

Paper Review Seminar # 1

(2022. 07. 12)

SimCLR & MoCo (v1 & v2)

통계데이터사이언스학과 통합과정 5학기 이승한

목차

1. Introduction of Contrastive Learning
- 2. SimCLR**
- 3. MoCo v1**
- 4. MoCo v2**
5. Reference

1. Introduction of Contrastive Learning

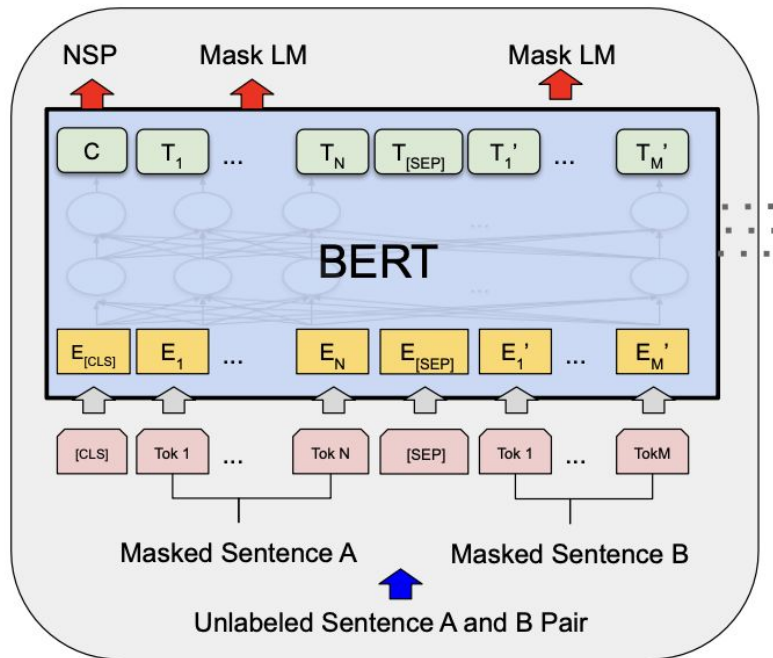
What is **Self-Supervised** Learning?

- Supervised Learning : label O
- Unsupervised Learning : label X
- Semi-supervised Learning : label O & X
- **Self-supervised** Learning : **get label from data itself !**

1. Introduction of Contrastive Learning

Self-Supervised Learning in NLP

- MLM (Masked Language Model)
- NSP (Next Sentence Prediction)



1. Introduction of Contrastive Learning

Self-Supervised Learning in CV

- Pretext Tasks
 - learn representation that could help downstream task
 - ex) Exemplar (2014), Context Prediction (2015), Jigsaw puzzle (2016), ...
- Contrastive Learning
 - make similar/dissimilar data close/far

1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Pretext Tasks

Exemplar (2014)

- crop 32x32 patch, where gradient is significant & data augmentation with this patch
- 1 instance = 1 class
- cons) unsuitable for large dataset

$$L(X) = \sum_{\mathbf{x}_i \in X} \sum_{T \in T_i} l(i, T\mathbf{x}_i)$$

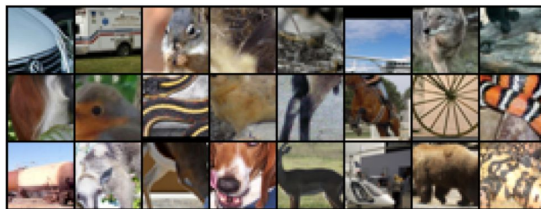


Fig. 1. Exemplary patches sampled from the STL unlabeled dataset which are later augmented by various transformations to obtain surrogate data for the CNN training.



Fig. 2. Several random transformations applied to one of the patches extracted from the STL unlabeled dataset. The original ('seed') patch is in the top left corner.

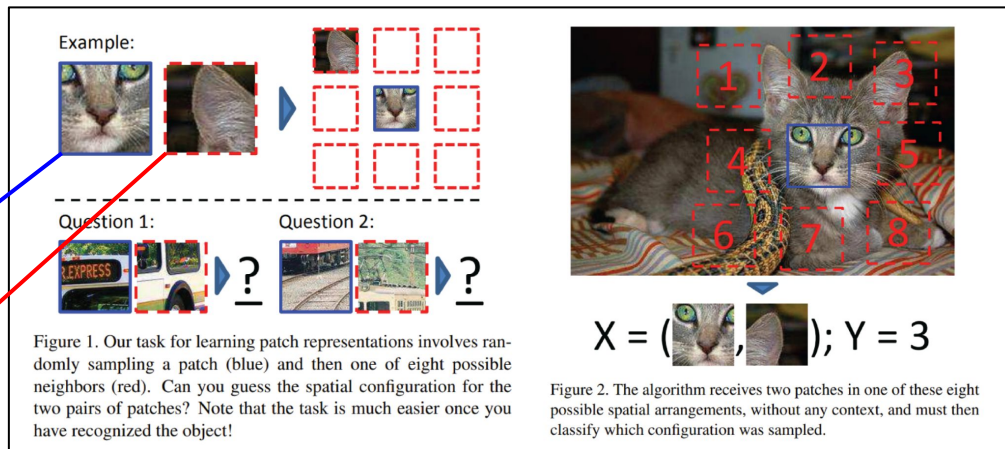
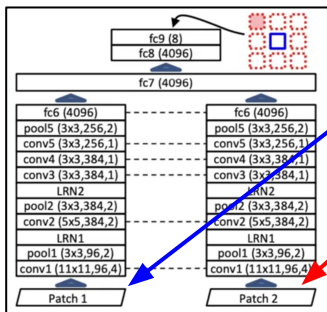
Discriminative Unsupervised Feature Learning with Exemplar Convolutional Neural Networks (Dosovitskiy, NIPS 2014)

1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Pretext Tasks

Context Prediction (2015)

- guess the relative location (1~8)
- cons) cheating with texture/boundary



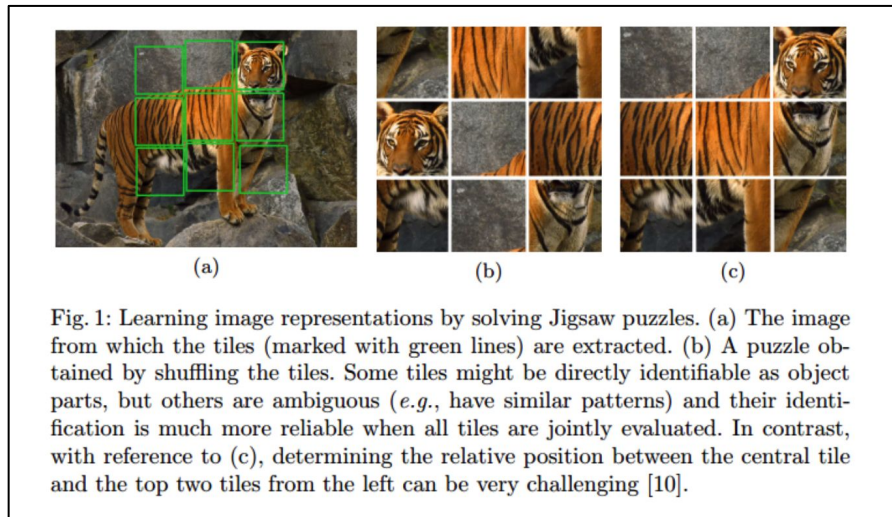
Unsupervised Visual Representation Learning by Context Prediction (Doersch et al., ICCV 2015)

1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Pretext Tasks

Jigsaw Puzzle (2016)

- solve jigsaw puzzle
- number of class :
 - original) $9! = 362,880$
 - proposed) 100
(remove similar permutation)



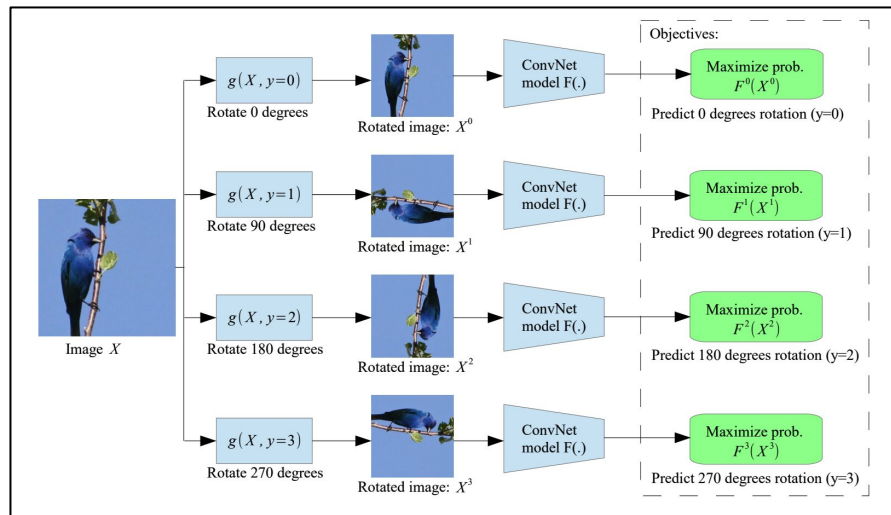
Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles (Noroozi and Favaro, ECCV 2016)

1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Pretext Tasks

Rotation Prediction (2016)

- rotate image & guess the rotation angle

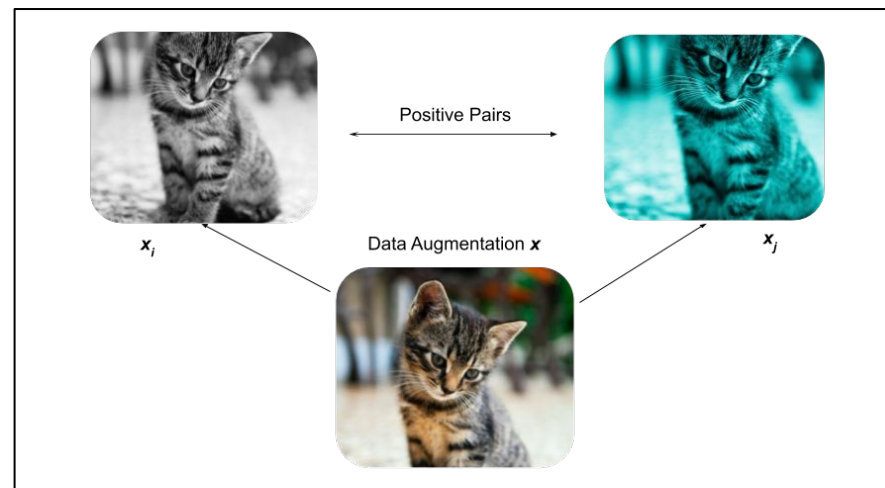
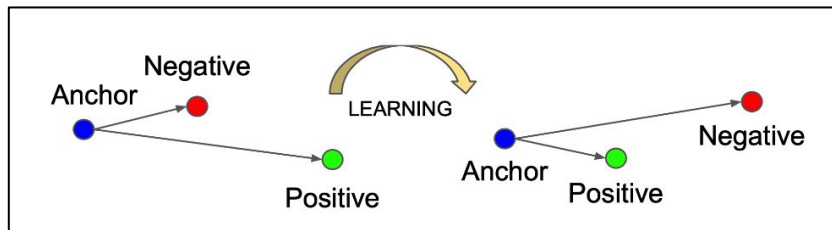


Unsupervised Representation Learning by Predicting Image Rotations (Gidaris et al., ICLR 2018)

1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Contrastive Learning

- make **similar** data **close** to each other
make **dissimilar** data **far** from each other
- various loss functions



1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Contrastive Learning

(1) Contrastive Loss

$$\mathcal{L}_{\text{cont}}(\mathbf{x}_i, \mathbf{x}_j, \theta) = \mathbb{1}[y_i = y_j] \|f_{\theta}(\mathbf{x}_i) - f_{\theta}(\mathbf{x}_j)\|_2^2 + \mathbb{1}[y_i \neq y_j] \max(0, \epsilon - \|f_{\theta}(\mathbf{x}_i) - f_{\theta}(\mathbf{x}_j)\|_2)^2$$

(2) Triplet Loss

$$\mathcal{L}_{\text{triplet}}(\mathbf{x}, \mathbf{x}^+, \mathbf{x}^-) = \sum_{\mathbf{x} \in \mathcal{X}} \max(0, \|f(\mathbf{x}) - f(\mathbf{x}^+)\|_2^2 - \|f(\mathbf{x}) - f(\mathbf{x}^-)\|_2^2 + \epsilon)$$

(3) Noise Contrastive Estimation

$$\mathcal{L}_{\text{NCE}} = -\frac{1}{N} \sum_{i=1}^N [\log \sigma(\ell_{\theta}(\mathbf{x}_i)) + \log(1 - \sigma(\ell_{\theta}(\tilde{\mathbf{x}}_i)))]$$

$$\text{where } \sigma(\ell) = \frac{1}{1 + \exp(-\ell)} = \frac{p_{\theta}}{p_{\theta} + q}$$

target sample $\sim P(\mathbf{x}|C = 1; \theta) = p_{\theta}(\mathbf{x})$

noise sample $\sim P(\tilde{\mathbf{x}}|C = 0) = q(\tilde{\mathbf{x}})$

1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Contrastive Learning

(1) Contrastive Loss

$$\mathcal{L}_{\text{cont}}(\mathbf{x}_i, \mathbf{x}_j, \theta) = \mathbb{1}[y_i = y_j] \boxed{\|f_{\theta}(\mathbf{x}_i) - f_{\theta}(\mathbf{x}_j)\|_2^2} + \mathbb{1}[y_i \neq y_j] \max(0, \epsilon - \boxed{\|f_{\theta}(\mathbf{x}_i) - f_{\theta}(\mathbf{x}_j)\|_2})^2$$

(2) Triplet Loss

$$\mathcal{L}_{\text{triplet}}(\mathbf{x}, \mathbf{x}^+, \mathbf{x}^-) = \sum_{\mathbf{x} \in \mathcal{X}} \max(0, \boxed{\|f(\mathbf{x}) - f(\mathbf{x}^+)\|_2^2} - \boxed{\|f(\mathbf{x}) - f(\mathbf{x}^-)\|_2^2} + \epsilon)$$

 : positive (= similar pair)
 : negative (= dissimilar pair)

(3) Noise Contrastive Estimation

$$\mathcal{L}_{\text{NCE}} = -\frac{1}{N} \sum_{i=1}^N [\log \sigma(\ell_{\theta}(\mathbf{x}_i)) + \log(1 - \sigma(\ell_{\theta}(\tilde{\mathbf{x}}_i)))]$$

$$\text{where } \sigma(\ell) = \frac{1}{1 + \exp(-\ell)} = \frac{p_{\theta}}{p_{\theta} + q}$$

target sample $\sim P(\mathbf{x}|C = 1; \theta) = p_{\theta}(\mathbf{x})$

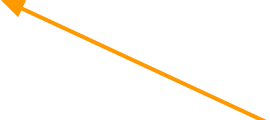
noise sample $\sim P(\tilde{\mathbf{x}}|C = 0) = q(\tilde{\mathbf{x}})$

1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Contrastive Learning

(4) InfoNCE

$$\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E} \left[\log \frac{f(\mathbf{x}, \mathbf{c})}{\sum_{\mathbf{x}' \in X} f(\mathbf{x}', \mathbf{c})} \right]$$


$$p(C = \text{pos} | X, \mathbf{c}) = \frac{p(x_{\text{pos}} | \mathbf{c}) \prod_{i=1, \dots, N; i \neq \text{pos}} p(\mathbf{x}_i)}{\sum_{j=1}^N [p(\mathbf{x}_j | \mathbf{c}) \prod_{i=1, \dots, N; i \neq j} p(\mathbf{x}_i)]} = \frac{\frac{p(\mathbf{x}_{\text{pos}} | \mathbf{c})}{p(\mathbf{x}_{\text{pos}})}}{\sum_{j=1}^N \frac{p(\mathbf{x}_j | \mathbf{c})}{p(\mathbf{x}_j)}} = \frac{f(\mathbf{x}_{\text{pos}}, \mathbf{c})}{\sum_{j=1}^N f(\mathbf{x}_j, \mathbf{c})}$$

1. Introduction of Contrastive Learning

Self-Supervised Learning in CV - Contrastive Learning

(4) InfoNCE

$$\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E} \left[\log \frac{f(\mathbf{x}, \mathbf{c})}{\sum_{\mathbf{x}' \in X} f(\mathbf{x}', \mathbf{c})} \right]$$

$$p(C = \text{pos} | X, \mathbf{c}) = \frac{p(x_{\text{pos}} | \mathbf{c}) \prod_{i=1, \dots, N; i \neq \text{pos}} p(\mathbf{x}_i)}{\sum_{j=1}^N [p(\mathbf{x}_j | \mathbf{c}) \prod_{i=1, \dots, N; i \neq j} p(\mathbf{x}_i)]} = \frac{\frac{p(\mathbf{x}_{\text{pos}} | \mathbf{c})}{p(\mathbf{x}_{\text{pos}})}}{\sum_{j=1}^N \frac{p(\mathbf{x}_j | \mathbf{c})}{p(\mathbf{x}_j)}} = \frac{f(\mathbf{x}_{\text{pos}}, \mathbf{c})}{\sum_{j=1}^N f(\mathbf{x}_j, \mathbf{c})}$$

1 positive

1 positive + (N-1) negatives

SimCLR (2020)

A Simple Framework for Contrastive Learning of Visual Representations

Ting Chen¹ Simon Kornblith¹ Mohammad Norouzi¹ Geoffrey Hinton¹

Abstract

This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently proposed contrastive self-supervised learning algorithms without requiring specialized architectures or a memory bank. In order to understand what enables the contrastive prediction tasks to learn useful representations, we systematically study the major components of our framework. We show that (1) composition of data augmentations plays a critical role in defining effective predictive tasks, (2) introducing a learnable nonlinear transformation between the representation and the contrastive loss substantially improves the quality of the learned representations, and (3) contrastive learning benefits from larger batch sizes and more training steps compared to supervised learning. By combining these findings,

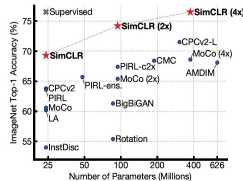


Figure 1. ImageNet Top-1 accuracy of linear classifiers trained on representations learned with different self-supervised methods (pretrained on ImageNet). Gray cross indicates supervised ResNet-50. Our method, SimCLR, is shown in bold.

<https://arxiv.org/abs/2002.05709>

A Simple Framework for Contrastive Learning of Visual ... - arXiv

by T Chen · 2020 · Cited by 5112 — Abstract: This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently proposed ...

Cite as: arXiv:2002.05709

MoCo v1 (2019)

Momentum Contrast for Unsupervised Visual Representation Learning

Kaiming He Haoqi Fan Yuxin Wu Saining Xie Ross Girshick

Facebook AI Research (FAIR)

Code: <https://github.com/facebookresearch/moco>

Abstract

We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning [29] as dictionary look-up, we build a dynamic dictionary with a queue and a moving-averaged encoder. This enables building a large and consistent dictionary on-the-fly that facilitates contrastive unsupervised learning. MoCo provides competitive results under the common linear protocol on ImageNet classification. More importantly, the representations learned by MoCo transfer well to downstream tasks. MoCo can outperform its supervised pre-training counterpart in 7 detection/segmentation tasks on PASCAL VOC, COCO, and other datasets, sometimes surpassing it by large margins. This suggests that the gap between unsupervised and supervised representation learning has been largely closed in many vision tasks.

1. Introduction

Unsupervised representation learning is highly successful in natural language processing, e.g., as shown by GPT [50, 51] and BERT [12]. But supervised pre-training is still

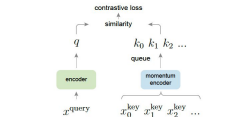


Figure 1. Momentum Contrast (MoCo) trains a visual representation encoder by matching an encoded query q to a dictionary of encoded keys by a contrastive loss. The dictionary keys $\{k_0, k_1, k_2, \dots\}$ are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

<https://arxiv.org/abs/1911.05722>

Momentum Contrast for Unsupervised Visual Representation ...

by K He · 2019 · Cited by 3735 — Abstract: We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning as ...

Cite as: arXiv:1911.05722

MoCo v2 (2020)

Improved Baselines with Momentum Contrastive Learning

Xinlei Chen Haoqi Fan Ross Girshick Kaiming He
Facebook AI Research (FAIR)

Abstract

Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR. In this note, we verify the effectiveness of two of SimCLR's design improvements by implementing them in the MoCo framework. With simple modifications to MoCo—namely, using an MLP projection head and more data augmentation—we establish stronger baselines that outperform SimCLR and do not require large training batches. We hope this will make state-of-the-art unsupervised learning research more accessible. Code will be made public.

1. Introduction

Recent studies on unsupervised representation learning from images [16, 13, 8, 17, 1, 9, 15, 6, 12, 2] are converging on a central concept known as *contrastive learning* [5]. The results are promising: e.g., Momentum Contrast (MoCo) [6] shows that unsupervised pre-training can surpass its ImageNet-supervised counterpart in multiple detection and segmentation tasks, and SimCLR [2] further reduces the gap in linear classifier performance between unsupervised and supervised pre-training representations.

This note establishes stronger and more feasible base-

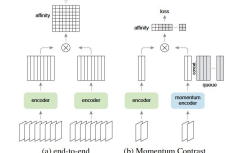


Figure 1. A batching perspective of two optimization mechanisms for contrastive learning. Images are encoded into a representation space, in which pairwise affinities are computed.

ffective contrastive loss function, called InfoNCE [13], is:

$$\mathcal{L}_{\text{InfoNCE}}(x, y) = -\log \frac{\exp(q(x, y)/r)}{\exp(q(x, y)/r) + \sum_{k \neq y} \exp(q(x, k)/r)} \quad (1)$$

<https://arxiv.org/abs/2003.04297>

Improved Baselines with Momentum Contrastive Learning

by X Chen · 2020 · Cited by 1088 — Abstract: Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR.

Cite as: arXiv:2003.04297

SimCLR (2020)

A Simple Framework for Contrastive Learning of Visual Representations

Ting Chen¹ Simon Kornblith¹ Mohammad Norouzi¹ Geoffrey Hinton¹

Abstract

This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently proposed contrastive self-supervised learning algorithms without requiring specialized architectures or a memory bank. In order to understand what enables the contrastive prediction tasks to learn useful representations, we systematically study the major components of our framework. We show that (1) composition of data augmentations plays a critical role in defining effective predictive tasks, (2) introducing a learnable nonlinear transformation between the representation and the contrastive loss substantially improves the quality of the learned representations, and (3) contrastive learning benefits from larger batch sizes and more training steps compared to supervised learning. By combining these findings,

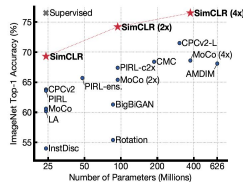


Figure 1. ImageNet Top-1 accuracy of linear classifiers trained on representations learned with different self-supervised methods (pretrained on ImageNet). Gray cross indicates supervised ResNet-50. Our method, SimCLR, is shown in bold.

<https://arxiv.org/abs/2002.05709>

A Simple Framework for Contrastive Learning of Visual ... - arXiv

by T Chen · 2020 · Cited by 5112 — Abstract: This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently proposed ...

Cite as: arXiv:2002.05709

MoCo v1 (2019)

Momentum Contrast for Unsupervised Visual Representation Learning

Kaiming He Haoqi Fan Yuxin Wu Saining Xie Ross Girshick

Facebook AI Research (FAIR)

Code: <https://github.com/facebookresearch/moco>

Abstract

We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning [29] as dictionary look-up, we build a dynamic dictionary with a queue and a moving-averaged encoder. This enables building a large and consistent dictionary on-the-fly that facilitates contrastive unsupervised learning. MoCo provides competitive results under the common linear protocol on ImageNet classification. More importantly, the representations learned by MoCo transfer well to downstream tasks. MoCo can outperform its supervised pre-training counterpart in 7 detection/segmentation tasks on PASCAL VOC, COCO, and other datasets, sometimes surpassing it by large margins. This suggests that the gap between unsupervised and supervised representation learning has been largely closed in many vision tasks.

1. Introduction

Unsupervised representation learning is highly successful in natural language processing, e.g., as shown by GPT [50, 51] and BERT [12]. But supervised pre-training is still

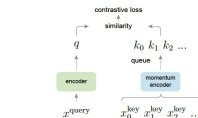


Figure 1. Momentum Contrast (MoCo) trains a visual representation of encoded keys using a contrastive loss. The dictionary keys $\{k_0, k_1, k_2, \dots\}$ are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

<https://arxiv.org/abs/1911.05722>

Momentum Contrast for Unsupervised Visual Representation ...

by K He · 2019 · Cited by 3735 — Abstract: We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning as ...

Cite as: arXiv:1911.05722

MoCo v2 (2020)

Improved Baselines with Momentum Contrastive Learning

Xinlei Chen Haoqi Fan Ross Girshick Kaiming He

Facebook AI Research (FAIR)

Abstract

Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR. In this note, we verify the effectiveness of two of SimCLR's design improvements by implementing them in the MoCo framework. With simple modifications to MoCo—namely, using an MLP projection head and more data augmentation—we establish stronger baselines that outperform SimCLR and do not require large training batches. We hope this will make state-of-the-art unsupervised learning research more accessible. Code will be made public.

1. Introduction

Recent studies on unsupervised representation learning from images [16, 13, 8, 17, 1, 9, 15, 6, 12, 2] are converging on a central concept known as *contrastive learning* [5]. The results are promising: e.g., Momentum Contrast (MoCo) [6] shows that unsupervised pre-training can surpass its ImageNet-supervised counterpart in multiple detection and segmentation tasks, and SimCLR [2] further reduces the gap in linear classifier performance between unsupervised and supervised pre-training representations.

This note establishes stronger and more feasible base-

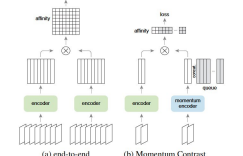


Figure 1. A batching perspective of two optimization mechanisms for contrastive learning. Images are encoded into a representation space, in which pairwise affinities are computed.

ffective contrastive loss function, called InfoNCE [13], is:

$$\mathcal{L}_{\text{InfoNCE}}(x, y) = -\log \frac{\exp(q \cdot k^+ / r)}{\exp(q \cdot k^+ / r) + \sum_{k^-} \exp(q \cdot k^- / r)}, \quad (1)$$

<https://arxiv.org/abs/2003.04297>

Improved Baselines with Momentum Contrastive Learning

by X Chen · 2020 · Cited by 1088 — Abstract: Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR.

Cite as: arXiv:2003.04297

2. SimCLR

SimCLR = **Sim**ple Framework for **C**ontrastive **L**earning of Visual **R**epresentation

- contrastive **SELF-SUPERVISED** learning algorithm

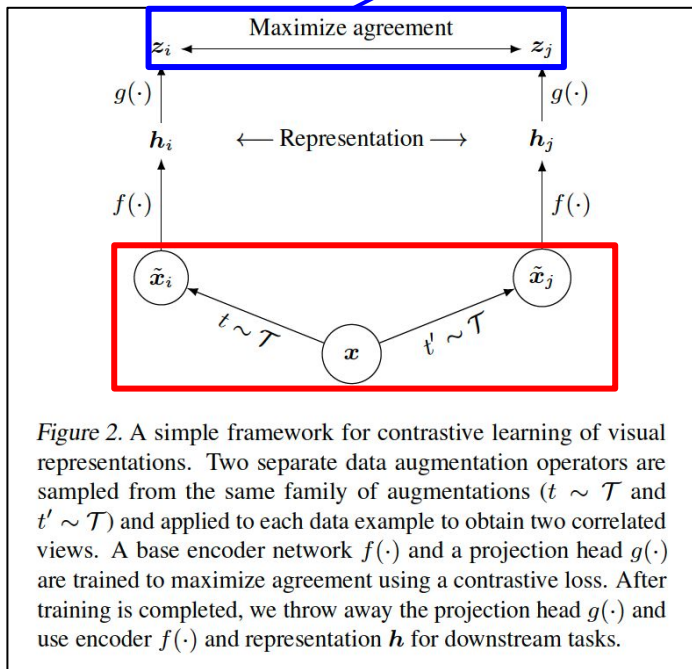
Three Findings

- (1) **composition of data augmentations**
- (2) **learnable non-linear transformation** between representation & contrastive loss
- (3) contrastive learning benefits from ..
 - **larger batch sizes**
 - **more training steps**

2. SimCLR

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

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



learns representation by

maximizing agreement between

differently augmented versions of **same data**,

with **contrastive loss**

Figure 2. A simple framework for contrastive learning of visual representations. Two separate data augmentation operators are sampled from the same family of augmentations ($t \sim \mathcal{T}$ and $t' \sim \mathcal{T}$) and applied to each data example to obtain two correlated views. A base encoder network $f(\cdot)$ and a projection head $g(\cdot)$ are trained to maximize agreement using a contrastive loss. After training is completed, we throw away the projection head $g(\cdot)$ and use encoder $f(\cdot)$ and representation \mathbf{h} for downstream tasks.

2. SimCLR

4 Major Components

1. Stochastic Data Augmentation
2. Base Encoder
3. Projection head
4. Contrastive Loss Function

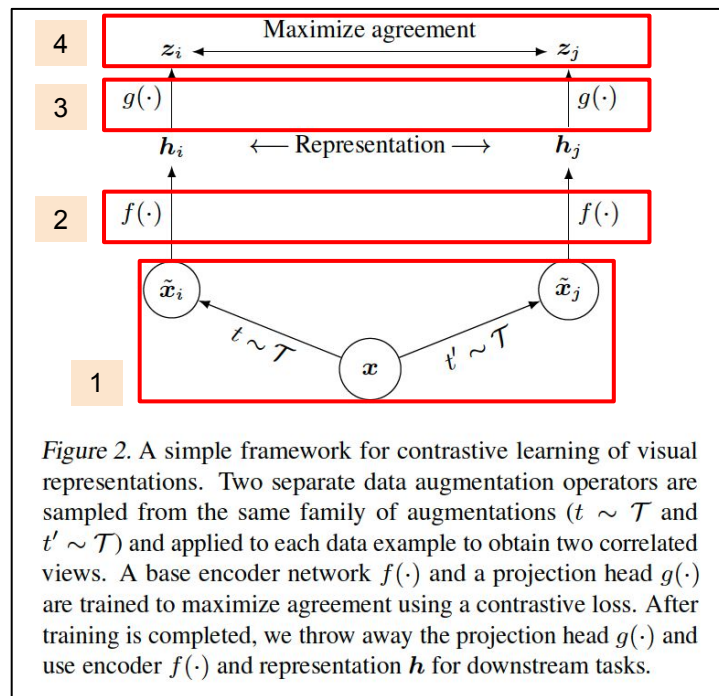


Figure 2. A simple framework for contrastive learning of visual representations. Two separate data augmentation operators are sampled from the same family of augmentations ($t \sim \mathcal{T}$ and $t' \sim \mathcal{T}$) and applied to each data example to obtain two correlated views. A base encoder network $f(\cdot)$ and a projection head $g(\cdot)$ are trained to maximize agreement using a contrastive loss. After training is completed, we throw away the projection head $g(\cdot)$ and use encoder $f(\cdot)$ and representation h for downstream tasks.

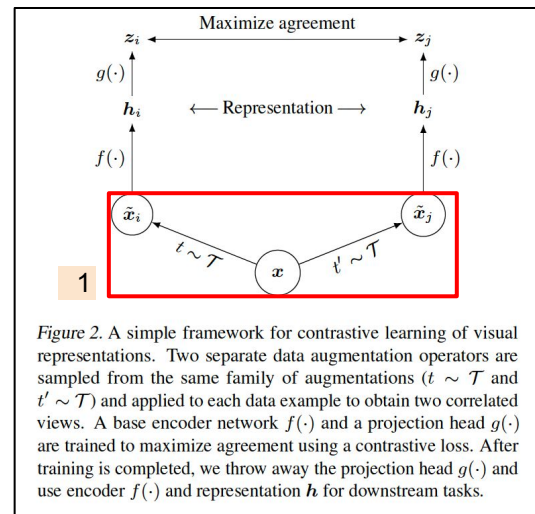
2. SimCLR

4 Major Components

1. Stochastic Data Augmentation

- positive pair \tilde{x}_i & \tilde{x}_j
- apply 3 simple augmentations
 - (1) random cropping
 - (2) random color distortions
 - (3) random Gaussian blur

→ (1) + (2) : good performance!

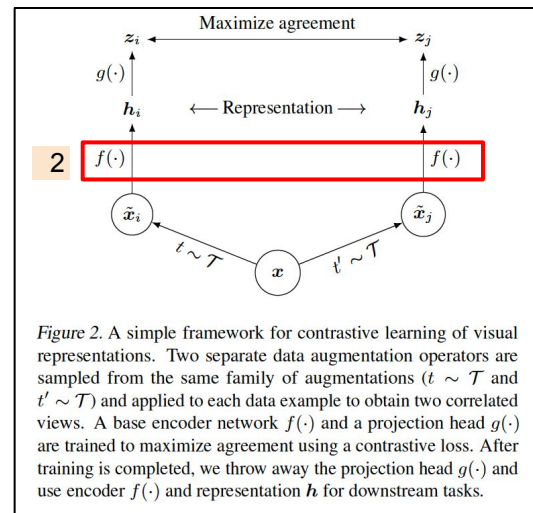


2. SimCLR

4 Major Components

2. Base Encoder $f(\cdot)$

- $\mathbf{h}_i = f(\tilde{\mathbf{x}}_i) = \text{ResNet}(\tilde{\mathbf{x}}_i)$.
 - where $\mathbf{h}_i \in \mathbb{R}^d$ is the output after the GAP
 - extract representations from **augmented data samples**
- use ResNet

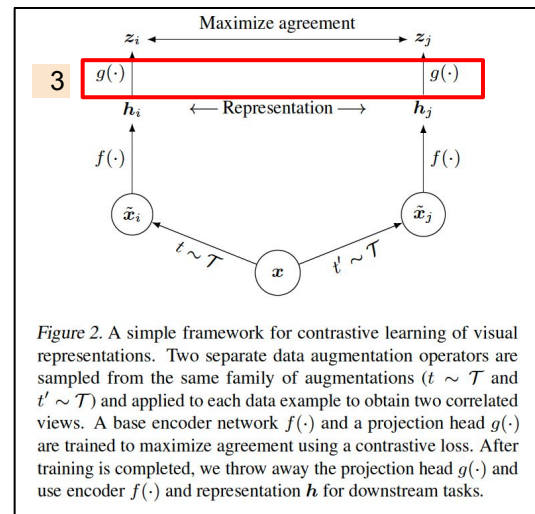


2. SimCLR

4 Major Components

3. Projection head $g(\cdot)$

- $z_i = g(h_i) = W^{(2)} \sigma(W^{(1)} h_i)$.
 - maps representations to the space **where contrastive loss is applied**
- use MLP with 1 hidden layer



2. SimCLR

4 Major Components

4. Contrastive Loss Function

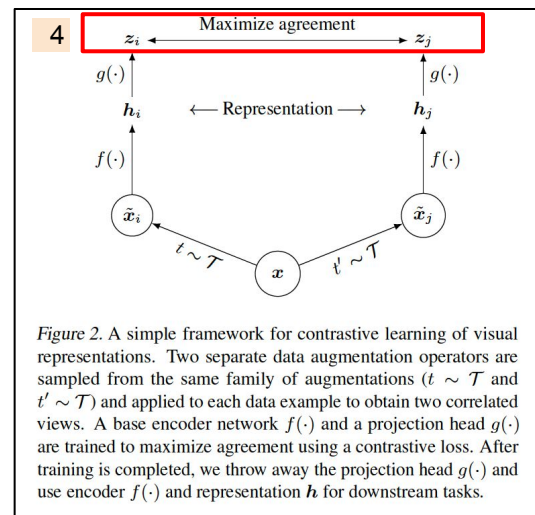
- Data : set $\{\tilde{\mathbf{x}}_k\}$ including a positive pair ($\tilde{\mathbf{x}}_i$ and $\tilde{\mathbf{x}}_j$)
- Task : aims to identify $\tilde{\mathbf{x}}_j$ in $\{\tilde{\mathbf{x}}_k\}_{k \neq i}$ for a given $\tilde{\mathbf{x}}_i$.

Sample mini batches of size N

→ 2 augmentations → $2N$ data points

(no negative samples ...only positive pairs)

- Just treat $2(N - 1)$ augmented samples within a mini-batch as negative examples.



2. SimCLR

4 Major Components

4. Contrastive Loss Function

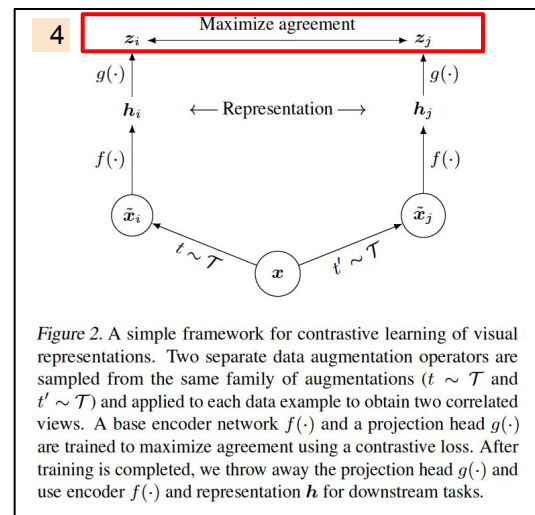
Loss Function for a positive pair of examples (NT-Xent)

(= Normalized Temperature scaled CE loss)

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

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

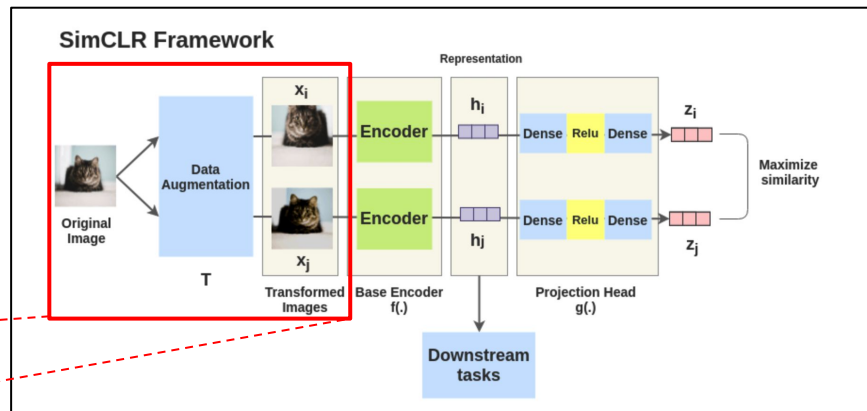
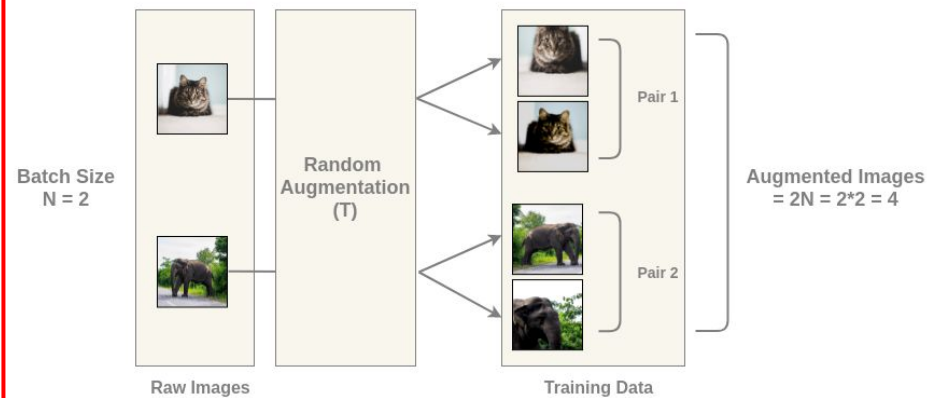
→ final loss : computed across all positive pairs



2. SimCLR

Example

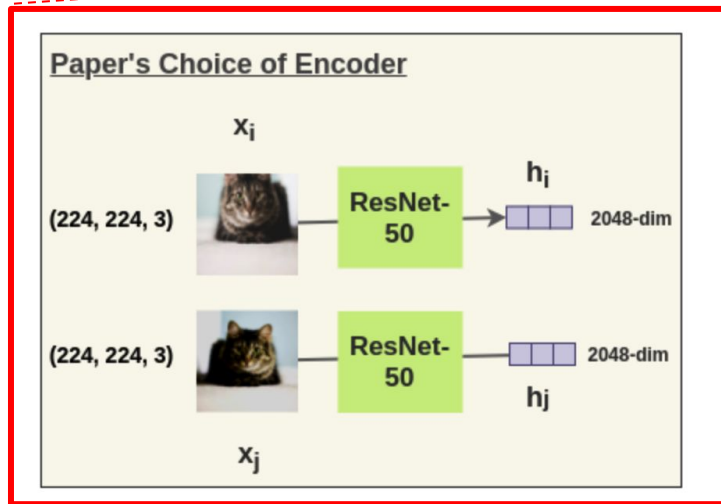
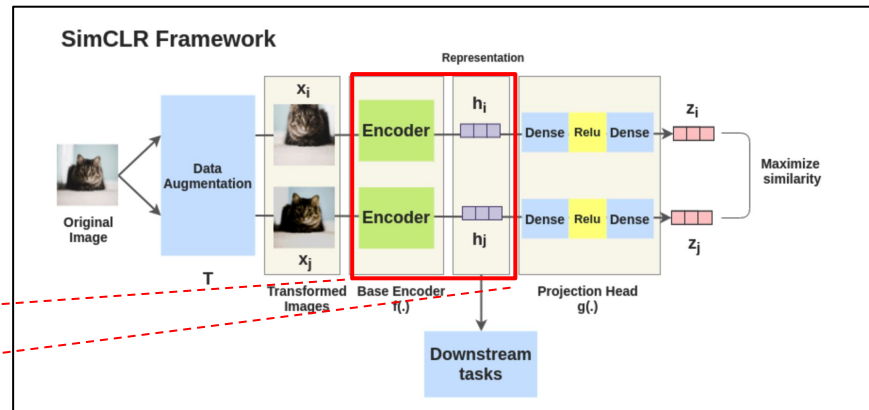
Preparing similar pairs in a batch



Step 1) Data Augmentation

2. SimCLR

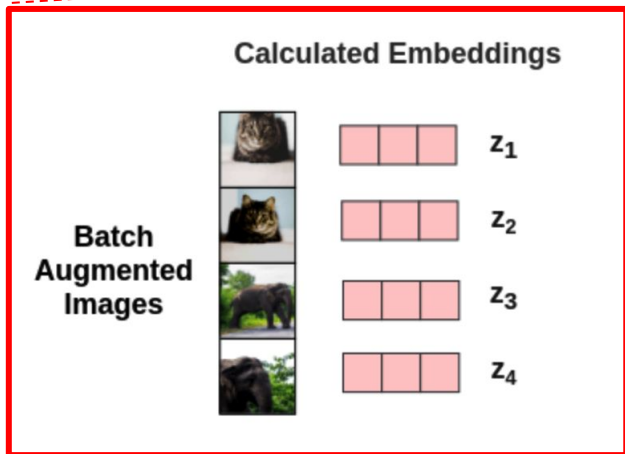
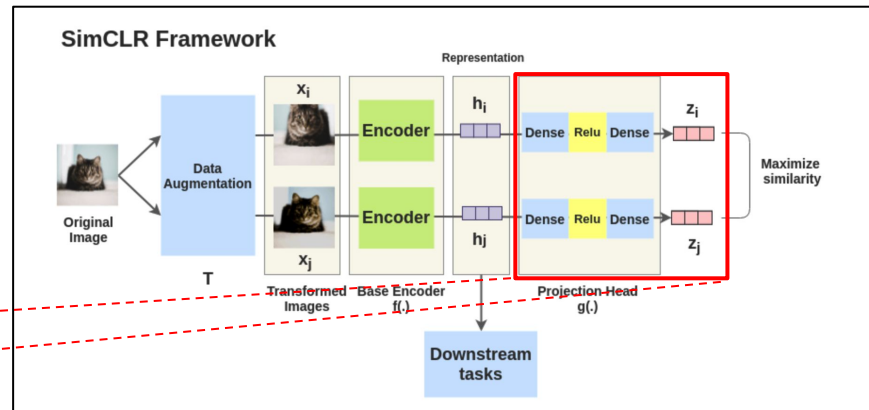
Example



Step 2) Embedding with Base Encoder (= ResNet)

2. SimCLR

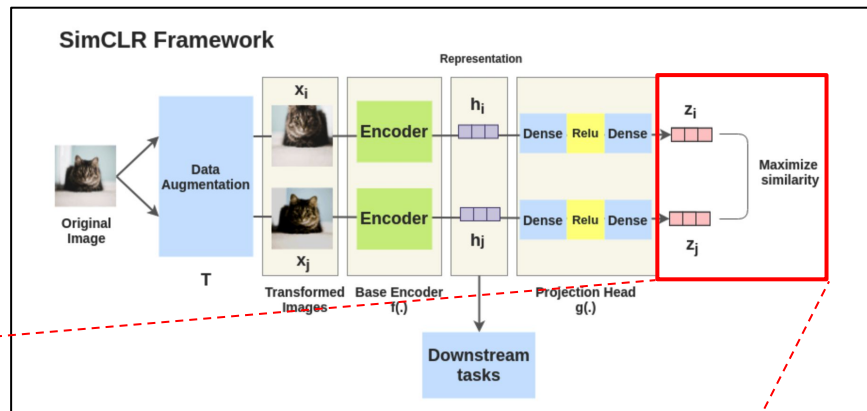
Example



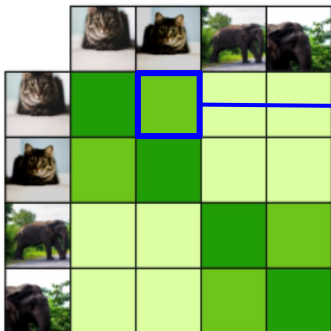
Step 3) Embedding with Projection Head

2. SimCLR

Example



Pairwise cosine similarity



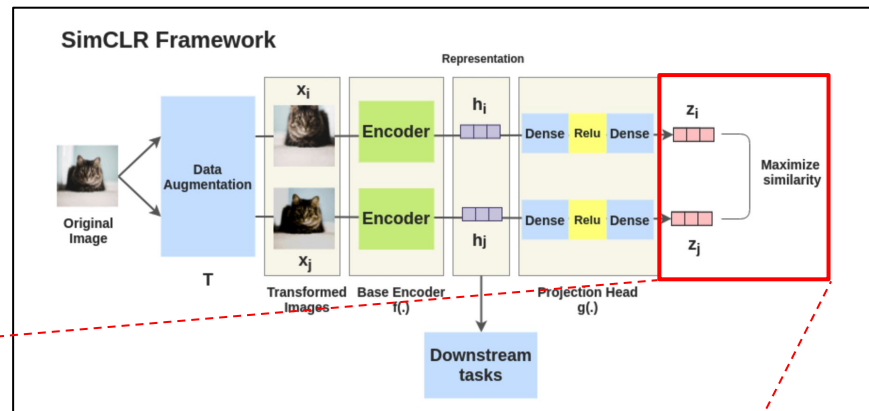
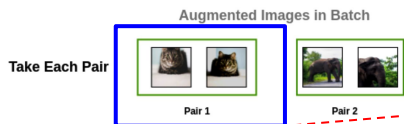
Similarity Calculation of Augmented Images

$$\text{similarity}(\begin{matrix} x_i \\ \text{cat} \end{matrix}, \begin{matrix} x_j \\ \text{cat} \end{matrix}) = \text{cosine similarity} \left(\begin{matrix} z_i \\ \text{cat} \end{matrix}, \begin{matrix} z_j \\ \text{cat} \end{matrix} \right)$$

$$s_{i,j} = \frac{z_i^T z_j}{(\tau \|z_i\| \|z_j\|)}$$

2. SimCLR

Example



$$l(i, j) = -\log \frac{\exp(s_{i,j})}{\sum_{k=1}^{2N} l_{[k \neq i]} \exp(s_{i,k})}$$

$$l(\text{cat}_1, \text{cat}_2) = -\log \left(\frac{\exp(\text{similarity}(\text{cat}_1, \text{cat}_2))}{\exp(\text{similarity}(\text{cat}_1, \text{cat}_2)) + \exp(\text{similarity}(\text{cat}_1, \text{elephant}_1)) + \exp(\text{similarity}(\text{cat}_1, \text{elephant}_2))} \right)$$

Interchanged

$$l(\text{cat}_2, \text{cat}_1) = -\log \left(\frac{\exp(\text{similarity}(\text{cat}_2, \text{cat}_1))}{\exp(\text{similarity}(\text{cat}_2, \text{cat}_1)) + \exp(\text{similarity}(\text{cat}_2, \text{elephant}_1)) + \exp(\text{similarity}(\text{cat}_2, \text{elephant}_2))} \right)$$

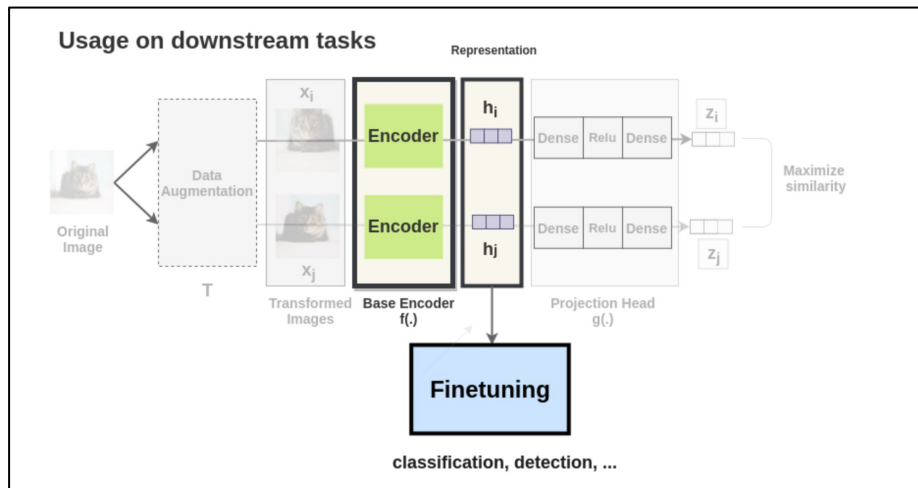
$$L = \frac{1}{2N} \sum_{k=1}^N [l(2k-1, 2k) + l(2k, 2k-1)]$$

Pair 1 Loss (k=1)

Pair 2 Loss (k=2)

$$L = \frac{[l(\text{cat}_1, \text{cat}_2) + l(\text{cat}_2, \text{cat}_1)] + [l(\text{elephant}_1, \text{elephant}_2) + l(\text{elephant}_2, \text{elephant}_1)]}{2 * 2}$$

2. SimCLR



Representation obtained from **base encoder (not projection head)** can be used for other tasks!
(will be discussed in the Experiment Section)

2. SimCLR

PseudoCode

Algorithm 1 SimCLR's main learning algorithm.

input: batch size N , constant τ , structure of f, g, \mathcal{T} .
for sampled minibatch $\{x_k\}_{k=1}^N$ **do**
 for all $k \in \{1, \dots, N\}$ **do**
 draw two augmentation functions $t \sim \mathcal{T}, t' \sim \mathcal{T}$
 # the first augmentation
 $\tilde{x}_{2k-1} = t(x_k)$
 $h_{2k-1} = f(\tilde{x}_{2k-1})$ # representation
 $z_{2k-1} = g(h_{2k-1})$ # projection
 # the second augmentation
 $\tilde{x}_{2k} = t'(x_k)$
 $h_{2k} = f(\tilde{x}_{2k})$ # representation
 $z_{2k} = g(h_{2k})$ # projection
 end for
 for all $i \in \{1, \dots, 2N\}$ and $j \in \{1, \dots, 2N\}$ **do**
 $s_{i,j} = z_i^\top z_j / (\|z_i\| \|z_j\|)$ # pairwise similarity
 end for
 define $\ell(i, j)$ **as** $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$
 $\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$
 update networks f and g to minimize \mathcal{L}
end for
return encoder network $f(\cdot)$, and throw away $g(\cdot)$

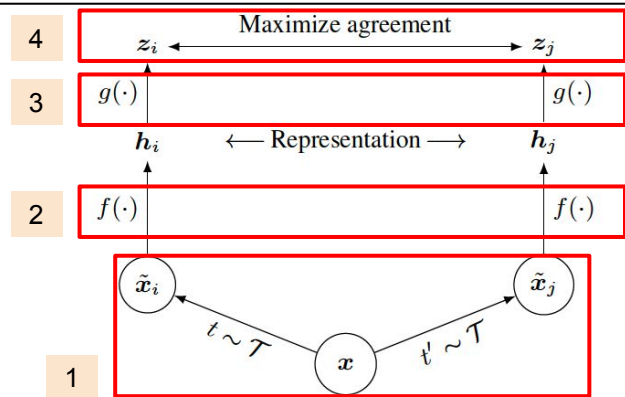
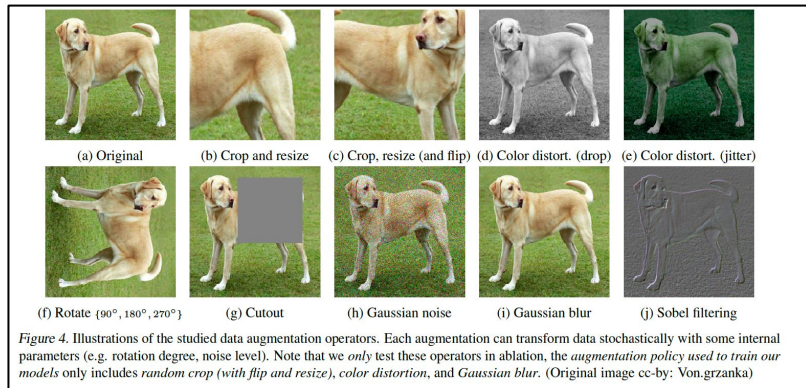


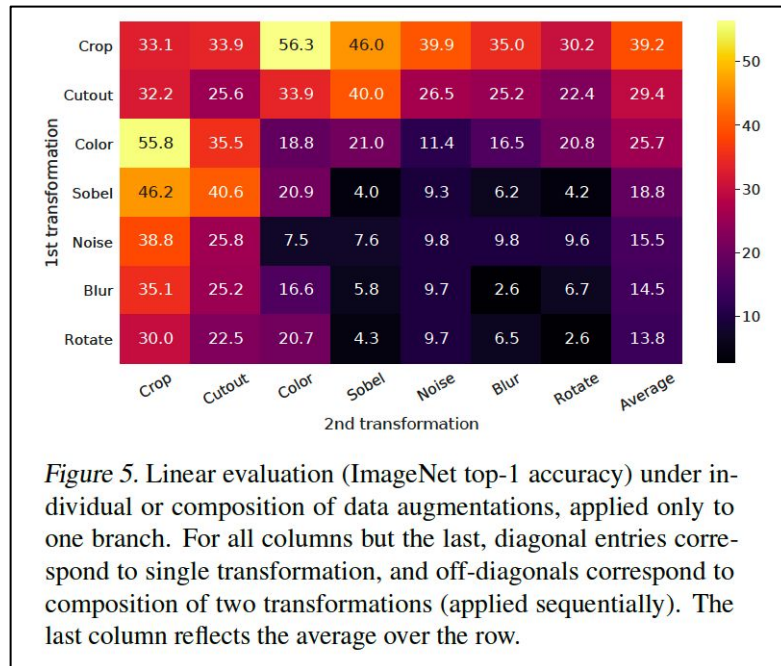
Figure 2. A simple framework for contrastive learning of visual representations. Two separate data augmentation operators are sampled from the same family of augmentations ($t \sim \mathcal{T}$ and $t' \sim \mathcal{T}$) and applied to each data example to obtain two correlated views. A base encoder network $f(\cdot)$ and a projection head $g(\cdot)$ are trained to maximize agreement using a contrastive loss. After training is completed, we throw away the projection head $g(\cdot)$ and use encoder $f(\cdot)$ and representation h for downstream tasks.

2. SimCLR

Experiments



Data Augmentations



Data Augmentations Pairs Results

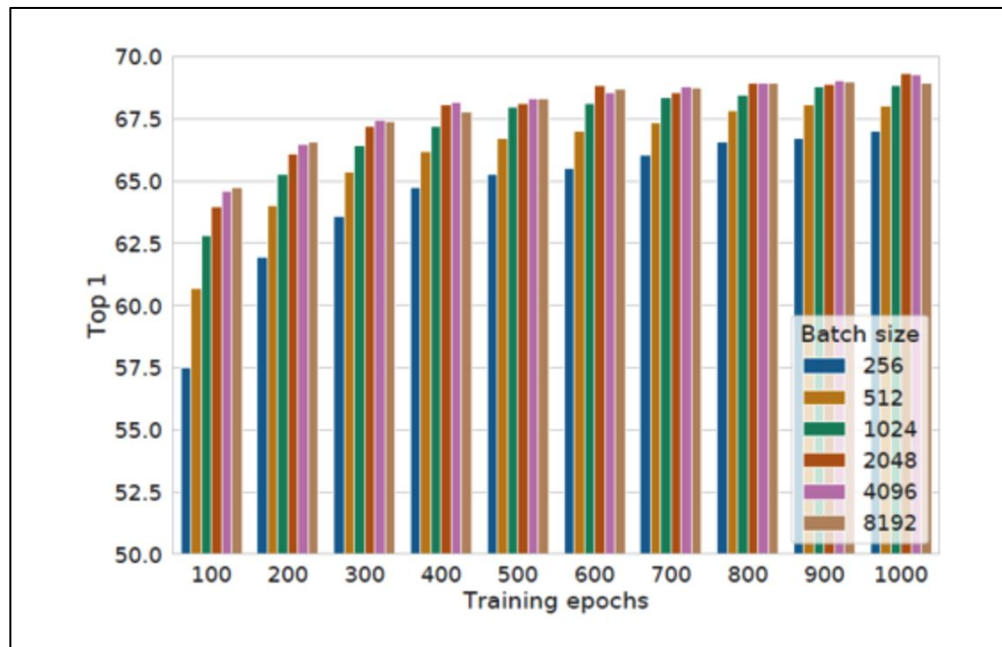
2. SimCLR

Experiments

contrastive learning benefits from ..

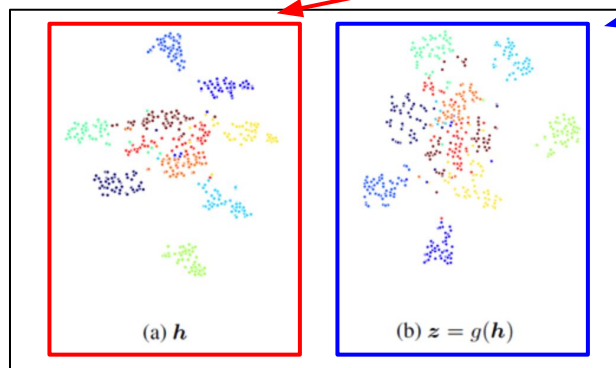
- **larger batch sizes**
- **more training steps**

than supervised learning

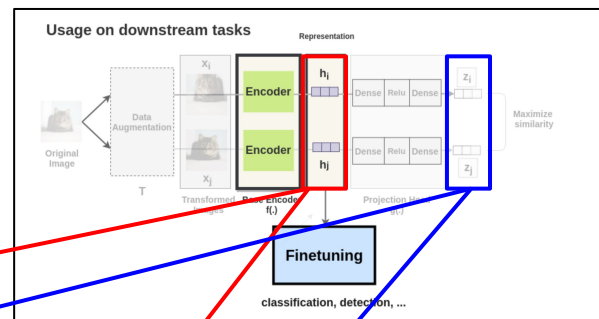


2. SimCLR

Experiments



What to predict?	Random guess	Representation	
		h	$g(h)$
Color vs grayscale	80	99.3	97.4
Rotation	25	67.6	25.6
Orig. vs corrupted	50	99.5	59.6
Orig. vs Sobel filtered	50	96.6	56.3



Representation obtained from **base encoder (not projection head)** can be used for other tasks!

SimCLR (2020)

A Simple Framework for Contrastive Learning of Visual Representations

Ting Chen¹ Simon Kornblith¹ Mohammad Norouzi¹ Geoffrey Hinton¹

Abstract

This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently proposed contrastive self-supervised learning algorithms without requiring specialized architectures or a memory bank. In order to understand what enables the contrastive prediction tasks to learn useful representations, we systematically study the major components of our framework. We show that (1) composition of data augmentations plays a critical role in defining effective predictive tasks, (2) introducing a learnable nonlinear transformation between the representation and the contrastive loss substantially improves the quality of the learned representations, and (3) contrastive learning benefits from larger batch sizes and more training steps compared to supervised learning. By combining these findings,

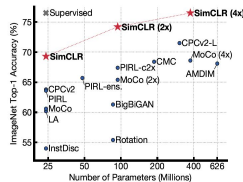


Figure 1. ImageNet Top-1 accuracy of linear classifiers trained on representations learned with different self-supervised methods (pretrained on ImageNet). Gray cross indicates supervised ResNet-50. Our method, SimCLR, is shown in bold.

<https://arxiv.org> > cs

A Simple Framework for Contrastive Learning of Visual ... - arXiv

by T Chen · 2020 · Cited by 5112 — Abstract: This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently proposed ...

Cite as: [arXiv:2002.05709](https://arxiv.org/abs/2002.05709)

MoCo v1 (2019)

Momentum Contrast for Unsupervised Visual Representation Learning

Kaiming He Haoqi Fan Yuxin Wu Saining Xie Ross Girshick

Facebook AI Research (FAIR)

Code: <https://github.com/facebookresearch/moco>

Abstract

We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning [29] as dictionary look-up, we build a dynamic dictionary with a queue and a moving-averaged encoder. This enables building a large and consistent dictionary on-the-fly that facilitates contrastive unsupervised learning. MoCo provides competitive results under the common linear protocol on ImageNet classification. More importantly, the representations learned by MoCo transfer well to downstream tasks. MoCo can outperform its supervised pre-training counterpart in 7 detection/segmentation tasks on PASCAL VOC, COCO, and other datasets, sometimes surpassing it by large margins. This suggests that the gap between unsupervised and supervised representation learning has been largely closed in many vision tasks.

1. Introduction

Unsupervised representation learning is highly successful in natural language processing, e.g., as shown by GPT [50, 51] and BERT [12]. But supervised pre-training is still

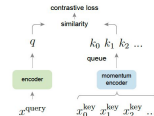


Figure 1. Momentum Contrast (MoCo) trains a visual representation of encoded keys by matching an encoded query q to a dictionary of encoded keys using a contrastive loss. The dictionary keys $\{k_0, k_1, k_2, \dots\}$ are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

<https://arxiv.org> > cs

Momentum Contrast for Unsupervised Visual Representation ...

by K He · 2019 · Cited by 3735 — Abstract: We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning as ...

Cite as: [arXiv:1911.05722](https://arxiv.org/abs/1911.05722)

MoCo v2 (2020)

Improved Baselines with Momentum Contrastive Learning

Xinlei Chen Haoqi Fan Ross Girshick Kaiming He

Facebook AI Research (FAIR)

Abstract

Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR. In this note, we verify the effectiveness of two of SimCLR's design improvements by implementing them in the MoCo framework. With simple modifications to MoCo—namely, using an MLP projection head and more data augmentation—we establish stronger baselines that outperform SimCLR and do not require large training batches. We hope this will make state-of-the-art unsupervised learning research more accessible. Code will be made public.

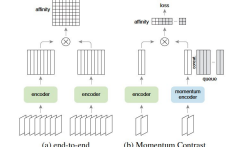


Figure 1. A batching perspective of two optimization mechanisms for contrastive learning. Images are encoded into a representation space, in which pairwise affinities are computed.

effective contrastive loss function, called InfoNCE [13], is:

$$\mathcal{L}_{\text{InfoNCE}}(x, y) = -\log \frac{\exp(q(x, y)/r)}{\exp(q(x, y)/r) + \sum_{k \neq y} \exp(q(x, k)/r)}. \quad (1)$$

<https://arxiv.org> > cs

Improved Baselines with Momentum Contrastive Learning

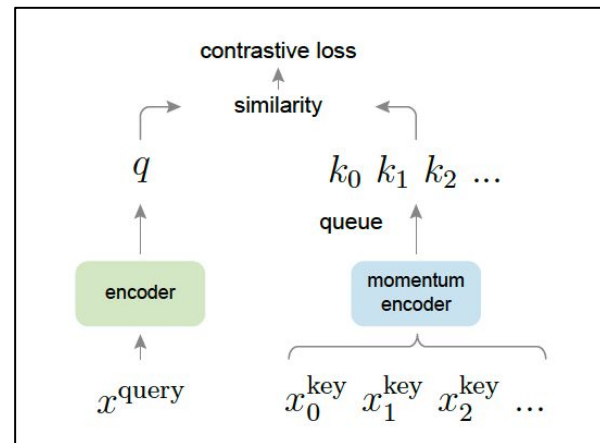
by X Chen · 2020 · Cited by 1088 — Abstract: Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR.

Cite as: [arXiv:2003.04297](https://arxiv.org/abs/2003.04297)

3. MoCo v1

MoCo = **M**omentum **C**ontrast

- **UN-SUPERVISED** visual representation learning algorithm
- **(1) Dictionary** look-up perspective
 - build a **Dynamic dictionary** (with a queue (FIFO))
- **(2) Moving Average Encoder**

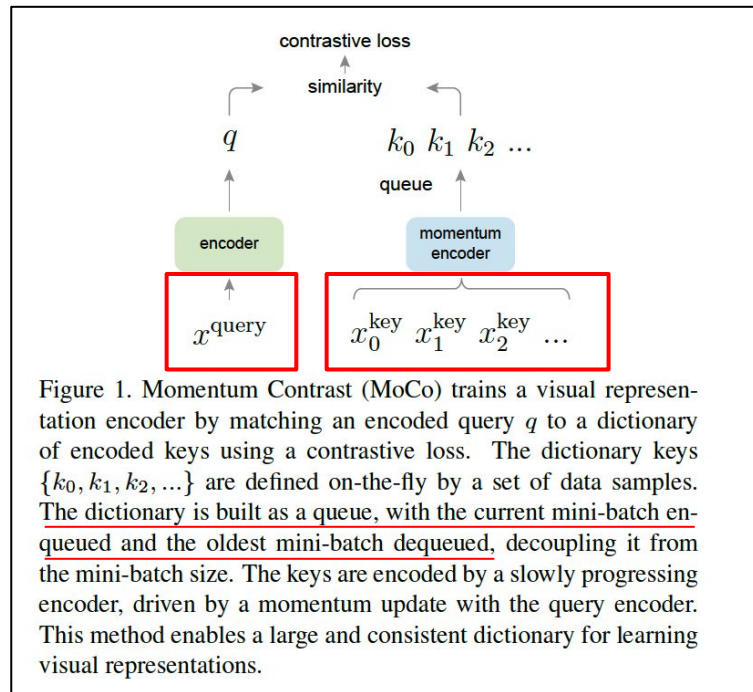
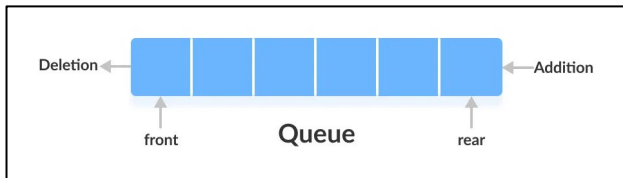


3. MoCo v1

(1) Dynamic Dictionary

- (SimCLR) positive pair & negative pair in ONE BATCH
- (MoCo) define a Dictionary
 - **match O** with query -> **positive** key
 - **match X** with query -> **negative** key

Use “**Dynamic**” Dictionary, by using **queue (FIFO)**



3. MoCo v1

(1) Dynamic Dictionary

Notation

- encoded query : q
- keys of a dictionary :
 - set of encoded samples : $\{k_0, k_1, k_2, \dots\}$
 - positive key : k_+
 - negative key : k_-

Contrastive Loss : low , when...

- q is similar to its positive key k_+
- dissimilar to other key (= negative keys)

(similarity : measured by **dot product**)

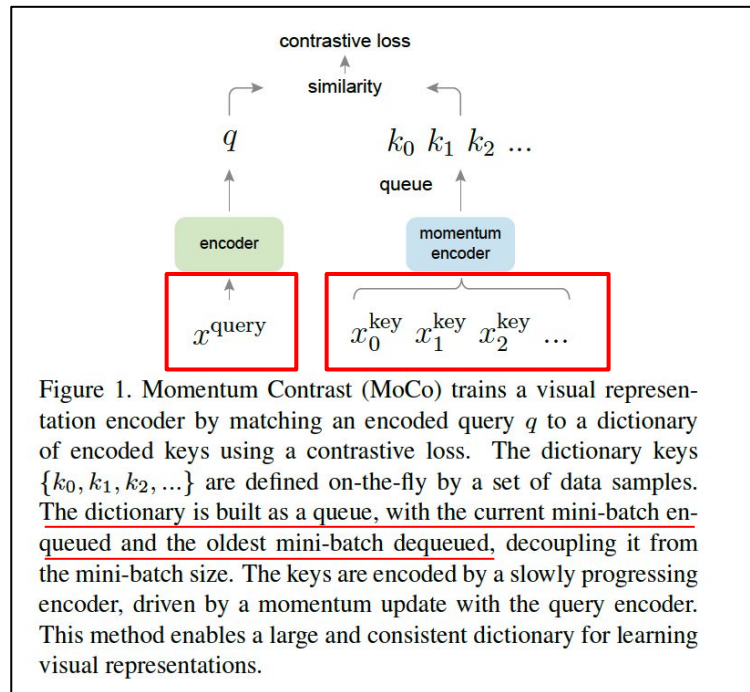


Figure 1. Momentum Contrast (MoCo) trains a visual representation encoder by matching an encoded query q to a dictionary of encoded keys using a contrastive loss. The dictionary keys $\{k_0, k_1, k_2, \dots\}$ are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

3. MoCo v1

(2) Loss Function : **InfoNCE**

$$\mathcal{L}_q = -\log \frac{\exp(q \cdot k_+ / \tau)}{\sum_{i=0}^K \exp(q \cdot k_i / \tau)}.$$

- τ : temperature

→ sum over **one** positive & **K** negatives

(= log loss of $(K + 1)$ way softmax classifier , that tries to **classify** q as k_+)

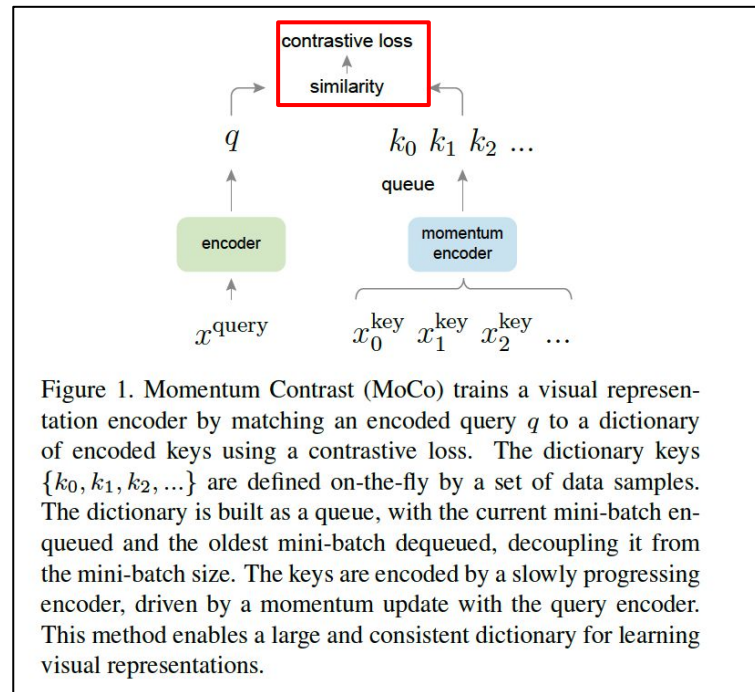
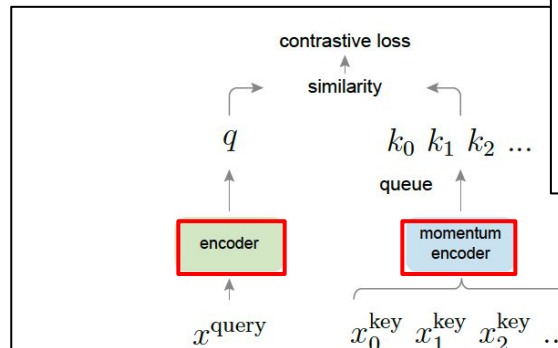
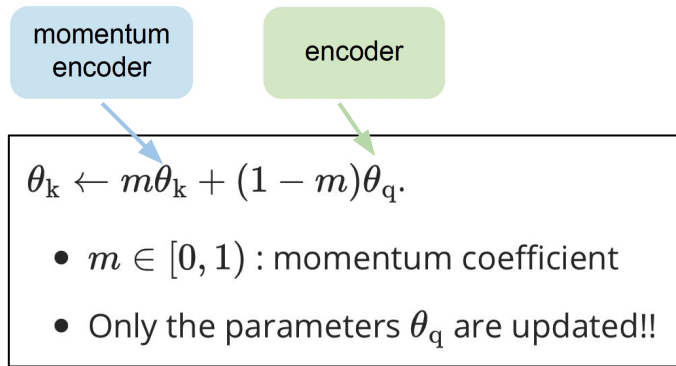


Figure 1. Momentum Contrast (MoCo) trains a visual representation encoder by matching an encoded query q to a dictionary of encoded keys using a contrastive loss. The dictionary keys $\{k_0, k_1, k_2, \dots\}$ are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

3. MoCo v1

(3) Moving Average Encoder

- **slowly progressing** key encoder
= **momentum-based MA of query encoder**
- to maintain “**consistency**”



Model Notation

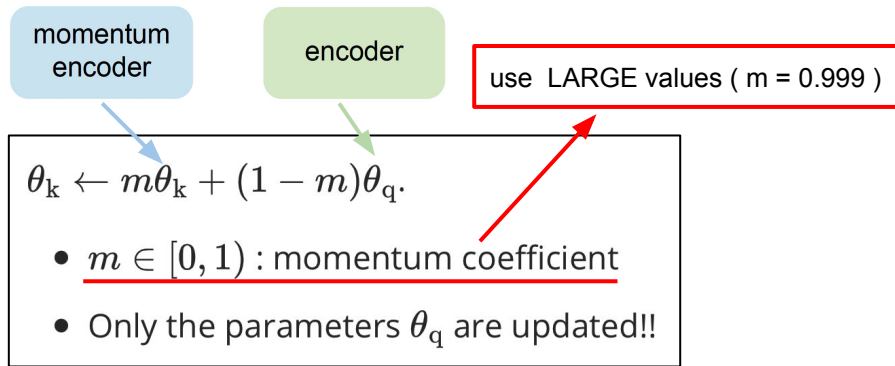
- query : $q = f_q(x^q)$
 - f_q : encoder network
 - x^q : query
- key : $k = f_k(x^k)$
 - f_k : encoder network
 - x^k : key

Figure 1. Momentum Contrast (MoCo) trains a visual representation encoder by matching an encoded query q to a dictionary of encoded keys using a contrastive loss. The dictionary keys $\{k_0, k_1, k_2, \dots\}$ are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

3. MoCo v1

(3) Moving Average Encoder

- **slowly progressing** key encoder
= **momentum-based MA of query encoder**
- to maintain “**consistency**”



Model Notation

- query : $q = f_q(x^q)$
 - f_q : encoder network
 - x^q : query
- key : $k = f_k(x^k)$
 - f_k : encoder network
 - x^k : key

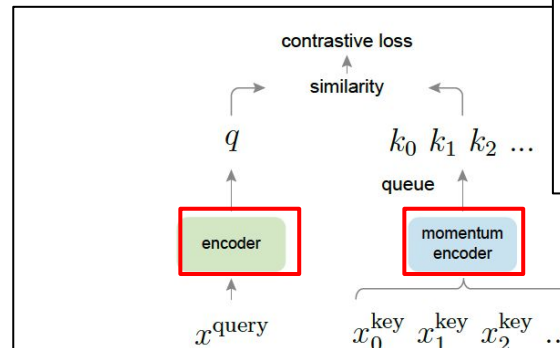


Figure 1. Momentum Contrast (MoCo) trains a visual representation encoder by matching an encoded query q to a dictionary of encoded keys using a contrastive loss. The dictionary keys $\{k_0, k_1, k_2, \dots\}$ are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

3. MoCo v1

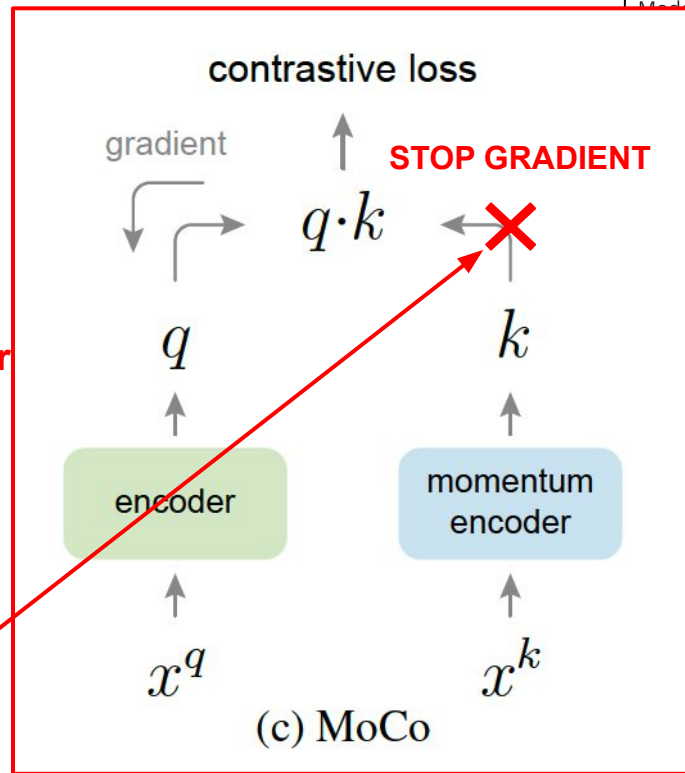
(3) Moving Average Encoder

- **slowly progressing** key encoder
= **momentum-based MA of query encoder**
- to maintain “**consistency**”



$$\theta_k \leftarrow m\theta_k + (1 - m)\theta_q.$$

- $m \in [0, 1)$: momentum coefficient
- Only the parameters θ_q are updated!!



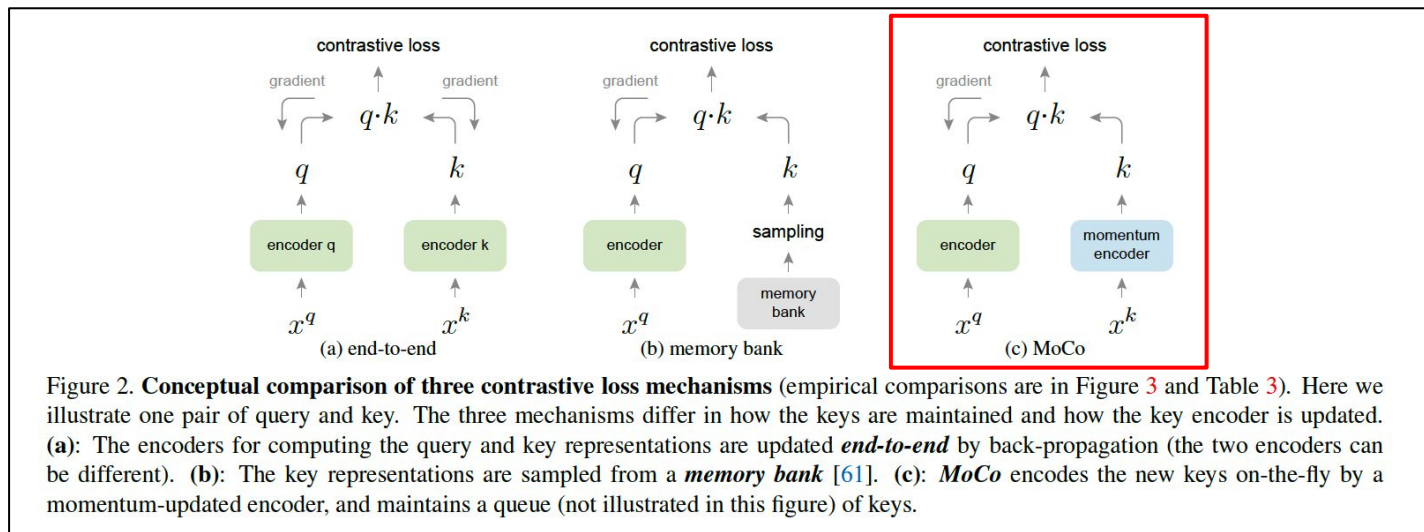
Model Notation

- query : $q = f_q(x^q)$
- f_q : encoder network
- x^q : query
- key : $k = f_k(x^k)$
- f_k : encoder network
- x^k : key

representational
query keys
samples.
match en-
it from
gressing
encoder.
learning

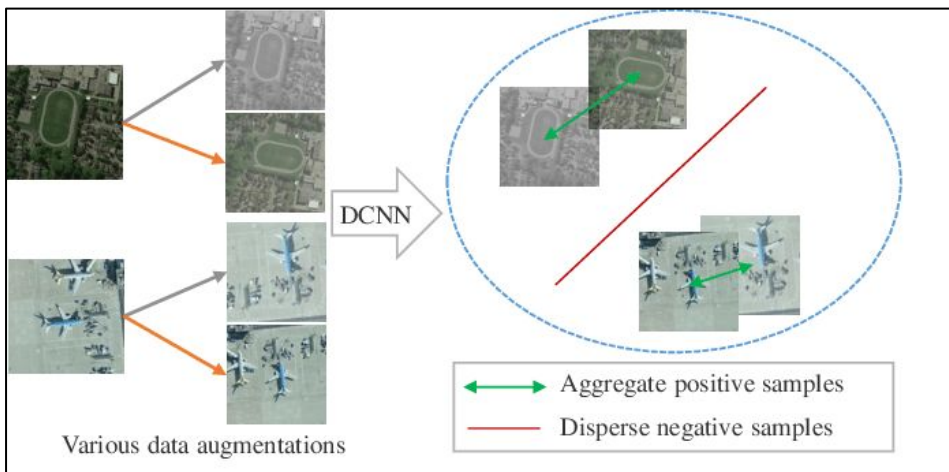
3. MoCo v1

Relations to Previous Mechanisms (end-to-end & memory bank)



3. MoCo v1

Pre-text Task : Instance Discrimination



3. MoCo v1

Pseudo-code

Algorithm 1 Pseudocode of MoCo in a PyTorch-like style.

```
# f_q, f_k: encoder networks for query and key
# queue: dictionary as a queue of K keys (CxK)
# m: momentum
# t: temperature

f_k.params = f_q.params # initialize
for x in loader: # load a minibatch x with N samples
    x_q = aug(x) # a randomly augmented version
    x_k = aug(x) # another randomly augmented version

    q = f_q.forward(x_q) # queries: NxK
    k = f_k.forward(x_k) # keys: NxK
    k = k.detach() # no gradient to keys

    # positive logits: Nx1
    l_pos = bmm(q.view(N,1,C), k.view(N,C,1))

    # negative logits: NxK
    l_neg = mm(q.view(N,C), queue.view(C,K))

    # logits: Nx(1+K)
    logits = cat([l_pos, l_neg], dim=1)

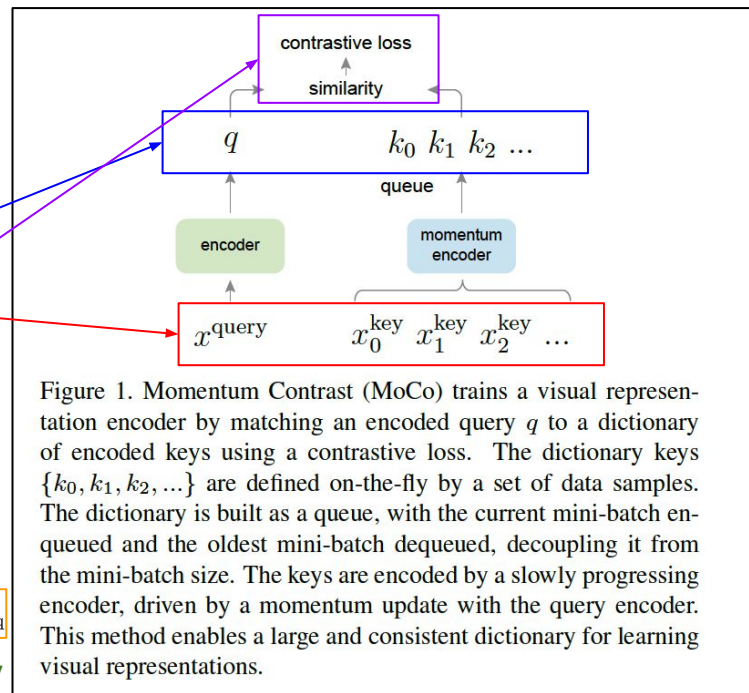
    # contrastive loss, Eqn. (1)
    labels = zeros(N) # positives are the 0-th
    loss = CrossEntropyLoss(logits/t, labels)

    # SGD update: query network
    loss.backward()
    update(f_q.params)

    # momentum update: key network
    f_k.params = m*f_k.params + (1-m)*f_q.params

    # update dictionary
    enqueue(queue, k) # enqueue the current minibatch
    dequeue(queue) # dequeue the earliest minibatch
```

bmm: batch matrix multiplication; mm: matrix multiplication; cat: concatenation.



SimCLR (2020)

A Simple Framework for Contrastive Learning of Visual Representations

Ting Chen¹ Simon Kornblith¹ Mohammad Norouzi¹ Geoffrey Hinton¹

Abstract

This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently proposed contrastive self-supervised learning algorithms without requiring specialized architectures or a memory bank. In order to understand what enables the contrastive prediction tasks to learn useful representations, we systematically study the major components of our framework. We show that (1) composition of data augmentations plays a critical role in defining effective predictive tasks, (2) introducing a learnable nonlinear transformation between the representation and the contrastive loss substantially improves the quality of the learned representations, and (3) contrastive learning benefits from larger batch sizes and more training steps compared to supervised learning. By combining these findings,

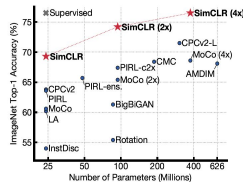


Figure 1. ImageNet Top-1 accuracy of linear classifiers trained on representations learned with different self-supervised methods (pretrained on ImageNet). Gray cross indicates supervised ResNet-50. Our method, SimCLR, is shown in bold.

<https://arxiv.org/abs/2002.05709>

A Simple Framework for Contrastive Learning of Visual ... - arXiv

by T Chen · 2020 · Cited by 5112 — Abstract: This paper presents SimCLR: a simple framework for contrastive learning of visual representations. We simplify recently proposed ...

Cite as: arXiv:2002.05709

MoCo v1 (2019)

Momentum Contrast for Unsupervised Visual Representation Learning

Kaiming He Haoqi Fan Yuxin Wu Saining Xie Ross Girshick

Facebook AI Research (FAIR)

Code: <https://github.com/facebookresearch/moco>

Abstract

We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning [29] as dictionary look-up, we build a dynamic dictionary with a queue and a moving-averaged encoder. This enables building a large and consistent dictionary on-the-fly that facilitates contrastive unsupervised learning. MoCo provides competitive results under the common linear protocol on ImageNet classification. More importantly, the representations learned by MoCo transfer well to downstream tasks. MoCo can outperform its supervised pre-training counterpart in 7 detection/segmentation tasks on PASCAL VOC, COCO, and other datasets, sometimes surpassing it by large margins. This suggests that the gap between unsupervised and supervised representation learning has been largely closed in many vision tasks.

1. Introduction

Unsupervised representation learning is highly successful in natural language processing, e.g., as shown by GPT [50, 51] and BERT [12]. But supervised pre-training is still

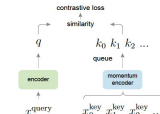


Figure 1. Momentum Contrast (MoCo) trains a visual representation of encoded keys using a contrastive loss. The dictionary keys $\{k_0, k_1, k_2, \dots\}$ are defined on-the-fly by a set of data samples. The dictionary is built as a queue, with the current mini-batch enqueued and the oldest mini-batch dequeued, decoupling it from the mini-batch size. The keys are encoded by a slowly progressing encoder, driven by a momentum update with the query encoder. This method enables a large and consistent dictionary for learning visual representations.

<https://arxiv.org/abs/1911.05722>

Momentum Contrast for Unsupervised Visual Representation ...

by K He · 2019 · Cited by 3735 — Abstract: We present Momentum Contrast (MoCo) for unsupervised visual representation learning. From a perspective on contrastive learning as ...

Cite as: arXiv:1911.05722

MoCo v2 (2020)

Improved Baselines with Momentum Contrastive Learning

Xinlei Chen Haoqi Fan Ross Girshick Kaiming He

Facebook AI Research (FAIR)

Abstract

Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR. In this note, we verify the effectiveness of two of SimCLR's design improvements by implementing them in the MoCo framework. With simple modifications to MoCo—namely, using an MLP projection head and more data augmentation—we establish stronger baselines that outperform SimCLR and do not require large training batches. We hope this will make state-of-the-art unsupervised learning research more accessible. Code will be made public.

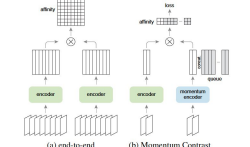


Figure 1. A batching perspective of two optimization mechanisms for contrastive learning. Images are encoded into a representation space, in which pairwise affinities are computed.

effective contrastive loss function, called InfoNCE [13], is:

$$\mathcal{L}_{\text{InfoNCE}}(x, y) = -\log \frac{\exp(q(x, y)/r)}{\exp(q(x, y)/r) + \sum_{k \neq y} \exp(q(x, k)/r)}. \quad (1)$$

<https://arxiv.org/abs/2003.04297>

Improved Baselines with Momentum Contrastive Learning

by X Chen · 2020 · Cited by 1088 — Abstract: Contrastive unsupervised learning has recently shown encouraging progress, e.g., in Momentum Contrast (MoCo) and SimCLR.

Cite as: arXiv:2003.04297

4. MoCo v2

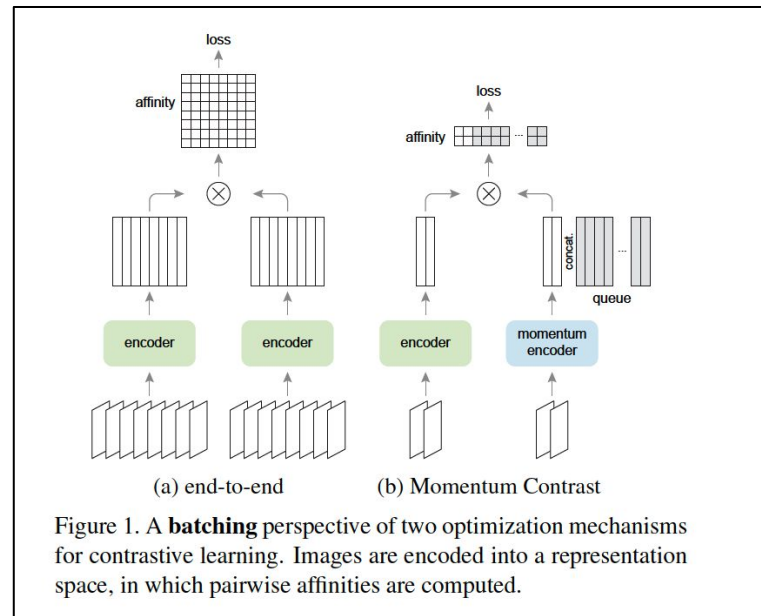
(1) **MoCo v1** :

- (a) Dynamic Dictionary
- (b) Moving-Averaged Encoder

(2) **SimCLR**

- (a) larger batch size for lots of negative samples
- (b) **stronger augmentation**
- (c) **MLP Projection head**

MoCo v2 = **MoCo v1** + **SimCLR** ((b) + (c))



5. Reference

- Discriminative Unsupervised Feature Learning with Exemplar Convolutional Neural Networks (Dosovitskiy, NIPS 2014)
- Unsupervised Visual Representation Learning by Context Prediction (Doersch et al., ICCV 2015)
- Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles (Noroozi and Favaro, ECCV 2016)
- Unsupervised Representation Learning by Predicting Image Rotations (Gidaris et al., ICLR 2018)
- A Simple Framework for Contrastive Learning of Visual Representations (Chen et al., ICML 2020)
- Momentum Contrast for Unsupervised Visual Representation Learning (He et al., CVPR 2020)
- Improved Baselines with Momentum Contrastive Learning (Chen et al., arxiv 2020)
- <https://lilianweng.github.io/posts/2021-05-31-contrastive/>
- <https://seunghan96.github.io/categories/cl/>

Thank You