  $x_i$ 's feature. A memory bank is a set of features  $V = \{\mathbf{v}_i\}, \forall i$  randomly initialized and updated with  $\mathbf{v}_i \leftarrow \mathbf{f}_i$  during each training iteration.



## InsDis

Non-Parametric Softmax Classifier

**NCE & Regularization**

Image Class Classifier

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

## Summary

Due to the large number of classes,  $n$  scales up to millions and computing Eq. 2 can be prohibitively expensive. Noise-contrastive estimation (NCE) is used to solve this problem here.

The multi-class classification task before is cast to a binary classification that discriminates *data samples* and *noise samples*. The probability of a representation  $\mathbf{v}$  corresponding to the  $i$ -th example in  $V$  is:

$$P(i|\mathbf{v}) = \frac{\exp(\mathbf{v}^\top \mathbf{f}_i / \tau)}{\mathbf{Z}_i}, \quad (3)$$

$$\mathbf{Z}_i = \sum_{j=1}^n \exp(\mathbf{v}_j^\top \mathbf{f}_i / \tau). \quad (4)$$



Uniform distribution  $U(n)$  is used to formulate *noise distribution*. Suppose noise samples are  $m$  times more frequent than data samples. The posterior probability of sample  $i$  with feature  $\mathbf{v}$  from it is:

$$h(i, \mathbf{v}) = \frac{P(i|\mathbf{v})}{P(i|\mathbf{v}) + mP_n(i)}.$$

The objective to minimize negative log-posterior distribution is:

$$J_{NCE}(\boldsymbol{\theta}) = -E_{P_d} [\log h(i, \mathbf{v})] - m \cdot E_{P_n} [\log(1 - h(i, \mathbf{v}'))].$$

$P_d$  is the actual data distribution while  $P_n$  is the sampled noise distribution.  $\mathbf{v}'$  is the feature of another image.



Compute  $\mathbf{Z}_i$  from Eq. 4 is expensive. It is approximated via Monte Carlo:

$$\mathbf{Z} \simeq \mathbf{Z}_i \simeq nE_j [\exp(\mathbf{v}_j^\top \mathbf{f}_i / \tau)] = \frac{n}{m} \sum_{k=1}^m \exp(\mathbf{v}_{j_k}^\top \mathbf{f}_i / \tau).$$

With the NCE approximation, the computational complexity reduces from  $O(n)$  to  $O(1)$  per sample.



During each iteration only one instance per class will be seen, and this causes fluctuation. An additional regularization term  $\lambda \left\| \mathbf{v}_i^{(t)} - \mathbf{v}_i^{(t-1)} \right\|$  is added for each positive sample.

The final objective function is:

$$J_{NCE}(\boldsymbol{\theta}) = - E_{P_d} \left[ \log h(i, \mathbf{v}) - \lambda \left\| \mathbf{v}_i^{(t)} - \mathbf{v}_i^{(t-1)} \right\| \right] \\ - m \cdot E_{P_n} [\log(1 - h(i, \mathbf{v}'))].$$



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

**Image Class Classifier**

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

## Summary



To classify test image  $\hat{x}$ , first compute  $\hat{\mathbf{f}} = f_{\theta}(\hat{x})$ . Then compare it with all representations in  $V$  with cosine similarity and find the  $k$ -nearest neighbors  $\mathcal{N}_k$ . The class with maximum weight from  $\mathcal{N}_k$  will be the predicted image class.



- ▶ Problem formation: mutual information maximization.
- ▶ Learning objective: a representation that aims to maximize mutual information between different views of the same scene.
- ▶ Novelty: contrastive learning with different views.

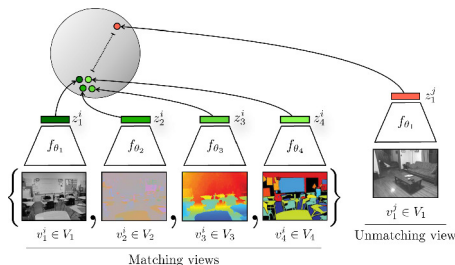


Figure 3: Contrastive multiview coding framework.



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

**Contrastive Learning with Two Views**

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

## Summary

The idea is learn a representation that separates samples from different distributions.

Given two views  $V_1$  and  $V_2$  of a collections of samples  $\{v_1^i, v_2^i\}_{i=1}^N$ , contrast congruent and incongruent pairs  $x \sim p(v_1, v_2)$  (*positives*) and  $y \sim p(v_1)p(v_2)$  (*negatives*).

A discrimination function  $h_{\theta}(\cdot)$  is trained to give high values for positives and low values for negatives. The contrastive loss is thus defined on a sample set  $S = \{x, y_1, y_2, \dots, y_k\}$  as:

$$\mathcal{L}_{contrast} = -E_S \left[ \log \frac{h_{\theta}(x)}{h_{\theta}(x) + \sum_{j=1}^k h_{\theta}(y_j)} \right].$$

To easily construct the sample set  $S$ , one view is fixed and the other enumerates positives and negatives:

$$\mathcal{L}_{contrast}^{V_1, V_2} = -E_{\{v_1, v_2^1, \dots, v_2^{k+1}\}} \left[ \log \frac{h_{\theta}(v_1^1, v_2^1)}{\sum_{j=1}^{k+1} h_{\theta}(v_1^1, v_2^j)} \right]. \quad (5)$$

The discrimination function  $h_{\theta}$  is a neural network. Specifically, two encoders  $f_{\theta_1}$  and  $f_{\theta_2}$  is used.

$$h_{\theta}(\{v_1, v_2\}) = \exp \left( \frac{1}{\tau} \cdot \frac{f_{\theta}(v_1) \cdot f_{\theta}(v_2)}{\|f_{\theta}(v_1)\| \cdot \|f_{\theta}(v_2)\|} \right).$$



Symmetrically,  $\mathcal{L}_{contrast}^{V_2, V_1}$  can be derived from Eq. 5.

Adding  $\mathcal{L}_{contrast}^{V_1, V_2}$  and  $\mathcal{L}_{contrast}^{V_2, V_1}$ , the two-view contrastive loss is:

$$\mathcal{L}(V_1, V_2) = \mathcal{L}_{contrast}^{V_1, V_2} + \mathcal{L}_{contrast}^{V_2, V_1}. \quad (6)$$



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

Contrastive Learning with Two Views

**Contrastive Learning with Multiple Views**

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

## Summary



Suppose there are  $M$  views  $V_1, V_2, \dots, V_M$ . The “core” view is the one to optimize. There are two paradigms that can be used.

- ▶ *Pair-wise*:  $\mathcal{L}_C = \sum_{j=2}^M \mathcal{L}(V_1, V_j)$ .
- ▶ *Full grap*:  $\mathcal{L}_F = \sum_{1 \leq i < j \leq M} \mathcal{L}(V_i, V_j)$ .



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

**Implementation Details**

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

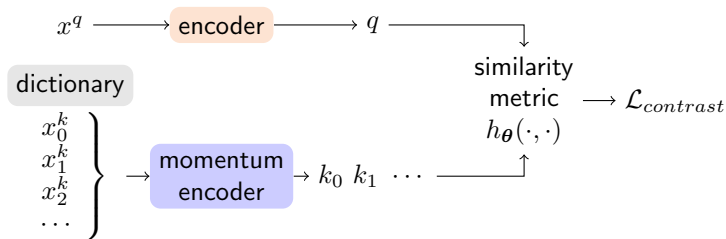
## Summary





In the last model, noise-contrast estimation is used to reduce number of classes. CMC uses *negative sampling* and formulate contrastive loss as a  $(m + 1)$ -way softmax classification. The idea of memory bank is also adopted here and it is dynamically updated.

- ▶ Problem formation: instance-level classification.
- ▶ Learning objective: moving-averaged encoder.
- ▶ Novelty: build a large and consistent dictionary on-the-fly that facilitates unsupervised contrastive learning.



**Figure 4:** Contrastive as dictionary look-up. The optimization goal is retrieve the positive key(s) from the dictionary.



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

## Summary

Consider an encoded query  $q$  and a set of encoded samples (*keys*)  $\{k_0, k_1, \dots\}$ . MoCo use the following loss.

### InfoNCE

$$\mathcal{L}_q = -\log \frac{\exp(q \cdot k_+ / \tau)}{\sum_{j=0}^K \exp(q \cdot k_j / \tau)} \quad (7)$$

Intuitively,  $\mathcal{L}_q$  is a  $K + 1$ -way softmax classification loss that classifies  $q$  as  $k_+$ . Two identical or distinct encoders  $f_{\theta_q}$  and  $f_{\theta_k}$  are used to encode  $q$  and  $k$ 's, respectively.



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

**Momentum Contrast**

## SupCon

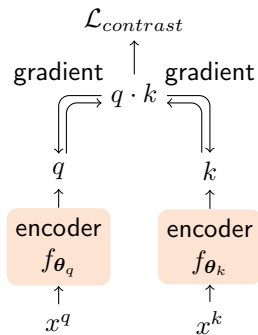
Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

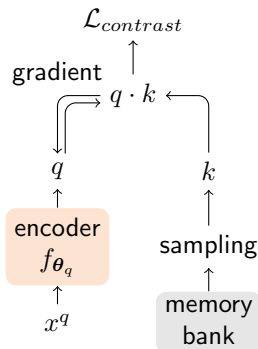
## Summary



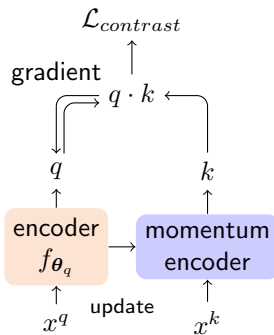
- ▶ *Dictionary as a queue*: maintains dictionary as a queue (first in, first out), decouples its size from batch size, and progressively updated with new batches and iterations.
- ▶ *Momentum update*: By maintaining a queue, the encoder  $f_{\theta_k}$  can not be updated via back-propagation. Momentum update is used instead to update it consistently:  $\theta_k \leftarrow m \cdot \theta_k + (1 - m) \cdot \theta_q$ .



(a) end-to-end



(b) memory bank  
instance-only: InsDic  
multiview: CMC



(c) MoCo

- ▶ Problem formation: supervised image classification
- ▶ Learning objective: a representation from supervised class labels.
- ▶ Novelty: extend self-supervised contrast to fully-supervised scenarios.

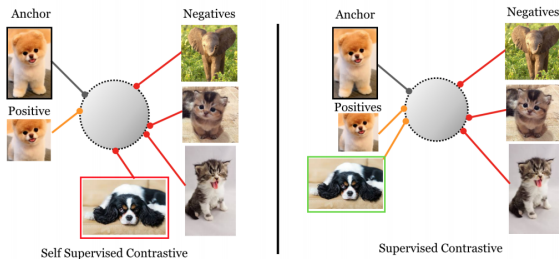


Figure 5: Self-supervised v.s. Supervised Contrast.





## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

**Supervised Contrastive Losses**

Choice of  $\mathcal{L}^{sup}$

## Summary

Self supervised contrastive losses 2, 6 and 7 do not consider additional label information. To include it, two kinds of supervised contrastive losses are proposed:

$$\mathcal{L}_{out}^{sup} = \sum_{i \in I} \mathcal{L}_{out,i}^{sup} = \sum_{i \in I} -\frac{1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_p / \tau)}{\sum_{a \in A(i)} \exp(\mathbf{z}_i \cdot \mathbf{z}_a / \tau)}$$

$$\mathcal{L}_{in}^{sup} = \sum_{i \in I} \mathcal{L}_{in,i}^{sup} = \sum_{i \in I} -\log \left\{ \frac{1}{|P(i)|} \sum_{p \in P(i)} \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_p / \tau)}{\sum_{a \in A(i)} \exp(\mathbf{z}_i \cdot \mathbf{z}_a / \tau)} \right\}$$

- ▶  $A(i)$  is all the samples except  $i$  itself.
- ▶  $P(i)$  is all the samples in  $A(i)$  that has the same class label as  $i$ .



- ▶ Generalize to an arbitrary number of positives.
- ▶ Capability increases with more negatives.
- ▶ Ability to perform hard positive/negative mining.



## InsDis

Non-Parametric Softmax Classifier

NCE & Regularization

Image Class Classifier

## CMC

Contrastive Learning with Two Views

Contrastive Learning with Multiple Views

Implementation Details

## MoCo

Contrastive Learning as Dictionary Look-up

Momentum Contrast

## SupCon

Supervised Contrastive Losses

Choice of  $\mathcal{L}^{sup}$

## Summary

The gradients for both  $\mathcal{L}_{out,i}^{sup}$  and  $\mathcal{L}_{in,i}^{sup}$  w.r.t.  $\mathbf{z}_i$  have the same form:

$$\frac{\partial \mathcal{L}_{\cdot,i}^{sup}}{\partial \mathbf{z}_i} = \frac{1}{\tau} \left\{ \sum_{p \in P(i)} \mathbf{z}_p (P_{ip} - X_{ip}) + \sum_{n \in N(i)} \mathbf{z}_n P_{in} \right\}$$

►  $N(i)$  is the complement of  $P(i)$  w.r.t.  $A(i)$ .

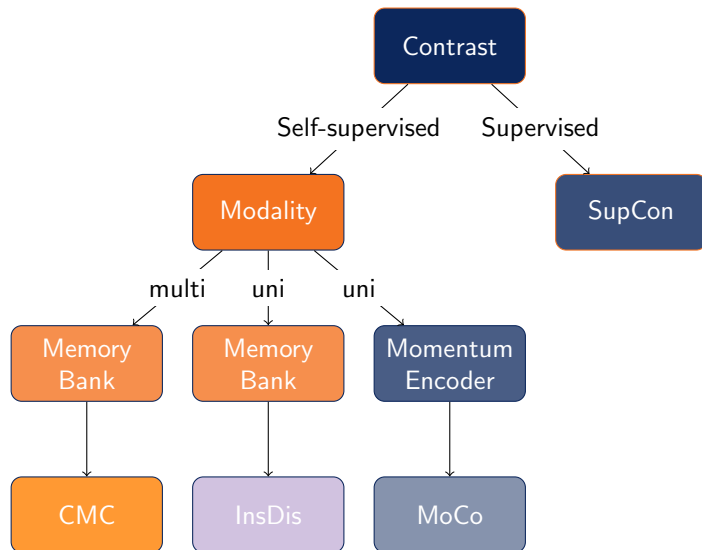
►  $P(ix) \equiv \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_x / \tau)}{\sum_{a \in A(i)} \exp(\mathbf{z}_i \cdot \mathbf{z}_a / \tau)}$ ,  $x \in \{p, n\}$ .

►  $X_{ip} = \begin{cases} \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_p / \tau)}{\sum_{p' \in P(i)} \exp(\mathbf{z}_i \cdot \mathbf{z}_{p'} / \tau)}, & \text{if } \mathcal{L}_{\cdot,i}^{sup} = \mathcal{L}_{in,i}^{sup} \\ \frac{1}{|P(i)|}, & \text{if } \mathcal{L}_{\cdot,i}^{sup} = \mathcal{L}_{out,i}^{sup} \end{cases}$ .

If  $z_p$  is substituted with the less biased mean representation  $\bar{z}$ ,  $X_{ip}^{in}$  will reduce to  $X_{ip}^{out}$ . This will lead to more stability in training and reduce impact of a single positive sample.

Supervised contrastive loss  $\mathcal{L}^{sup}$

$$\mathcal{L}_{out}^{sup} = \sum_{i \in I} \mathcal{L}_{out,i}^{sup} = \sum_{i \in I} -\frac{1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(z_i \cdot z_p / \tau)}{\sum_{a \in A(i)} \exp(z_i \cdot z_a / \tau)}$$



# Thank You for Your Attention



Q & A



- [1] Y. Tian, D. Krishnan, and P. Isola, “Contrastive multiview coding,” in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XI*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds., ser. Lecture Notes in Computer Science, vol. 12356, Springer, 2020, pp. 776–794. DOI: 10.1007/978-3-030-58621-8\\_45. [Online]. Available: [https://doi.org/10.1007/978-3-030-58621-8%5C\\_45](https://doi.org/10.1007/978-3-030-58621-8%5C_45).
- [2] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, “Unsupervised feature learning via non-parametric instance discrimination,” in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, IEEE Computer Society, 2018, pp. 3733–3742. DOI: 10.1109/CVPR.2018.00393. [Online]. Available: [http://openaccess.thecvf.com/content%5C\\_cvpr%5C\\_2018/html/Wu%5C\\_Unsupervised%5C\\_Feature%5C\\_Learning%5C\\_CVPR%5C\\_2018%5C\\_paper.html](http://openaccess.thecvf.com/content%5C_cvpr%5C_2018/html/Wu%5C_Unsupervised%5C_Feature%5C_Learning%5C_CVPR%5C_2018%5C_paper.html).



- [3] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, "Momentum contrast for unsupervised visual representation learning," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, IEEE, 2020, pp. 9726–9735. DOI: [10.1109/CVPR42600.2020.00975](https://doi.org/10.1109/CVPR42600.2020.00975). [Online]. Available: <https://doi.org/10.1109/CVPR42600.2020.00975>.
- [4] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M.-F. Balcan, and H.-T. Lin, Eds., 2020. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/hash/d89a66c7c80a29b1bdbab0f2a1a94af8-Abstract.html>.