

# **Pretext tasks**

## **3. SimCLR**

1	2	SELF-PREDICTION	INNATE RELATIONSHIP (Context-based)	1. ROTATION 2. RELATIVE POSITION	IMAGE
3		CONTRASTIVE LEARNING	INTER-SAMPLE CLASSIFICATION	1. Instance Discrimination 2. SimCLR [Contrastive Loss] 3. Theory – Guarantees / Bounds	IMAGE
4		CONTRASTIVE LEARNING	INTER-SAMPLE CLASSIFICATION	Contrastive Predictive Coding (CPC), [NCE, InfoNCE Loss]	AUDIO/ SPEECH
5		SELF-PREDICTION	GENERATIVE (VAE)	1. AE – Variational Bayes 2. VQ-VAE + AR	IMAGE  AUDIO/ SPEECH
6		SELF-PREDICTION	GENERATIVE (AR)	1. AR-LM – GPT 2. Masked-LM – BERT	LANGUAGE
7		SELF-PREDICTION	MASKED-GEN (Masked LM for ASR)	1. Wav2Vec / 2.0 2. HuBERT	AUDIO/ SPEECH

# A Simple Framework for Contrastive Learning of Visual Representations

Ting Chen, Simon Kornblith, Mohammad Norouzi, Geoffrey Hinton

# A Simple Framework for Contrastive Learning of Visual Representations

Ting Chen<sup>1</sup> Simon Kornblith<sup>1</sup> Mohammad Norouzi<sup>1</sup> Geoffrey Hinton<sup>1</sup>

<sup>1</sup>Google Research, Brain Team. Correspondence to: Ting Chen <iamtingchen@google.com>.

*Proceedings of the 37<sup>th</sup> International Conference on Machine Learning*, Vienna, Austria, PMLR 119, 2020. Copyright 2020 by the author(s).

<sup>1</sup>Code available at <https://github.com/google-research/simclr>.

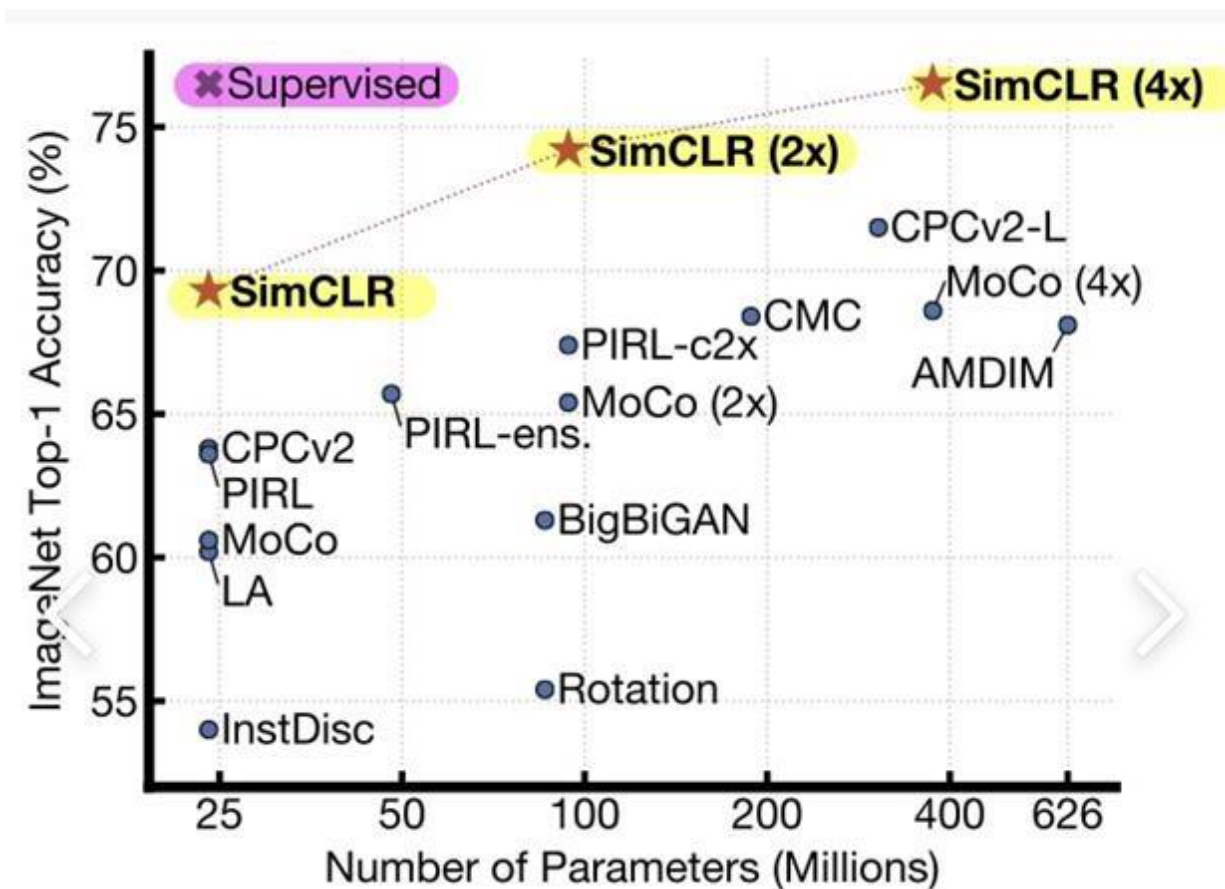
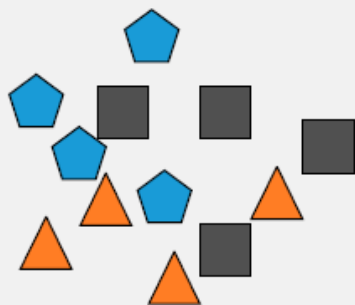


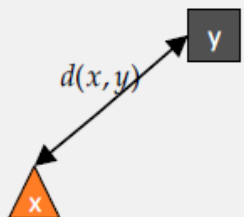
Figure 1. ImageNet Top-1 accuracy of linear classifiers trained on representations learned with different self-supervised methods (pretrained on ImageNet). Gray cross indicates supervised ResNet-50. Our method, SimCLR, is shown in bold.

# Deep Metric Learning

a) Original data space



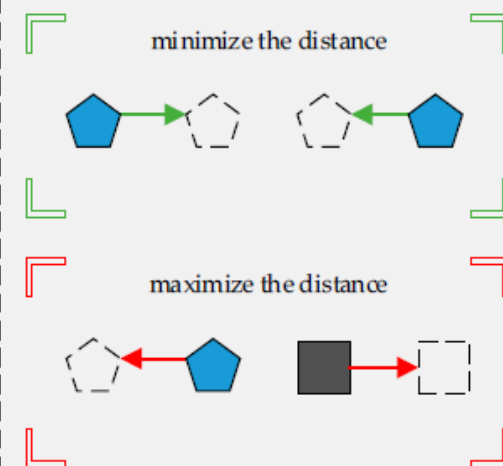
b) Euclidean metric



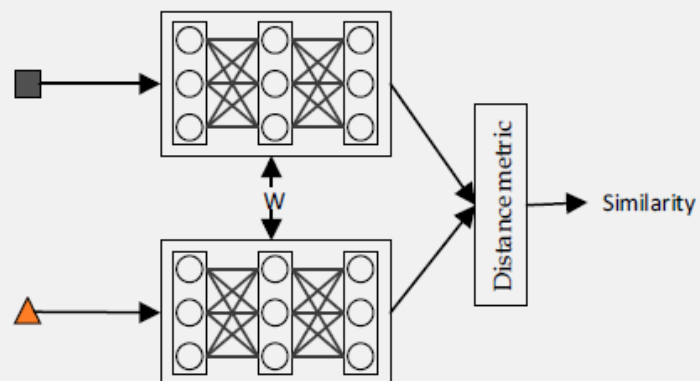
A diagram showing two points,  $x$  (an orange triangle) and  $y$  (a dark gray square), with a double-headed arrow between them labeled  $d(x, y)$ .

$$d(x, y) = \sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

c) Purpose of metric learning

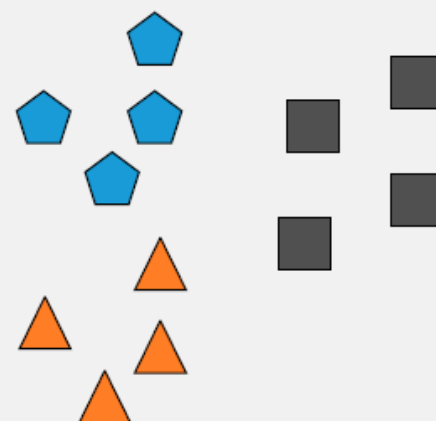


d) Deep metric learning example\*

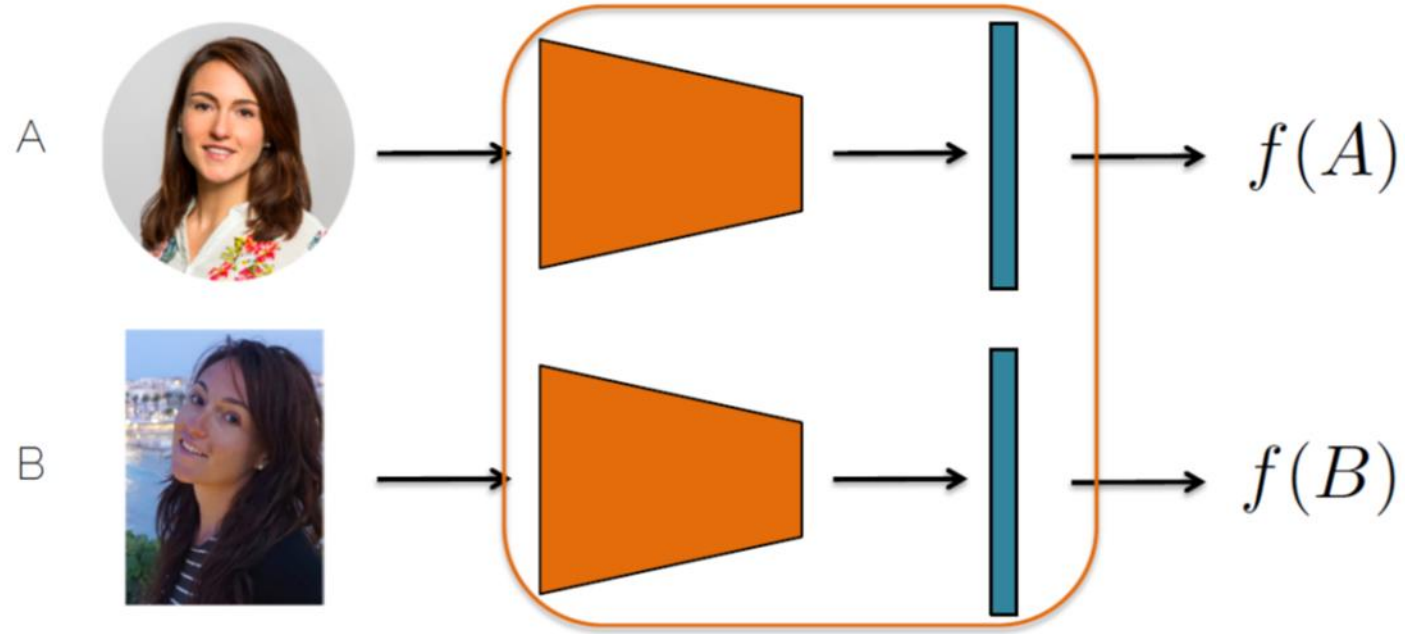


\*Siamese Network

e) Transformed data space



- Siamese network = shared weights



- Distance function  $d(A, B) = ||f(A) - f(B)||^2$
- Training: learn the parameter such that
  - If  $A$  and  $B$  depict the same person,  $d(A, B)$  is small
  - If  $A$  and  $B$  depict a different person,  $d(A, B)$  is large

- Contrastive loss:

$$\mathcal{L}(A, B) = y^* ||f(A) - f(B)||^2 + (1 - y^*) \max(0, m^2 - ||f(A) - f(B)||^2)$$



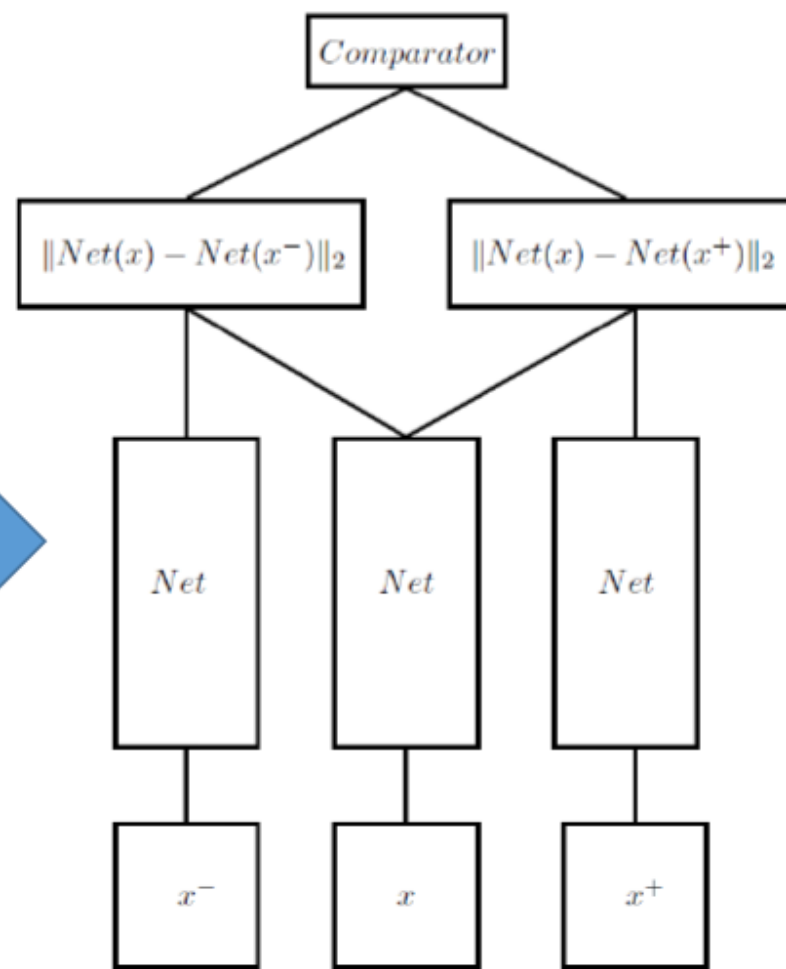
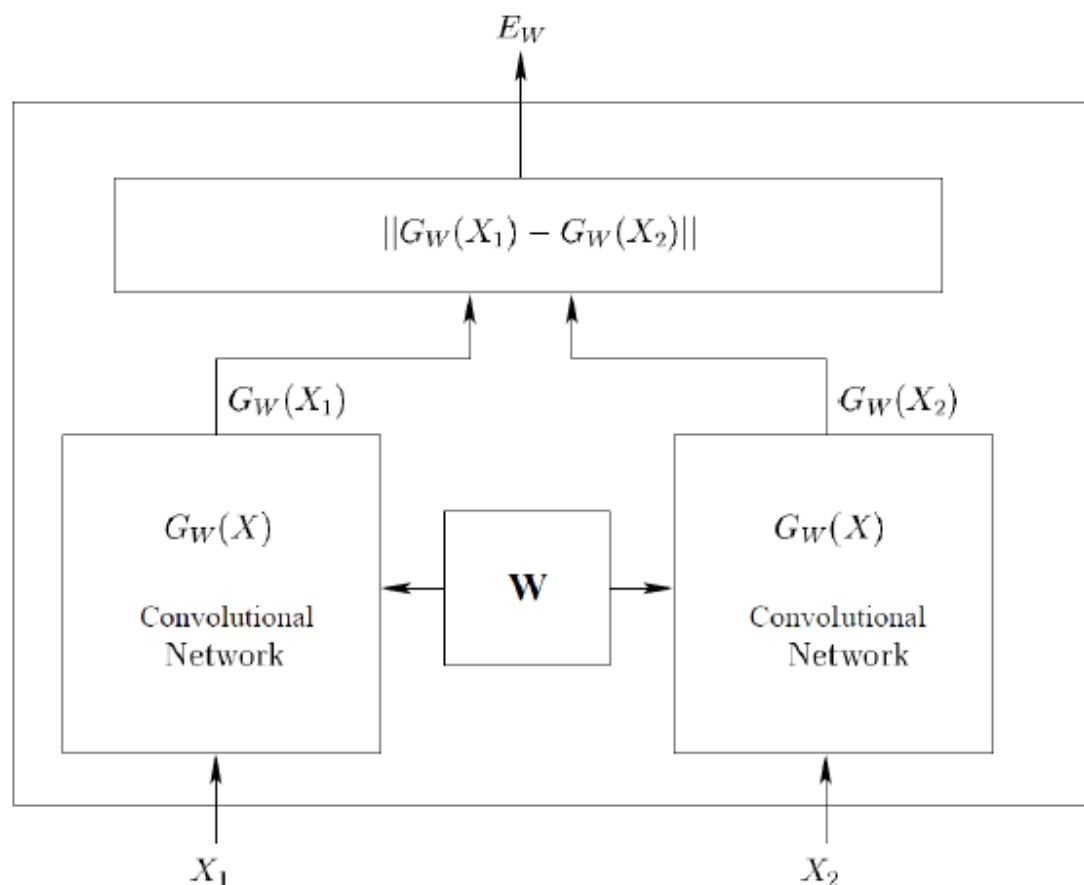
Positive pair,  
reduce the distance  
between the  
elements



Negative pair,  
brings the elements  
further apart up to a  
margin

- Training the siamese networks
  - You can update the weights for each channel independently and then average them
- This loss function allows us to learn to bring positive pairs together and negative pairs apart

# Triplet Network



From Siamese to Triplet Network



# Triplet loss

- Triplet loss allows us to learn a ranking



Anchor (A)



Positive (P)



Negative (N)

We want:  $\|f(A) - f(P)\|^2 < \|f(A) - f(N)\|^2$

- Triplet loss allows us to learn a ranking

$$\|f(A) - f(P)\|^2 < \|f(A) - f(N)\|^2$$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 < 0$$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + m < 0$$

  
margin

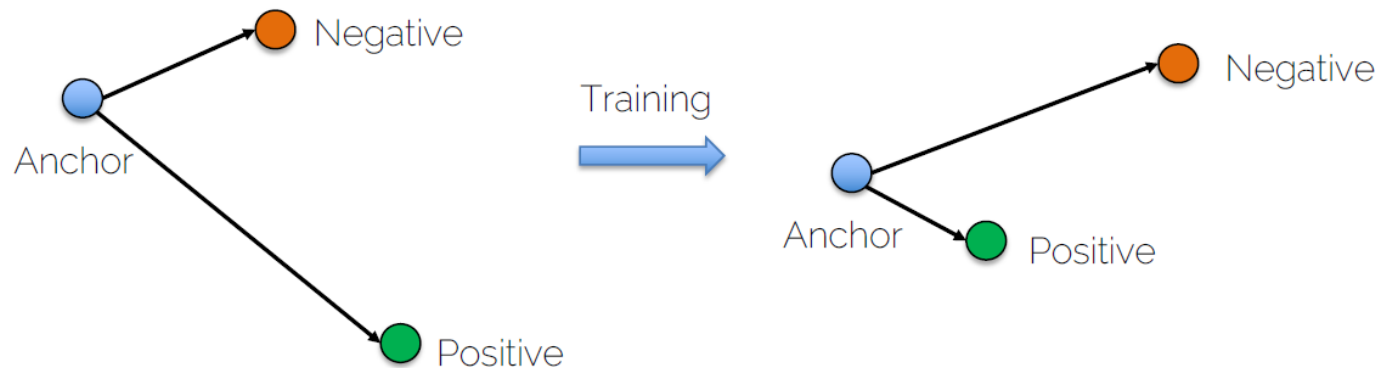
- Triplet loss allows us to learn a ranking

$$||f(A) - f(P)||^2 < ||f(A) - f(N)||^2$$

$$||f(A) - f(P)||^2 - ||f(A) - f(N)||^2 < 0$$

$$||f(A) - f(P)||^2 - ||f(A) - f(N)||^2 + m < 0$$

$$\mathcal{L}(A, P, N) = \max(0, ||f(A) - f(P)||^2 - ||f(A) - f(N)||^2 + m)$$



# The proposed SimCLR framework

A simple idea: maximizing the agreement of representations under data transformation, using a contrastive loss in the latent/feature space.

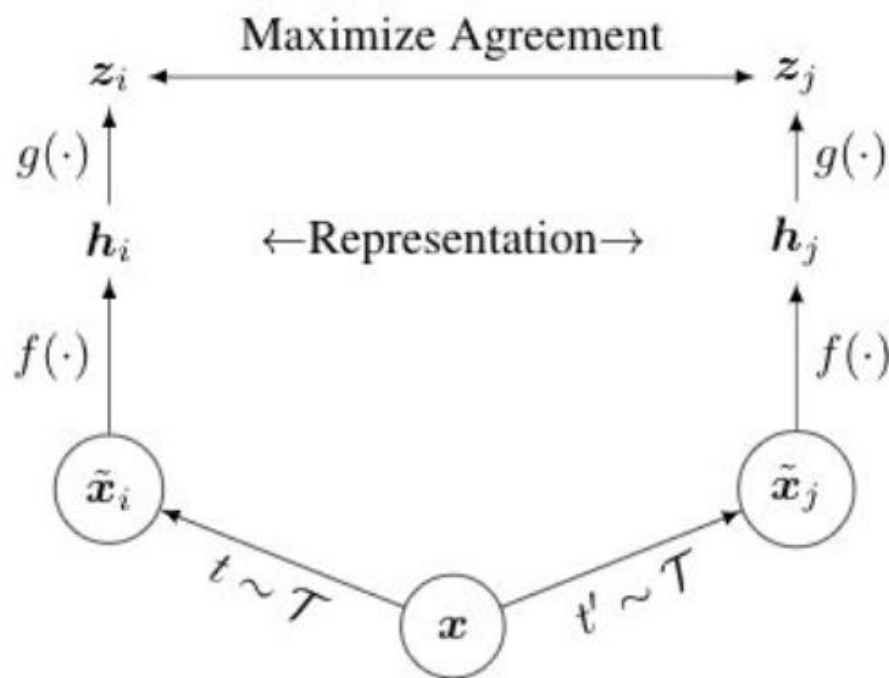
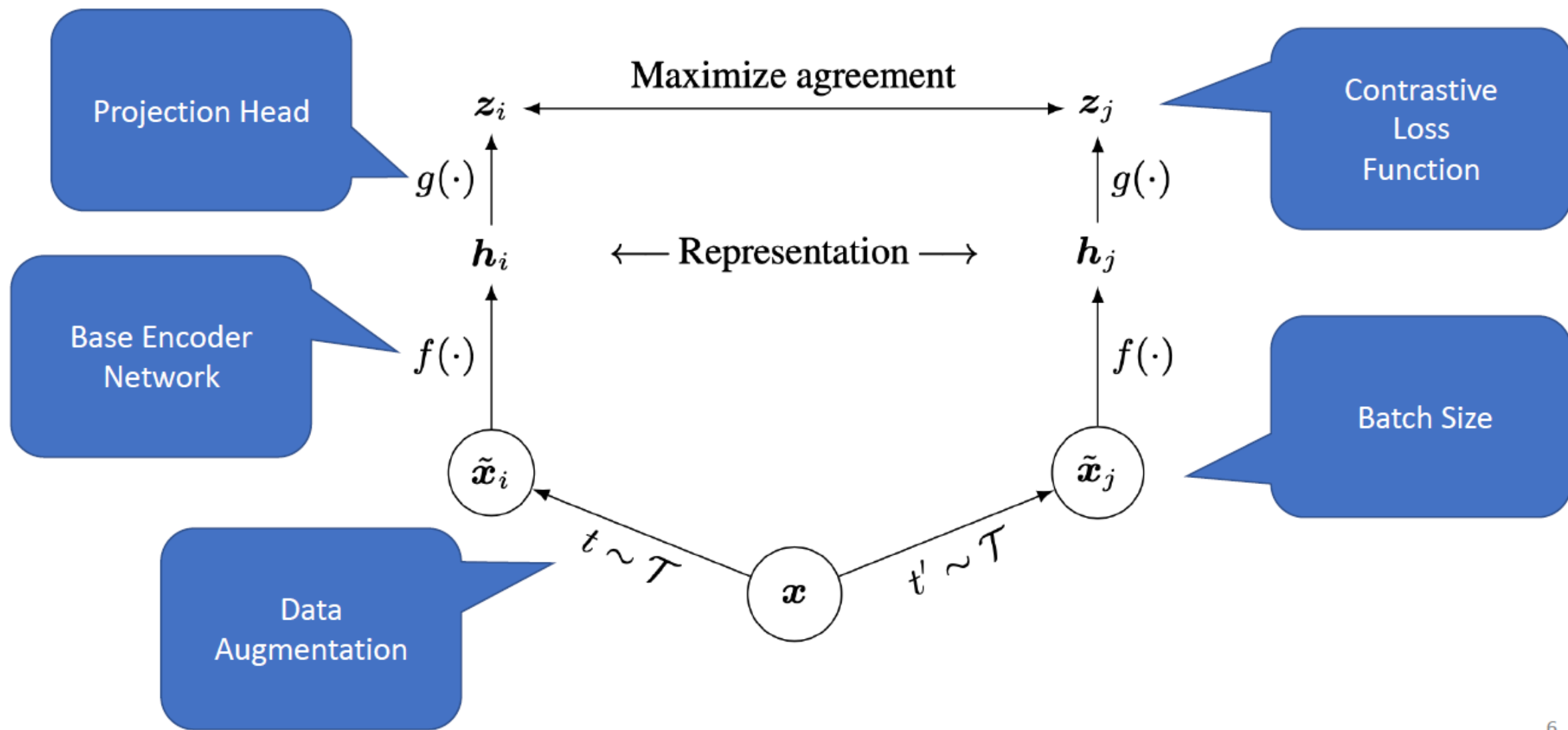


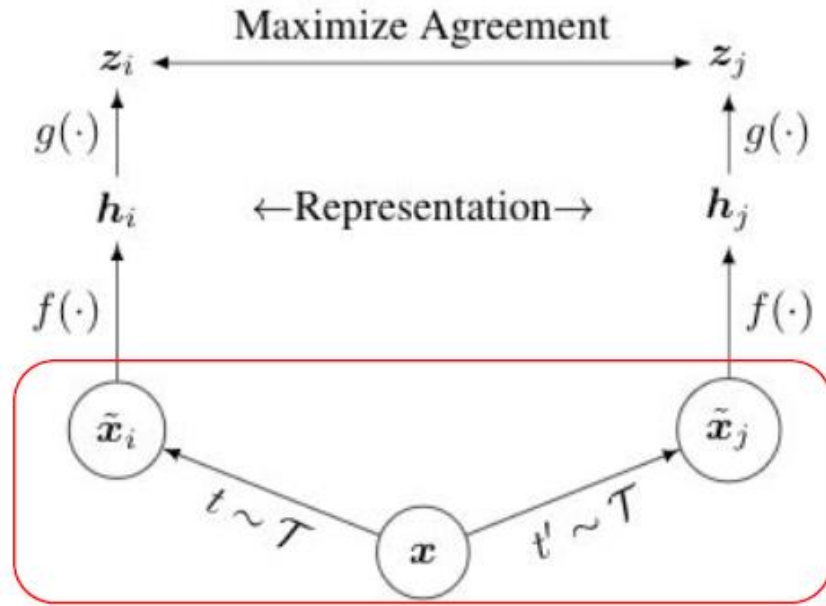
Figure 2. A framework for contrastive representation learning. Two separate stochastic data augmentations  $t, t' \sim \mathcal{T}$  are applied to each example to obtain two correlated views. A base encoder network  $f(\cdot)$  with a projection head  $g(\cdot)$  is trained to maximize agreement in *latent representations* via a contrastive loss.

# Framework



We use random crop and color distortion for augmentation.

Examples of augmentation applied to the left most images:



Three augmentations applied sequentially

- Random cropping

- Random color distortions

- Random Gaussian blur



# Systematically study a set of augmentation



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



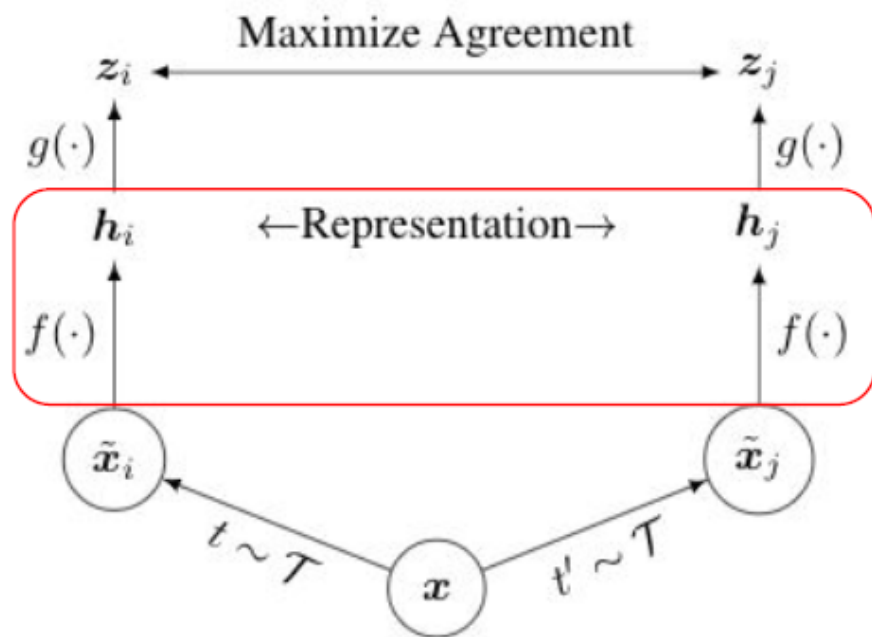
(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

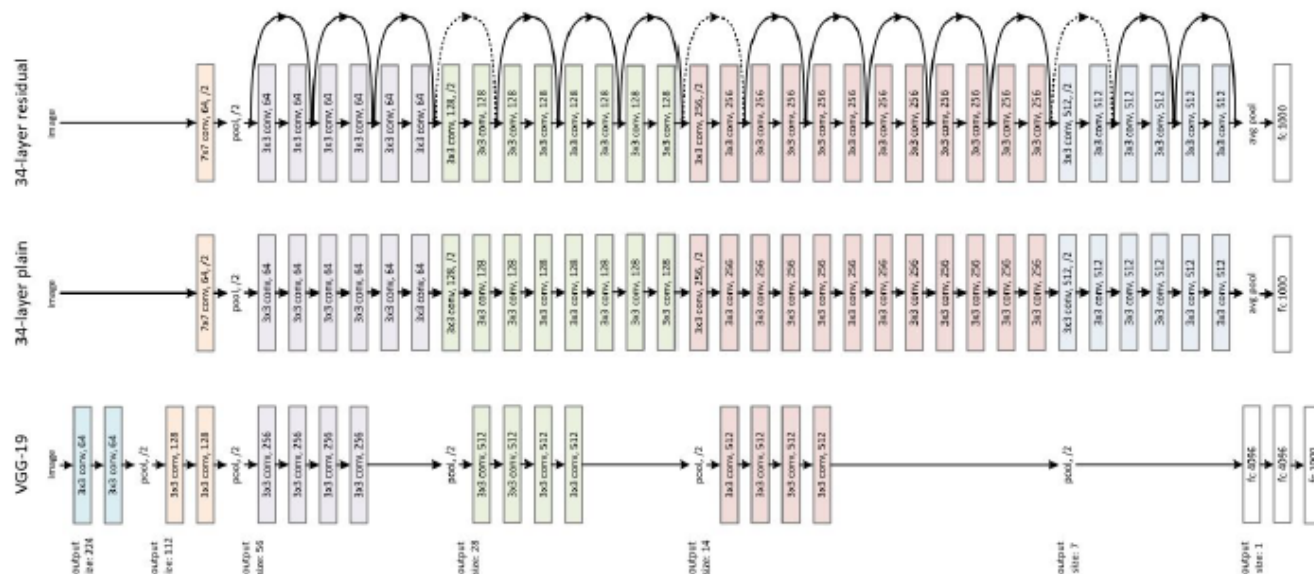


SimCLR chooses ResNet

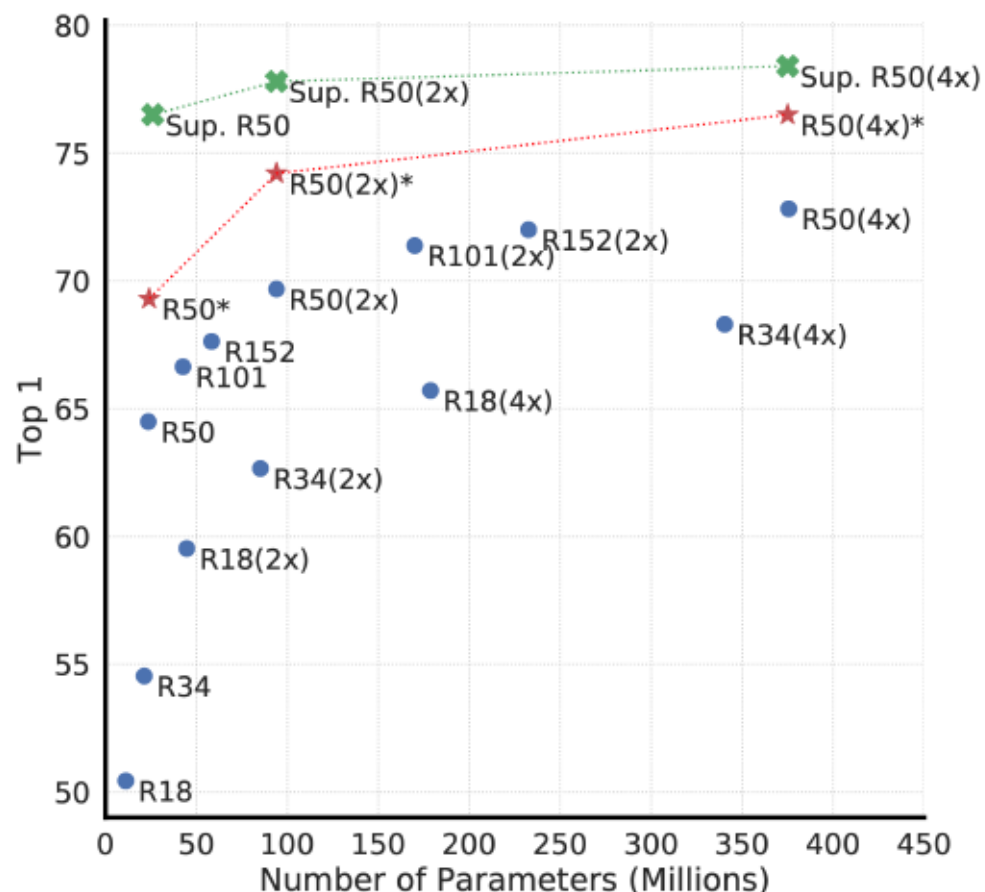
$$h_i = f(\tilde{x}_i) = \text{ResNet}(\tilde{x}_i)$$

$f(x)$  is the base network that computes internal representation.

We use (unconstrained) ResNet in this work. However, it can be other networks.

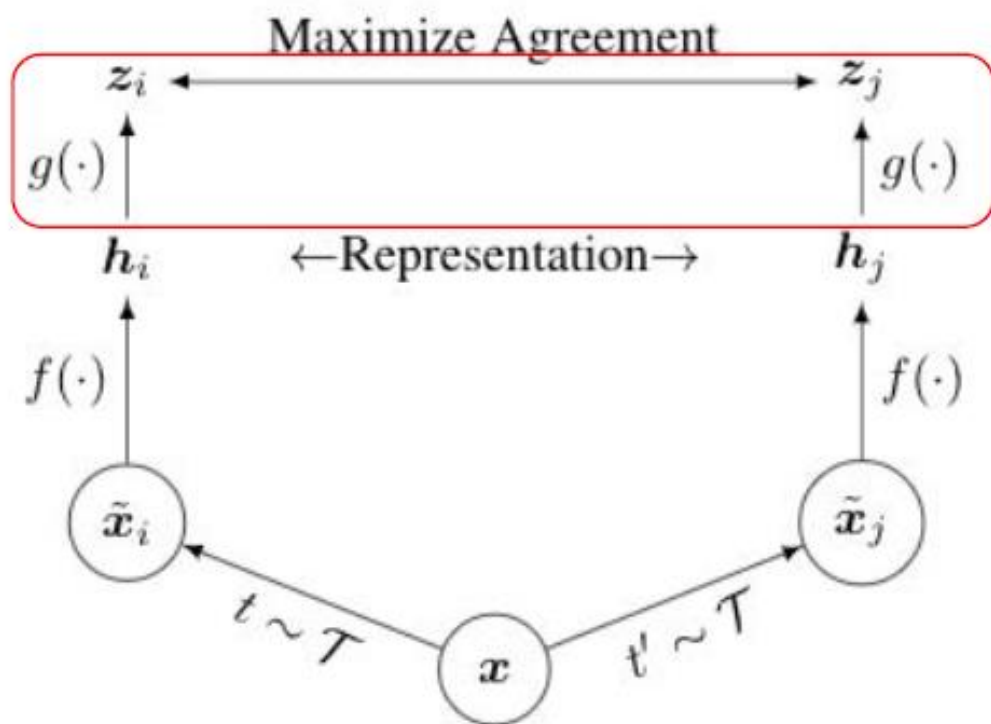


# Base Encoder



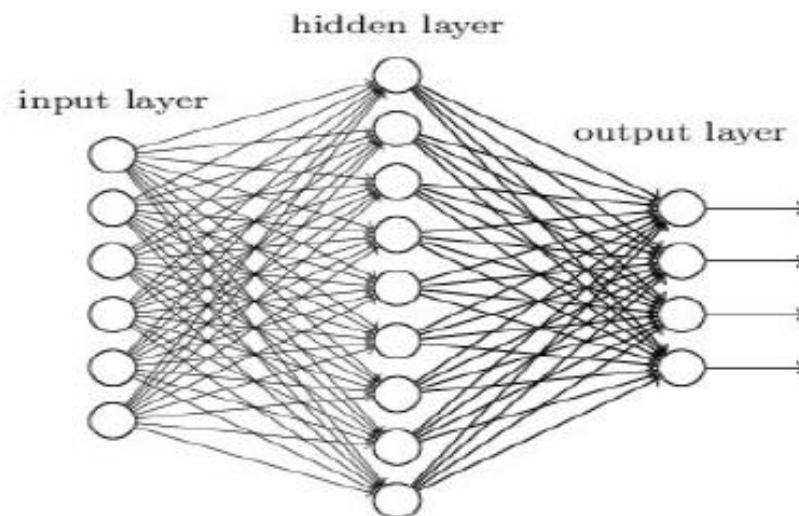
Performance gap shrinks as model size increases  
Unsupervised learning benefits more from bigger models





$g(h)$  is a projection network that project representation to a latent space.

We use a 2-layer non-linear MLP (fully connected net).



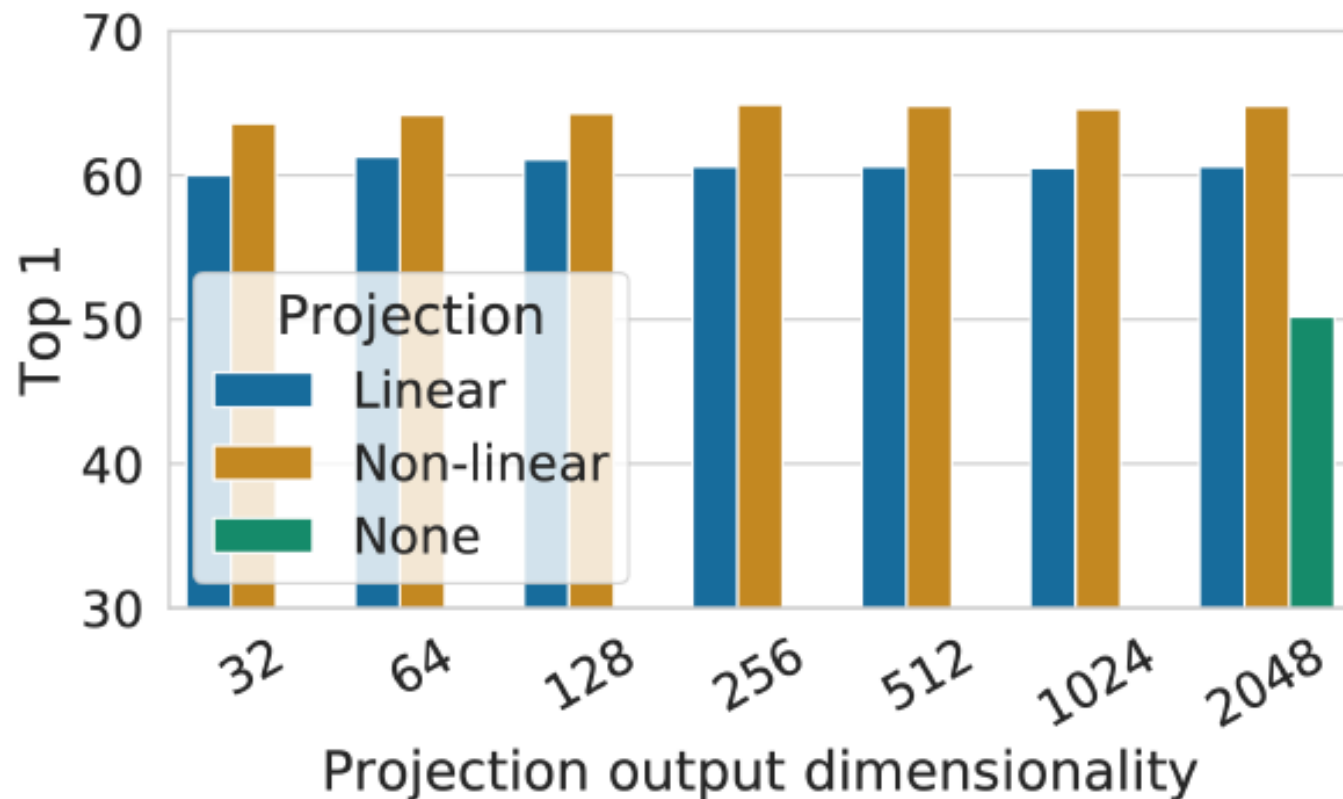
A small neural network

Multilayer Perceptron (MLP)

$$z_i = g(h_i) = W^{(2)} \sigma(W^{(1)} h_i)$$

$\sigma$  is ReLU (non-linearity)

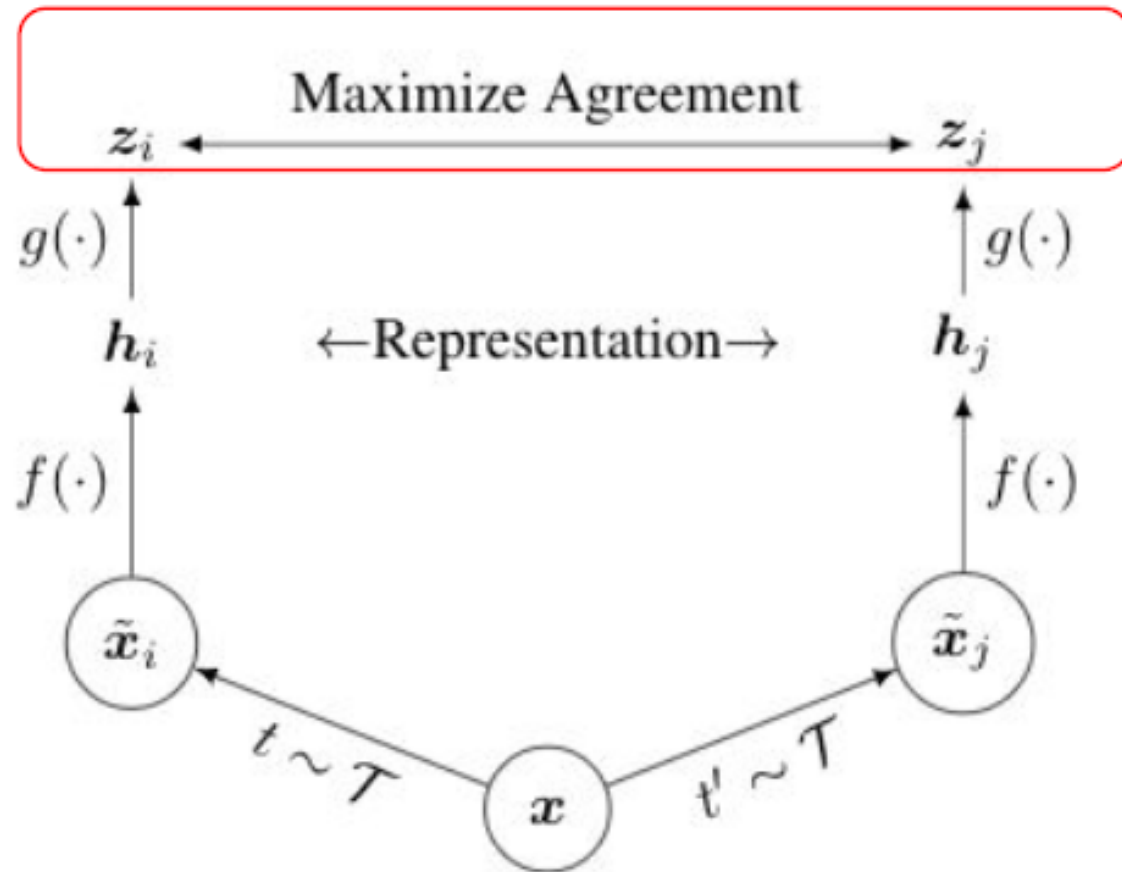
# Projection Head



Non-Linear > Linear >> None

Maximize agreement using a contrastive task:

Given  $\{x_k\}$  where two different examples  $x_i$  and  $x_j$  are a positive pair, identify  $x_j$  in  $\{x_k\}_{k \neq i}$  for  $x_i$ .



Loss function:

$$\text{Let } \text{sim}(\mathbf{u}, \mathbf{v}) = \mathbf{u}^\top \mathbf{v} / \|\mathbf{u}\| \|\mathbf{v}\|$$

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j) / \tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k) / \tau)}$$

# SimCLR pseudo code

---

**Algorithm 1** SimCLR's main learning algorithm.

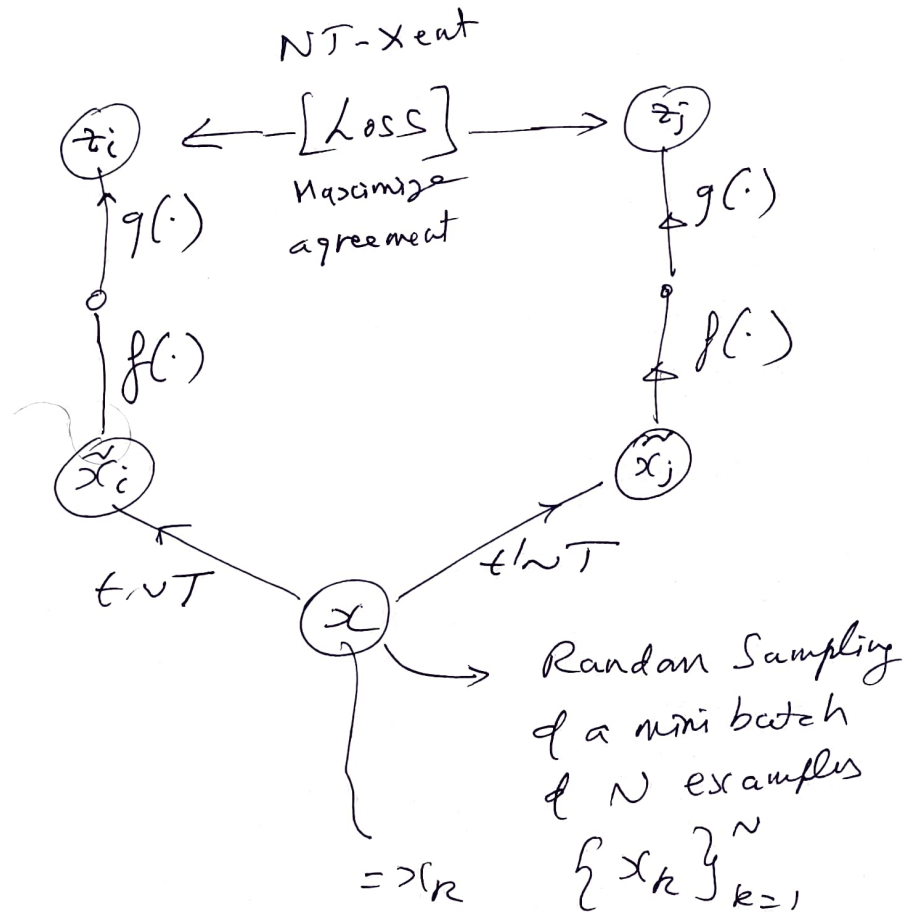
---

**input:** batch size  $N$ , temperature  $\tau$ , form of  $f$ ,  $g$ ,  $\mathcal{T}$ .  
**for** sampled mini-batch  $\{\mathbf{x}_k\}_{k=1}^N$  **do**  
  **for all**  $k \in \{1, \dots, N\}$  **do**  
    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$   
    # the first augmentation  
     $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$   
     $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation  
     $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection  
    # the second augmentation  
     $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$   
     $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation  
     $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection  
  **end for**  
  **for all**  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  **do**  
     $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\tau \|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity  
  **end for**  
  **define**  $\ell(i, j)$  **as**  $-s_{i,j} + \log \sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k})$   
   $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$   
  update networks  $f$  and  $g$  to minimize  $\mathcal{L}$   
**end for**  
**return** encoder network  $f$

---

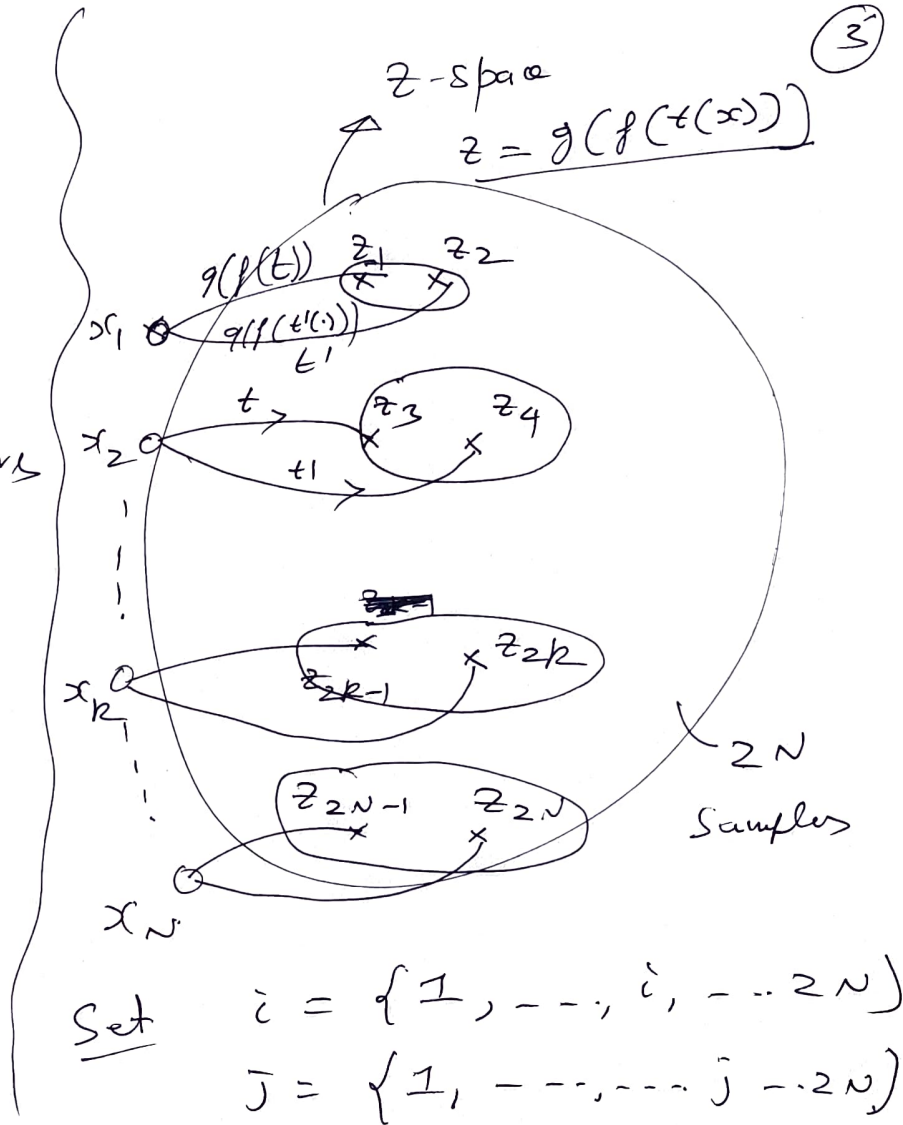
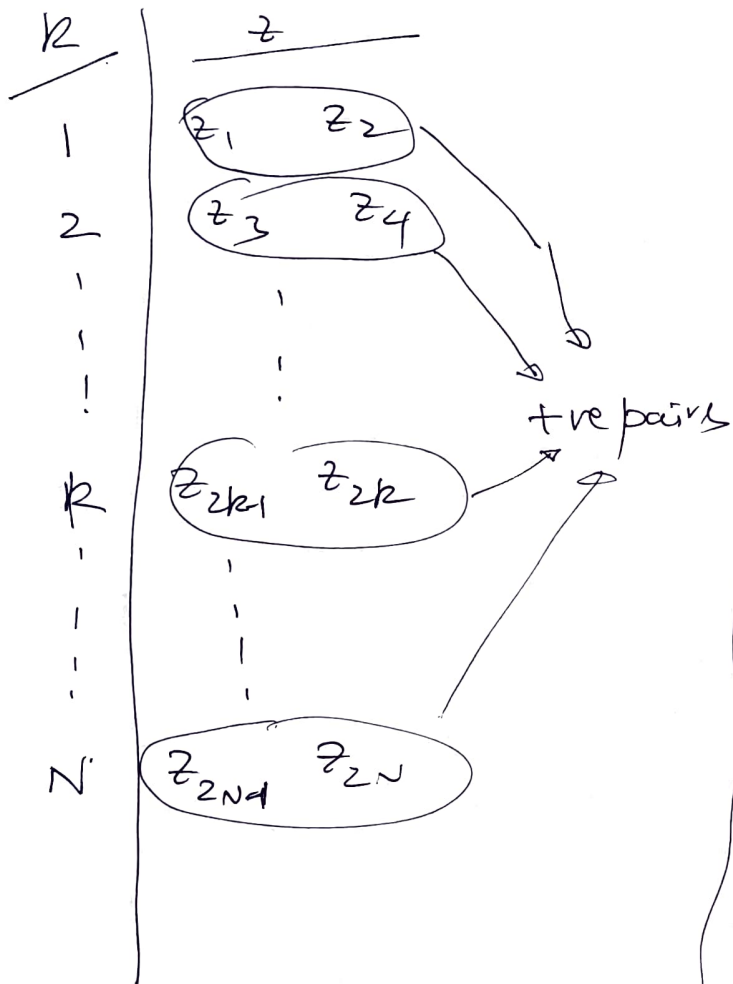
# SimCLR

①



$k$	$x_k$	$t(x_k)$	$t'(x_k)$	$z = g(f(\tilde{x}))$	
1	$x_1$	$\tilde{x}_1$	$\tilde{x}_2$	$z_1$	$z_2$
2	$x_2$	$\tilde{x}_3$	$\tilde{x}_4$	$z_3$	$z_4$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$x_k$	$\tilde{x}_{2k-1}$	$\tilde{x}_{2k}$	$z_{2k-1}$	$z_{2k}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$N$	$x_N$	$\tilde{x}_{2N-1}$	$\tilde{x}_{2N}$	$z_{2N-1}$	$z_{2N}$

Transform  
 $\boxed{t, t'} \rightarrow$   
 Augment  
 $g(f(\cdot))$





$\forall i = \{1, \dots, i, \dots, 2N\}$  &  $j = \{1, \dots, j, \dots, 2N\}$  ④  
 Define & Compute Cosine Similarity [Dot product between  $L_2$  normalized  $z_i, z_j$ ]

$$S_{ij} = \frac{z_i^T z_j}{\|z_i\| \|z_j\|}$$

$$\forall i, j \in [1, \dots, 2N]$$

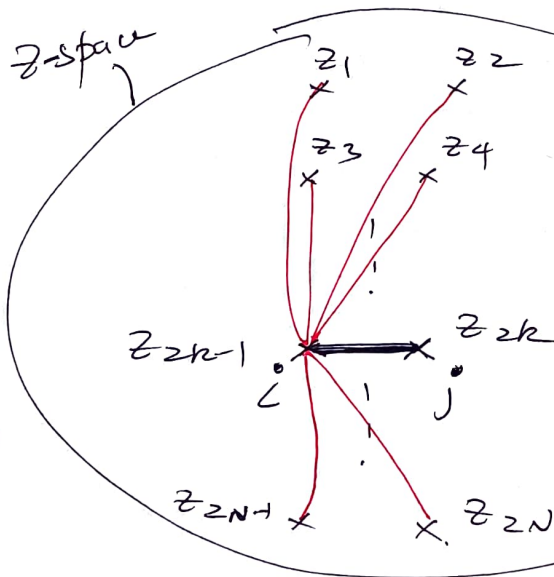
True pairs

$$i = 2k-1$$

$$j = 2k$$

Contrastive loss

- Reduce ~~————~~
- Increase ————

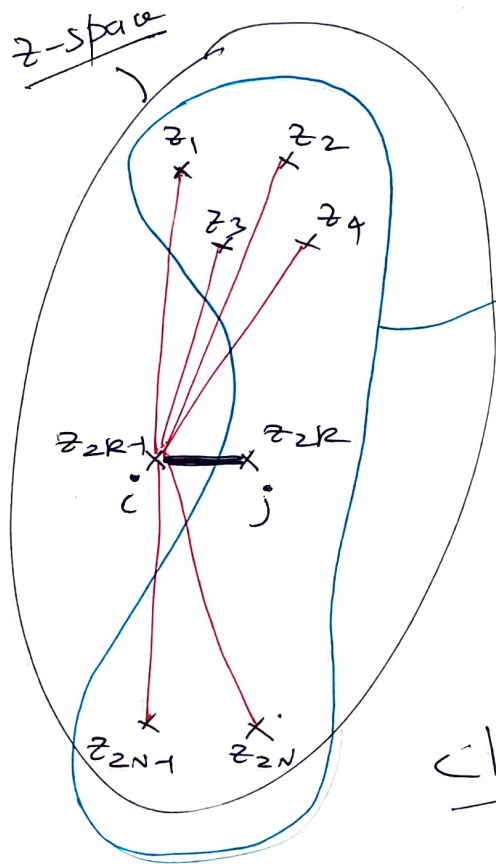


Ensure  $\tilde{x}_j$  is closer to  $\tilde{x}_i$  than all  $\tilde{x}_k, k \neq i, k \neq j$

Identify  $\tilde{x}_j$  in  $\{\tilde{x}_k\}$  for a given  $\tilde{x}_i$



⑤



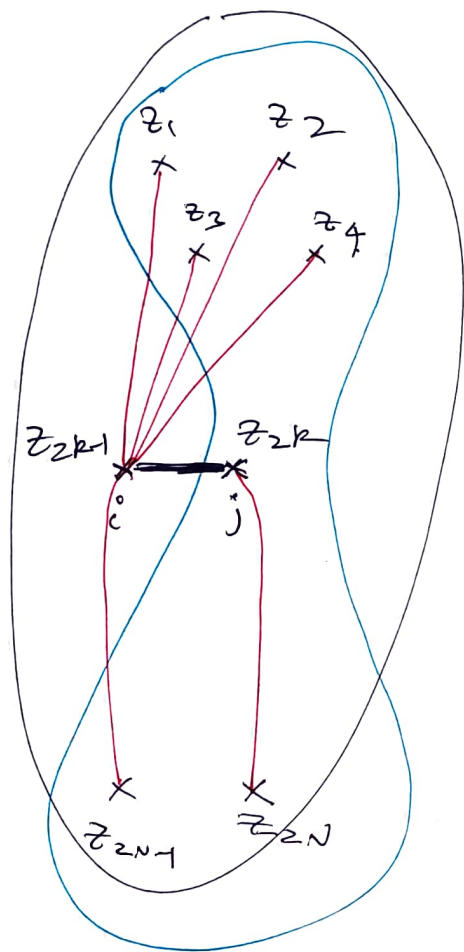
$$l(i, j) = -\log \frac{e^{\frac{\sum_{i,j}}{\tau}}}{\sum_{\substack{m=1 \\ m \neq i}}^{2N} e^{\frac{\sum_{i,m}}{\tau}}}$$

$T$ : Temperature

$\ell(\cdot, \cdot)$ : Normalized Temperature-Scaled  
Cross-Entropy Loss

closer  $i \leftrightarrow j$  or  $j \leftrightarrow i$   $\delta_{ij} \neq 0 \rightarrow 1$

$$\begin{array}{c} \xrightarrow{\quad} e^{\delta_{ij}} \\ \text{Father } i \leftrightarrow m \\ \text{Sim} \downarrow \rightarrow \begin{smallmatrix} 0 \\ -1 \end{smallmatrix} \text{ or } \rightarrow e^{\delta_{im}} \end{array}$$



$$l(i, j) =$$

$$-\log \frac{e^{\frac{\delta_{ij}}{\tau}}}{e^{\frac{\delta_{ij}}{\tau}} + \sum_{\substack{m=1 \\ m \neq j \\ m \neq i}}^{2N-2} e^{\frac{\delta_{im}}{\tau}}}$$

⑥

+ve pair CS

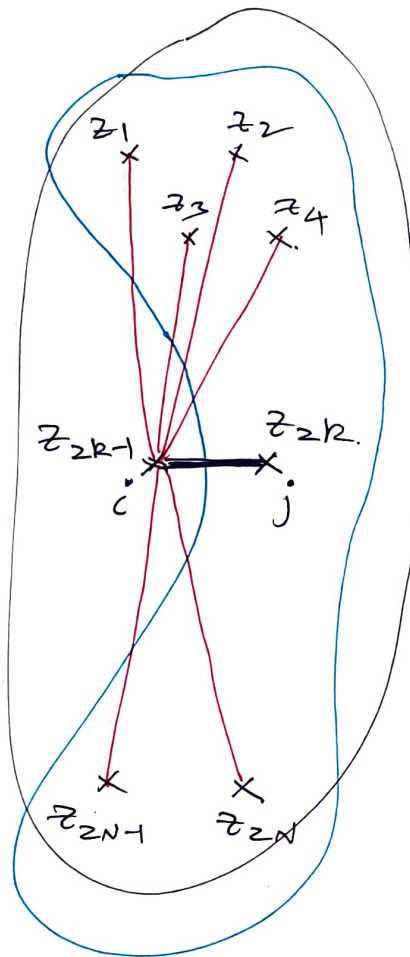
+ve pair CS

-ve pair CS

$$l(i, j) = -\log \frac{CS^+}{CS^+ + CS^-}$$

Sum of  
2N-2  
terms

$$= \frac{\text{Black-line}}{\text{Black-line} + \sum_{2N-2} \text{Red-lines}}$$



Identify  $\tilde{x}_j$  (ie  $z_{2k}$ ) for a given  $\tilde{x}_i \rightarrow z_{2k-1}$  in  $\{\tilde{x}_k\}$  ⑦

Ensure  $\tilde{x}_j$  is closer to  $\tilde{x}_i$  than all  $\tilde{x}_k, k \neq i, k \neq j$

How

1.  $\{x_i, x_j\}$  are the pairs  $x$  closer than  $\{x_i, x_m\}$

$$\frac{CS^+ \text{ (Black-Line)}}{CS^+ + CS^- \text{ (Red-Lines)}} =$$

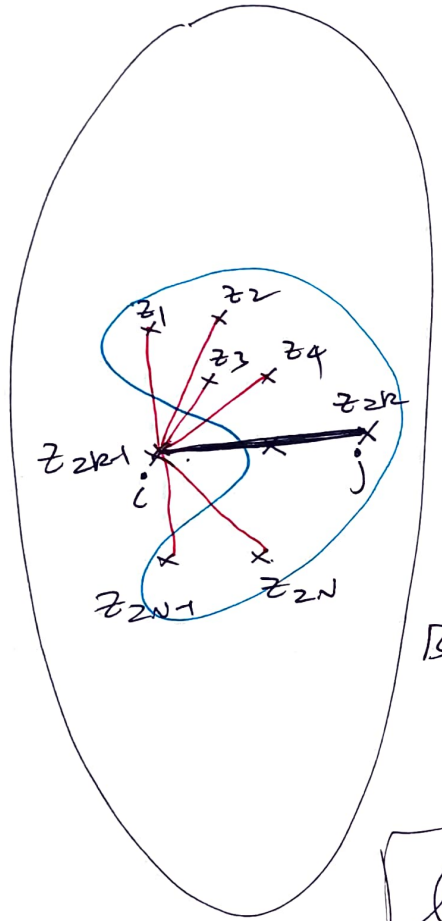
$$\frac{\text{High}}{\text{High} + \text{Low}} < 1 \rightarrow 1$$



⑧

2. when  $\{x_i, x_j\}$  are pairs  
are not yet close

[when  $g(f(\cdot))$  is not effective yet]

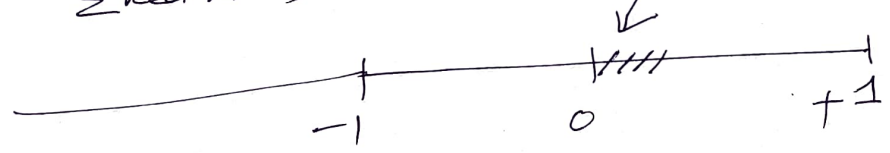


$$\frac{CS^+}{CS^+ + CS^-}$$

/                      \

Black-line                   $\leq$  Red-lines

$$= \frac{\text{Low}}{\text{Low} + \text{Hish.}} \sim 0 \rightarrow 0$$



$$l(i, j) = -\log \frac{CS^+}{CS^+ + CS^-}$$

# Contrastive

+ve pairs w.r.t -ve pairs

Good embeddings of  $f(\cdot)$  eg  $(\cdot)$   
= learning with epochs.

Poor embeddings / earnings  
of  $f(\cdot)$  eg  $(\cdot)$   
Low  
Low + High.

Higher the better

High  
High + Low



Loss

$$l(i, j) = -\log$$

$$\left[ \frac{e^{cs^+/\tau}}{e^{cs^+/\tau} + \sum_{\text{Red lines}} e^{\frac{cs^-}{\tau}}} \right]$$

(2N-2) terms

$$= -\log \left[ \frac{\text{Black-line}}{\text{Black-line} + \sum_{\text{Red lines}} 2N-2 \text{ terms}} \right]$$

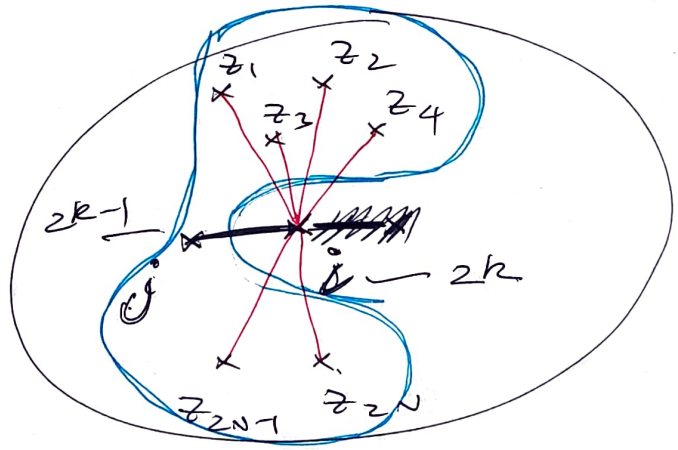
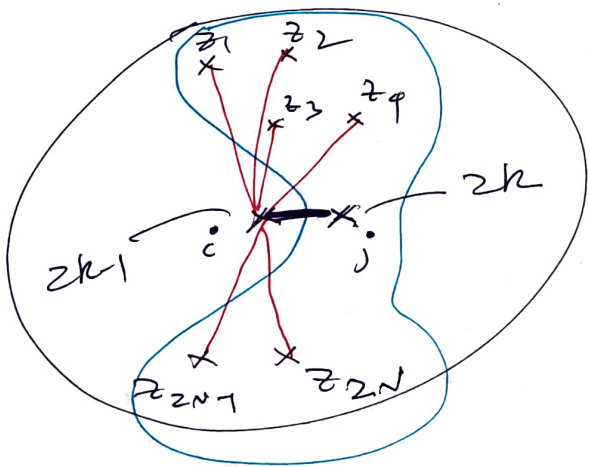
Lower the Better

Since

$$\underline{l(i, j) \neq l(j, i)}$$

Note The Denominator difference on the  $\sum$  red lines

$$L = \frac{1}{2N} \sum_{k=1}^N \left[ \underset{i \quad j}{l(2k-1, 2k)} + \underset{i \quad j}{l(2k, 2k-1)} \right]$$



# SimCLR pseudo code

---

**Algorithm 1** SimCLR's main learning algorithm.

---

**input:** batch size  $N$ , temperature  $\tau$ , form of  $f$ ,  $g$ ,  $\mathcal{T}$ .  
**for** sampled mini-batch  $\{\mathbf{x}_k\}_{k=1}^N$  **do**  
  **for all**  $k \in \{1, \dots, N\}$  **do**  
    draw two augmentation functions  $t \sim \mathcal{T}, t' \sim \mathcal{T}$   
    # the first augmentation  
     $\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$   
     $\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$  # representation  
     $\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$  # projection  
    # the second augmentation  
     $\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$   
     $\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$  # representation  
     $\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$  # projection  
  **end for**  
  **for all**  $i \in \{1, \dots, 2N\}$  and  $j \in \{1, \dots, 2N\}$  **do**  
     $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\tau \|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # pairwise similarity  
  **end for**  
  **define**  $\ell(i, j)$  **as**  $-s_{i,j} + \log \sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k})$   
   $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$   
  update networks  $f$  and  $g$  to minimize  $\mathcal{L}$   
**end for**  
**return** encoder network  $f$

---

Thank you !!