

# DEEP LEARNING

## Lecture 12: Self-Supervised Learning

Dr. Yang Lu

Department of Computer Science and Technology

[luyang@xmu.edu.cn](mailto:luyang@xmu.edu.cn)



# Self-Supervised Learning



Lecun

Yann LeCun and Yoshua Bengio say  
at ICLR 2020:

*Self-supervised learning could lead  
to the creation of artificial  
intelligence (AI) programs that are  
more humanlike in their reasoning.*



Bengio



# Supervised and Unsupervised Learning

- Given a task and enough labels, supervised learning can solve it really well.
- However, good performance usually requires **a decent amount of labels**, but collecting manual labels is expensive (i.e. ImageNet) and hard to be scaled up.

# Supervised and Unsupervised Learning

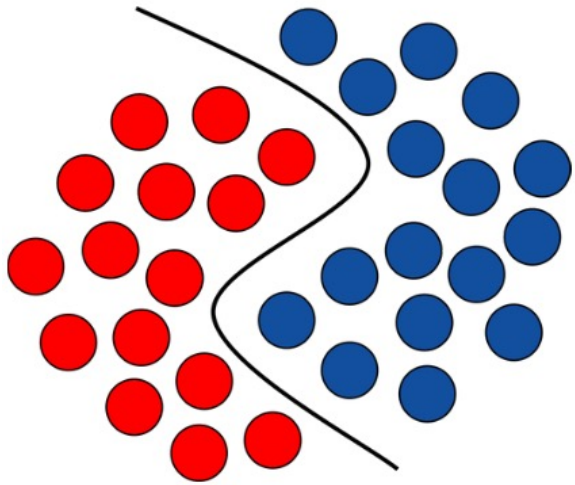
- Unlabeled data (e.g. free text, all the images on the Internet) is substantially more than a limited number of human curated labelled datasets,
  - It is kind of wasteful not to use them.
- However, unsupervised learning is not easy and usually works much less efficiently than supervised learning.



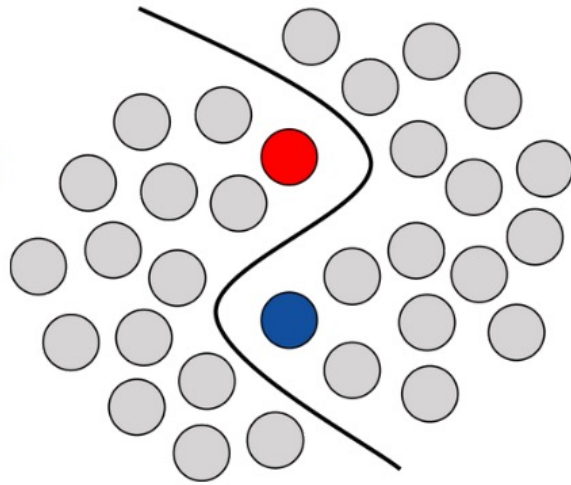
# Self-Supervised Learning

- What if we can **automatically generate labels** by some rules for unlabeled data and train unsupervised dataset in a supervised manner?
  - E.g. use a part of the data to predict the rest. The partition can be generated by rules, rather than human annotation.
- In this way, all the information needed, both inputs and labels, has been provided. This is known as **self-supervised learning**.
- The main purpose of self-supervised learning is to pre-train representations that can be transferred to downstream tasks by fine-tuning.

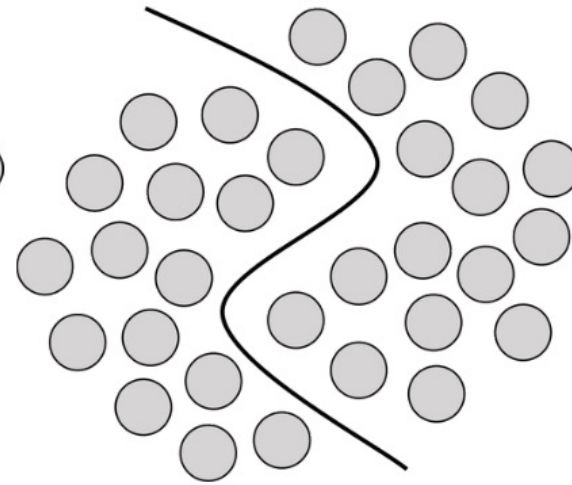
# Self-Supervised Learning



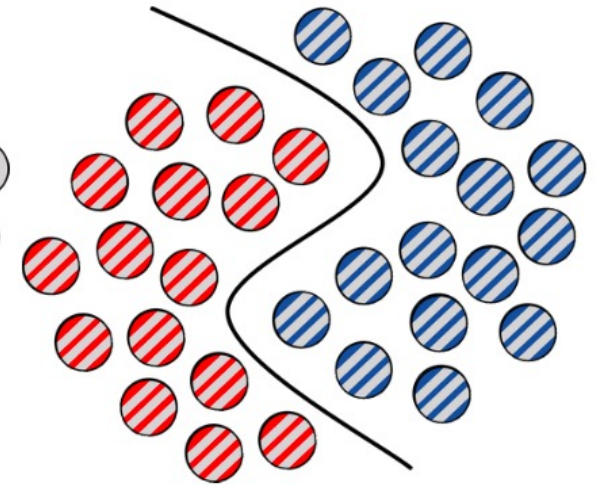
(a) Supervised



(b) Semi-supervised



(c) Unsupervised



(d) Self-supervised



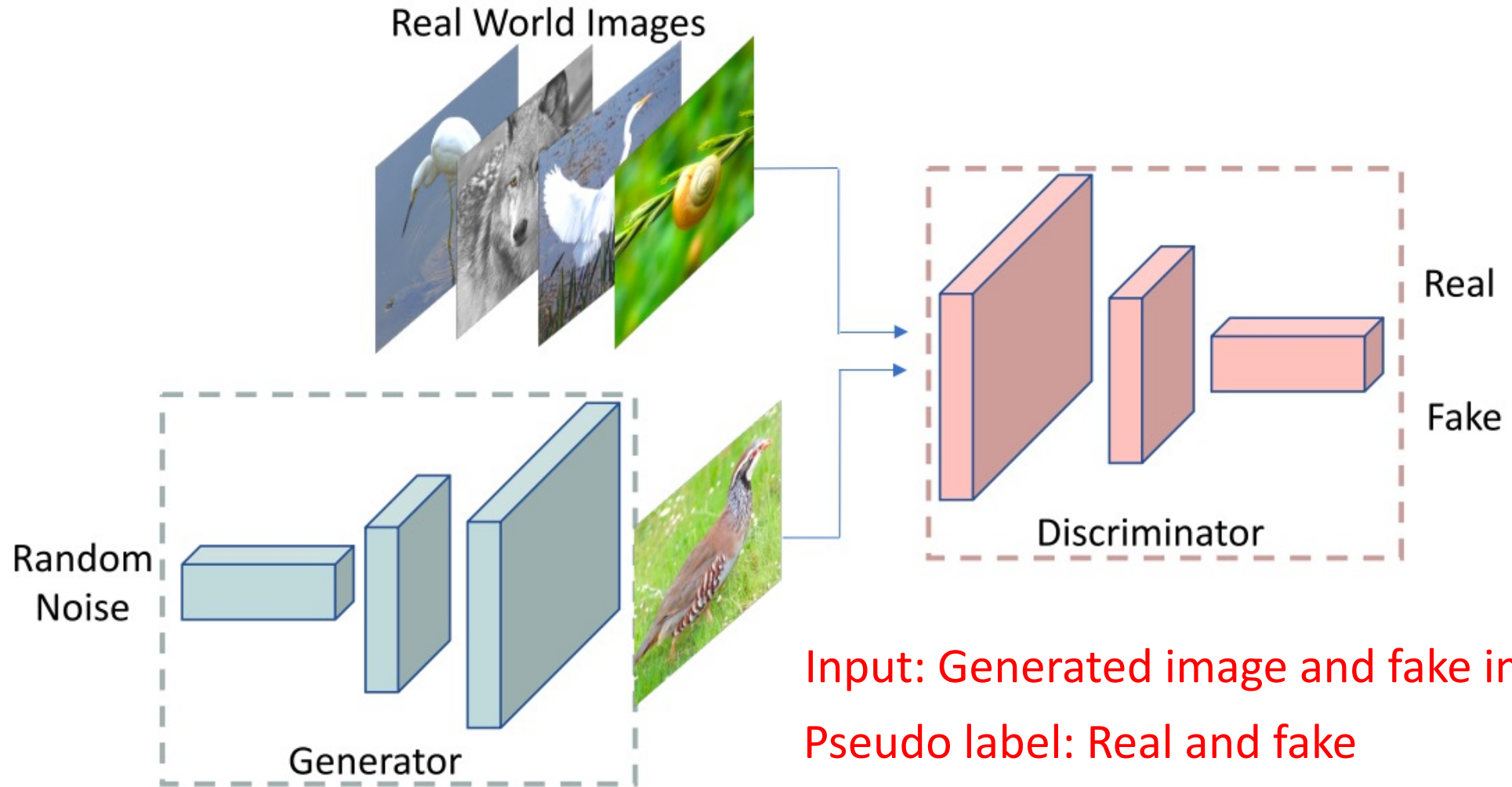
# Self-Supervised Learning

We have seen examples of self-supervised learning.

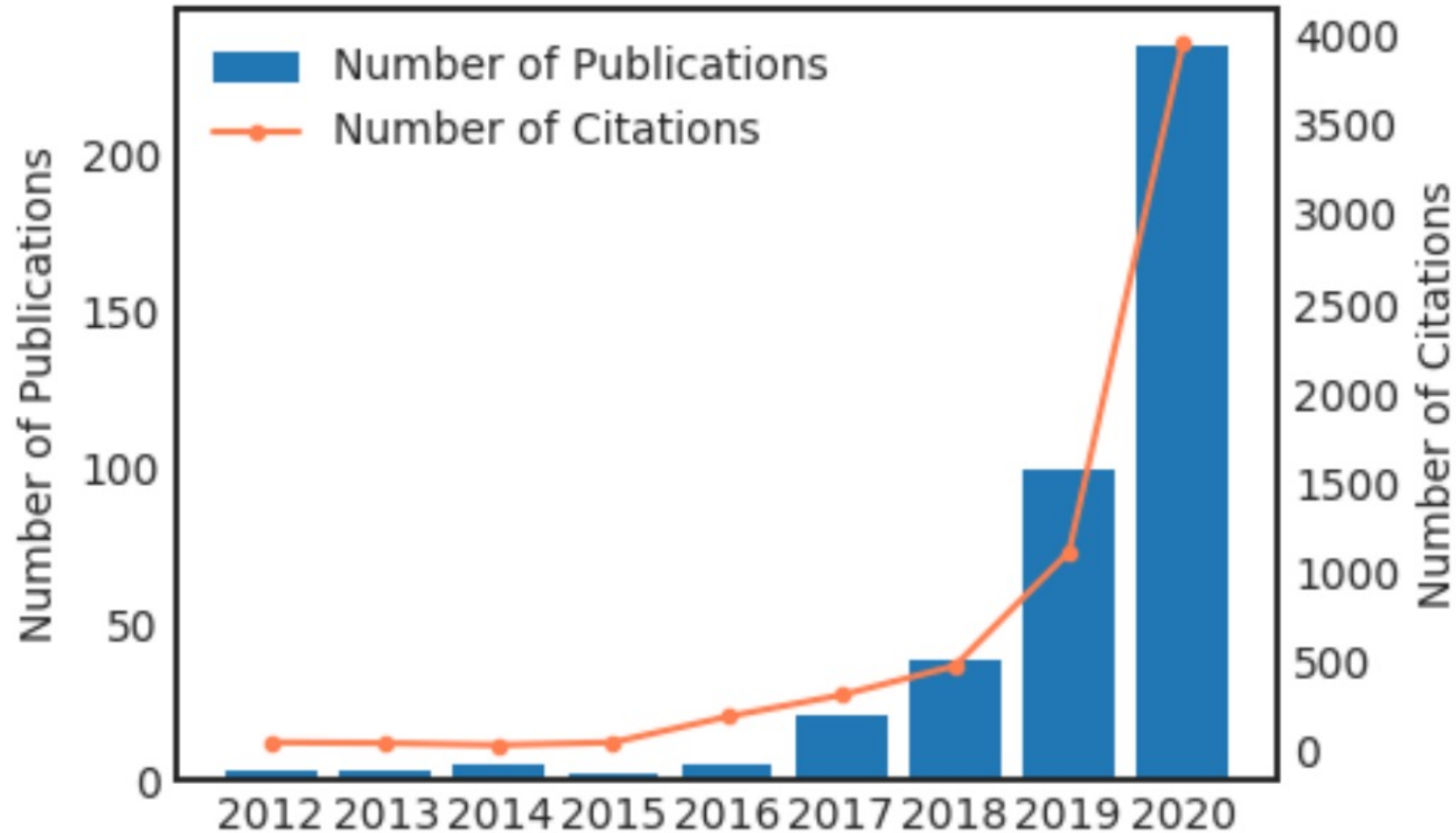
- Word2vec uses center word to predict context words.
  - The label (context words) is generated by sliding window.
- BERT has two tasks:
  - Use mask token to predict the missing word.
  - Concat two sentences to predict their order.
- GAN uses real images and fake images as labels.
- Graph embedding uses neighbors as labels.

**All the labels are automatically generated without human annotation for supervised learning task.**

# GAN



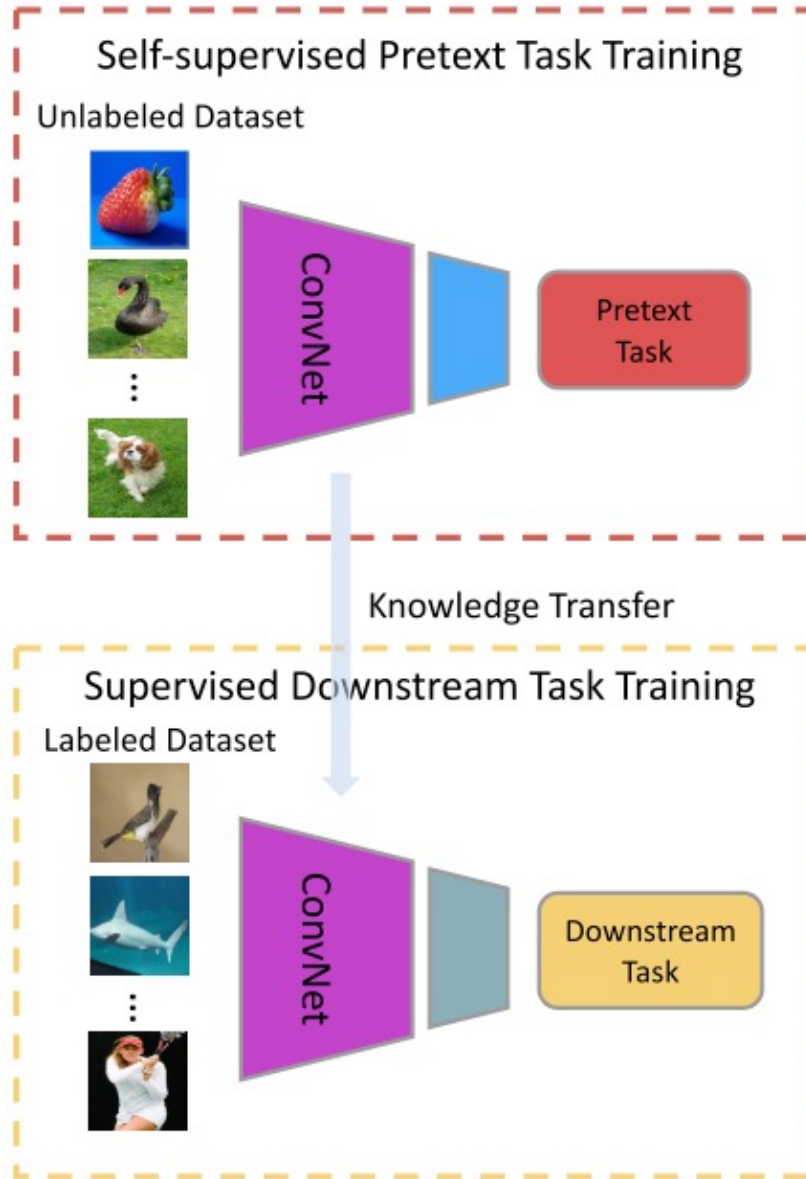
# Self-Supervised Learning



## Term Definition

- **Pretext Task**: Pre-designed tasks for networks to solve, in order to learn features as a pre-trained model.
- **Downstream Task**: Applications that are used to evaluate the quality of features learned by self-supervised learning.
- **Human-annotated label**: Labels of data that are manually annotated by human workers.
- **Pseudo label**: **Automatically generated labels** based on data attributes for pretext tasks.





# Outlines

- Generation-Based Methods
- Context-Based Methods
- Free Semantic Label-Based Methods
- Cross Modal-Based Methods
- Contrastive Learning



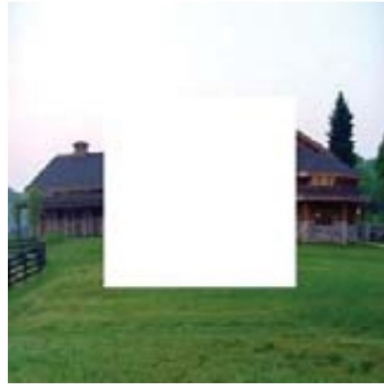
# GENERATION-BASED METHODS



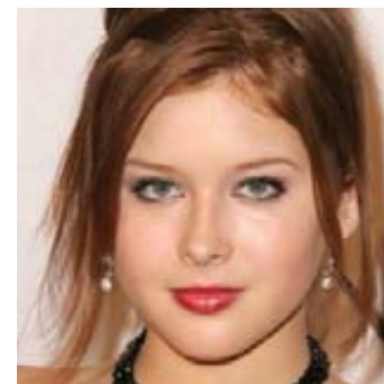
# Generation-Based Methods

- Idea: Use modified image to generate original image.
- The generator is able to learn image features by the loss between generated image and original image.
- The pseudo label is usually the original image.

# Image Generation with Inpainting



Input:  
Image with missing region



Pseudo label:  
Original Image



# Image Generation with Super Resolution



Input:  
Low resolution image

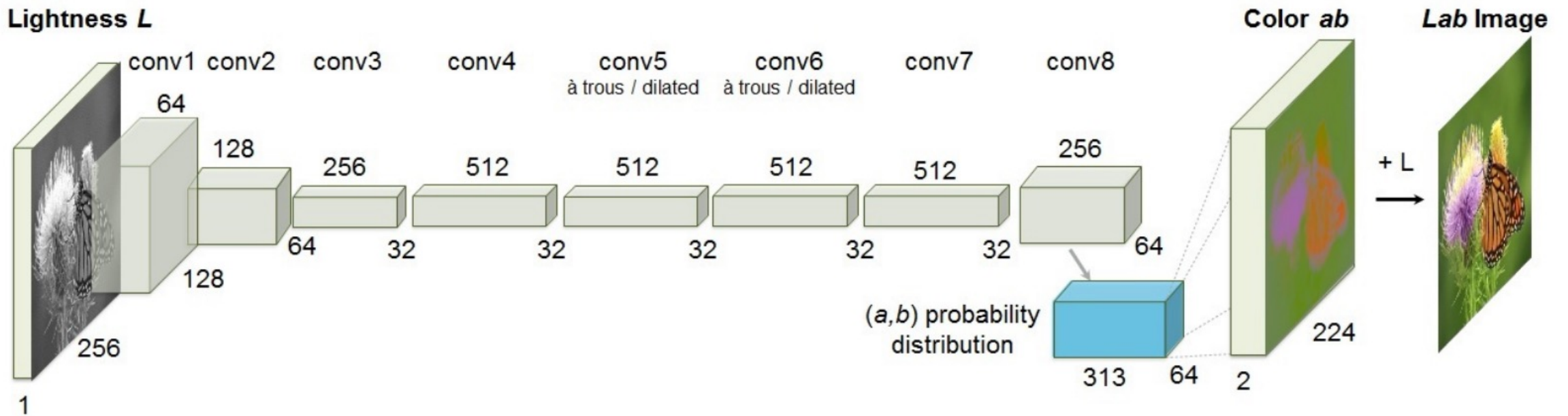


Pseudo label:  
High resolution image





# Image Generation with Colorization

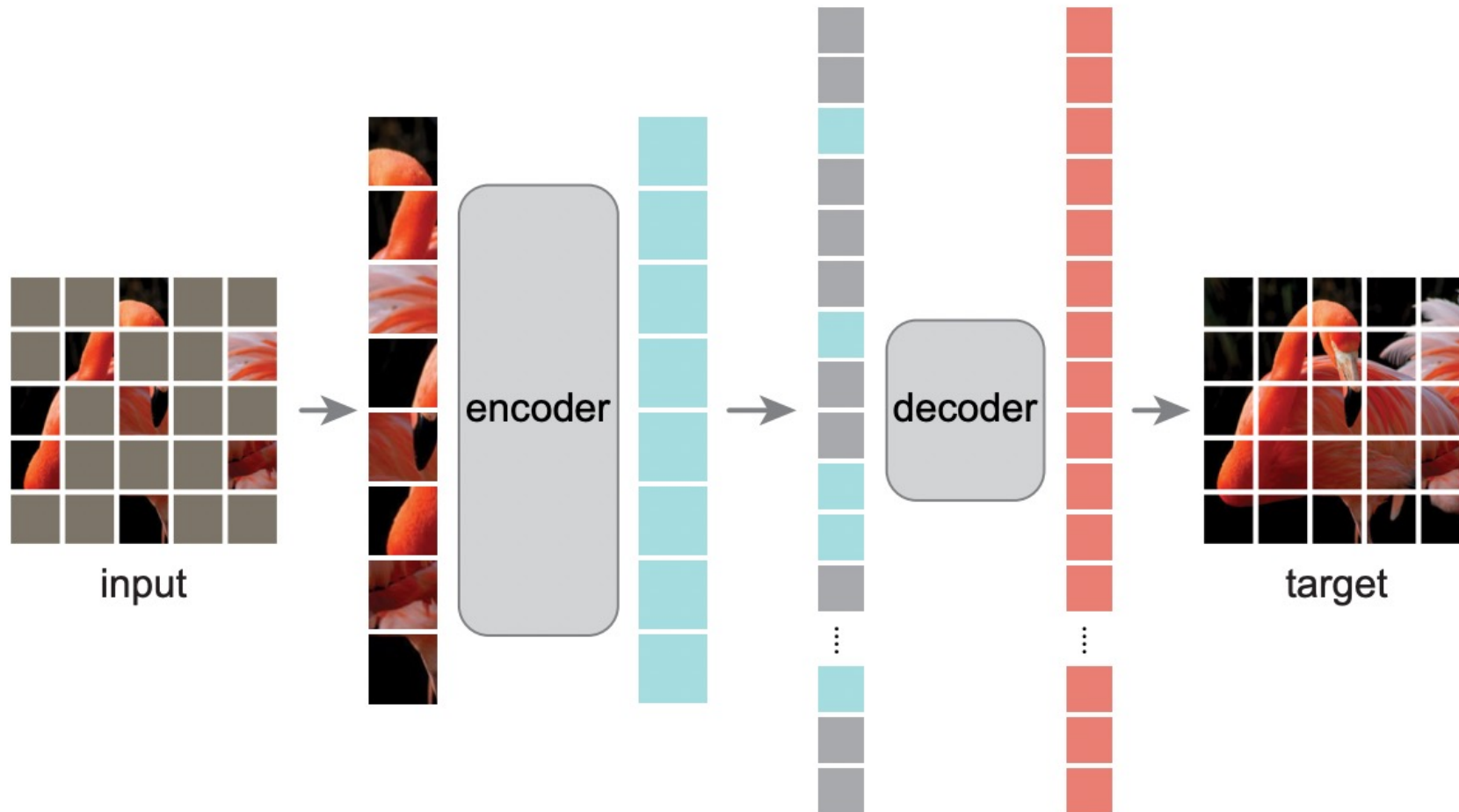


Input:  
Transformed grey level image

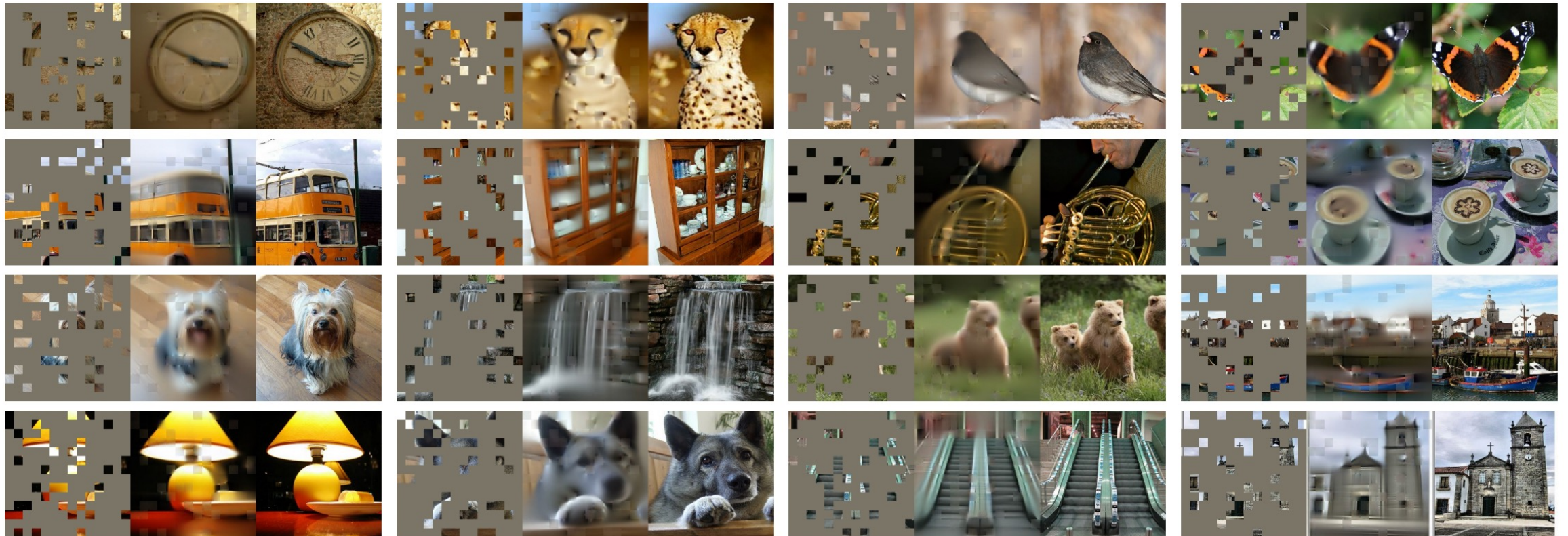
Pseudo label:  
Original colorful mage



# MAE



# MAE





# Video Generation with Colorization



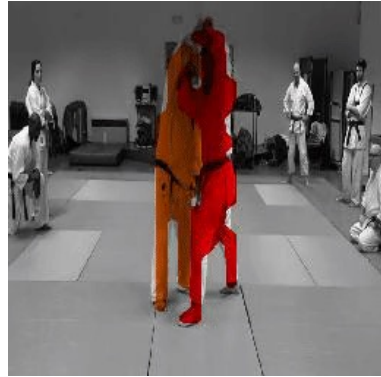
Reference  
colored frame

Input video

Predicted  
colorized video



# Video Generation with Colorization



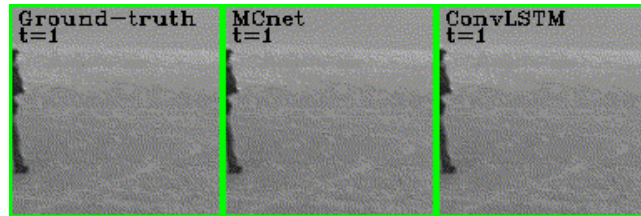
After learning to colorize videos, a mechanism for tracking automatically emerges without supervision.



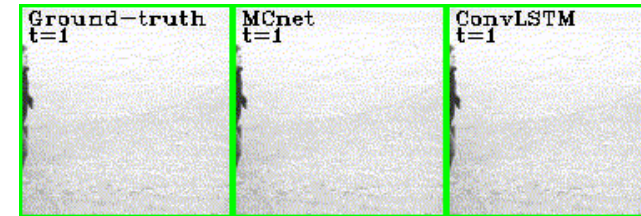


# Video Prediction

## Running



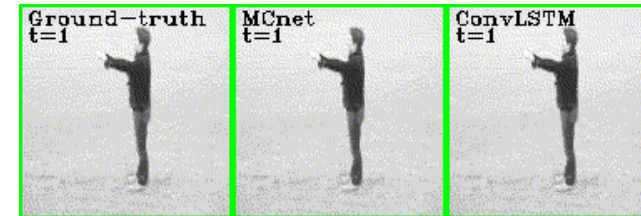
## Jogging



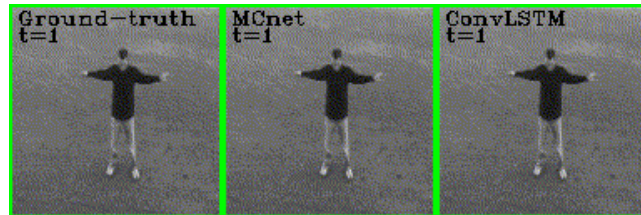
## Walking



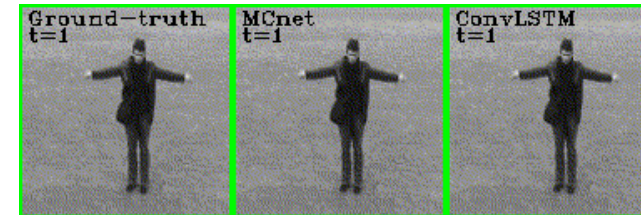
## Boxing



## Handclapping



## Handwaving



All models are trained to observe 10 frames (green) and predict 10 frames (red)





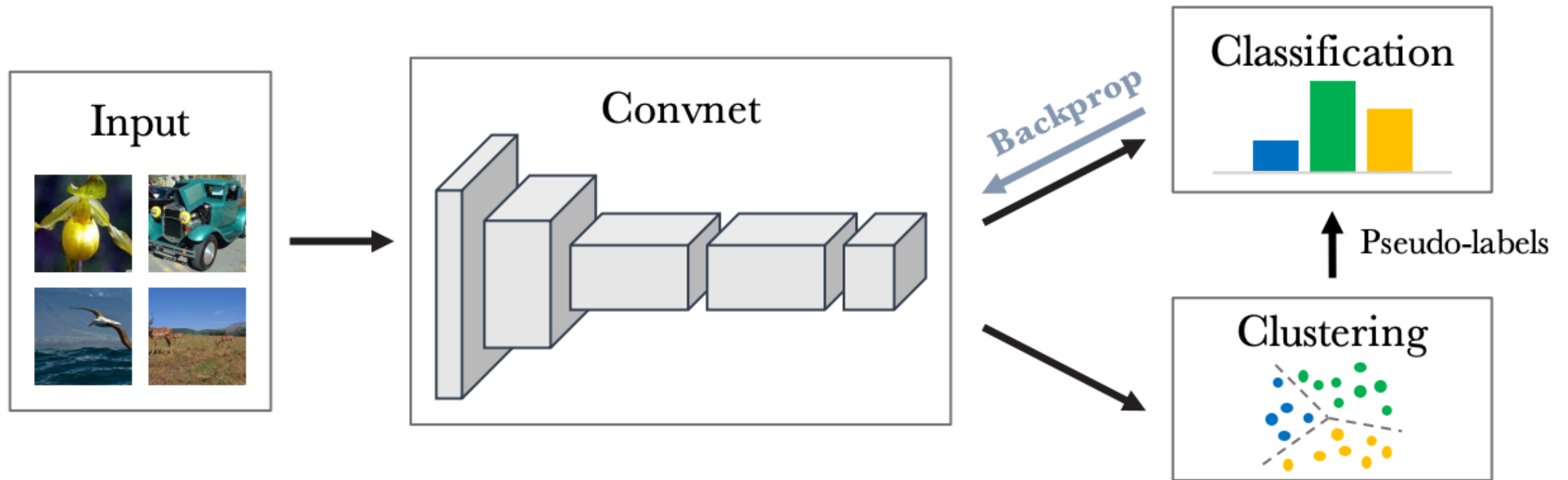


# CONTEXT-BASED METHODS

# Context-Based Methods

- The context-based pretext tasks mainly employ the context features of images as the supervision signal, including
  - context similarity;
  - spatial structure;
  - temporal structure;
  - ...

# Context Similarity

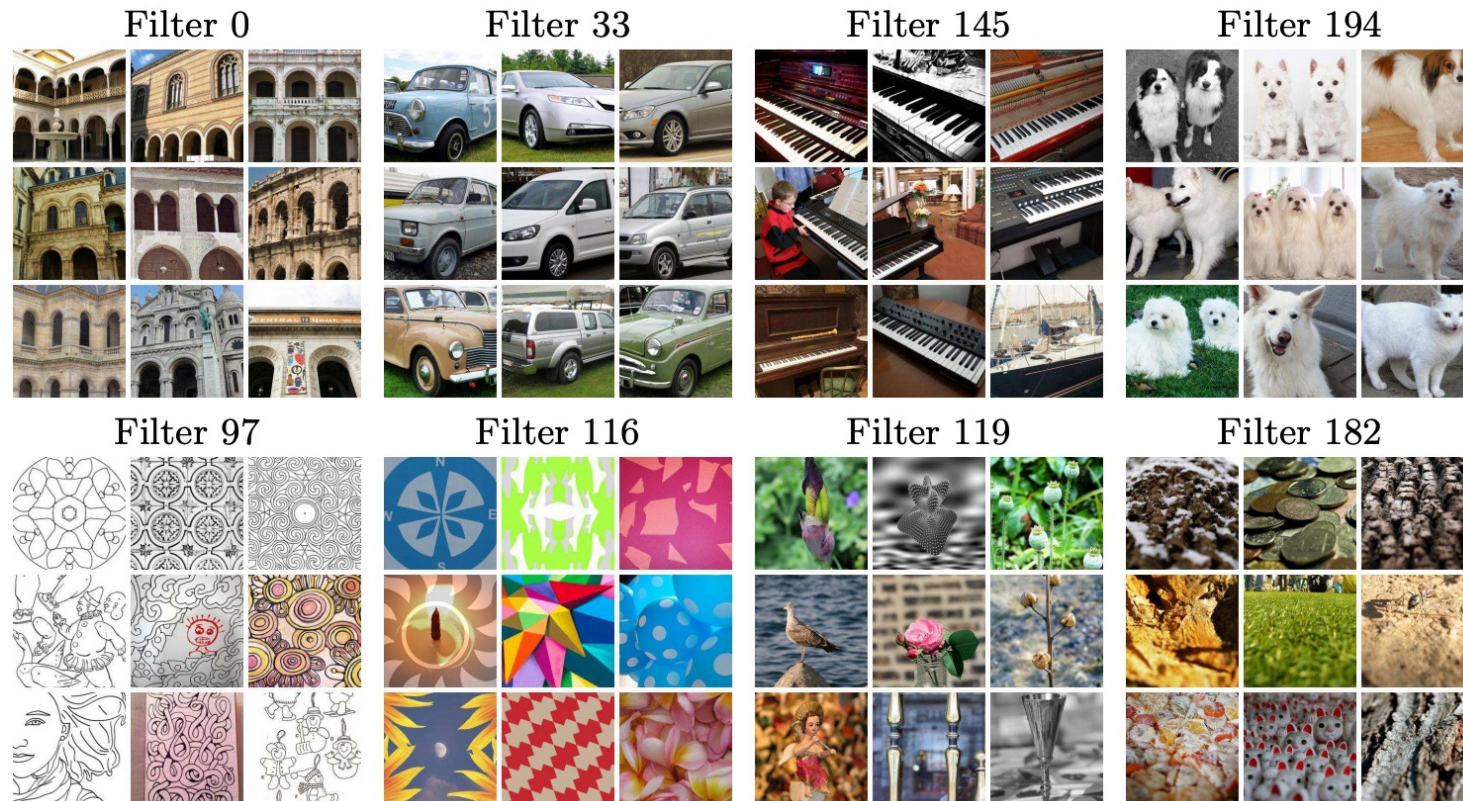


Iteratively cluster deep features and use the cluster assignments as pseudo-labels to learn the parameters of CNN



# Context Similarity

## ■ Clustering



Sensitive to  
specific objects

Sensitive to  
stylistic effect

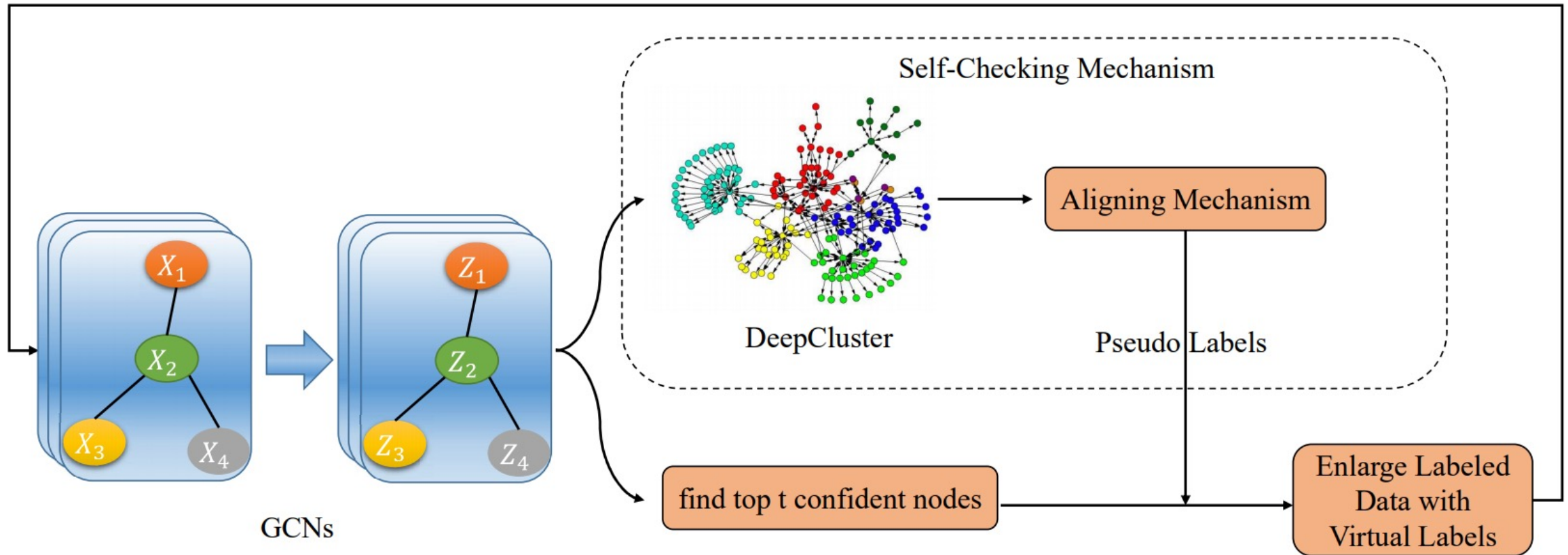
Top 9 activated images from a random subset of 10 millions images from YFCC100M for target filters in the last convolutional layer.



# Context Similarity

## ■ Clustering

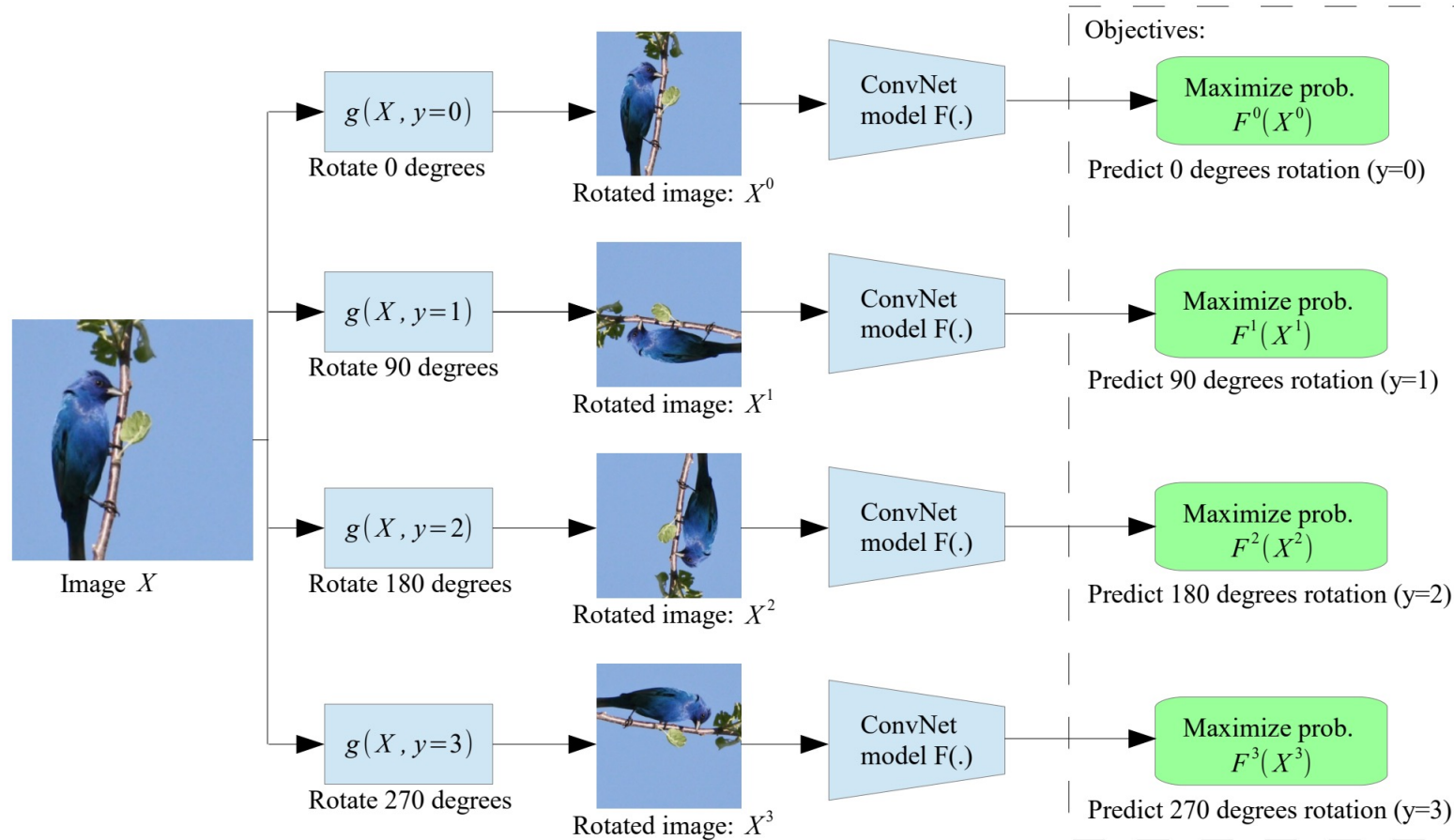
MultiStage Self-Training Framework





# Spatial Context Structure

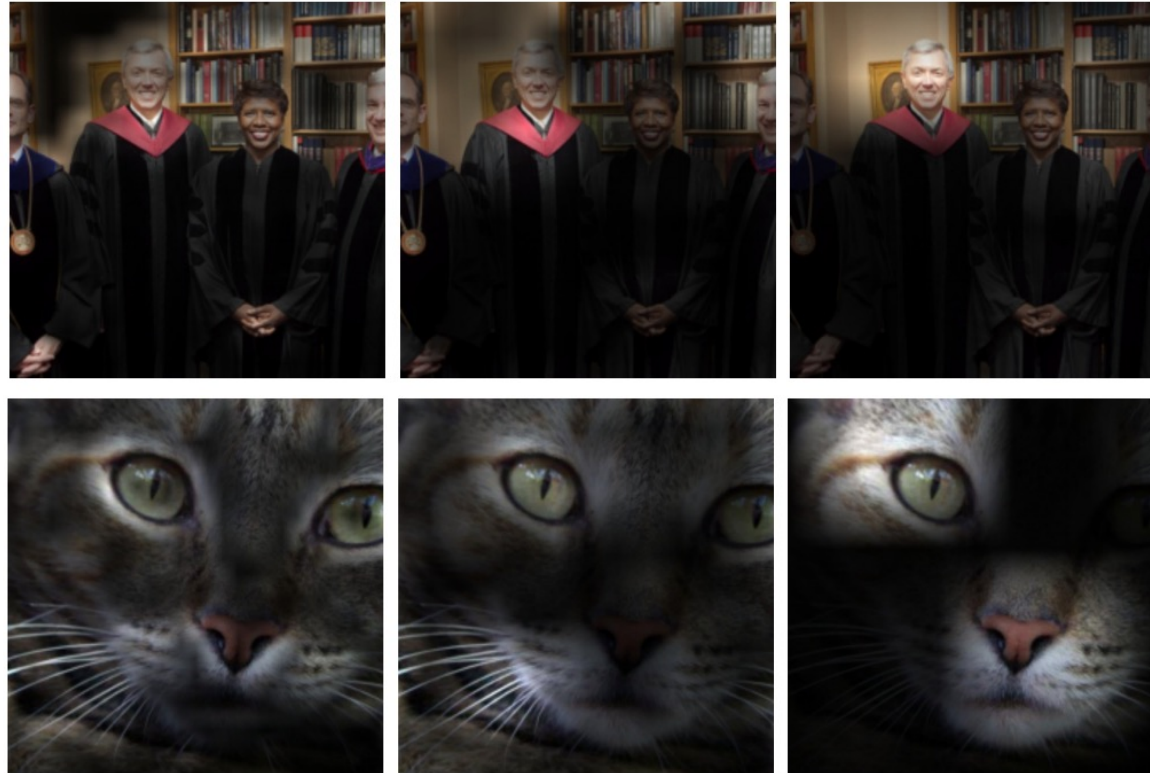
## ■ Rotation prediction





# Spatial Context Structure

## ■ Rotation prediction



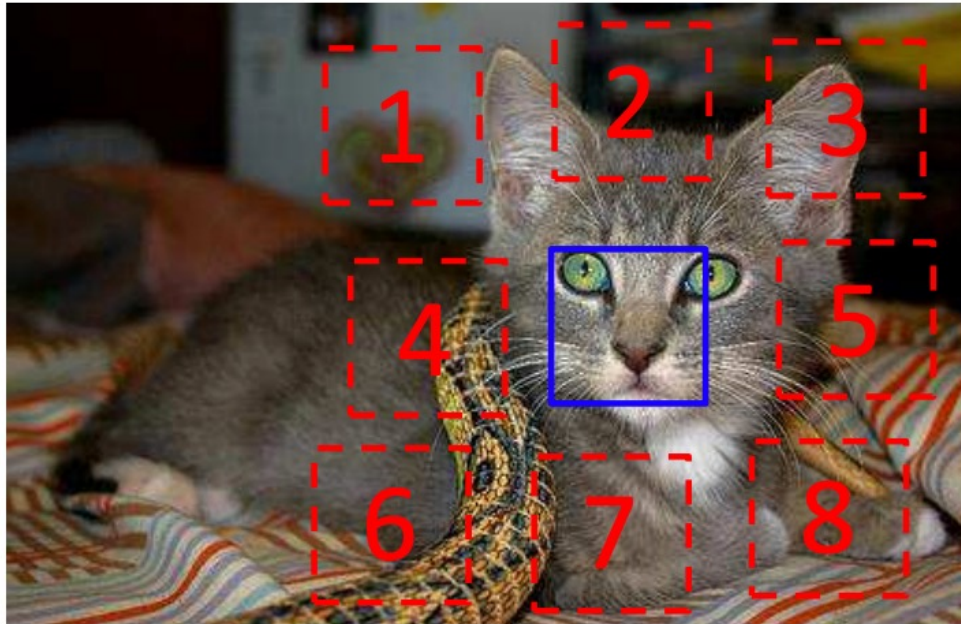
Conv1  $27 \times 27$  Conv3  $13 \times 13$  Conv5  $6 \times 6$

Attention maps of self-supervised model



# Spatial Context Structure

- Relative position prediction



$$X = (\text{cat face}, \text{cat ear}); Y = 3$$

Question 1:



Question 2:





# Spatial Context Structure

## Jigsaw puzzle

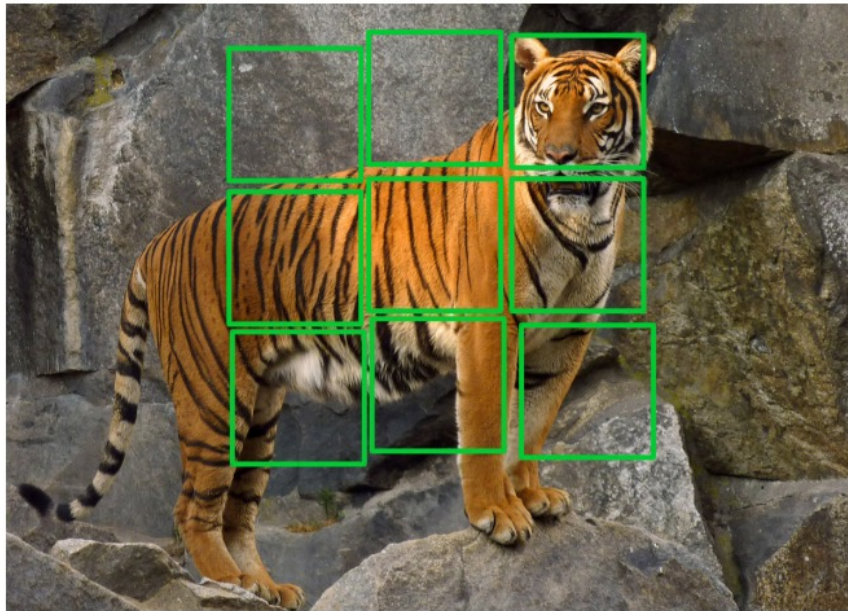
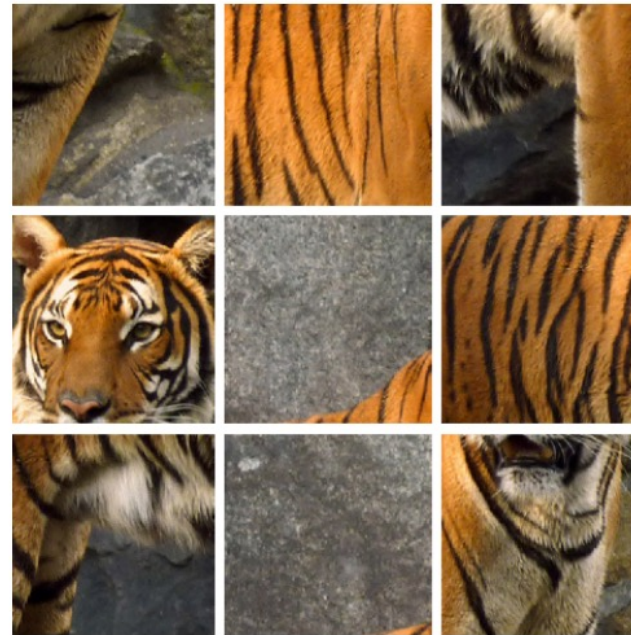


Image with 9 sampled image patches



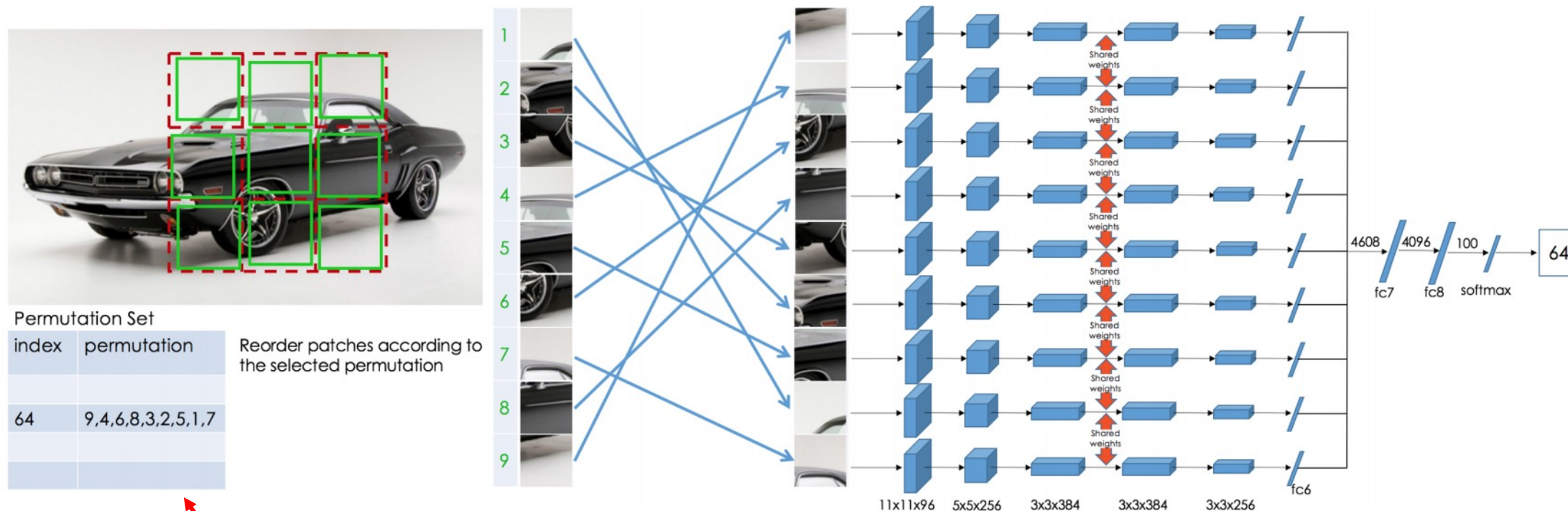
Shuffled image patches



Correct order of the sampled 9 patches

# Spatial Context Structure

## Jigsaw puzzle



How many classes?  $9! = 362,880$ .

# Spatial Context Structure

- Impossible to iterate over all possible permutation.
  - Similar permutation is somehow redundant.
- Choose permutations with max Hamming distance.

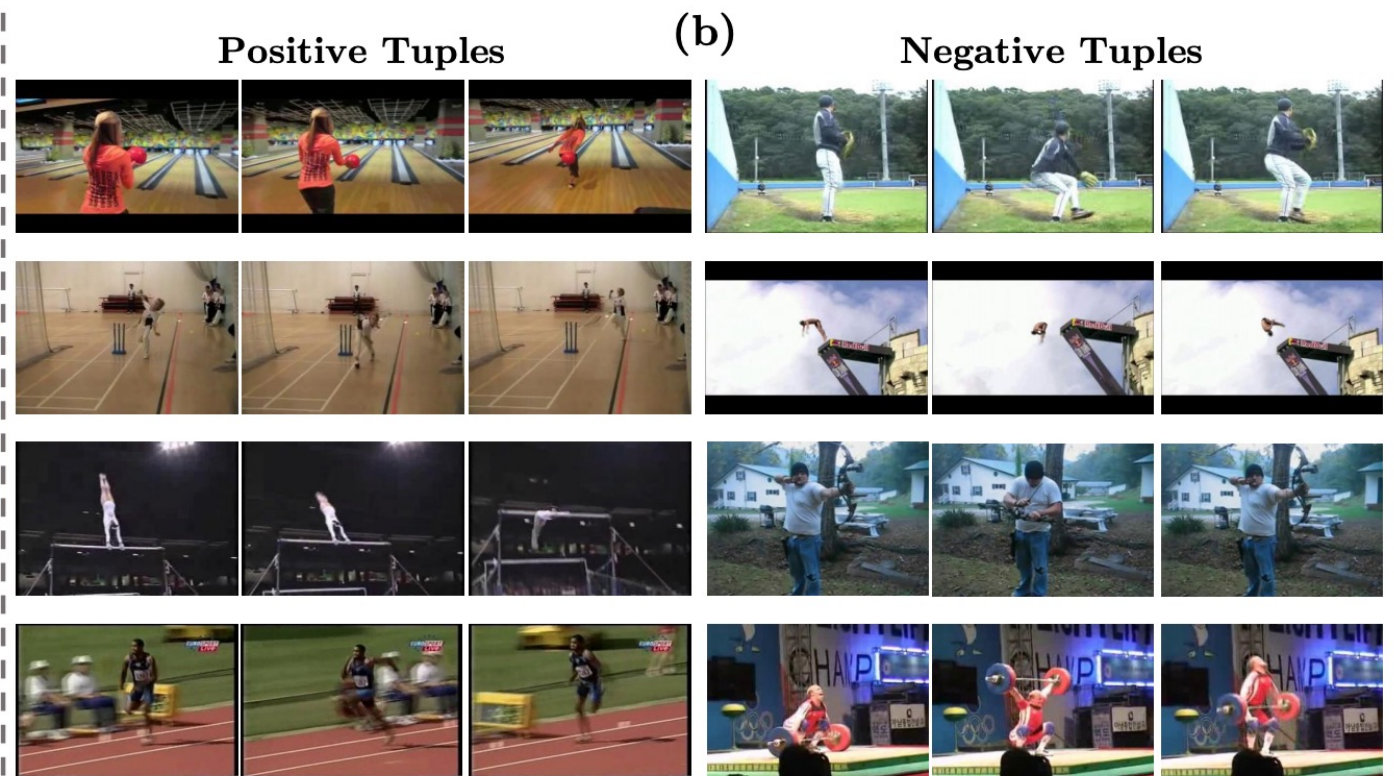
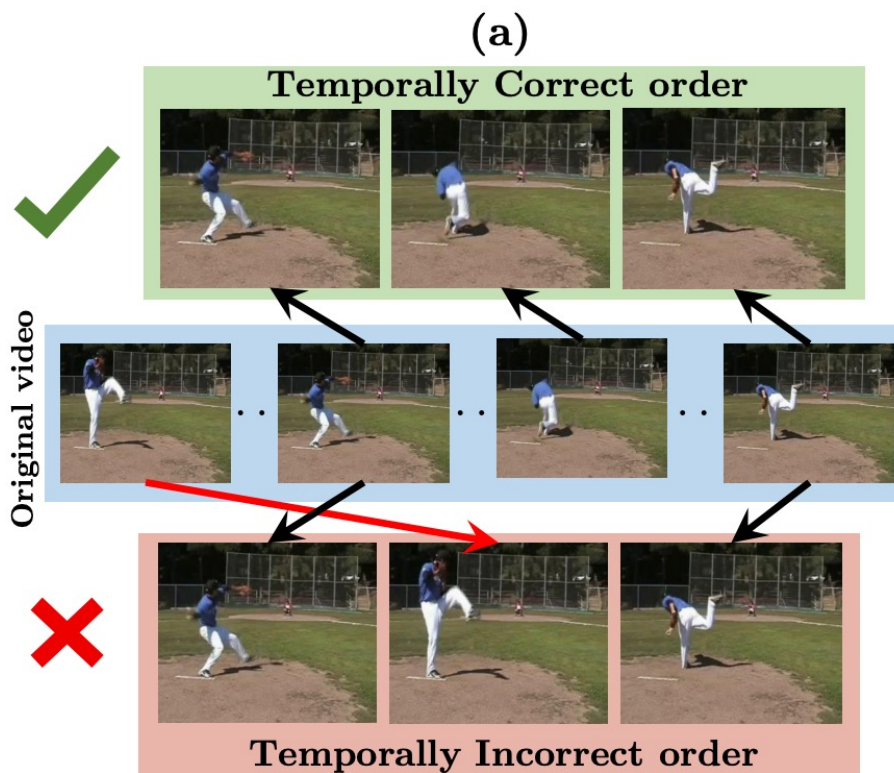
Number of permutations	Average hamming distance	Minimum hamming distance	Jigsaw task accuracy	Detection performance
1000	8.00	2	71	<b>53.2</b>
1000	6.35	2	62	51.3
1000	3.99	2	54	50.2





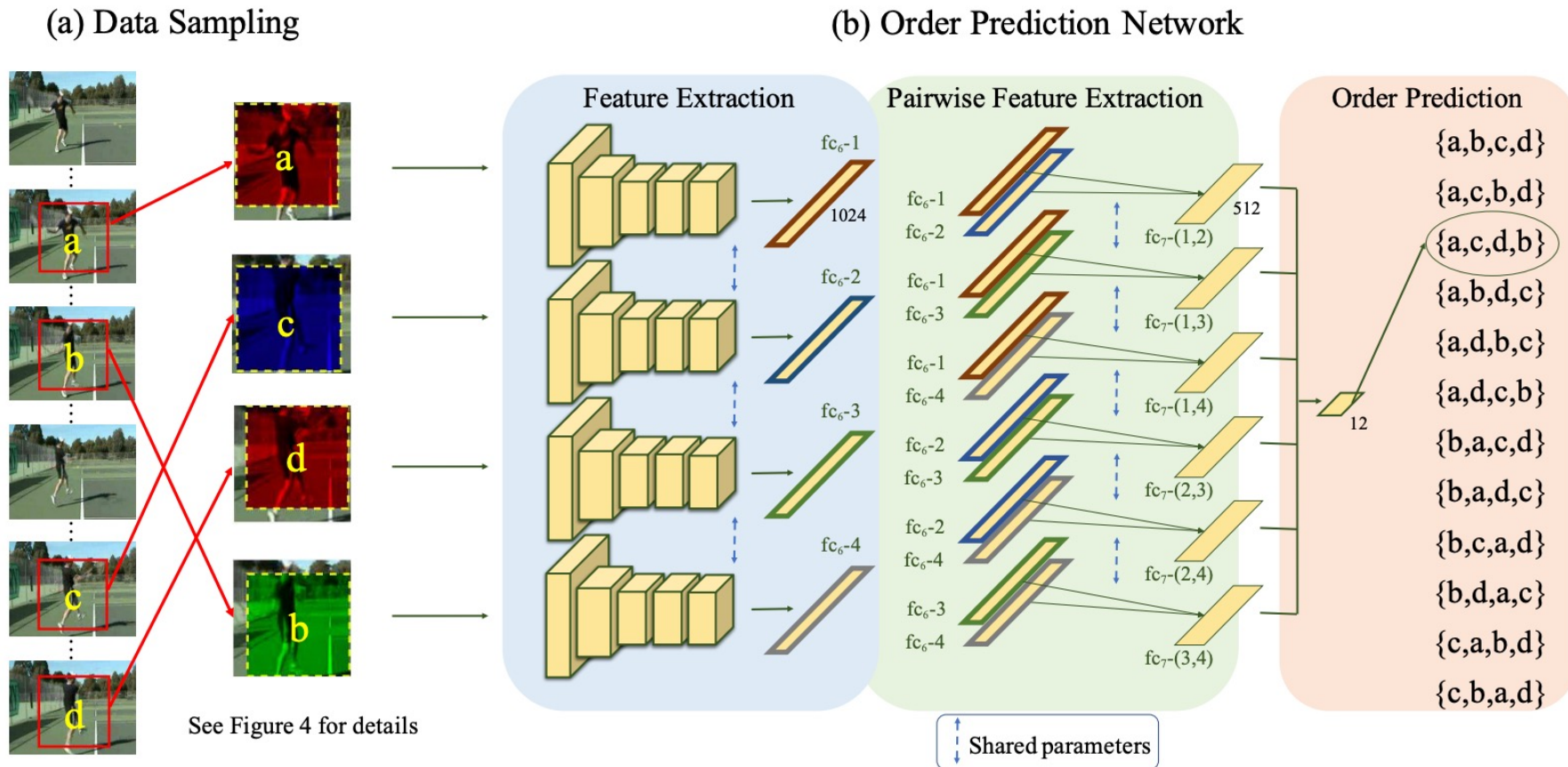
# Temporal Context Structure

## Temporal order verification



# Temporal Context Structure

## Temporal order prediction





# Sentence Context Structure

- Emoji prediction



# Sentence Context Structure

- Sentence permutation and rotation

I did X. Then I did Y. Finally I did Z.

I am going outside. I will be back in the evening.

original text



# Sentence Context Structure

- Gap sentence generation

TRANSFORMER



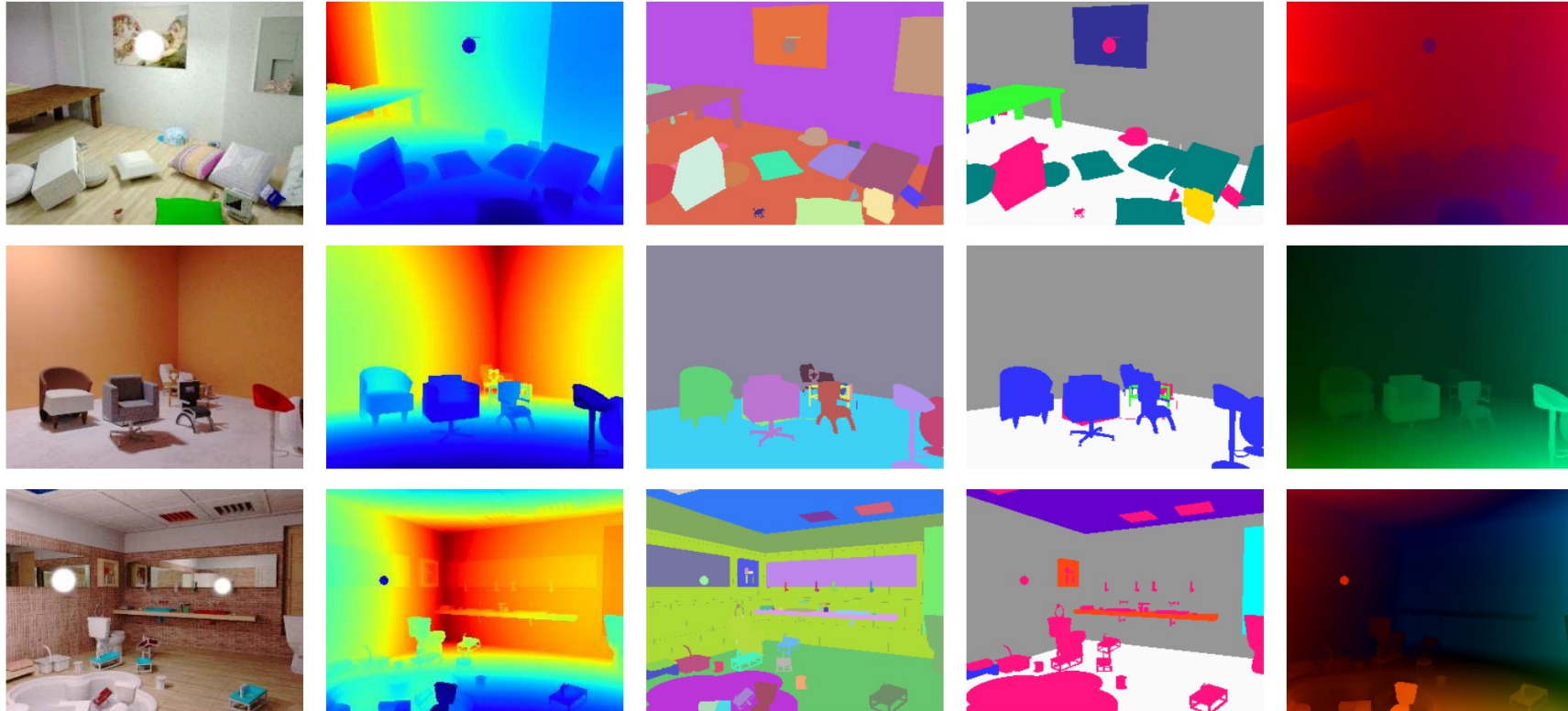


# FREE SEMANTIC LABEL-BASED METHODS

# Free Semantic Label-Based Methods

- Self-supervised learning requires no human annotations.
- Alternatively, we may obtain some semantic information as labels by.
  - Game engines: generate realistic images with accurate pixel-level labels with very low cost.
  - Auxiliary automatic annotators: generate saliency, foreground masks, contours, depth for images and videos.

# Game Engines



Synthetic  
image

Depth

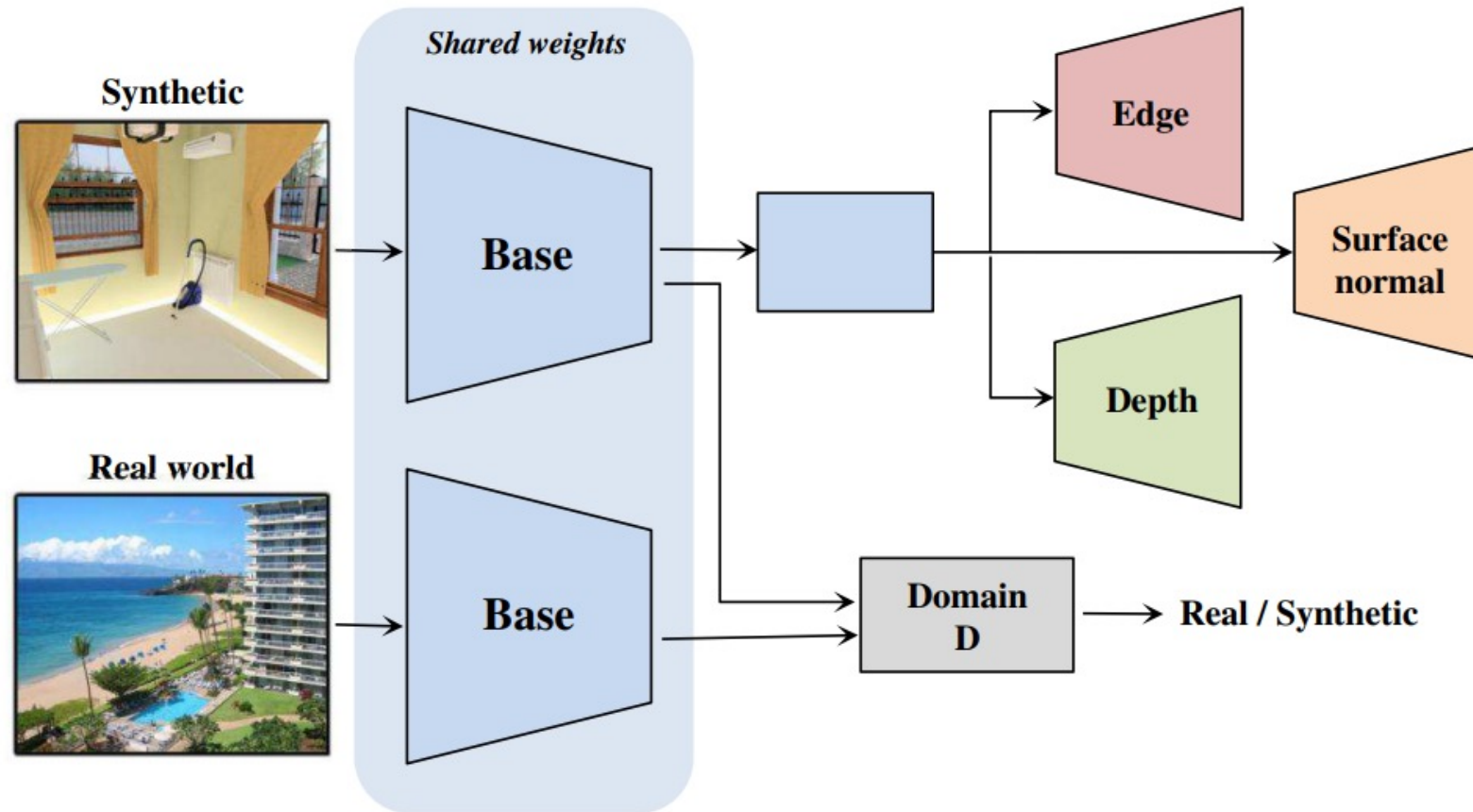
Segmentation

Semantic  
segmentation

Optical  
flow



# Game Engines

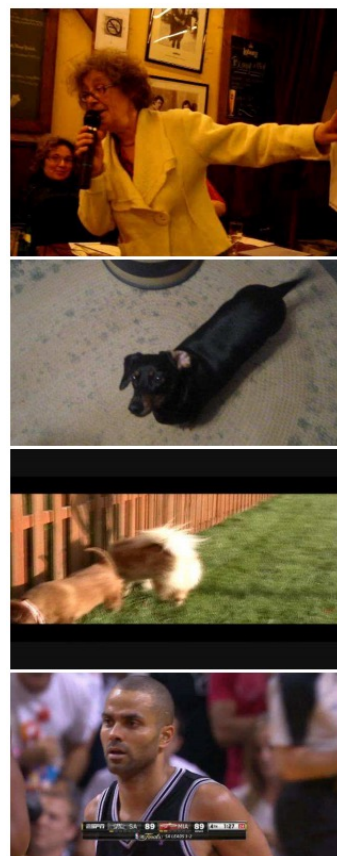


A discriminator network  $D$  is employed to minimize the difference of feature space domains between real-world and synthetic data

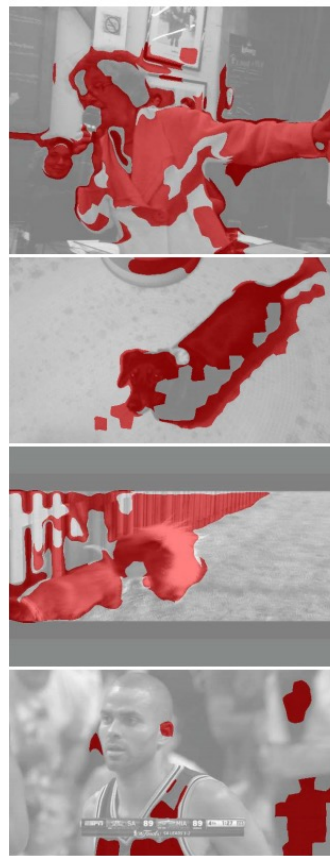




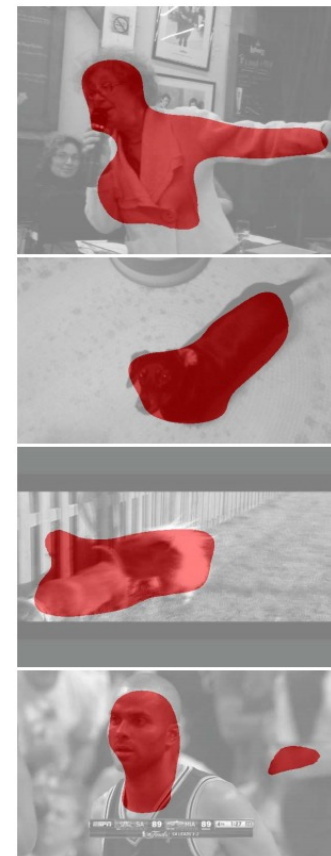
# Auxiliary Automatic Annotators



A video frame



Auxiliary motion detector



Trained detector



# Auxiliary Automatic Annotators



- Top: input image.
- Middle: relative depth image computed using a formula.
- Bottom: Predicted depth maps using our trained model.





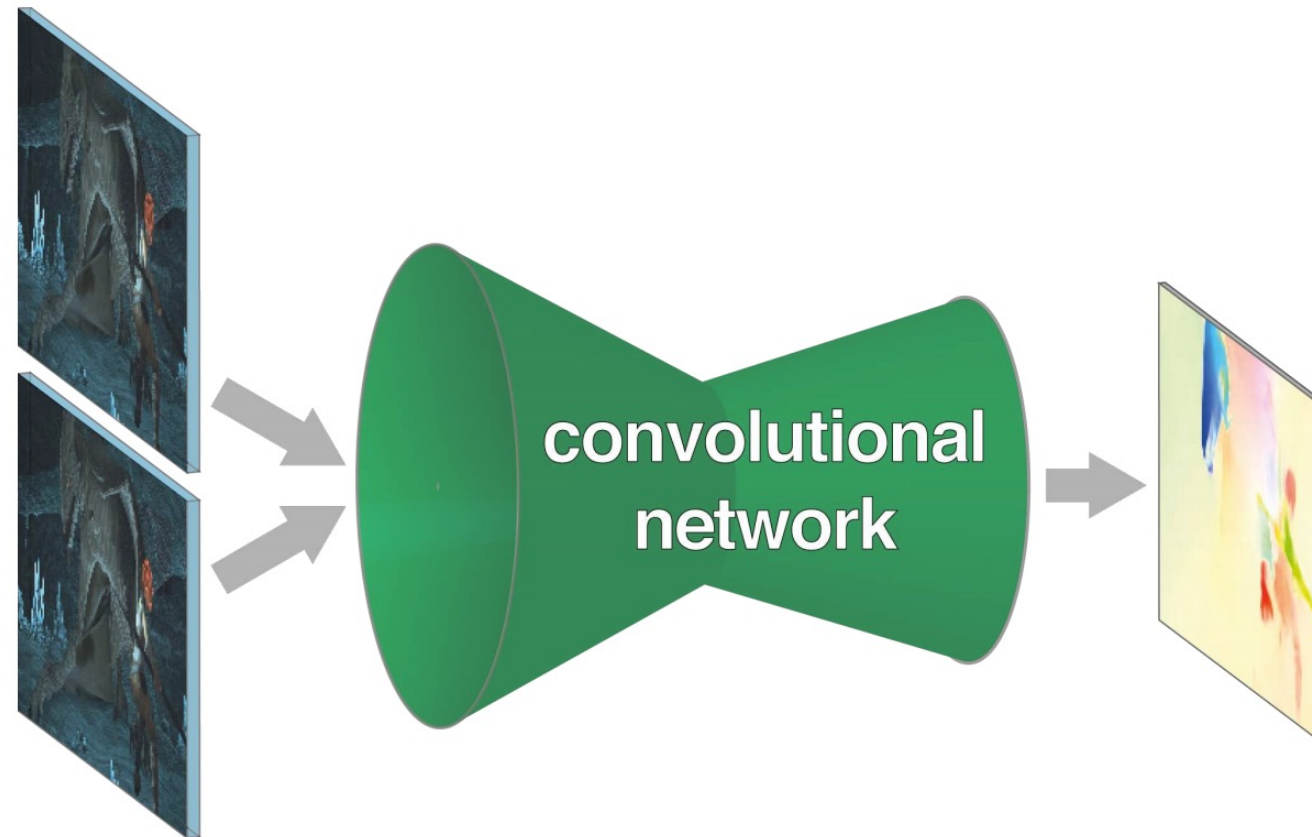
# CROSS MODAL-BASED METHODS

# Cross Modal-based Learning

- Use different modal as pseudo label.
  - Optical flow;
  - Audio;
  - Text;
  - Camera poses...



# RGB-Flow Correspondence

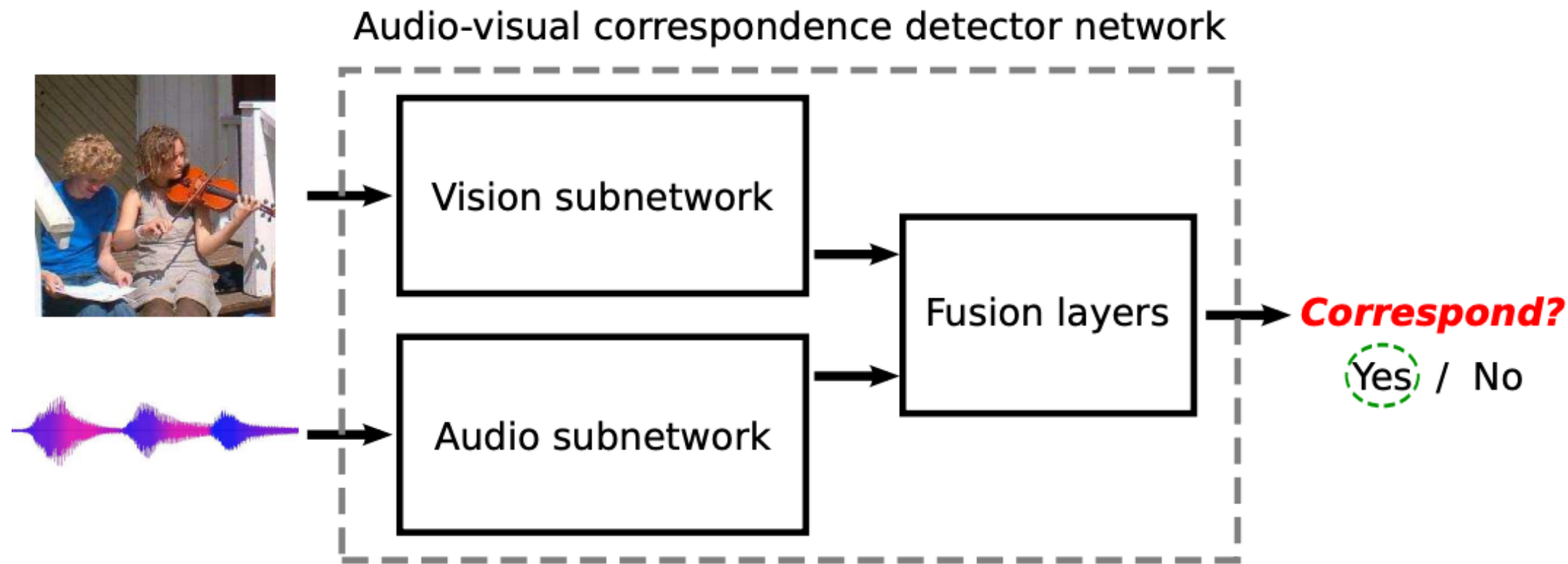


Large scale optical flow modal is also hard to obtain. It can also be generated by some auxiliary algorithm.





# Visual-Audio Correspondence

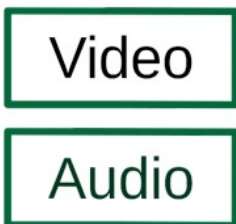


Learn to determine whether a pair of video and audio clip correspond to each other or not

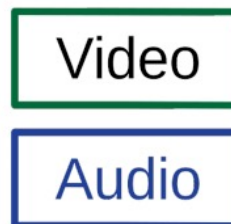


# Visual-Audio Correspondence

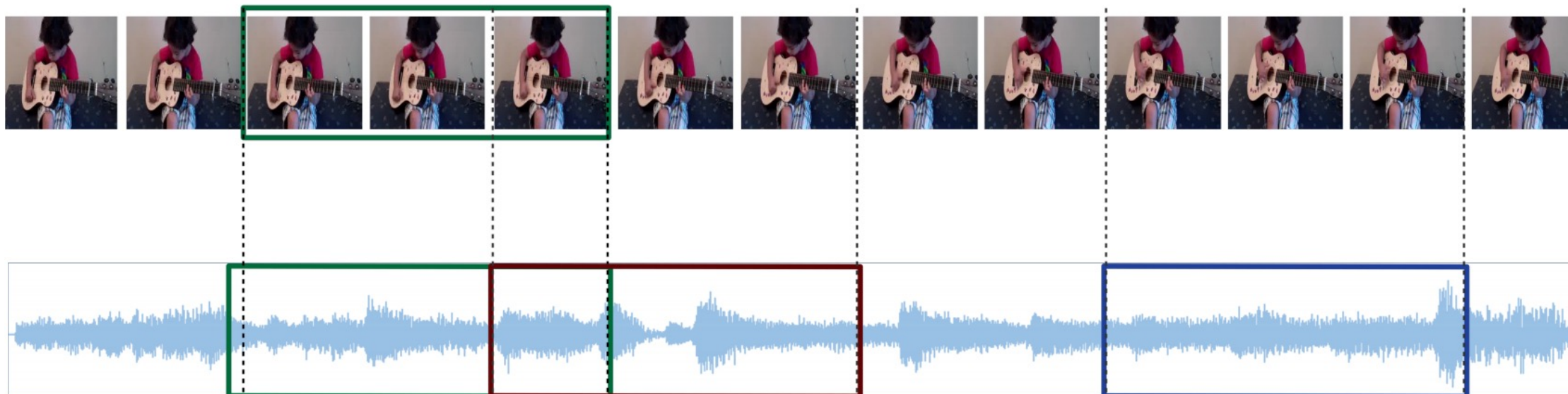
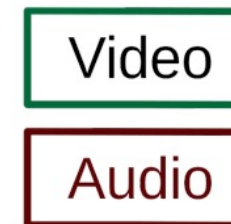
Positive pair



Hard negative pair



Super hard negative pair



# Visual-Audio Correspondence

## Objects that Sound

Relja Arandjelović<sup>1</sup>, Andrew Zisserman<sup>1,2</sup>  
<sup>1</sup>DeepMind <sup>2</sup>University of Oxford

Frames are processed completely independently, motion information is not used, and there is no temporal smoothing

Input single frame



Frame/  
Localization  
overlaid



Localization



## Localizing objects that sound



廈門大學信息學院 (特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)



廈門大學 计算机科学与技术系

Department of Computer Science and Technology, Xiamen University

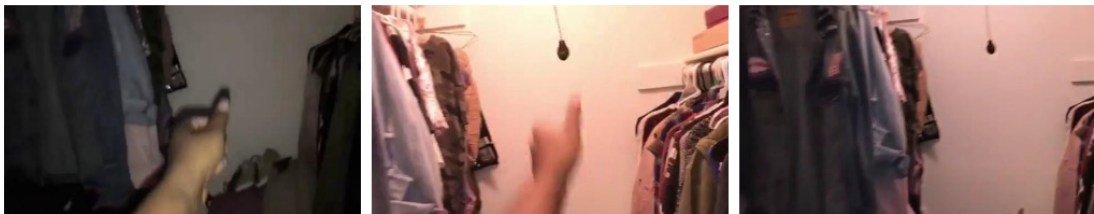
Reference: Arandjelovic, Relja, and Andrew Zisserman. "Objects that sound." In Proceedings of the European Conference on Computer Vision (ECCV), pp. 435-451. 2018.

# Visual-Text Correspondence

## Strongly related pairs

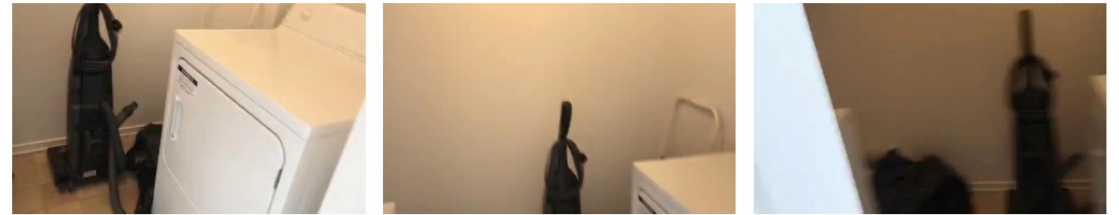


Subtitle: Over here is my bike. I love my bike.

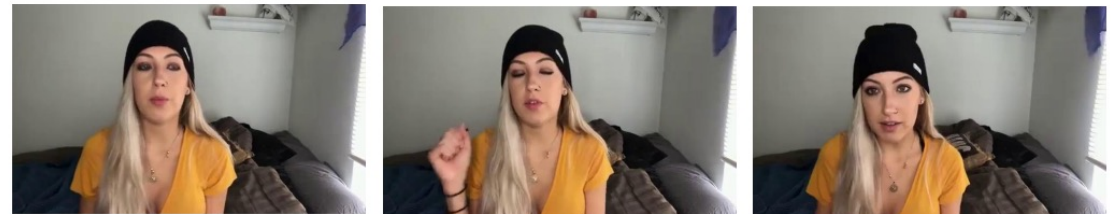


Subtitle: Let me grab the light. It's my closet.

## Weakly or not related pairs



Subtitle: It's just a mess in here right now.



Subtitle: My sister's going back to school.

Use subtitle as supervision. But subtitles usually contain large number of noises.



# CONTRASTIVE LEARNING



# Contrastive Learning

Generate  $\longleftrightarrow$  Distinguish  
?



Do we have to able to draw cash, in order to distinguish cash?



## Generative / Predictive



## Contrastive



# Contrastive Learning

- For any data point  $\mathbf{x}$ , which is commonly referred to as an “anchor” data point, contrastive methods aim to learn a feature mapping  $f$  such that:

$$\text{score}(f(\mathbf{x}), f(\mathbf{x}^+)) \gg \text{score}(f(\mathbf{x}), f(\mathbf{x}^-)).$$

- $\mathbf{x}^+$  is a data point similar to  $\mathbf{x}$ , referred to as a positive sample.
- $\mathbf{x}^-$  is a data point dissimilar to  $\mathbf{x}$ , referred to as a negative sample.
- the score function is a metric that measures the similarity between two features.



# Contrastive Learning

- To optimize for this property, we can construct a softmax classifier that classifies positive and negative samples correctly:

$$\mathcal{L} = -\mathbb{E}_{\mathcal{X}} \left[ \log \frac{\exp(f(\mathbf{x})^T f(\mathbf{x}^+))}{\exp(f(\mathbf{x})^T f(\mathbf{x}^+)) + \sum_{j=1}^{N-1} \exp(f(\mathbf{x})^T f(\mathbf{x}_j))} \right]$$

- It is commonly called the **InfoNCE loss** in the contrastive learning literature.
- But the key problem is:

How do we know data similarity?



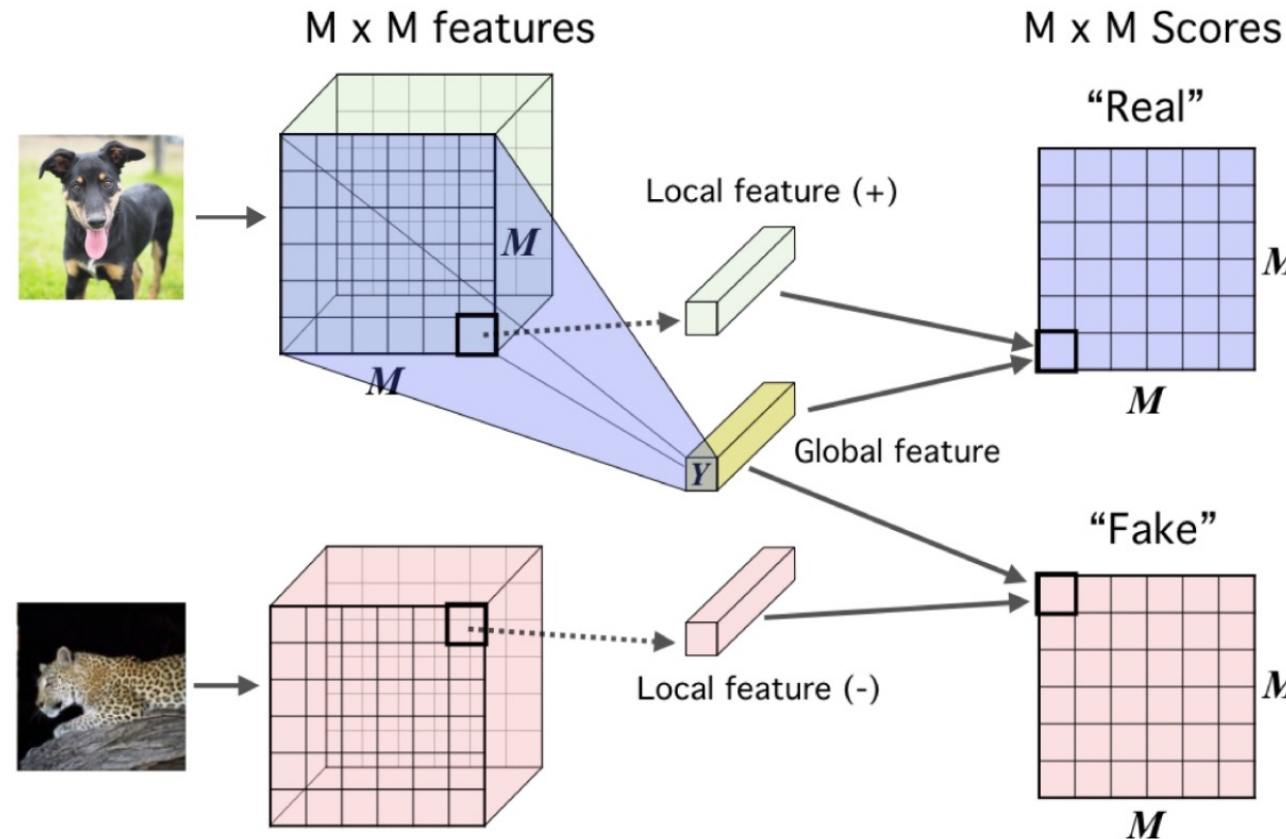
# Deep Infomax

Learning deep representations by mutual information estimation and maximization

RD Hjelm, A Fedorov, S Lavoie-Marchildon... - arXiv preprint arXiv ..., 2018 - arxiv.org

In this work, we perform unsupervised learning of representations by maximizing mutual information between an input and the output of a deep neural network encoder. Importantly, we show that structure matters: incorporating knowledge about locality of the input to the objective can greatly influence a representation's suitability for downstream tasks. We further control characteristics of the representation by matching to a prior distribution adversarially. Our method, which we call Deep InfoMax (DIM), outperforms a number of popular ...

☆ 被引用 483 相关文章 所有 4 版本



Classify whether a pair of global features and local features are from the same image or not.





# Deep Infomax

- Global features  $E_\psi(X)$  are the final output of a convolutional encoder.
- Local features  $C_\psi^{(i)}(X)$  are the output of an intermediate layer in the encoder (an  $M \times M$  feature map).
  - Each local feature map has a limited receptive field.
- We want to maximize the mutual information between local and global features of the same image:

$$\operatorname{argmax}_{\omega, \psi} \frac{1}{M^2} \sum_{i=1}^{M^2} I_{\omega, \psi} \left( C_\psi^{(i)}(X); E_\psi(X) \right)$$

and minimize it for different image.

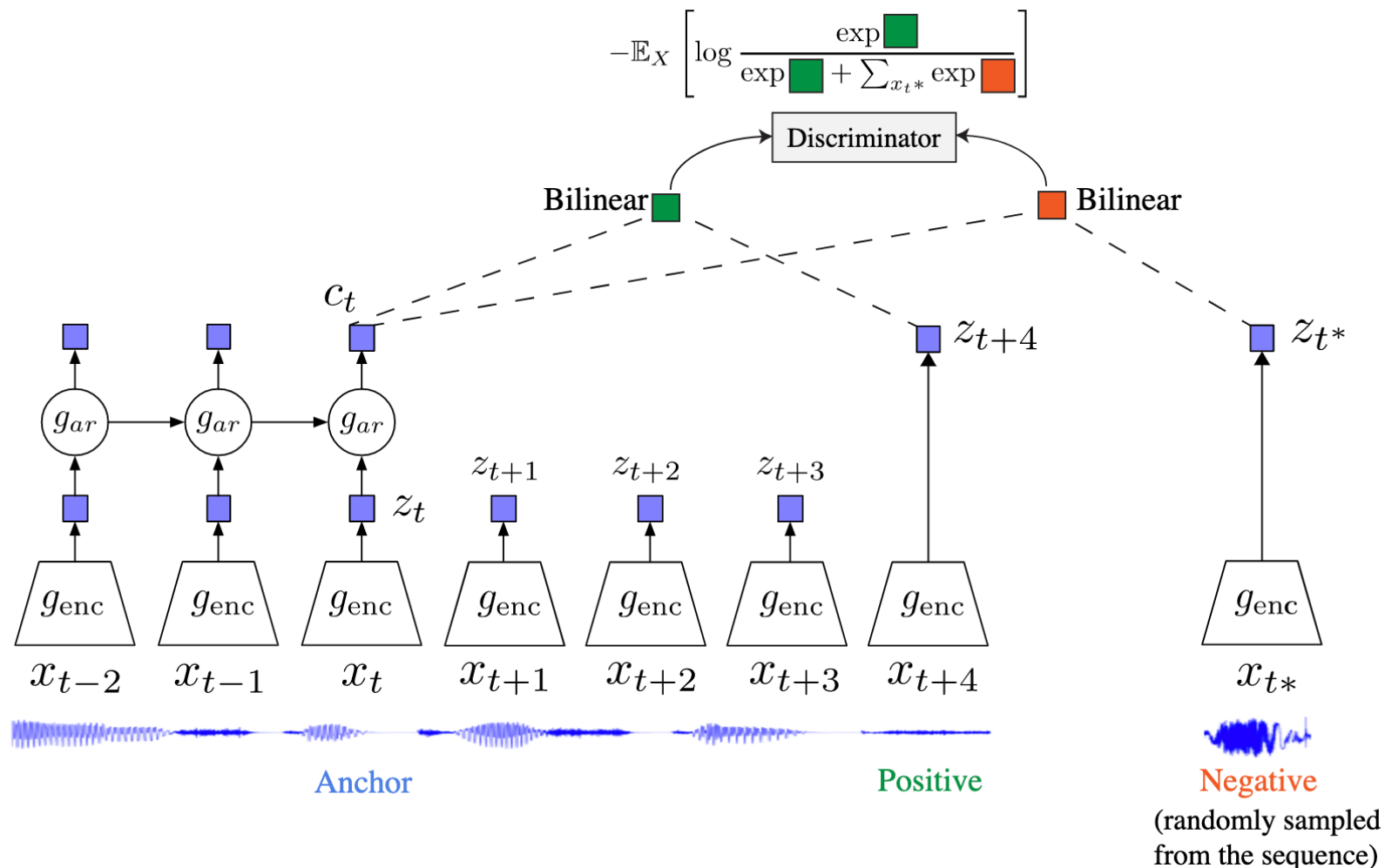
# Contrastive Predictive Coding

## Representation learning with contrastive predictive coding

A Oord, Y Li, O Vinyals - arXiv preprint arXiv:1807.03748, 2018 - arxiv.org

... learning approach to extract useful **representations** from high-dimensional data, which we call **Contrastive** ... The key insight of our model is to **learn** such **representations** by predicting the ...

☆ Save 📄 Cite Cited by 6707 Related articles All 4 versions 🔗



# Contrastive Multiview Coding

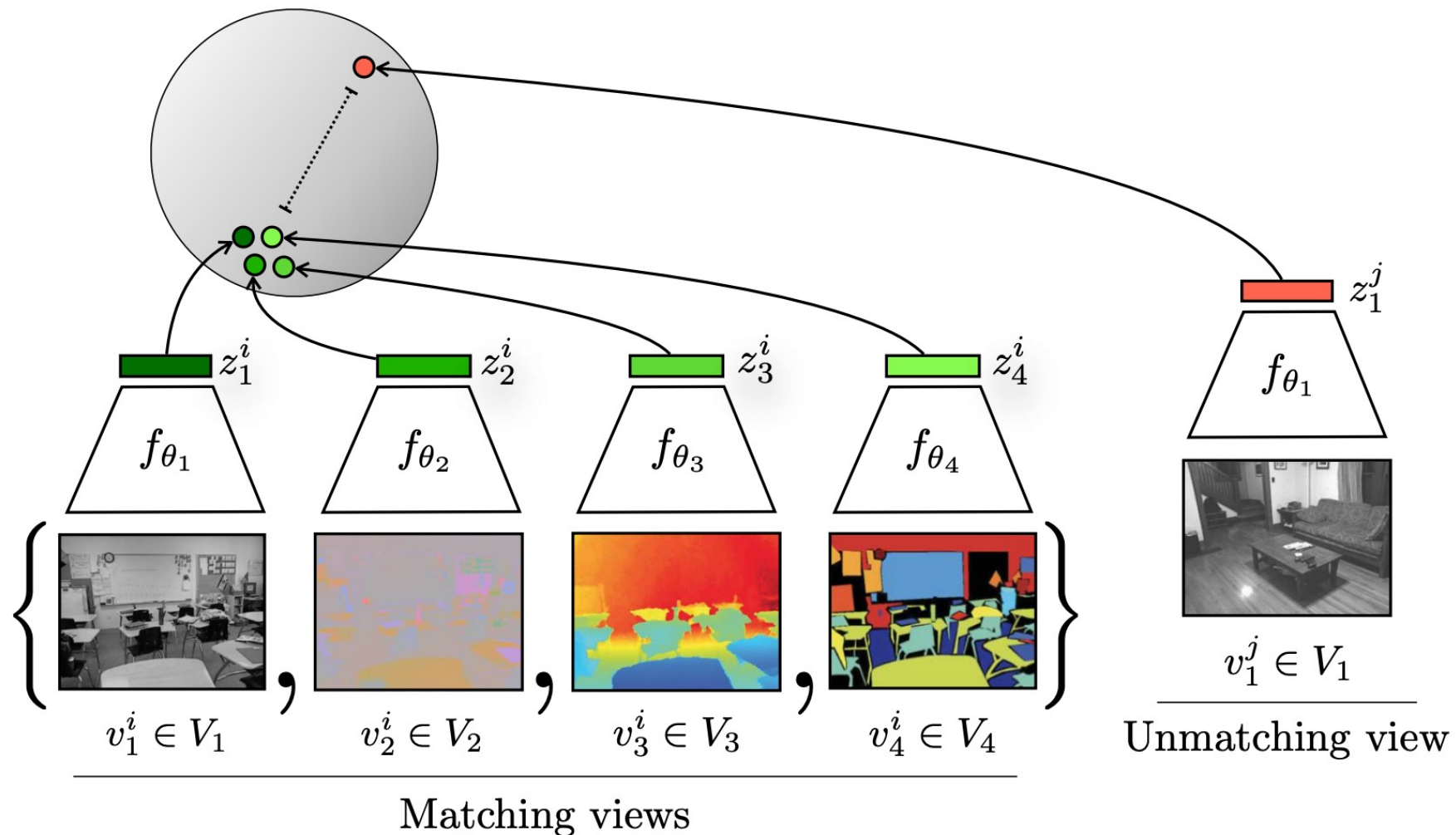
## Contrastive multiview coding

Y Tian, D Krishnan, P Isola - Computer Vision—ECCV 2020: 16th European ..., 2020 - Springer

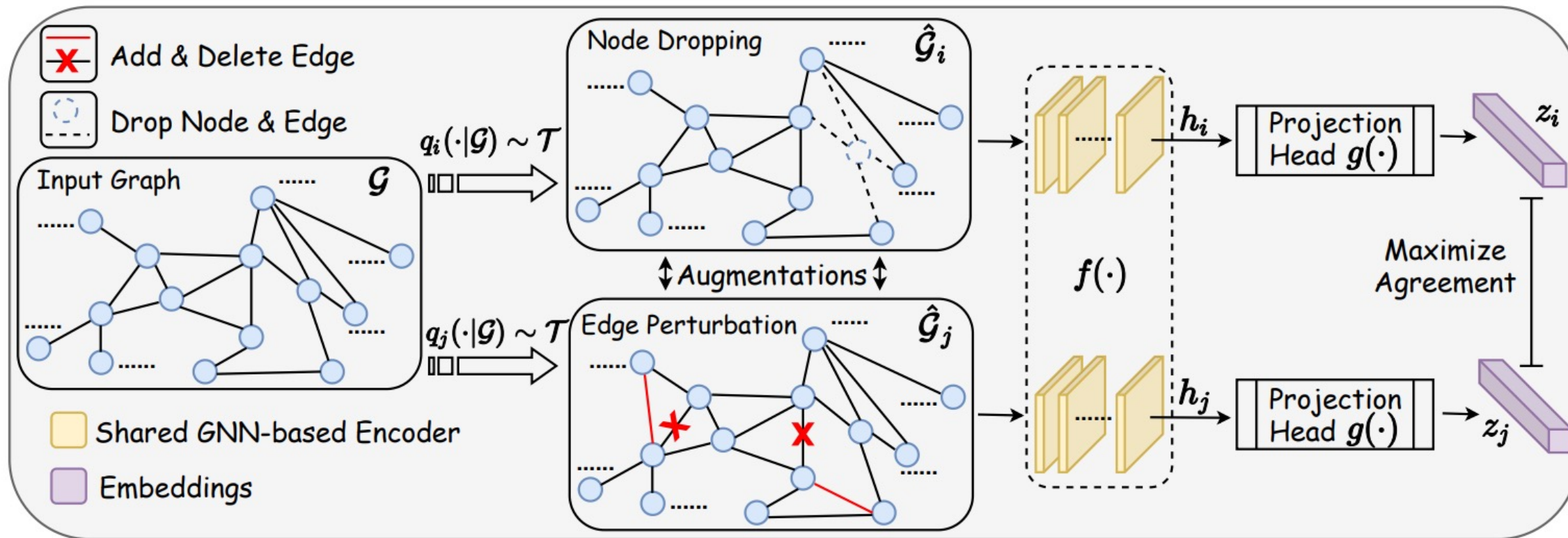
... Finally, we compare the **contrastive** formulation of **multiview** learning to the recently ...

**contrastive** approach learns stronger representations. The core ideas that we build on: **contrastive** ...

☆ Save 引 Cite Cited by 2082 Related articles All 11 versions

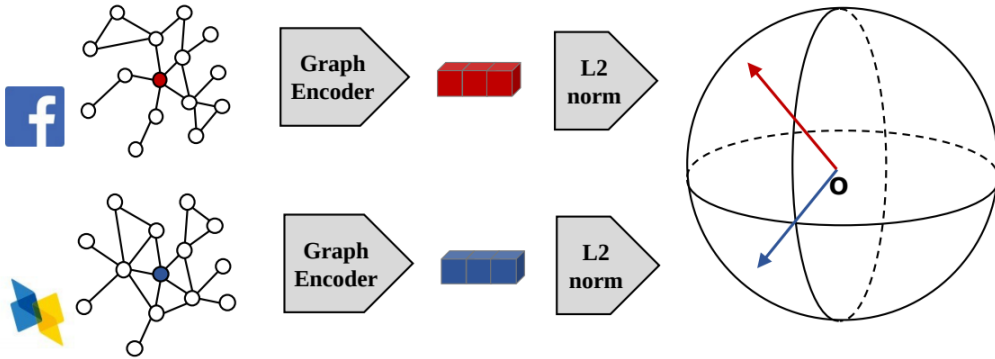


# Graph Contrastive Learning

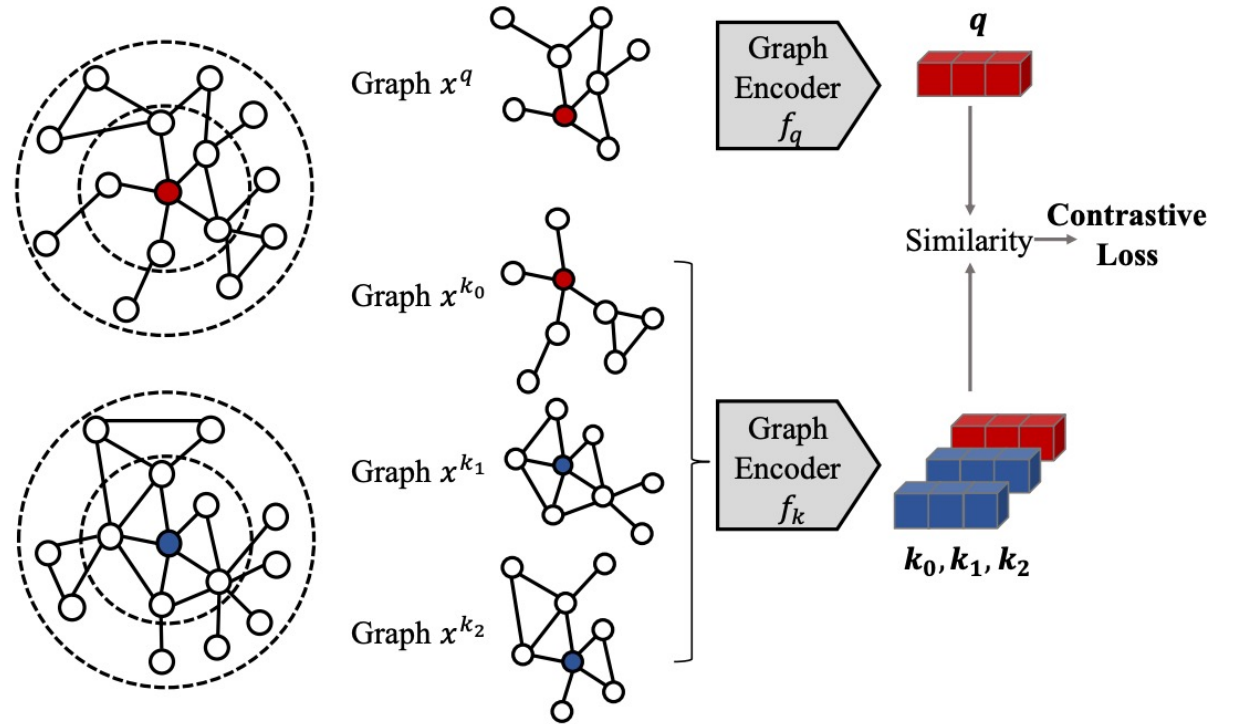


# Graph Contrastive Coding

## Facebook social graph



## DBLP co-authorship graph

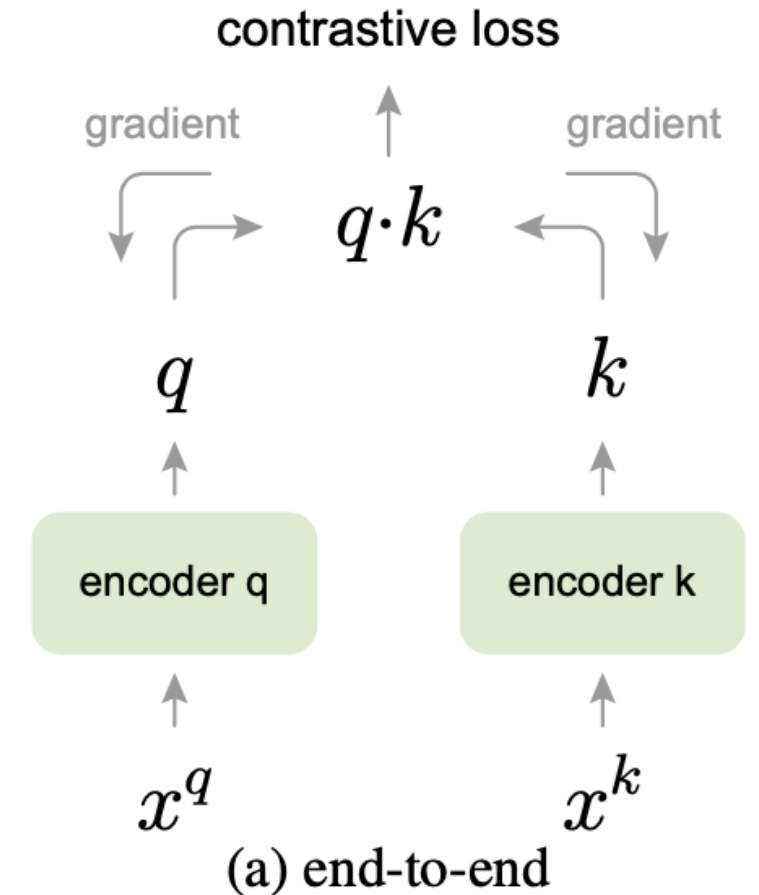


Capture the universal network topological properties across multiple networks



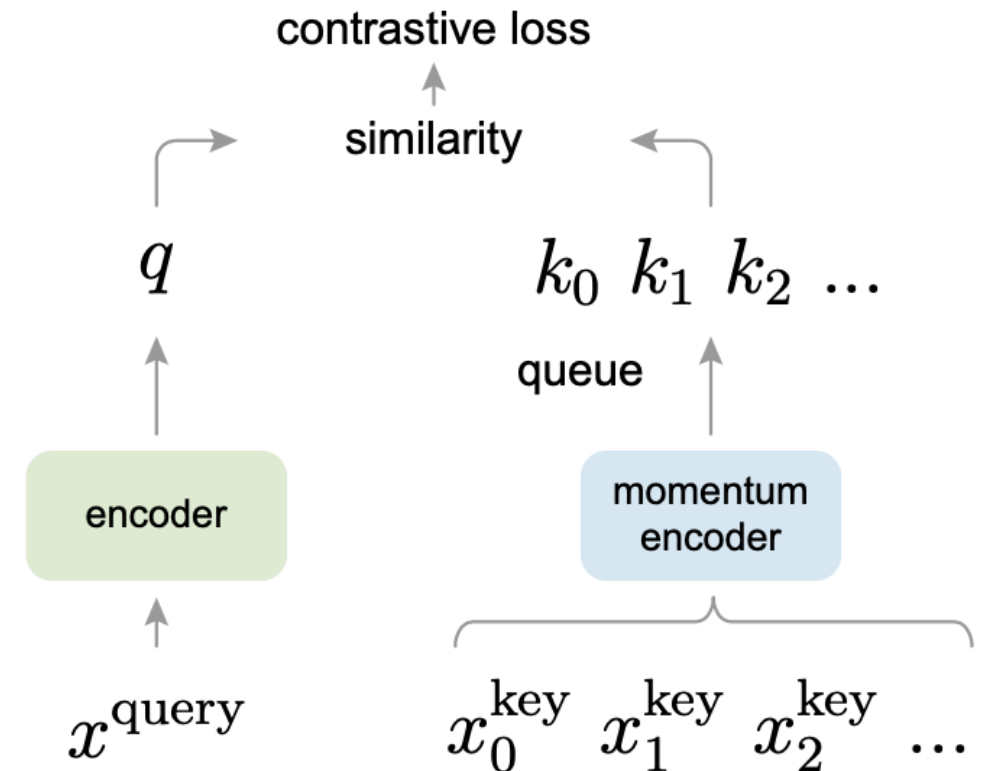


- A general framework for contrastive learning optimization.
- Motivation: number of negative samples should be **large** to make contrast.
  - However, the size is usually limited by batch size and GPU memory size.
- Idea: Reuse the representations of negative samples.



# MoCo

- Contrastive learning can be thought of as training an encoder for a **dictionary look-up** task.
- Query is the anchor sample.  $N$  keys contains 1 positive sample and  $N - 1$  negative samples.
- $q$  and  $k_0, k_1, \dots$  are encoded samples. InfoNCE is calculated on them.



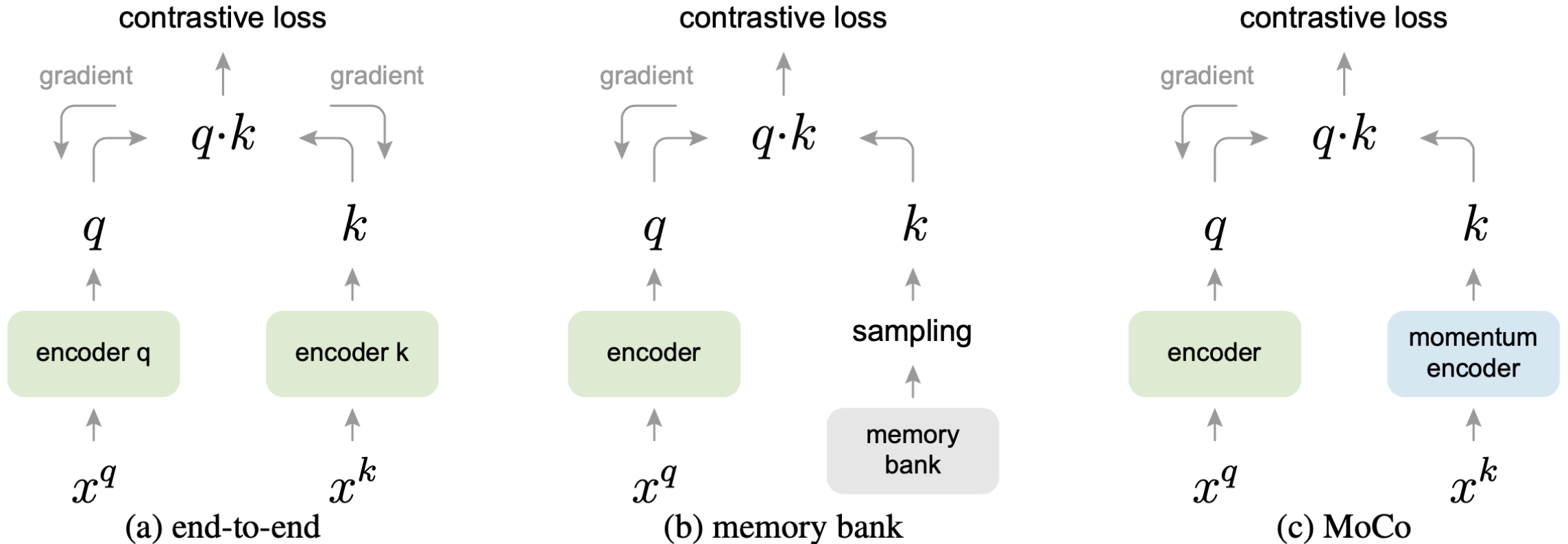
- Use a queue to store encoded negative samples for reuse.
- The queue is dynamic updated during training.
- Momentum update is adopted to slow down the frequency of key encoder:

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

- $m$  is usually set very close to 1 (e.g. 0.999).



# MoCo



### 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
```

No parameter  
update for keys

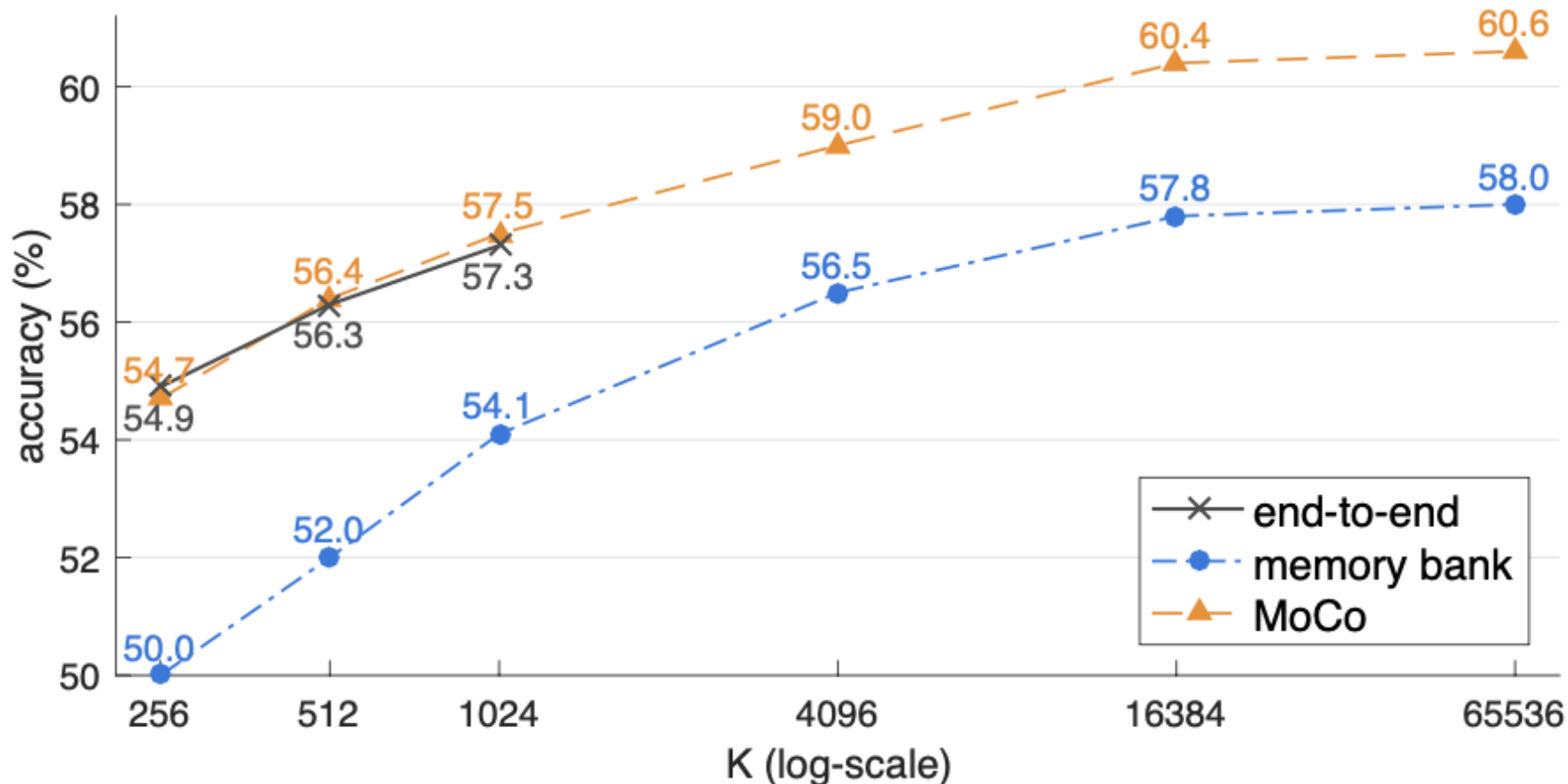
Use all negative  
samples in the queue

Momentum update

Enqueue negative samples



# MoCo



Comparison of three contrastive loss mechanisms under the ImageNet linear classification protocol



## Contributions:

- Use data augmentations.
- Introduce a learnable nonlinear transformation between the representation and the contrastive loss.
- Contrastive learning benefits from larger batch sizes and more training steps compared to supervised learning.
  - Batch size 8192 with 128 TPU v3 cores...

# SimCLR



廈門大學信息學院(特色化示范性软件学院)

School of Informatics Xiamen University (National Characteristic Demonstration Software School)

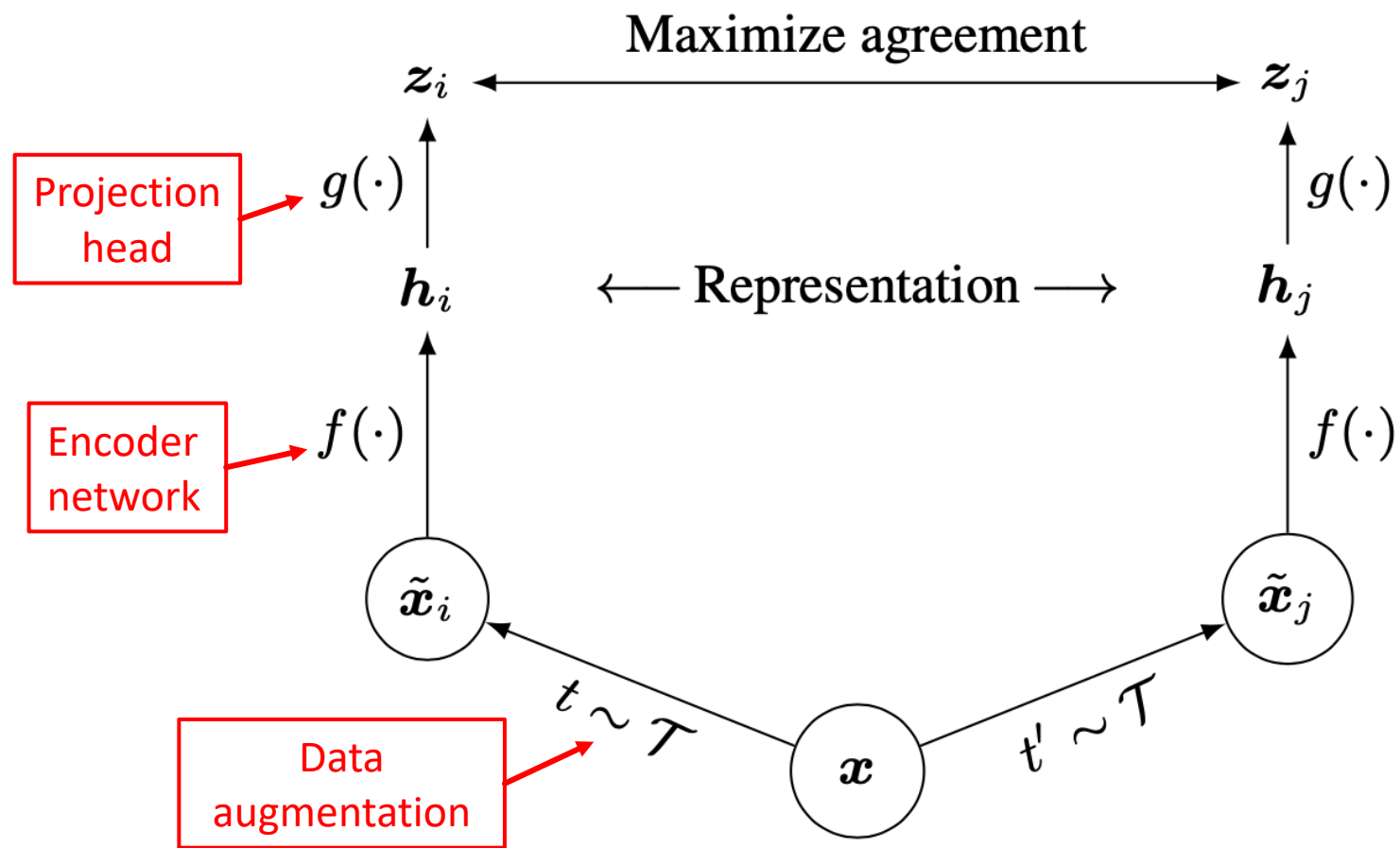


廈門大學 计算机科学与技术系

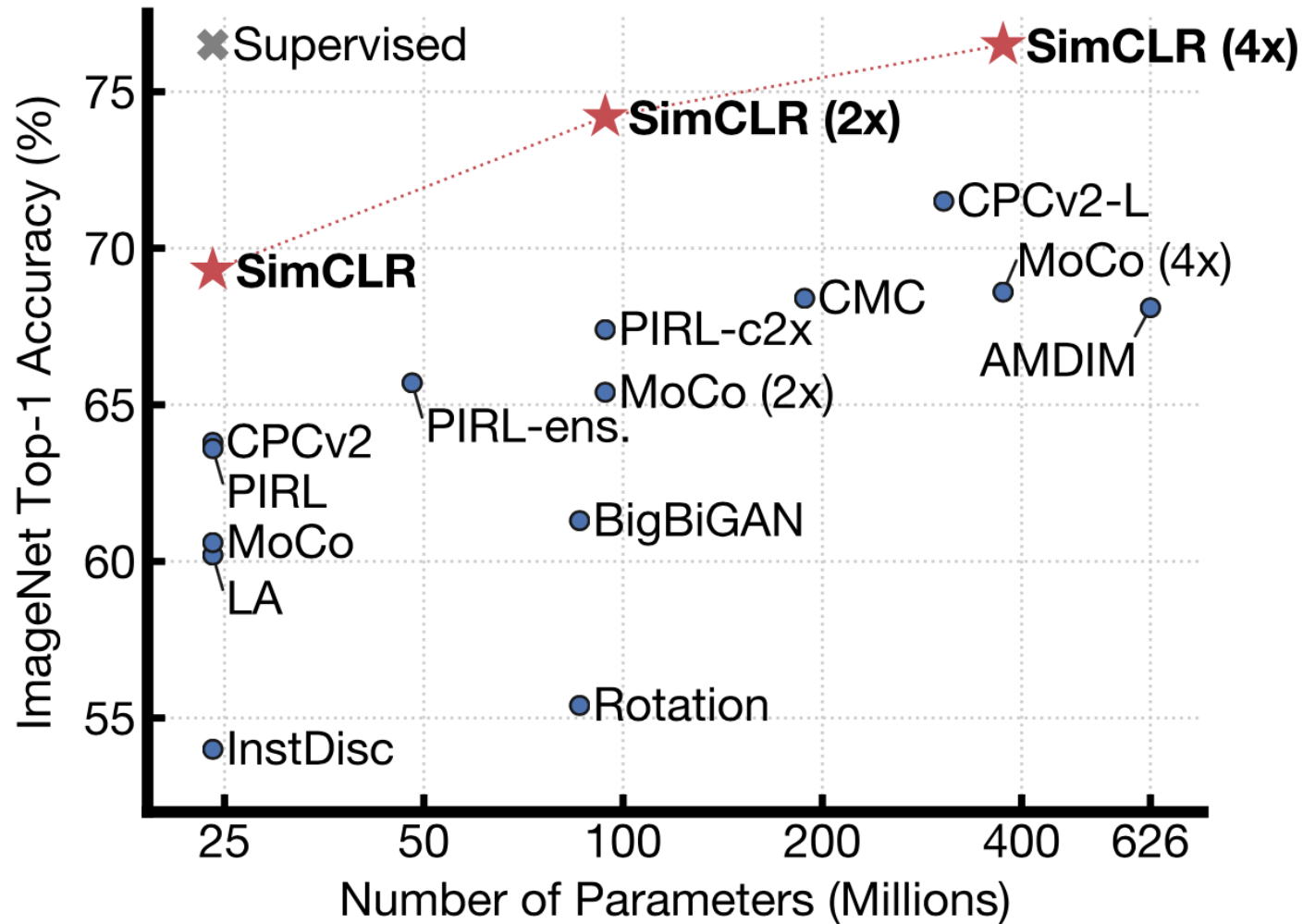
Department of Computer Science and Technology, Xiamen University

Image source: <https://ai.googleblog.com/2020/04/advancing-self-supervised-and-semi.html>

# SimCLR



# SimCLR





This 2-page short paper declares:

- Two design improvements used in SimCLR, namely, an MLP projection head and stronger data augmentation, are **orthogonal** to the frameworks of MoCo and SimCLR, and when used with MoCo they lead to better image classification and object detection transfer learning results.
- In contrast to SimCLR’s large 4k~8k batches, which require TPU support, our “MoCo v2” baselines can run on a typical 8-GPU machine and achieve better results than SimCLR.

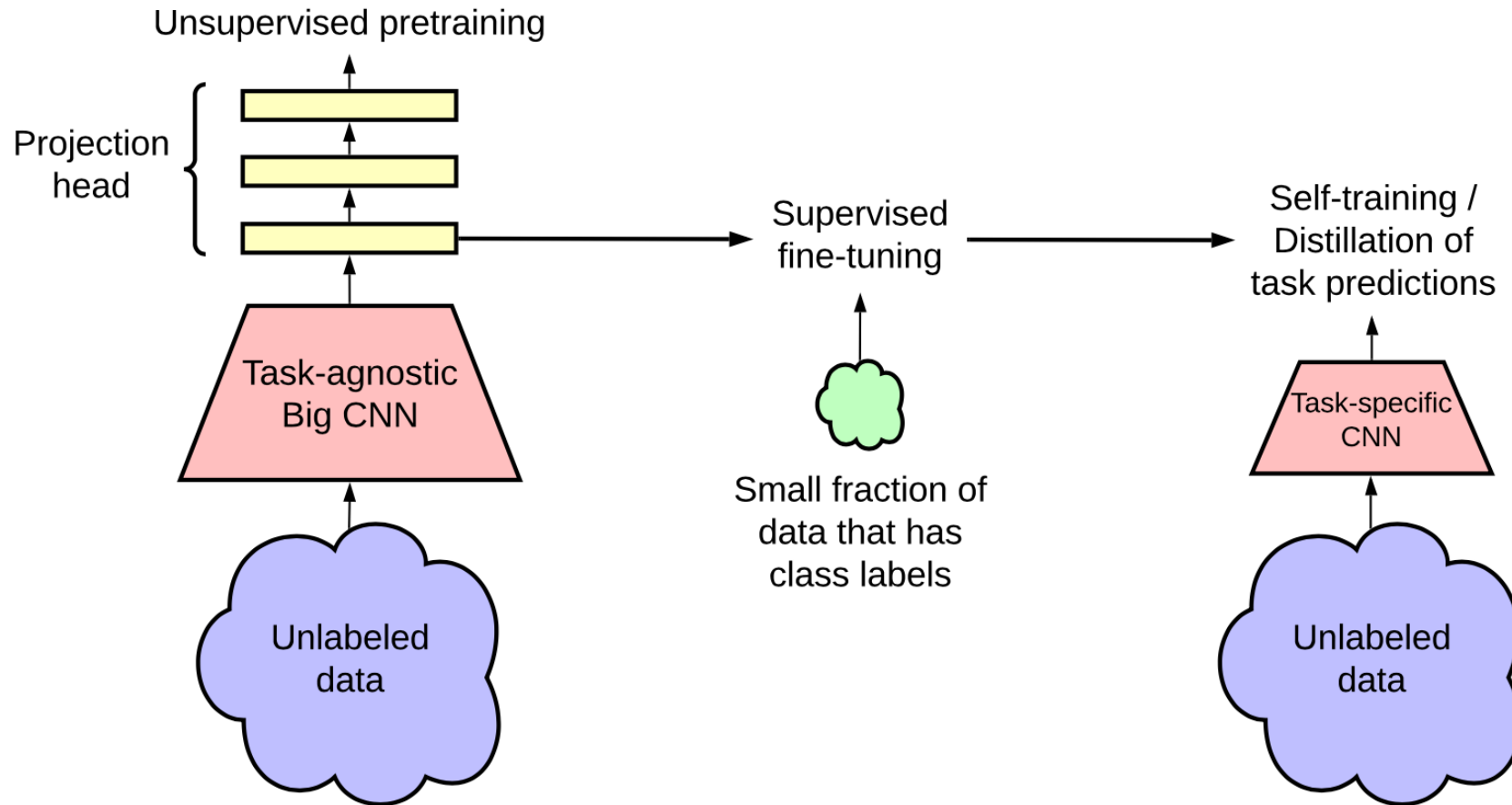
# SimCLR v2

Big self-supervised models are strong semi-supervised learners

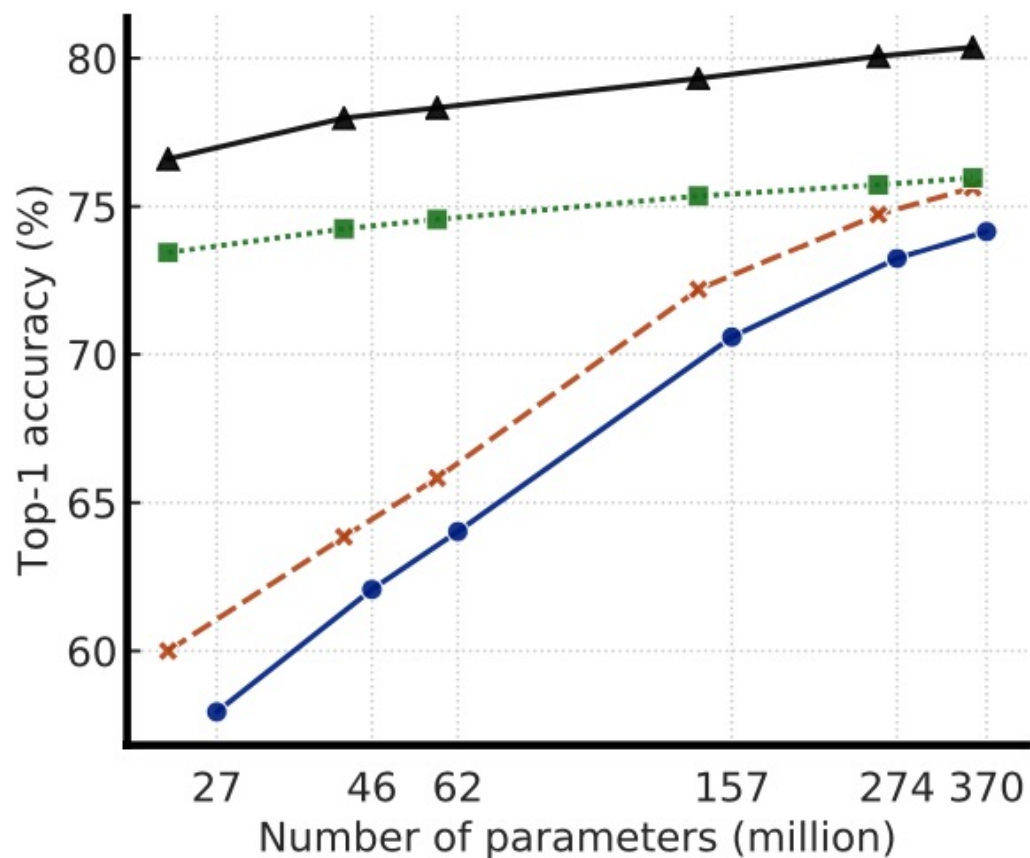
T. Chen, S. Kornblith, K. Swersky... - Advances in neural ..., 2020 - proceedings.neurips.cc

..., supervised fine-tune" paradigm for semi-supervised learning on ImageNet [21]. During self-supervised ... : Using a big (deep and wide) neural network for self-supervised pretraining and ...

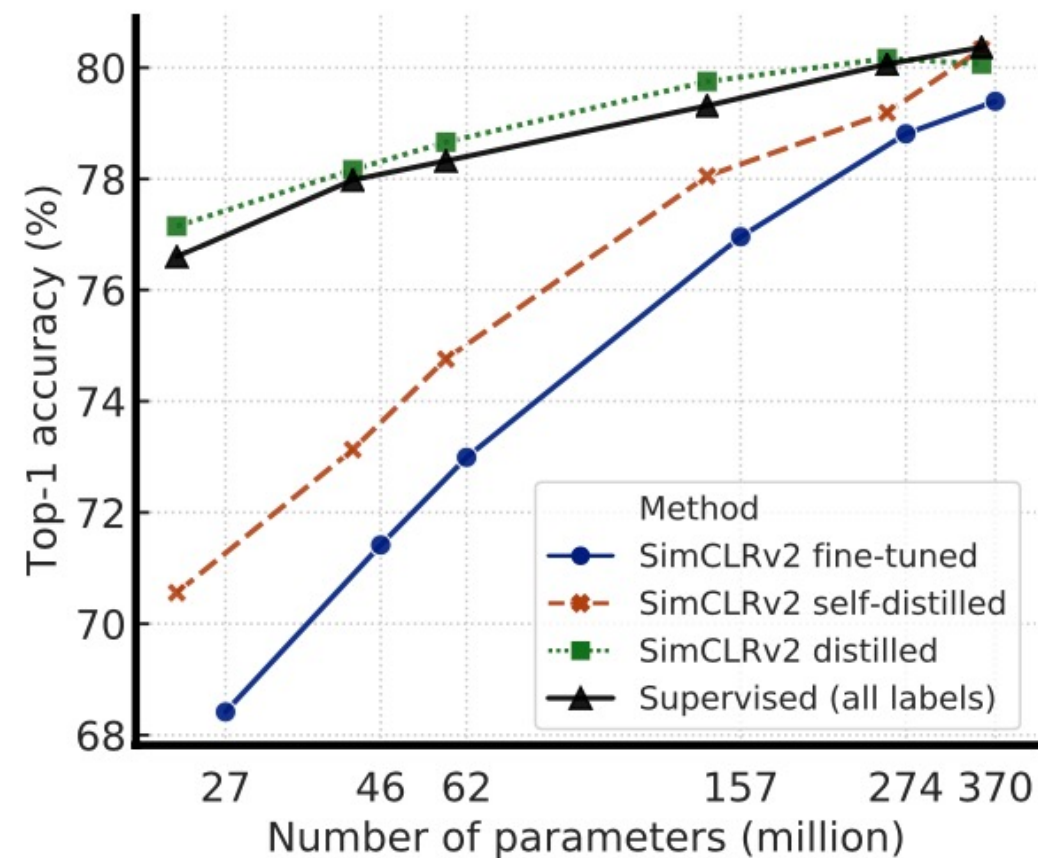
☆ Save 📄 Cite Cited by 1901 Related articles All 13 versions 🔗



# SimCLR v2



(a) Label fraction 1%

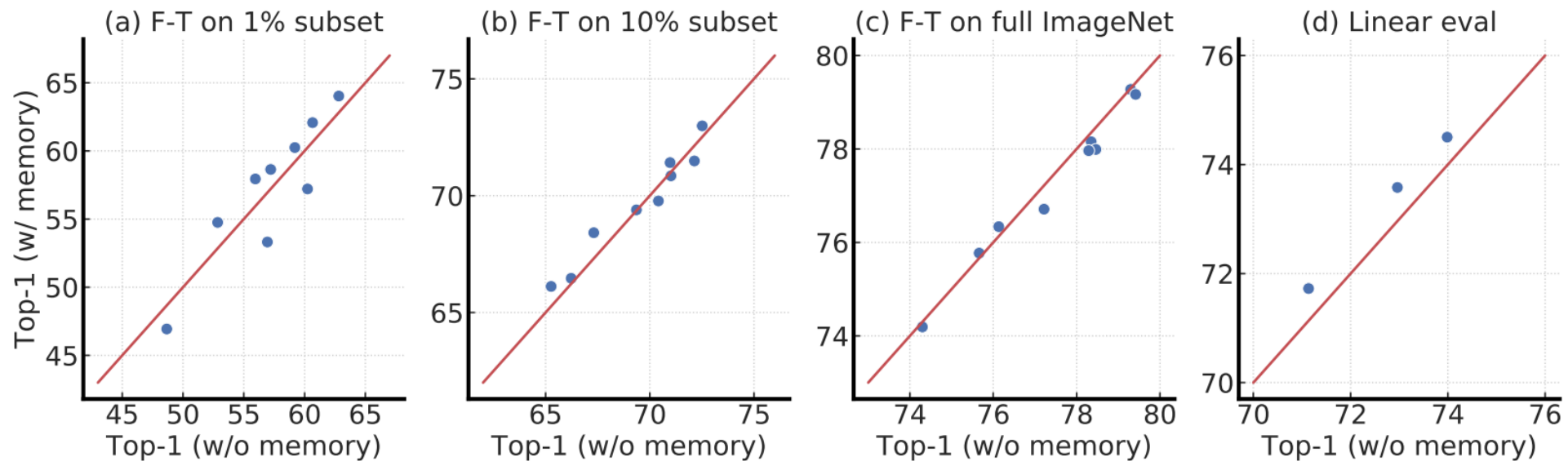


(b) Label fraction 10%



# SimCLR v2

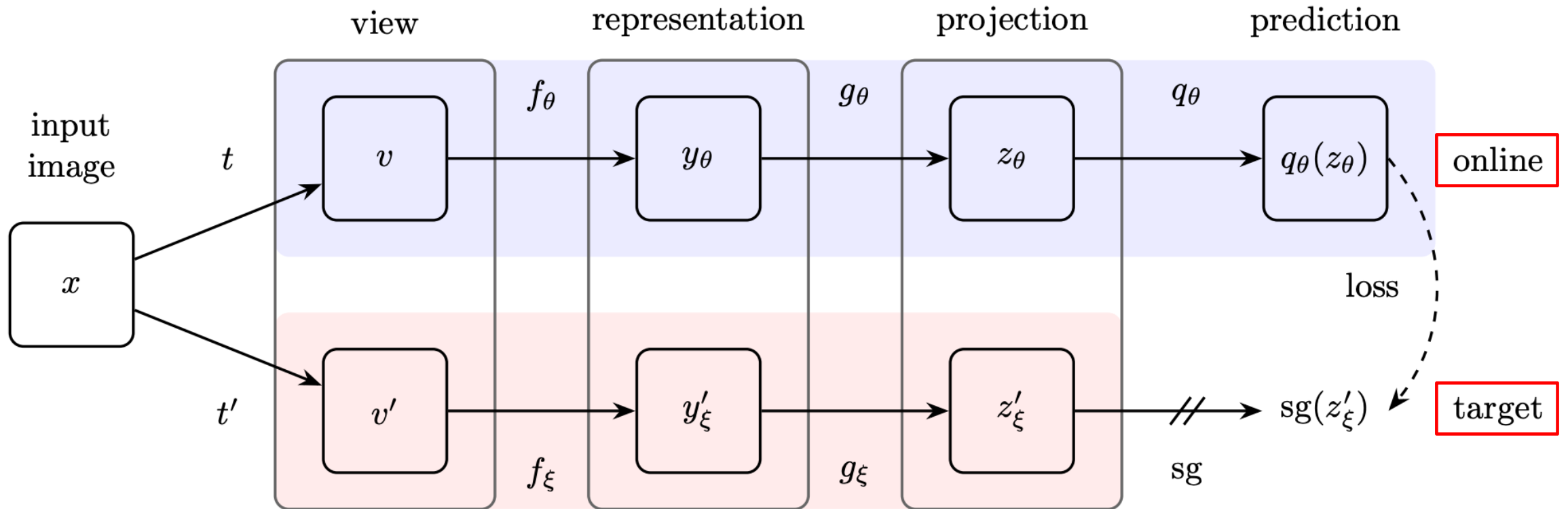
- “Memory provides **modest advantages** in terms of linear evaluation and fine-tuning with 1% of the labels; the improvement is around 1%.”
- “We believe the reason that memory only provides **marginal improvement** is that we already use a big batch size (i.e. 4096).”



Top-1 results of ResNet-50, ResNet-101, and ResNet-152 trained with or without memory.



## ■ Are negative samples necessary for contrastive learning?





# SimSiam

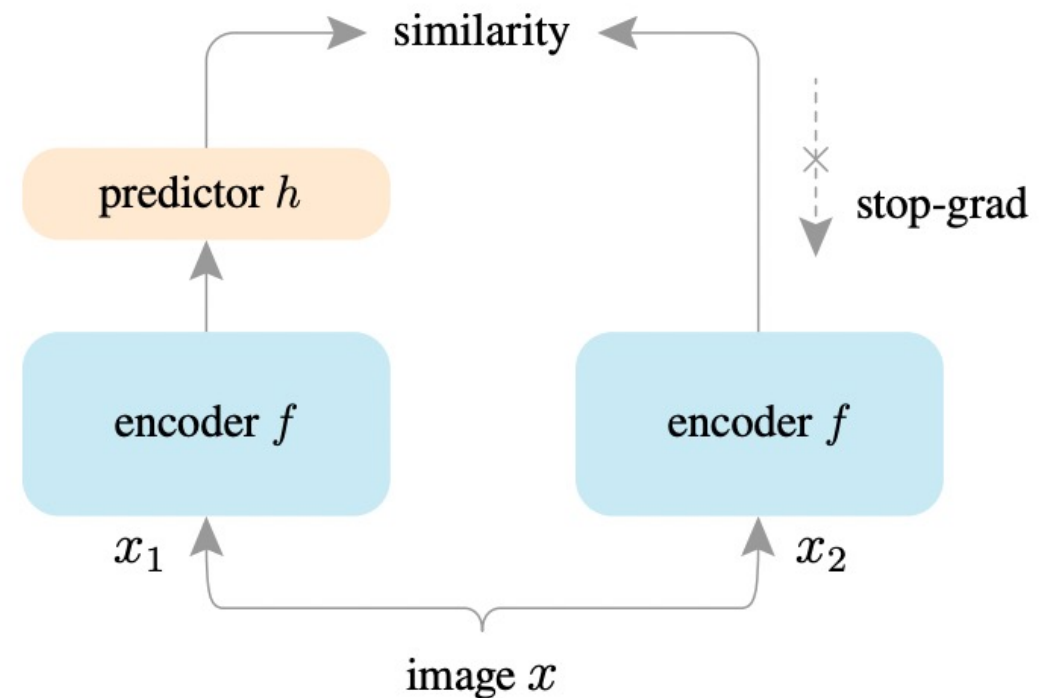
## Exploring simple siamese representation learning

X Chen, K He - Proceedings of the IEEE/CVF conference on ..., 2021 - openaccess.thecvf.com

... In this paper, we report surprising empirical results that **simple** Siamese networks can learn meaningful ... We hope this **simple** baseline will motivate people to rethink the roles of Siamese ...

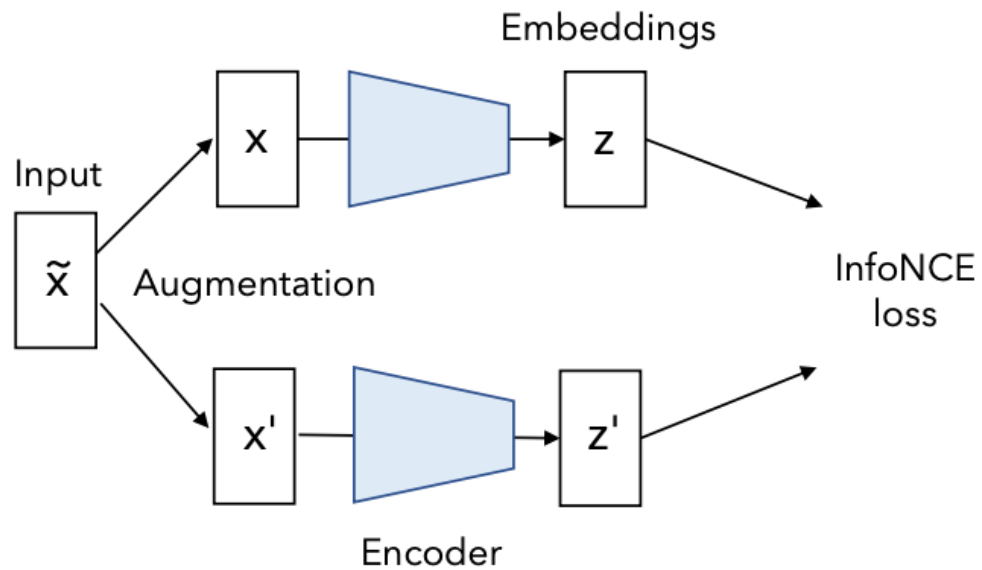
☆ Save 7 Cite Cited by 2986 Related articles All 6 versions ⇨

- Surprising empirical results that simple Siamese networks can learn meaningful representations, even using none of the following:
  - negative sample pairs,
  - large batches,
  - momentum encoders.

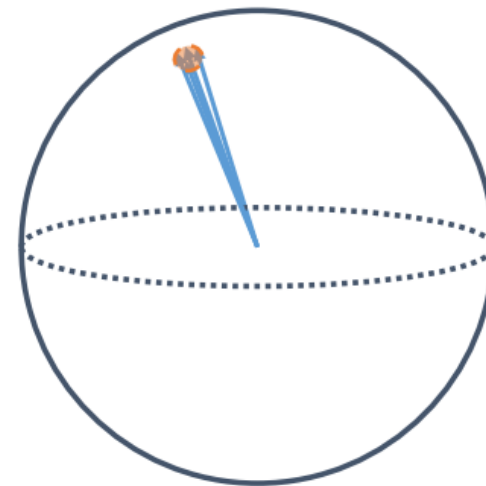


# Collapse

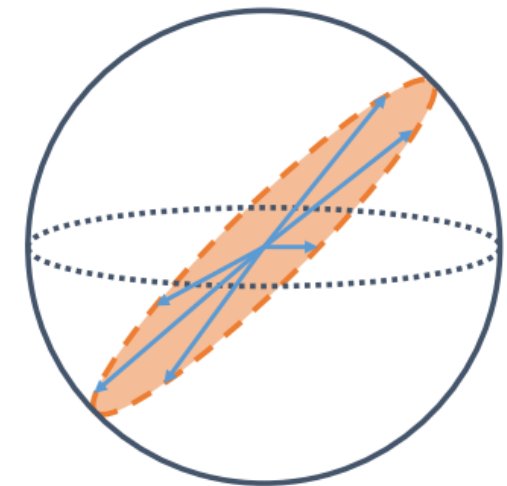
- Collapse: All representations tend to be the same.



(a) embedding space



(b) complete collapse



(c) dimensional collapse



# Downstream Tasks for Evaluation

- To compare different self-supervised learning methods, there are some commonly used downstream tasks for evaluation.

CV:

- Semantic segmentation
- Object detection
- Image classification
- Human action recognition
- ...

NLP:

- Question answering
- Named entity recognition
- Sentiment classification
- Natural language inference
- ...



# Conclusion

After this lecture, you should know:

- What is the difference between supervised and self-supervised learning.
- What is pretext task and pseudo label?
- How can we generate pseudo label?
- What is contrastive learning?

## Suggested Reading

- Self-supervised Visual Feature Learning with Deep Neural Networks: A Survey
- Self-supervised Learning: Generative or Contrastive
- Awesome Self-Supervised Learning
- Contrastive Self-Supervised Learning
- 对比学习(Contrastive Learning)相关进展梳理





# Thank you!

- Any question?
- Don't hesitate to send email to me for asking questions and discussion. 😊