

# Self-Supervised Learning

## Tretext Tasks

Klaudia Balcer

Advanced Data Mining



Uniwersytet  
Wrocławski

# Contents

- 1 A **Simple Contrastive Learning** of Visual **Representations**
- 2 **Momentum Contrast**
- 3 **Bootstrap your own Latent**
- 4 Self-**distillation** with **no** Labels
- 5 Graph Contrastive Learning with Augmentations
- 6 **Context-aware Diffusion-based Contrastive Learning** for Sequential **Recommendation**

# A Simple Contrastive Learning of Visual Representations

# Image Augmentations



(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

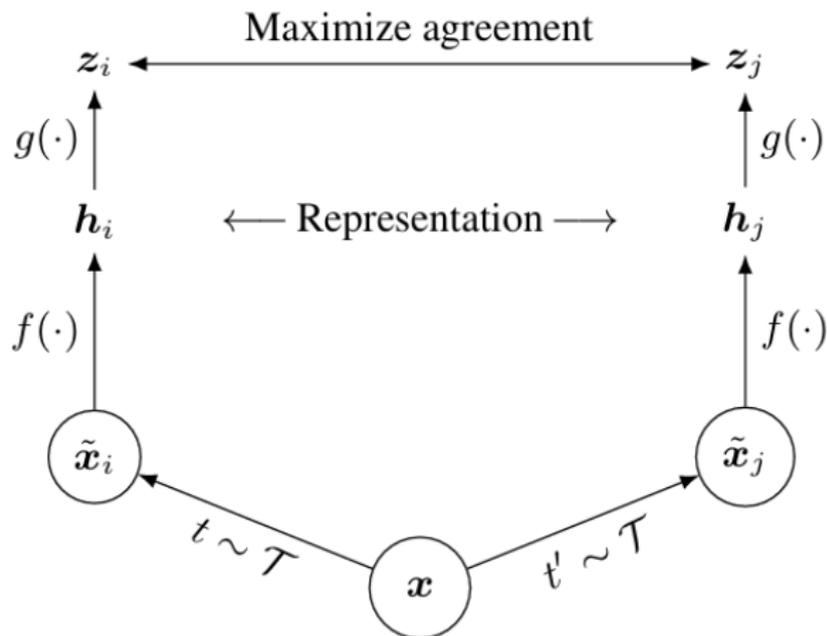


(i) Gaussian blur



(j) Sobel filtering

## SimCLR



<https://dl.acm.org/doi/10.5555/3524938.3525087>

# SimCLR

- stochastic data augmentation  $t$ :
  - random cropping
  - random color distortions
  - Gaussian blur
- base encoder  $f(\cdot)$  – ResNet (<https://arxiv.org/abs/1512.03385>)
- projection head  $g(\cdot)$  – MLP
- contrastive loss:

$$\mathcal{L} = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(z_i, z_k)/\tau)},$$

where  $\text{sim}$  is cosine similarity and as negatives we treat both views of all other images;

– NT-Xent (the normalized temperature-scaled cross entropy loss).

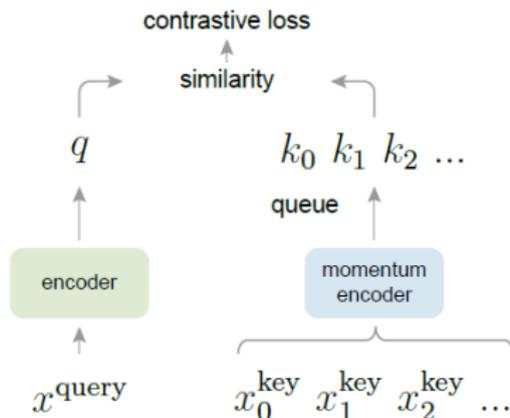
- huge batch size – from 256 to 8192, LARS optimizer, global batch normalization

# Momentum Contrast

## MoCo

- Contrastive learning as dynamic dictionary:
  - keys – instances from the data
  - values / queries – encoded images
- InfoNCE loss:

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



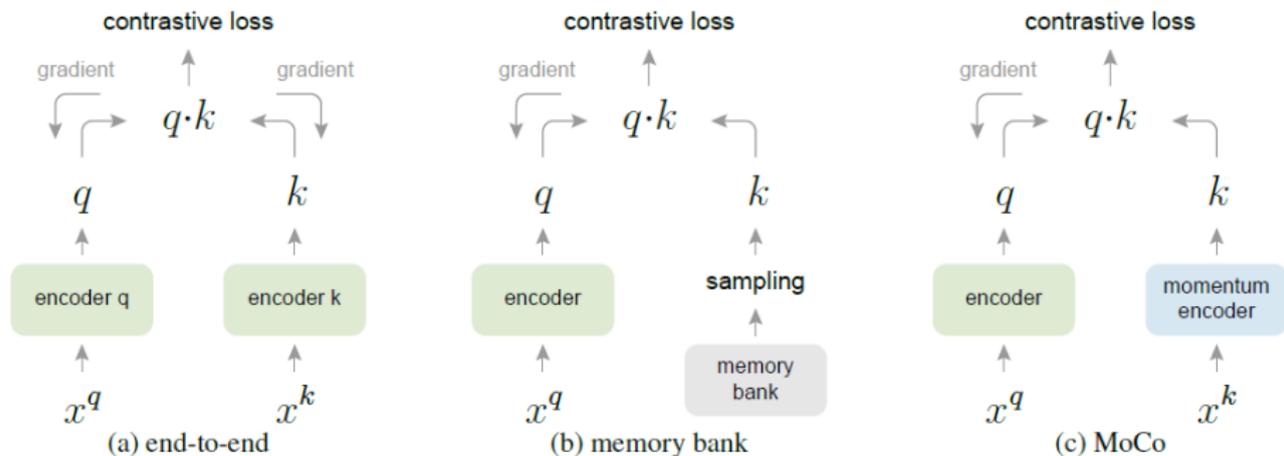
# MoCo

- Queue:
  - consist of embeddings of images from **multiple mini-batches**
  - progressive replacement: when a new mini-batch is added, the oldest is removed
- Momentum encoder:
  - when we do have *keys* from outside the current mini-batch, we cannot do back-propagation on the key encoder
  - using the query encoder to the key encoder leads to rapid changes and inconsistency in the queued embeddings
  - we use a momentum update rule:

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

- a slowly evolving key encoder is a core to making use of a queue
- Look at the pseudo-code from the paper (page 4)

## MoCo



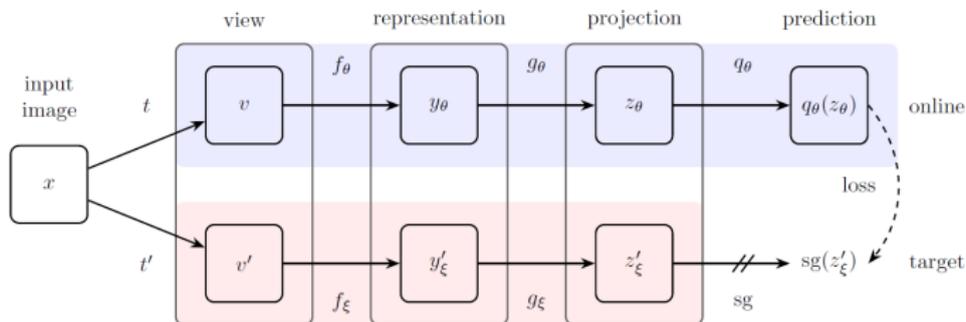
Technical details: ResNet encoder, temperature  $\tau = 0.07$ , augmentations: a  $224 \times 224$ -pixel crop is taken from a randomly resized image, and then undergoes random color jittering, random horizontal flip, and random grayscale conversion, all available in PyTorch's torchvision package, per device batch normalization (details in the paper), SGD optimizer.

[https://openaccess.thecvf.com/content\\_CVPR\\_2020/papers/He\\_Momentum\\_Contrast\\_for\\_Unsupervised\\_Visual\\_Representation\\_Learning\\_CVPR\\_2020\\_pa](https://openaccess.thecvf.com/content_CVPR_2020/papers/He_Momentum_Contrast_for_Unsupervised_Visual_Representation_Learning_CVPR_2020_paper.pdf)

# Bootstrap your own Latent

# BYOL

- experimental fining: when learning to mimic a random network, we can obtain much better representations than a random initialization (18.8% vs 1.4% on classification downstream task)
- two neural networks (online and target networks) that interact and learn from each other
- BYOL trains its online network to predict the target network's representation of another augmented view of the same image
- the effectiveness of BYOL depends less on the augmentations used than on previous approaches



## BYOL

- online network is defined by a set of weights  $\theta$  and is comprised of three stages: an encoder  $f_\theta$ , a projector  $g_\theta$  and a predictor  $q_\theta$ ,
- The target network has the same architecture as the online network, but uses a different set of weights  $\xi$
- we update the target network with moving average (momentum encoder):

$$\xi \leftarrow \tau \xi + (1 - \tau) \theta$$

- the prediction is made to map one representation to the other
- loss function (one summand, second summand with augmentations swap):

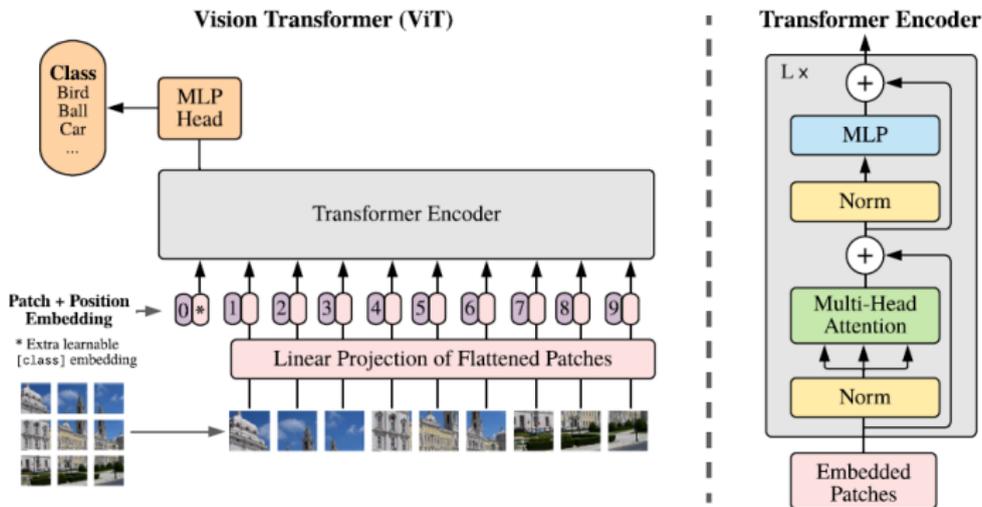
$$\mathcal{L}_{\theta, \xi} \triangleq \left\| \overline{q_\theta(z_\theta)} - \overline{z'_\xi} \right\|_2^2 = 2 - 2 \cdot \frac{\langle q_\theta(z_\theta), z'_\xi \rangle}{\|q_\theta(z_\theta)\|_2 \cdot \|z'_\xi\|_2}$$

- online network is optimized using a gradient based optimizer (LARS)
- architecture: ResNet50 (convnet)

# Self-distillation with **no** Labels

# ViT – recap

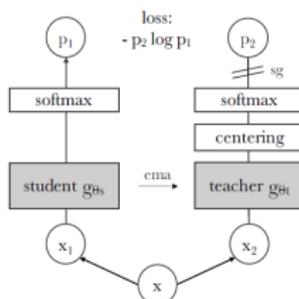
- show their power on large datasets (14M-300M images)
- applying a standard Transformer directly to images, with the fewest possible modifications



<https://openreview.net/pdf?id=YicbFdNTTy>

## DINO

- One of the components of the success of Transformers in NLP was self-supervised pretraining; ViT is trained in a supervised manner



- Both networks have the same architecture but different parameters.
- We apply a stop-gradient (sg) operator on the teacher to propagate gradients only through the student.
- The teacher parameters are updated with an exponential moving average (ema) of the student parameters.

## DINO

- Knowledge distillation without any true labels (with soft pseudo-labels), and the teacher is built during the training, somewhat related to co-distillation (similar to mutual training)
- *Global* augmentations (containing over 50% of the image), *Local* augmentations (containing less than 50% of the image),  $V$  – set of all views

$$\min_{\theta_s} \sum_{x \in \{x_1^g, x_2^g\}} \sum_{\substack{x' \in V \\ x' \neq x}} H(P_t(x), P_s(x'))$$

- teacher is updated with students weights in momentum encoder:

$$\theta_t \leftarrow \lambda \theta_t + (1 - \lambda) \theta_s$$

with  $\lambda$  from cosine schedule from 0.996 to 1, no batch-normalization

- avoiding collapsing: sharpening (temperature in softmax) and centering:

$$g_t(x) \leftarrow g_t(x) + c$$

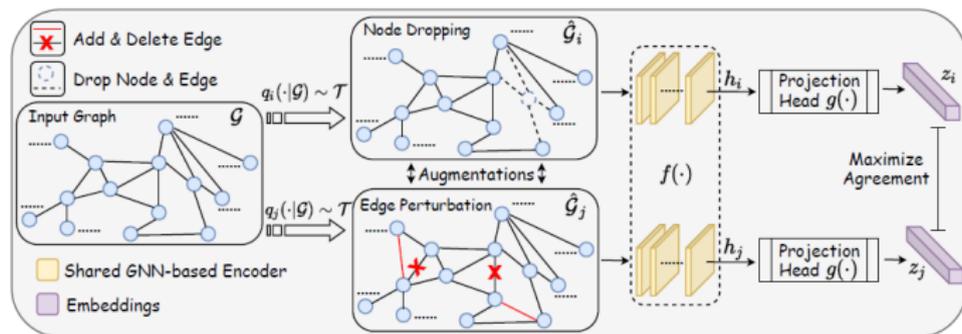
$$c \leftarrow mc = (1 - m) \frac{1}{B} \sum_{i=1}^B g_{\theta_t}(x_i)$$

# Graph Contrastive Learning with Augmentations

# GraphCL

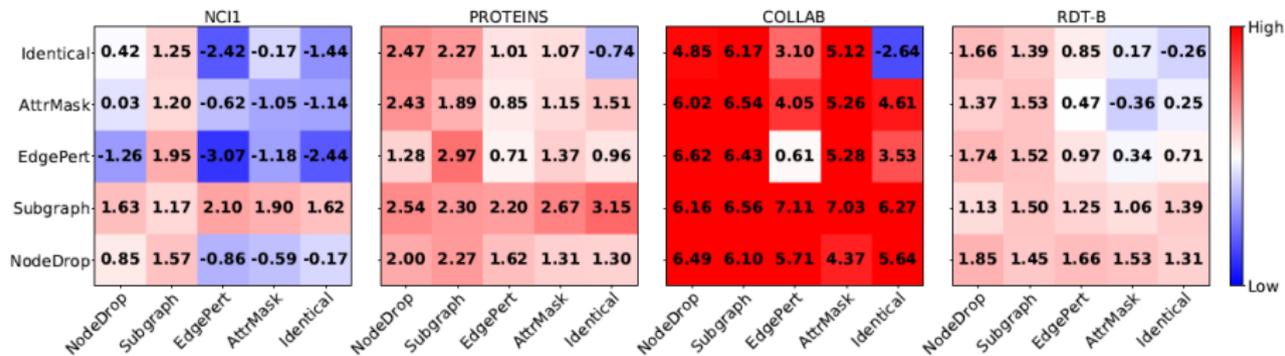
- Augmentations: node dropout, edge perturbation, attribute masking, subgraph selection (via random walk). The selection depends on the domain.
- NT-Xent loss:

$$\ell_n = -\log \frac{\exp(\text{sim}(z_{n,i}, z_{n,j}) / \tau)}{\sum_{n'=1, n' \neq n}^N \exp(\text{sim}(z_{n,i}, z_{n',j}) / \tau)}$$



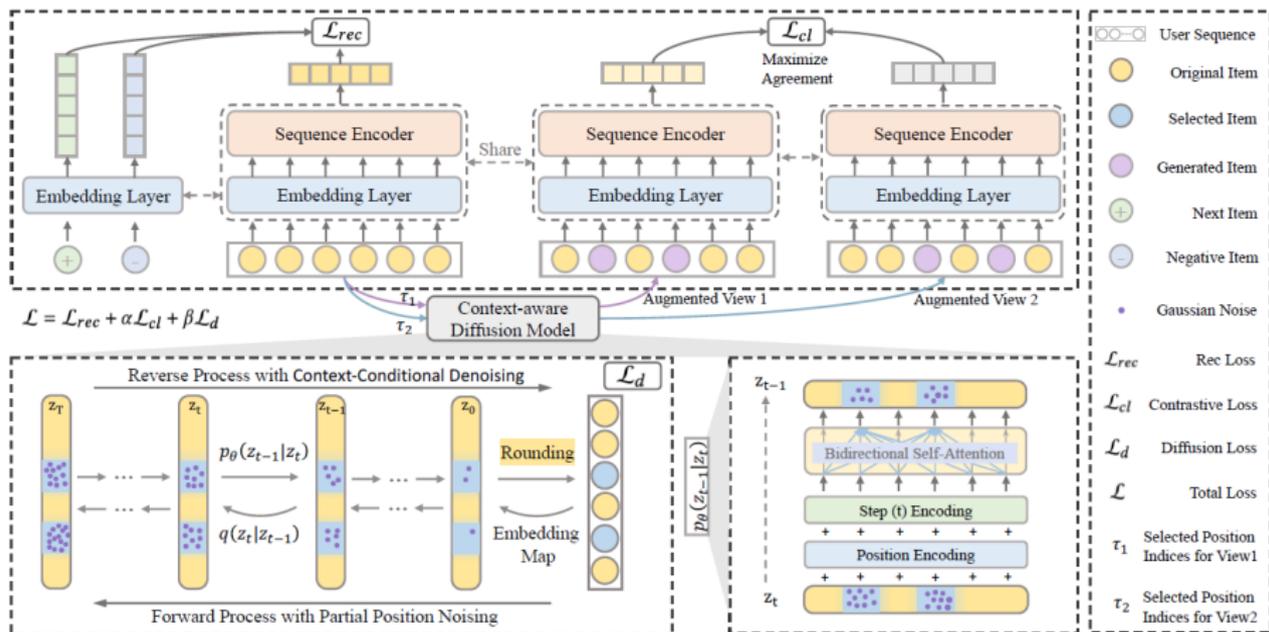
<https://proceedings.nips.cc/paper/2020/file/3fe230348e9a12c13120749e3f9fa4cd-Paper.pdf>

## GraphCL



# Context-aware **D**iffusion-based Contrastive Learning for Sequential **R**ecommendation

## CaDiRec



<https://arxiv.org/html/2405.09369v5>