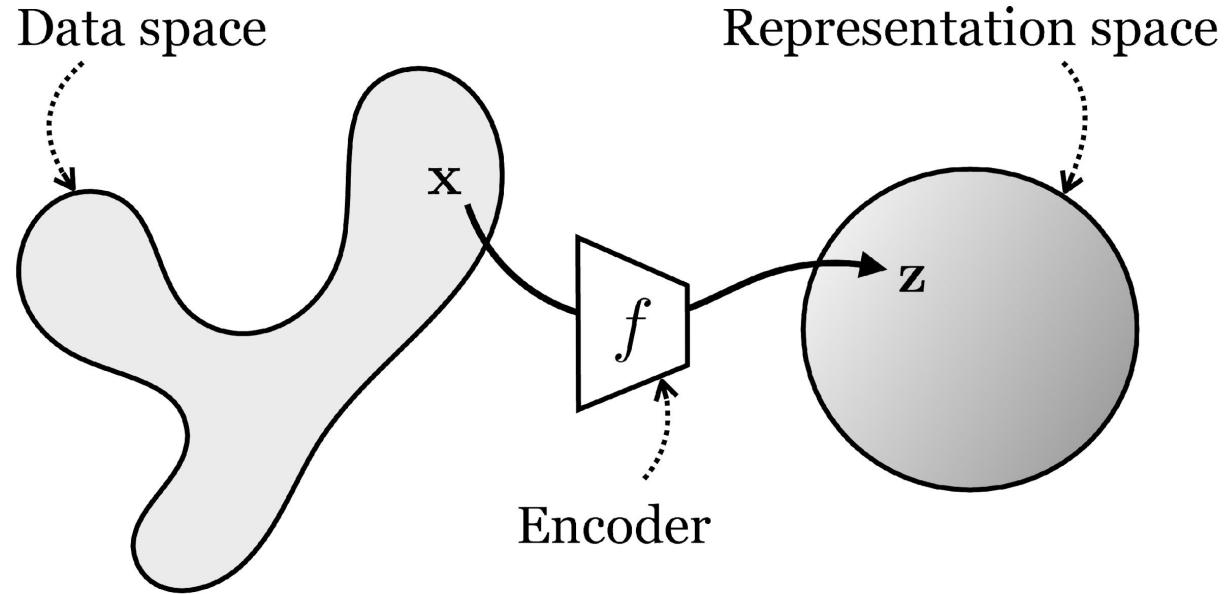


# A Unifying Framework for Representation Learning

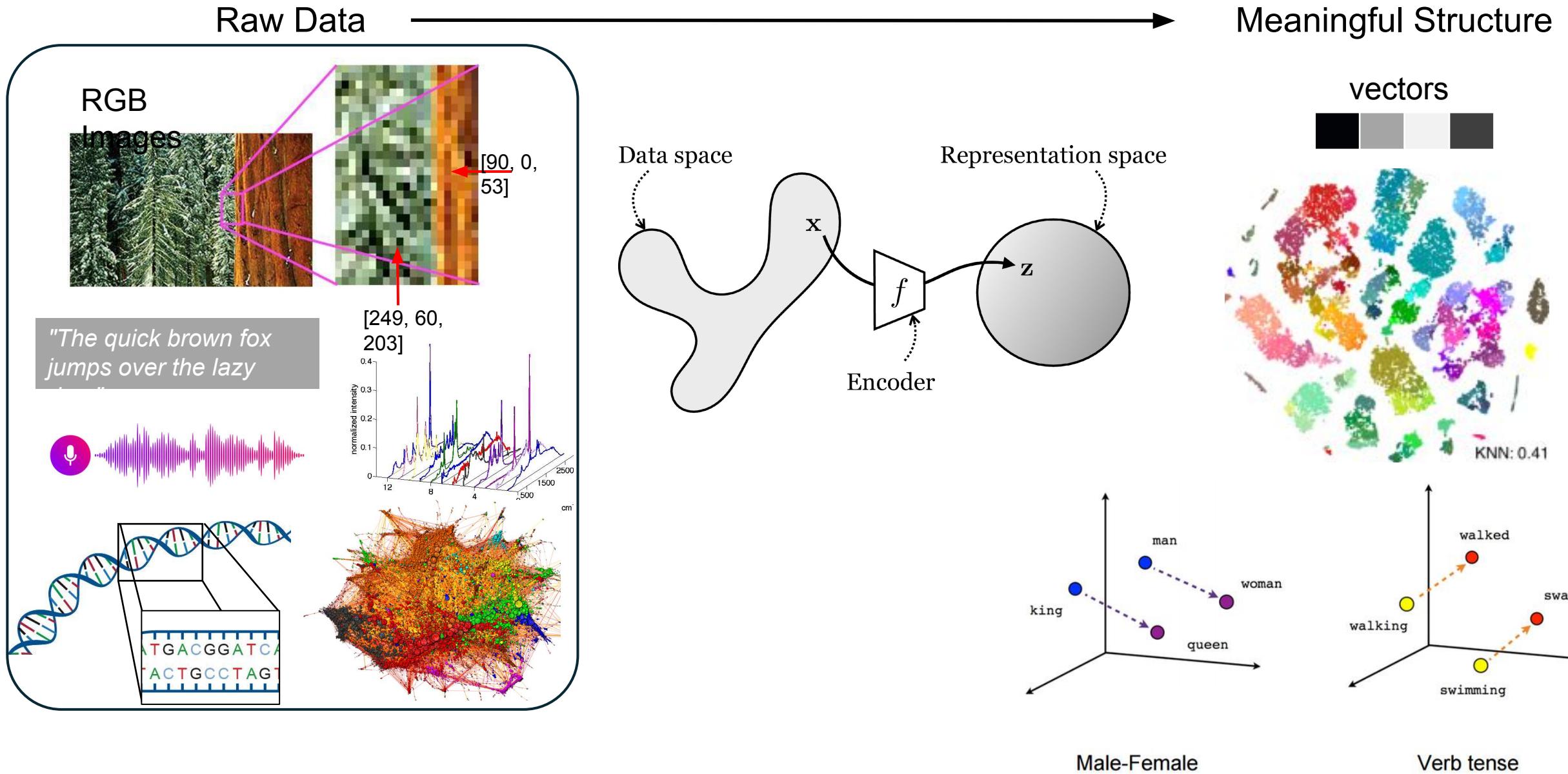
Shaden Alshammary

# What's Representation Learning?

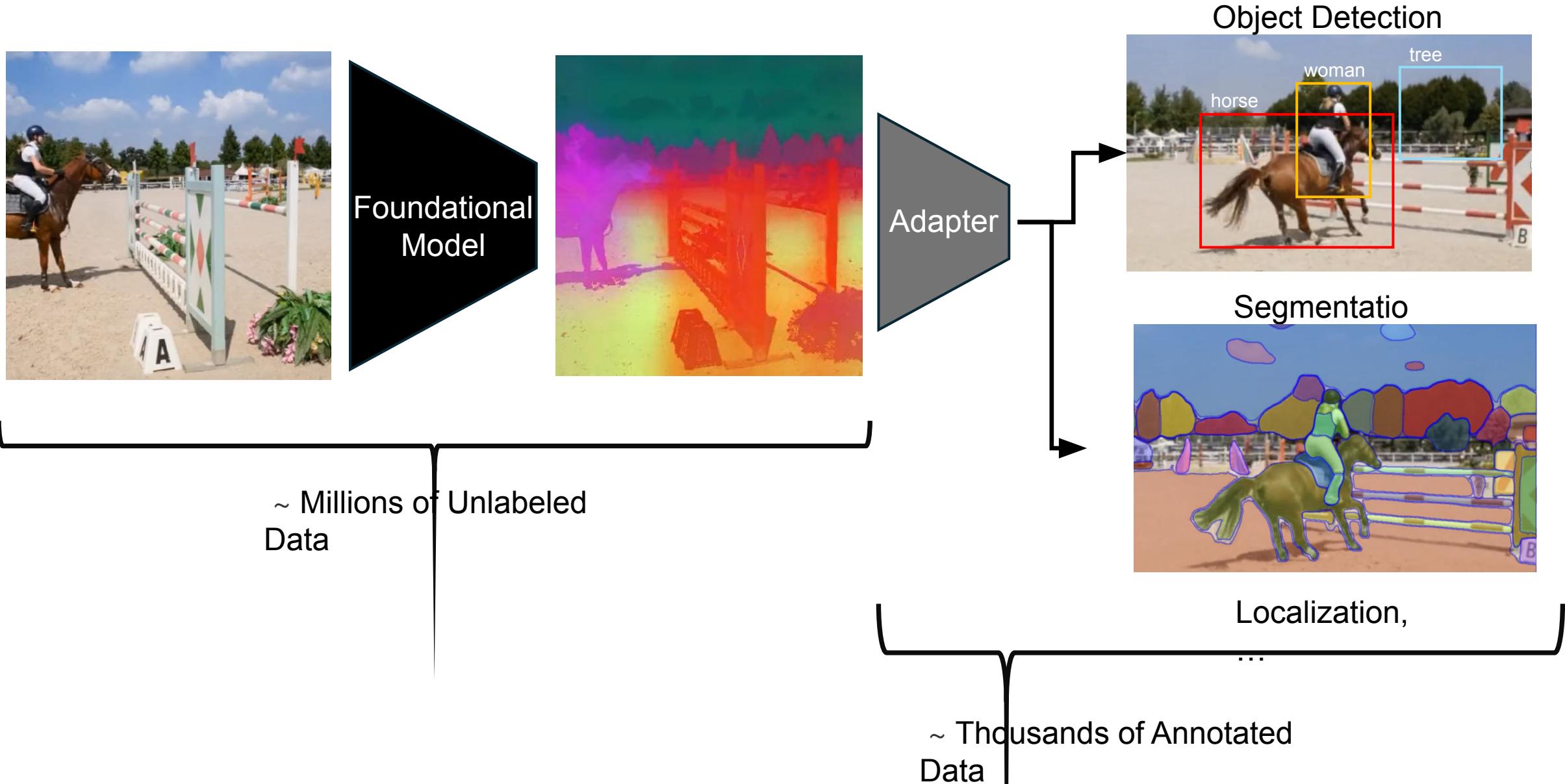


from raw, messy, high dimensional data to meaningful structure

# What's Representation Learning?



# Why Learning Good Representations Matters?



# How Much Information is the Machine Given during Learning?

## ► “Pure” Reinforcement Learning (**cherry**)

- The machine predicts a scalar reward given once in a while.

## ► **A few bits for some samples**



## ► Supervised Learning (**icing**)

- The machine predicts a category or a few numbers for each input
- Predicting human-supplied data
- **10→10,000 bits per sample**

## ► Self-Supervised Learning (**cake génoise**)

- The machine predicts any part of its input for any observed part.
- Predicts future frames in videos
- **Millions of bits per sample**

[Slide Credit: Yann LeCun]

# Supervised Learning

Classification

Object Detection,  
Segmentation etc

# Unsupervised Learning

Dimensionality  
Reduction

PCA, t-SNE

Clustering

K-Means

# Self-Supervised Learning (SSL)

Compressive

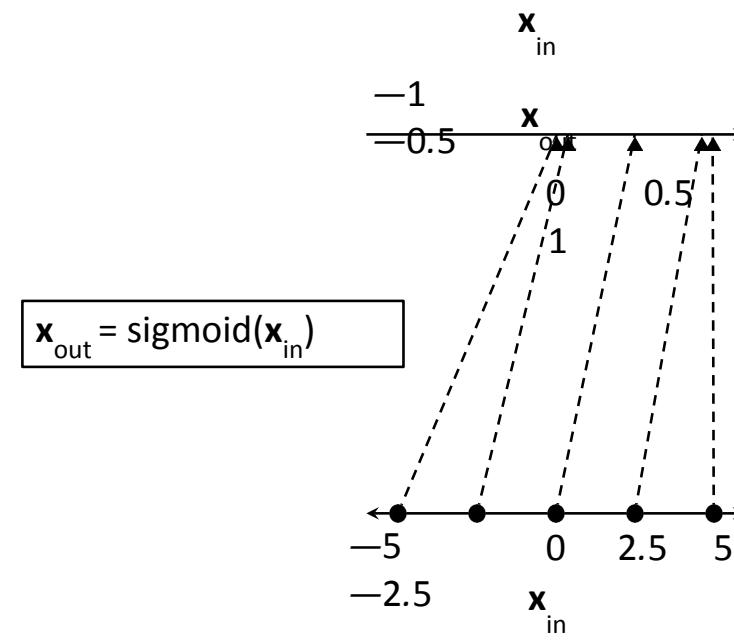
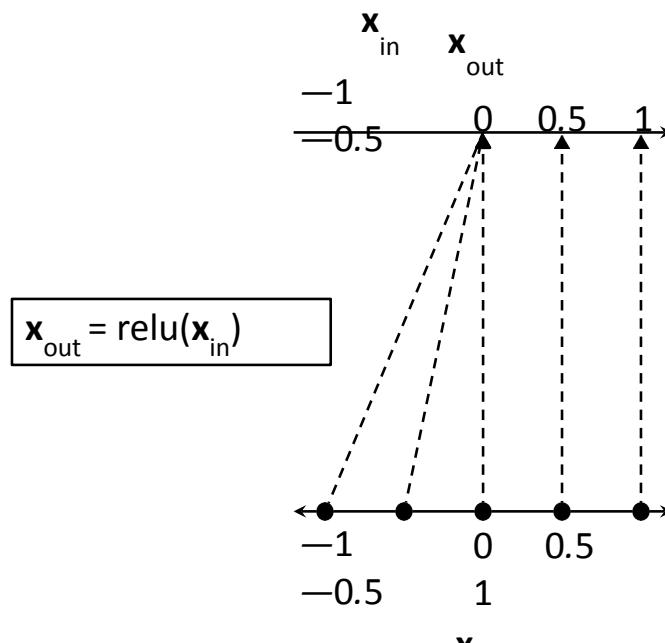
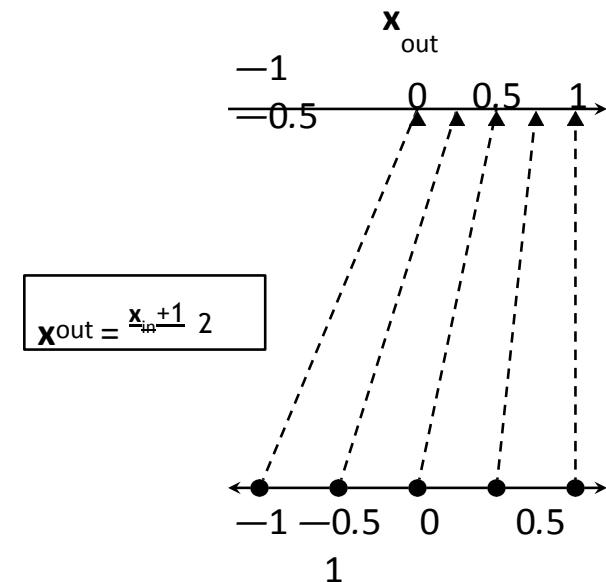
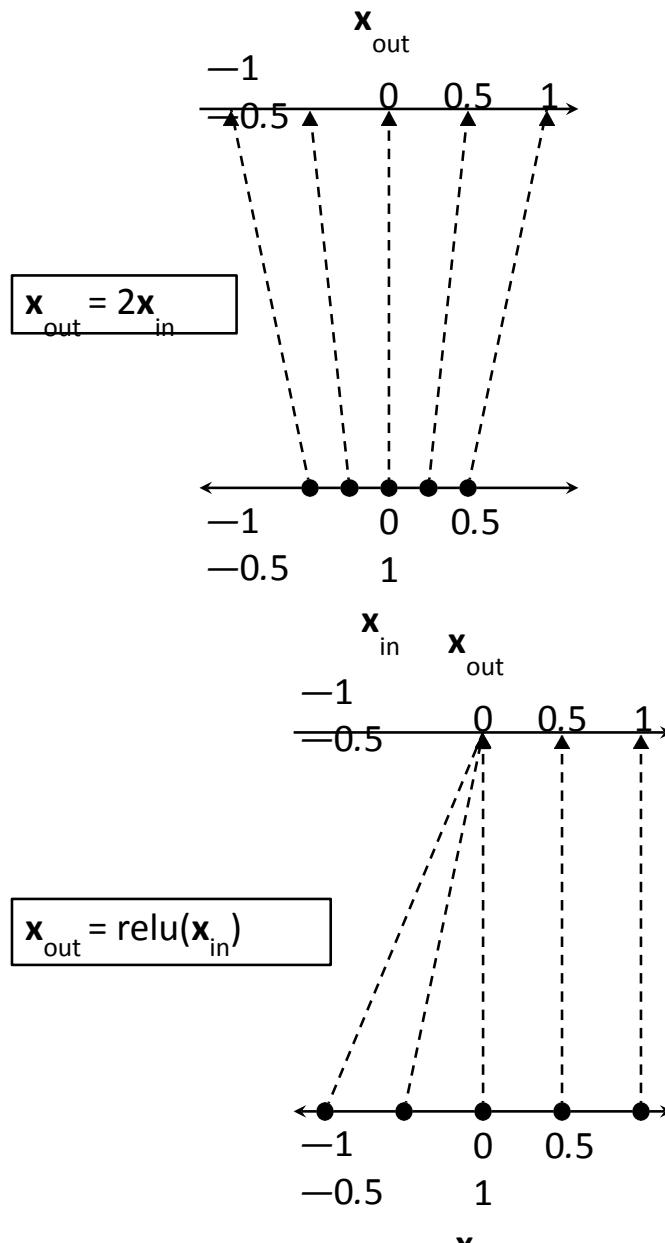
Predictive

Contrastive

# What Do Deep Networks Learn?

What Do ~~Deep~~ Networks Learn?  
**Shallow**

# Data transformations for a variety of neural net layers



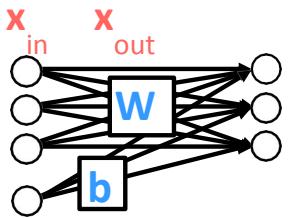
Wiring graph

Equation

Activations

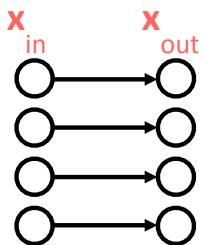
Parameters

linear



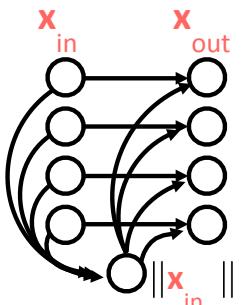
$$\mathbf{x}_{\text{out}} = \mathbf{W}\mathbf{x}_{\text{in}} + \mathbf{b}$$

relu



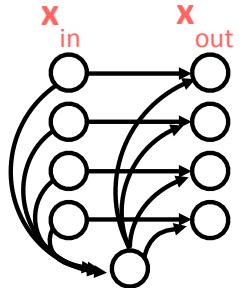
$$x_{\text{out}}[i] = \max(x_{\text{in}}[i], 0)$$

L2-norm



$$x_{\text{out}}[i] = \frac{x_{\text{in}}[i]}{\|\mathbf{x}_{\text{in}}\|_2}$$

softmax



$$x_{\text{out}}[i] = \frac{e^{-\tau x_{\text{in}}[i]}}{\sum_{k=1}^K e^{-\tau x_{\text{in}}[k]}}$$

Wiring graph

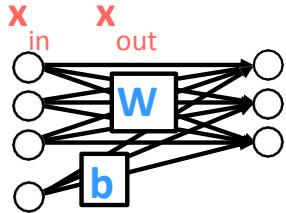
Equation

Mapping

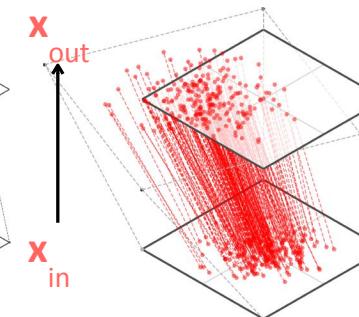
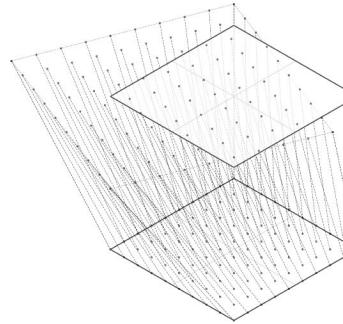
Activations

Parameters

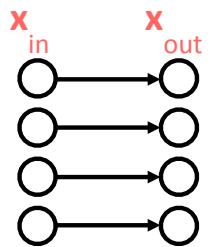
linear



$$x_{out} = Wx_{in} + b$$

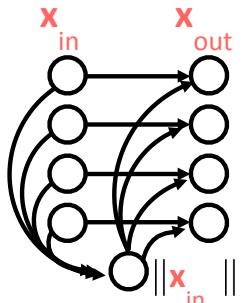


relu



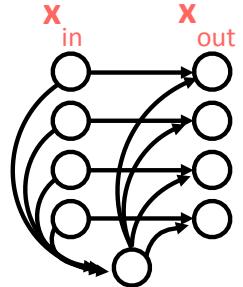
$$x_{out}[i] = \max(x_{in}[i], 0)$$

L2-norm



$$x_{out}[i] = \frac{x_{in}[i]}{\| x_{in} \|_2}$$

softmax



$$x_{out}[i] = \frac{e^{-\tau x_{in}[i]}}{\sum_{k=1}^K e^{-\tau x_{in}[k]}}$$

Wiring graph

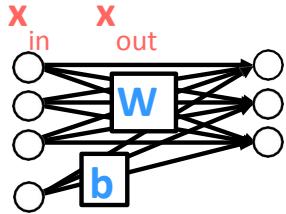
Equation

Mapping

Activations

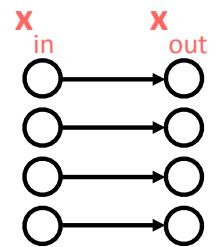
Parameters

linear

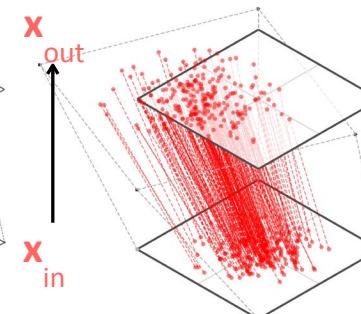
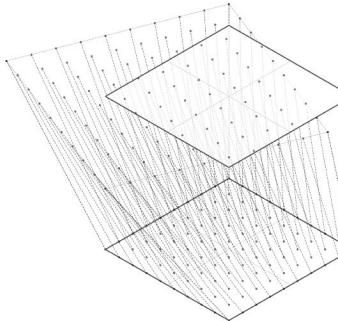


$$\mathbf{x}_{\text{out}} = \mathbf{W}\mathbf{x}_{\text{in}} + \mathbf{b}$$

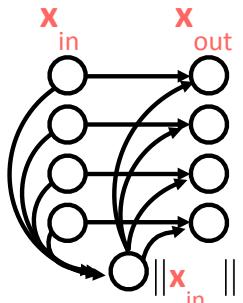
relu



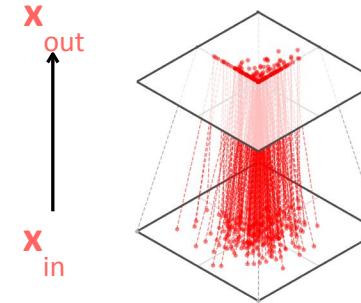
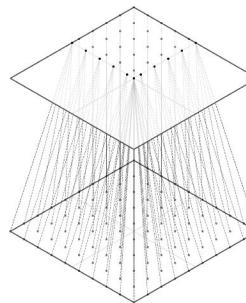
$$x_{\text{out}}[i] = \max(x_{\text{in}}[i], 0)$$



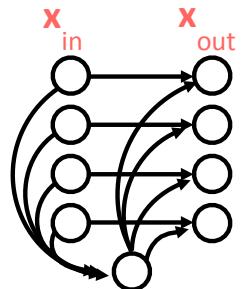
L2-norm



$$x_{\text{out}}[i] = \frac{x_{\text{in}}[i]}{\|x_{\text{in}}\|_2}$$



softmax



$$x_{\text{out}}[i] = \frac{e^{-\tau x_{\text{in}}[i]}}{\sum_{k=1}^K e^{-\tau x_{\text{in}}[k]}}$$

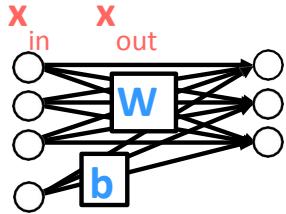
Wiring graph

Equation

Mapping

Activations

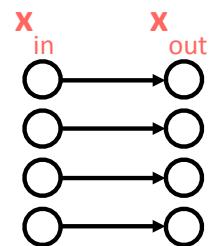
linear



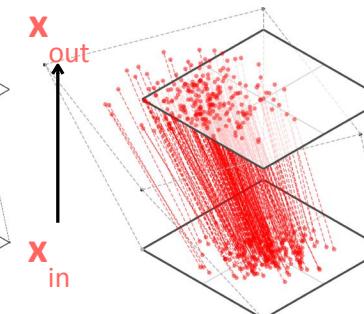
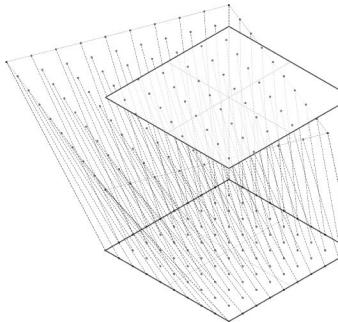
$$x_{\text{out}} = Wx_{\text{in}} + b$$

Parameters

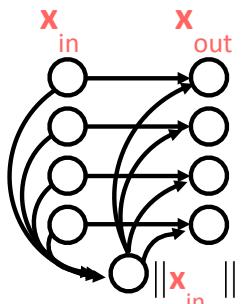
relu



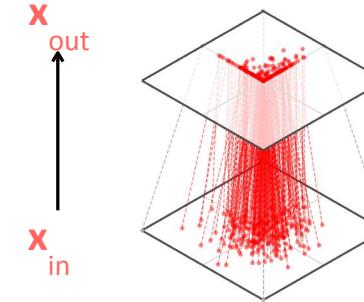
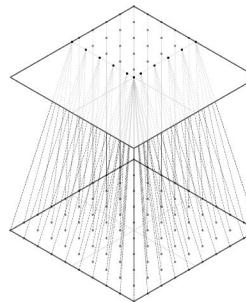
$$x_{\text{out}}[i] = \max(x_{\text{in}}[i], 0)$$



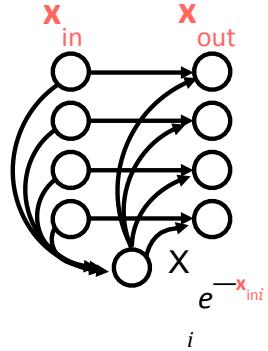
L2-norm



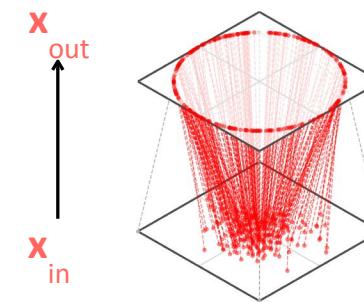
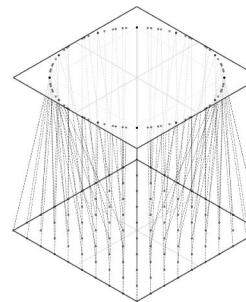
$$x_{\text{out}}[i] = \frac{x_{\text{in}}[i]}{\|x_{\text{in}}\|_2}$$



softmax



$$x_{\text{out}}[i] = \frac{e^{-\tau x_{\text{in}}[i]}}{\sum_{k=1}^K e^{-\tau x_{\text{in}}[k]}}$$



## Wiring graph

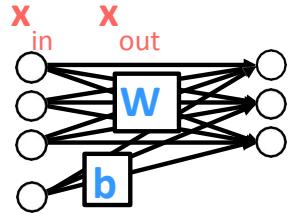
## Equation

## Mapping

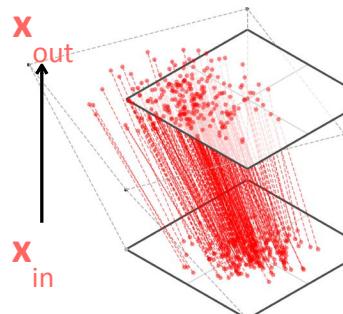
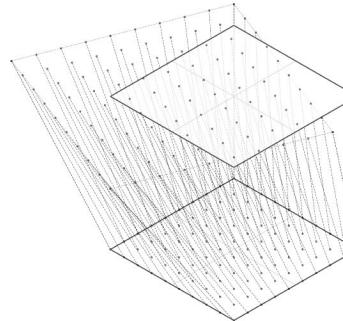
Activations

Parameters

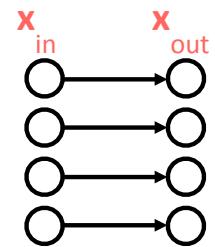
linear



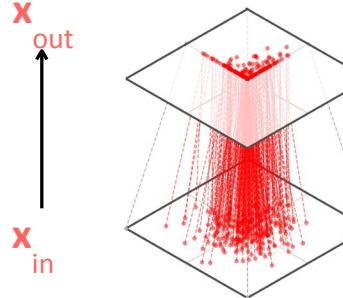
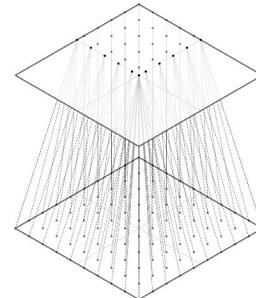
$$x_{out} = Wx_{in} + b$$



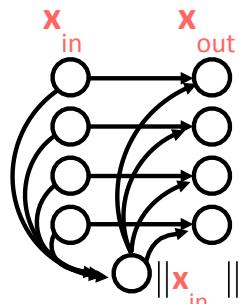
relu



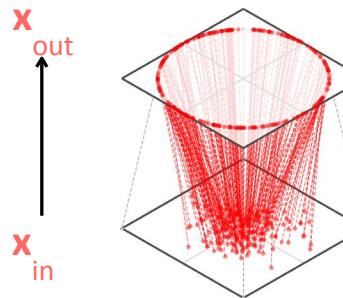
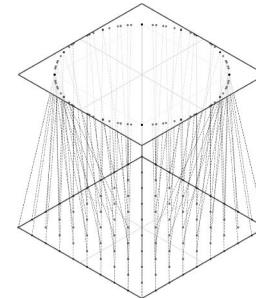
$$x_{out}[i] = \max(x_{in}[i], 0)$$



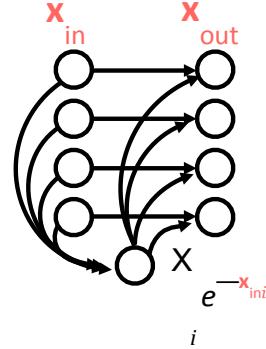
L2-norm



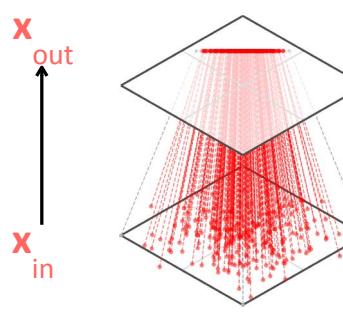
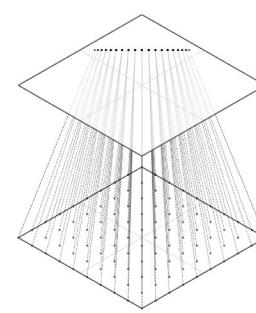
$$x_{out}[i] = \frac{x_{in}[i]}{\|x_{in}\|_2}$$



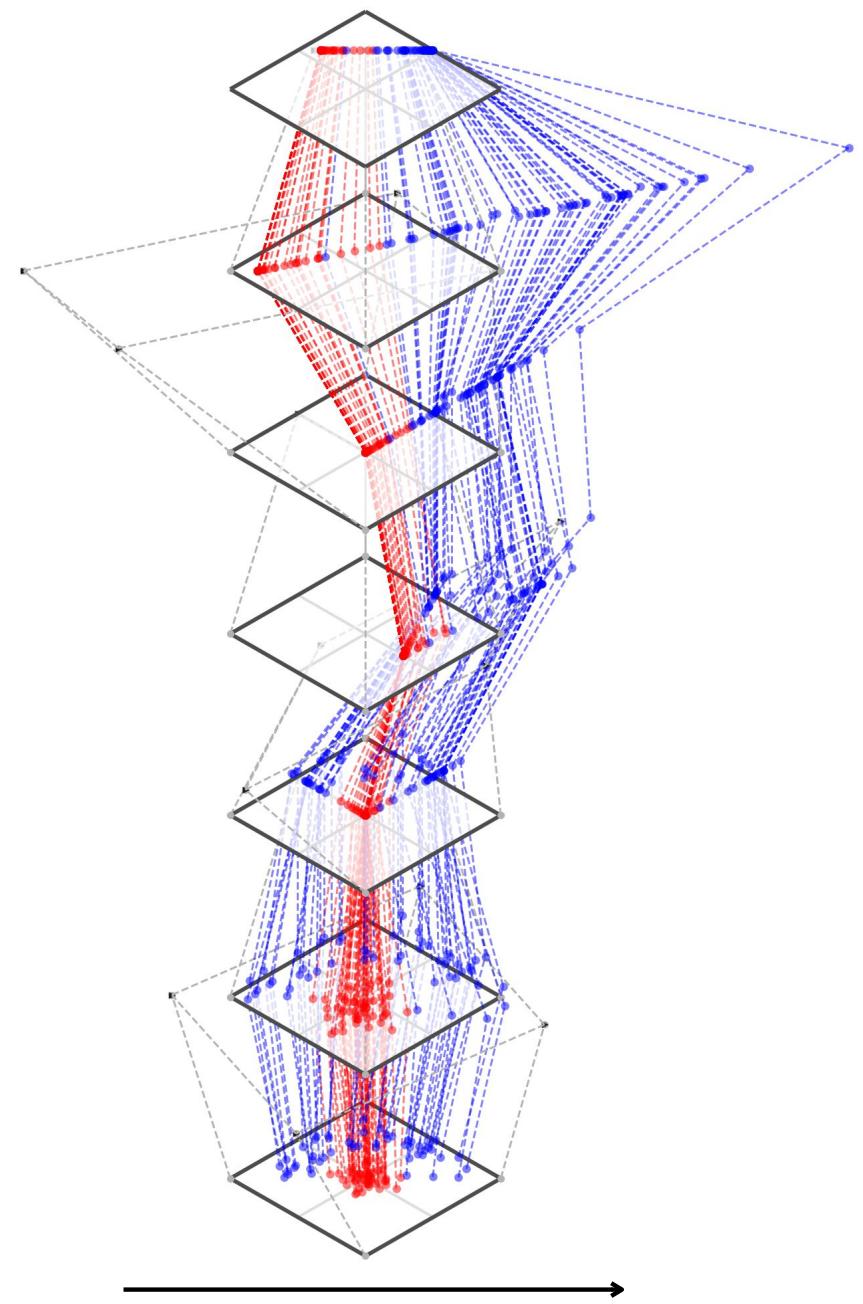
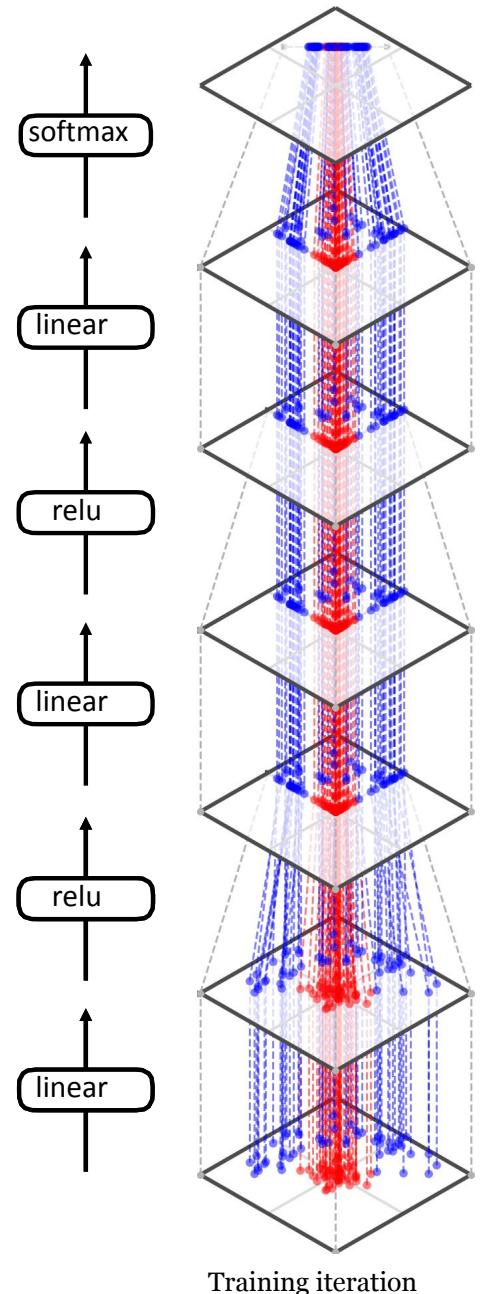
softmax



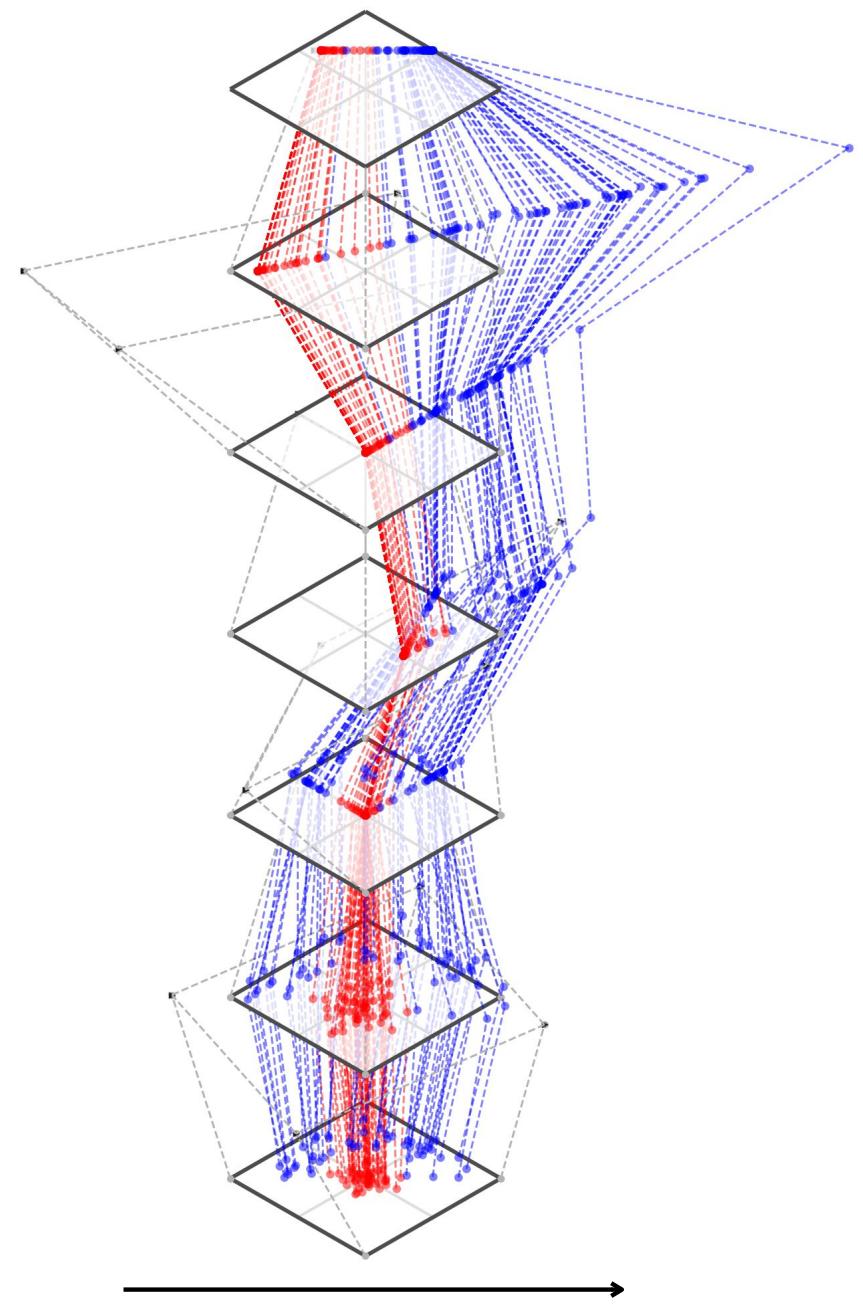
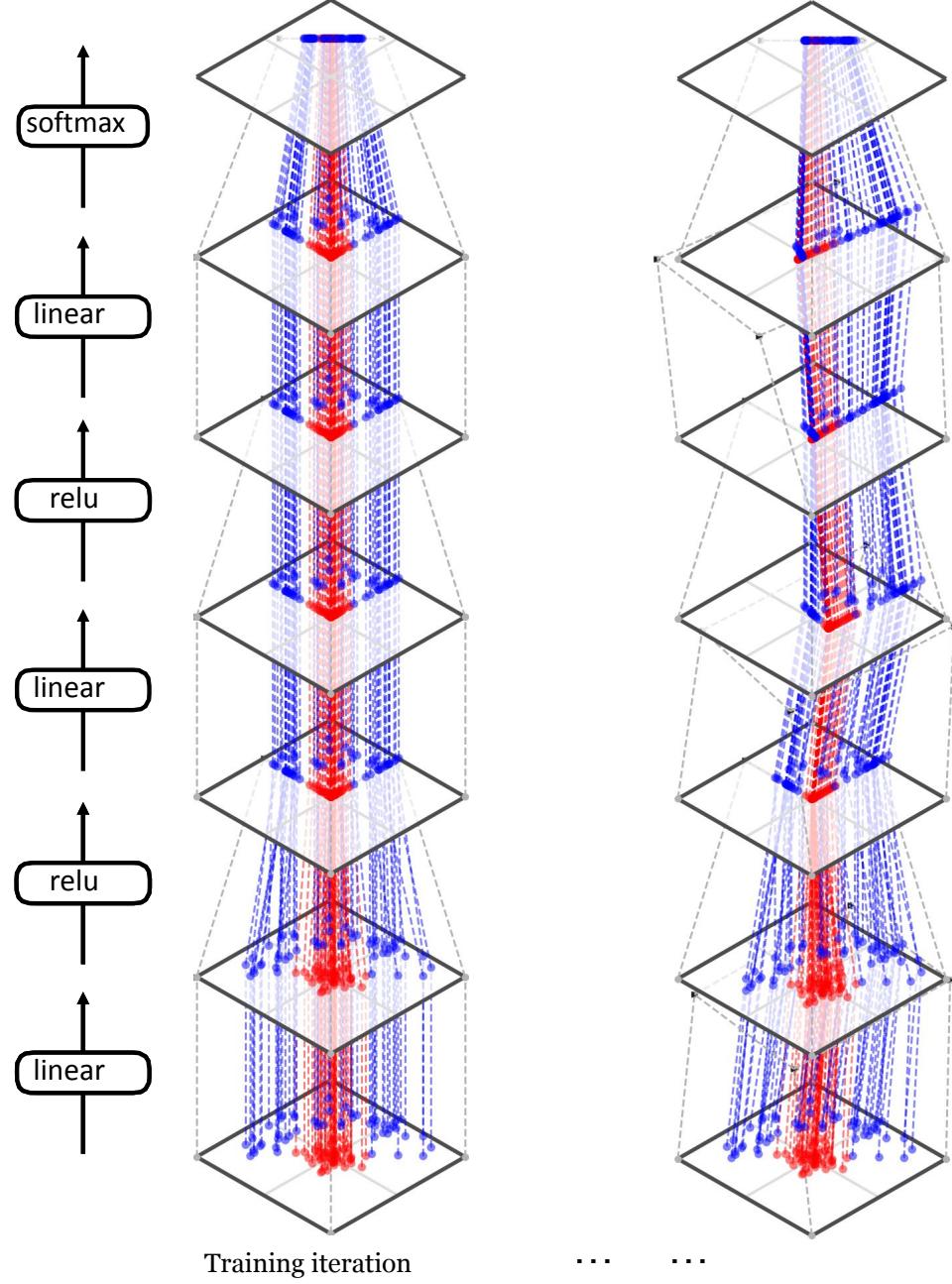
$$x_{out}[i] = \frac{e^{-\tau x_{in}[i]}}{\sum_{k=1}^K e^{-\tau x_{in}[k]}}$$



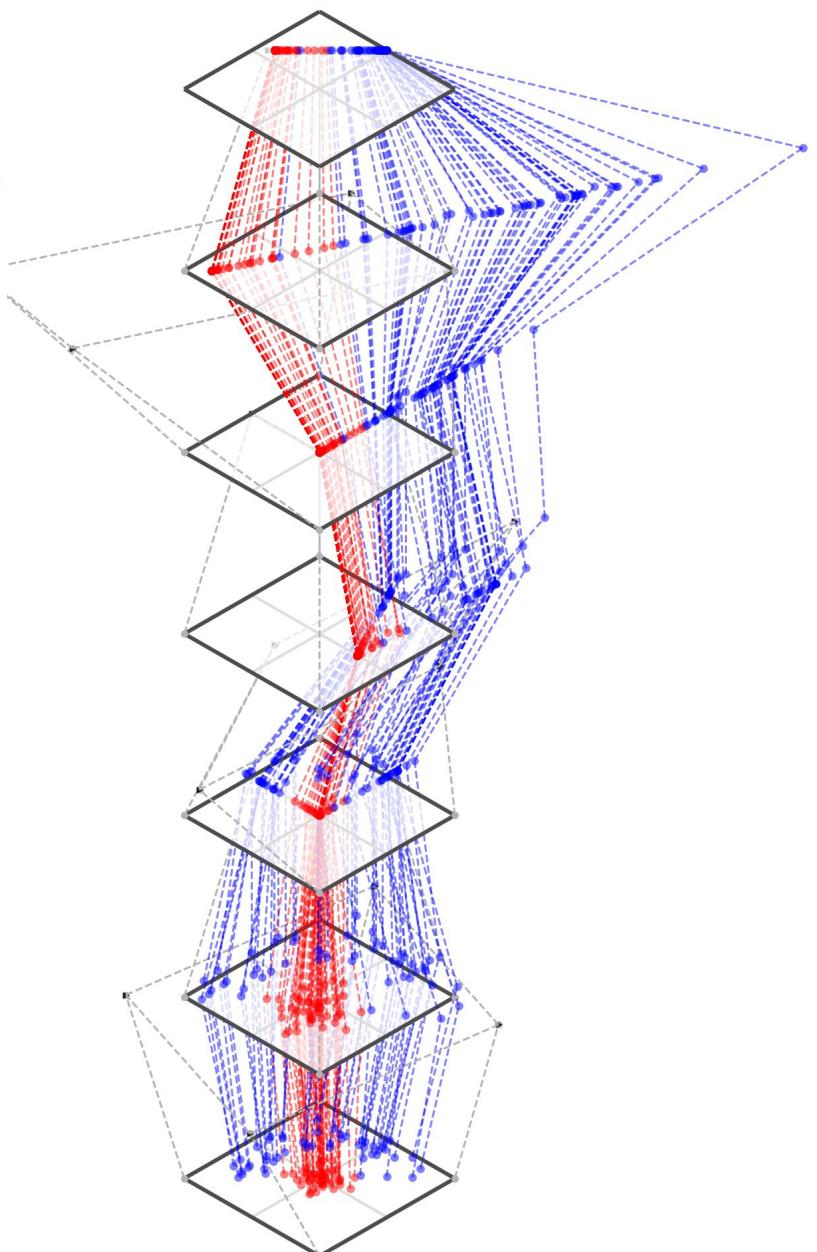
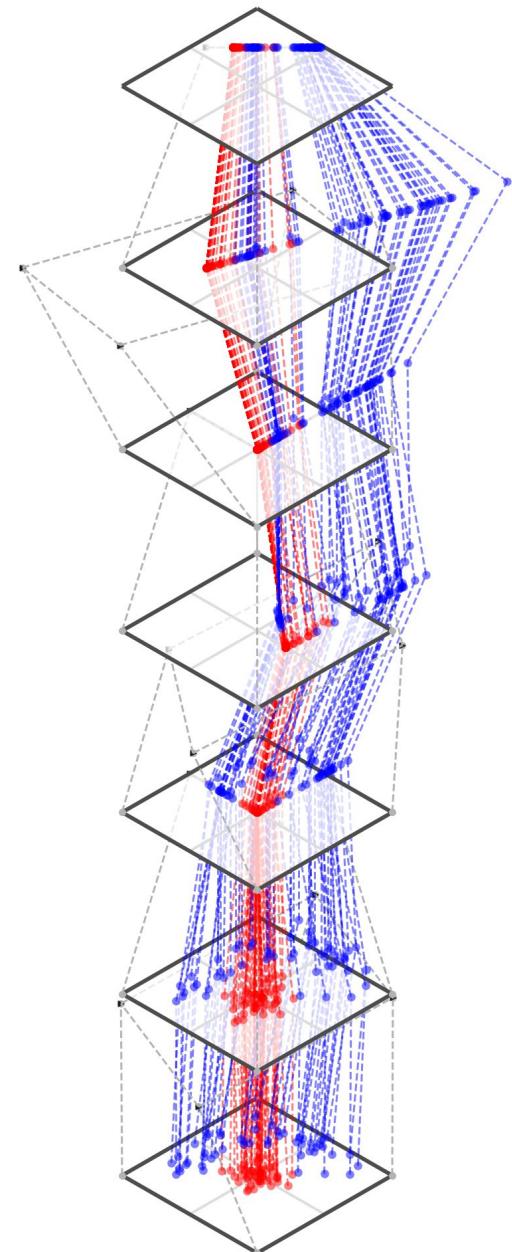
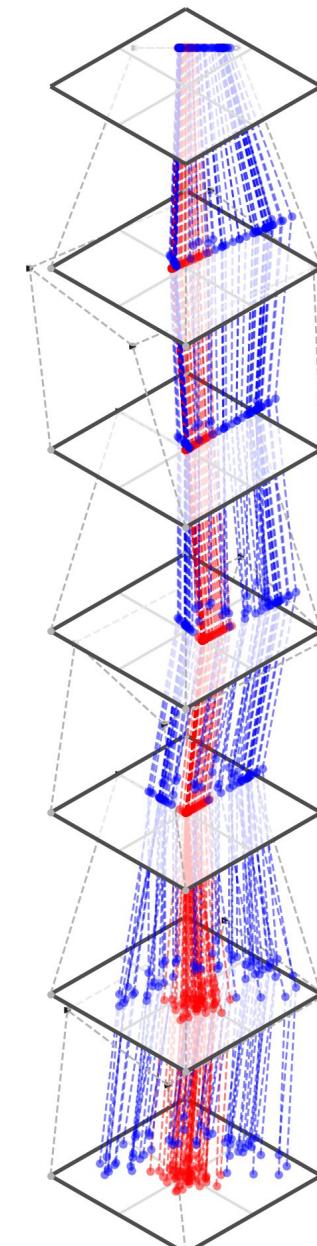
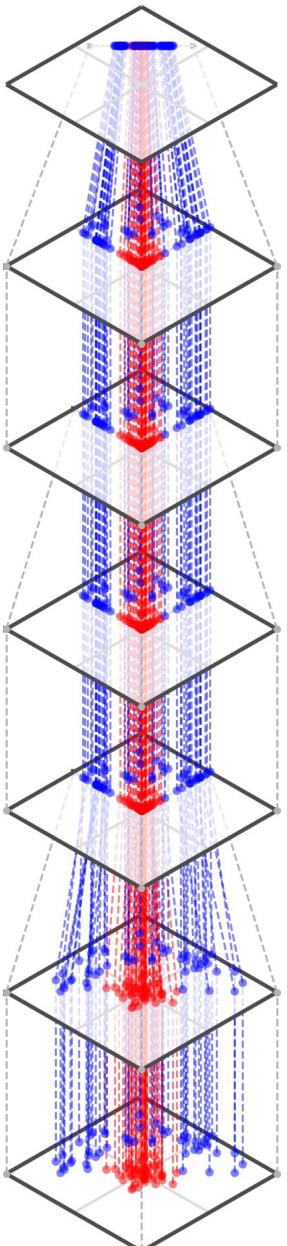
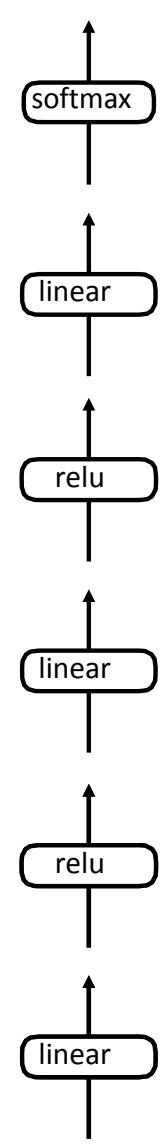
# MLP



# MLP



# MLP

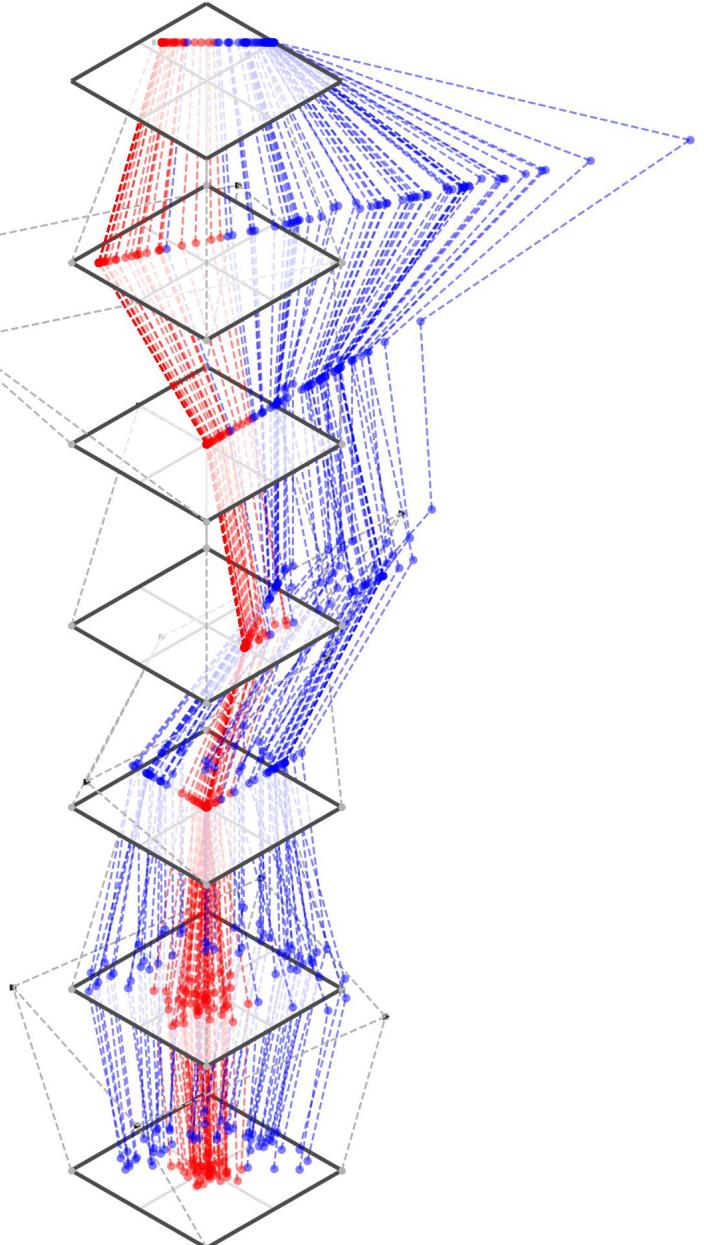
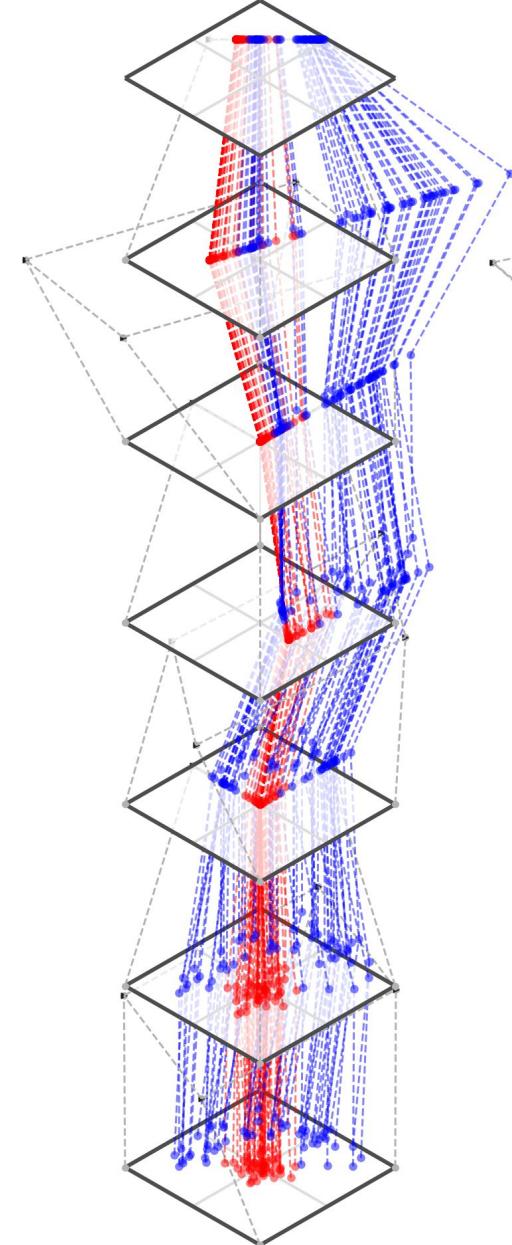
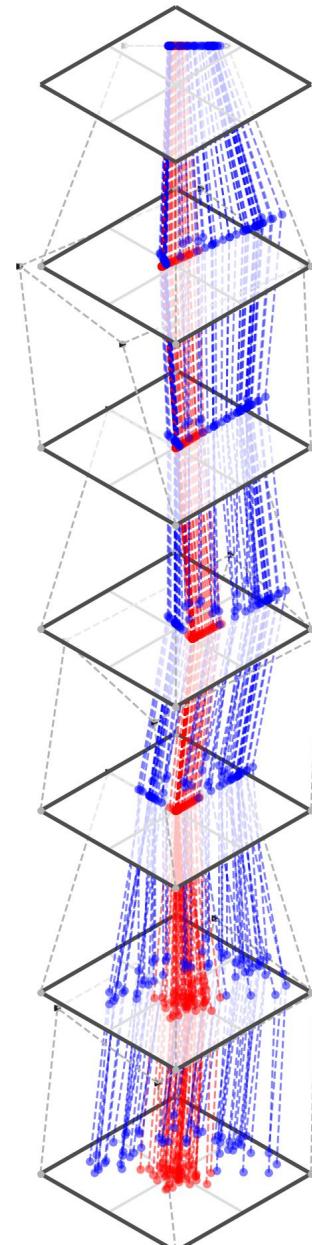
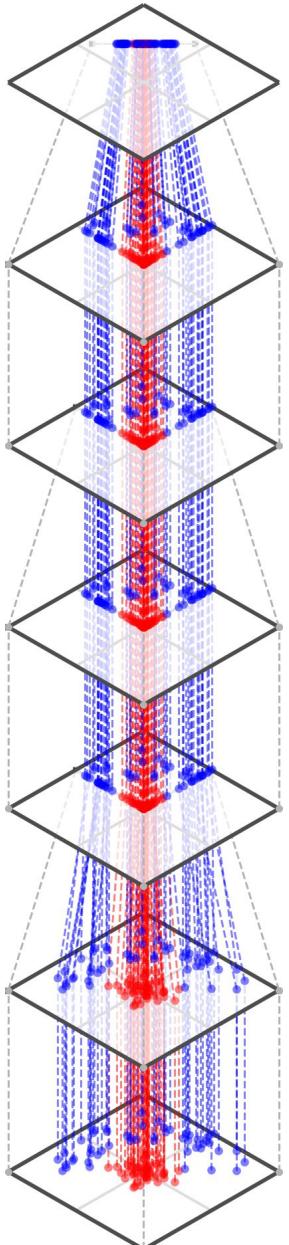
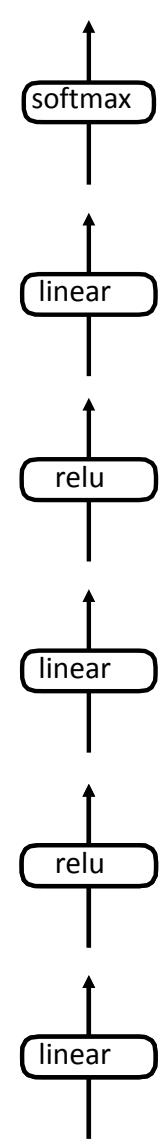


Training iteration

... ...

→

# MLP

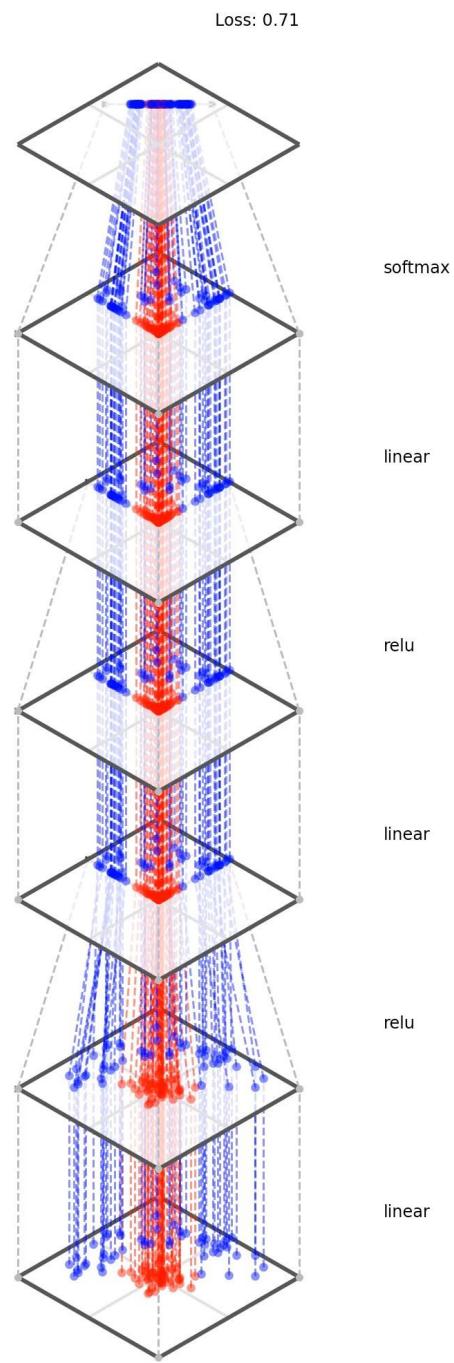


Training iteration

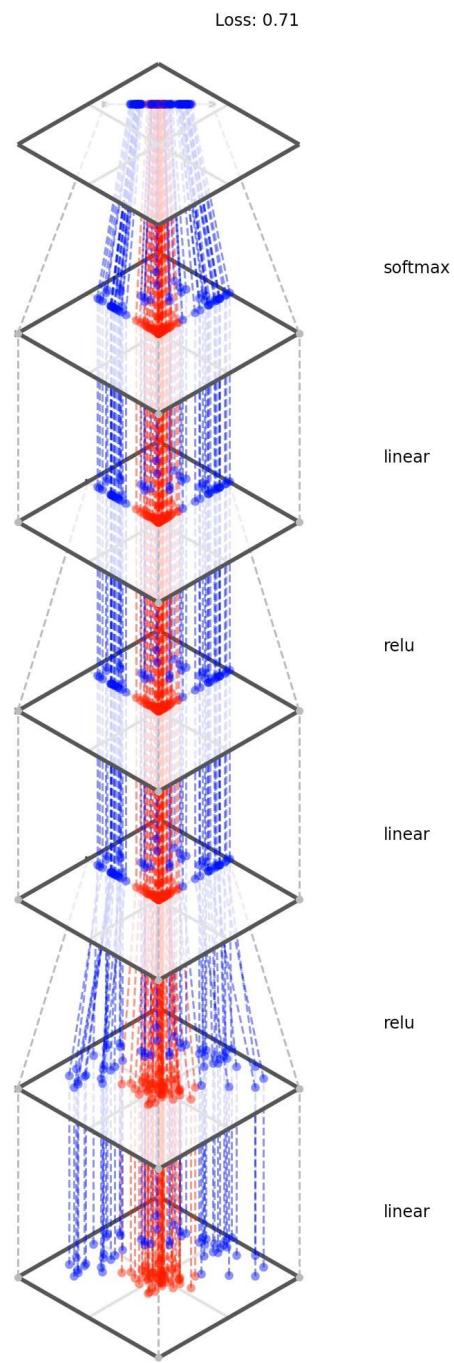
... ...

→

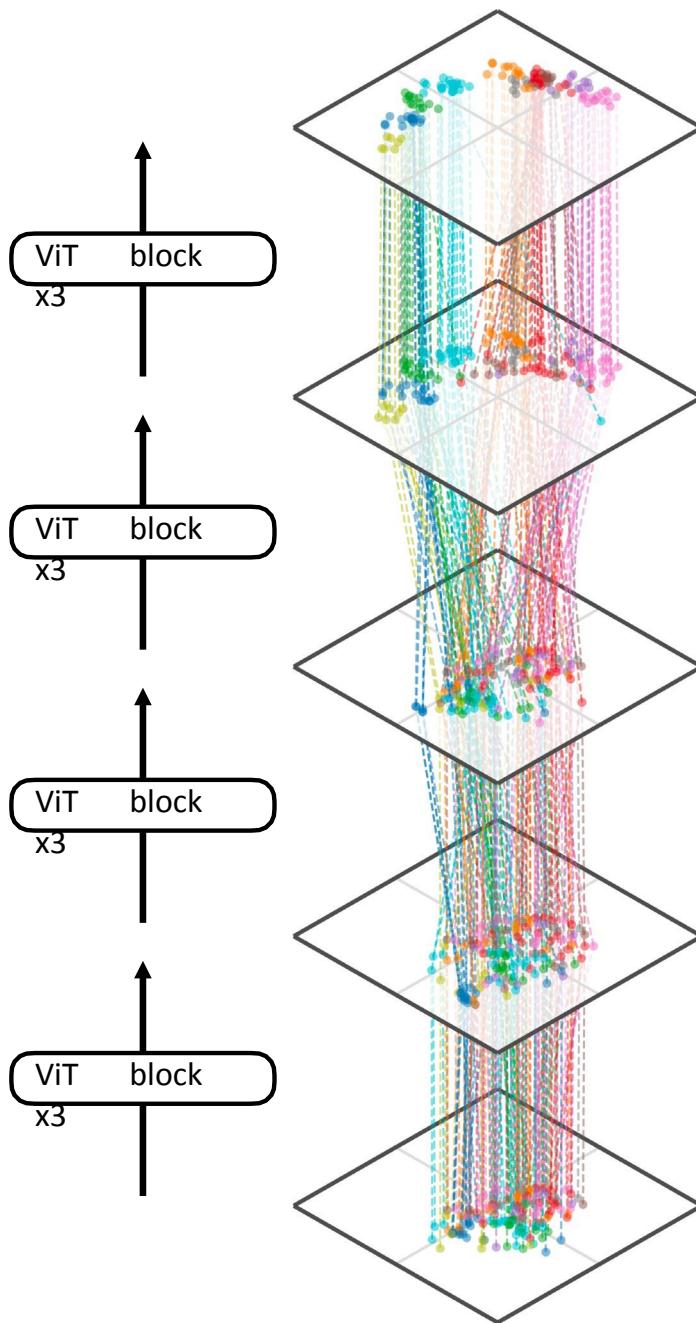
ML  
P



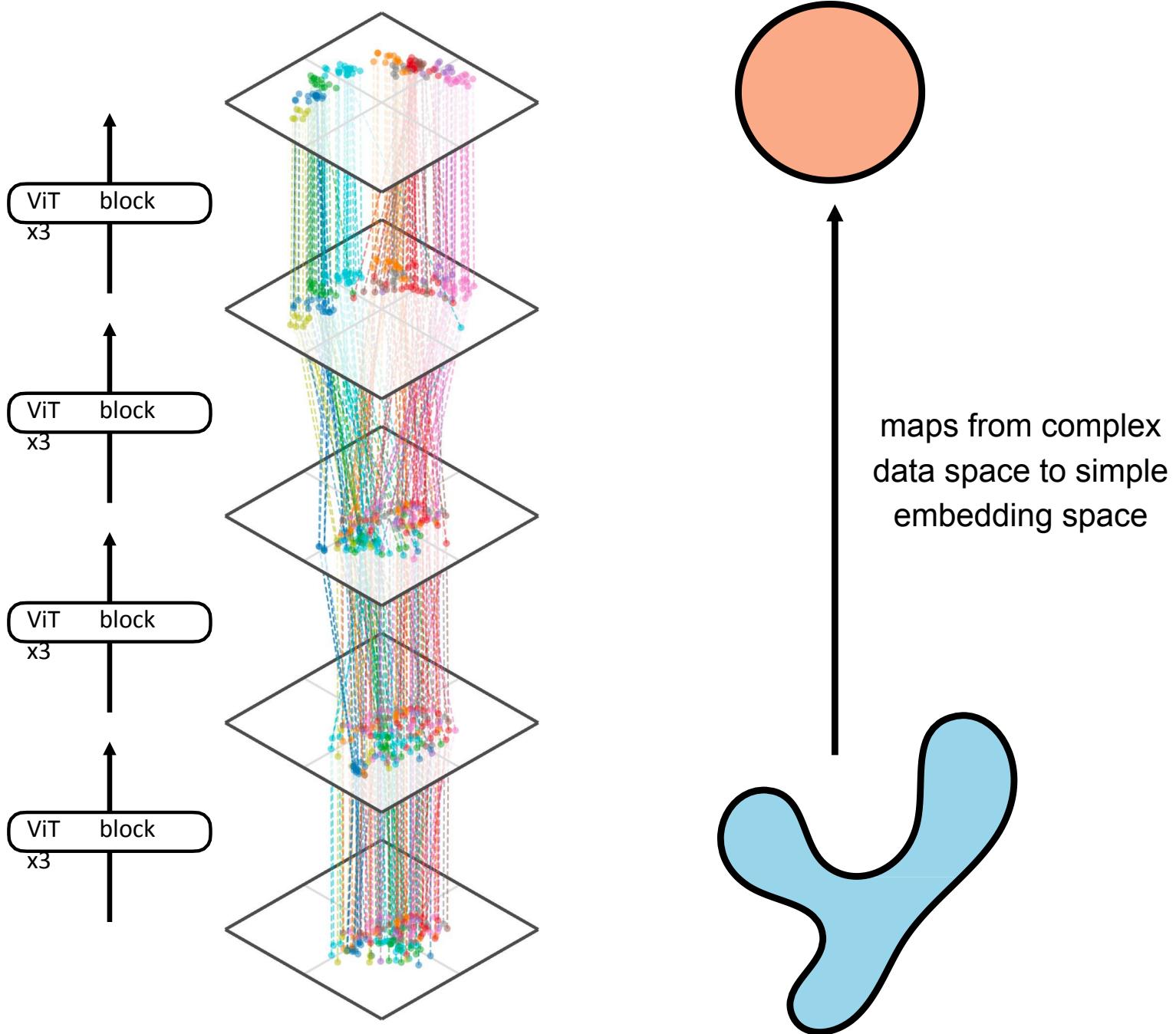
ML  
P



# CLIP



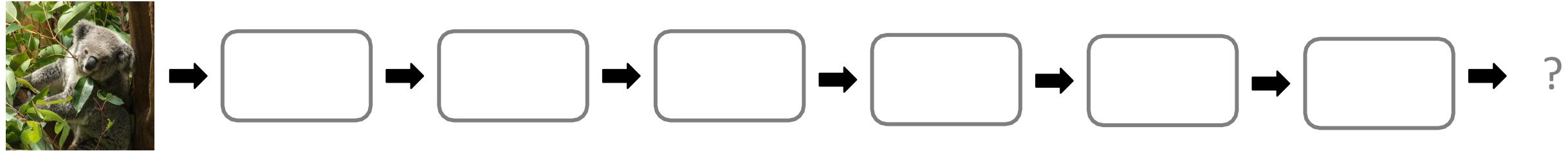
# CLIP



# What Do Deep Networks Learn?

# How to Represent Images w/ Deep Learning?

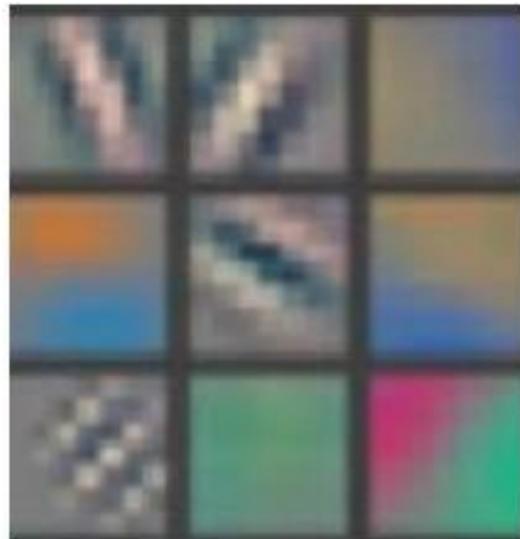
general modules (instead of specialized features)



**compose simple modules into complex functions**

- build multiple levels of abstractions
- learn by back-prop
- learn from data
- reduce domain knowledge and feature engineering

# Multiple Levels of Representations



features



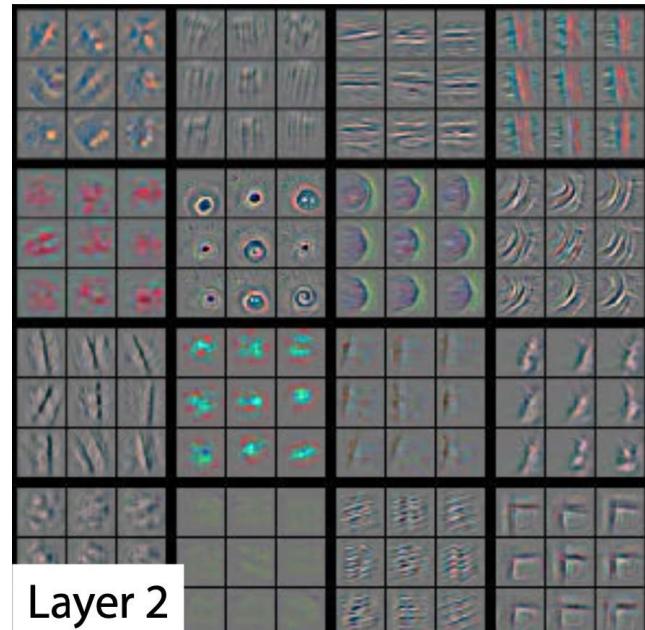
stimuli

(patches with the highest  
1-hot activations)

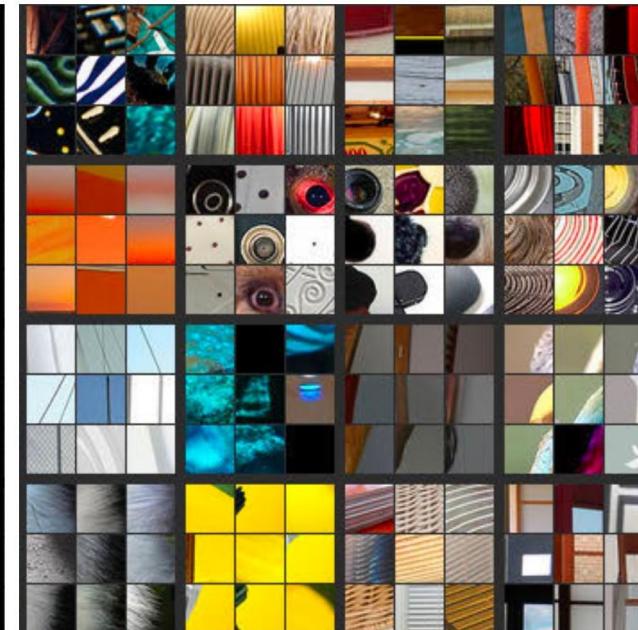
# Multiple Levels of Representations



# features



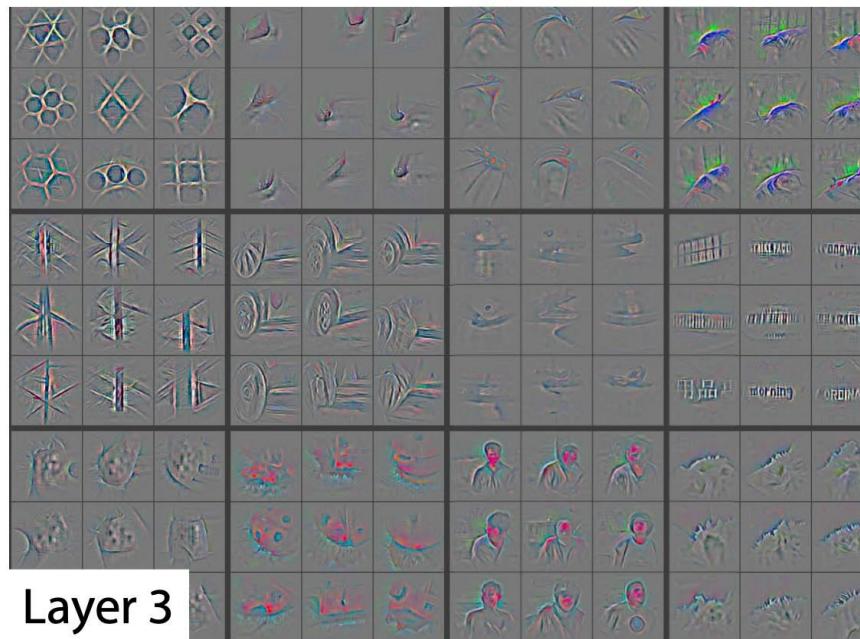
## stimuli



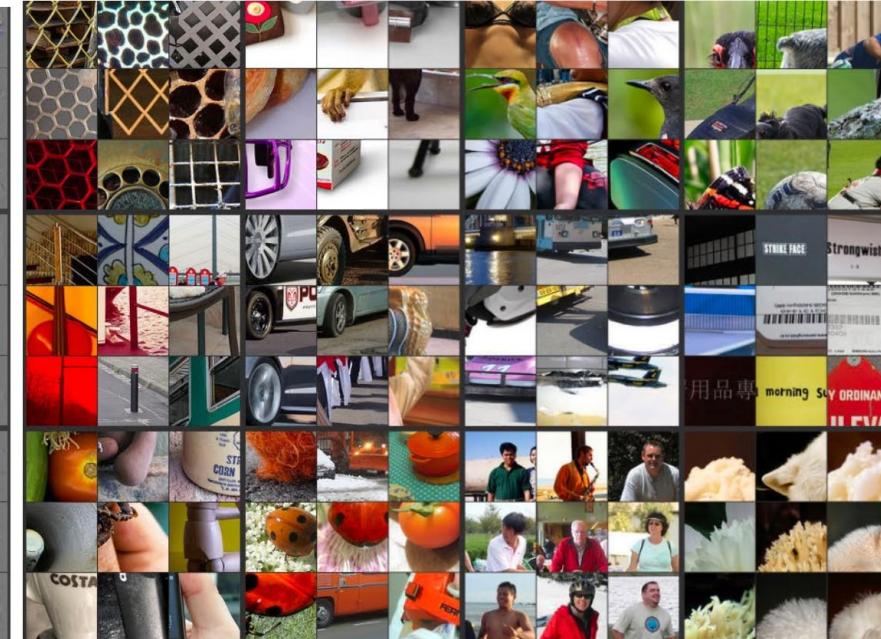
# Multiple Levels of Representations



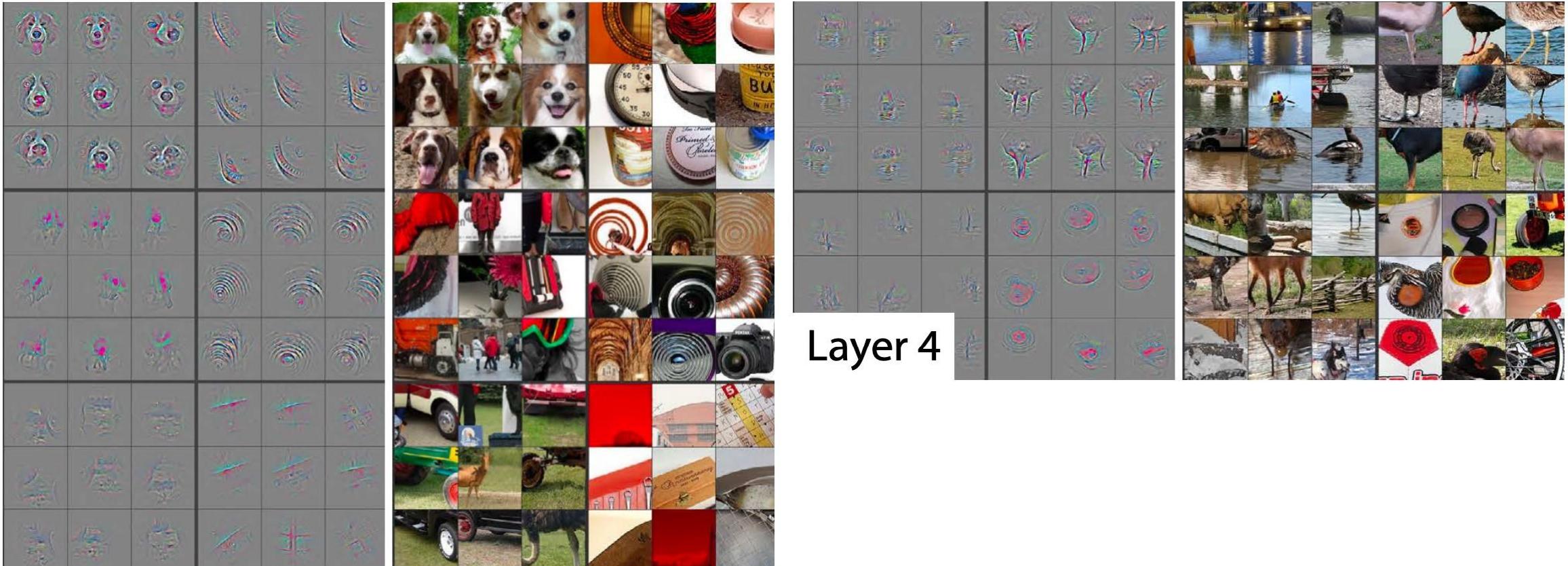
# features



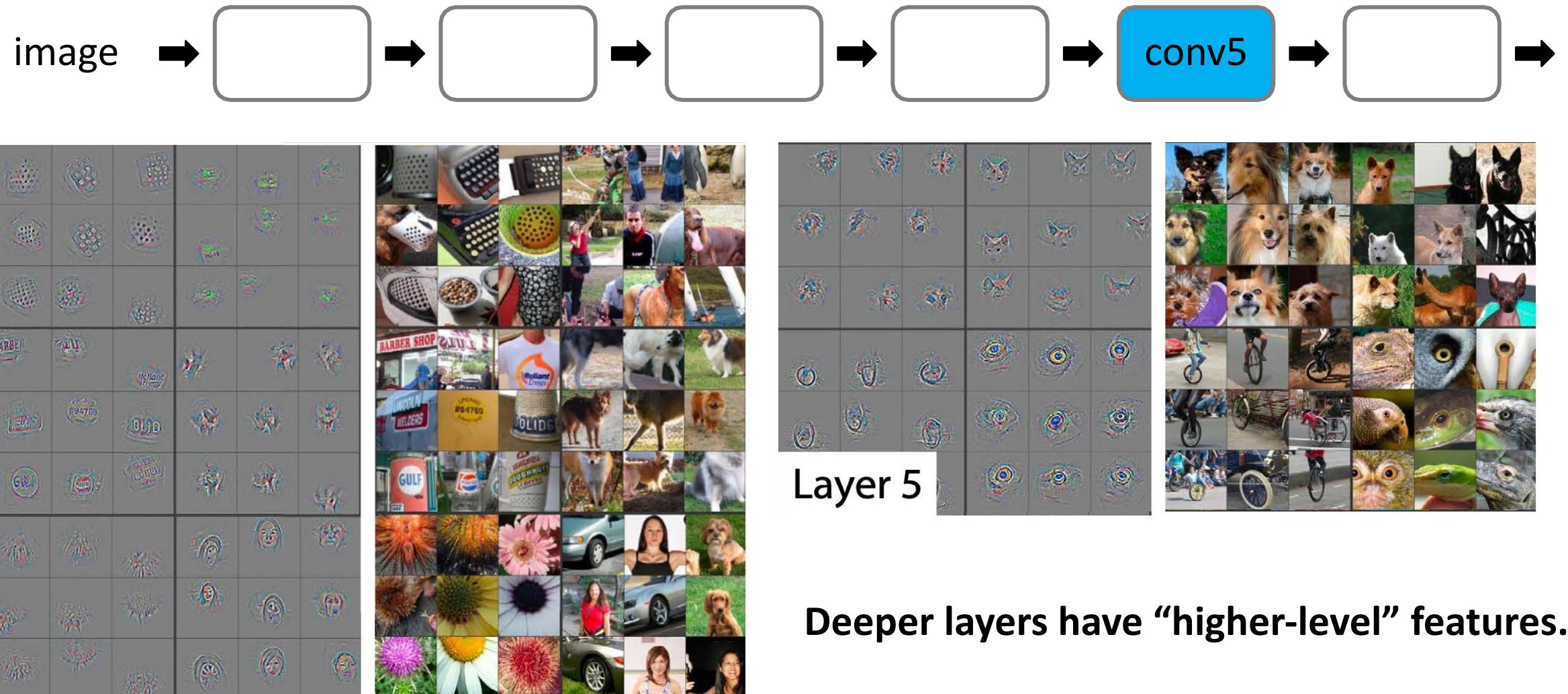
## stimuli



# Multiple Levels of Representations

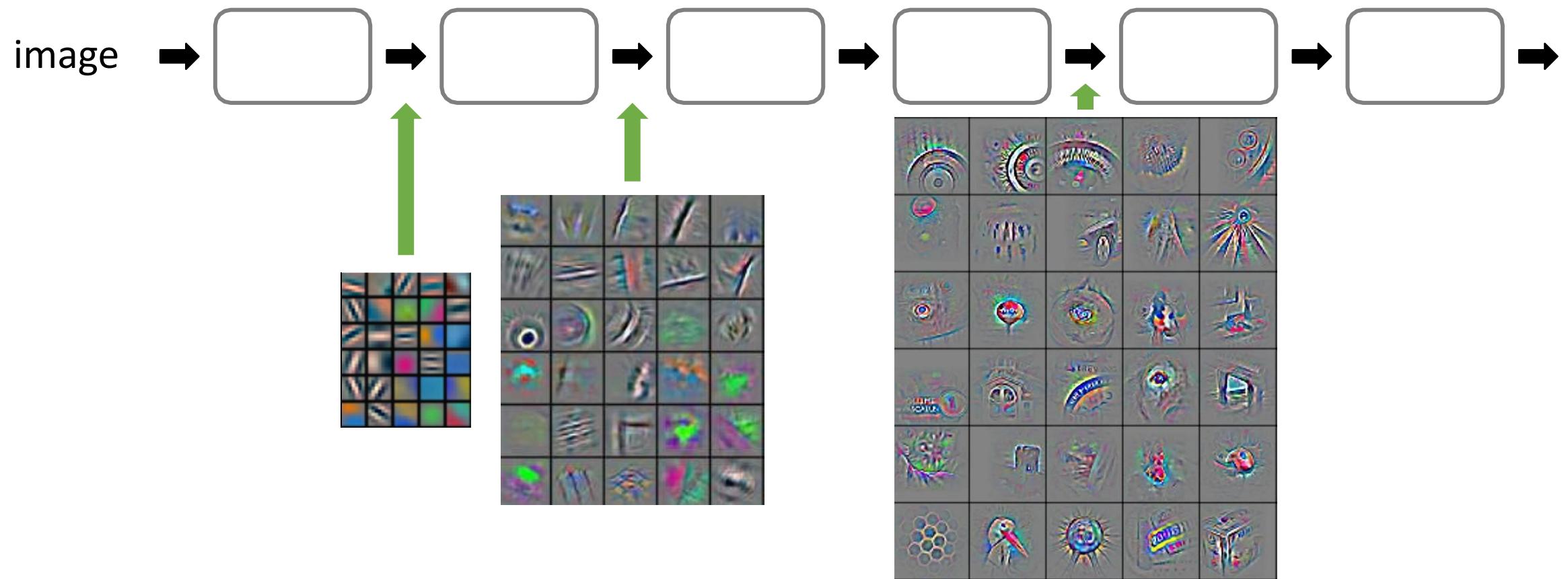


# Multiple Levels of Representations



Deeper layers have “higher-level” features.

# Multiple Levels of Representations



**Deeper layers have “higher-level” features.**

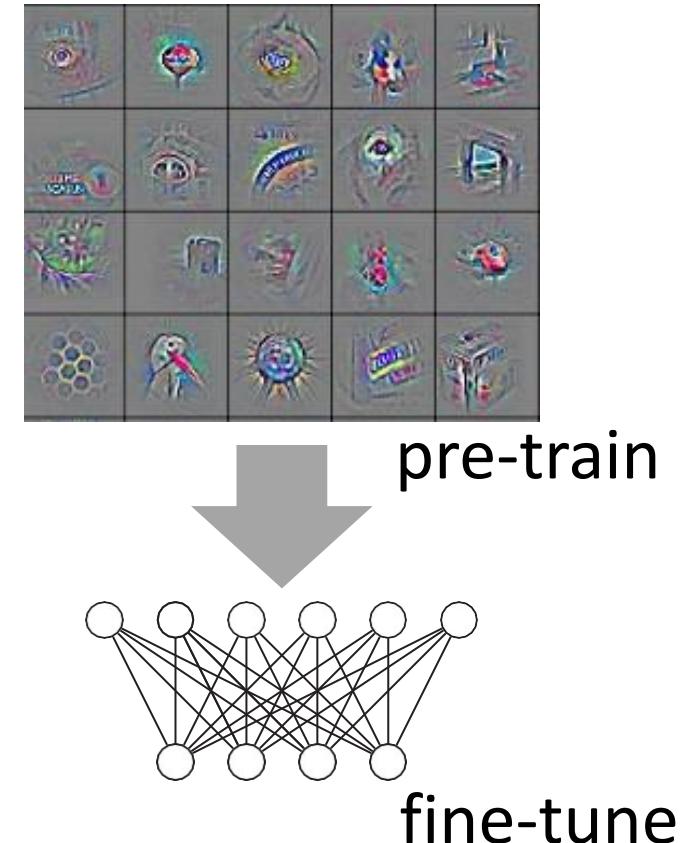
# Deep Representations are

## Transferable!

The single most important discovery in DL revolution

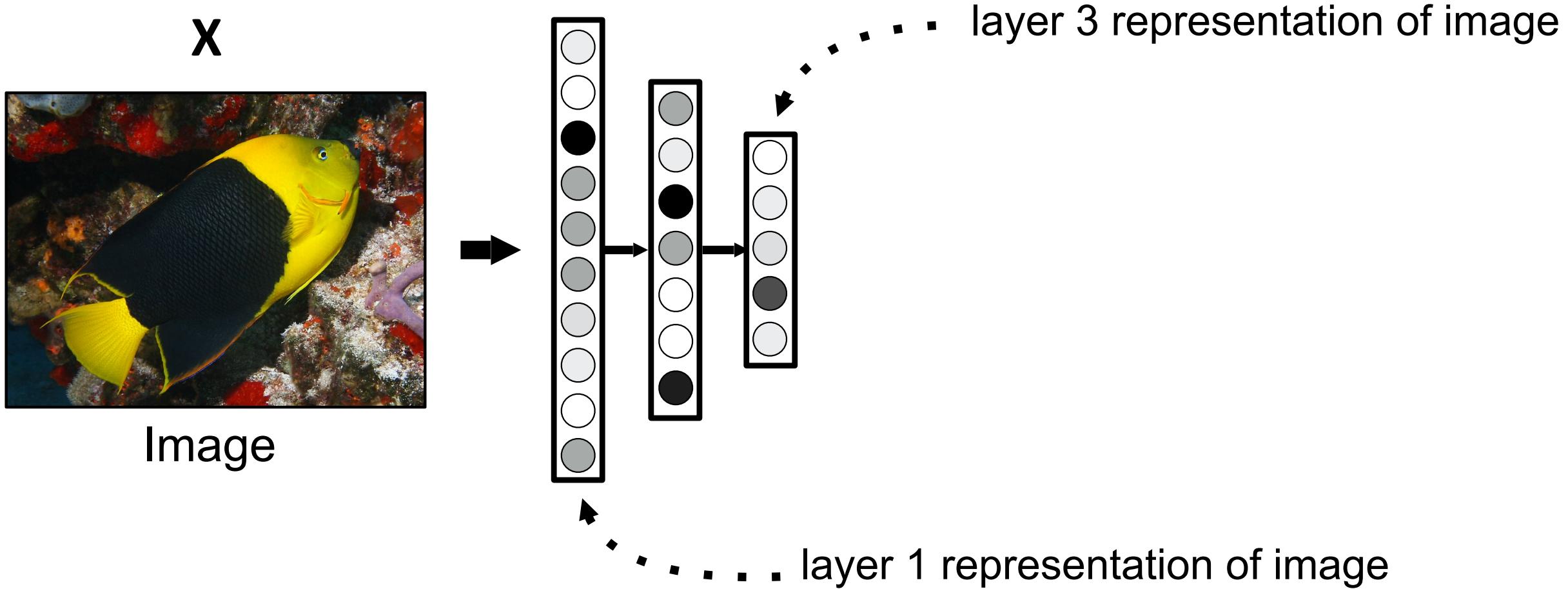
Transfer learning:

- pre-train on large-scale data
- fine-tune on small-scale data
- enable DL for small datasets
- revolutionize computer vision
- data: engine for general representation
- GPT: a similar principle



"DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", Donahue et al. arXiv 2013  
"Visualizing and Understanding Convolutional Networks", Zeiler & Fergus. arXiv 2013  
"CNN Features off-the-shelf: an Astounding Baseline for Recognition", Razavian. arXiv 2014

# Classification Networks Produce Good Features



Represent data as a neural **embedding** — a vector/tensor of neural activations  
(perhaps representing a vector of detected texture patterns or object parts)

# What we covered so far?

1. Neural Networks naturally learn representations
2. Good representations are transferable/useful for downstream tasks
3. We have huge unlabeled datasets, small labeled ones

HOW CAN WE LEARN  
WITHOUT LABEELS?!!!



# Supervised Learning

Classification

Object Detection,  
Segmentation etc

# Unsupervised Learning

Dimensionality  
Reduction

Clustering

PCA, t-SNE

K-Means

No Labels

# Self-Supervised Learning (SSL)

Compressive

Predictive

Contrastive

# Self-Supervised Learning (SSL)

Compressive

Predictive

Contrastive

# Learning via compression

x



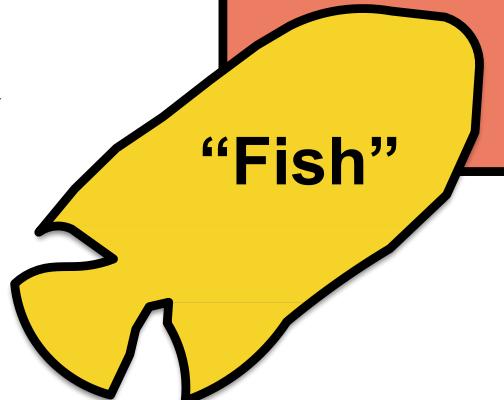
Image

# Learning via compression

X



Image



Compact  
mental  
representation

# Learning via compression

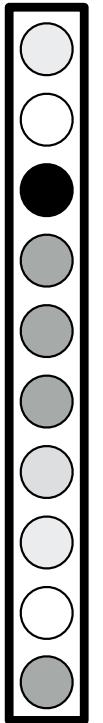
x



Image

# Learning via compression

X



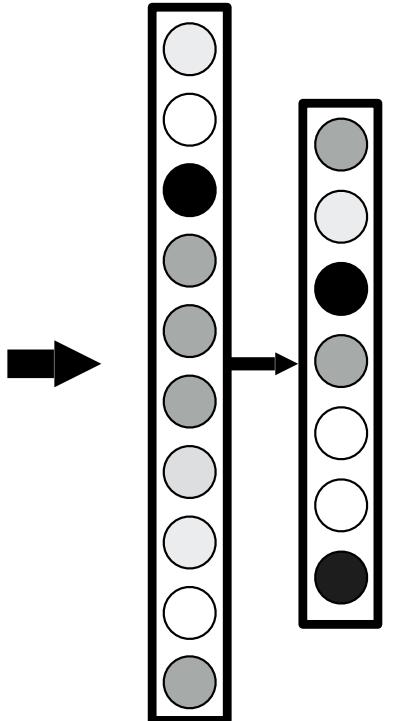
Image

# Learning via compression

X



Image

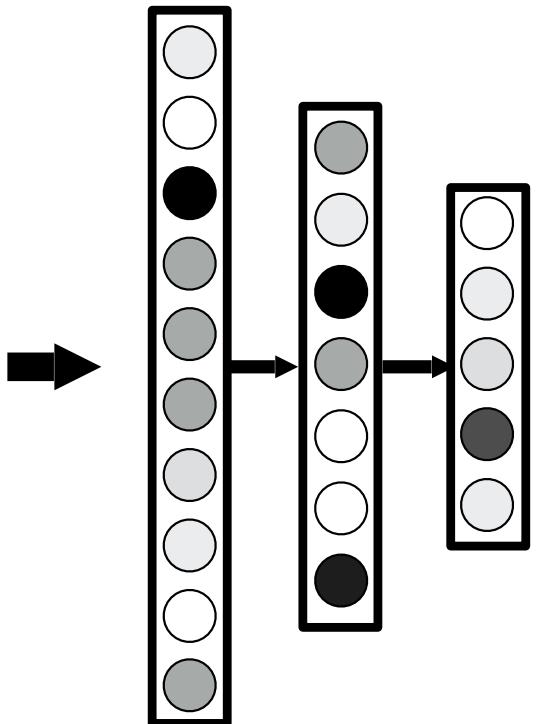


# Learning via compression

X



Image

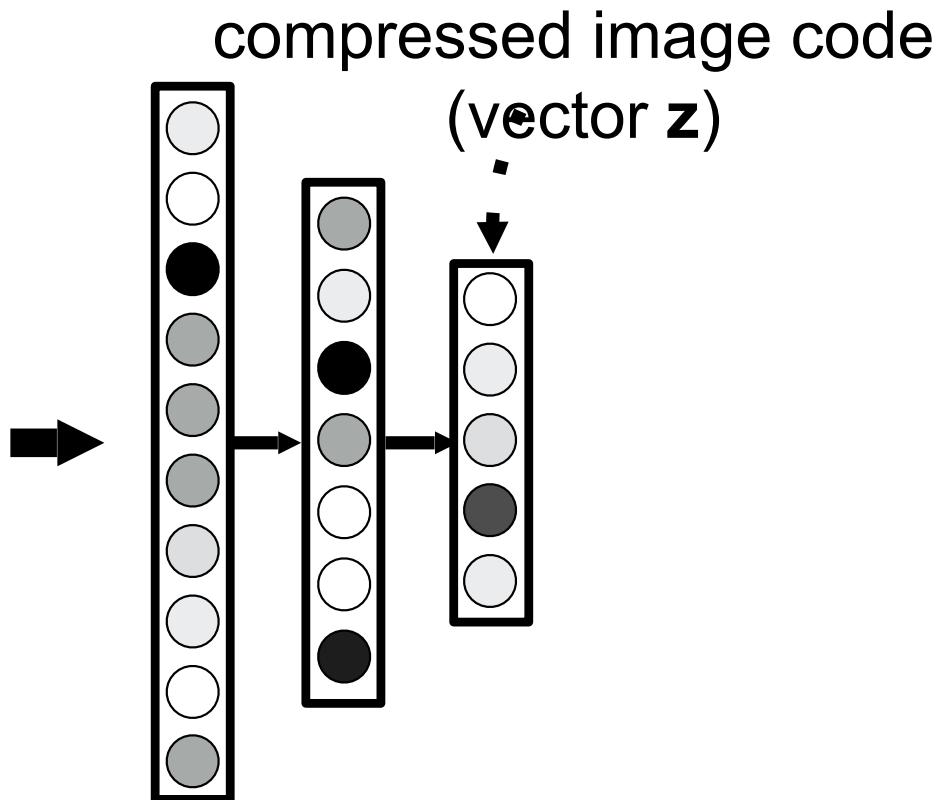


# Learning via compression

**X**

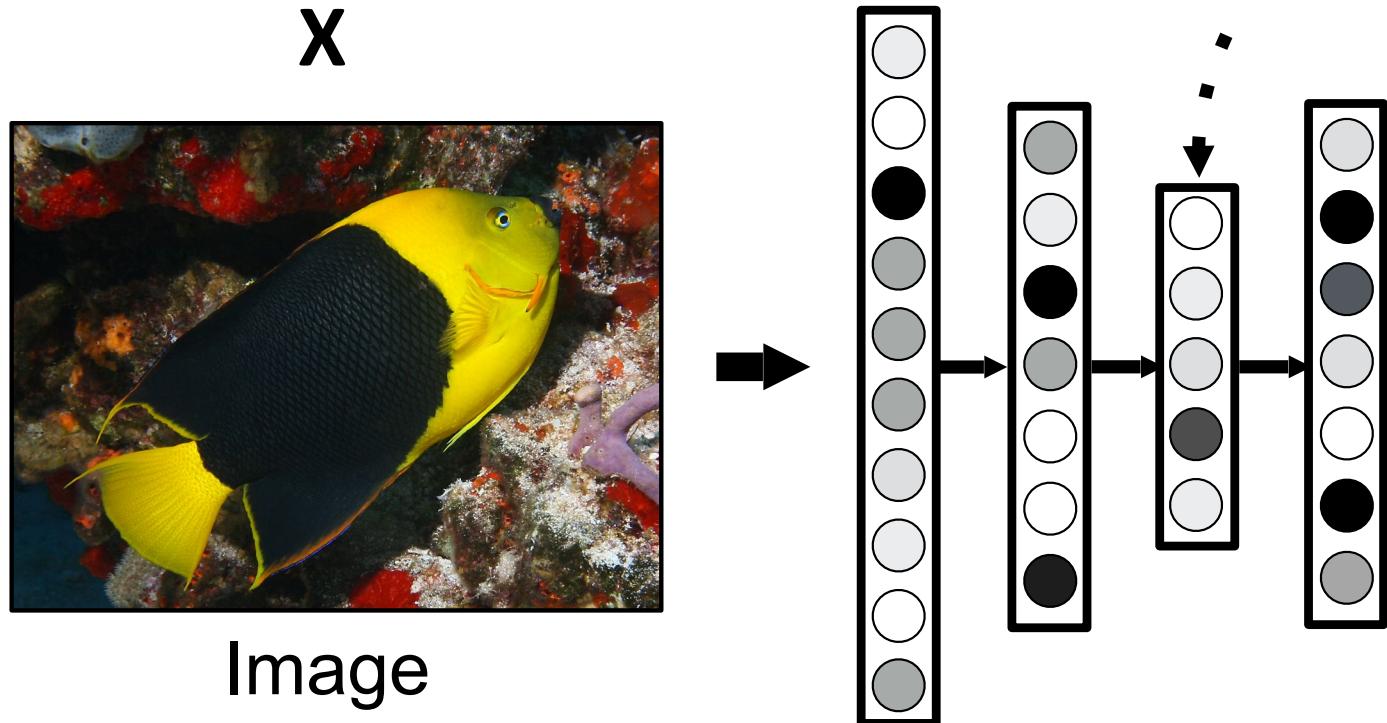


Image



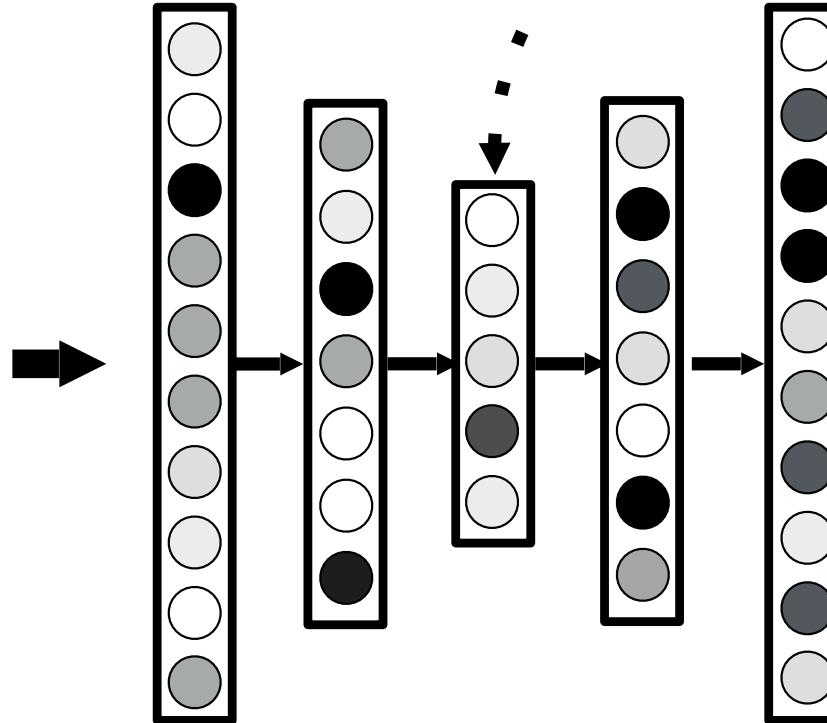
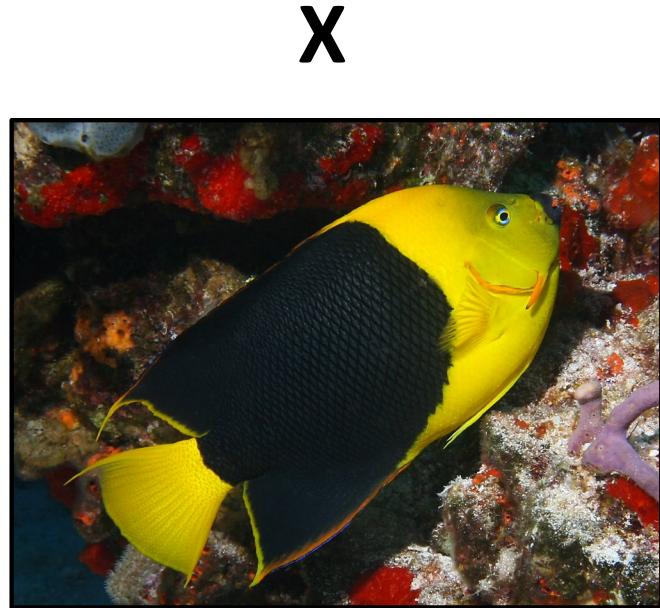
# Learning via compression

compressed image code  
(vector  $z$ )



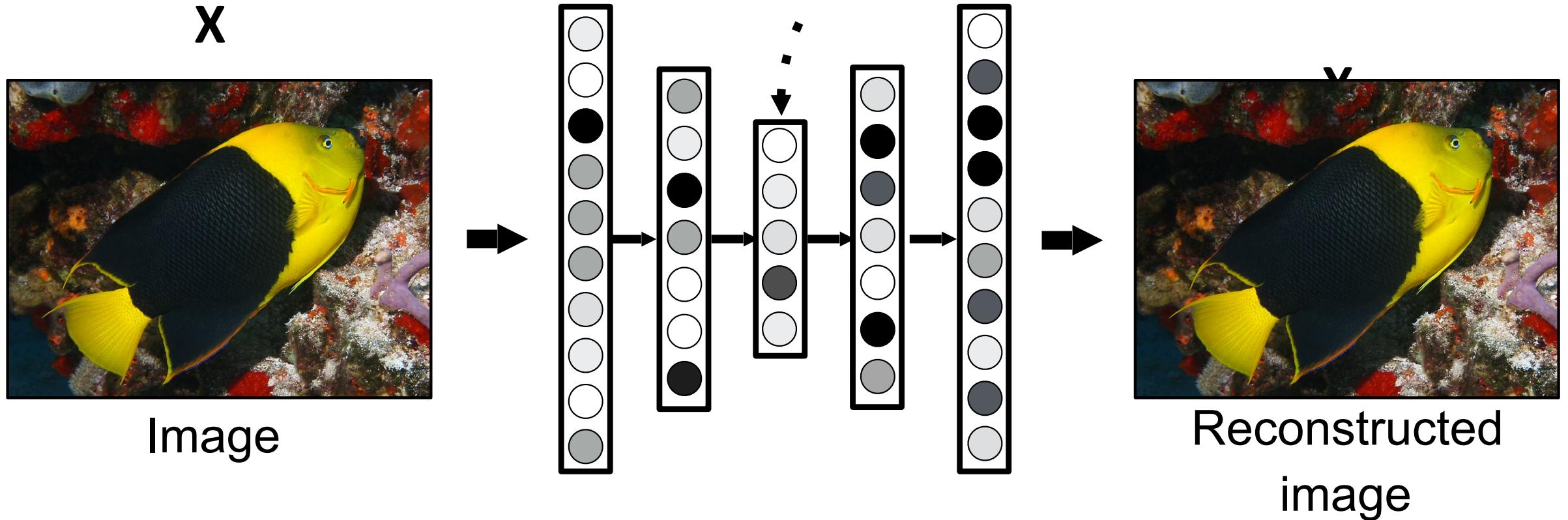
# Learning via compression

compressed image code  
(vector  $z$ )



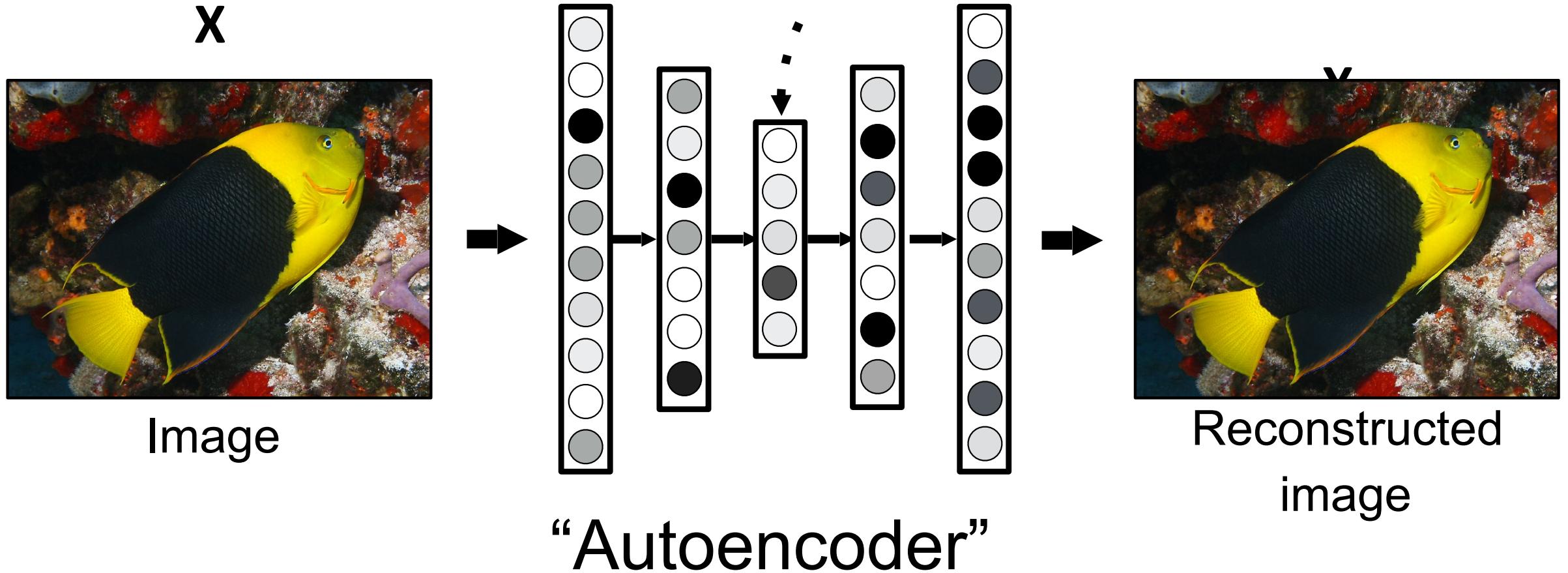
# Learning via compression

compressed image code  
(vector  $z$ )



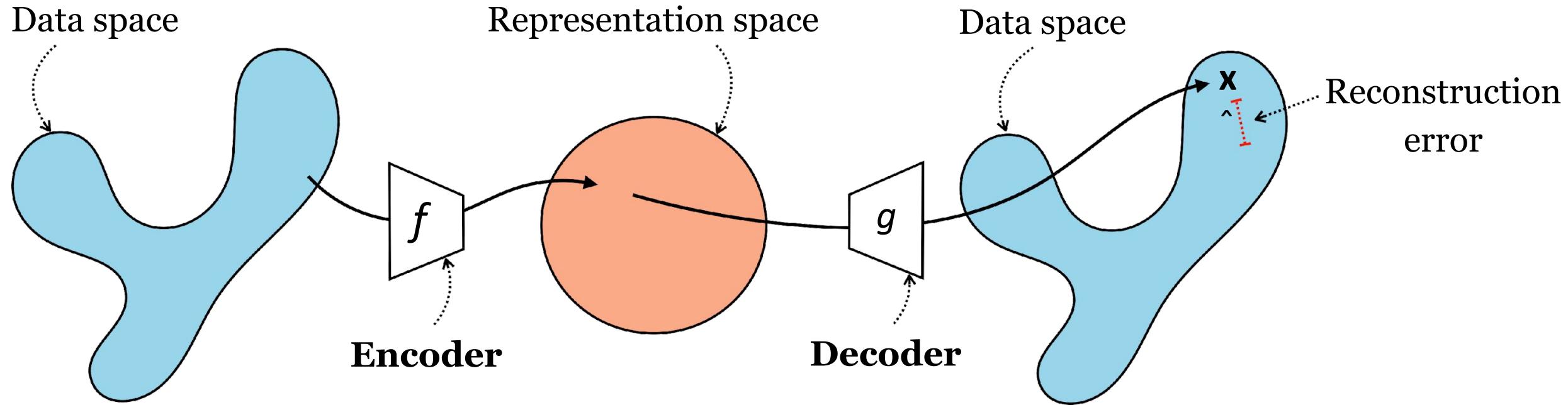
# Learning via compression

compressed image code  
(vector  $z$ )



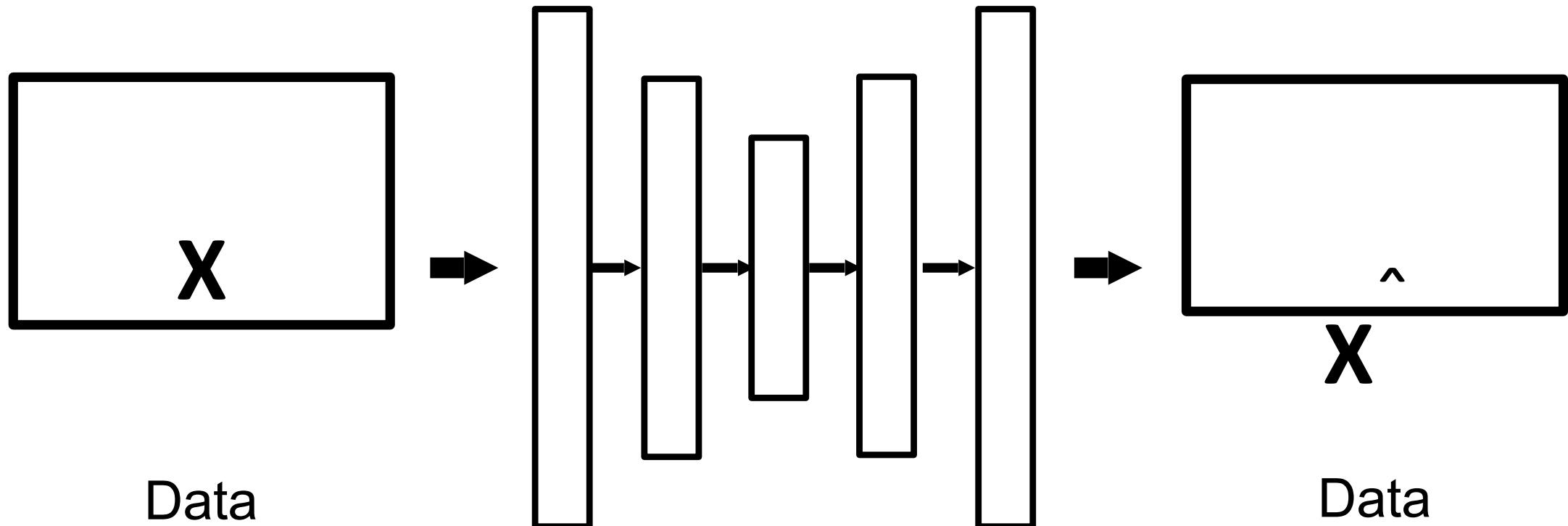
[e.g., Hinton & Salakhutdinov, Science 2006]

# Autoencoder



$$f^*, g^* = \arg \min_{f,g} E_x \| \mathbf{x} - g(f(\mathbf{x})) \|_2^2$$

# Datacompression



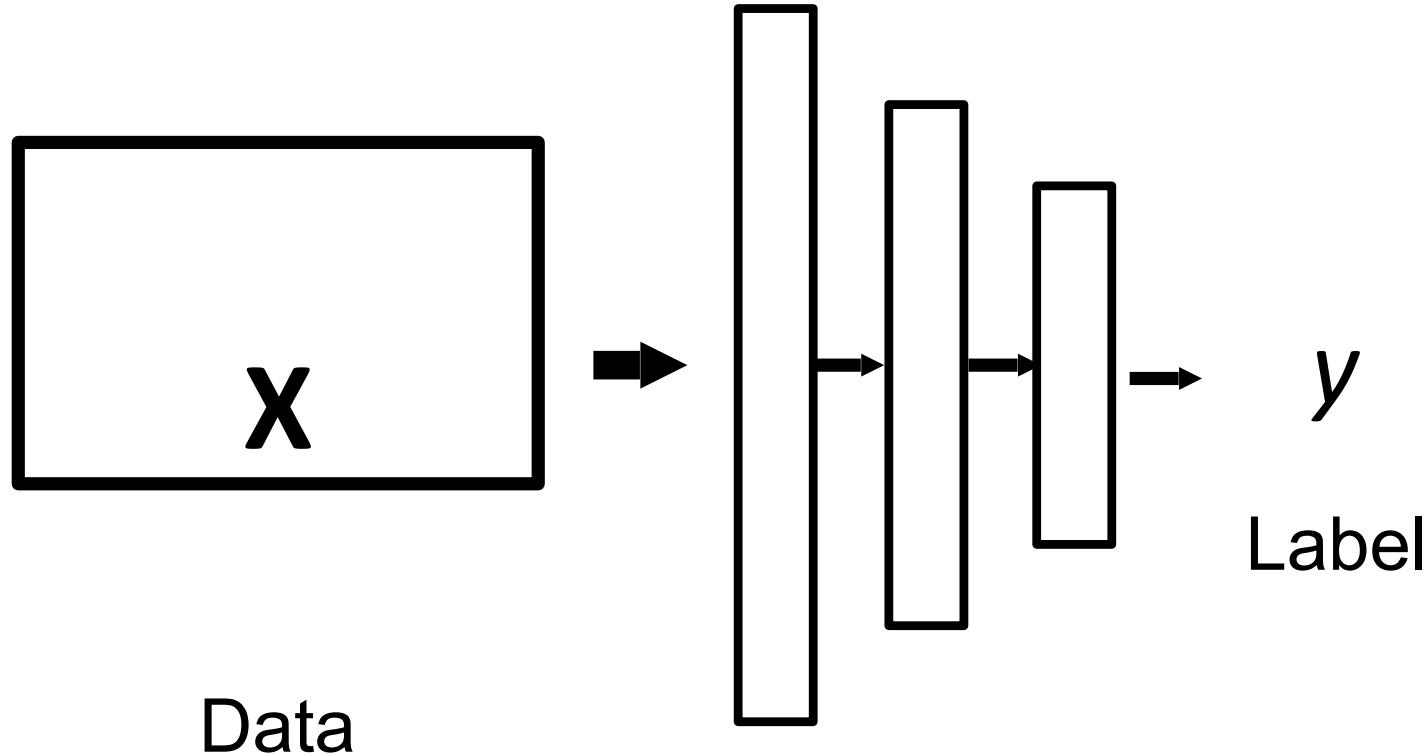
# Self-Supervised Learning (SSL)

Compressive

Predictive

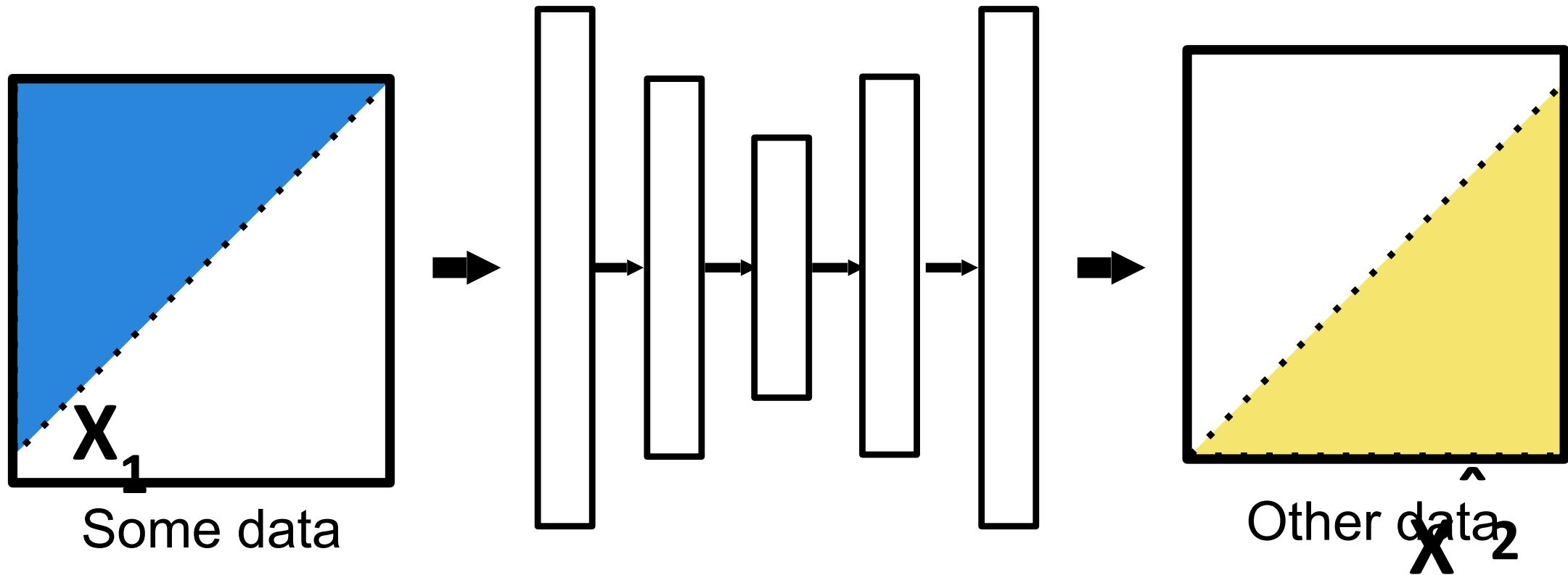
Contrastive

# Label prediction



# Dataprediction

aka “self-supervised learning”

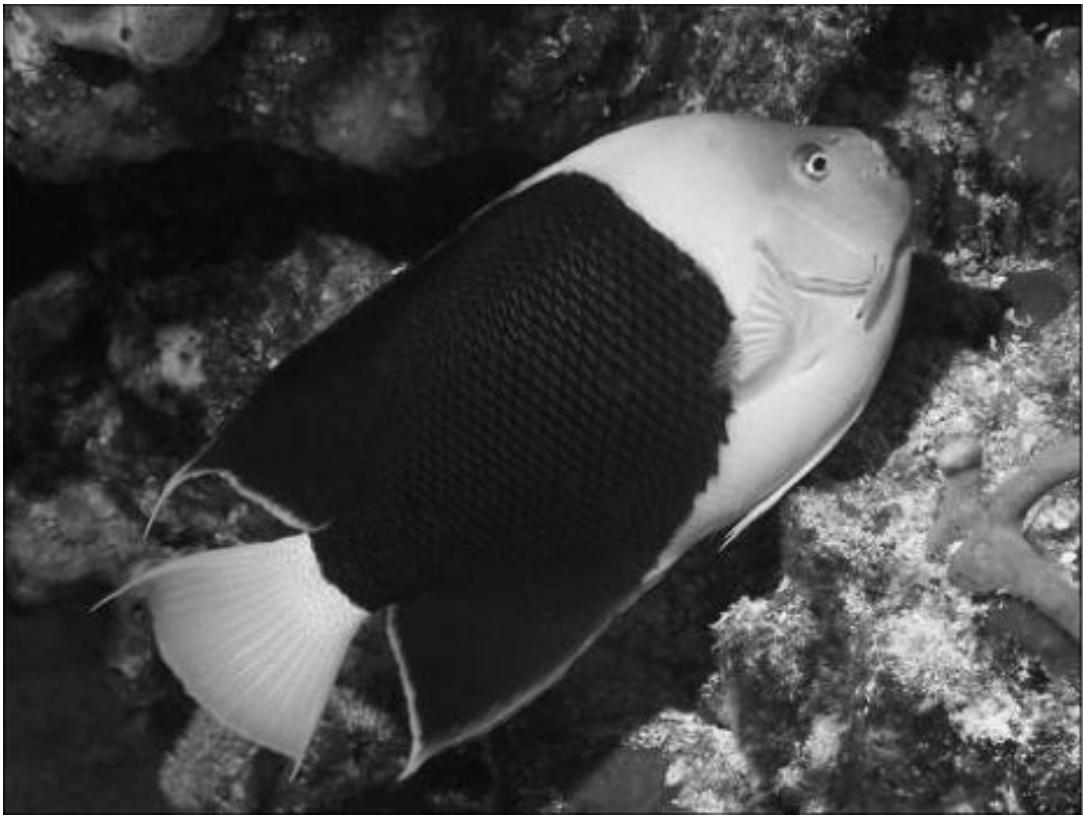




Grayscale image: L  
 $\text{chan}_X \in \mathbb{R}^{H \times W \times 1}$



[Zhang, Isola, Efros, ECCV  
2016]

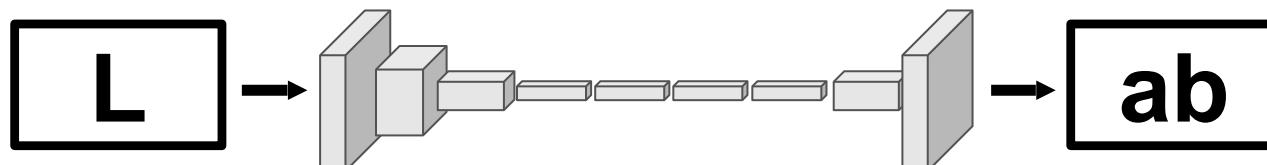


$$\xrightarrow{\mathcal{F}}$$



Grayscale image: L  
chan $\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$

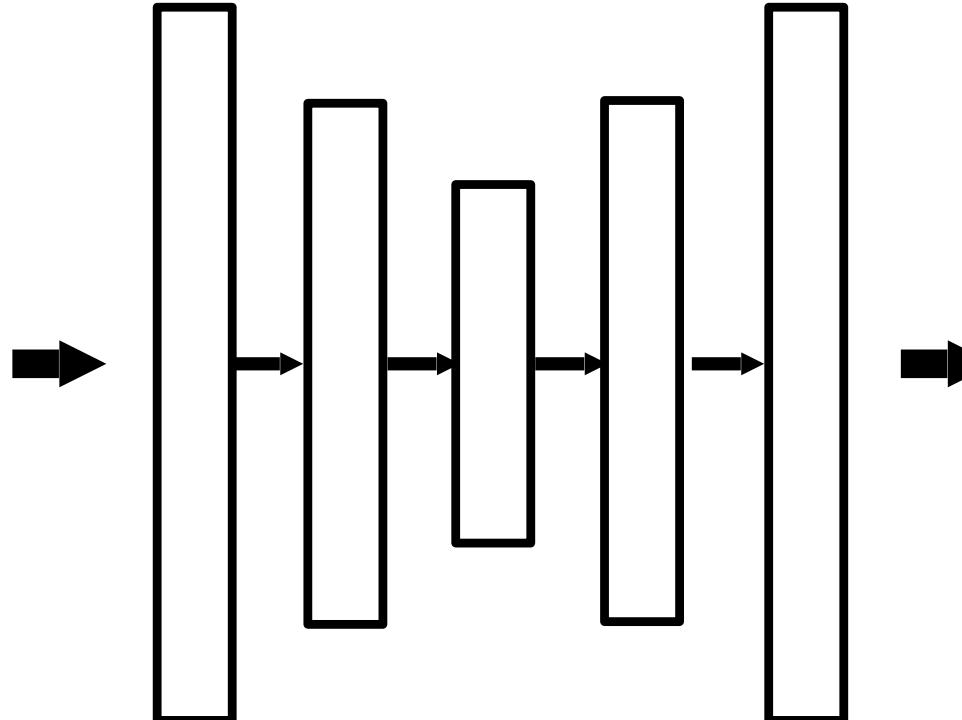
Color information: ab channels  
 $\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$



[Zhang, Isola, Efros, ECCV  
2016]

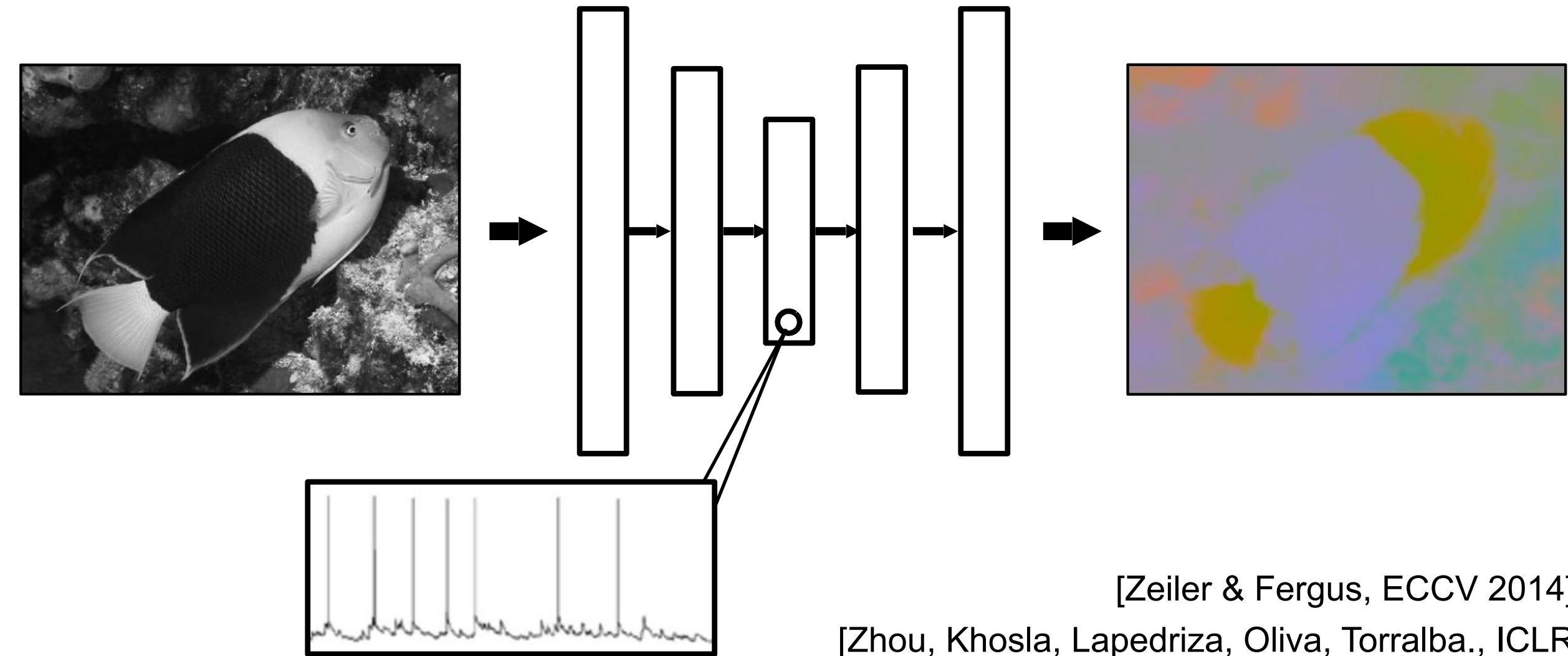
# Deep Net “Electrophysiolog

y”



[Zeiler & Fergus, ECCV 2014]  
[Zhou, Khosla, Lapedriza, Oliva, Torralba., ICLR  
2015]

# Deep Net



# Stimuli that drive selected neurons (conv5 layer)

face  
s



# Stimuli that drive selected neurons (conv5 layer)

face  
s



dog  
face  
s



# Stimuli that drive selected neurons (conv5 layer)

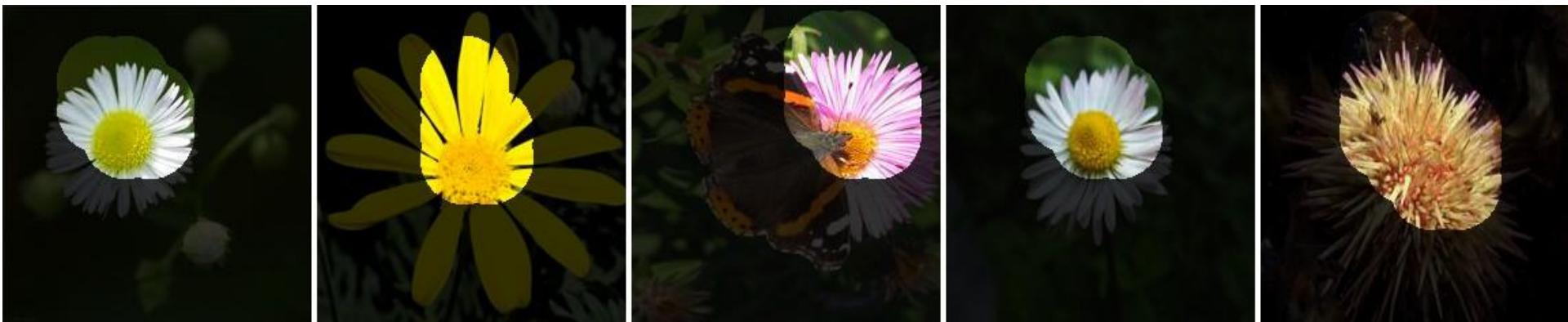
face  
s



dog  
face  
s



flowers



# Predictive Learning: Language Models

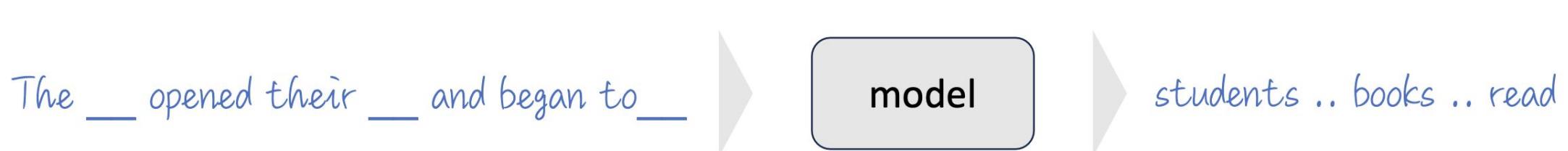
## Next word prediction (GPT)

- Predict the next word (token) given a prefix



## Masked language modeling (BERT)

- Predict the masked words (tokens) in a text



Radford, et al., "Improving Language Understanding by Generative Pre-Training", 201

Devlin, et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding", 201

# Predictive Learning: Computer Vision

## Masked image modeling (Context Encoders)

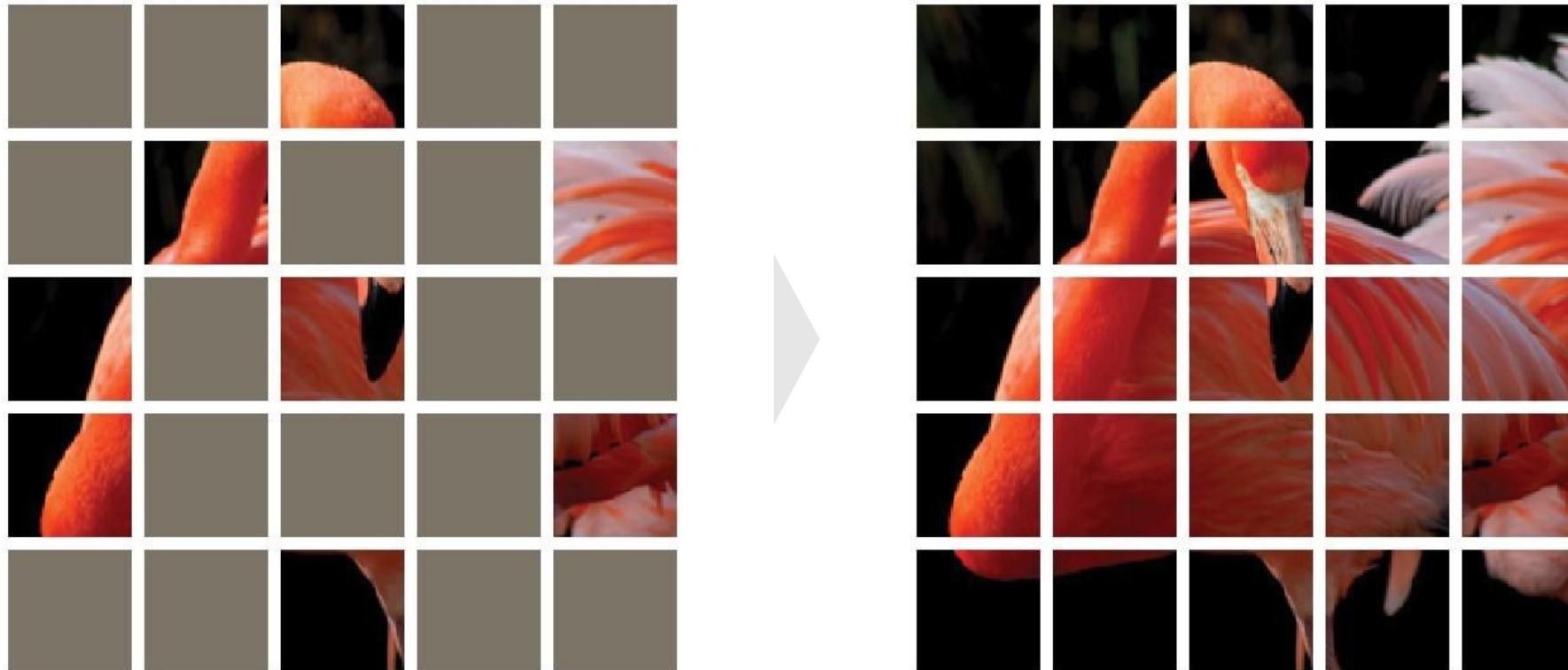
- Predict the masked regions using ConvNets



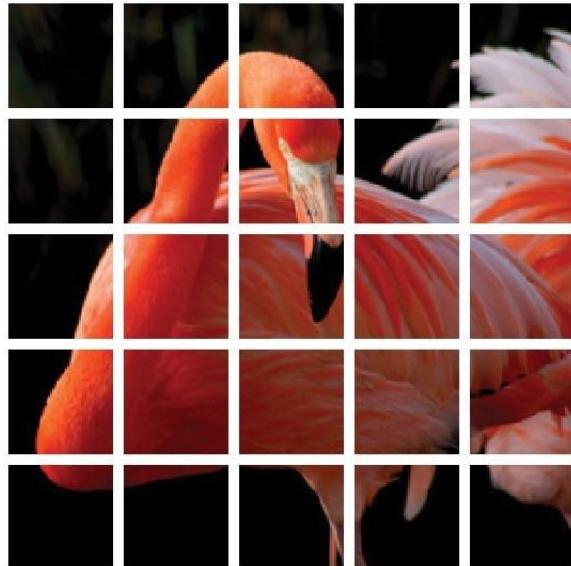
# Predictive Learning: Computer Vision

## Masked image modeling (Masked Autoencoder)

- Predict the masked patches using  
Transformers



# Masked Autoencoder (MAE)



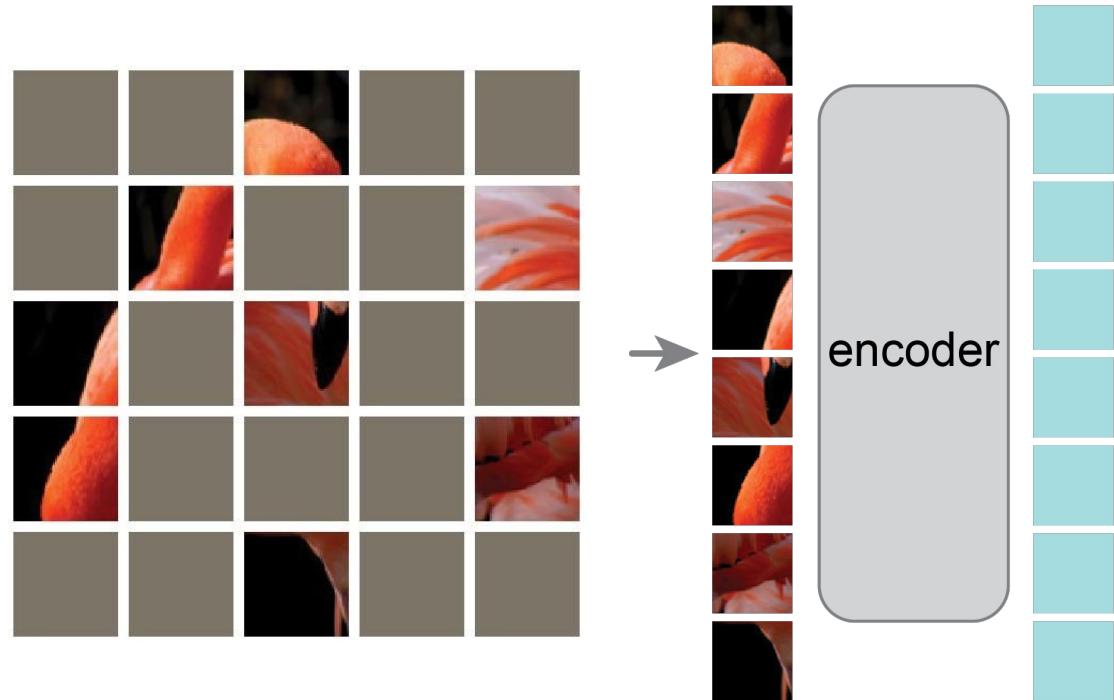
patches as visual  
tokens (Vision  
Transformer)

# Masked Autoencoder (MAE)



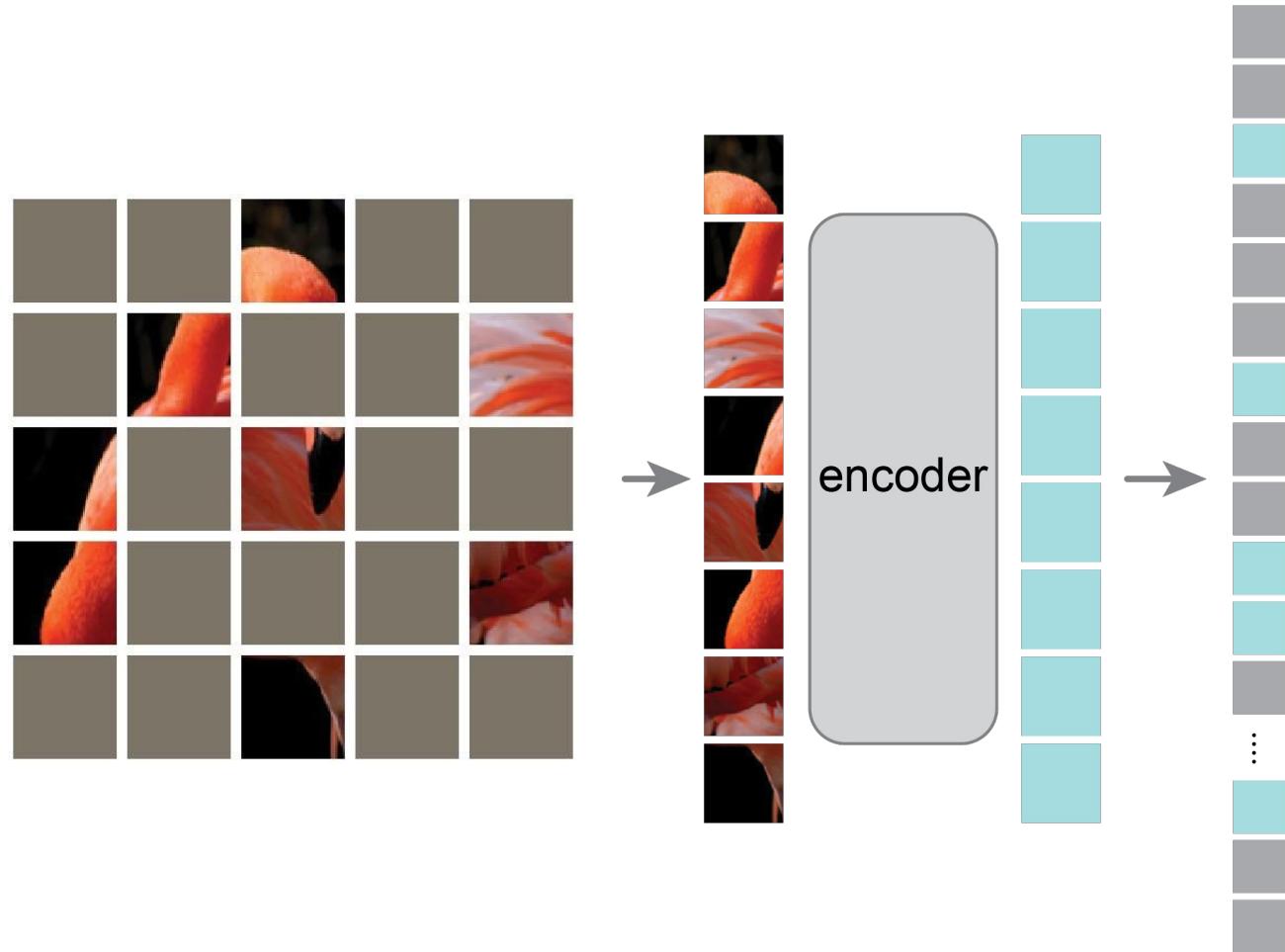
random masking

# Masked Autoencoder (MAE)



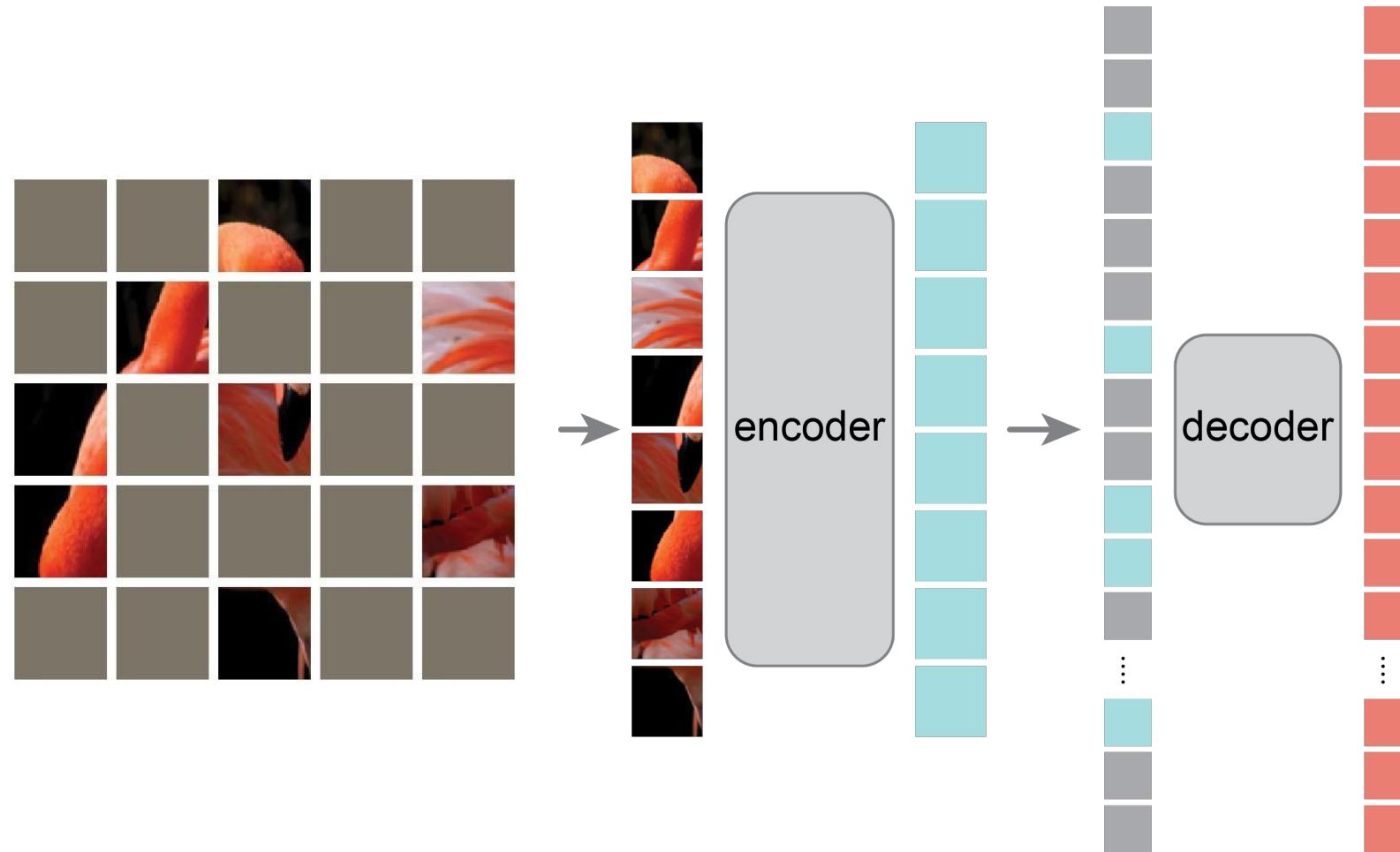
encode visible  
patches w/  
Transformer

# Masked Autoencoder (MAE)



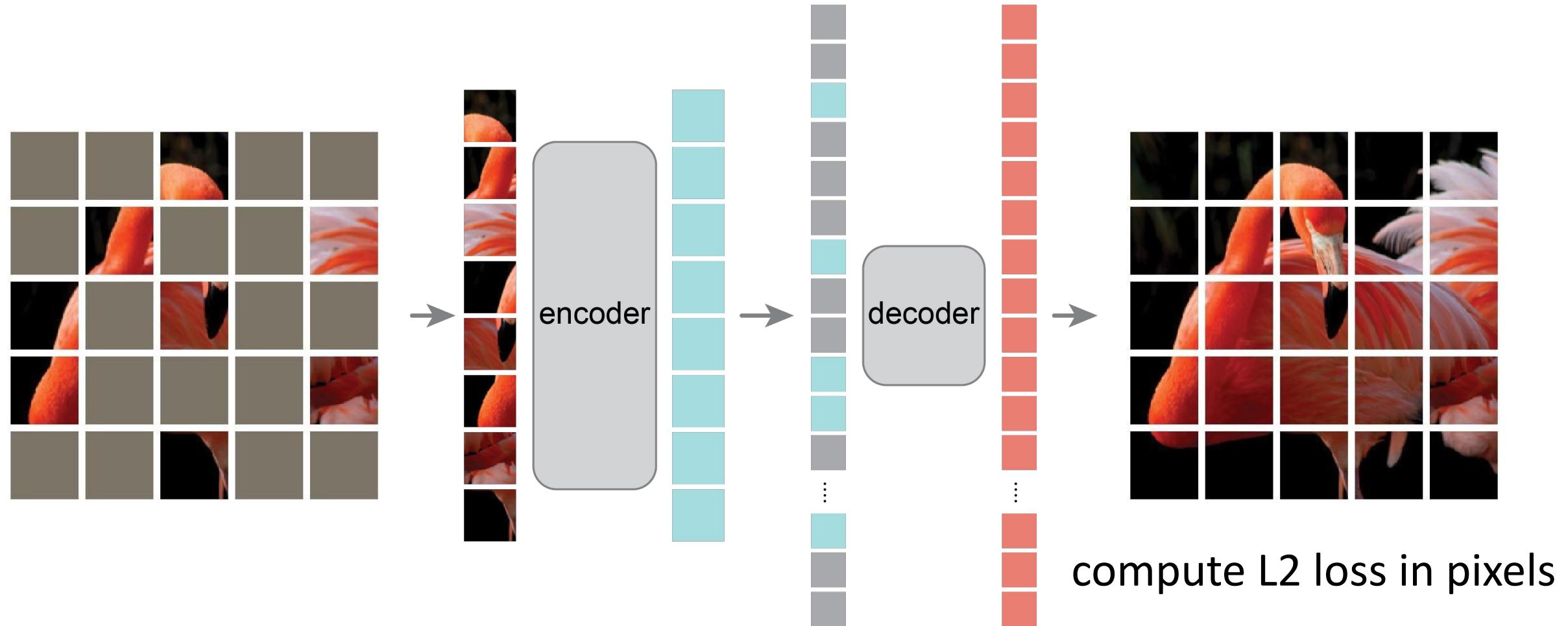
- expand with “mask tokens” (specify where to predict)

# Masked Autoencoder (MAE)

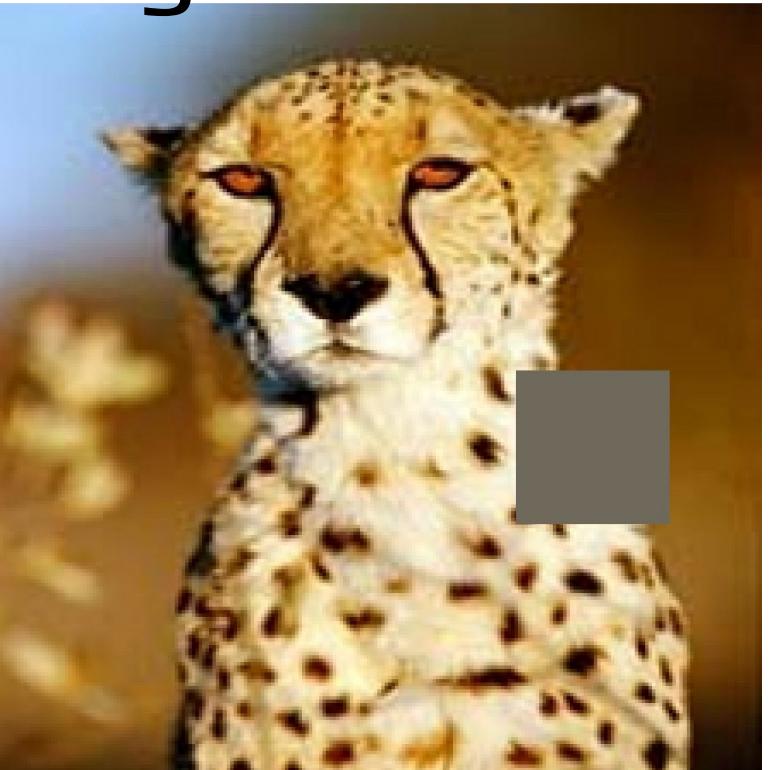


predict the unknown

# Masked Autoencoder (MAE)



# How to learn good representations by predicting?



- predicting a small portion may not require high-level understanding
- predicting a large portion of unknown patches encourages to learn semantic features

# How to learn good representations by predicting?



input



MAE prediction



original



- **The learning process:** the network gradually makes sense of the semantic patterns



- **The learning process:** the network gradually makes sense of the semantic patterns

# Self-Supervised Learning (SSL)

Compressive

Predictive

Contrastive

# Contrastive Learning

Motivation: metric learning



# Contrastive Learning

Motivation: metric learning

Distances in Pixel Space



Distances in Human Representation



# How do we get positive pairs without labels?



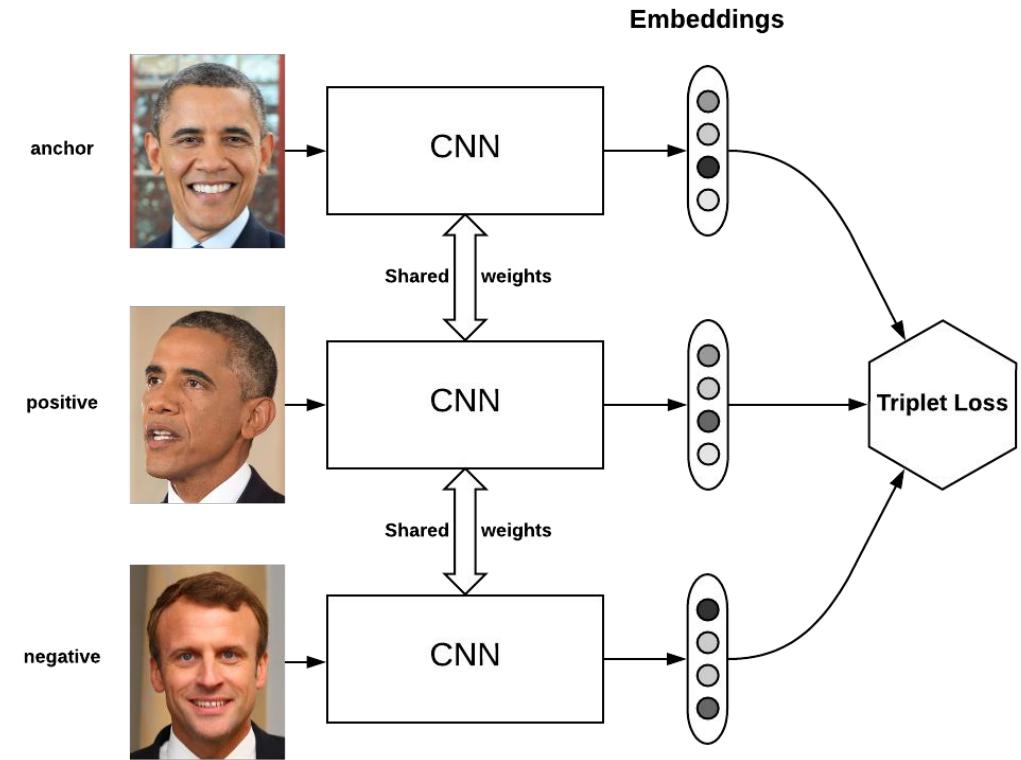
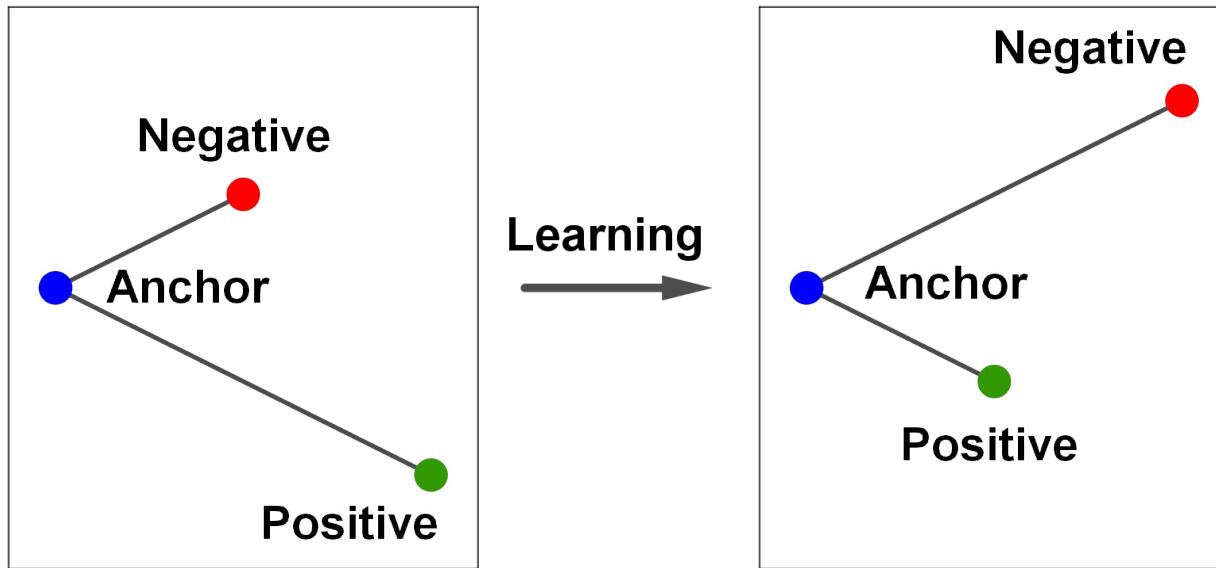
crops



augmentations



# Triplet Loss



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

# SimCLR: Simple Framework for Contrastive Learning



...



# Common setup

- Encoder maps data onto a hypersphere:  $f: \mathcal{X} \rightarrow \mathbb{S}^{d-1}$
- Cross-entropy for softmax “classifier” to discriminate “classes” defined by similarities

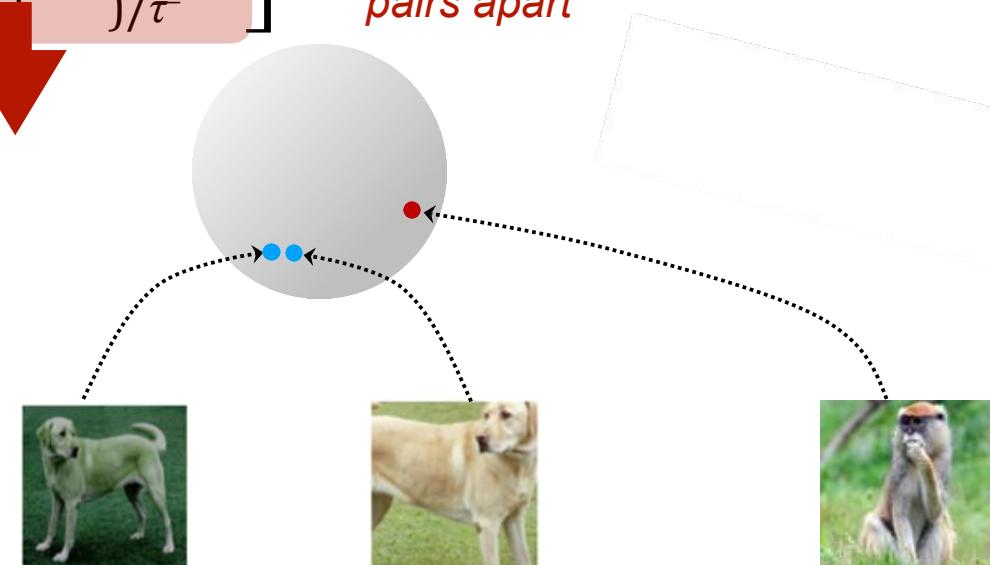
$$\min_f \mathbb{E}_{\substack{(\mathbf{x}, \mathbf{x}^+) \sim p_{\text{pos}} \\ \{\mathbf{x}_i\}_i^N \sim p_{\text{data}}}} \left[ -\log \frac{e^{f(\mathbf{x})^\top f(\mathbf{x}^+)/\tau}}{e^{f(\mathbf{x})^\top f(\mathbf{x}^+)/\tau} + \sum_{i=1}^N e^{f(\mathbf{x})^\top f(\mathbf{x}_i)/\tau}} \right]$$

*pull positive pair together*

*push negative pairs apart*

Symmetry:  $\forall \mathbf{x}, \mathbf{x}^+, p_{\text{pos}}(\mathbf{x}, \mathbf{x}^+) = p_{\text{pos}}(\mathbf{x}^+, \mathbf{x})$

Matching marginal:  $\forall \mathbf{x}, \int p_{\text{pos}}(\mathbf{x}, \mathbf{x}^+) d\mathbf{x}^+ = p_{\text{data}}(\mathbf{x})$



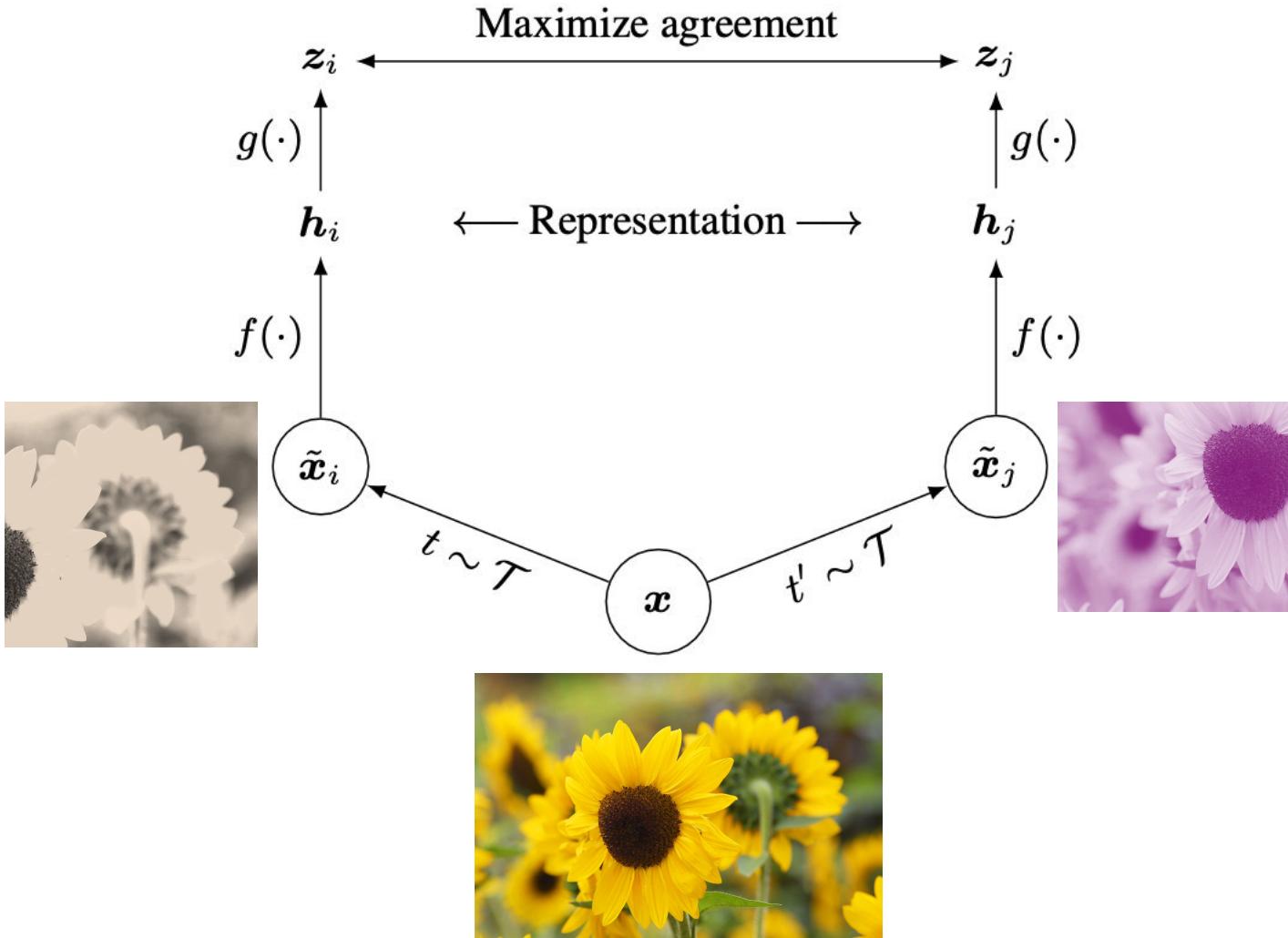
# SimCLR: the implementation details

Algorithm 1 SimCLR's main learning algorithm.

```

input: batch size  $N$ , constant  $\tau$ , structure of  $f, g, \mathcal{T}$ .
for sampled minibatch  $\{\mathbf{x}_k\}_{k=1}^N$  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
         $s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$  # 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  $g(\cdot)$ 

```



Algorithm 1 summarizes the proposed method.

## A Simple Framework for Contrastive Learning of Visual Representations

---



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



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



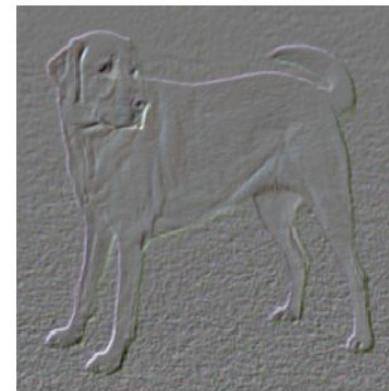
(g) Cutout



(h) Gaussian noise



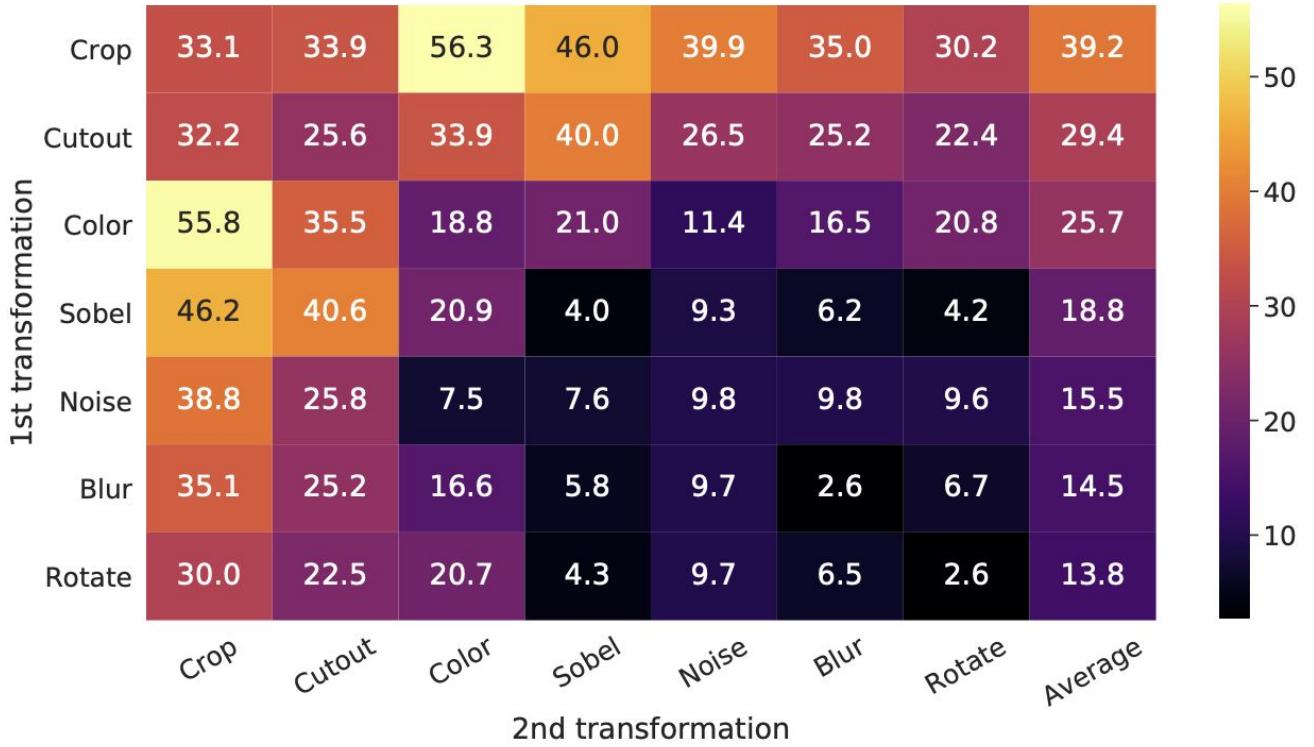
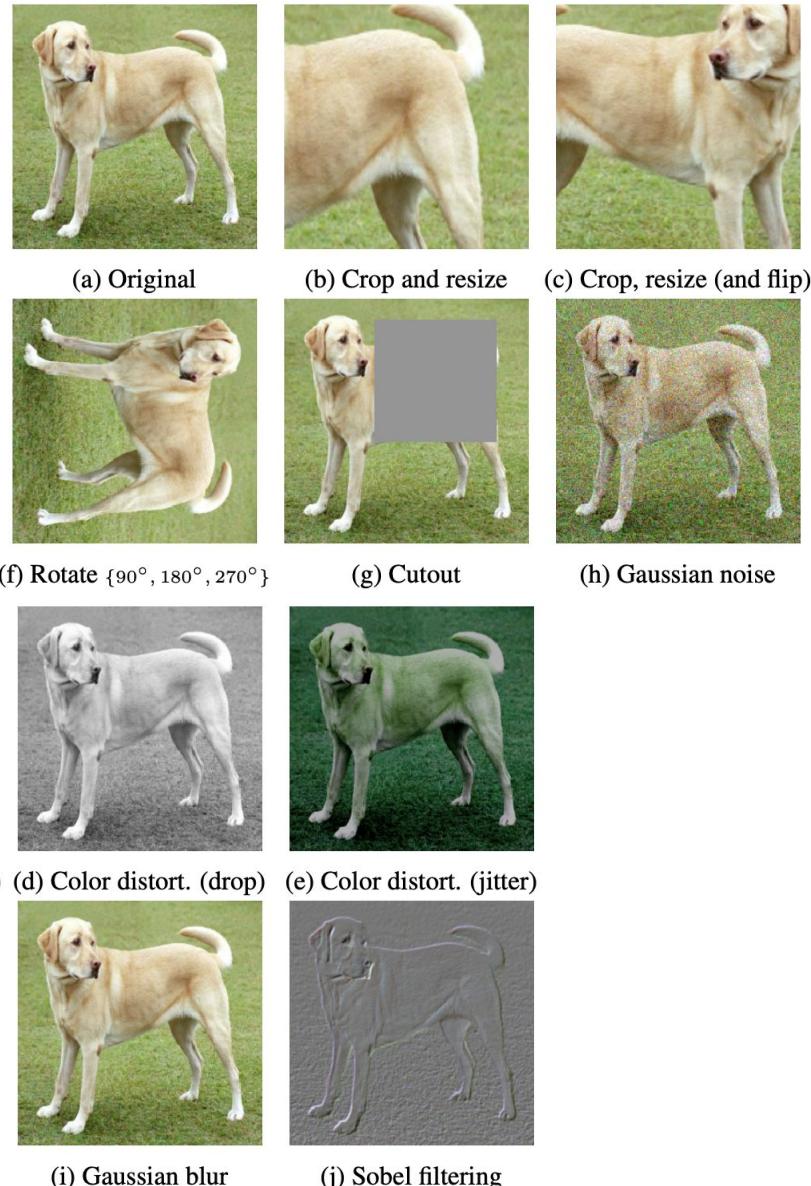
(i) Gaussian blur



(j) Sobel filtering

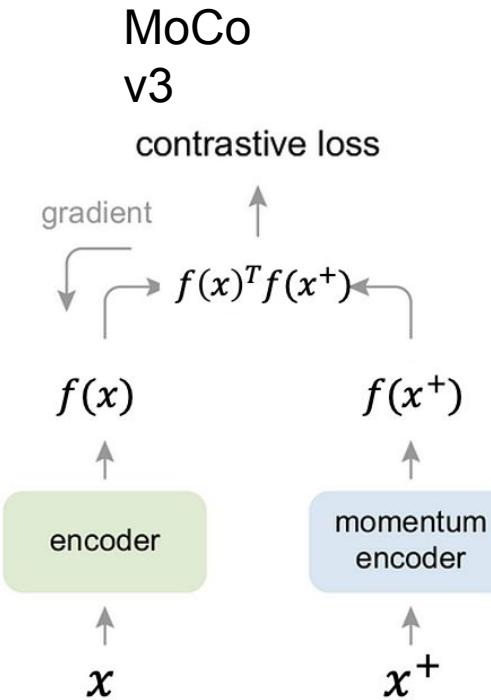
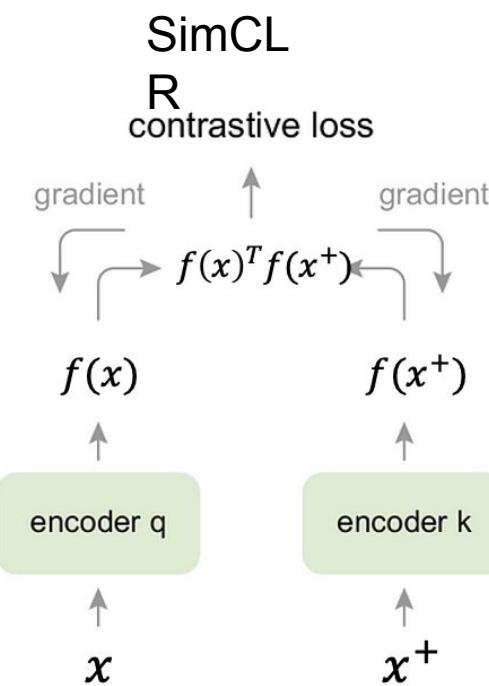
Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy used to train our models* only includes *random crop (with flip and resize)*, *color distortion*, and *Gaussian blur*. (Original image cc-by: Von.grzanka)

# A deeper look into augmentations



*Figure 5.* Linear evaluation (ImageNet top-1 accuracy) under individual or composition of data augmentations, applied only to one branch. For all columns but the last, diagonal entries correspond to single transformation, and off-diagonals correspond to composition of two transformations (applied sequentially). The last column reflects the average over the row.

# MoCo v3: Momentum Contrast



**InfoNCE Loss**

$$\min_f \mathbb{E}_{(x, x^+) \sim p_{pos}, \{x_i^-\}_{i=1}^N \sim p_{data}} \left[ -\log \frac{e^{f(x)^T f(x^+)/\tau}}{e^{f(x)^T f(x^+)/\tau} + \sum_{i=1}^N e^{f(x)^T f(x^-_i)/\tau}} \right]$$

Annotations:

- green arrow pointing up: pull positive pair together
- red arrow pointing down: push negative pairs apart

---

## Algorithm 1 MoCo v3: PyTorch-like Pseudocode

---

```

# f_q: encoder: backbone + proj mlp + pred mlp
# f_k: momentum encoder: backbone + proj mlp
# m: momentum coefficient
# tau: temperature

for x in loader: # load a minibatch x with N samples
    x1, x2 = aug(x), aug(x) # augmentation
    q1, q2 = f_q(x1), f_q(x2) # queries: [N, C] each
    k1, k2 = f_k(x1), f_k(x2) # keys: [N, C] each

    loss = ctr(q1, k2) + ctr(q2, k1) # symmetrized
    loss.backward()

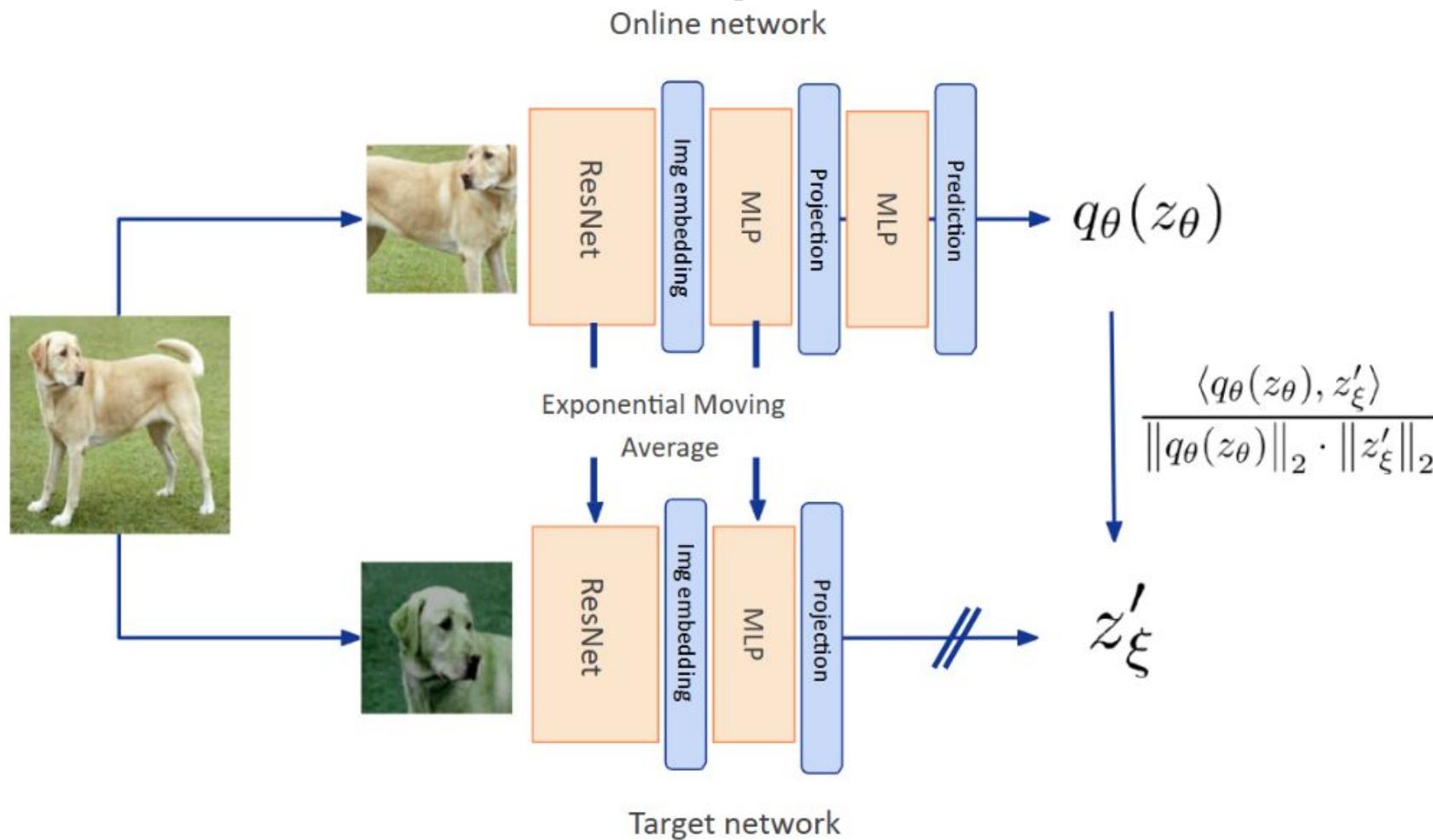
    update(f_q) # optimizer update: f_q
    f_k = m*f_k + (1-m)*f_q # momentum update: f_k

# contrastive loss
def ctr(q, k):
    logits = mm(q, k.t()) # [N, N] pairs
    labels = range(N) # positives are in diagonal
    loss = CrossEntropyLoss(logits/tau, labels)
    return 2 * tau * loss
  
```

---

**Notes:** mm is matrix multiplication. k.t() is k's transpose. The prediction head is excluded from f\_k (and thus the momentum update).

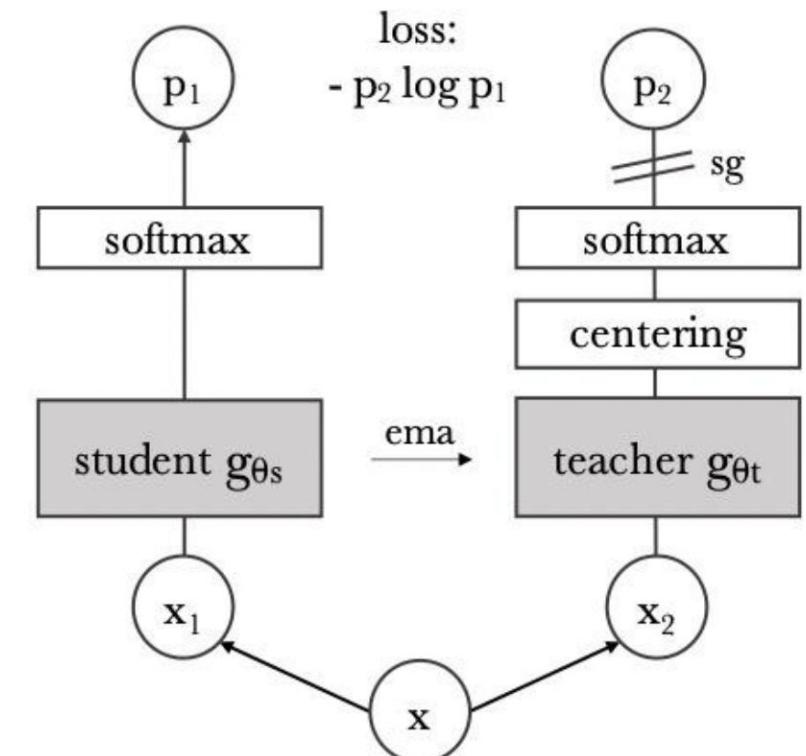
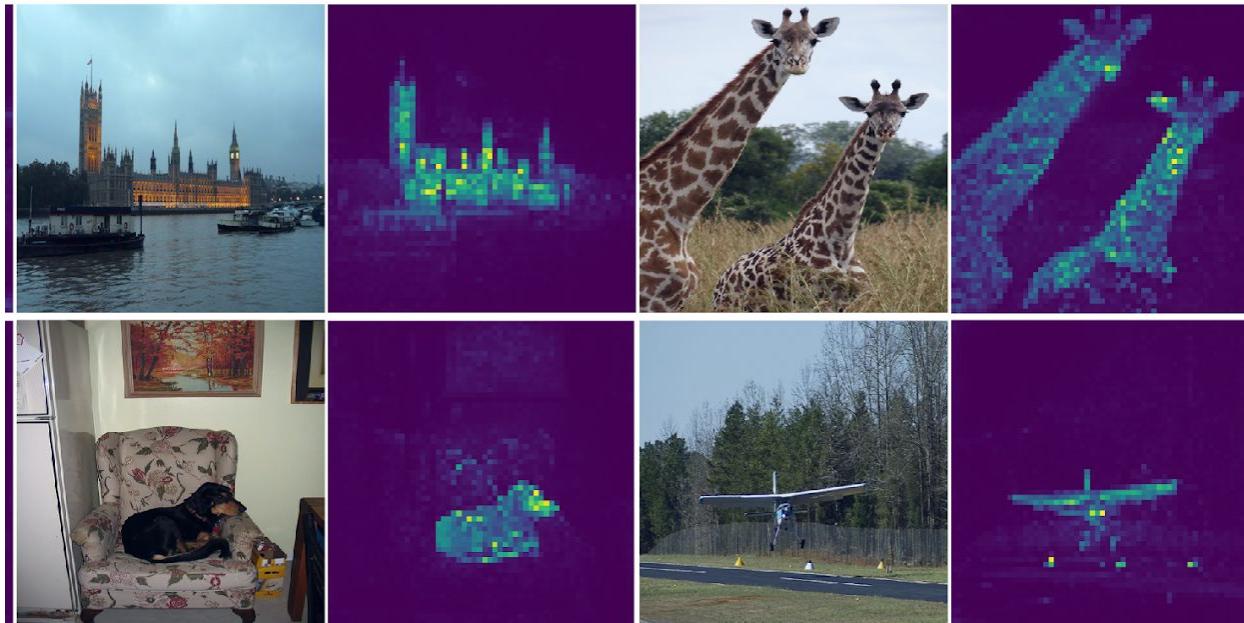
# BYOL: Bootstrap Your Own Latent



## How does BYOL work?

- ▶ There are two networks: an **online network** and a **target network**. Both start out the same.
- ▶ Each network has an encoder (think of it as a feature extractor).
- ▶ The online network has an extra part called the **predictor head**.
- ▶ The target network's parameters are updated slowly using the online network's parameters (using something called an exponential moving average).
- ▶ The goal is simple: make the online network's prediction as close as possible to the target network's output, using mean squared error between their normalized outputs.

# DiNO: DIstillation with NO labels



# **DiNO: DIstillation with NO labels**

## How does DINO work?

- ▶ **Student and Teacher Networks:**

There are two networks with the same design. The teacher is just a slowly updated version of the student (using EMA).

- ▶ **Multi-crop Strategy:**

The model looks at the same image in different ways—two large views and several small crops—to learn features that work at different scales.

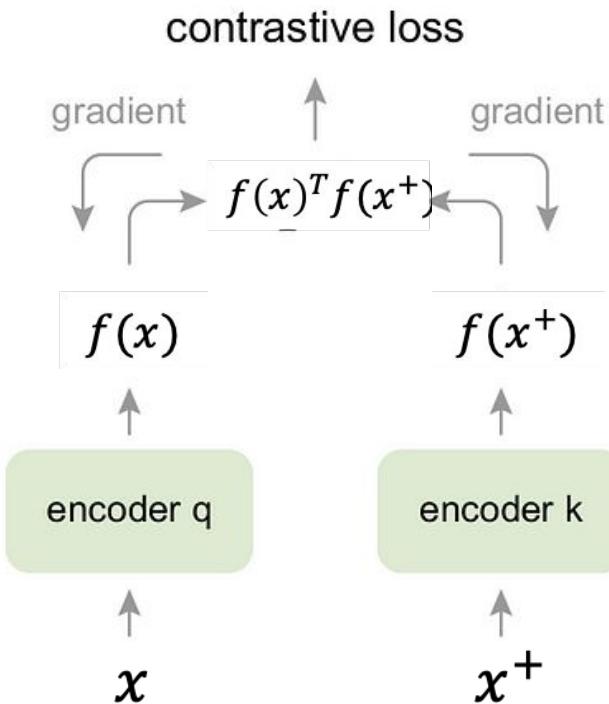
- ▶ **No Negatives Needed:**

Instead of comparing with negative samples, DINO uses a cross-entropy loss on soft similarity scores between the student and teacher outputs.

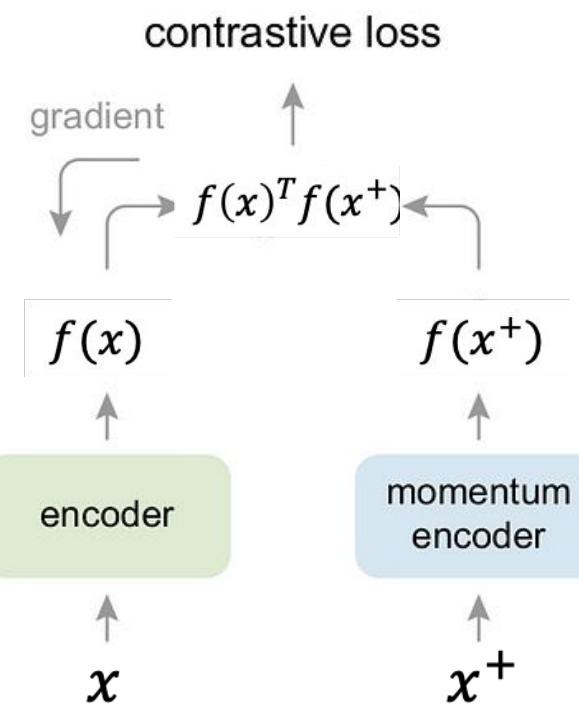
- ▶ **Avoiding Collapse:**

To make sure the model doesn't just output the same thing for every image, DINO centers and sharpens the teacher's outputs.

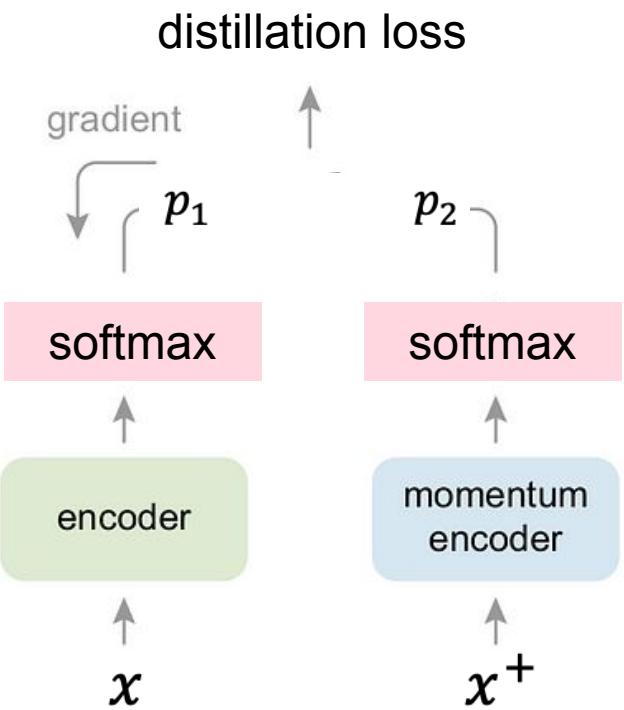
# SimCLR



# MoCo v3



# DiNo



InfoNCE Loss

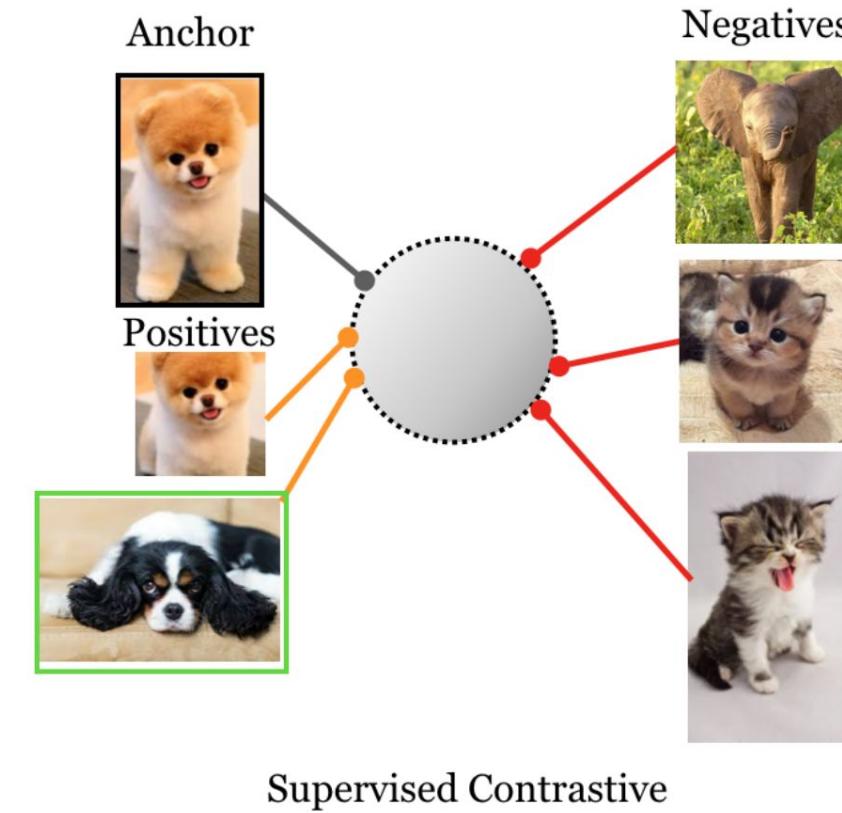
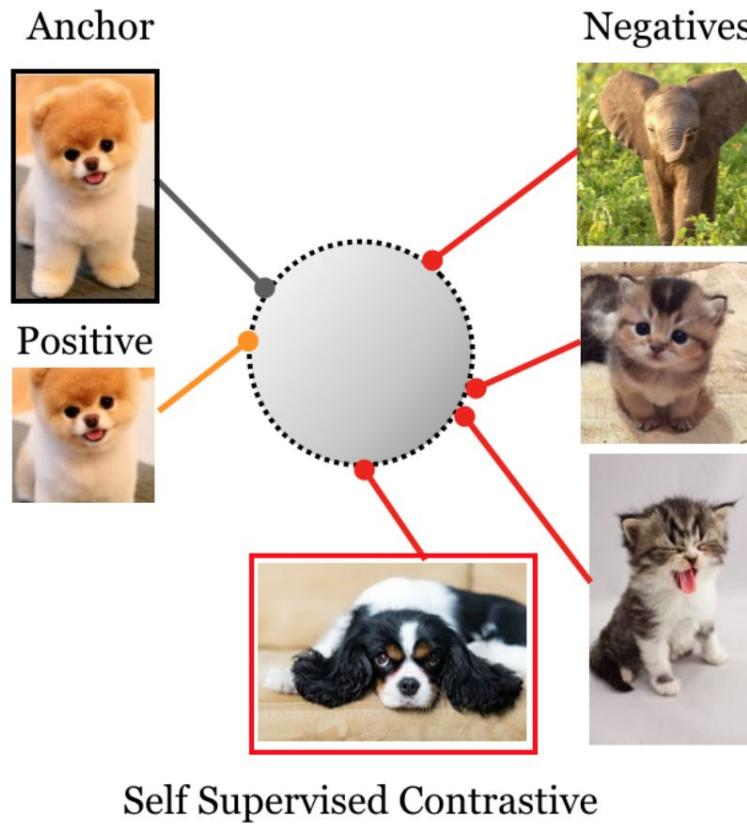
$$\min_f \mathbb{E}_{(x, x^+) \sim p_{pos}, \{x_i^-\}_{i=1}^N \sim p_{data}} \left[ -\log \frac{e^{f(x)^T f(x^+)/\tau}}{e^{f(x)^T f(x^+)/\tau} + \sum_{i=1}^N e^{f(x)^T f(x^-_i)/\tau}} \right]$$

↑ pull positive pair together  
↓ push negative pairs apart

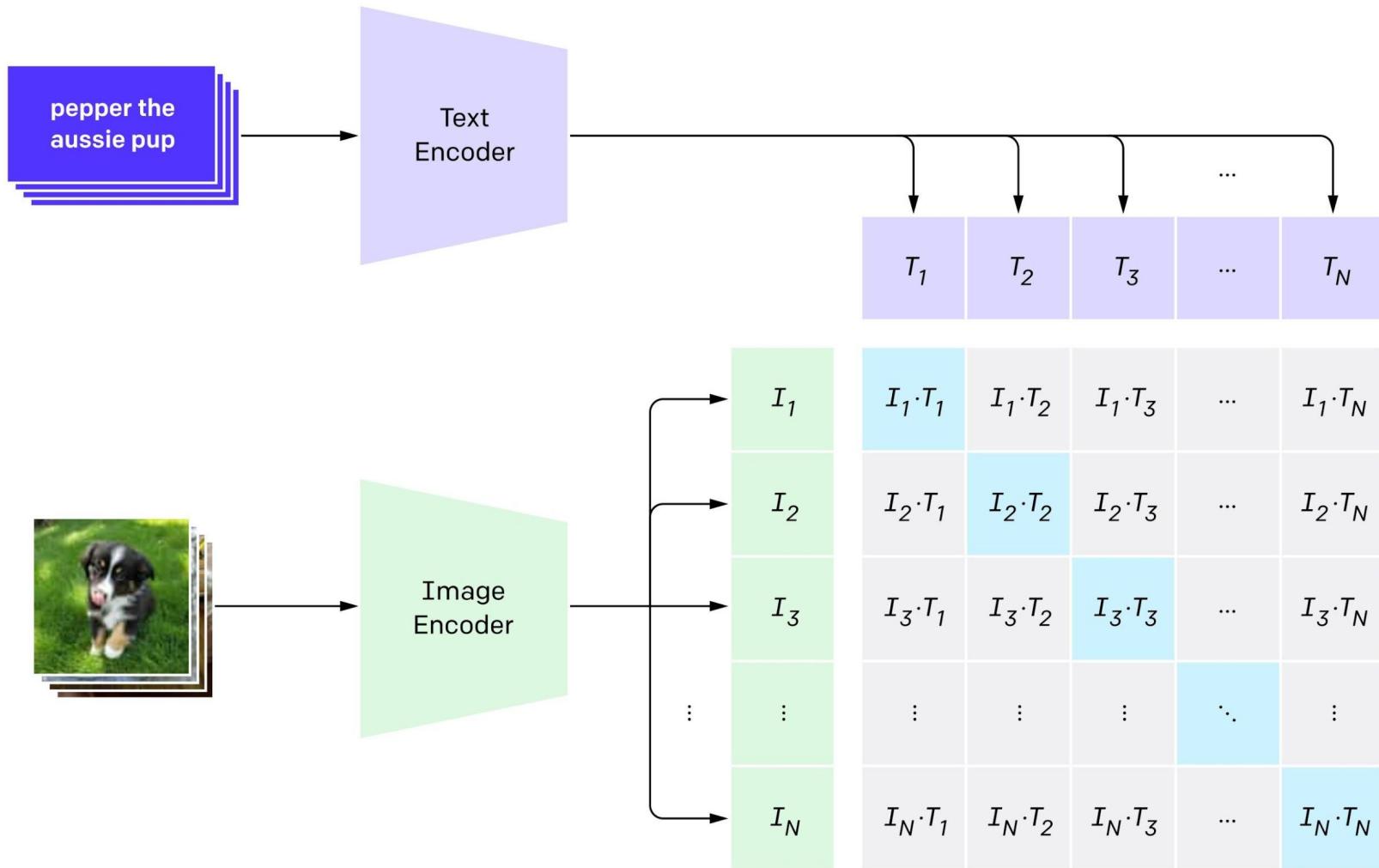
Model Distillation Loss

$$-p_2 \log p_1$$

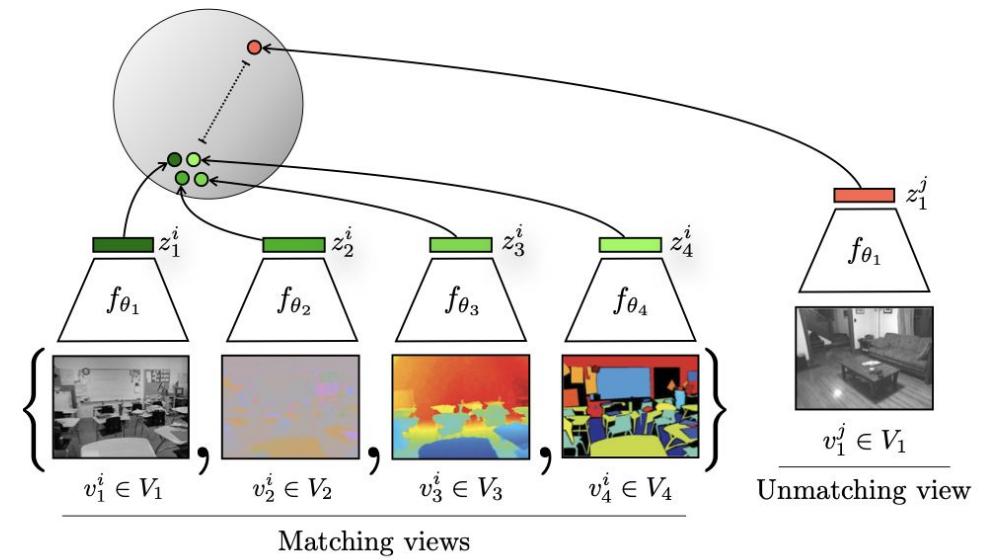
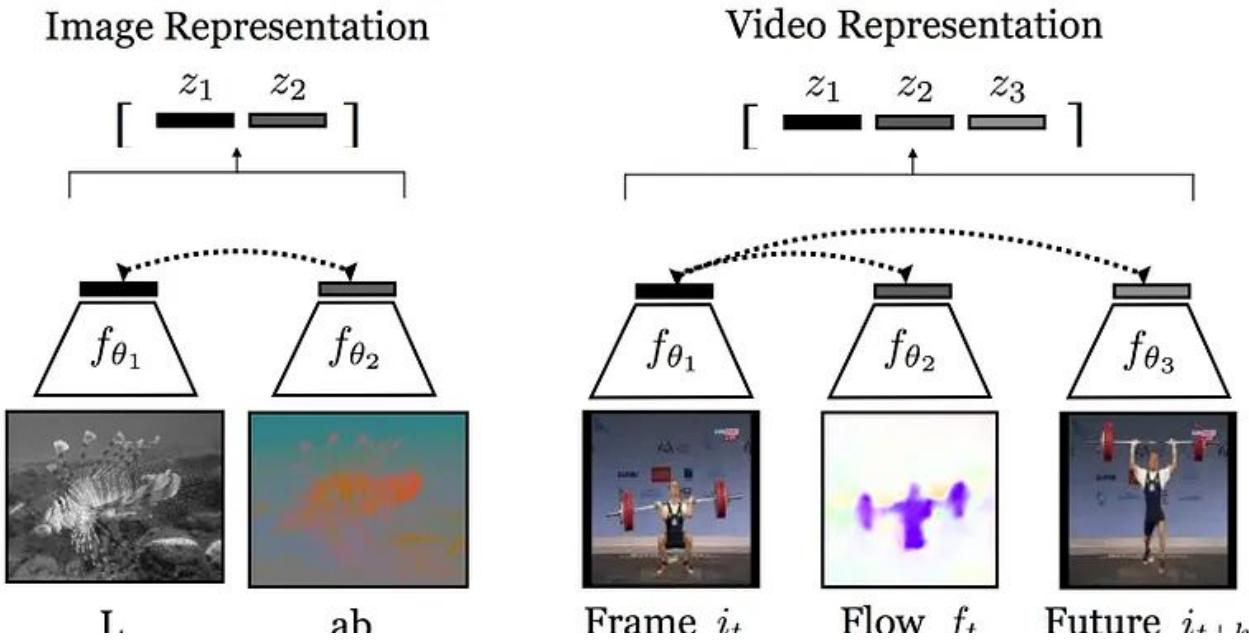
# SupCon: Supervised Contrastive Learning



# CLIP: Contrastive Language Image Pretraining



# CMC: Contrastive Multiview Coding



**Figure 1:** Given a set of sensory views, a deep representation is learnt by bringing views of the *same* scene together in embedding space, while pushing views of *different* scenes apart. Here we show an example of a 4-view dataset (NYU RGBD [53]) and its learned representation. The encodings for each view may be concatenated to form the full representation of a scene.

# LGSimCLR: Language SimCLR

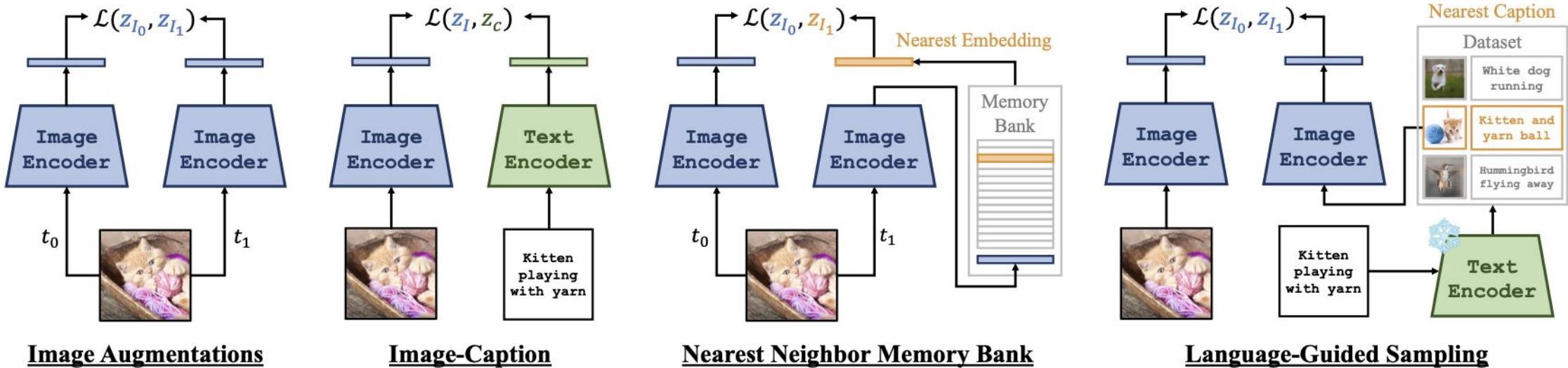


Figure 3. **Contrasting Contrastive Formulations.** While image-only and image-text contrastive learning directly extract views from the instance, nearest-neighbor methods rely on a memory bank of previously extracted features for training. In contrast, our approach samples nearest neighbors in caption embedding space using a pretrained language model and use the associated image for contrastive learning.



# A Unified Framework for Representation Learning

Stochastic Neighbor Embedding

A Simple Framework for Contrastive Learning of Visual Representations

Normalized Cuts and Image Segmentation

Published as a conference paper at ICLR 2023

## Learning Transferable Visual Models From Natural Language Supervision

Geoffrey Hinton  
Department of Computer Science  
10 King's College Road  
[www.cs.toronto.edu/~hinton/roweis](http://www.cs.toronto.edu/~hinton/roweis)

### Abstract

We describe a probabilistic approach by high-dimensional latent space to predict a fixed set of predetermined objects. This approach is restricted by supervision limits their generality and usability since additional labeled data is needed to represent any visual concept. Learning directly from raw text about images is a promising alternative which leverages a much broader source of supervision. We demonstrate that the simple pre-training task of predicting which caption goes with which image is an efficient and scalable way to learn SOTA models from scratch on a dataset of 400 million (image, text) pairs collected from the internet. After pre-training, natural language is used to learn visual concepts (or describe new ones) enabling zero-shot transfer of the model to downstream tasks. We study

Alec Radford<sup>1</sup> Jong Wook Kim<sup>1</sup> Chris Hallacy<sup>1</sup> Aditya Ramesh<sup>1</sup> Gabriel Goh<sup>1</sup> Sandhini Agarwal<sup>1</sup>  
Girish Sastry<sup>1</sup> Amanda Askell<sup>1</sup> Pamela Mishkin<sup>1</sup> Jack Clark<sup>1</sup> Gretchen Krueger<sup>1</sup> Ilya Sutskever<sup>1</sup>

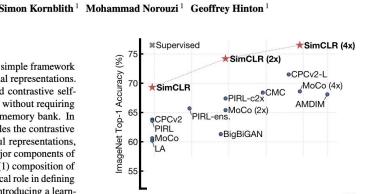


Figure 1. Comparison of ImageNet Top-1 Accuracy (%) on re-ports (pre-ResNet-18) and datasets.

Task-agnostic objectives such as autoregressive and masked language modeling have scaled across many orders of magnitude in compute, model capacity, and data, steadily improving capabilities. The development of “text-to-text” as a standardized input-output interface (McCann et al., 2018; Radford et al., 2019; Raffel et al., 2019) has enabled task-agnostic architectures to zero-shot transfer to downstream datasets removing the need for specialized output heads or dataset specific customization. Flagship systems like GPT-3 (Brown et al., 2020) are now competitive across many tasks with bespoke models while requiring little to no dataset-specific training data.

These results suggest that the aggregate supervision accessible to modern pre-training methods within web-scale collections of text surpasses that of high-quality crowd-labeled NLP datasets. However, in most such cases as captioning, vision is still standard practice to train models on crowd-labeled datasets such as ImageNet (Deng et al., 2009).

Pre-training methods which learn directly from text are a similar breakthrough in computer vision as encouraging.

Mori et al. (1999) explored improving image quality by using a local feature descriptor in documents paired with images. Srihari (2012) explored deep representations for multimodal Deep Boltzmann of low-level image and text tag features, modernizing this line of work and demonstrating its effectiveness.

Our approach scales to any number of views, and is view-agnostic. We analyze key properties of the approach that make it work, finding that the contrastive loss outperforms w/ pre-training and after the resulting semantics. Our open-source code is released at: <https://github.com/yonglongtian/ContrastiveMultiviewCoding>.

### Abstract

Humans view the world through many sensory channels, e.g., the long-wavelength light channel, viewed by the left eye, or the high-frequency channels, heard by the right ear. Each view is noisy and incomplete, but important factors, such as physics, geometry, and semantics, tend to be shared between all views (e.g., a “dog” can be seen, heard, and felt). We investigate the classic hypothesis that a powerful representation is one that models view-invariant factors. We study this hypothesis under the framework of multi-view contrastive learning, where we learn a representation that aims to maximize mutual information between different views of the same scene but is otherwise compact. Our approach scales to any number of views, and is view-agnostic. We analyze key properties of the approach that make it work, finding that the contrastive loss outperforms w/ pre-training and after the resulting semantics. Our open-source code is released at: <https://github.com/yonglongtian/ContrastiveMultiviewCoding>.

## Clustering by Low-Rank Doubly Stochastic Matrix Decomposition

Zhirong Yang [ZHIRONG.YANG@AALTO.FI](mailto:ZHIRONG.YANG@AALTO.FI)  
Department of Information and Computer Science, Aalto University, 00076, Finland

Erkki Oja [ERKKI.OJA@AALTO.FI](mailto:ERKKI.OJA@AALTO.FI)  
Department of Information and Computer Science, Aalto University, 00076, Finland

### Abstract

Contrastive learning, especially self-supervised contrastive learning (SSCL), has achieved great success in extracting powerful features from unlabeled data. In this work, we contribute to the theoretical understanding of SSCL and propose a generalization to the classic data visualization method, stochastic neighbor embedding (SNE (Hinton & Roweis, 2002)), whose goal is to preserve pairwise distances. From the perspective of preserving neighboring information, SSCL can be viewed as a special case of SNE with the input space pairwise similarities specified by data augmentation. The established correspondence facilitates deeper theoretical understanding of learned features of SSCL, as well as methodological guidelines for practical improvement. Specifically, through the lens of SNE, we provide novel analysis on domain-agnostic augmentations, implicit bias and robustness of learned features. To illustrate the practical advantage, we demonstrate that the modifications from SNE to t-SNE (Van de Maaten & Hinton, 2008) can also be adopted in the SSCL setting, achieving significant improvement in both in-distribution and out-of-distribution generalization.

## ABSTRACT

Contrastive learning, especially self-supervised contrastive learning (SSCL), has achieved great success in extracting powerful features from unlabeled data. In this work, we contribute to the theoretical understanding of SSCL and propose a generalization to the classic data visualization method, stochastic neighbor embedding (SNE (Hinton & Roweis, 2002)), whose goal is to preserve pairwise distances. From the perspective of preserving neighboring information, SSCL can be viewed as a special case of SNE with the input space pairwise similarities specified by data augmentation. The established correspondence facilitates deeper theoretical understanding of learned features of SSCL, as well as methodological guidelines for practical improvement. Specifically, through the lens of SNE, we provide novel analysis on domain-agnostic augmentations, implicit bias and robustness of learned features. To illustrate the practical advantage, we demonstrate that the modifications from SNE to t-SNE (Van de Maaten & Hinton, 2008) can also be adopted in the SSCL setting, achieving significant improvement in both in-distribution and out-of-distribution generalization.

## CLUSTERING AND A OF MULTIVARIATE OBSERVATIONS

J. MACQUEEN  
UNIVERSITY OF CALIFORNIA, LOS ANGELES

### 1. Introduction

The main purpose of this paper is to describe a  $p$ -dimensional population into  $k$  sets on the basis of which is called “ $k$ -means,” appears to give partition in the sense of within-class variance. That is, if function for this population,  $S = \{S_1, S_2, \dots, S_k\}$  is a  $i = 1, 2, \dots, k$ , is the conditional mean of  $p$  over  $\Omega$ .  $\sum_{j=1}^k \|S_j - u_i\|^2 dp_j$  tends to be low for the part of the part method. We say “tends to be low,” primarily because of corroborated to some extent by mathematical analytical experience. Also, the  $k$ -means procedure is computationally economical, so that it is feasible to program on a digital computer. Possible applications include grouping, nonlinear prediction, approximating multivariate parametric tests for independence among several variables. In addition to suggesting practical classification method has proved to be theoretically interesting. The  $k$ -means generalization of the ordinary sample mean, and one is naturally led to study the pertinent asymptotic behavior, the object being to establish some sort of law of large numbers for the  $k$ -means. This problem is sufficiently interesting, in fact, for us to devote a good portion of this paper to it. The  $k$ -means are defined in section 2.1, and the main results which have been obtained on the asymptotic behavior are given there. The rest of section 2 is devoted to the proofs of these results. Section 3 describes several specific possible applications, and reports some preliminary results from computer experiments conducted to explore the properties inherent in the  $k$ -means idea. The extension to general metric spaces

negativity constraint, together with various low-rank matrix approximation objectives, has widely been used for the relaxation purpose in the past decade.

The most popular nonnegative low-rank approximation method is Nonnegative Matrix Factorization (NMF). It finds a matrix that approximates the similarities and can be factorized into several nonnegative low-rank matrices. NMF was originally applied to vectorial data, where Ding et al. (2010) have shown that NMF is equivalent to the classical  $k$ -means method. Later NMF was applied to the (weighted) graph given by the pairwise similarities. For example, Ding et al. (2008) presented Nonnegative Spectral Cuts by using a multiplicative algorithm; Arora et al. (2011) proposed Latent Stochastic Decomposition that approximates a similarity matrix based on Euclidean distance and a left-stochastic matrix. Another stream in the same direction is topic modeling. Hofmann (1999) gave a generative model in Probabilistic Latent Semantic Interpreting (PLSI) for counting data, which is essentially equivalent to NMF using Kullback-Leibler (KL) divergence and Tri-factorizations. Bayesian treatment of PLSI by using Dirichlet prior was introduced by Blei et al. (2001). Symmetric PLSI with the same Bayesian treatment is called Interaction Component Model (ICM) (Sinkkonen et al., 2008).

Despite remarkable progress, the above relaxation approaches are still not fully satisfactory in all of the fol-

## Normalized Cuts and Image Segmentation

Jianbo Shi and Jitendra Malik, Member, IEEE

Abstract—We propose a novel approach to solving the perceptual grouping problem in vision. Rather than focusing on local features and their connections in the image data, our approach aims at extracting the global impression of an image. We treat image segments as a group of pixels with global connections, the normalized cut, for segmenting the graph. The normalized cut measures both the total dissimilarity between the different groups as well as the total similarity within the groups. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

encouraging. We have applied this approach to segmenting static or

# A Unified Framework for Representation Learning

- **I-Con:** A single information-theoretic loss that unifies several major classes of representation learning.
- **Periodic Table of Methods:** I-Con enables a structured organization of diverse learning
- **Filling the Gaps:** I-Con helps derive a new method that boosts unsupervised ImageNet-1K accuracy by **+8%**.

Supervisory Signal

	Gaussian	Student-T	Identity	Graph Kernel Weights	Uniform over K-Neighbors	Uniform over Positive Pairs	Cross-Modal Pairs	Uniform over Classes	Data-Label Pairs
Learned Representations	Gaussian	SNE [Hinton 2002]	Dual t-SNE	SNE Graph Embeddings	SNE with Uniform Affinities	InfoNCE [Bachman 2019]	CLIP [Radford 2021]	SupCon [Khosla 2020]	Cross Entropy [Good 1963]
	X-Sample CL [Sobal 2025]				LGSimCLR [El Banani 2023]	SimCLR [Chen 2020]			
	Gaussian $\sigma \rightarrow \infty$			PCA [Pearson 1901]		MoCoV3 [Chen 2021]	CMC [Tian 2020]		
	Gaussian $\sigma \rightarrow 0$					VI-Reg [Bardes 2021]	Average Margin CLIP	Average Margin SupCon	
	Student-T	t-SNE [Van der Maaten 2008]	Doubly t-SNE	t-SNE Graph Embedding	t-SNE with Uniform Affinities	Triplet Loss [Schroff 2015]	Triplet CLIP	Triplet SupCon	Error rate
Dimensionality Reduction Clusters	Dimensionality Reduction Clusters [MacQueen 1967]	Cluster K-Means Learning	Unimodal Normalization [ShSSL10]	Multimodal DCL [Yao 2012]	NCE [Our]	t-CLIP	t-SupCon	Harmonic Loss [Baek 2025]	
Interpretation of Gaps								Supervised Learning	

# A Unified Framework for Representation Learning

- **I-Con:** A single information-theoretic loss that unifies several major classes of representation learning.
- **Periodic Table of Methods:** I-Con enables a structured organization of diverse learning
- **Filling the Gaps:** I-Con helps derive a new method that boosts unsupervised ImageNet-1K accuracy by **+8%**.

Supervisory Signal

	Gaussian	Student-T	Identity	Graph Kernel Weights	Uniform over K-Neighbors	Uniform over Positive Pairs	Cross-Modal Pairs	Uniform over Classes	Data-Label Pairs
Learned Representations	Gaussian	SNE [Hinton 2002]	Dual t-SNE	SNE Graph Embeddings	SNE with Uniform Affinities	InfoNCE [Bachman 2019]	CLIP [Radford 2021]	SupCon [Khosla 2020]	Cross Entropy [Good 1963]
		X-Sample CL [Sobal 2025]			LGSimCLR [El Banani 2023]	SimCLR [Chen 2020]	CMC [Tian 2020]		
	Gaussian $\sigma \rightarrow \infty$			PCA [Pearson 1901]		VI-Reg [Bardes 2021]	Average Margin CLIP	Average Margin SupCon	
	Gaussian $\sigma \rightarrow 0$					Triplet Loss [Schroff 2015]	Triplet CLIP	Triplet SupCon	Error rate
Dimensionality Reduction Clusters	Student-T	t-SNE [Van der Maaten 2008]	Doubly t-SNE	t-SNE Graph Embedding	t-SNE with Uniform Affinities	t-SimCNE [Böhm 2023]	t-CLIP	t-SupCon	Harmonic Loss [Baek 2025]
						t-SimCLR [Hu 2023]			
Legend:									
Dimensionality Reduction Clusters			Cluster K-Means Learning		Normalizing	Unimodal [ShSSL10]	DCL	Multimodal NCE	Supervised Learning
Interpretation of Clusters			t-SNE Clustering [Ours]		SGCL Clustering		Supervised Clustering		Interpretation of Gaps

# A Unified Framework for Representation Learning

- **I-Con:** A single information-theoretic loss that unifies several major classes of representation learning.
- **Periodic Table of Methods:** I-Con enables a structured organization of diverse learning
- **Filling the Gaps:** I-Con helps derive a new method that boosts unsupervised ImageNet-1K accuracy by +8%.

Supervisory Signal

	Gaussian	Student-T	Identity	Graph Kernel Weights	Uniform over K-Neighbors	Uniform over Positive Pairs	Cross-Modal Pairs	Uniform over Classes	Data-Label Pairs
Learner	SNE [Hinton 2002]	X-Sample CL [Sobal 2025]	Dual t-SNE	SNE Graph Embeddings	SNE with Uniform Affinities	InfoNCE [Bachman 2019]	CLIP [Radford 2021]	SupCon [Khosla 2020]	Cross Entropy [Good 1963]
ended					LGSimCLR [El Banani 2023]	SimCLR [Chen 2020]	CMC [Tian 2020]		
representation					VILoss [Barde 2021]	Average Margin CLIP	Average Margin SupCon		
Dimensionality Reduction						Triplet Loss [Schroff 2015]	Triplet CLIP	Triplet SupCon	Error rate
Clusters	t-SNE [Van der Maaten 2008]	Doubly t-SNE		t-SNE Graph Embedding	t-SNE with Uniform Affinities	t-SimCNE [Böhm 2023]	t-CLIP	t-SupCon	Harmonic Loss [Baek 2025]
Interpretation of Concepts						t-SimCLR [Hu 2023]			

$\mathcal{L} = \int_i KL(p(i | j) || q(i | j)) di$

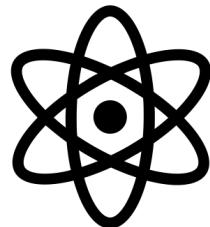
~Rows      ~Columns

**Legend:**

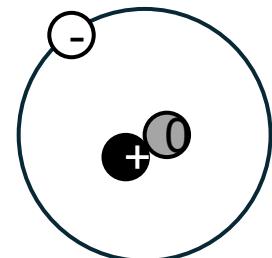
- Dimensionality Reduction Clusters [MacQueen 1967]
- Cluster Learning t K-Means [MacQueen 1967]
- Unimodal Clustering [SHCS 19]
- Multimodal NCE [Yao 2012]
- Supervised Learning [SupCon]
- Supervised Clustering [CLIP]
- Interpretation of Concepts [Goodfellow et al. 2015]

# Chemical Periodic Table

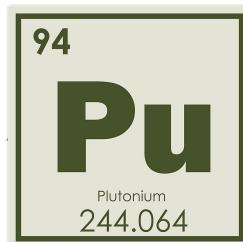
Element



Protons and Neutrons

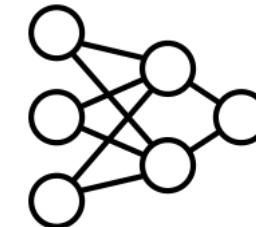


Gaps in the Periodic Table

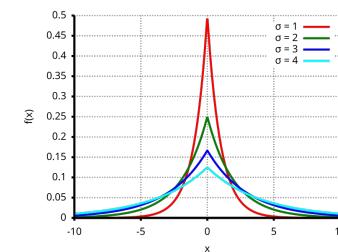
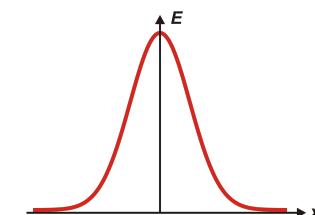


# Machine Learning Periodic Table

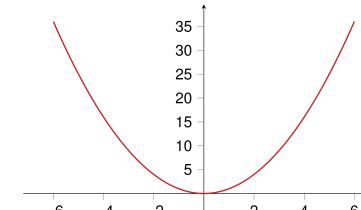
Algorithm



Distributions



New Loss Functions



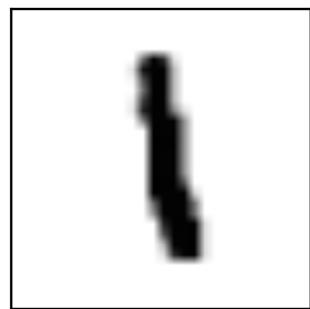
# Overview

- Intro
  - What's Representation Learning?
  - The Periodic Table
- Methods:
  - **Exploring some Examples**
    - Generalizing ML Methods with a single Equation
    - Building the Periodic Table
- Experiments
  - Building an Image Classifier that doesn't need human labels
- Future Directions

# Example 1: Dimensionality Reduction

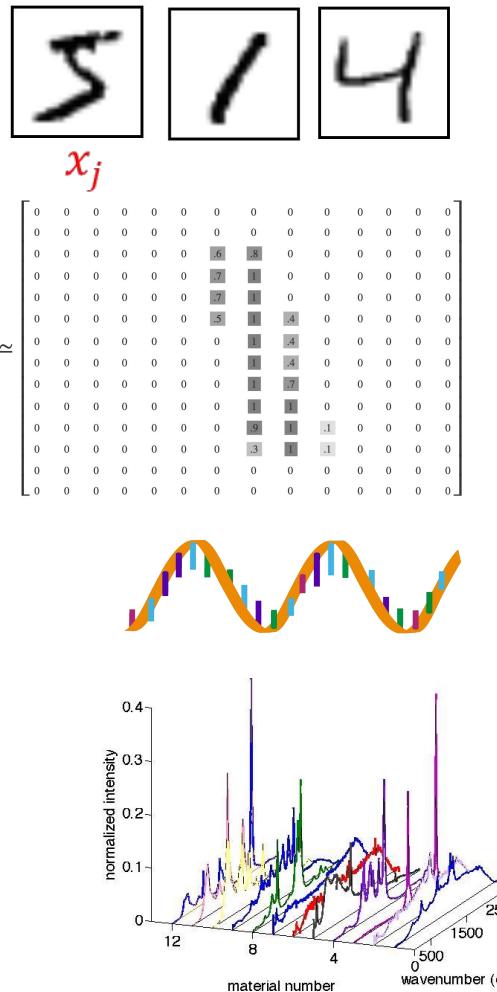
# High Dimensional Data

# RGB Images



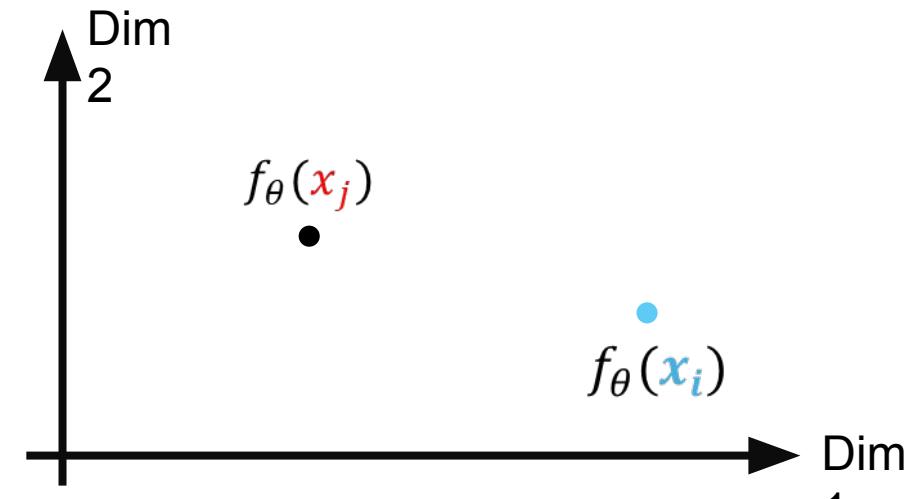
$x_i$   
mRNA-se  
q

# Spectroscopy Data



A blue rectangular box representing a neural network layer. The text "Mapper/LookUp" is centered inside the box. Below it, the mathematical expression  $f_\theta$  is written. Two black arrows, one on each side, point towards the center of the box.

# Low Dimensional Data



PCA

MDS

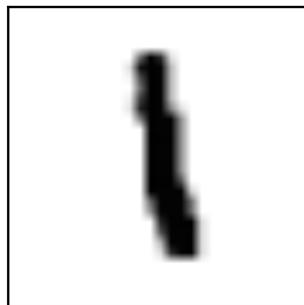
# tSNE

# Example 1: Dimensionality Reduction

PCA

# High Dimensional Data

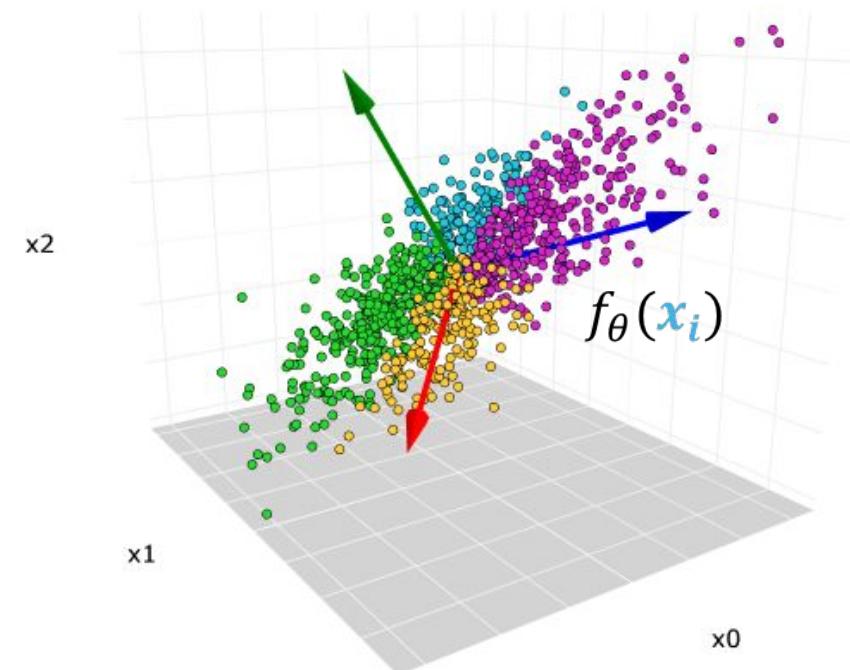
# RGB Images



$x_i$

# Projection Map

# Low Dimensional Data



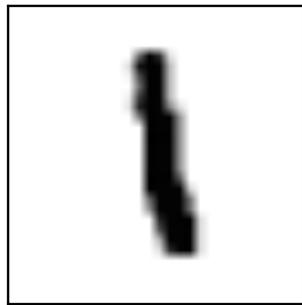
$$\mathcal{L} = -Var(f_\theta(X)) = -\sum_i \left( f_\theta(x_i) - f_\theta(x_j) \right)^2$$

where  $f_\theta$  is a linear projection map

# Example 1: Dimensionality Reduction

# High Dimensional Data

# RGB Images



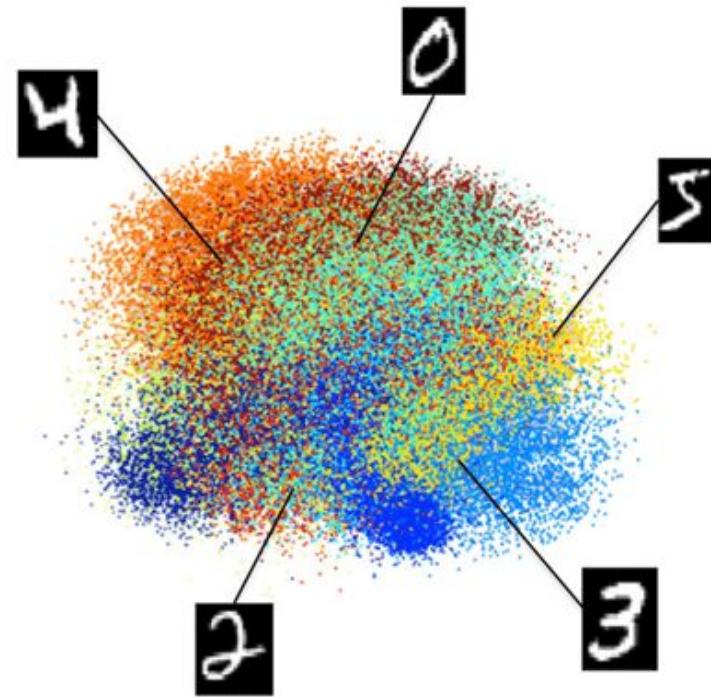
$x_i$

## Projection Map

Map

$$f_\theta$$

# Low Dimensional Data



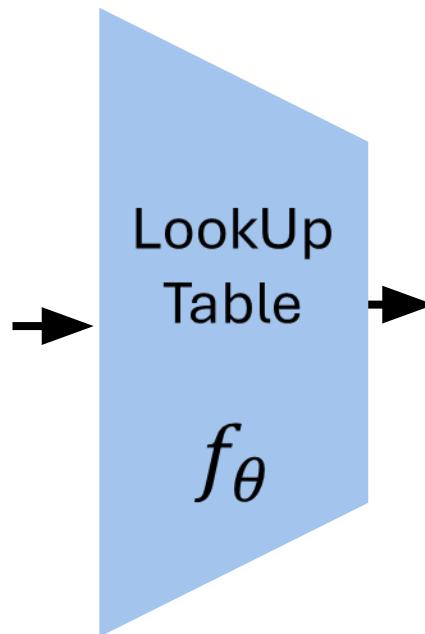
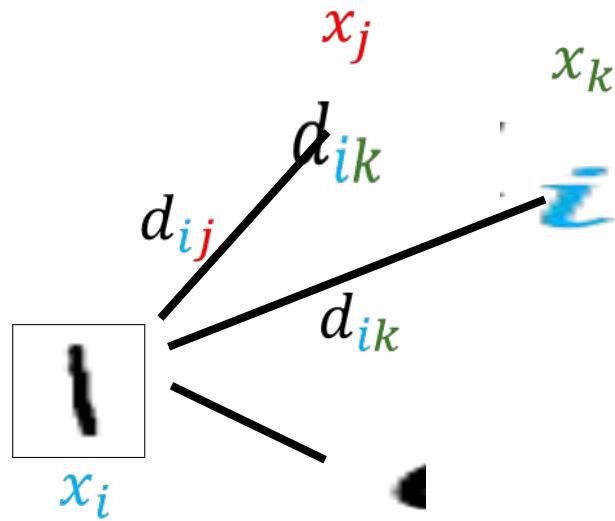
$$\mathcal{L} = -Var(f_\theta(X)) = -\sum_i \left( f_\theta(x_i) - f_\theta(x_j) \right)^2$$

where  $f_\theta$  is a linear projection map

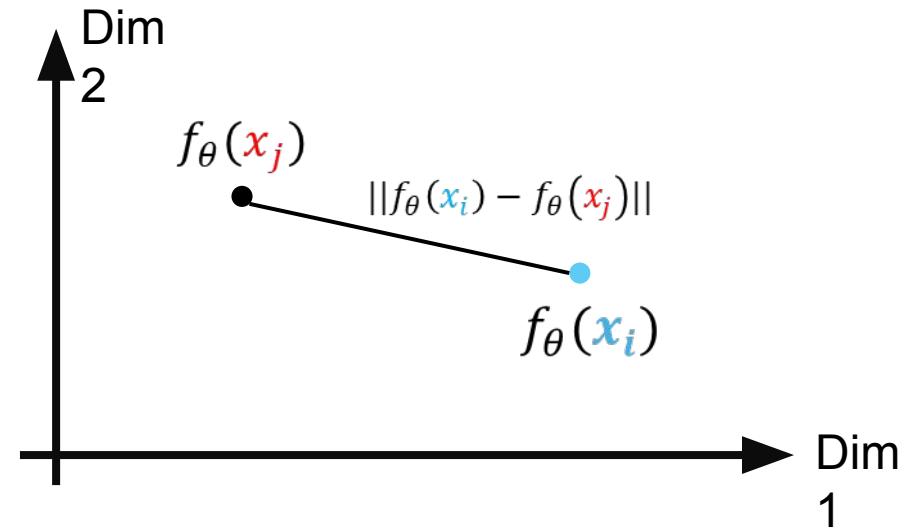
# Example 1: Dimensionality Reduction

MDS

High Dimensional Data



Low Dimensional Data



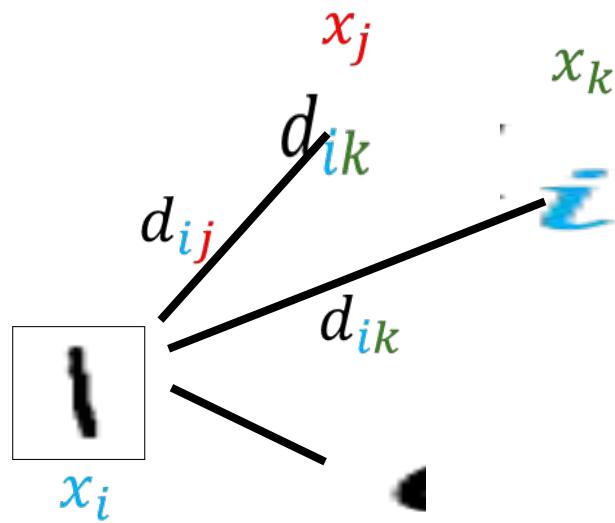
$$\mathcal{L} = \sum_i ( \|f_\theta(x_i) - f_\theta(x_j)\| - d_{ij} )^2$$

where  $f_\theta$  is a lookup table

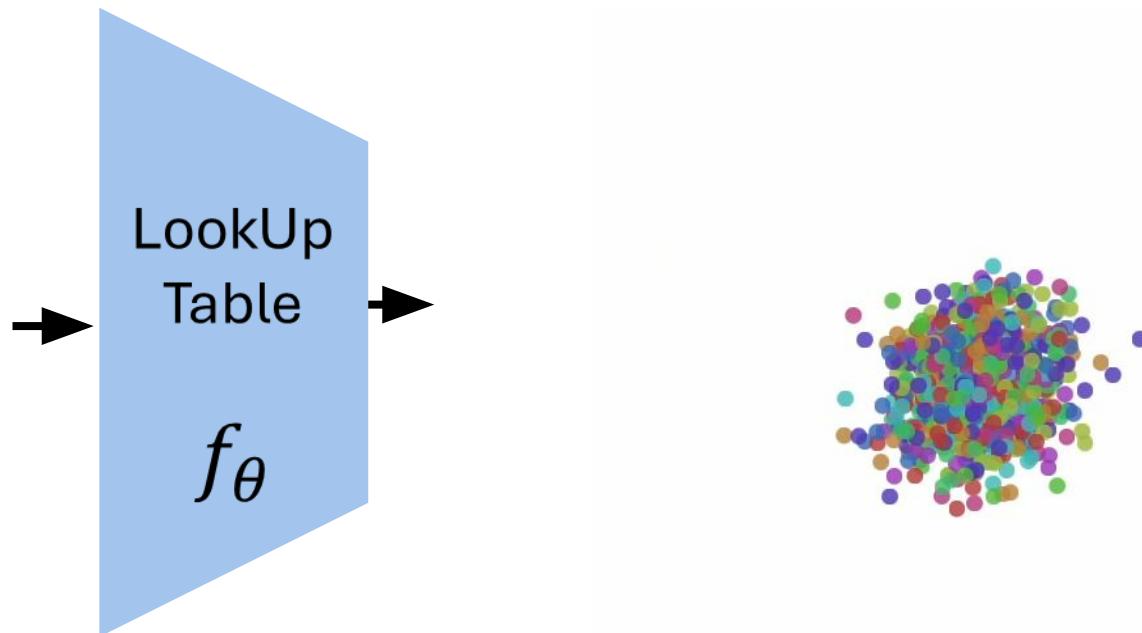
# Example 1: Dimensionality Reduction

MDS

High Dimensional Data



Low Dimensional Data



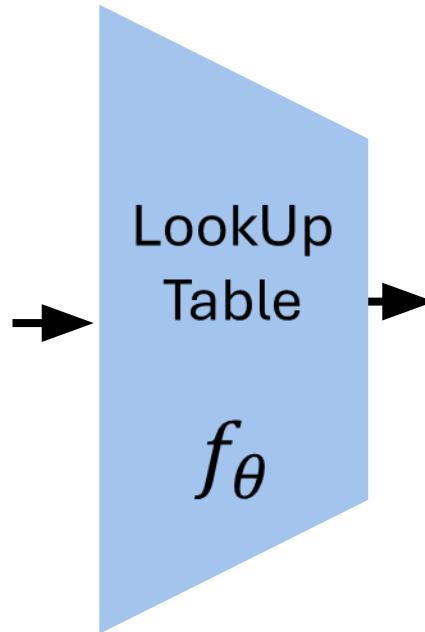
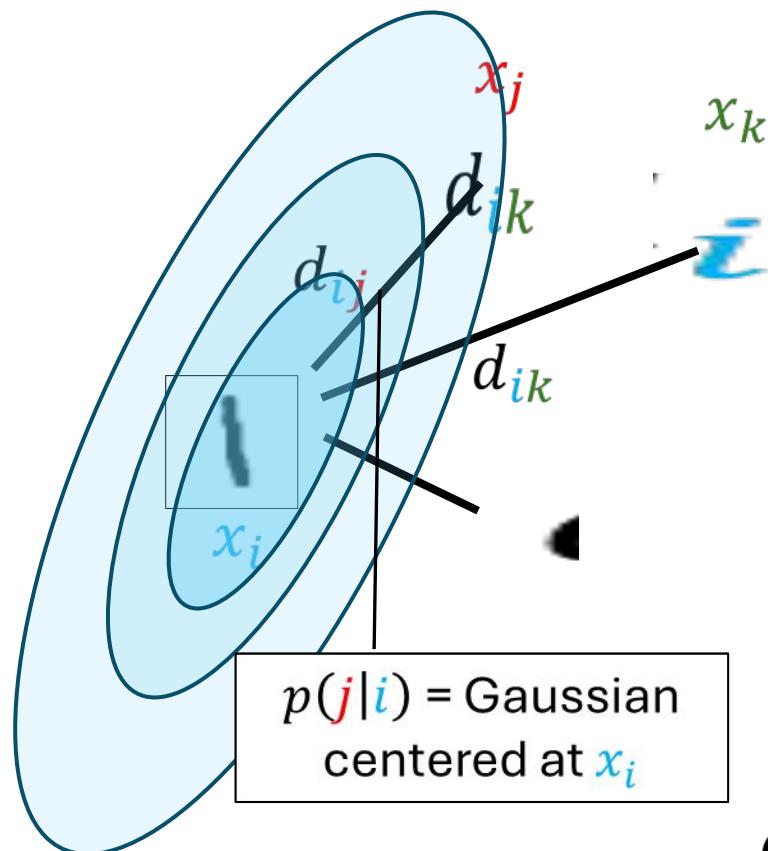
$$\mathcal{L} = \sum_i (||f_\theta(x_i) - f_\theta(x_j)|| - d_{ij})^2$$

where  $f_\theta$  is a lookup table

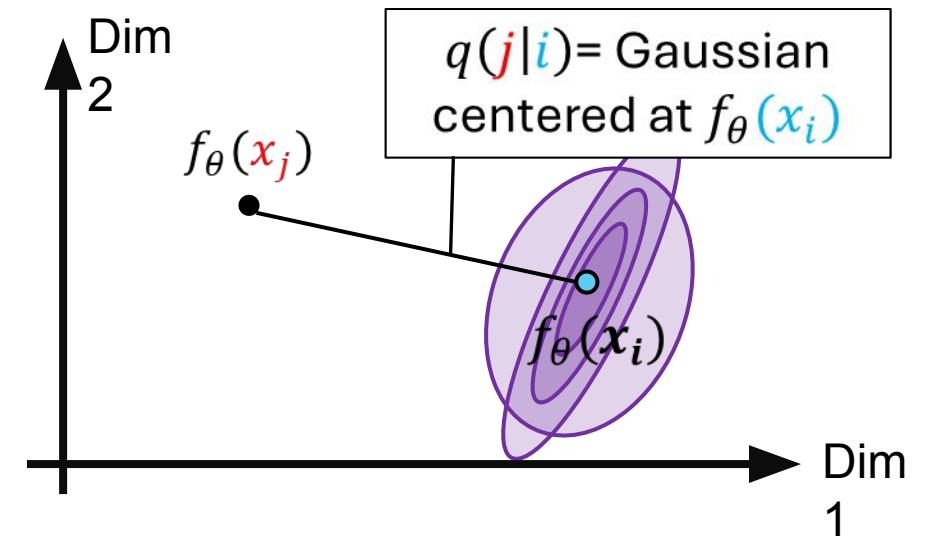
# Example 1: Dimensionality Reduction

SNE

High Dimensional Data



Low Dimensional Data

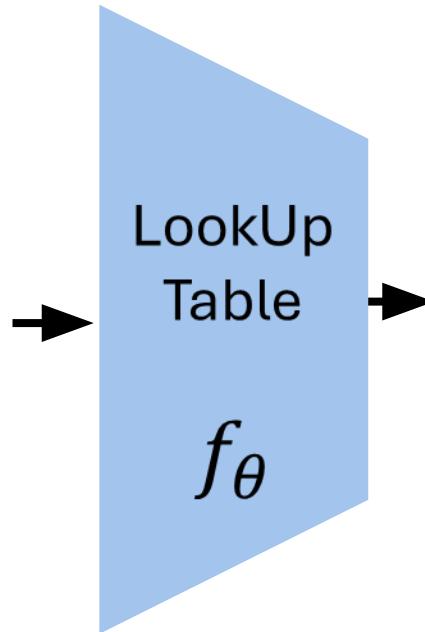
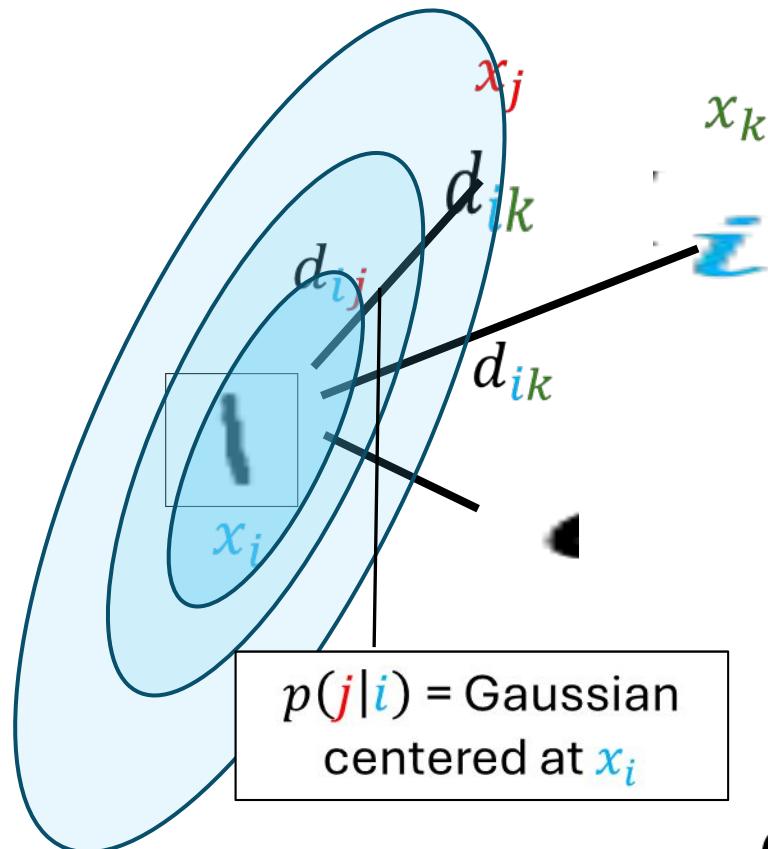


$$\mathcal{L} = \sum_i D_{KL}(p(\cdot | i) \parallel q(\cdot | i))$$

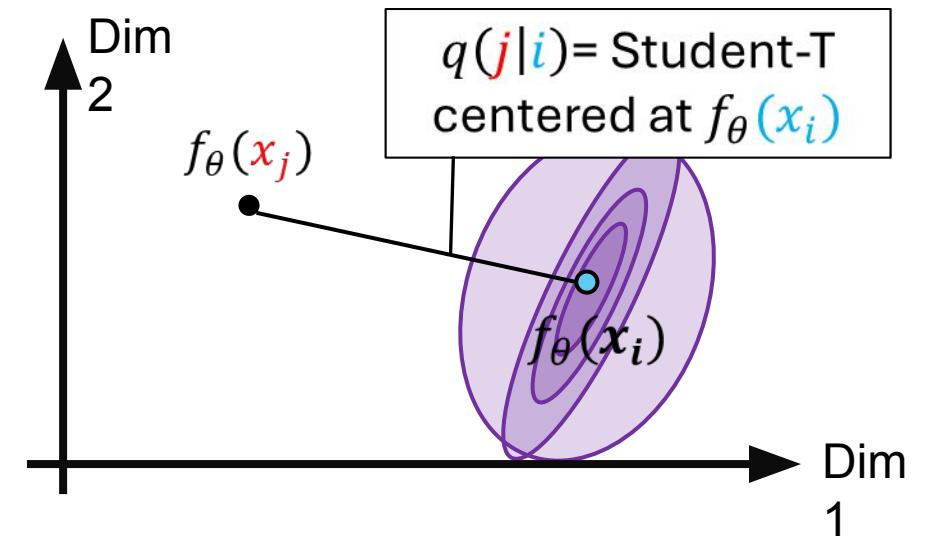
# Example 1: Dimensionality Reduction

t-SNE

High Dimensional Data



Low Dimensional Data

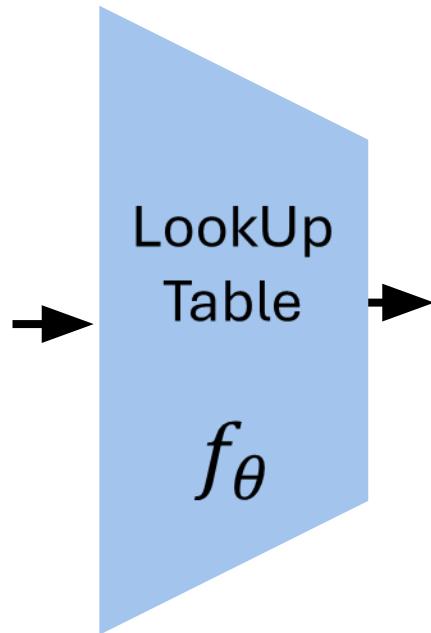
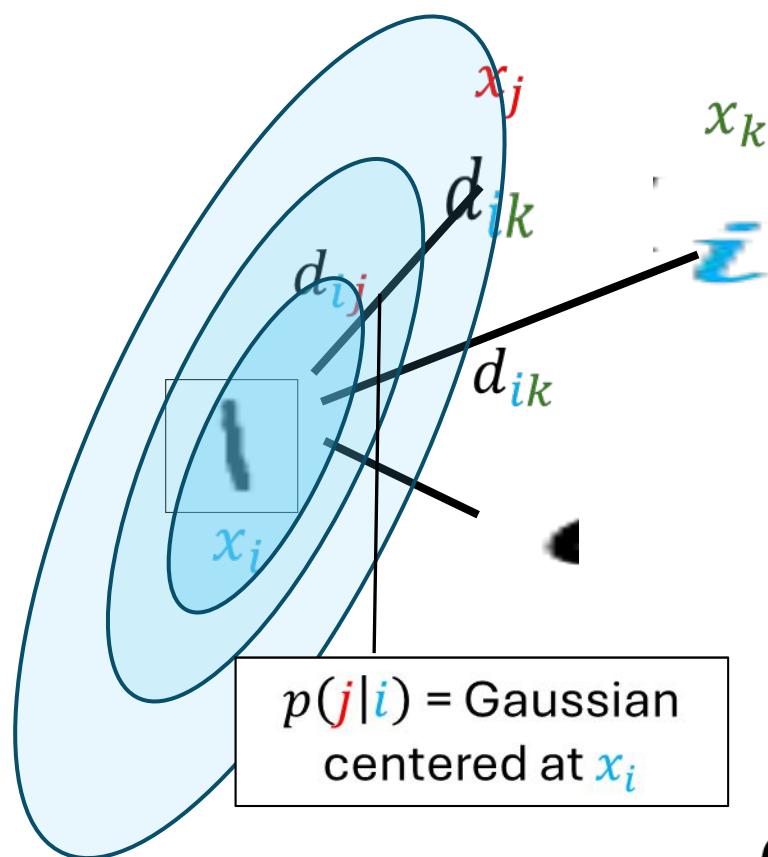


$$\mathcal{L} = \sum_i D_{KL}(p(\cdot | i) \parallel q(\cdot | i))$$

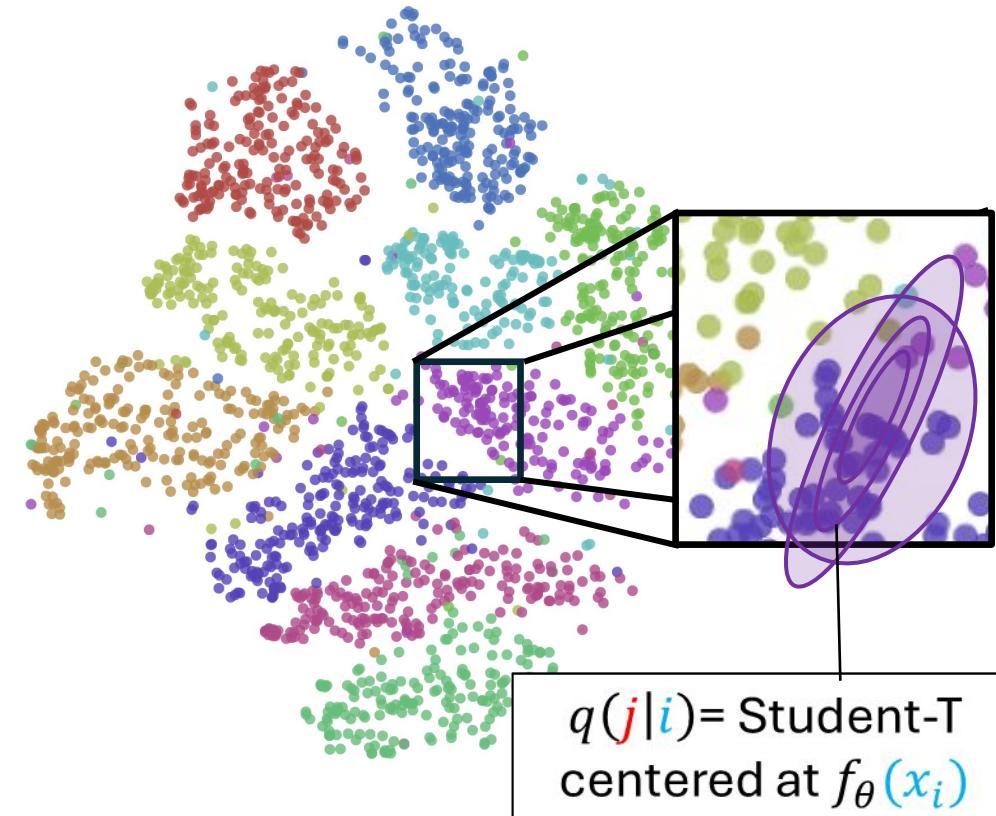
# Example 1: Dimensionality Reduction

t-SNE

High Dimensional Data



Low Dimensional Data



$$\mathcal{L} = \sum_i D_{KL}(p(\cdot | i) \parallel q(\cdot | i))$$

# Example 2: Contrastive Learning

SimCLR

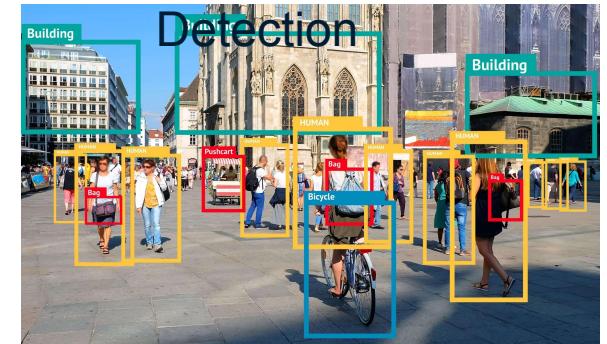


Downstream Tasks

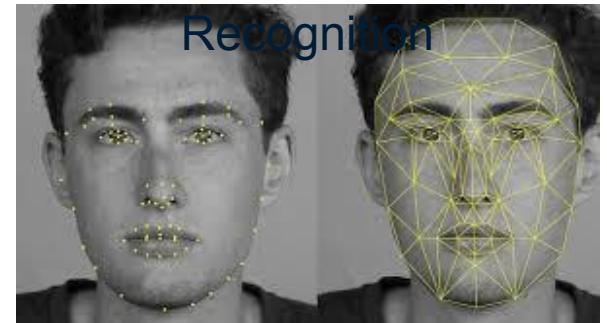
Segmentation



Object Detection



Face Recognition



# Example 2: Contrastive Learning

SimCLR

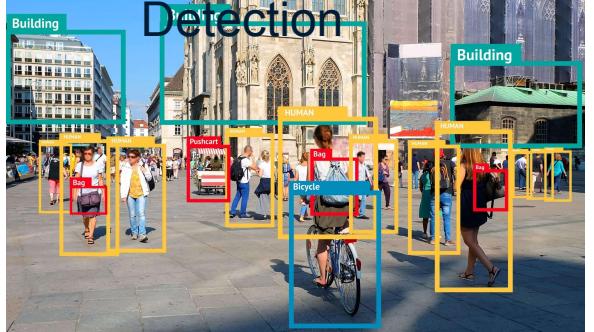
Downstream Tasks



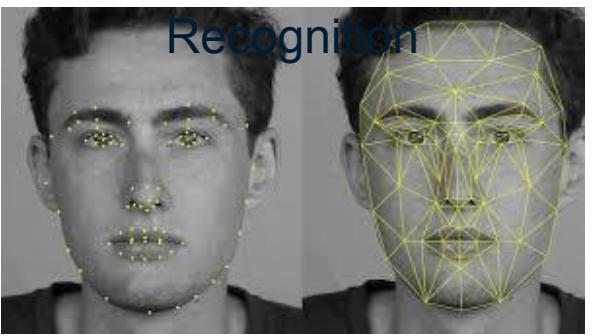
...



Segmentation



Object  
Detection

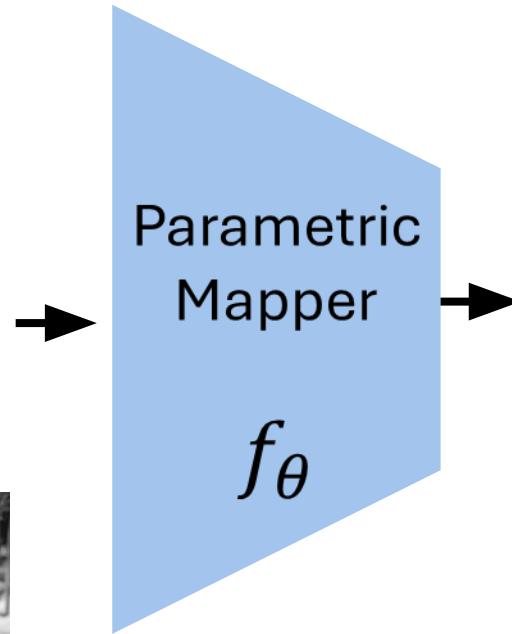
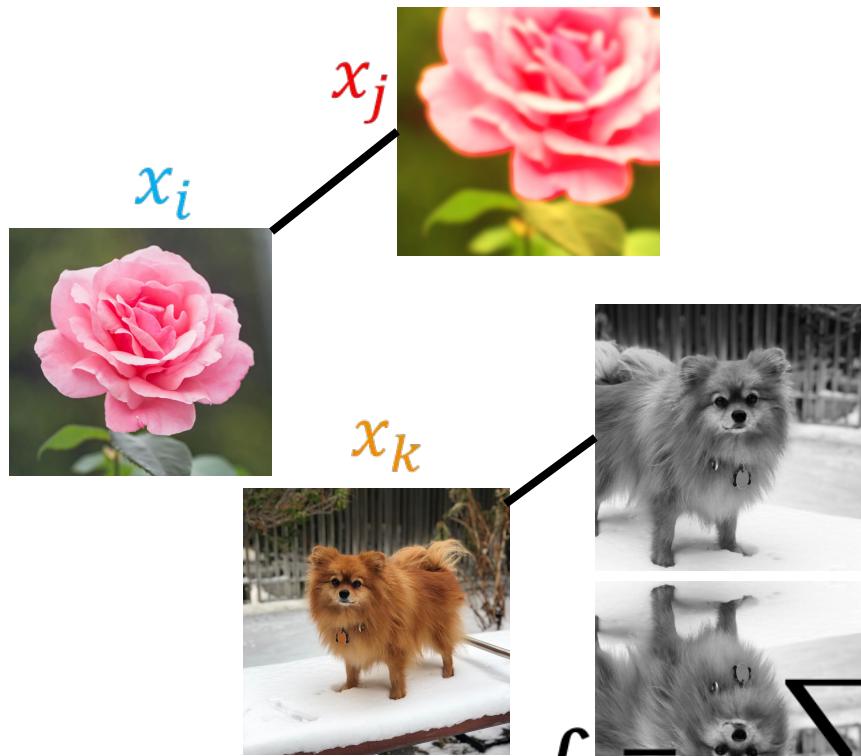


Face  
Recognition

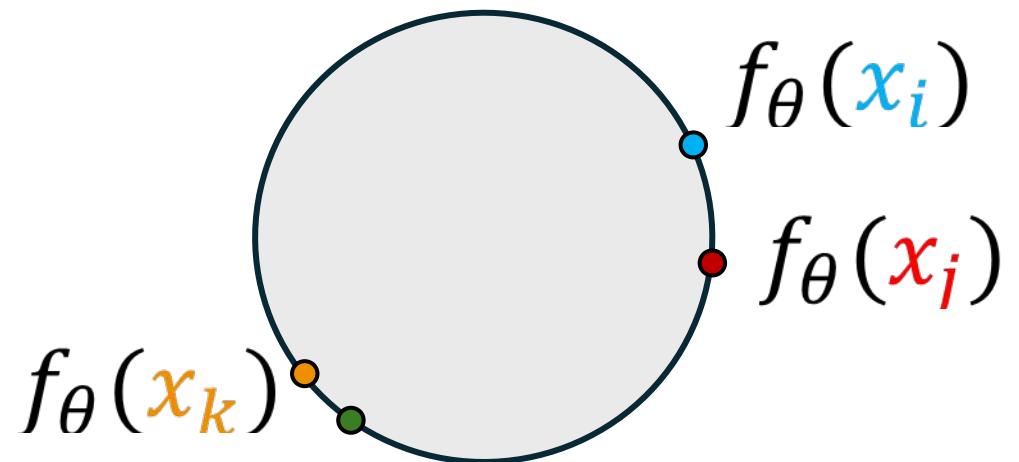
# Example 2: Contrastive Learning

SimCLR

RGB Images



Embeddings



$$\mathcal{L} = -\sum_i \log \frac{\exp(f_\theta(x_i) \cdot f_\theta(x_j)/t)}{\sum_k \exp(f_\theta(x_i) \cdot f_\theta(x_k)/t)}$$

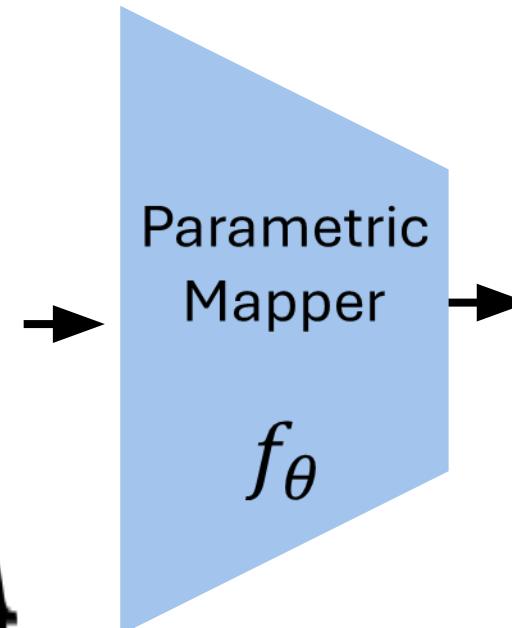
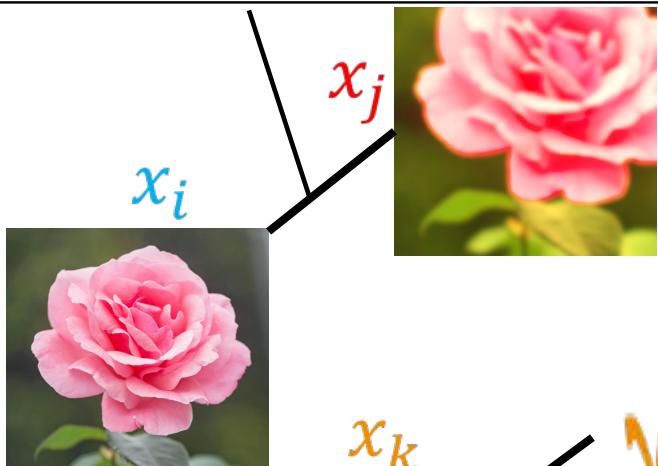
InfoNCE Loss

# Example 2: Contrastive Learning

SimCLR

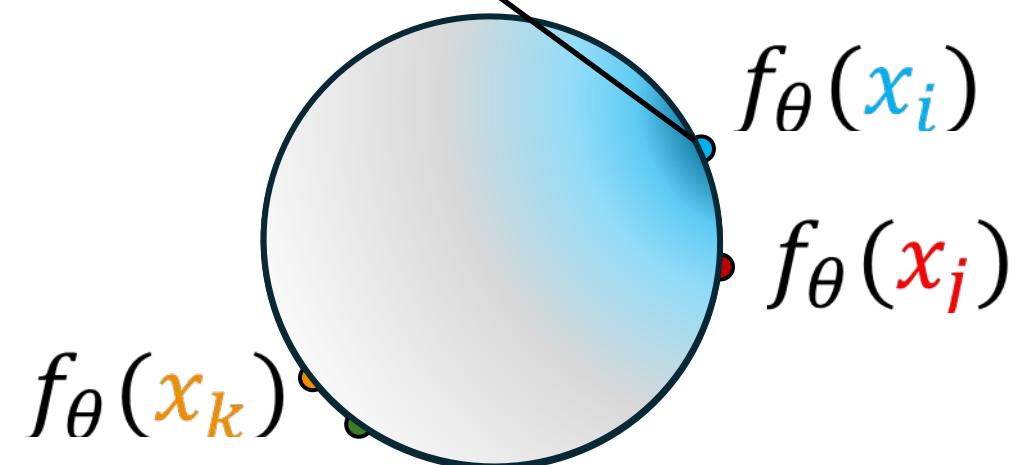
RGB Images

$p(j|i) = \mathbb{1}[i \text{ and } j \text{ are augmentations of the same image}]$



Embeddings

$q(j|i) = \text{Gaussian centered at } f_\theta(x_i)$

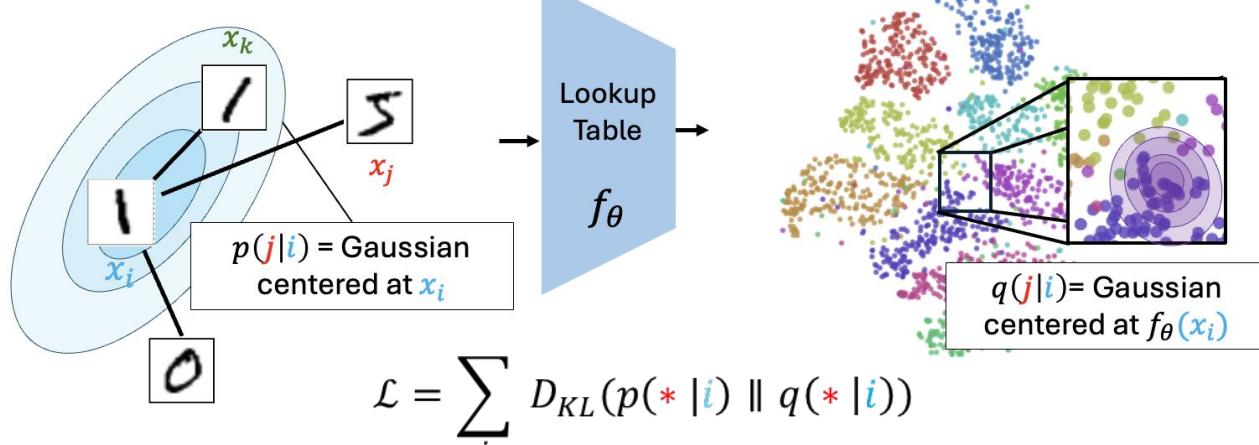


$$\mathcal{L} = \sum_i D_{KL}(p(\cdot | i) \parallel q(\cdot | i))$$

## Strong Common Pattern

SNE

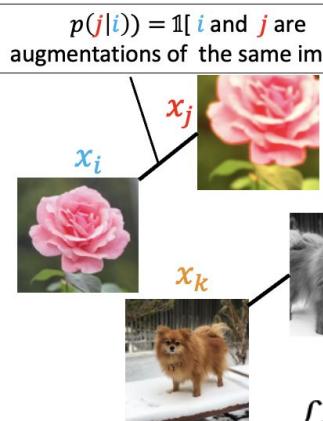
High Dimensional Data



SimCLR

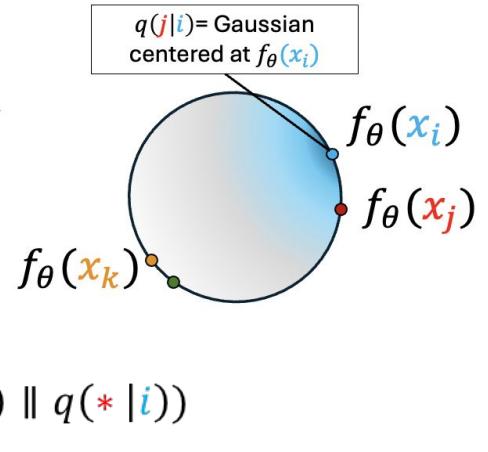
RGB Images

$p(j|i) = 1[i \text{ and } j \text{ are augmentations of the same image}]$



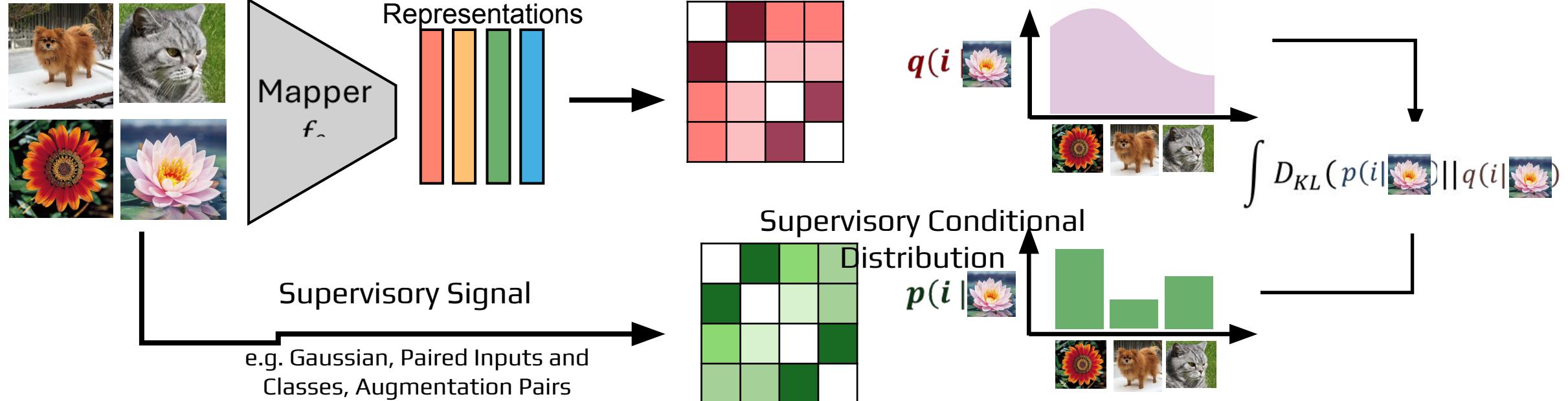
Parametric Mapper  $f_\theta$

Embeddings



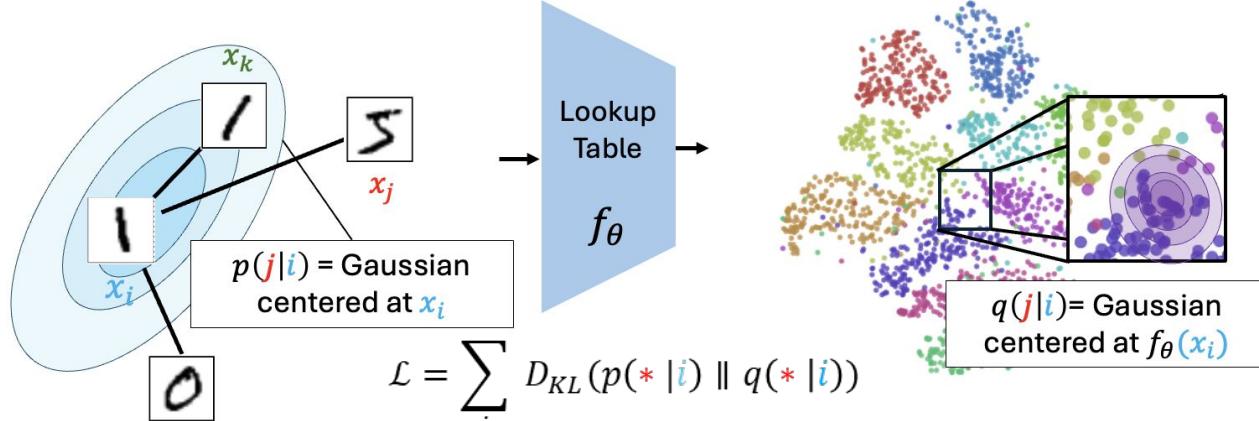
# I-Con Framework

## Overview



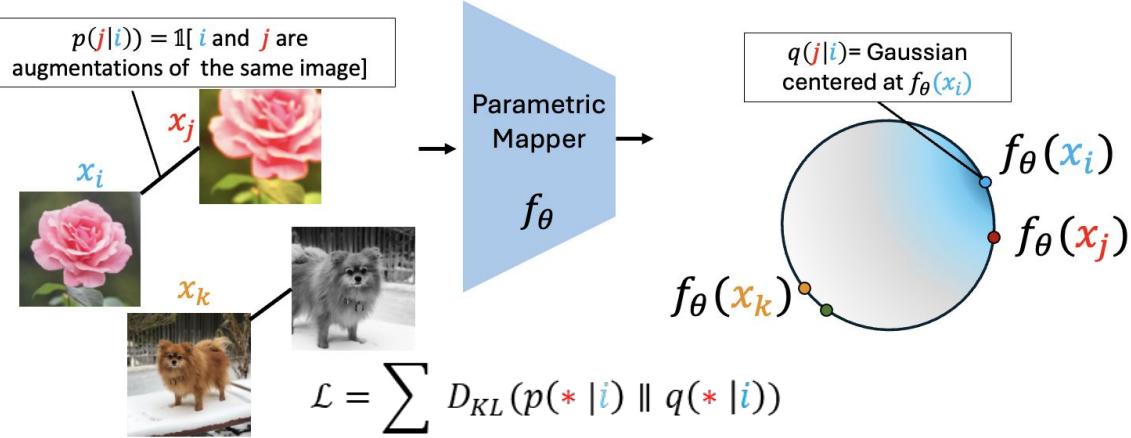
## SNE

High Dimensional Data

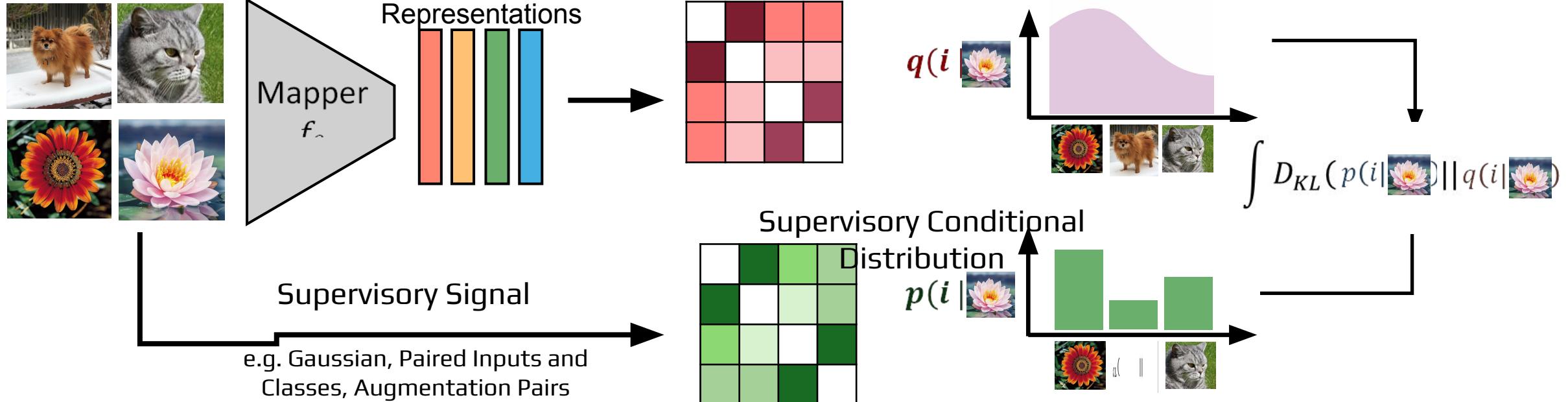


## SimCLR

RGB Images

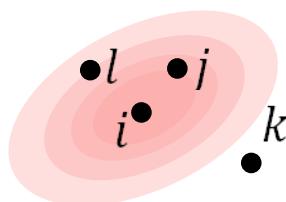


# Framework Overview



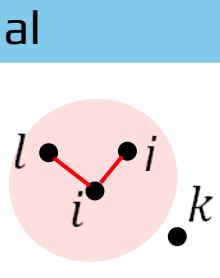
## Examples for types of $p$ and $q$

A) Spatial



e.g. Gaussian,  
Student-T  
Neighbors are based on distance

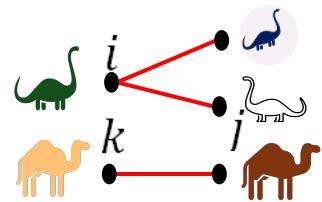
k-NN



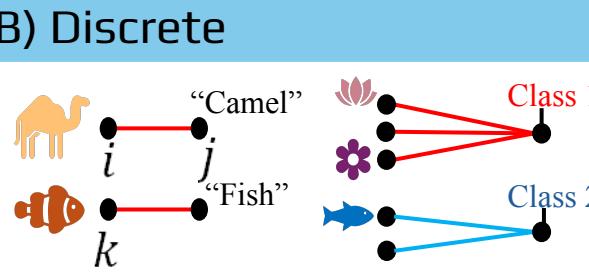
Positive Pairs

e.g., data augmentations or same-label pairs

B) Discrete

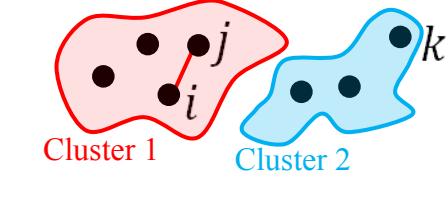


Cross-Modal  
e.g. image-text



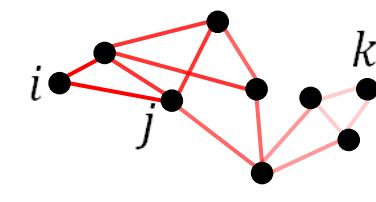
Labels  
e.g., images linked to category prototypes

C) Cluster



Neighbors grouped by shared cluster membership

D) Graph



Neighbors determined by graph connectivity

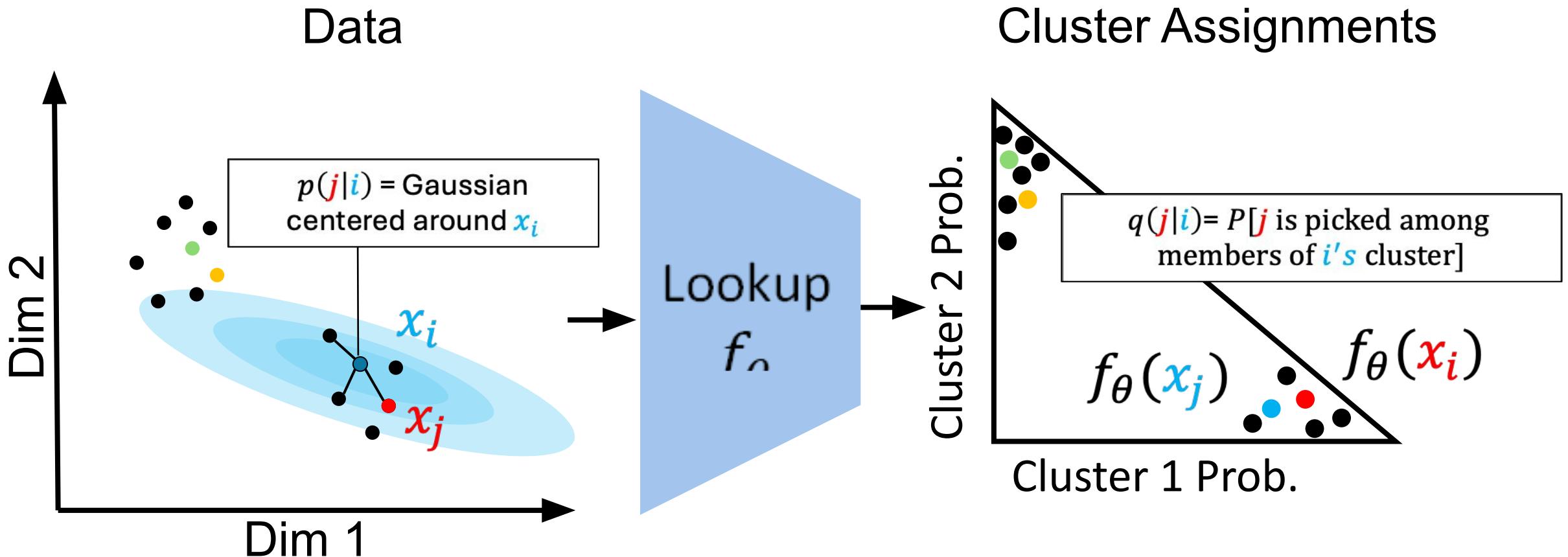
# Overview

- Intro
  - Welcome to the Zoo
  - The Periodic Table
- Methods:
  - Exploring some Examples
  - Generalizing ML Methods with a single Equation
  - **Building the Periodic Table**
- Experiments
  - Building an Image Classifier that doesn't need human labels
- Future Directions

## Supervisory Signal

# Example 3: Clustering

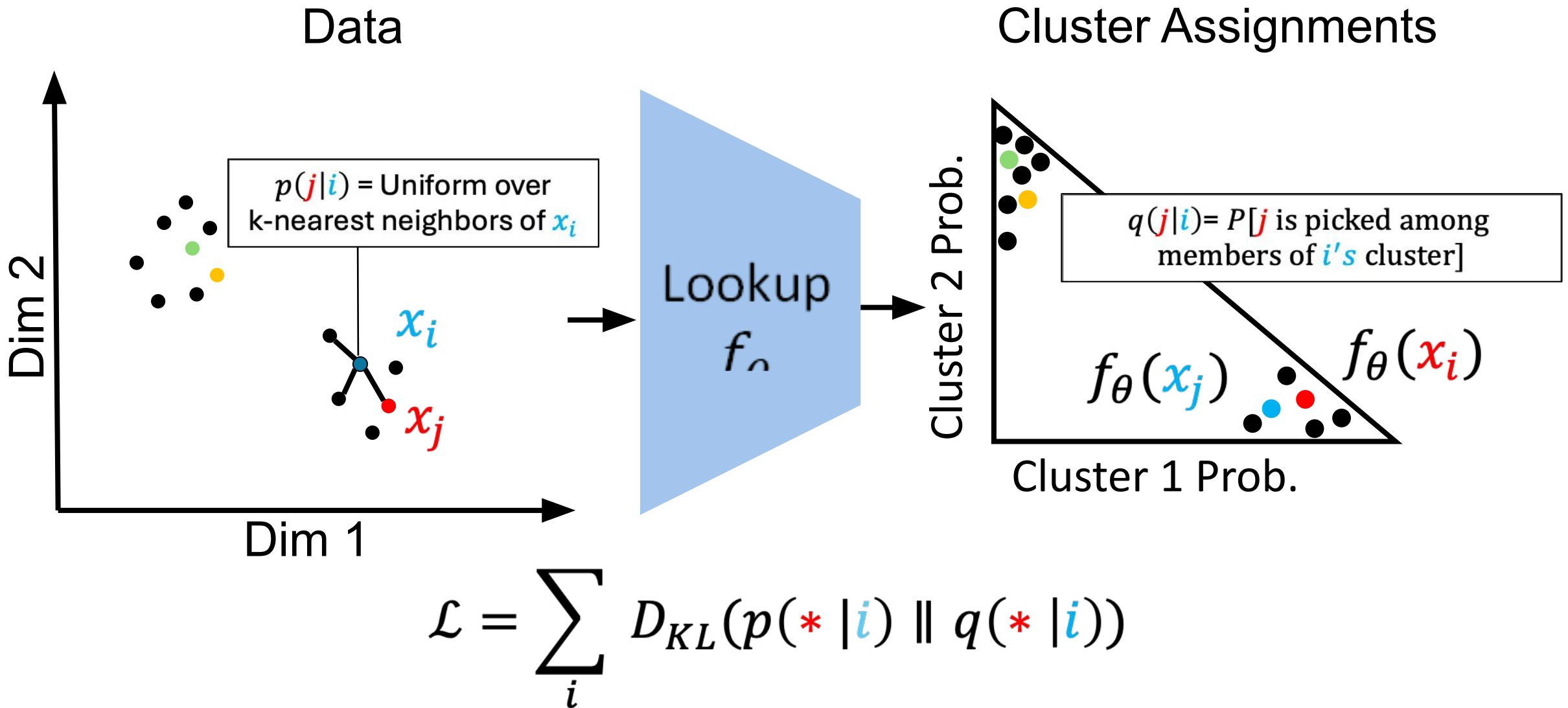
K-Means



$$\mathcal{L} = \sum_i D_{KL}(p(\cdot | i) \| q(\cdot | i))$$

# Example 3: Clustering

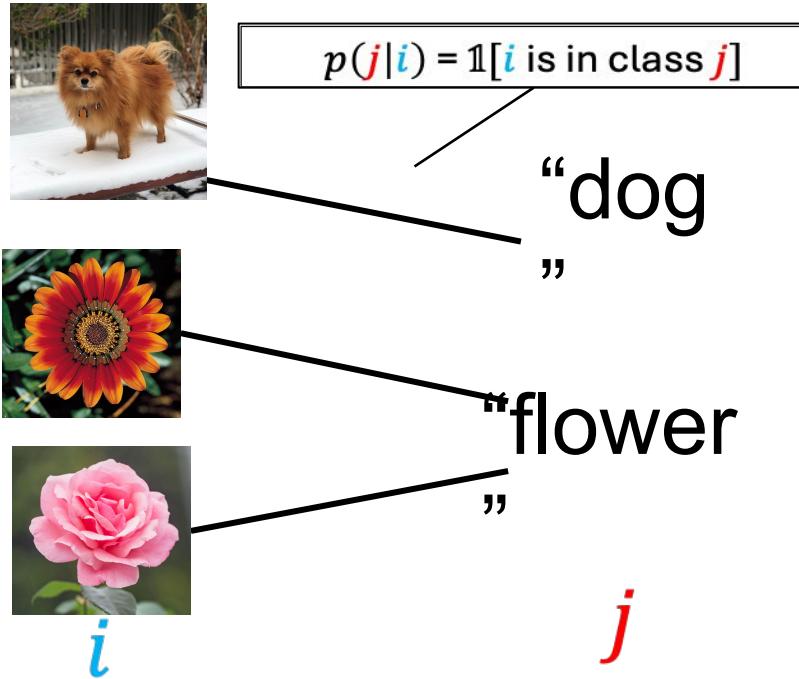
DCD



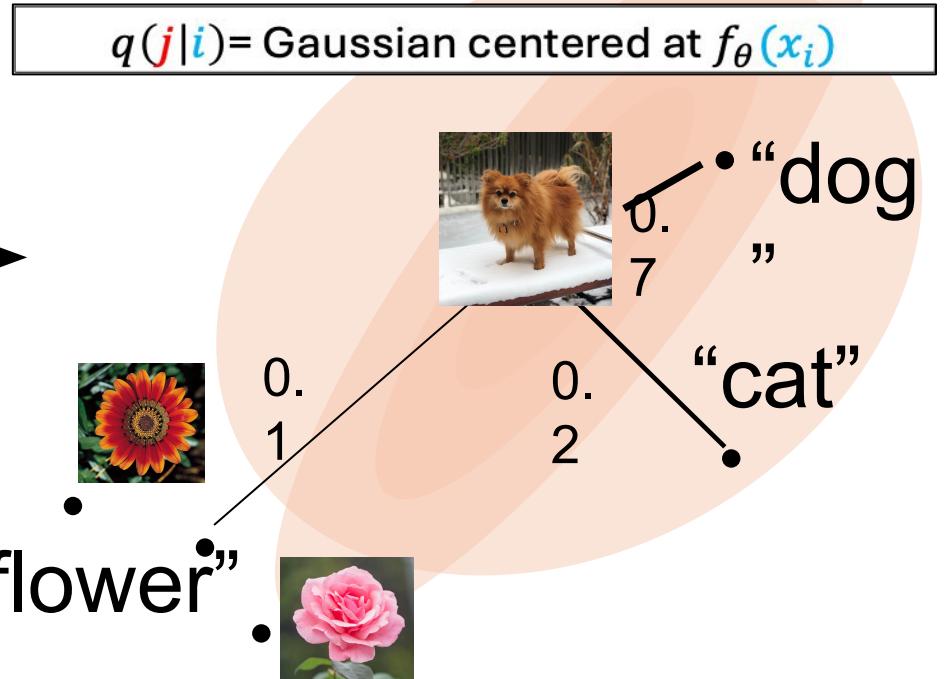
# Example 4: Supervised Learning

CrossEntropy

Data-Label Pairs



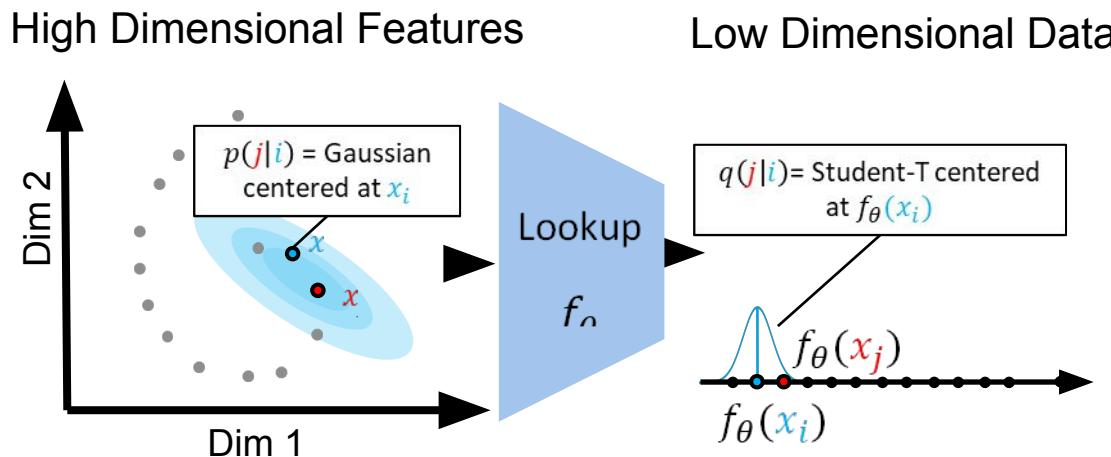
Data and Class Embeddings



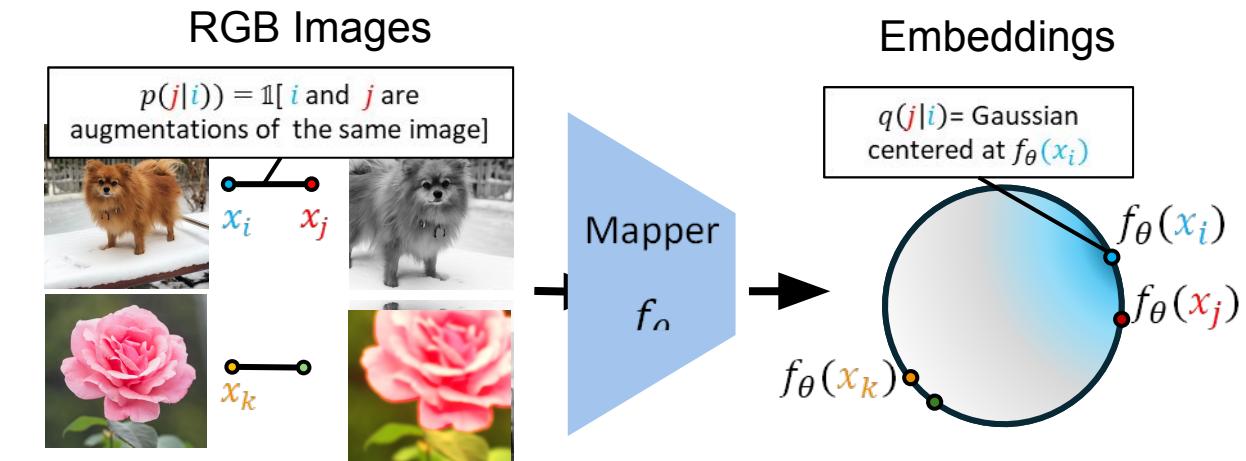
$$\mathcal{L} = \sum_i D_{KL}(p(\cdot | i) \| q(\cdot | i))$$

# Examples of methods as special cases of I-Con via different choices of $p$ , $q$ , and $f_\theta$

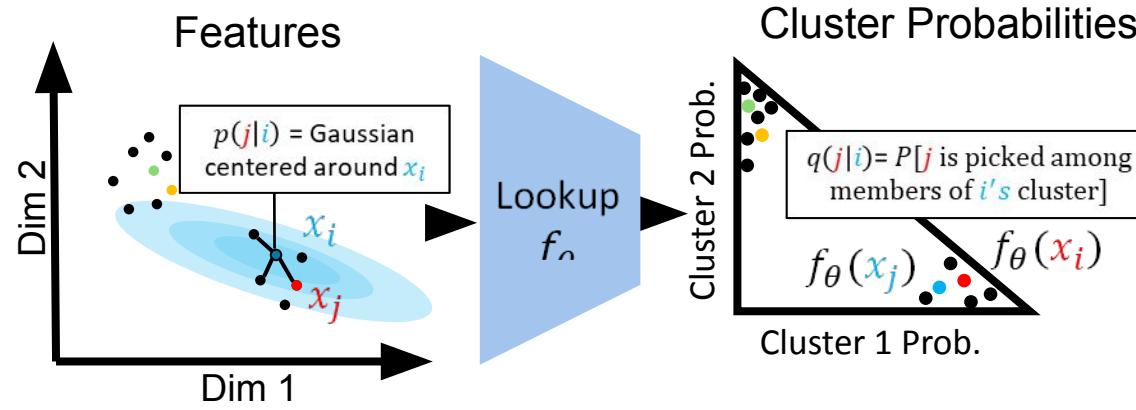
## t-SNE (Dimensionality Reduction)



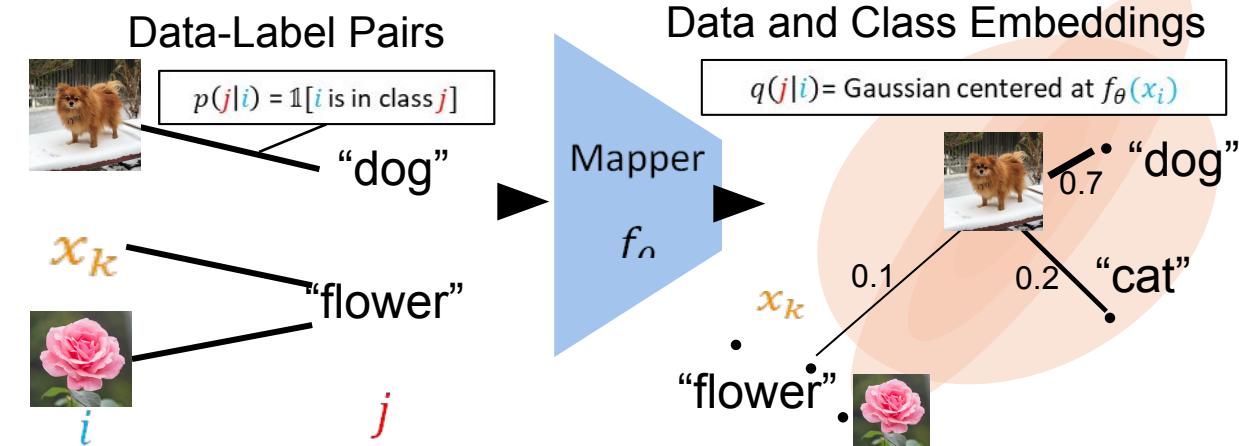
## SimCLR (Self-Supervised Learning)



## K-Means (Clustering)

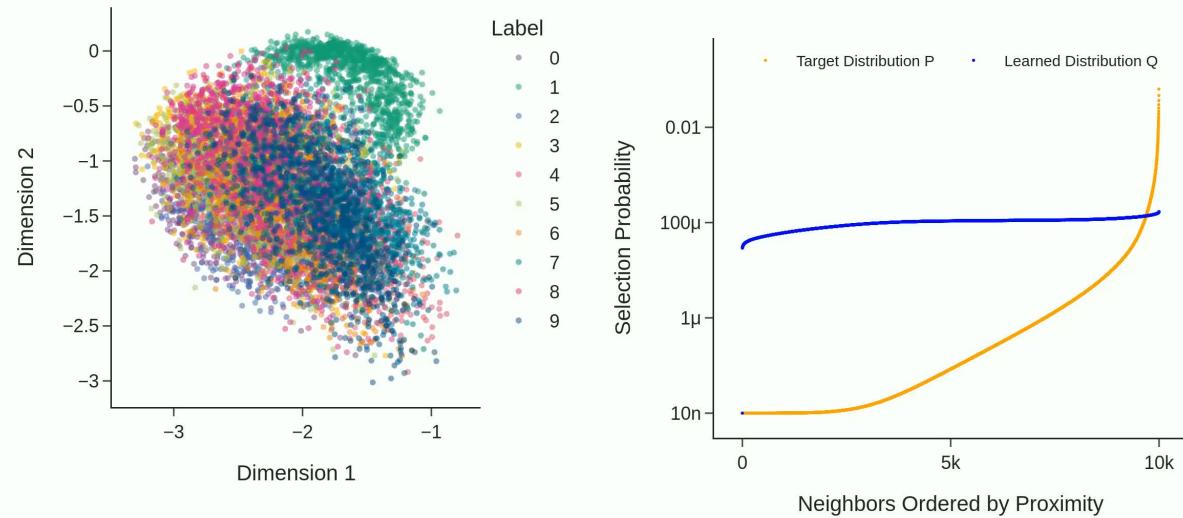


## Supervised Cross Entropy (Classification)

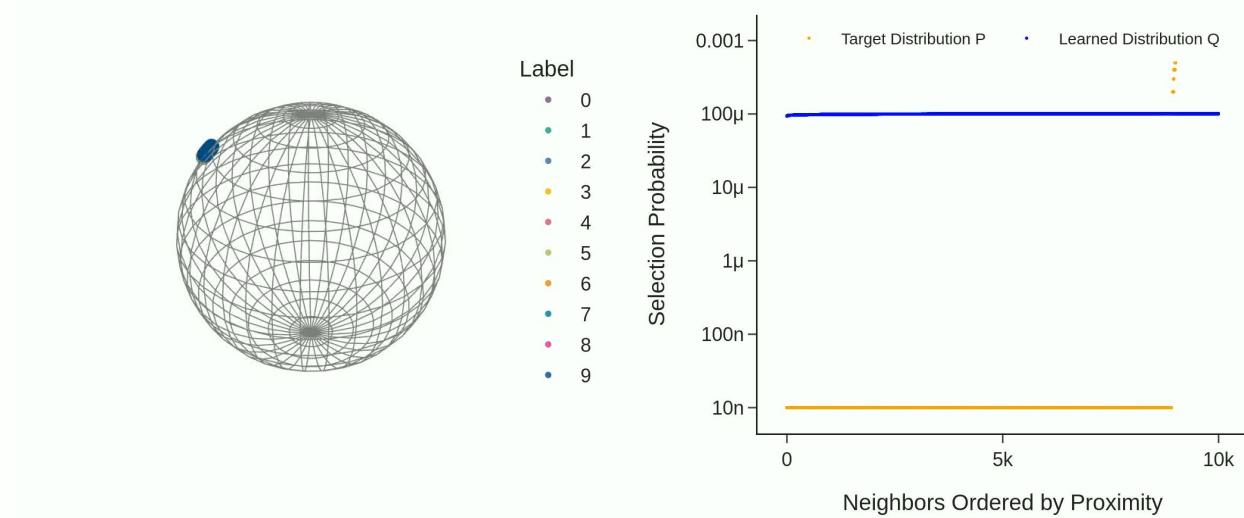


# Examples of methods as special cases of I-Con via different choices of $p$ , $q$ , and $f_\theta$ on MNIST

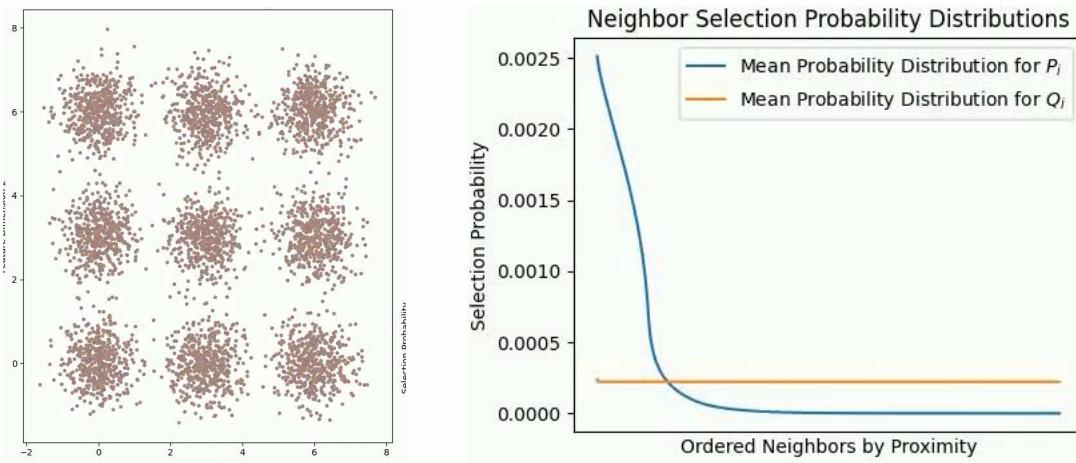
## t-SNE (Dimensionality Reduction)



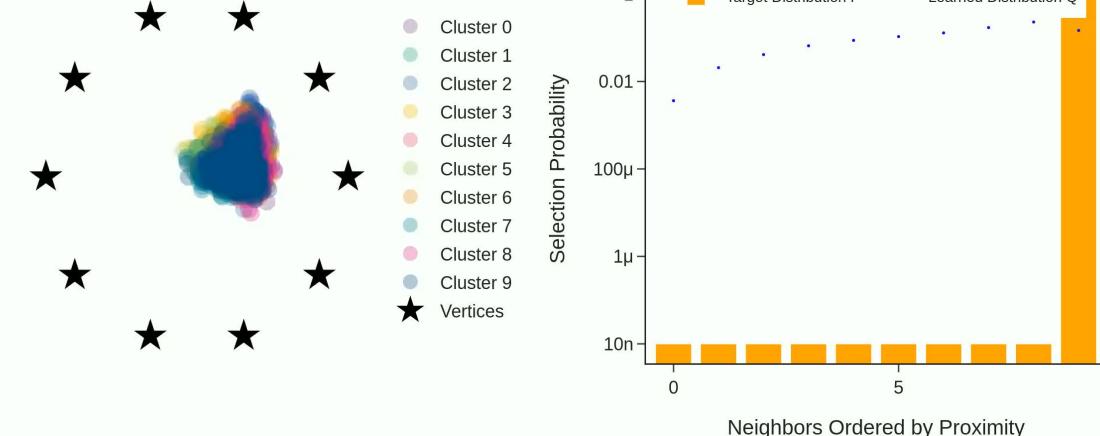
## Supervised Contrastive Learning



## K-Means or DCD (Clustering)

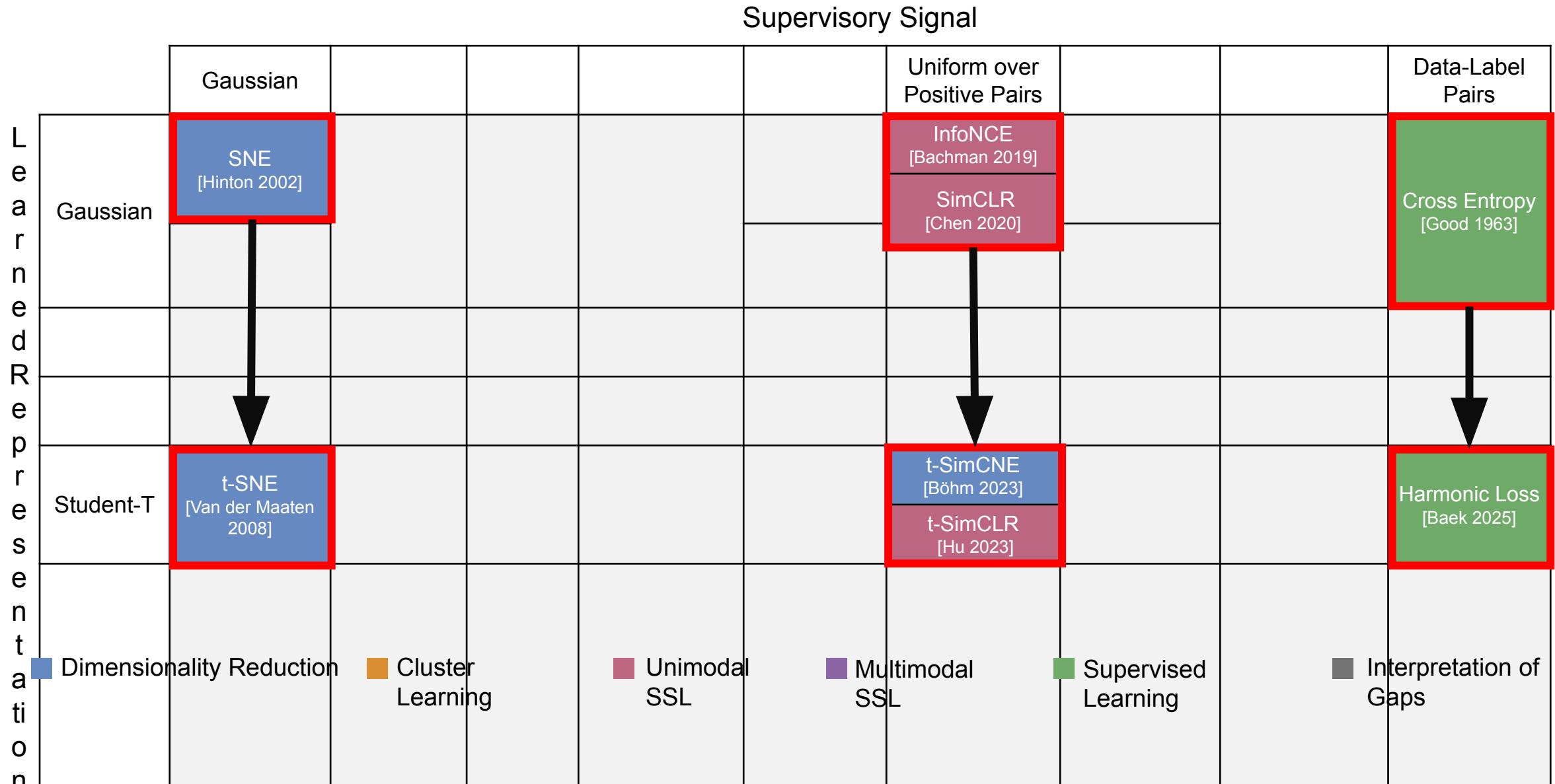


## Supervised Cross Entropy (Classification)



## Supervisory Signal

	Gaussian					Uniform over Positive Pairs			Data-Label Pairs
Learner	Gaussian	SNE [Hinton 2002]				InfoNCE [Bachman 2019]			Cross Entropy [Good 1963]
ended						SimCLR [Chen 2020]			
representation									
Dimensionality Reduction	Cluster Learning	Unimodal SSL	Multimodal SSL	Supervised Learning	Interpretation of Gaps				



	Gaussian		Identity		Uniform over K-Neighbors	Uniform over Positive Pairs			Data-Label Pairs
Learned Representations	Gaussian	SNE [Hinton 2002]				InfoNCE [Bachman 2019] SimCLR [Chen 2020]			Cross Entropy [Good 1963]
	Gaussian $\sigma \rightarrow \infty$			PCA [Pearson 1901]					
	Student-T	t-SNE [Van der Maaten 2008]				t-SimCNE [Böhm 2023] t-SimCLR [Hu 2023]			Harmonic Loss [Baek 2025]
	Dimensionality Reduction Clusters	X-Means [Macqueen 1967]		Cluster Learning		DCL [Yang 2012]	Multimodal NCE SSL Clustering [Ours]	Supervised Learning	Interpretation of Gaps

## Supervisory Signal

	Gaussian	Student-T	Identity	Graph Kernel Weights	Uniform over K-Neighbors	Uniform over Positive Pairs	Cross-Modal Pairs	Uniform over Classes	Data-Label Pairs
Learned Representations	Gaussian	SNE [Hinton 2002]			SNE with Uniform Affinities	InfoNCE [Bachman 2019] SimCLR [Chen 2020]	CLIP [Radford 2021]	SupCon [Khosla 2020]	Cross Entropy [Good 1963]
	Gaussian $\sigma \rightarrow \infty$	X-Sample CL [Sobal 2025]		PCA [Pearson 1901]		MoCoV3 [Chen 2021]	CMC [Tian 2020]		
	Gaussian $\sigma \rightarrow 0$					VI-Reg [Bardes 2021]			
	Student-T	t-SNE [Van der Maaten 2008]		t-SNE Graph Embedding	t-SNE with Uniform Affinities	t-SimCNE [Böhm 2023] t-SimCLR [Hu 2023]			Harmonic Loss [Baek 2025]
Dimensionality Reduction	Clusters	X-Means [Macqueen 1967]	Cluster Learning	Noisy Unimodal Clusters [Shrivastava 2010]	DCL [Yang 2012]	Multimodal NCE SSL Clustering [Ours]	Supervised Learning		Interpretation of Gaps

## Supervisory Signal

	Gaussian	Student-T	Identity	Graph Kernel Weights	Uniform over K-Neighbors	Uniform over Positive Pairs	Cross-Modal Pairs	Uniform over Classes	Data-Label Pairs
Learned Representations	Gaussian	SNE [Hinton 2002]	Dual t-SNE	SNE Graph Embeddings	SNE with Uniform Affinities	InfoNCE [Bachman 2019]	CLIP [Radford 2021]	SupCon [Khosla 2020]	Cross Entropy [Good 1963]
	X-Sample CL [Sobal 2025]				LGSimCLR [El Banani 2023]	SimCLR [Chen 2020]			
	Gaussian $\sigma \rightarrow \infty$					MoCoV3 [Chen 2021]	CMC [Tian 2020]		
	Gaussian $\sigma \rightarrow 0$					VI-Reg [Bardes 2021]	Average Margin CLIP	Average Margin SupCon	
Dimensionality Reduction	Student-T	t-SNE [Van der Maaten 2008]	Doubly t-SNE	t-SNE Graph Embedding	t-SNE with Uniform Affinities	t-SimCNE [Böhm 2023]	t-CLIP	Triplet CLIP	Triplet SupCon
						t-SimCLR [Hu 2023]			Error rate
							t-SupCon		Harmonic Loss [Baek 2025]
Clustering	K-Means [Macqueen 1967]	t K-Means Learning		NoSSL	Unimodal SSL [Shen 2010]	DCL	Multimodal NCE [Yang 2012]	Supervised Clustering	Interpretation of Gaps

# Overview

- Intro
  - Welcome to the Zoo
  - The Periodic Table
- Methods:
  - Exploring some Examples
  - Generalizing ML Methods with a single Equation
  - Building the Periodic Table
- Experiments
  - **Building an Image Classifier that doesn't need human labels**
- Future Directions

# Building New Representation Learners

Given the I-Con framework, we developed state-of-the-art clustering methods by:

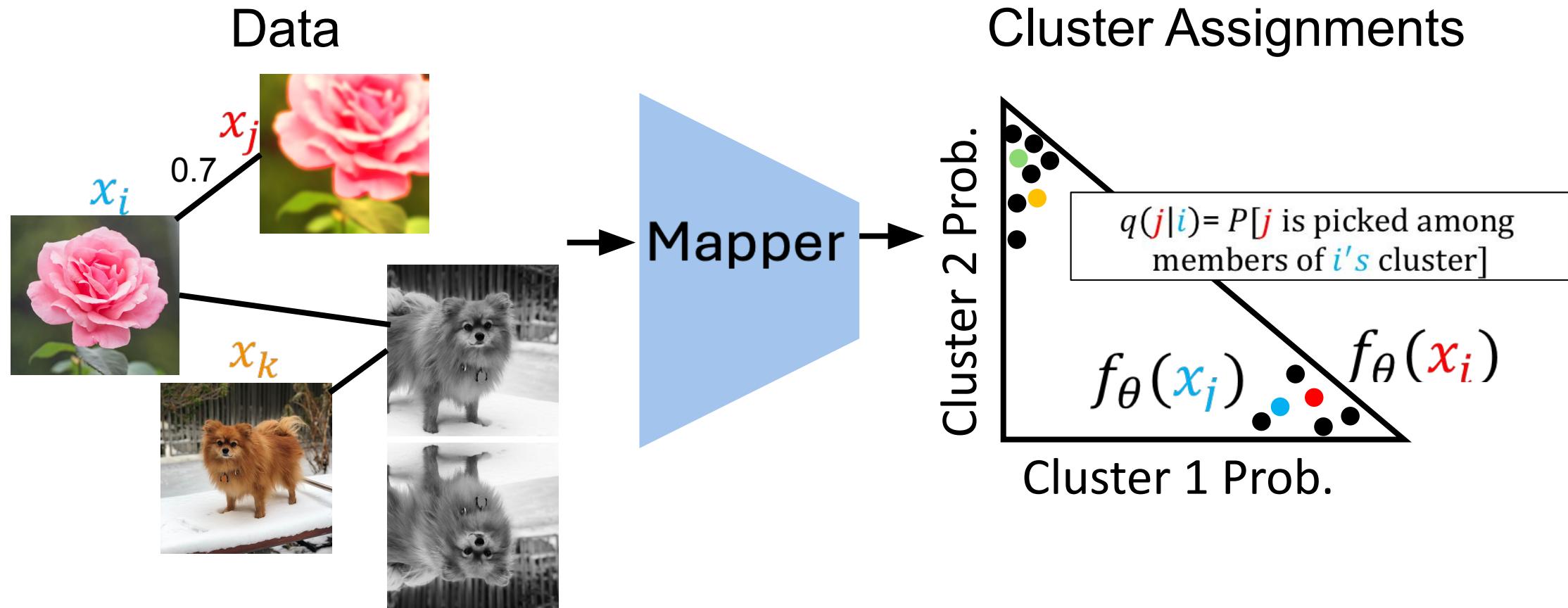
- Filling a gap in the table
- Transferring ideas from other fields
- Adapting the neighborhood width

		Supervisory Signal								
		Gaussian	Student-T	Identity	Graph Kernel Weights	Uniform over K-Neighbors	Uniform over Positive Pairs	Cross-Modal Pairs	Uniform over Classes	Data-Label Pairs
Learned Representant	Gaussian	SNE [Hinton 2002]	Dual t-SNE		SNE Graph Embeddings	SNE with Uniform Affinities	InfoNCE [Bachman 2019]	CLIP [Radford 2021]	SupCon [Khosla 2020]	Cross Entropy [Good 1963]
	X-Sample CL [Sobal 2025]						SimCLR [Chen 2020]			
	Gaussian $\sigma \rightarrow \infty$			PCA [Pearson 1901]			MoCoV3 [Chen 2021]			
	Gaussian $\sigma \rightarrow 0$						VI-Reg [Bardes 2021]	Average Margin CLIP	Average Margin SupCon	
Learned Representant	Student-T	t-SNE [Van der Maaten 2008]	Doubly t-SNE		t-SNE Graph Embedding	t-SNE with Uniform Affinities	Triplet Loss [Schroff 2015]	Triplet CLIP	Triplet SupCon	Error rate
	U-nif						t-SimCNE [Böhmm 2023]	t-CLIP	t-SupCon	Harmonic Loss [Baek 2025]
							t-SimCLR [Hu 2023]			

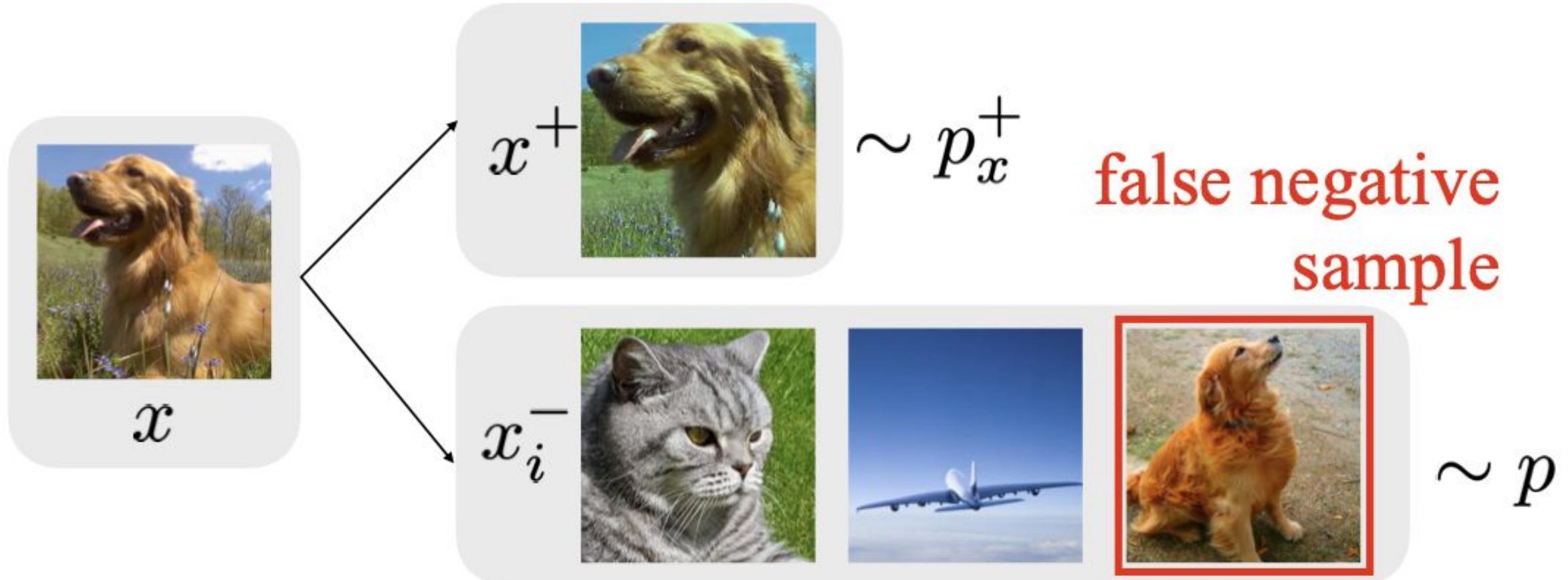
# Building New Representation Learners

Given the I-Con framework, we developed state-of-the-art clustering methods by:

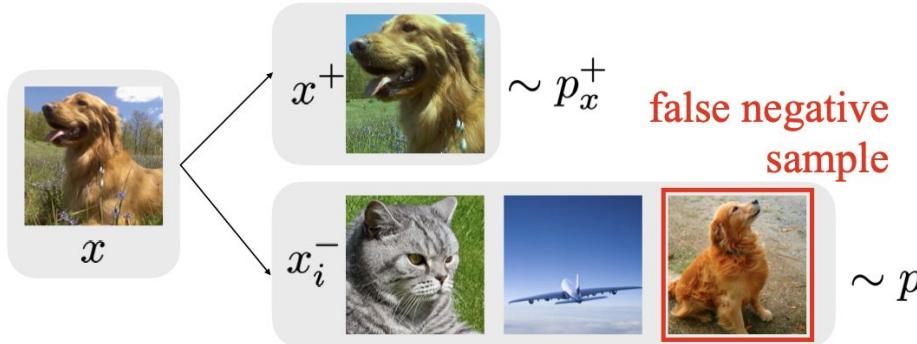
- Filling a gap in the table
- Transferring ideas from other fields
- Adapting the neighborhood width



# Debiased Contrastive Learning (Chuang et al, NeurIPS 2020)



# Debiased Contrastive Learning (Chuang et al, NeurIPS 2020)



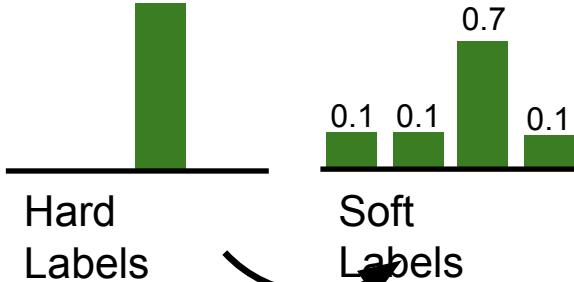
$$L_{\text{Debiased}}^{N,M}(f) = \mathbb{E}_{\substack{x \sim p; x^+ \sim p_x^+ \\ \{u_i\}_{i=1}^N \sim p \\ \{v_i\}_{i=1}^M \sim p_x^{+M}}} \left[ -\log \frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + Ng(x, \{u_i\}_{i=1}^N, \{v_i\}_{i=1}^M)} \right]$$

where

$$g(x, \{u_i\}_{i=1}^N, \{v_i\}_{i=1}^M) = \max \left\{ \frac{1}{\tau^-} \left( \frac{1}{N} \sum_{i=1}^N e^{f(x)^T f(u_i)} - \tau^+ \frac{1}{M} \sum_{i=1}^M e^{f(x)^T f(v_i)} \right), e^{-1/t} \right\}$$

doesn't represent a conditional probability

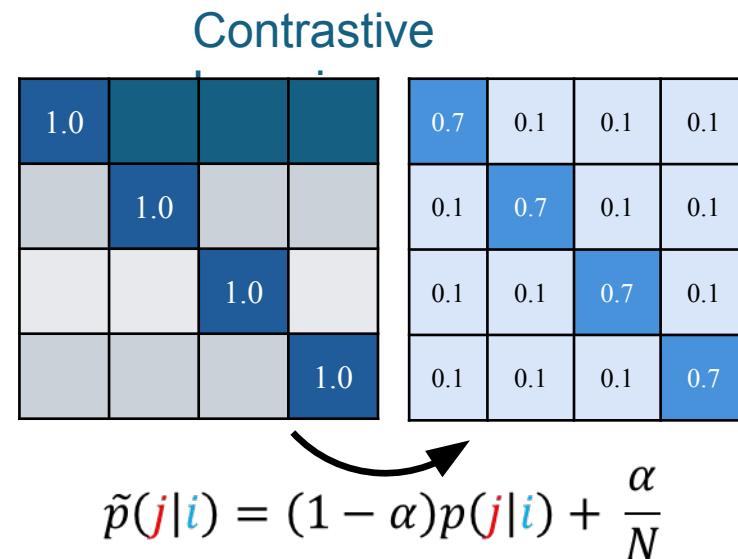
A more general approach for debiasing  
Supervised Learning



$$\tilde{p}(c|i) = (1 - \alpha)p(c|i) + \frac{\alpha}{N}$$

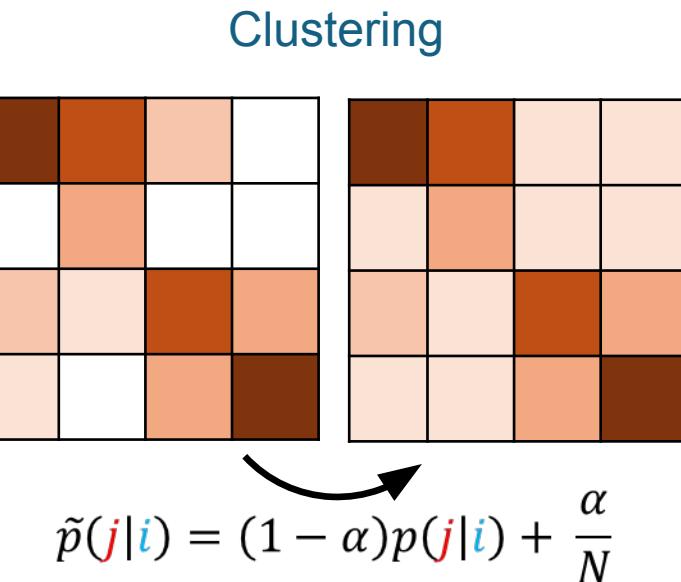
Reflect the uncertainty in the labels, not the softmax

$$\mathcal{L} = \sum \tilde{p}(c|i) \log q(c|i)$$



Soften target distribution, but learned distribution unchanged

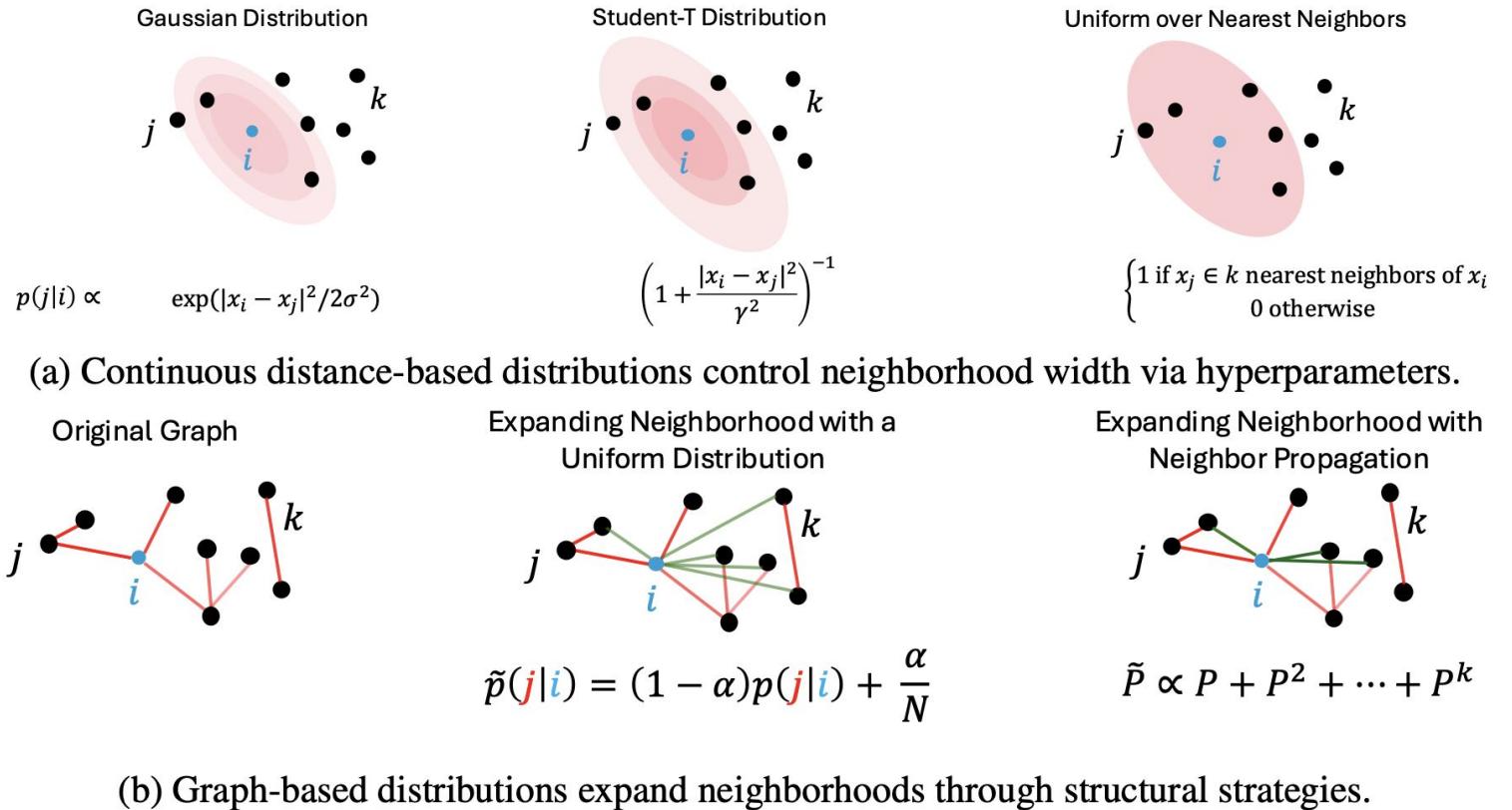
$$\mathcal{L} = \mathbb{E}_{\tilde{p}} \left[ -\log \frac{\exp(f_\theta(i) \cdot f_\theta(j))}{\sum_k \exp(f_\theta(i) \cdot f_\theta(k))} \right]$$



Learned distribution unchanged

$$\mathcal{L} = \sum_i KL(\tilde{p}(j|i) \| q(f_\theta(x_j) | f_\theta(x_i)))$$

# Adapting the neighborhood width

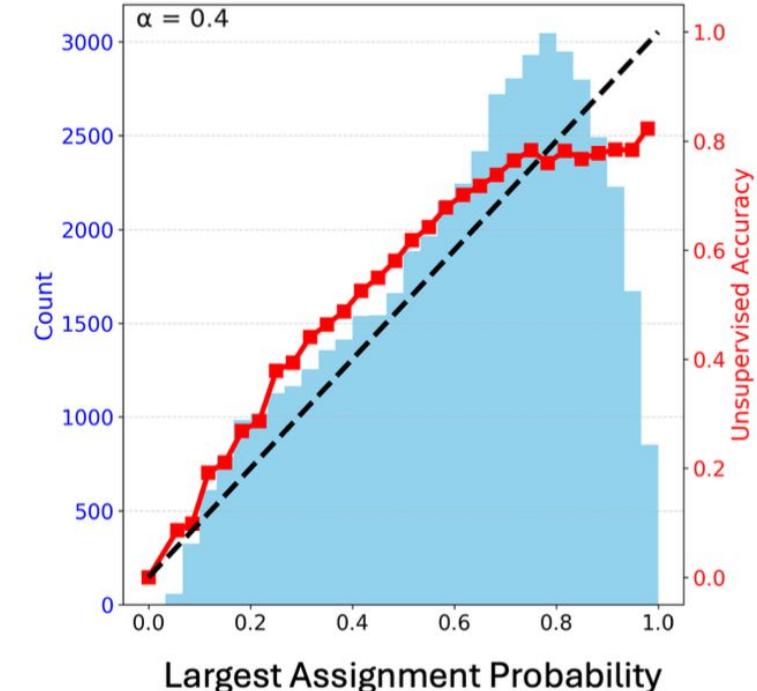
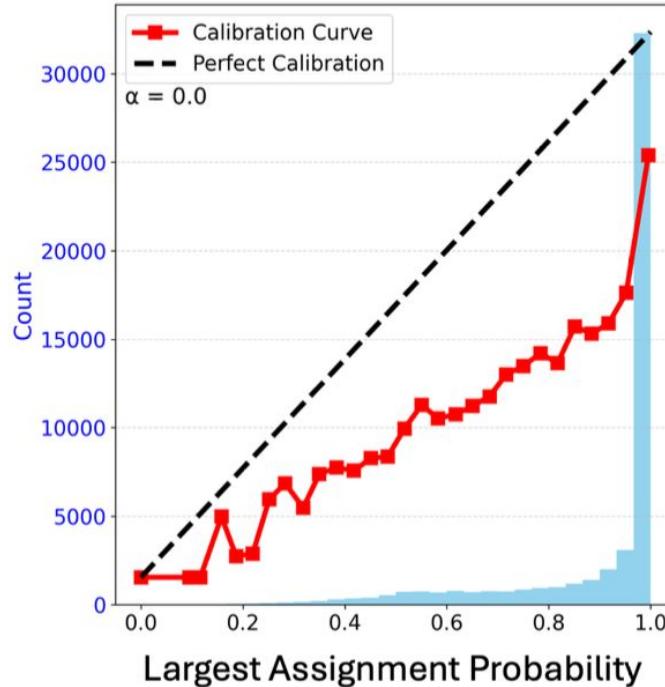


Method	DiNO ViT-S/14	DiNO ViT-B/14	DiNO ViT-L/14
Baseline	55.51	63.03	65.72
+ KNNs	56.43	64.26	65.70
+ 1-walks on KNN	<b>58.09</b>	<b>64.29</b>	65.97
+ 2-walks on KNN	57.84	64.27	<b>67.26</b>
+ 3-walks on KNN	57.82	64.15	67.02

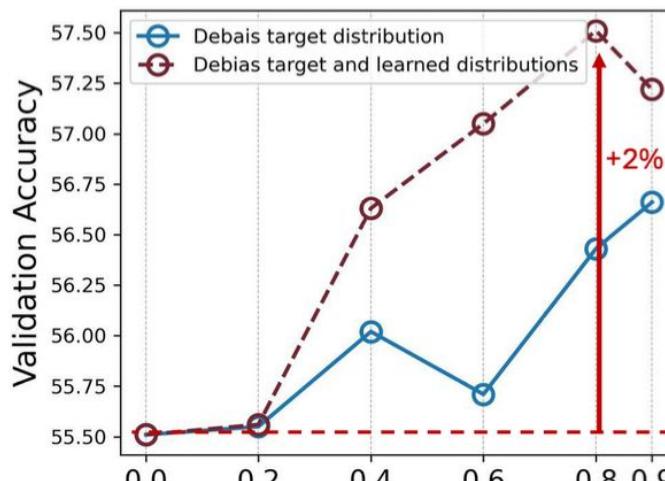
# Results for our Debiasing in Cluster Learning

$$\mathcal{L} = \sum_i KL(\tilde{p}(j|i) \parallel q(f_\theta(x_j)|f_\theta(x_i)))$$

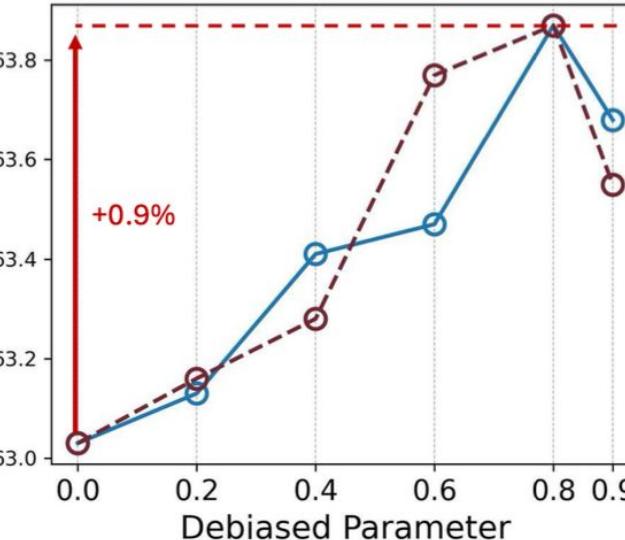
where  $\tilde{p}(j|i) = (1 - \alpha)p(j|i) + \frac{\alpha}{N}$



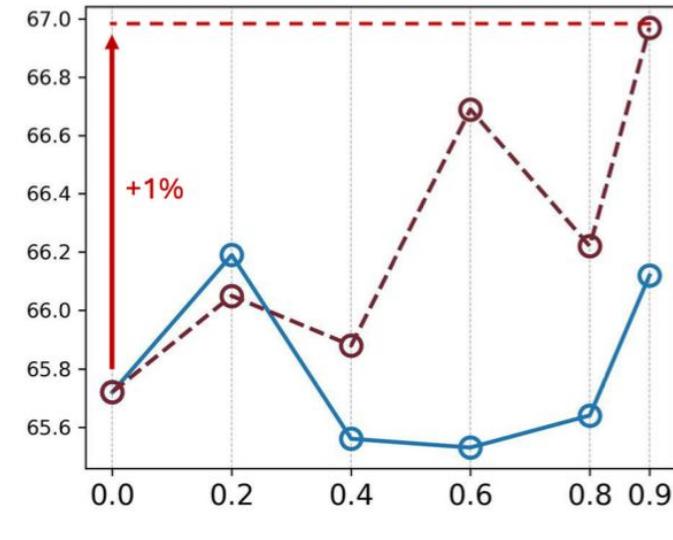
DiNO ViT-S/14



DiNO ViT-B/14



DiNO ViT-L/14



# More Results in Cluster Learning

Method	DiNO ViT-S/14	DiNO ViT-B/14	DiNO ViT-L/14
Baseline	55.51	63.03	65.70
+ Debiasing	57.27 ± 0.07	63.72 ± 0.09	66.87 ± 0.07
+ KNN Propagation	<b>58.45</b> ± 0.23	64.87 ± 0.19	67.25 ± 0.21
+ EMA	57.8 ± 0.26	<b>64.75</b> ± 0.18	<b>67.52</b> ± 0.28

Method	DiNO ViT-S/14	DiNO ViT-B/14	DiNO ViT-L/14
k-Means	51.84	52.26	53.36
Contrastive Clustering	47.35	55.64	59.84
SCAN	49.20	55.60	60.15
TEMI	56.84	58.62	—
<b>Debiased InfoNCE Clustering (Ours)</b>	<b>57.8</b> ± 0.26	<b>64.75</b> ± 0.18	<b>67.52</b> ± 0.28

# Slides Credits

- 6.S898 Deep Learning, MIT EECS  
<https://phillipi.github.io/6.s898/>
- 6.8300/6.8301: Advances in Computer Vision, MIT CSAIL,  
<https://advances-in-vision.github.io>