

# Lecture 12

## Representation Learning

# Overview

- Why should we care about Representation Learning?
- Transfer learning
- Unsupervised learning
- Self-supervised learning

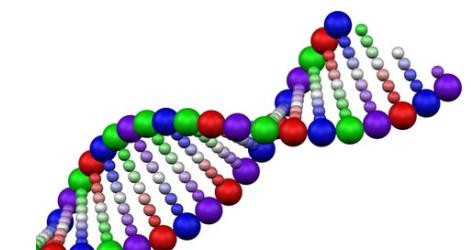
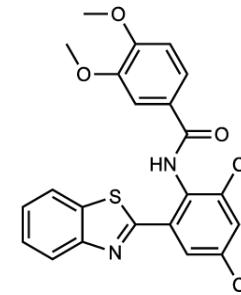
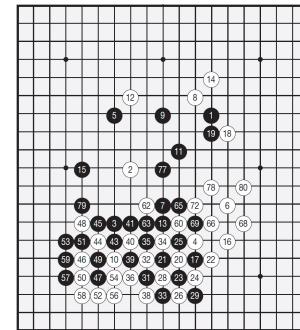
**Why should we care about  
Representation Learning?**

# Deep Learning is Representation Learning

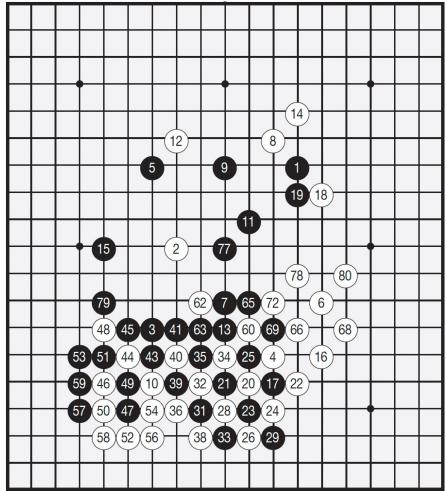
- Representation Learning: worth a conference name :)
  - International Conference on Learning Representations (ICLR)
- Represent raw data to solve complex problems
  - compression, abstraction, conceptualization
- Raw data in different forms
  - pixels, words, waves, states, molecules, DNA, ...



Current approaches to object recognition make essential use of machine learning methods. To improve their performance, we can collect larger datasets, learn more powerful models, and use better techniques for preventing overfitting. Until recently, datasets of labeled images were relatively small compared to those of natural language, NLP datasets such as Penn Treebank and CHiRA-10/100 [12]. Simple recognition tasks can be solved quite well with datasets of this size, especially if they are augmented with label-preserving transformations. For example, the current best performing neural network for image classification appears to be Inception v3 [14]. But objects in realistic settings exhibit considerable variability, so to learn to recognize them it is necessary to use much larger training sets. And indeed, the shortcomings of small image datasets have been well-known for some time. Most notably, it is difficult to learn to recognize objects without collecting datasets with millions of images. The new larger datasets include LabelMe [23], which consists of hundreds of thousands of fully-segmented images, and ImageNet [6], which consists of over 15 million labeled images and images in over 22,000 categories.



# Deep Learning is Representation Learning

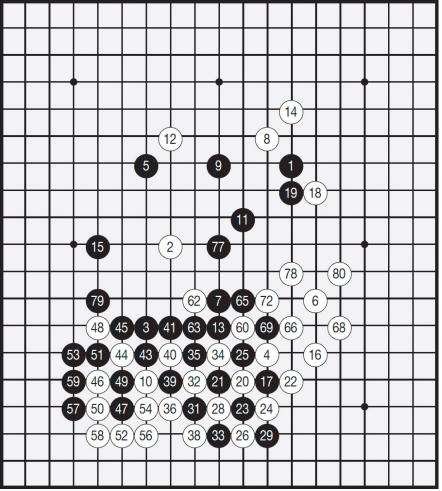


$3^{361}$  states?

## Game of Go

- an exponentially large number of states?
- infeasible to enumerate, memorize, or search

# Deep Learning is Representation Learning



$3^{361}$  states?

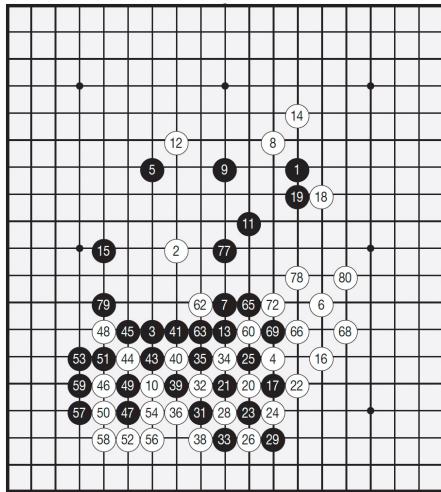
*Bad  
representation*



$256^{3 \times 500 \times 500}$ ?

- Image space has exponentially more states than Go.

# Deep Learning is Representation Learning

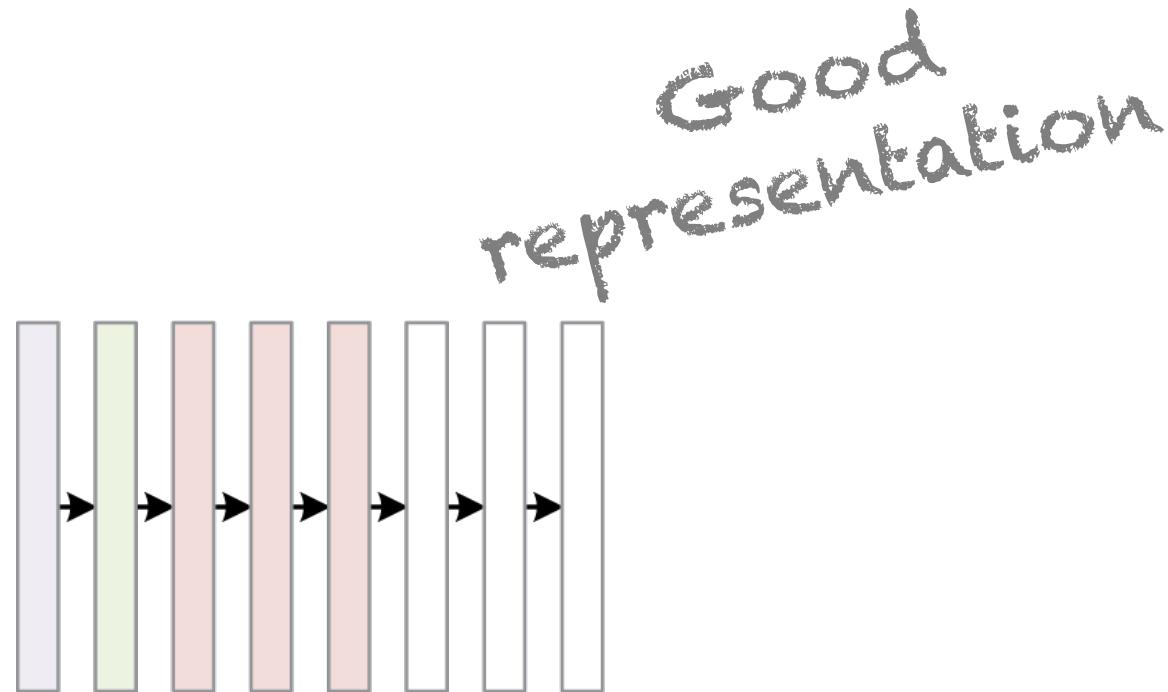


$3^{361}$  states?

- Image recognition is solved in representation space.

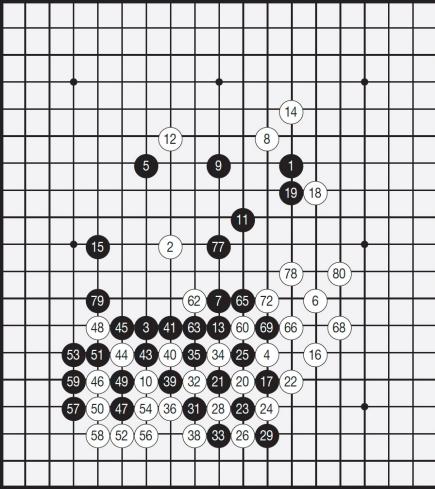


$256^{3 \times 500 \times 500}$ ?

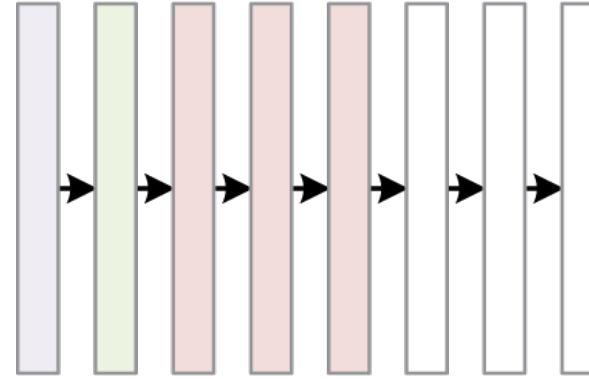


# Deep Learning is Representation Learning

- Go playing can be solved in representation space.



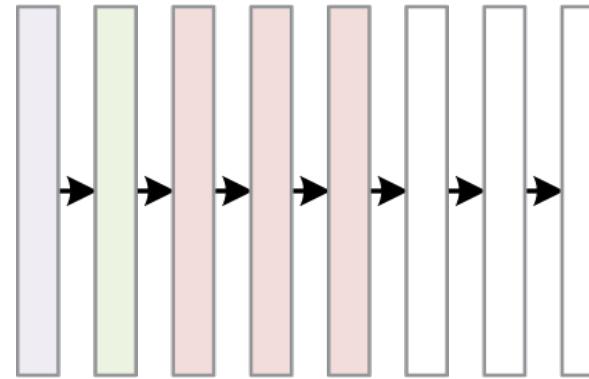
$3^{361}$  states?



Good  
representation



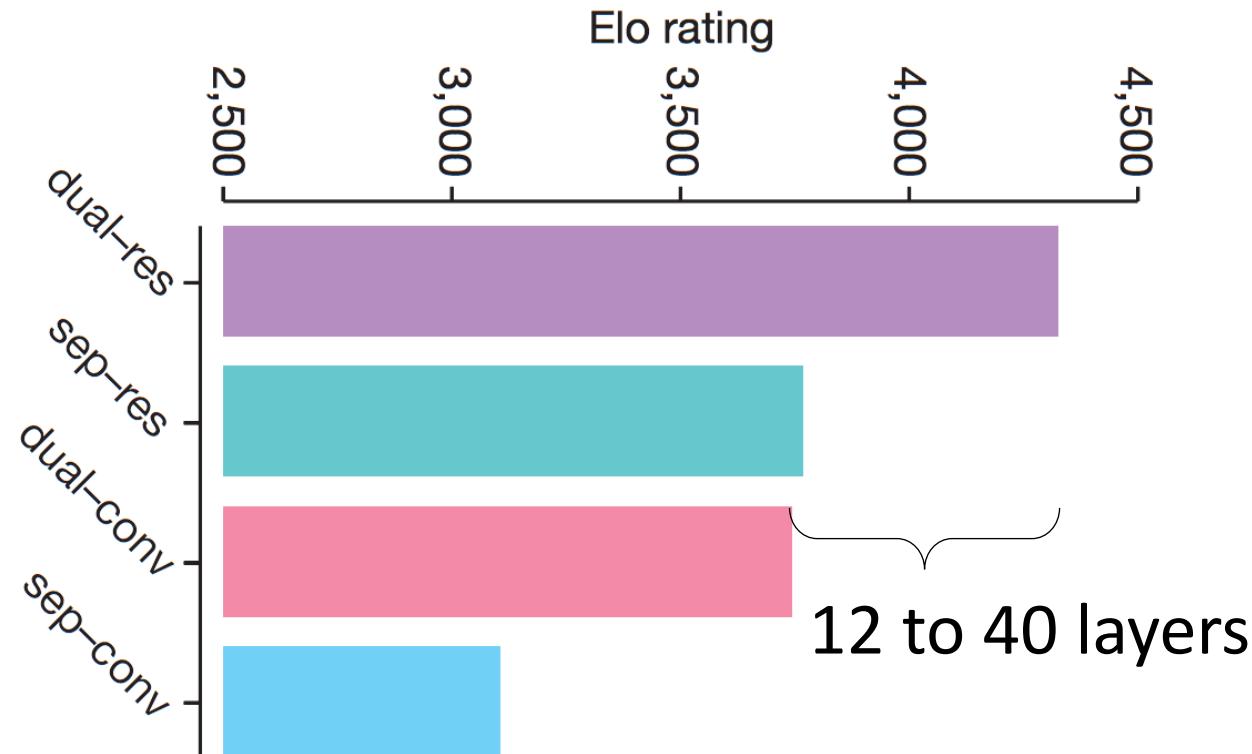
$256^{3 \times 500 \times 500}$ ?



# Deep Learning is Representation Learning

**AlphaGo:** better representation leads to higher rating

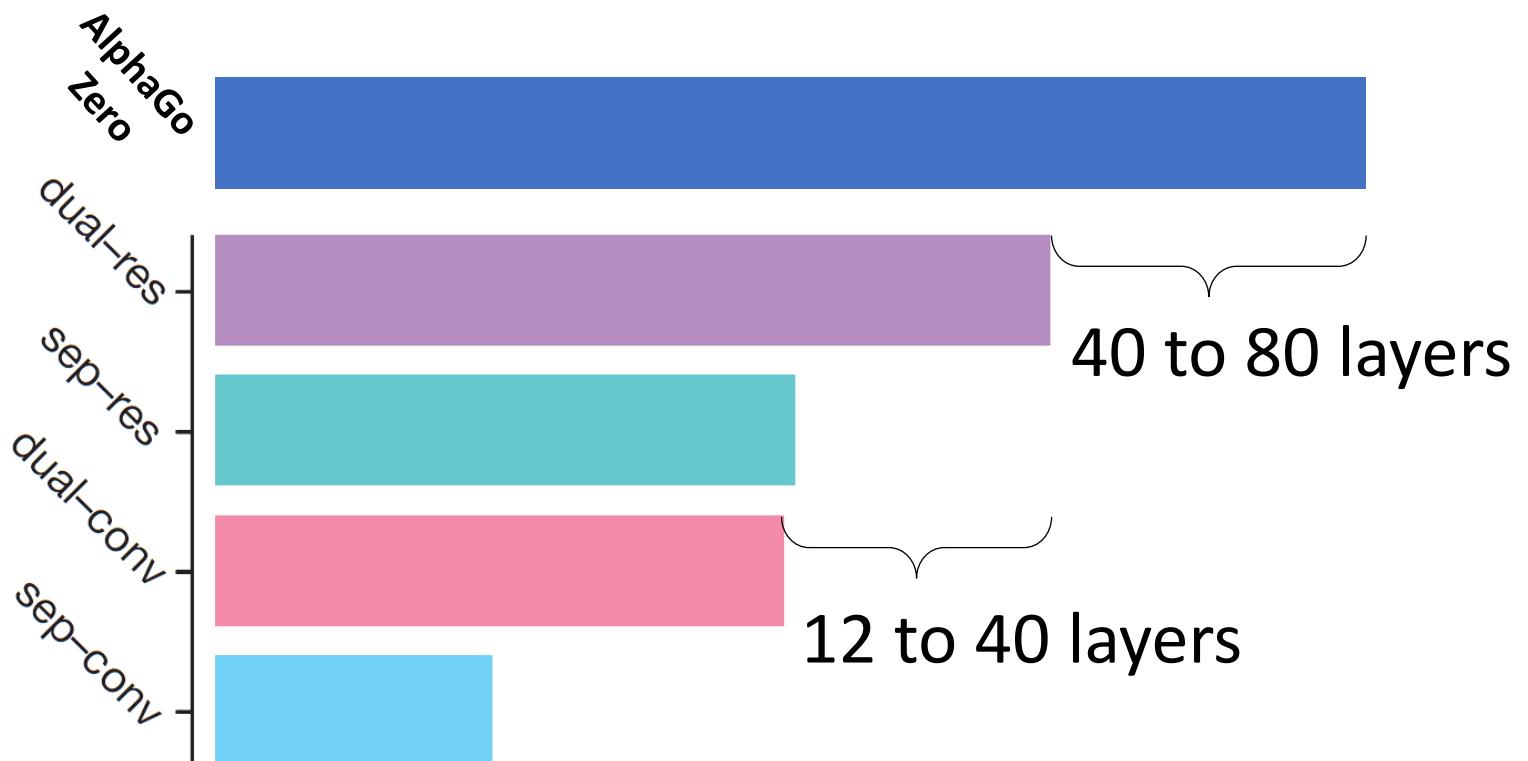
- outperform best human players
- without human knowledge



# Deep Learning is Representation Learning

**AlphaGo:** better representation leads to higher rating

- outperform best human players
- without human knowledge



# How to Represent Images?



→ class



→ edge → class

- Domain knowledge required
- OK-ish for low-level representations
- Extremely difficult to design high-level representations



→ edge → orientation → class

deeper



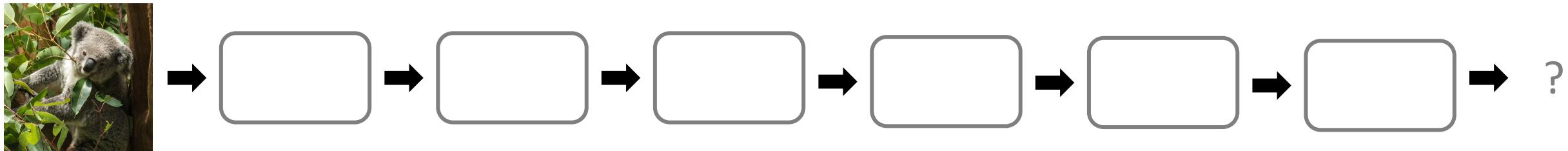
→ edge → orientation → histogram → class



→ edge → orientation → histogram → clusters → class

# How to Represent Images w/ Deep Learning?

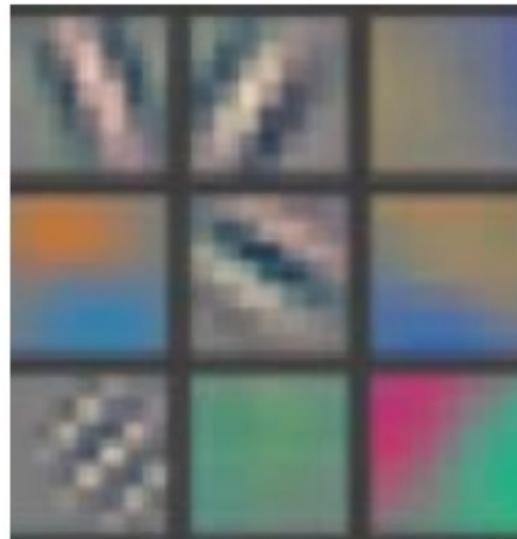
general modules (instead of specialized features)



**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



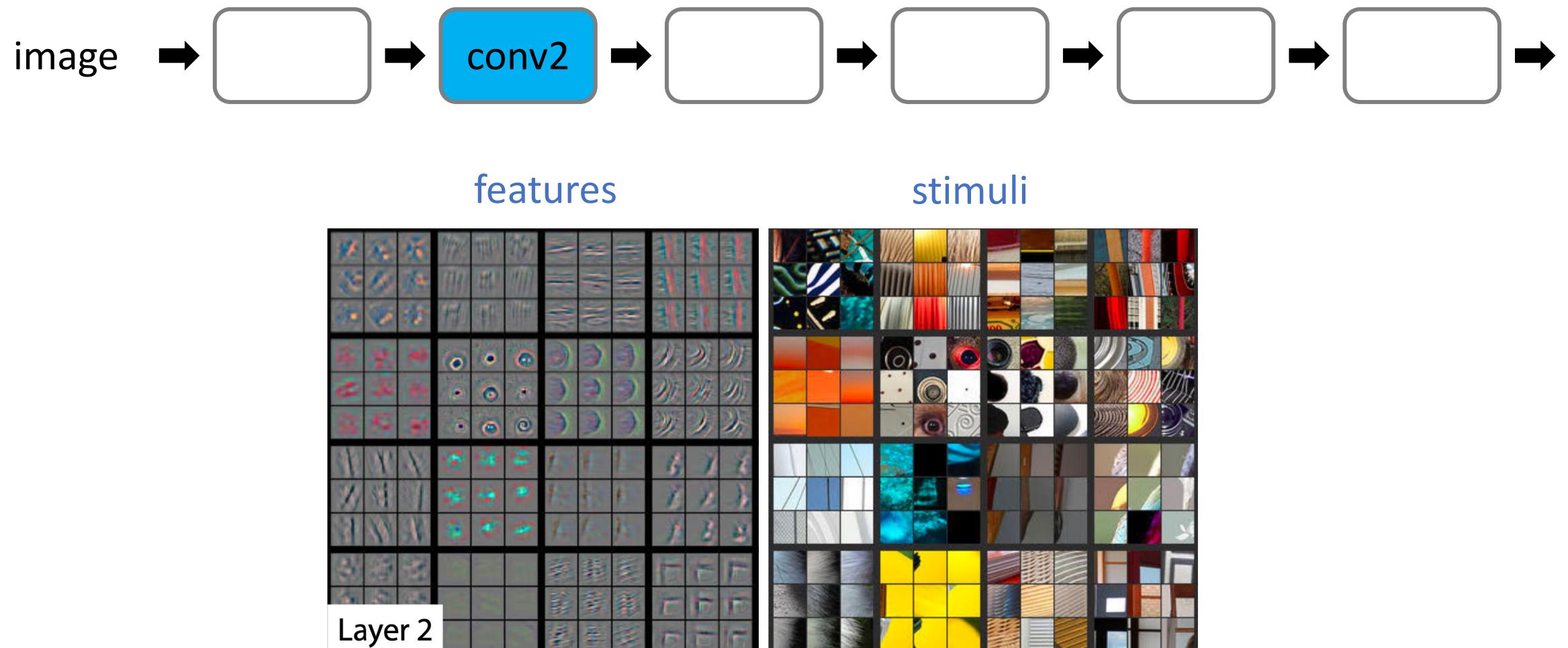
features



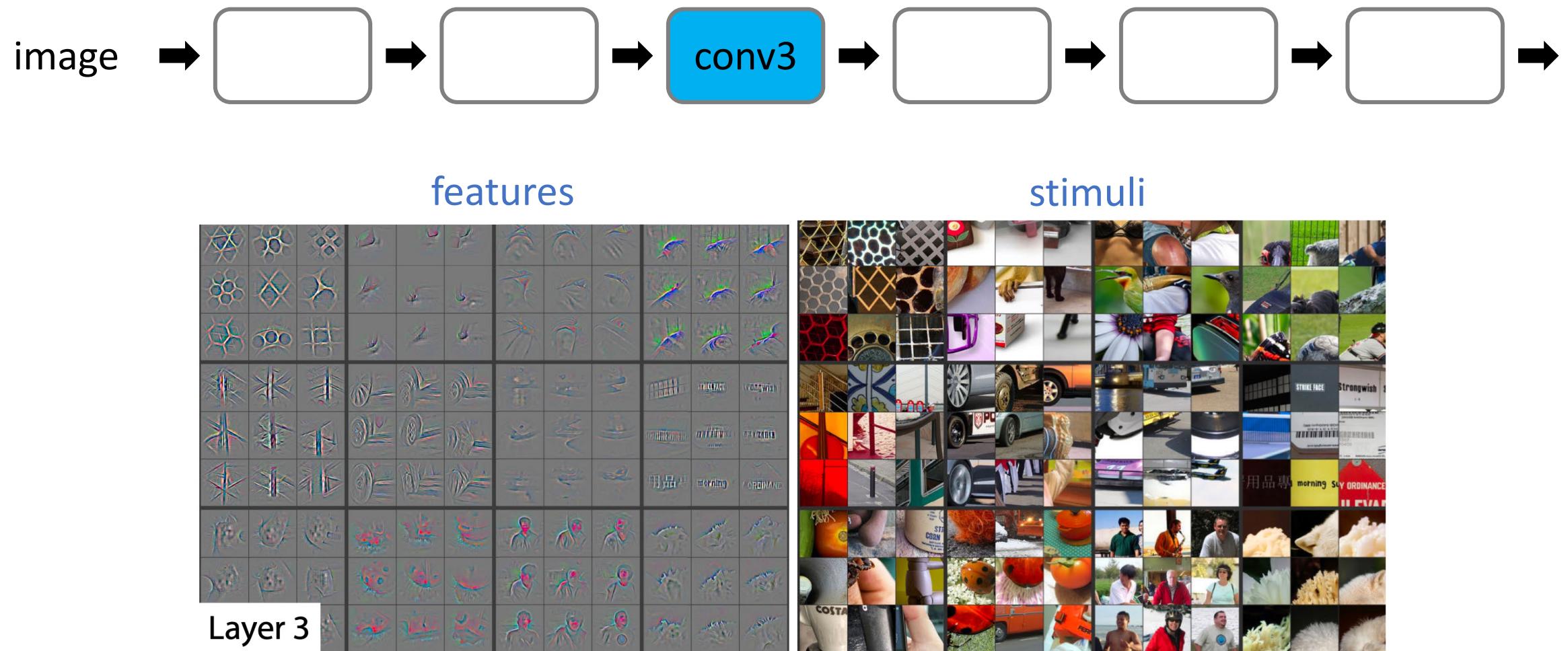
stimuli

(patches with the highest  
1-hot activations)

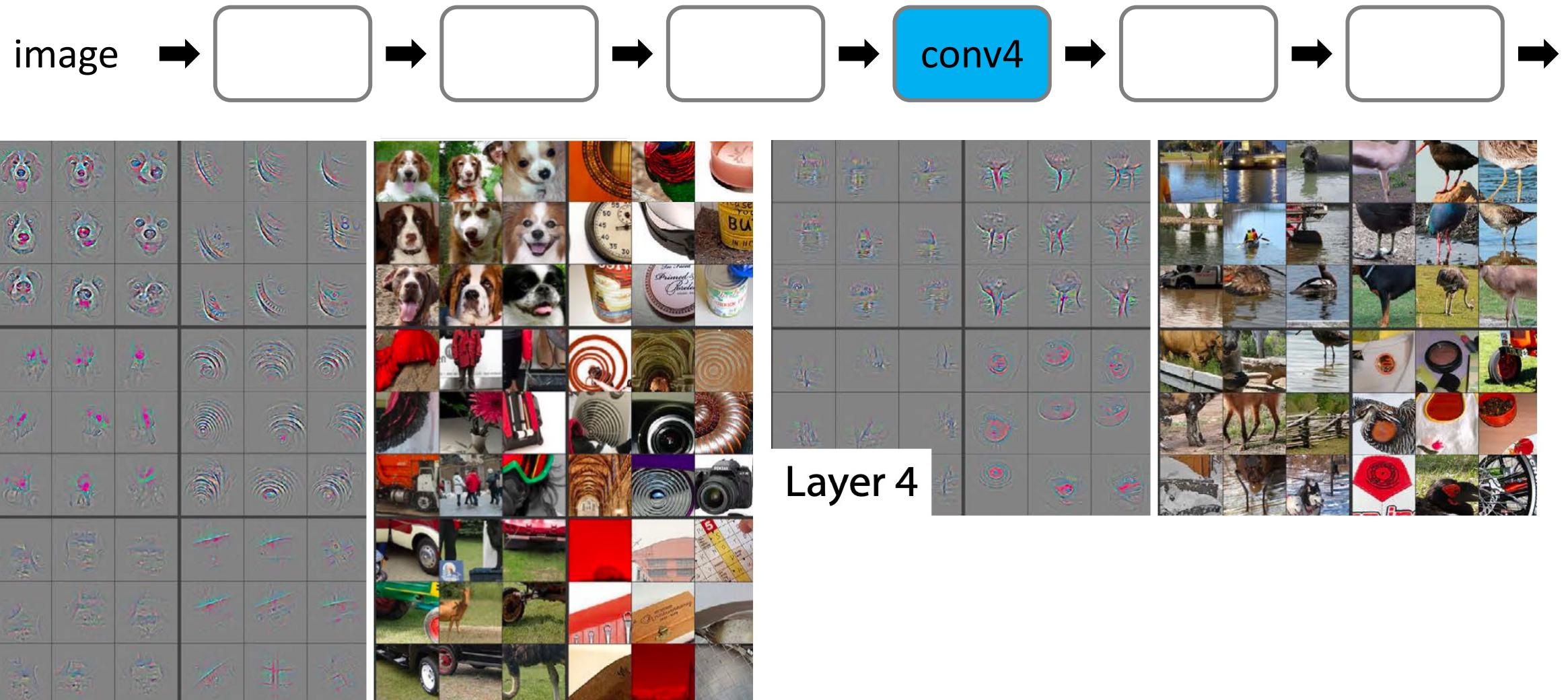
# Multiple Levels of Representations



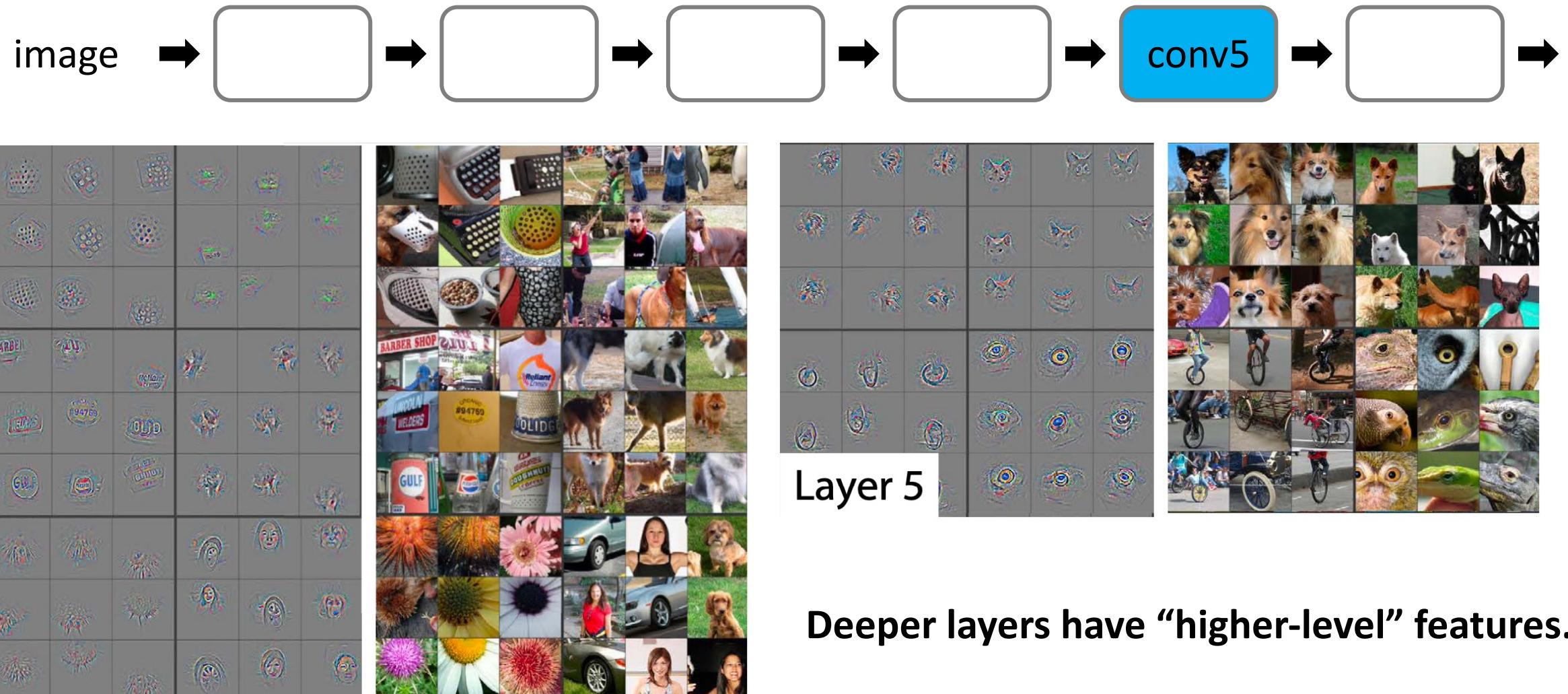
# Multiple Levels of Representations



# Multiple Levels of Representations

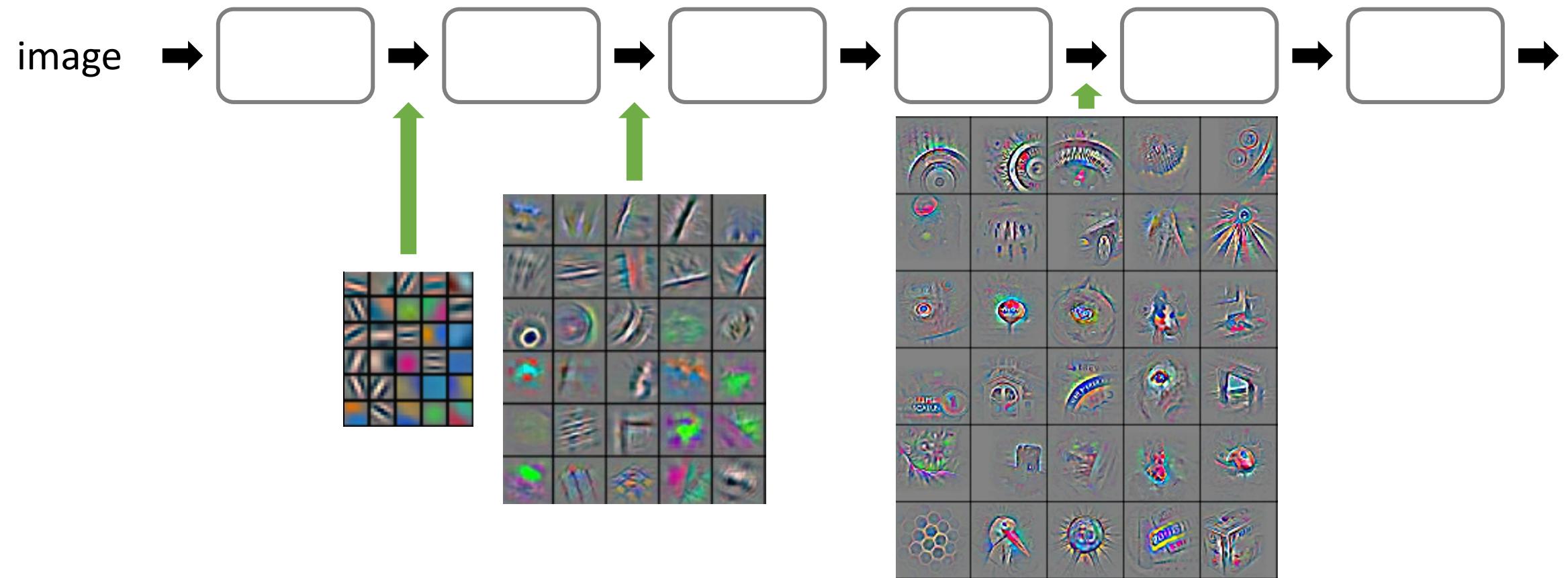


# Multiple Levels of Representations



Deeper layers have “higher-level” features.

# Multiple Levels of Representations



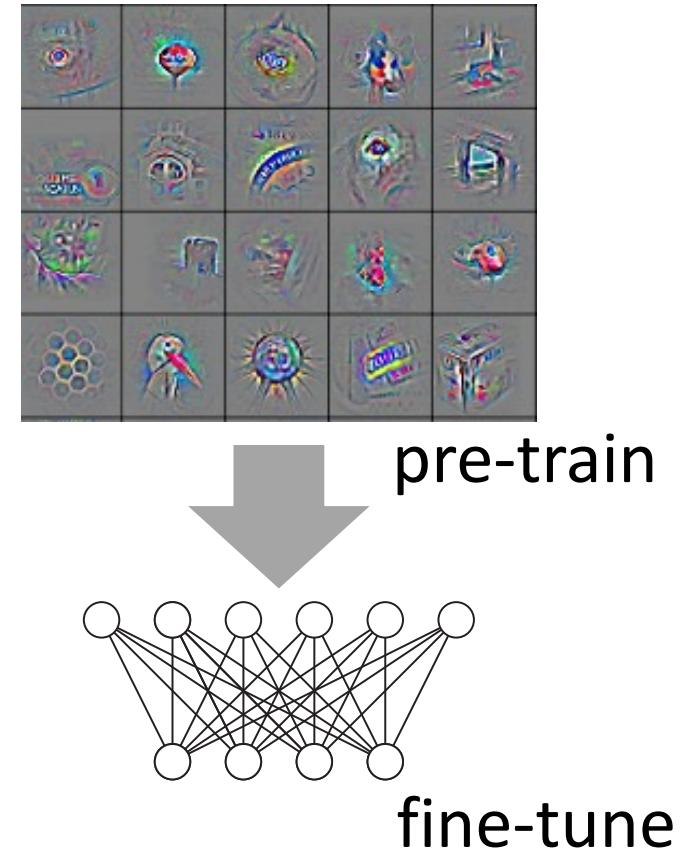
Deeper layers have “higher-level” features.

# Deep Representations are Transferrable!

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

# **Transfer learning**

# Transfer learning: Pre-training & Fine-tuning

## Pre-training



## Pre-training:

- to learn **general** representations
- on **large-scale** data
- train for a **long** time
- with **large** models

# Transfer learning: Pre-training & Fine-tuning

## Pre-training



## Fine-tuning



### Fine-tuning:

- transfer weights to **specific** tasks
- on **small-scale** data
- train for a **short** time, **lower** learning rate
- enable **large** models with lower risk of overfitting

# Transfer learning: Pre-training & Fine-tuning

## Pre-training



## Fine-tuning



## Partial transfer

- pre-train and target domains may differ
- highest-level features are too adapted to pre-training
- randomly initialize new layers

# Transfer learning: Pre-training & Fine-tuning

## Pre-training



## Fine-tuning



## Frozen weights

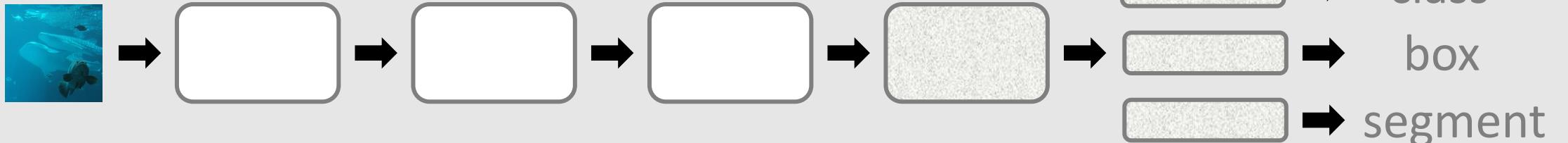
- freeze some/all pre-trained weights
- reduce overfitting if data is too little
- save memory, speed up training

# Transfer learning: Pre-training & Fine-tuning

## Pre-training



## Fine-tuning



## Network surgery

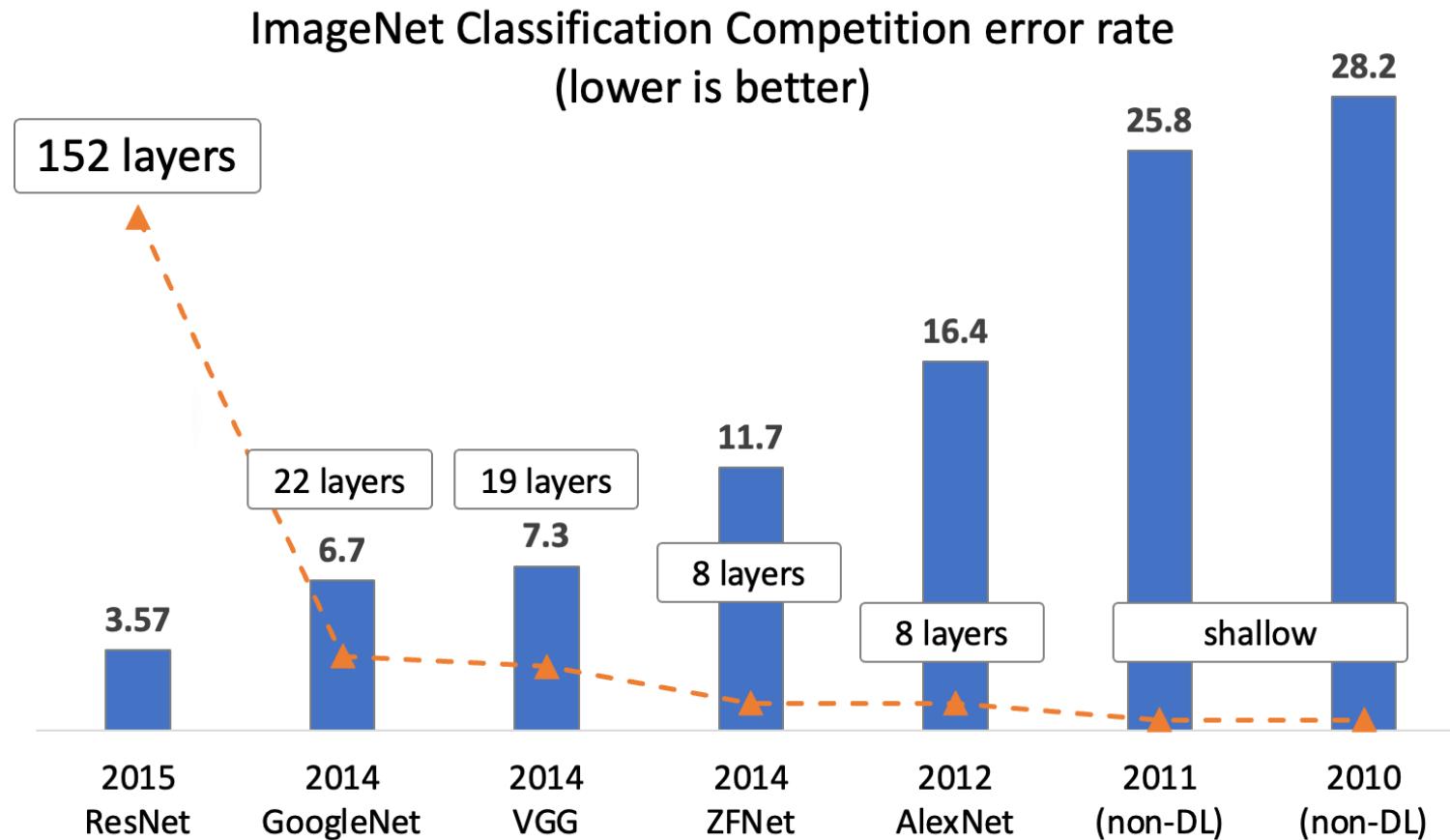
- re-purpose the model for other tasks (detect, segment)
- general features + task-specific predictions

# How can we (pre-)train good representations?

- **Model** (network architecture)
  - scaling: deep, wide, large

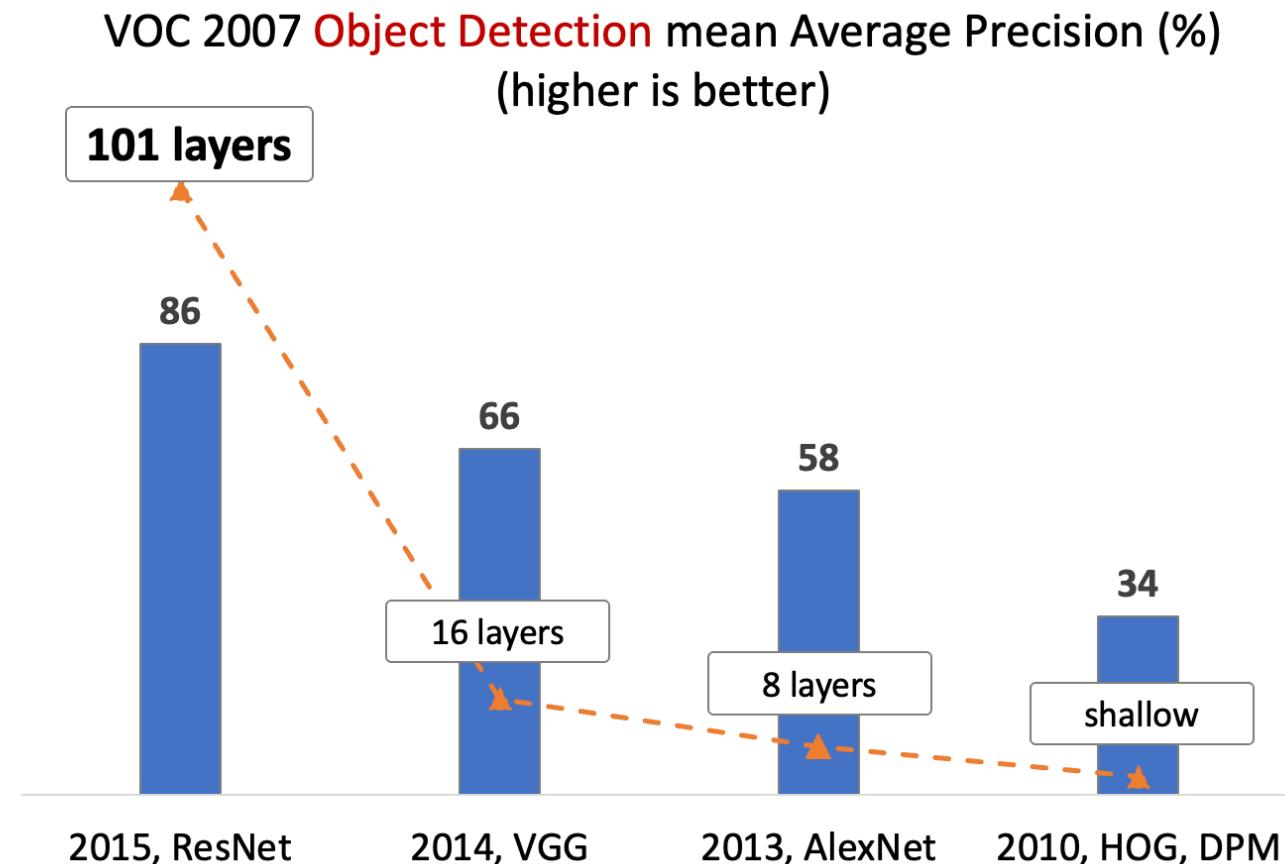
# How can we (pre-)train good representations?

- **Model** (network architecture)
  - scaling: deep, wide, large



# How can we (pre-)train good representations?

- **Model** (network architecture)
  - scaling: deep, wide, large



# How can we (pre-)train good representations?

- **Model** (network architecture)
  - scaling: deep, wide, large
  - inductive bias: convolution, recurrency, attention

# How can we (pre-)train good representations?

- **Model** (network architecture)
  - scaling: deep, wide, large
  - inductive bias: convolution, recurrency, attention
- **Data**
  - scaling
  - curating, cleaning, filtering, ...

# How can we (pre-)train good representations?

- **Model** (network architecture)
  - scaling: deep, wide, large
  - inductive bias: convolution, recurrency, attention
- **Data**
  - scaling
  - curating, cleaning, filtering, ...
- **Learning objective**
  - supervised
  - unsupervised
  - self-supervised

# Pre-training Objective

- **Supervised Learning**
  - w/ human annotated labels (e.g., classes)
- **Unsupervised Learning**
  - w/o human annotated labels
- **Self-supervised Learning**
  - labels induced by data “itself”
  - modern “unsupervised learning” in camouflage



remaining of  
this lecture

# **Unsupervised Learning**

# Unsupervised Learning

- Conceptually, “**unsupervised learning**” involves any learning method **without** using human labels.
- In the modern literature, “**unsupervised learning**” often refers to a family of classical methods for:
  - **clustering**
  - **dimension reduction**
- As we will show, even classical “**unsupervised learning**” methods are kind of “**self-supervised**”.

# Unsupervised Learning

Classical (but very useful) unsupervised learning methods:

- K-means clustering
- Principal Component Analysis (PCA)
- Autoencoder

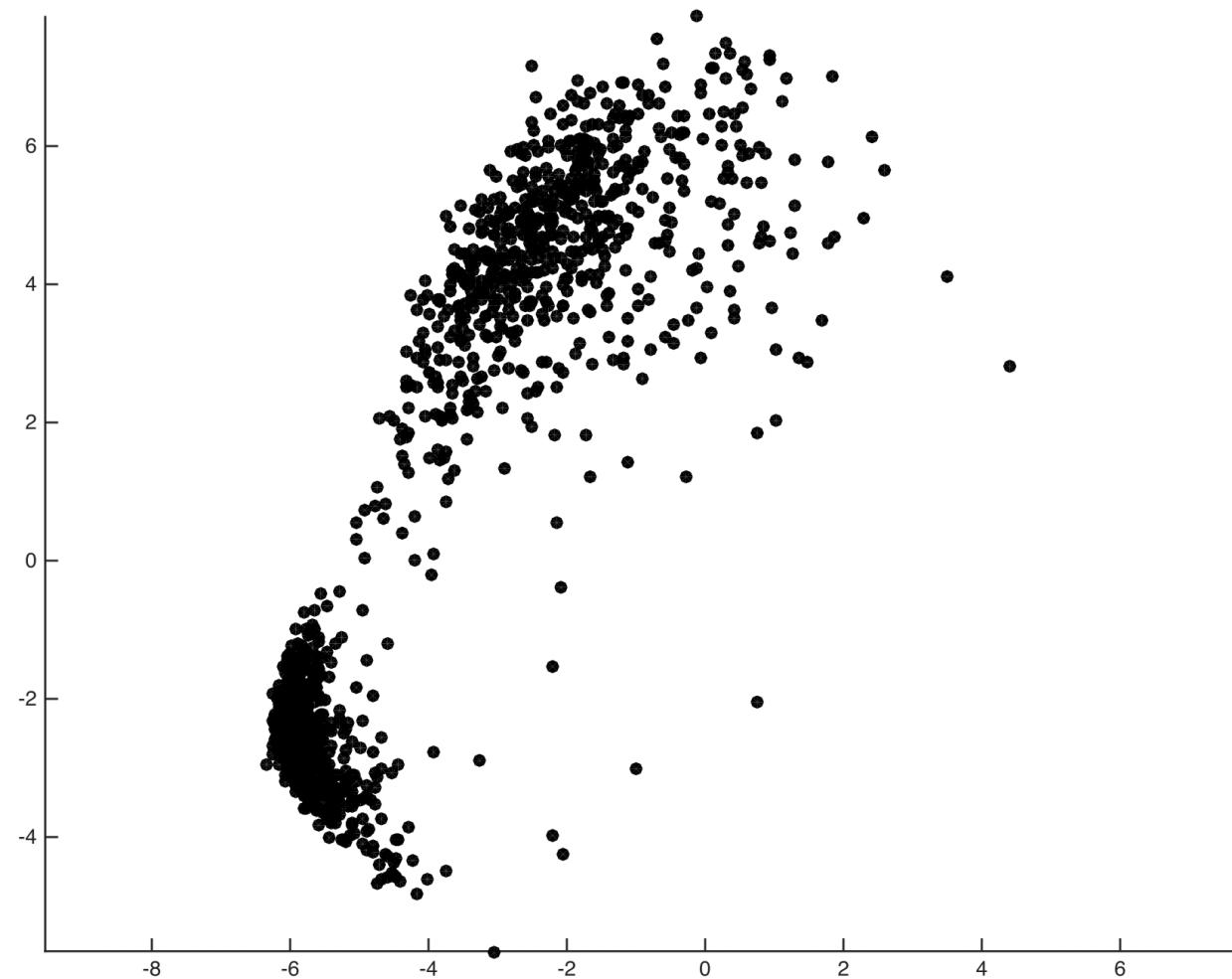
# Unsupervised Learning

Classical (but very useful) unsupervised learning methods:

- **K-means clustering**
- Principal Component Analysis (PCA)
- Autoencoder

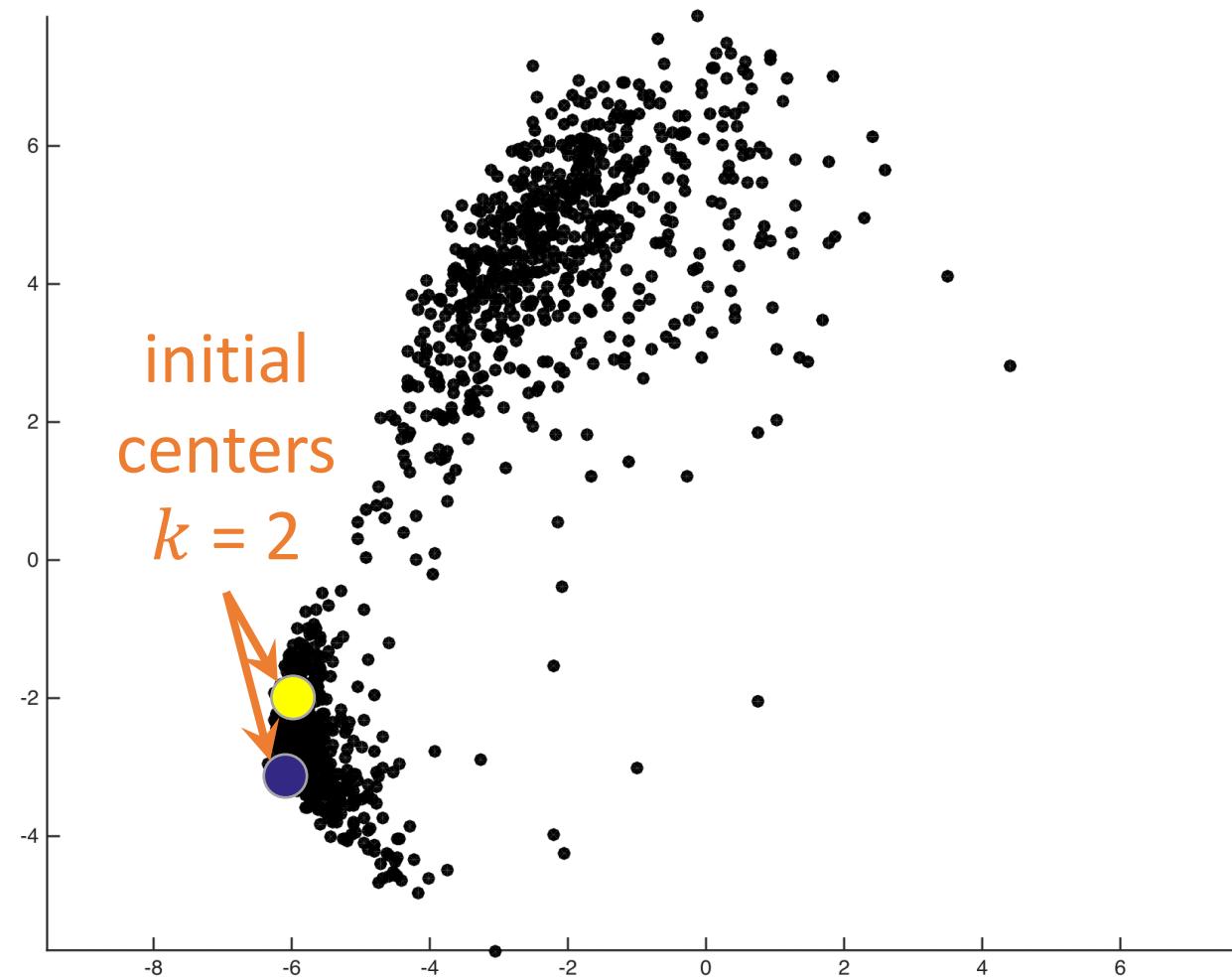
# K-means clustering

Represent data by  $k$  centers



# K-means clustering

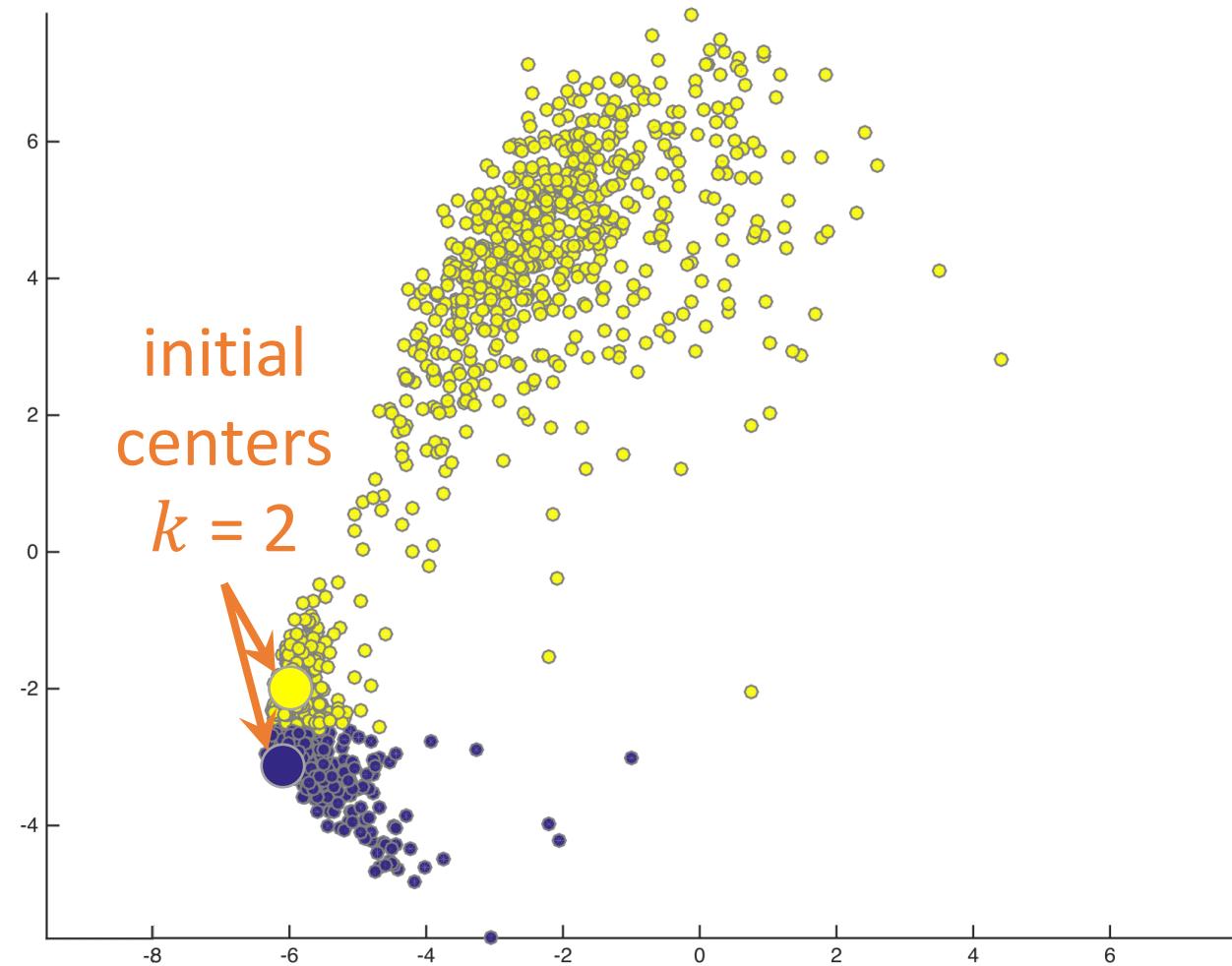
Represent data by  $k$  centers



# K-means clustering

Represent data by  $k$  centers

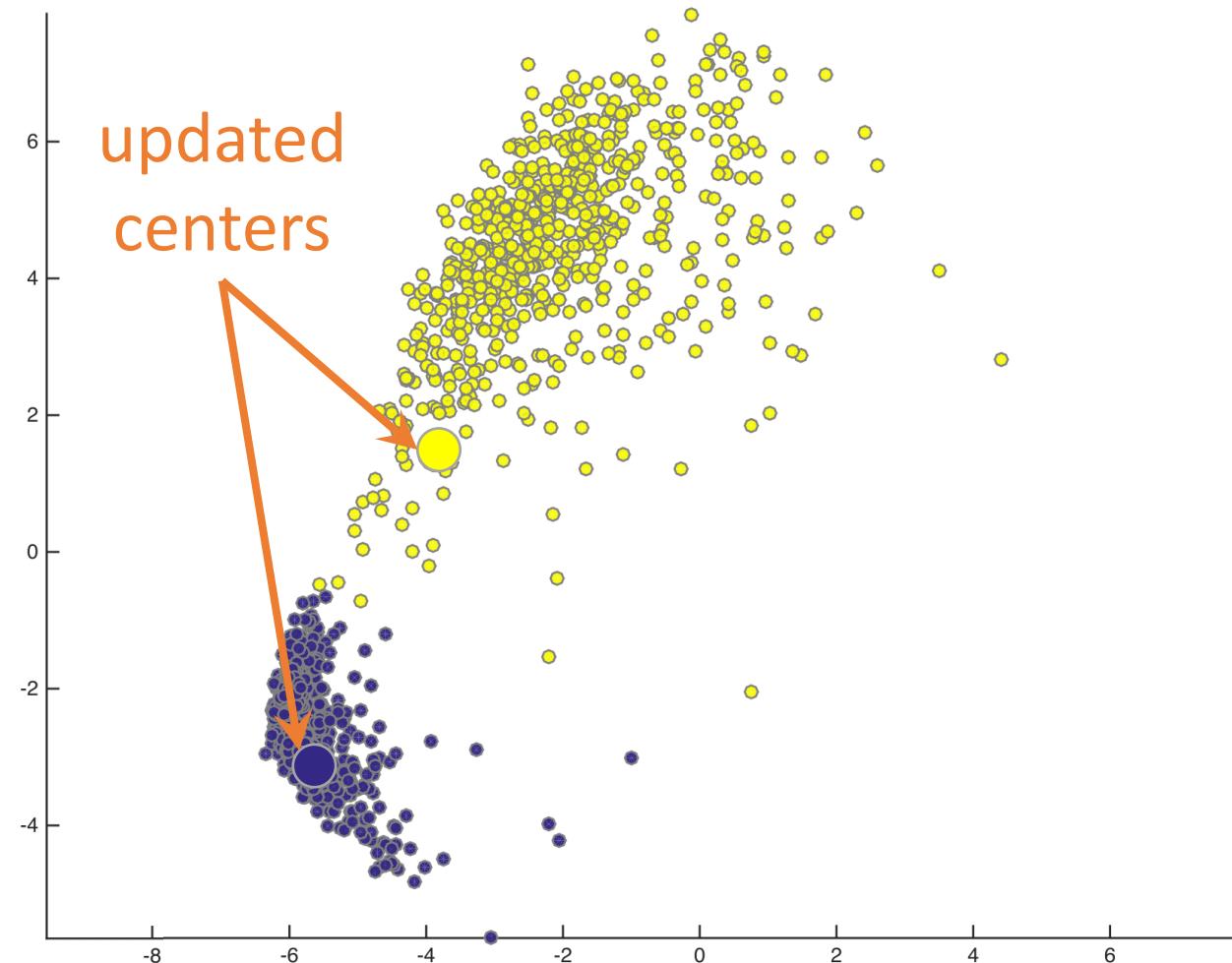
- **Assignment:** each sample is assigned to the nearest center



# K-means clustering

Represent data by  $k$  centers

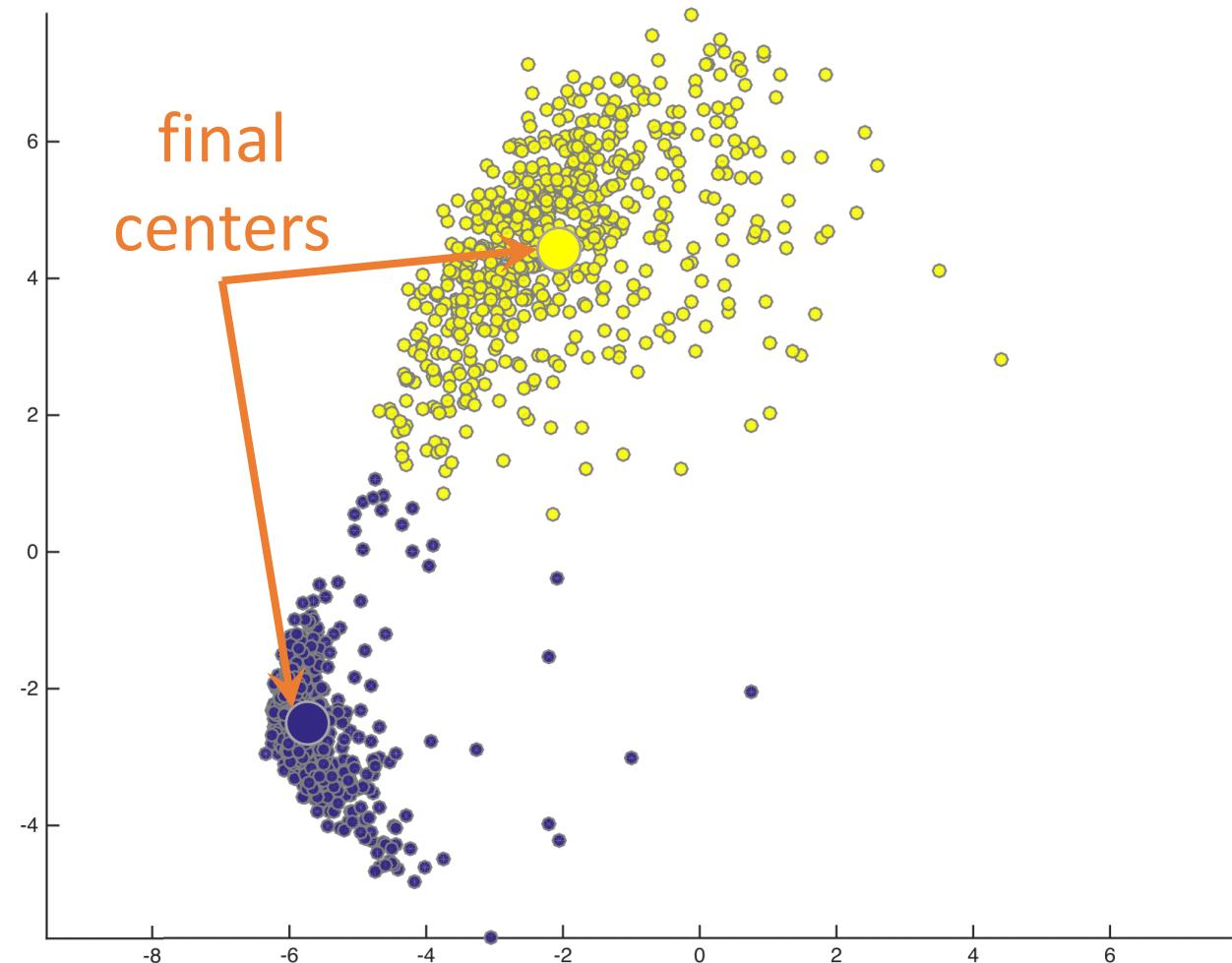
- **Assignment:** each sample is assigned to the nearest center
- **Update:** each center is updated by the mean of samples assigned to it



# K-means clustering

Represent data by  $k$  centers

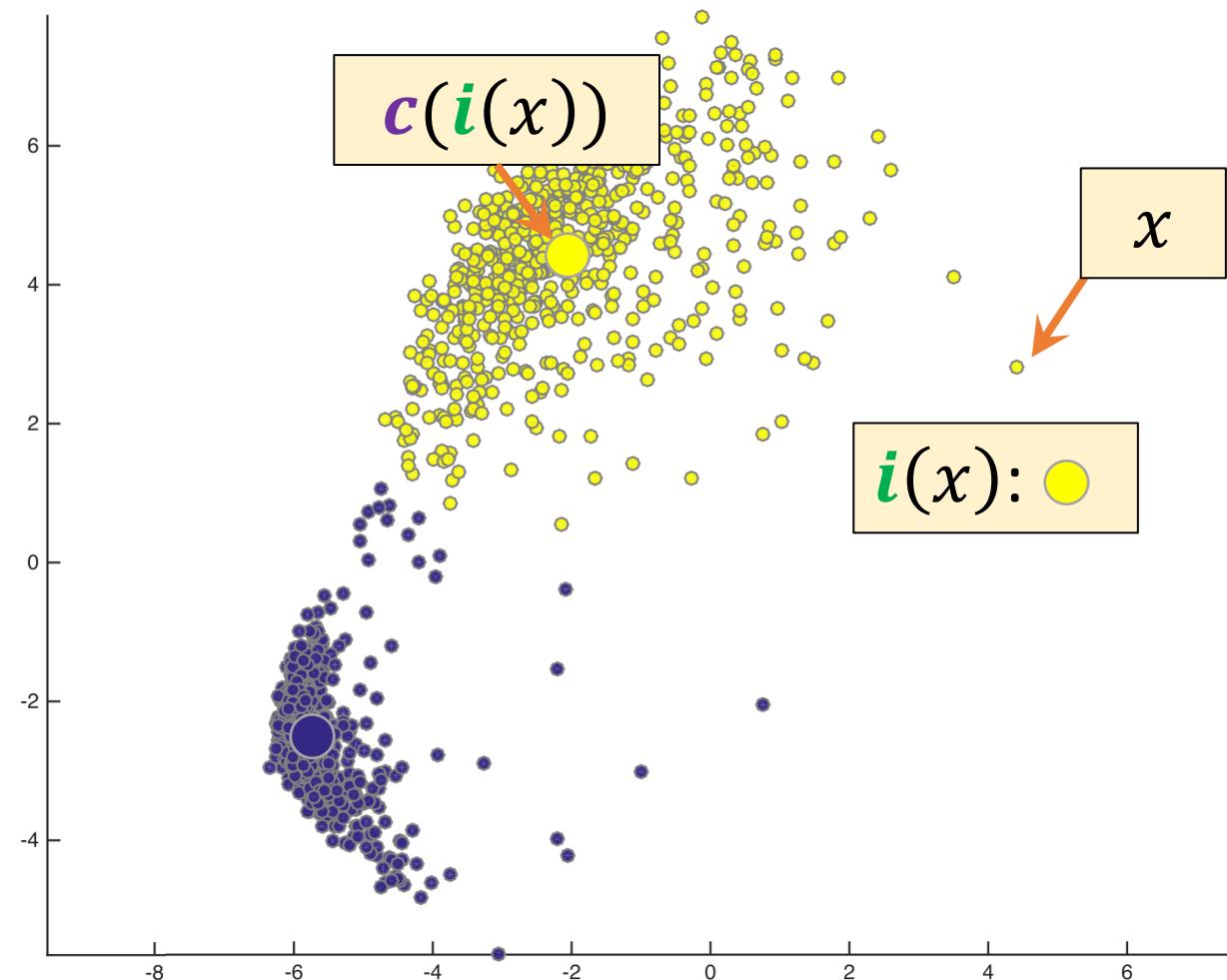
- **Assignment:** each sample is assigned to the nearest center
- **Update:** each center is updated by the mean of samples assigned to it
- Repeat until converge



# K-means clustering

$i(x)$ : the cluster index of sample  $x$   
 $c(i)$ : the  $i$ -th center

Loss function:

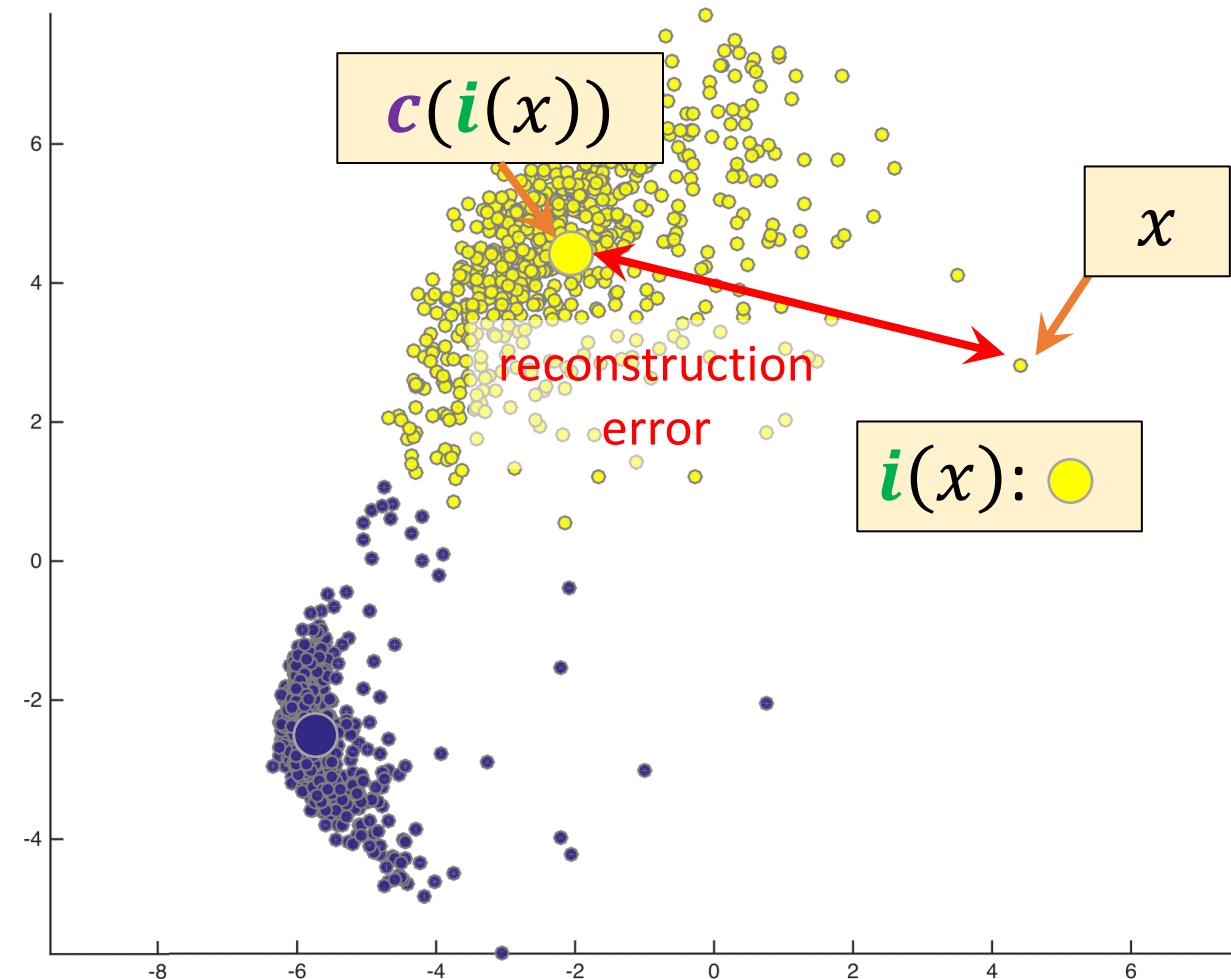


# K-means clustering

Loss function:

$$\min_{\mathbf{c}, \mathbf{i}} \sum_x \|x - \mathbf{c}(\mathbf{i}(x))\|^2$$

$\mathbf{i}(x)$ : the cluster index of sample  $x$   
 $\mathbf{c}(i)$ : the  $i$ -th center



# K-means clustering

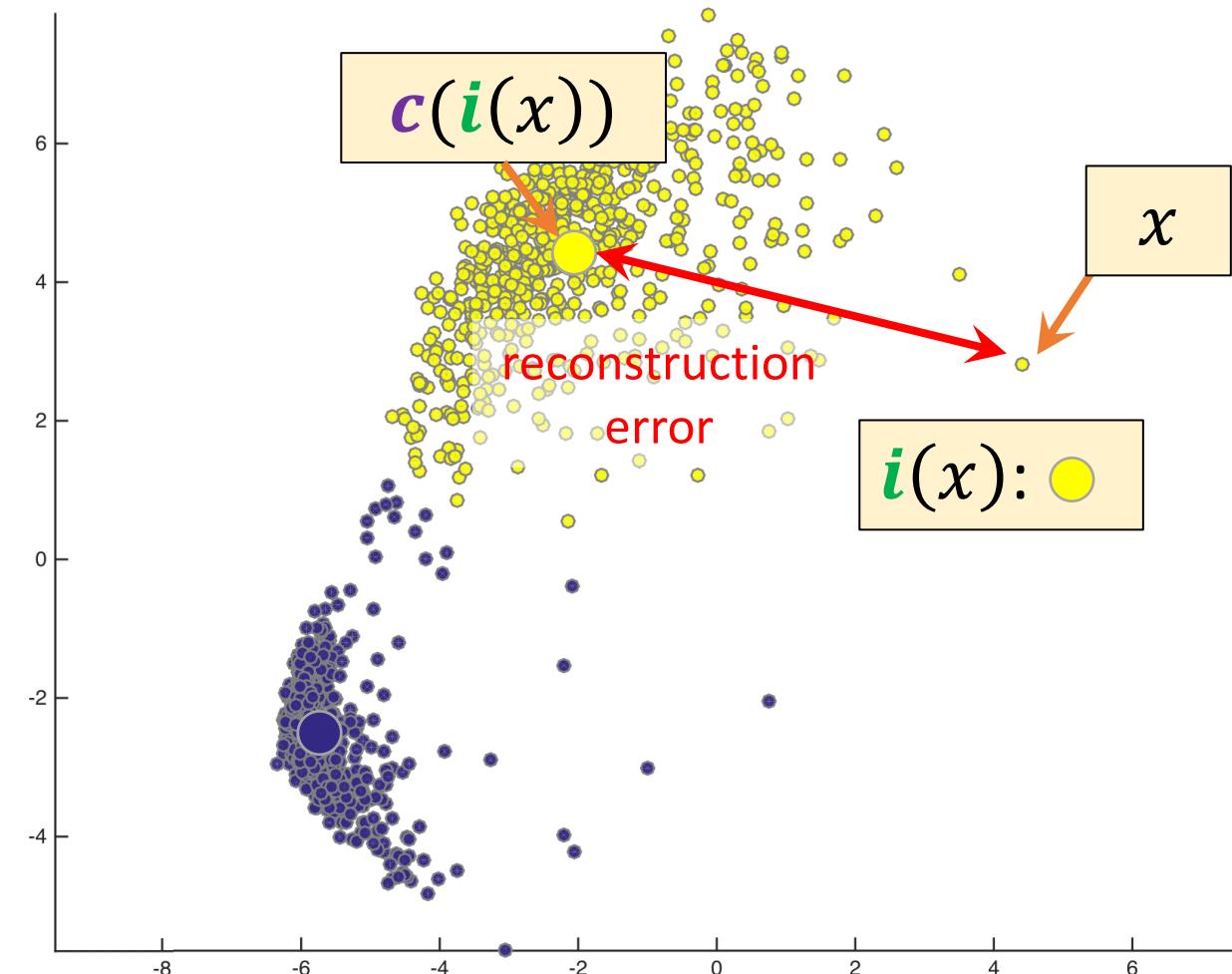
Loss function:

$$\min_{\mathbf{c}, \mathbf{i}} \sum_x \|x - \mathbf{c}(\mathbf{i}(x))\|^2$$

Assignment: fixed  $\mathbf{c}$ , optimize  $\mathbf{i}$

$$\mathbf{i}(x) = \operatorname{argmin}_{\mathbf{i}} \|x - \mathbf{c}(\mathbf{i}(x))\|^2$$

$\mathbf{i}(x)$ : the cluster index of sample  $x$   
 $\mathbf{c}(i)$ : the  $i$ -th center



# K-means clustering

**Loss function:**

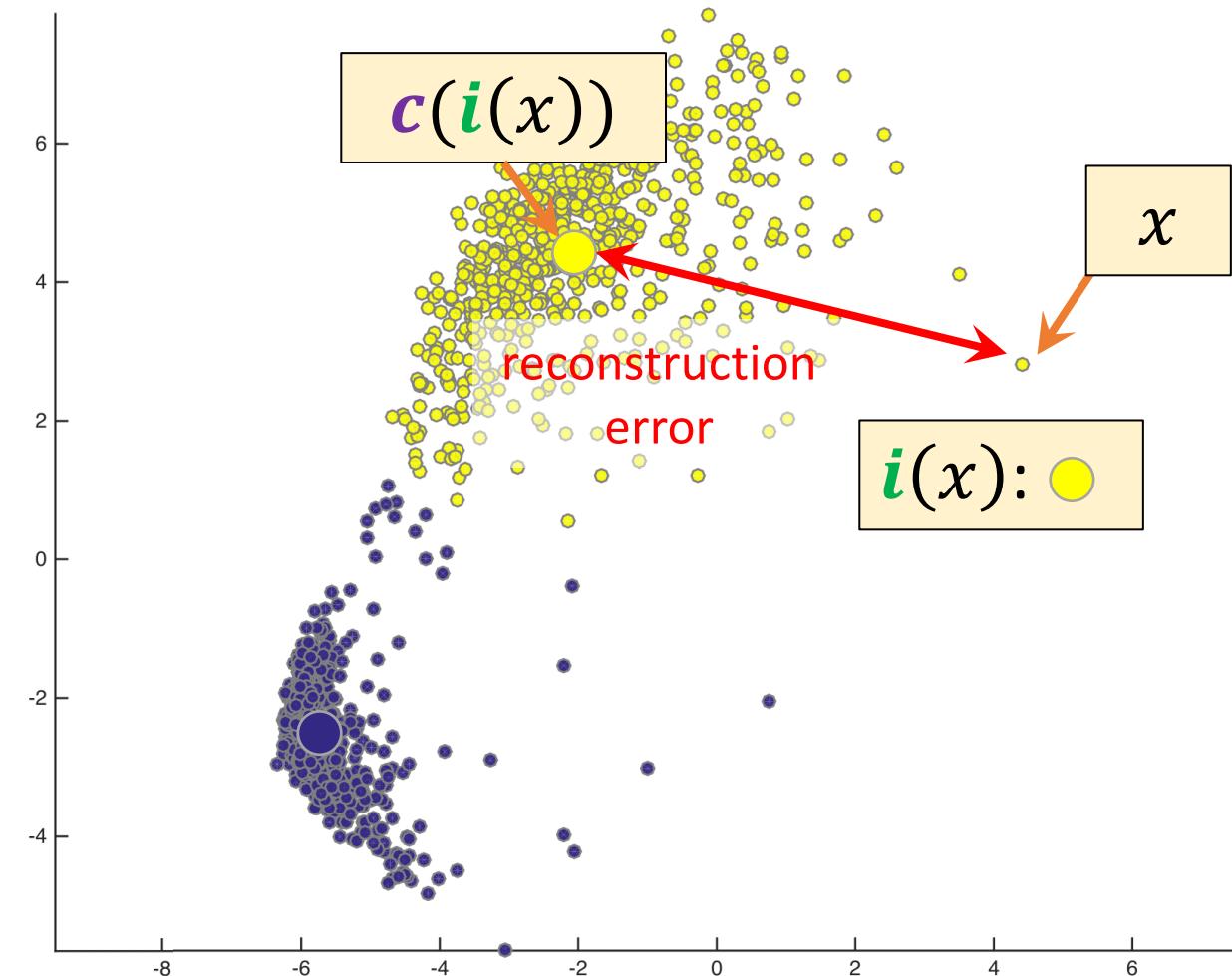
$$\min_{\mathbf{c}, \mathbf{i}} \sum_x \|x - \mathbf{c}(\mathbf{i}(x))\|^2$$

**Assignment:** fixed  $\mathbf{c}$ , optimize  $\mathbf{i}$

**Update:** fixed  $\mathbf{i}$ , optimize  $\mathbf{c}$

$$\mathbf{c}(j) = \operatorname{argmin}_{\mathbf{c}(j)} \sum_{x: \mathbf{i}(x)=j} \|x - \mathbf{c}(j)\|^2$$

$\mathbf{i}(x)$ : the cluster index of sample  $x$   
 $\mathbf{c}(i)$ : the  $i$ -th center



# K-means clustering

**Loss function:**

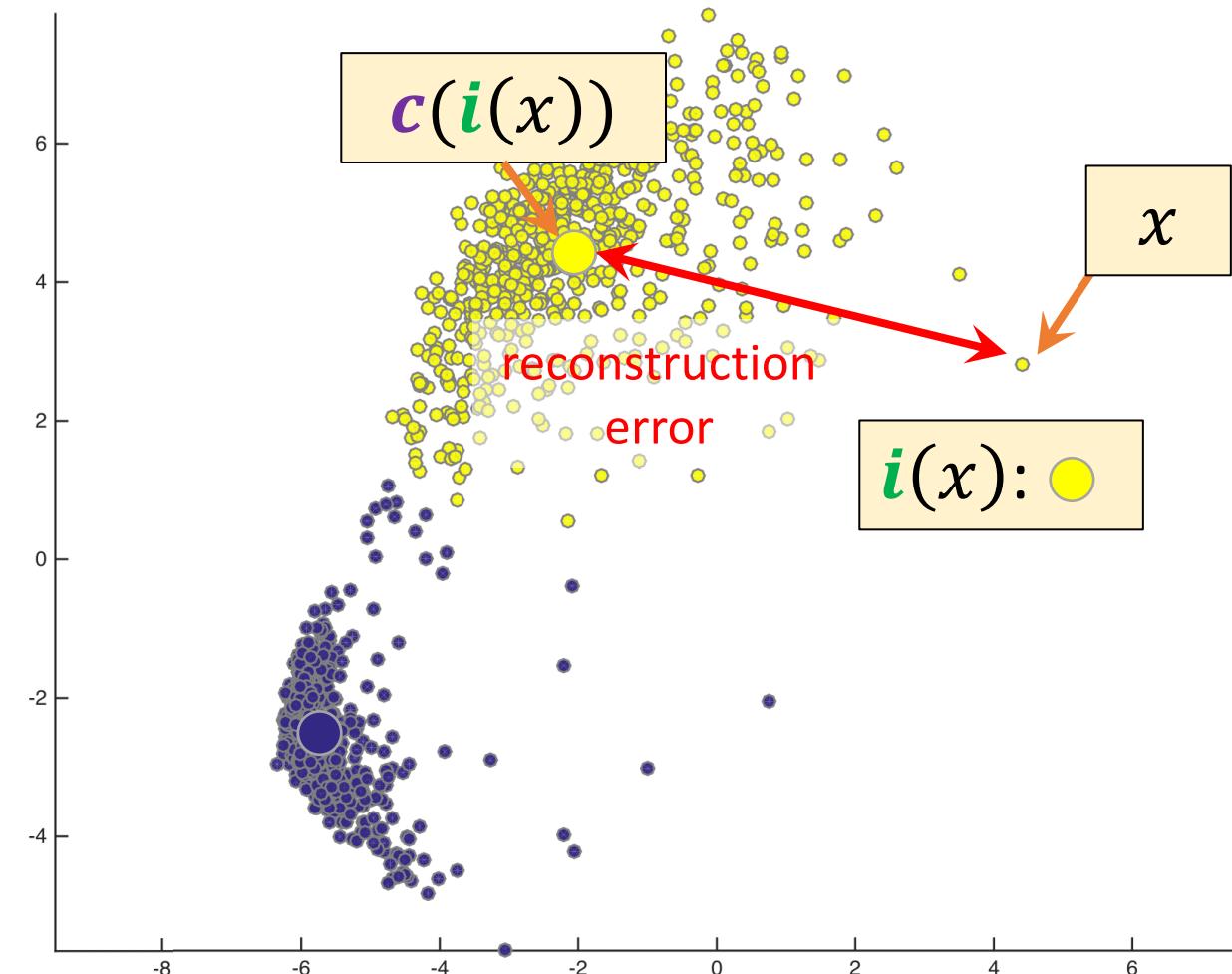
$$\min_{\mathbf{c}, \mathbf{i}} \sum_x \|x - \mathbf{c}(i(x))\|^2$$

**Assignment:** fixed  $\mathbf{c}$ , optimize  $\mathbf{i}$

**Update:** fixed  $\mathbf{i}$ , optimize  $\mathbf{c}$

Repeat until converge

$i(x)$ : the cluster index of sample  $x$   
 $c(i)$ : the  $i$ -th center



# K-means clustering

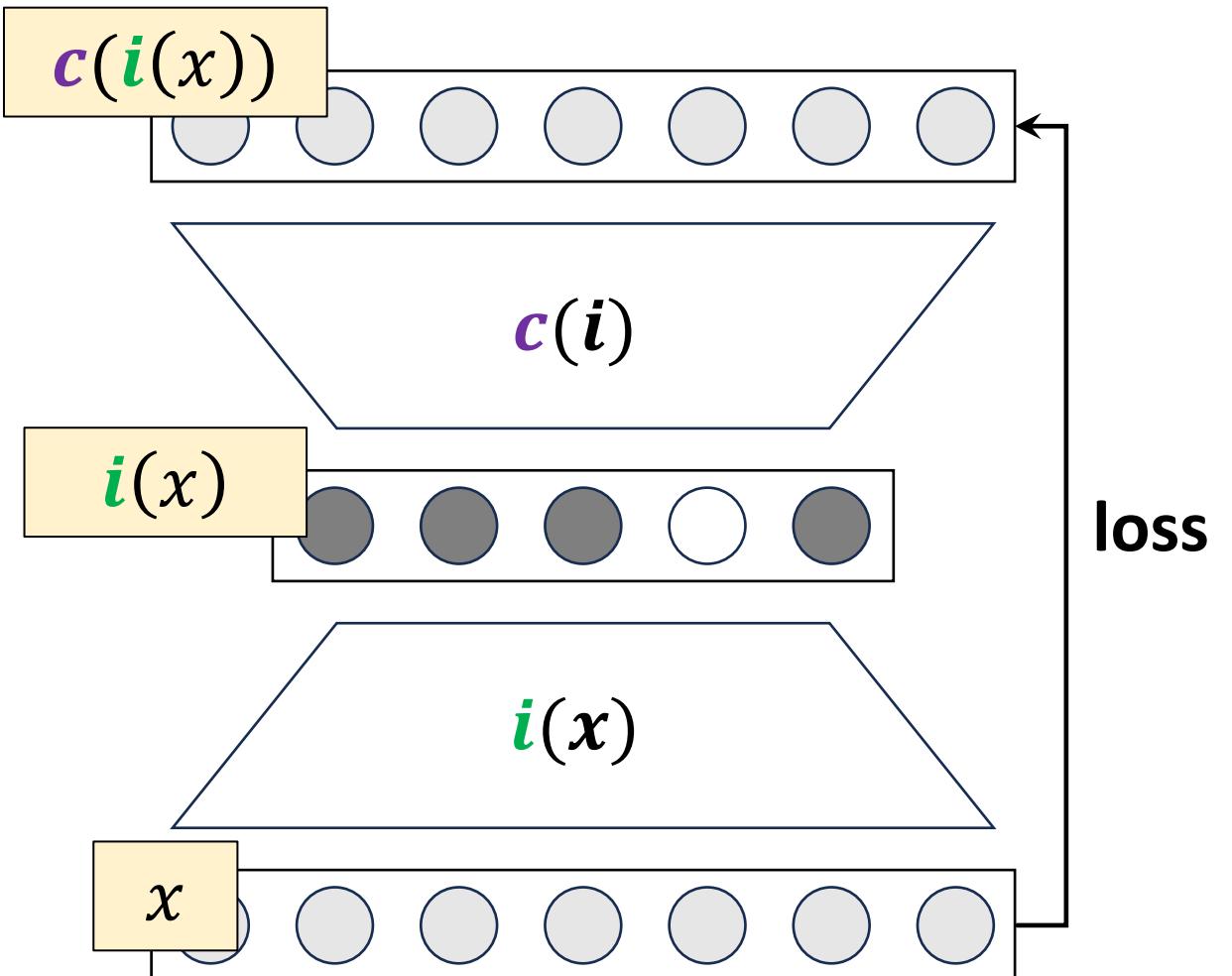
Loss function:

$$\min_{\mathbf{c}, \mathbf{i}} \sum_x \|x - \mathbf{c}(i(x))\|^2$$

Assignment: fixed  $\mathbf{c}$ , optimize  $\mathbf{i}$

Update: fixed  $\mathbf{i}$ , optimize  $\mathbf{c}$

Repeat until converge



# K-means clustering

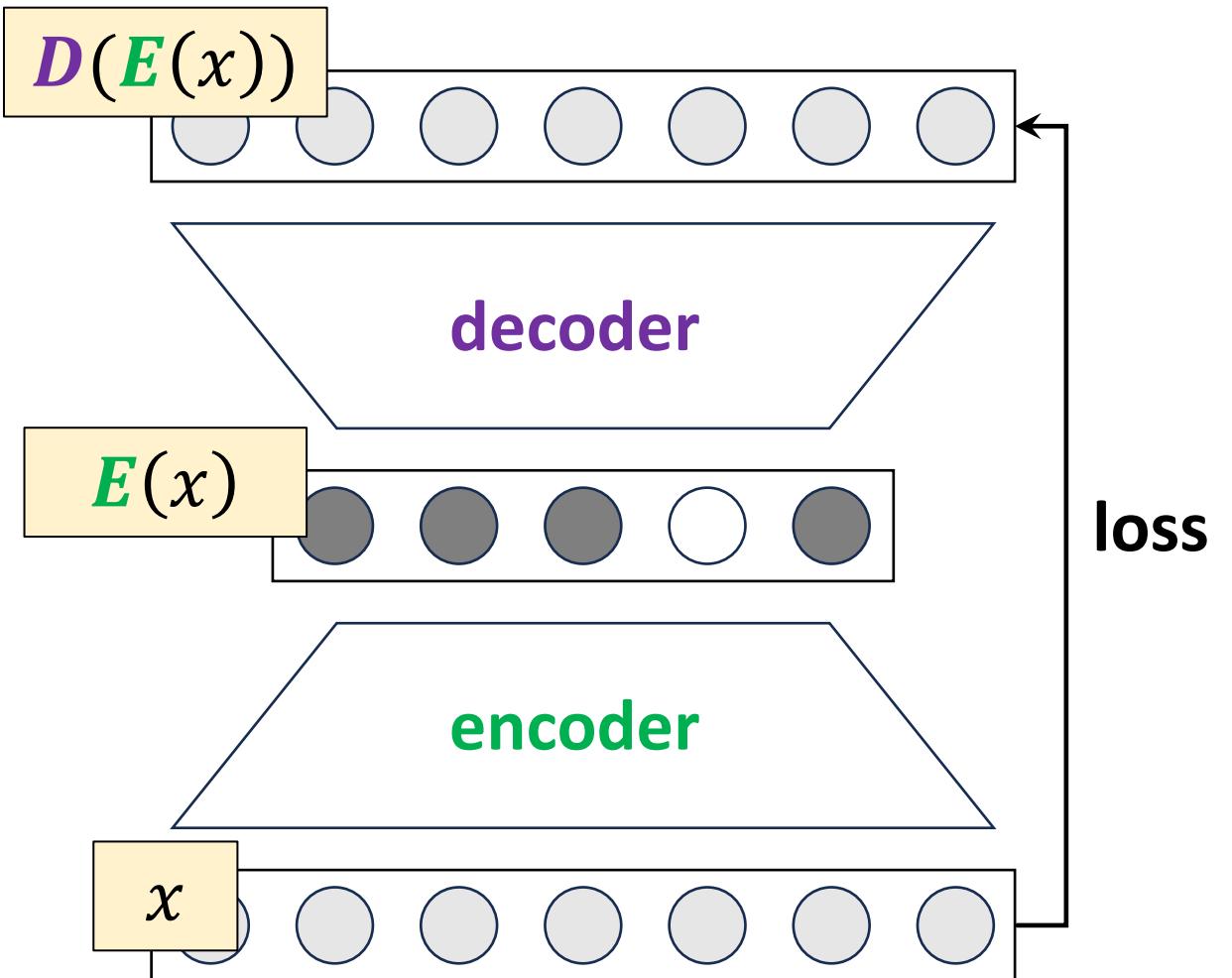
Loss function:

$$\min_{\mathbf{D}, \mathbf{E}} \sum_x \|x - \mathbf{D}(\mathbf{E}(x))\|^2$$

encoder  $\mathbf{E}$ : map  $x$  to a (one-hot) code

decoder  $\mathbf{D}$ : map a code to a center

K-means is like Autoencoding.

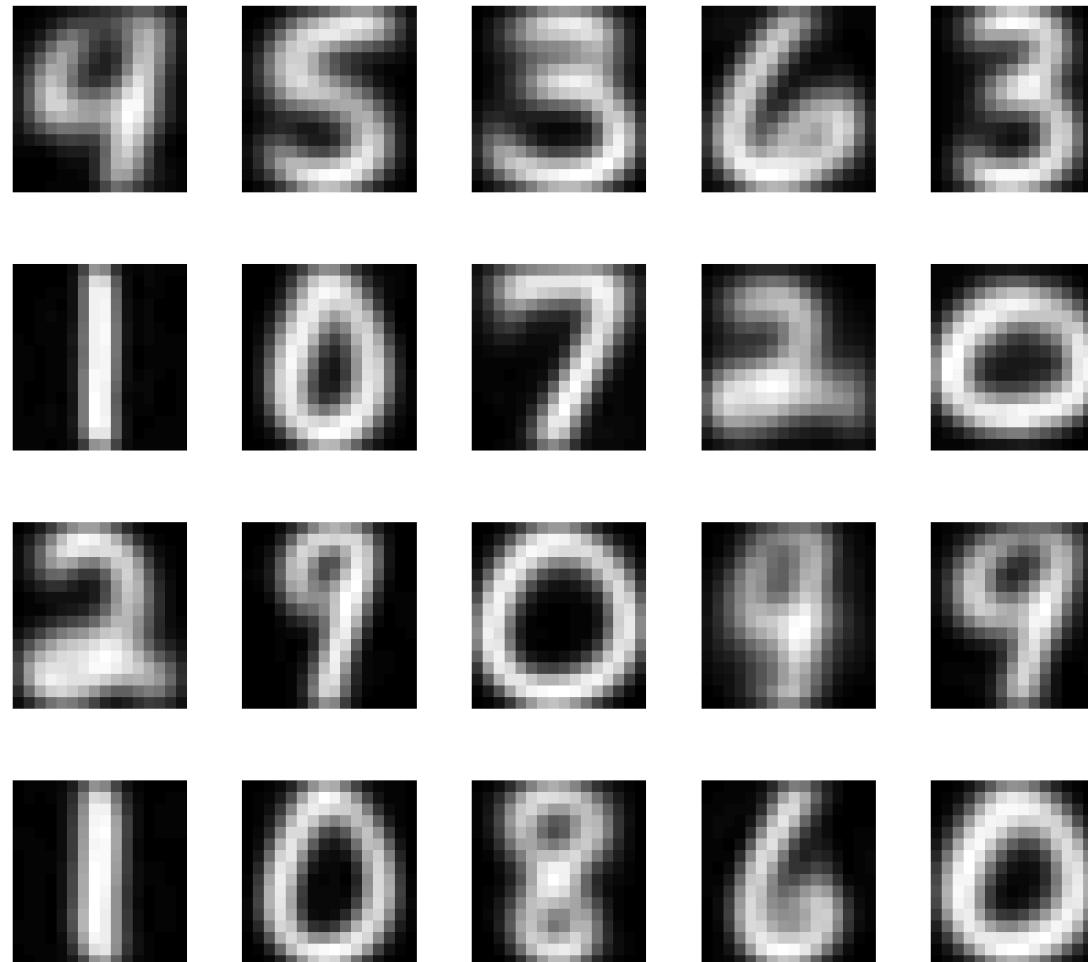


# Example: K-means on MNIST digits

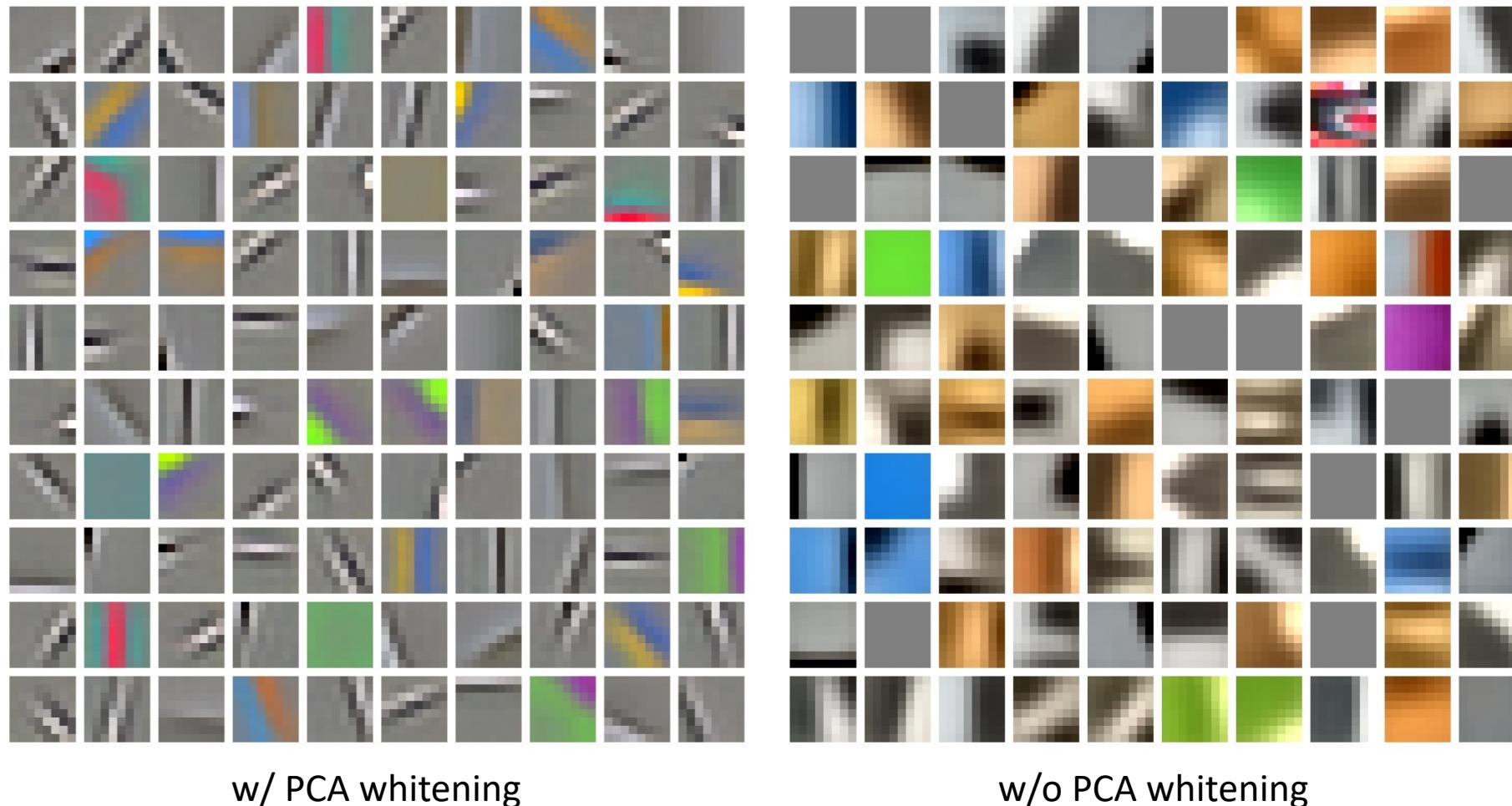
Visualized centers with  $k = 20$

These are not actual digit images.

These are the learned “representations”.



# Example: K-means on 8x8 image patches

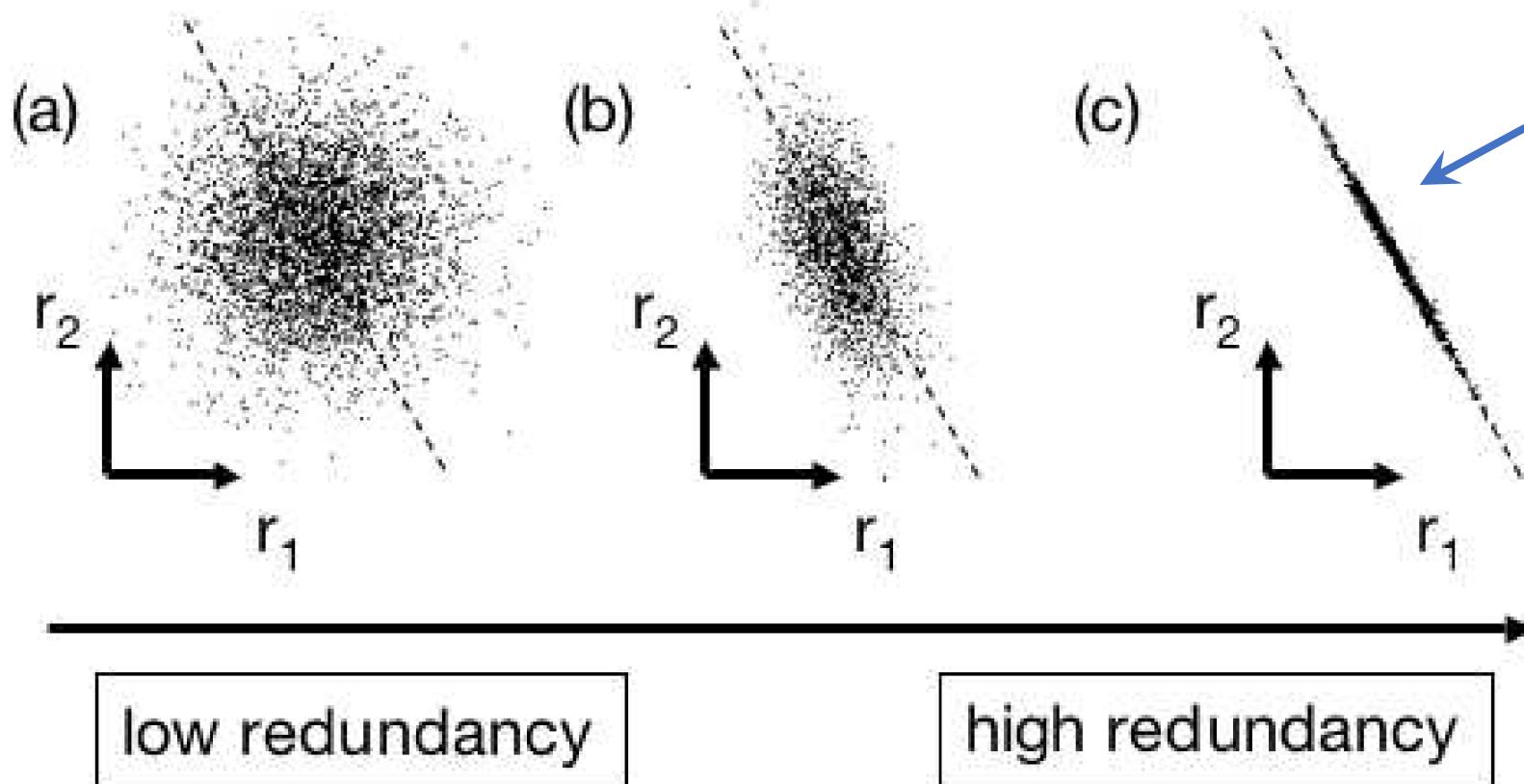


# Unsupervised Learning

Classical (but very useful) unsupervised learning methods:

- K-means clustering
- **Principal Component Analysis (PCA)**
- Autoencoder

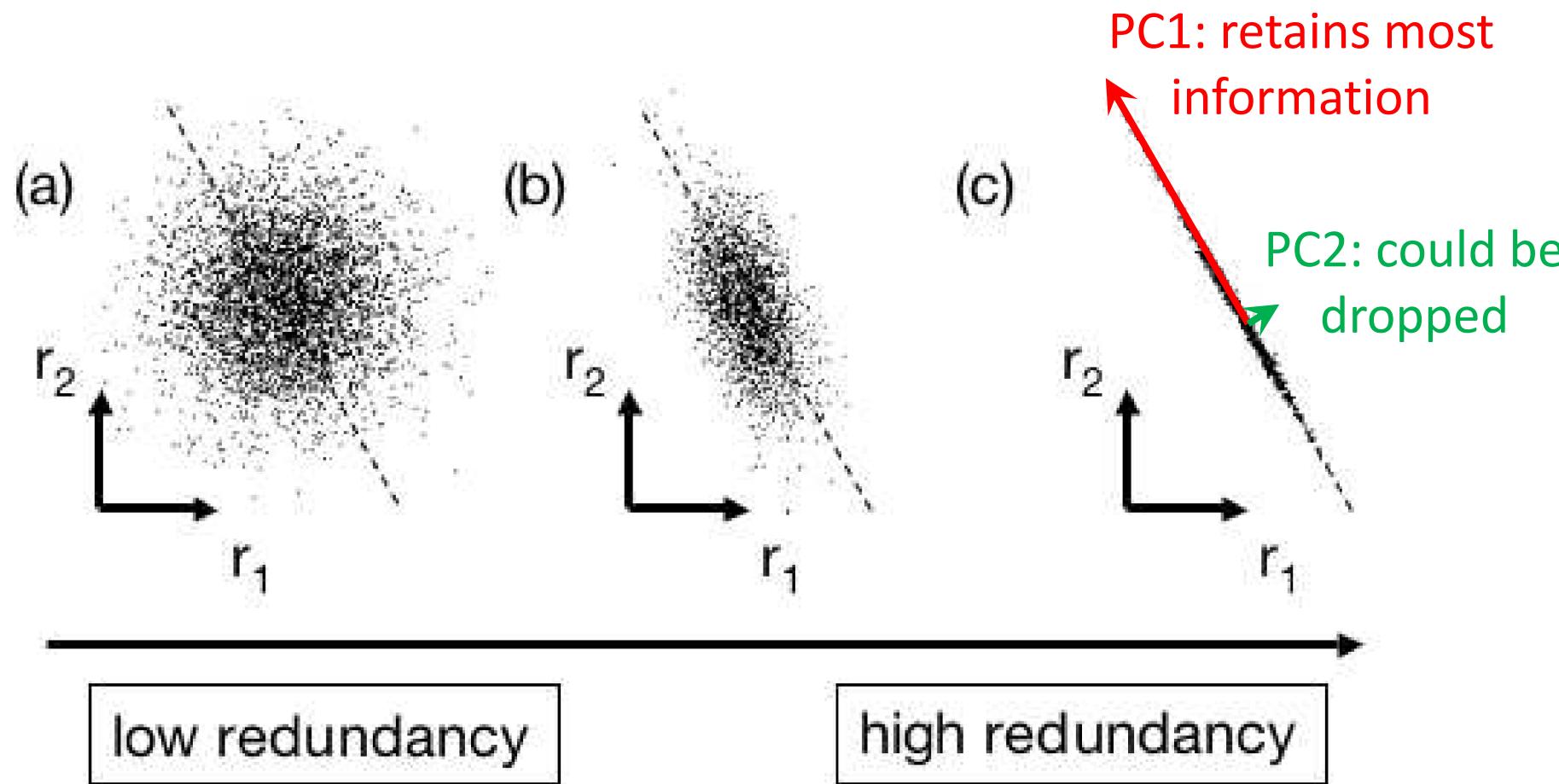
# Principal Component Analysis (PCA)



In this case, one can use  $r_1$  to represent  $r_2$  approximately.

So,  $(r_1, r_2)$  is not a good representation.

# Principal Component Analysis (PCA)



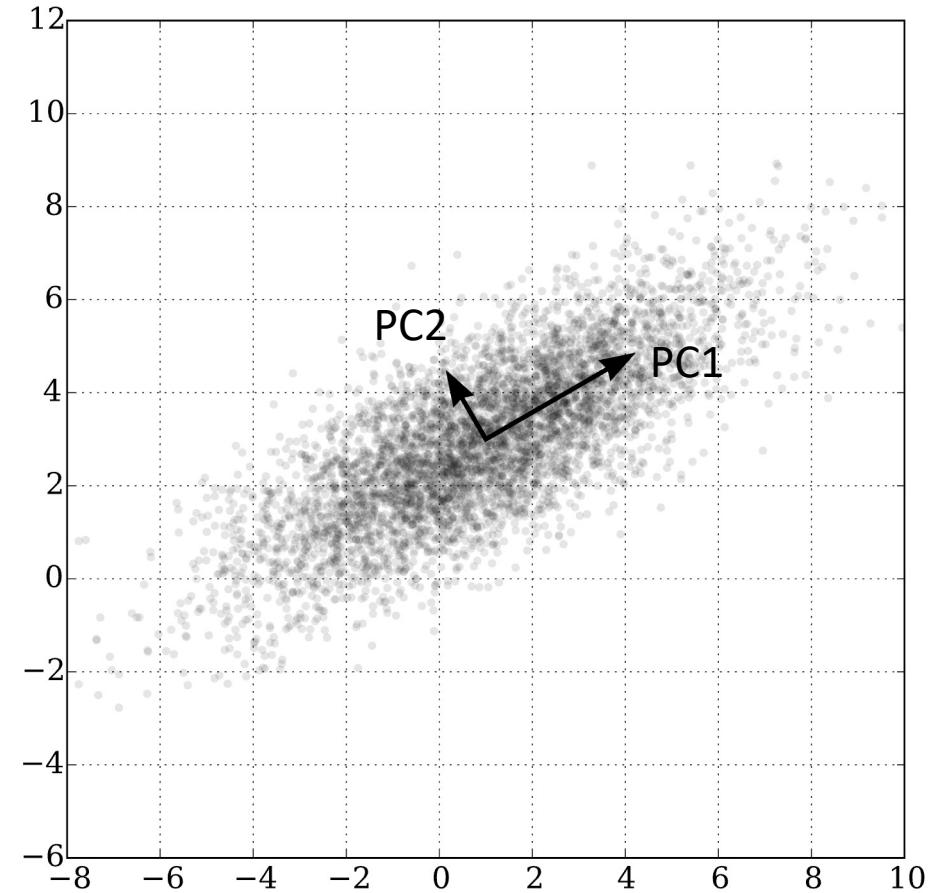
# Principal Component Analysis (PCA)

Represent data in the space of the largest  $k$  Principal Components (PCs)

$$\begin{aligned} \min_W \sum_x \|x - W^T W x\|^2 \\ s.t. \quad W W^T = I_{k \times k} \end{aligned}$$

$x$ :  $d$ -dimensional vector

$W$ :  $k \times d$  matrix



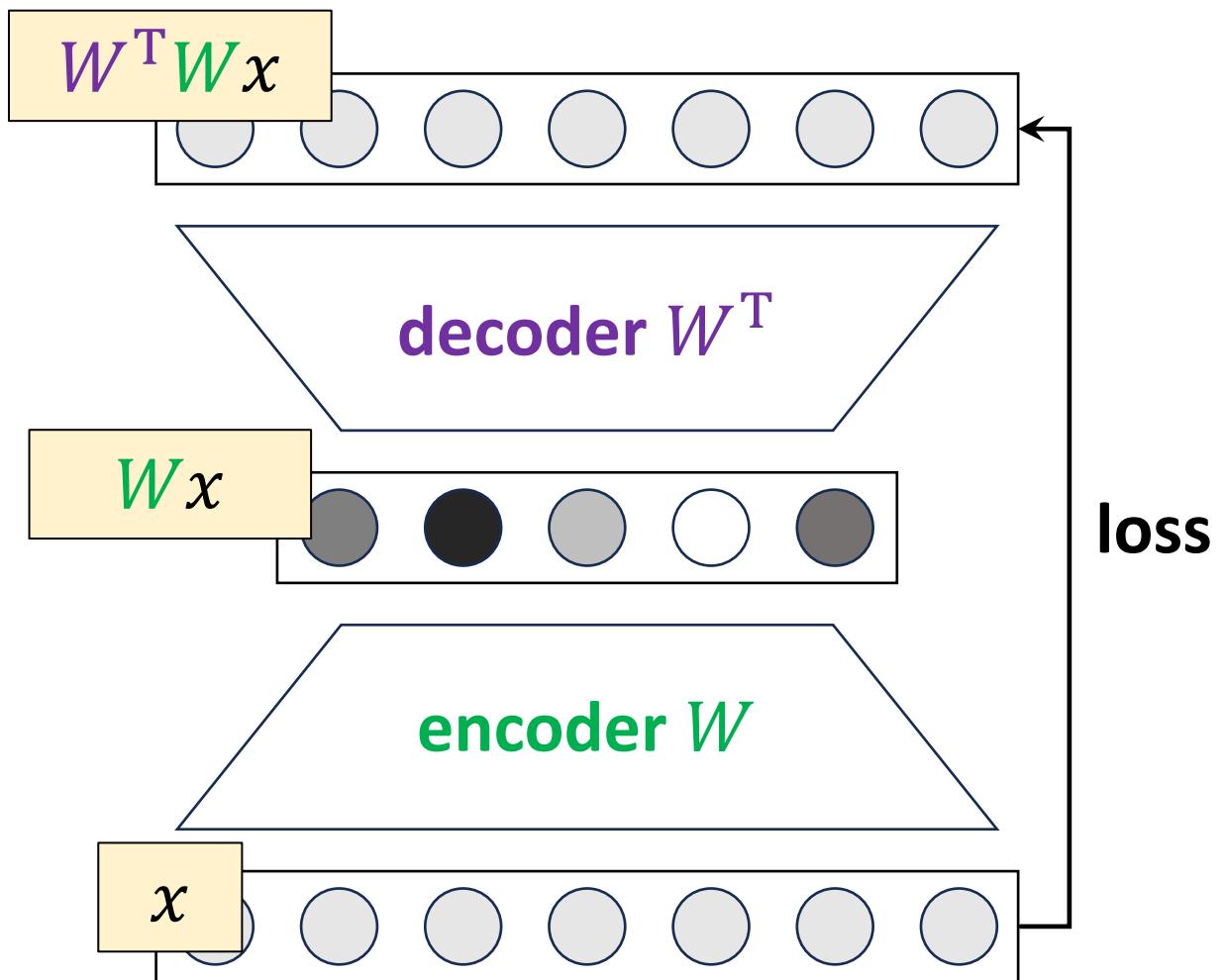
# Principal Component Analysis (PCA)

Represent data in the space of the largest  $k$  Principal Components (PCs)

$$\begin{aligned} \min_W \sum_x \|x - W^T W x\|^2 \\ s.t. \quad W W^T = I_{k \times k} \end{aligned}$$

**encoder  $W$** : project onto PCs  
**decoder  $W^T$** : project back

PCA is like Autoencoding.  
• w/ linear encoder & decoder  
• w/ orthogonal constraint



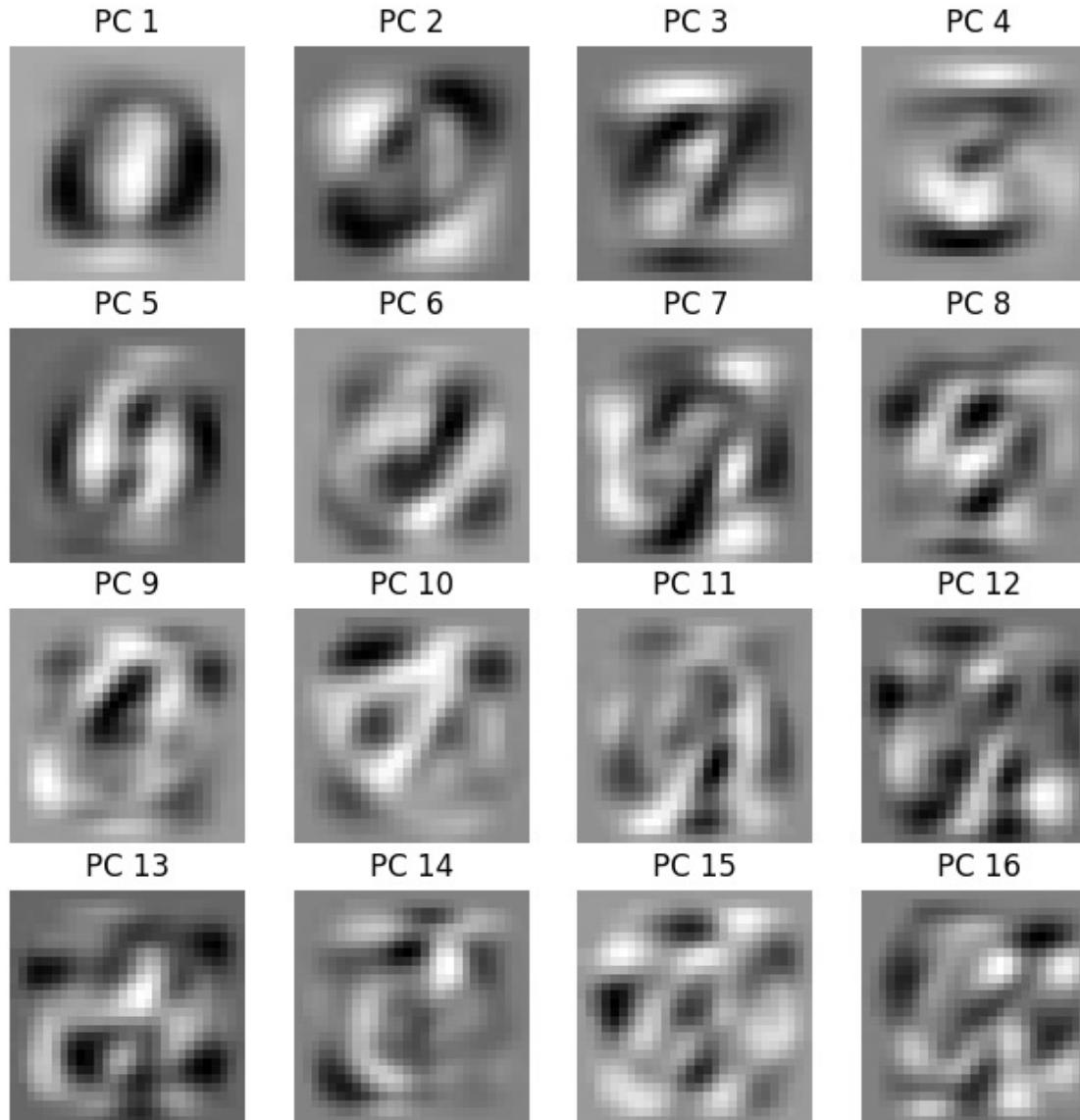
# Example: PCA on face data (Eigenfaces)

These are not actual  
face images.

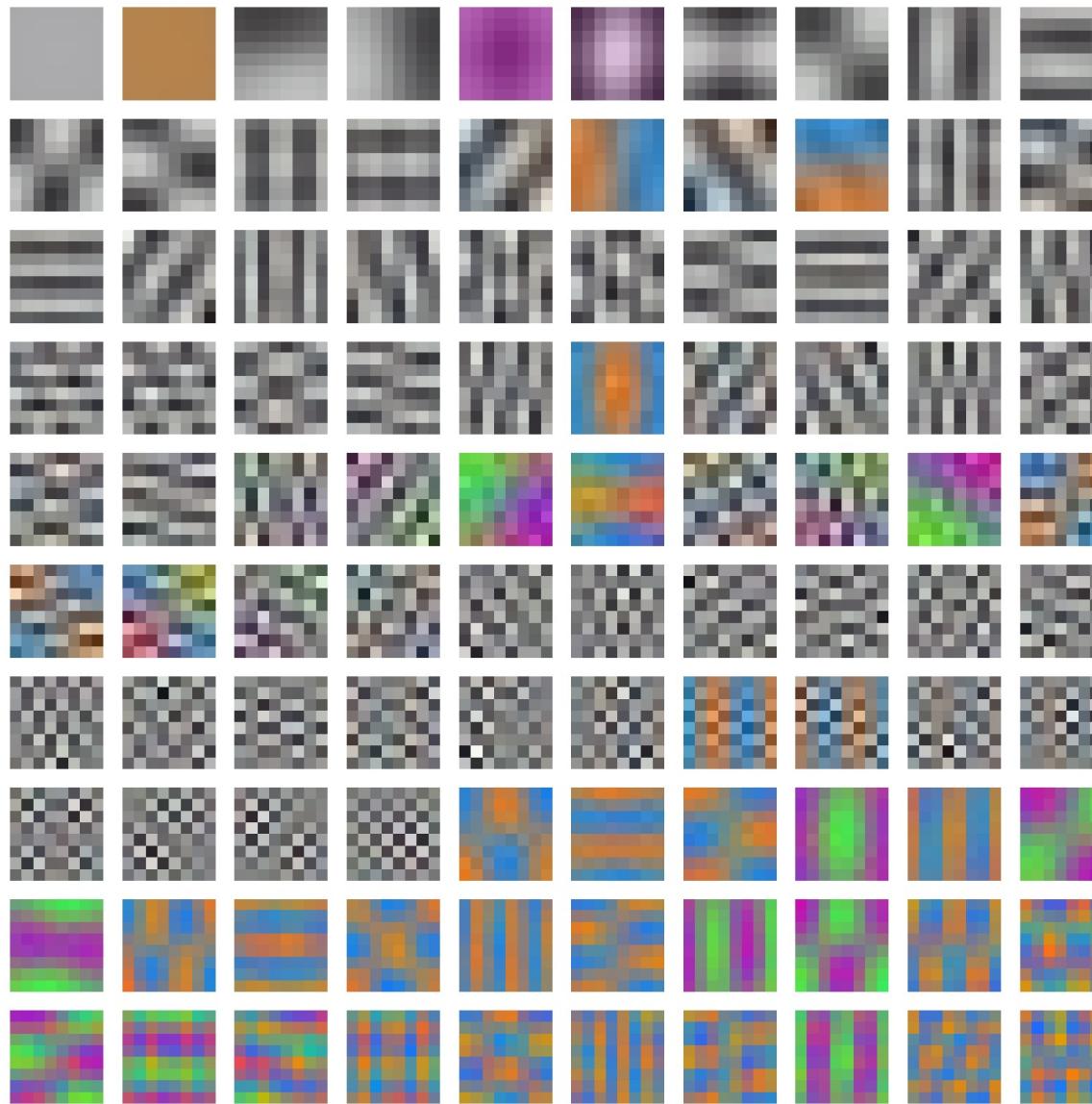
These are the learned  
“representations”.



# Example: PCA on MNIST digits



# Example: PCA on 8x8 image patches



# Unsupervised Learning

Classical (but very useful) unsupervised learning methods:

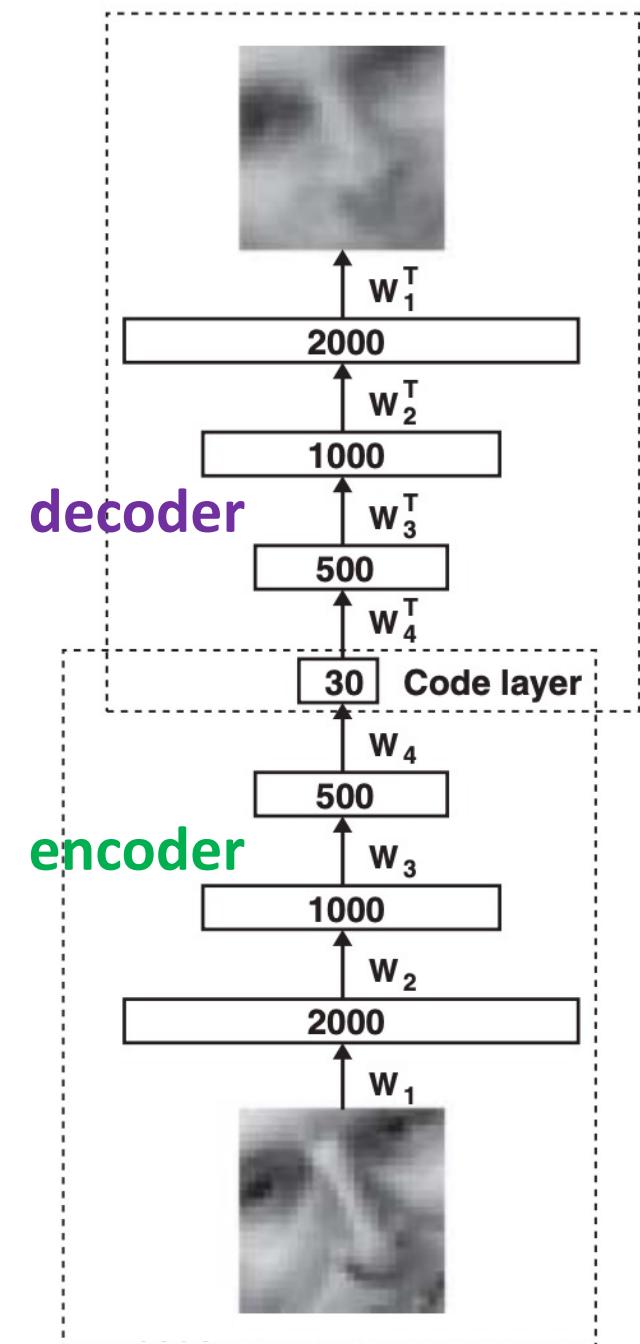
- K-means clustering
- Principal Component Analysis (PCA)
- **Autoencoder**

# Autoencoder

Represent data by compressing it into a low-dimensional latent space

- **encoder & decoder**: deep neural nets
- optimized by BackProp
- a lower-dimensional latent as a “bottleneck”
- reconstruct the input

$$\min_{\mathcal{D}, \mathcal{E}} \sum_x \|\mathcal{D}(\mathcal{E}(x)) - x\|^2$$

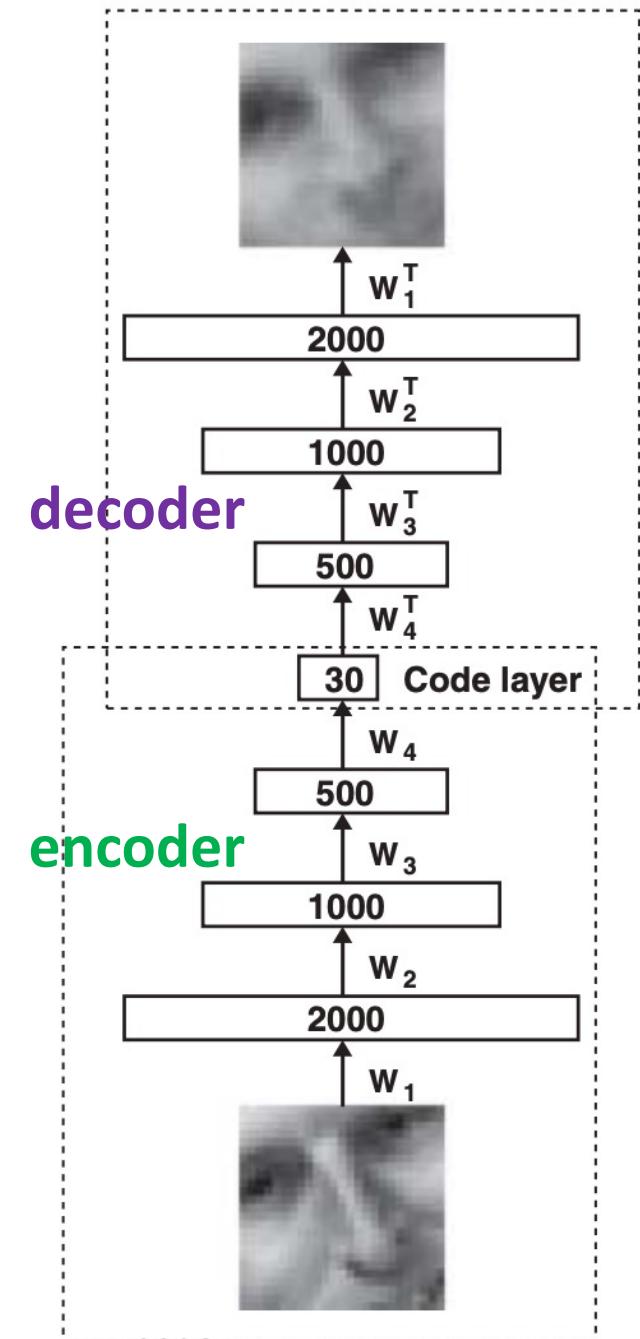


# Discussion

- K-means and PCA are like Autoencoding
- Autoencoding is “self-supervised learning”
  - “auto-” == “self-”
  - “encoding”  $\approx$  “learning”

Using data “itself”  
as the supervision

$$\min_{\mathcal{D}, \mathcal{E}} \sum_x \|\mathcal{D}(\mathcal{E}(x)) - x\|^2$$



# **Self-supervised Learning**

# Self-supervised Learning

- **Self-supervised learning:** Using the data itself as supervision
- **Self-supervised learning** is a form of “unsupervised learning”
- In general, the goal is to:
  - Learn representations from **large-scale** unlabeled data
  - Learn **high-level** semantic representations

# Self-supervised Learning

Modern self-supervised learning methods:

- Contrastive Learning
- Predictive Learning

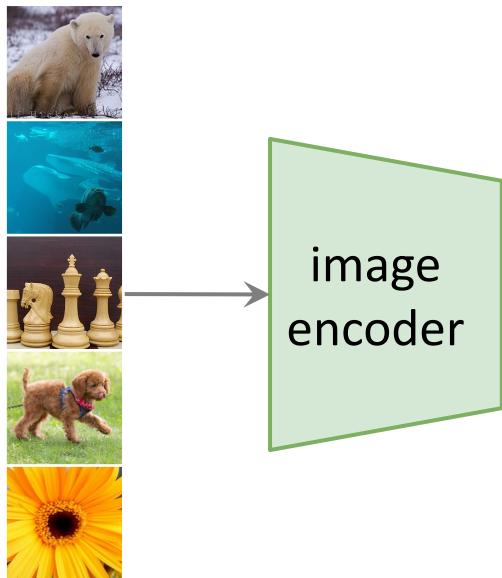
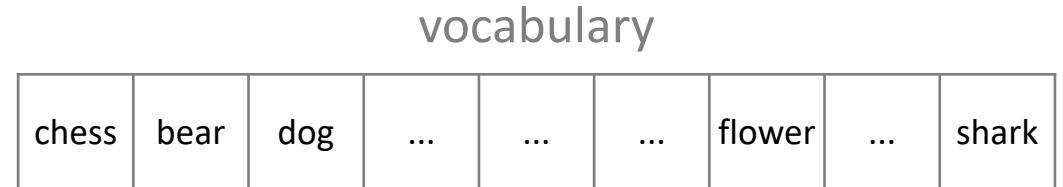
# Self-supervised Learning

Modern self-supervised learning methods:

- **Contrastive Learning**
- Predictive Learning

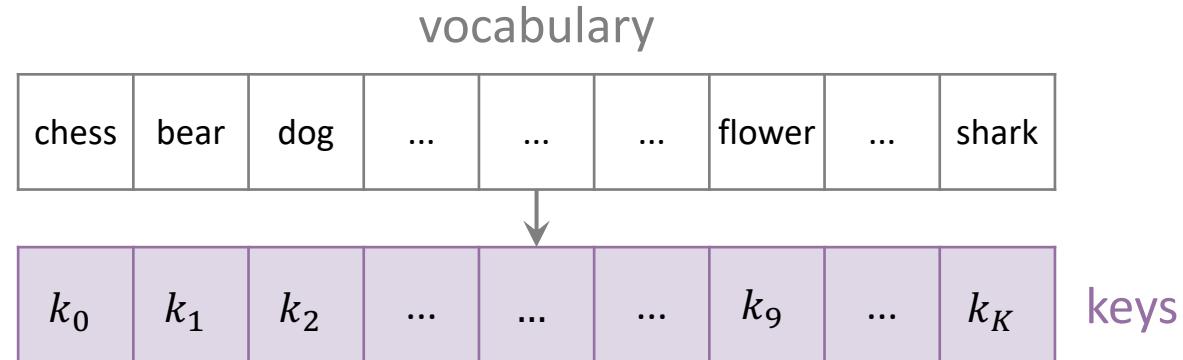
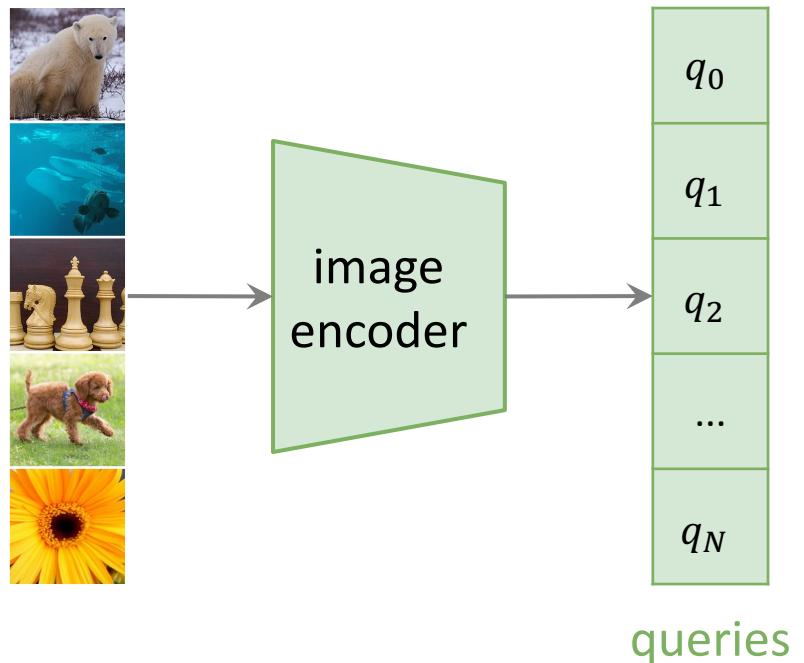
# Contrastive Learning

- **Supervised classification** is Contrastive Learning



# Contrastive Learning

- Supervised classification is Contrastive Learning

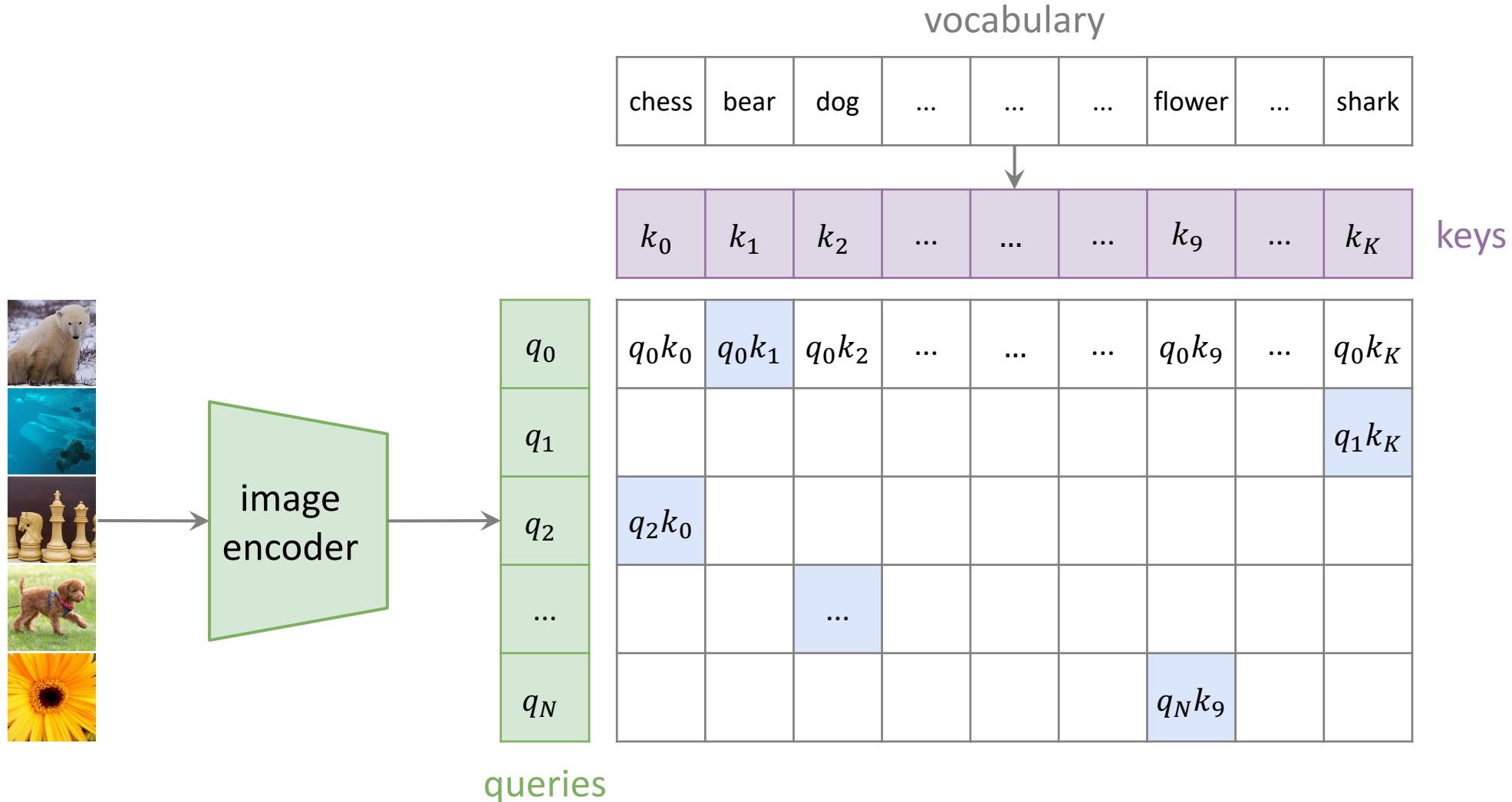


These are the weights in  
the last fc layer

These are the  
representations before  
the last fc layer

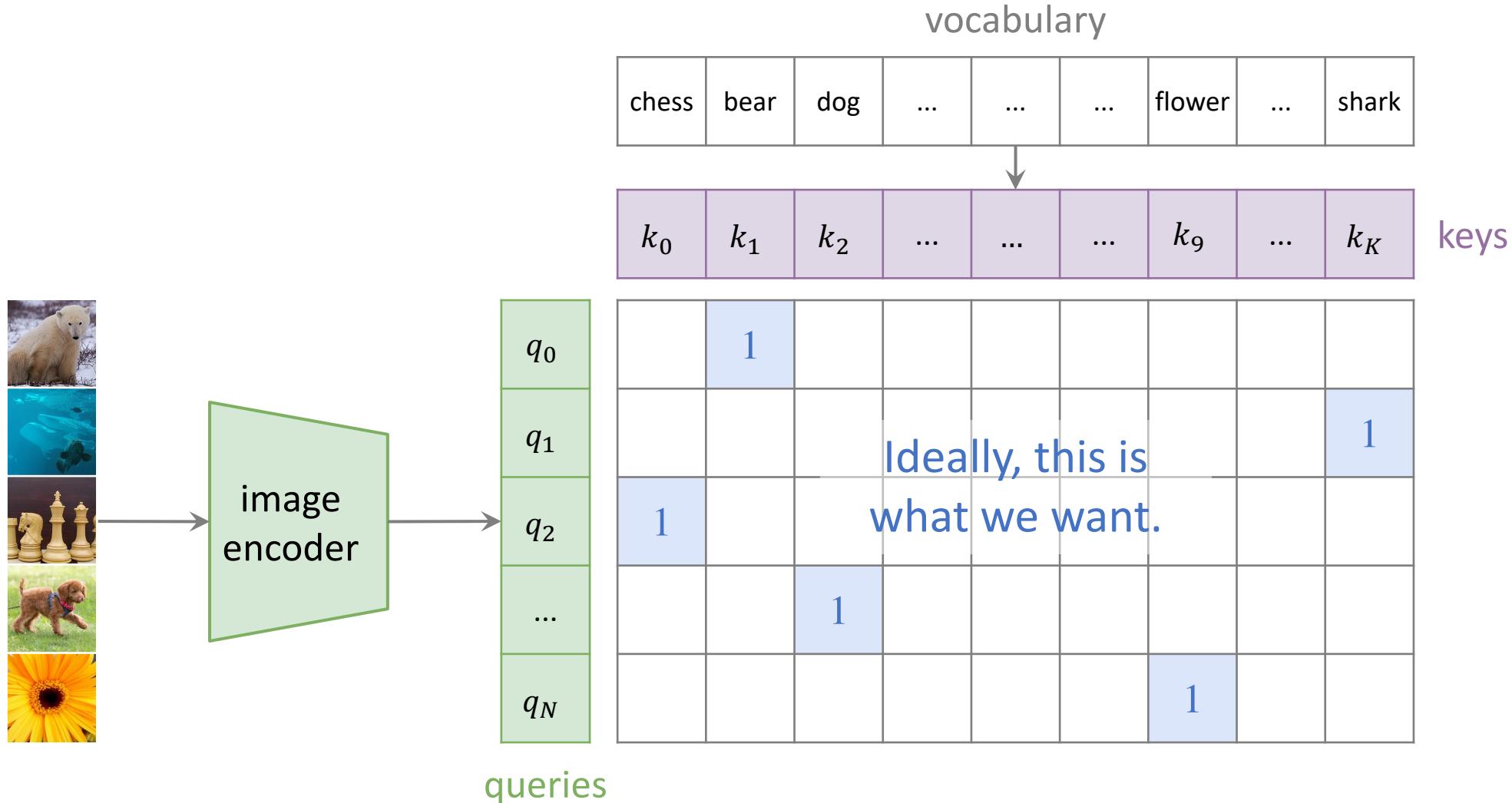
# Contrastive Learning

- Supervised classification is Contrastive Learning



# Contrastive Learning

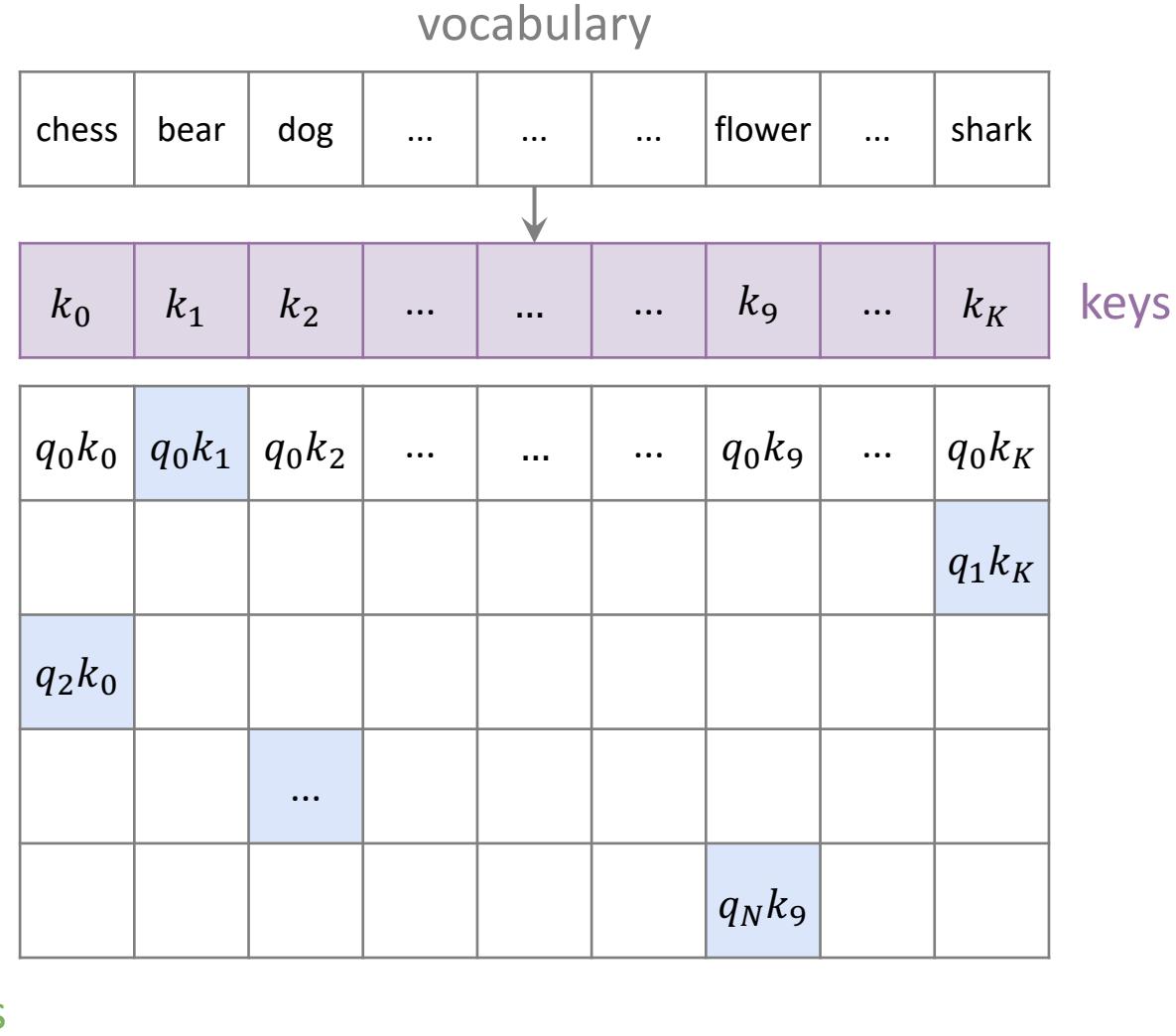
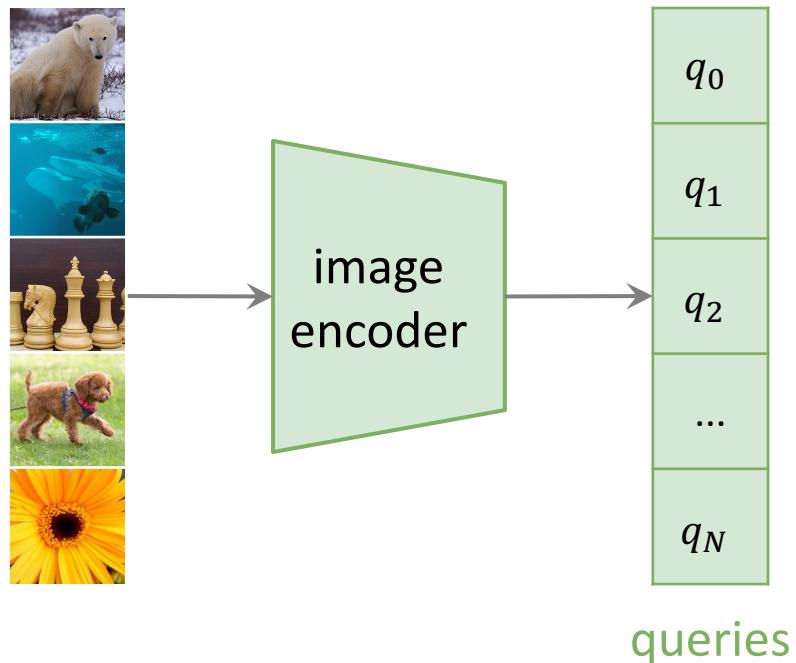
- Supervised classification is Contrastive Learning



# Contrastive Learning

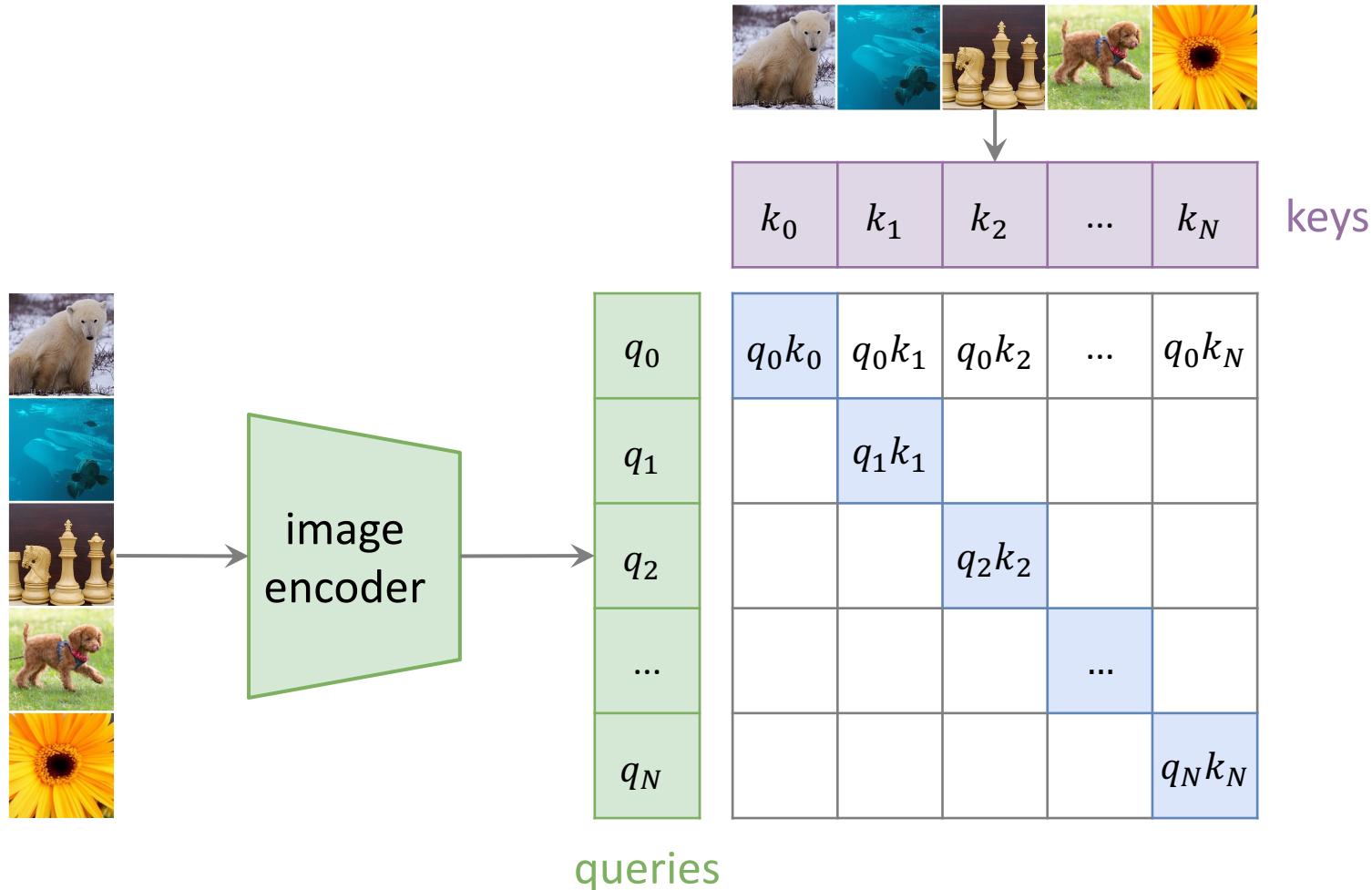
- Supervised classification is Contrastive Learning

- maximize similarities of positive pairs
- minimize similarities of negative pairs



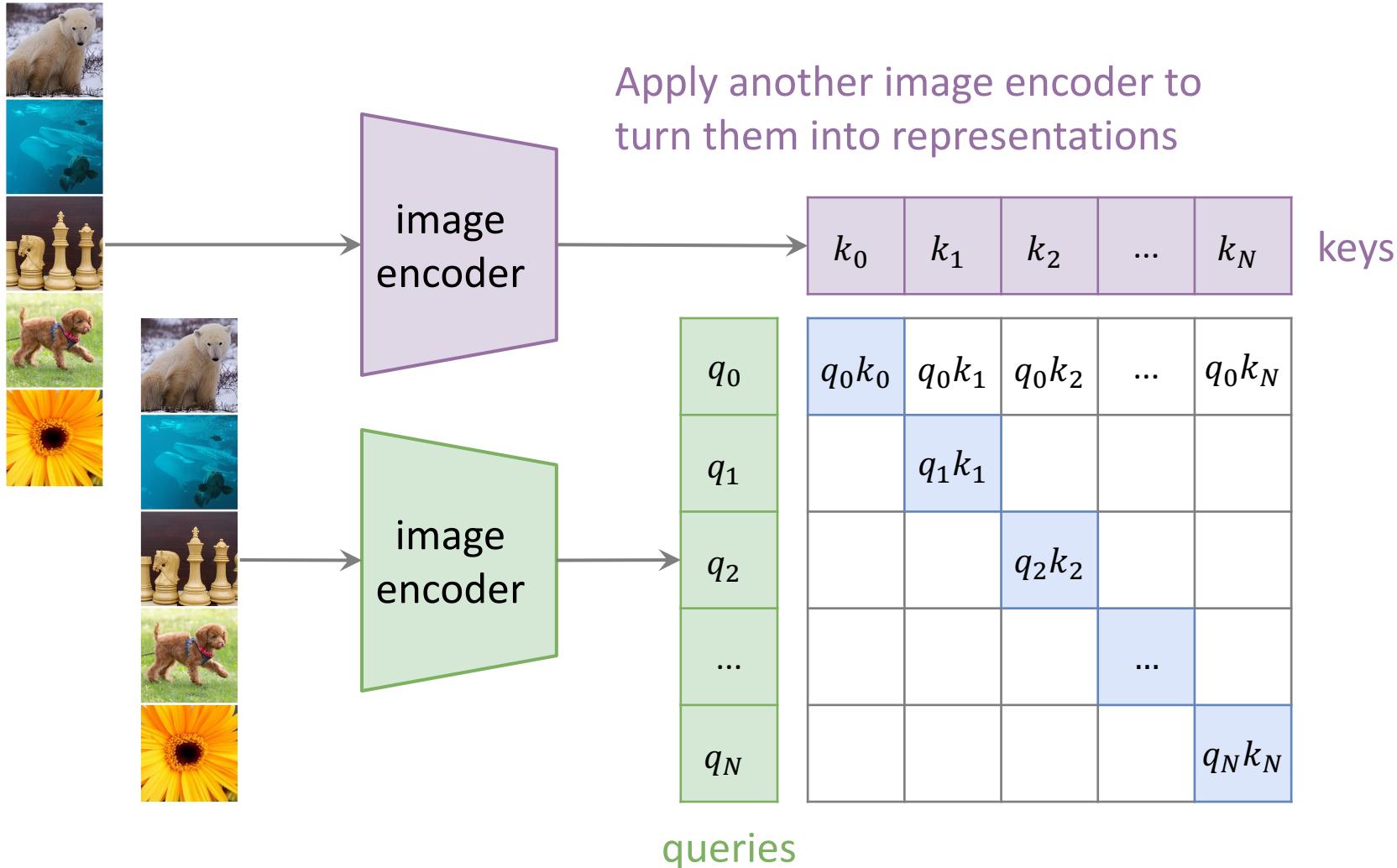
# Contrastive Learning

- What if we don't have the vocabulary? Form a vocabulary w/ data **itself**



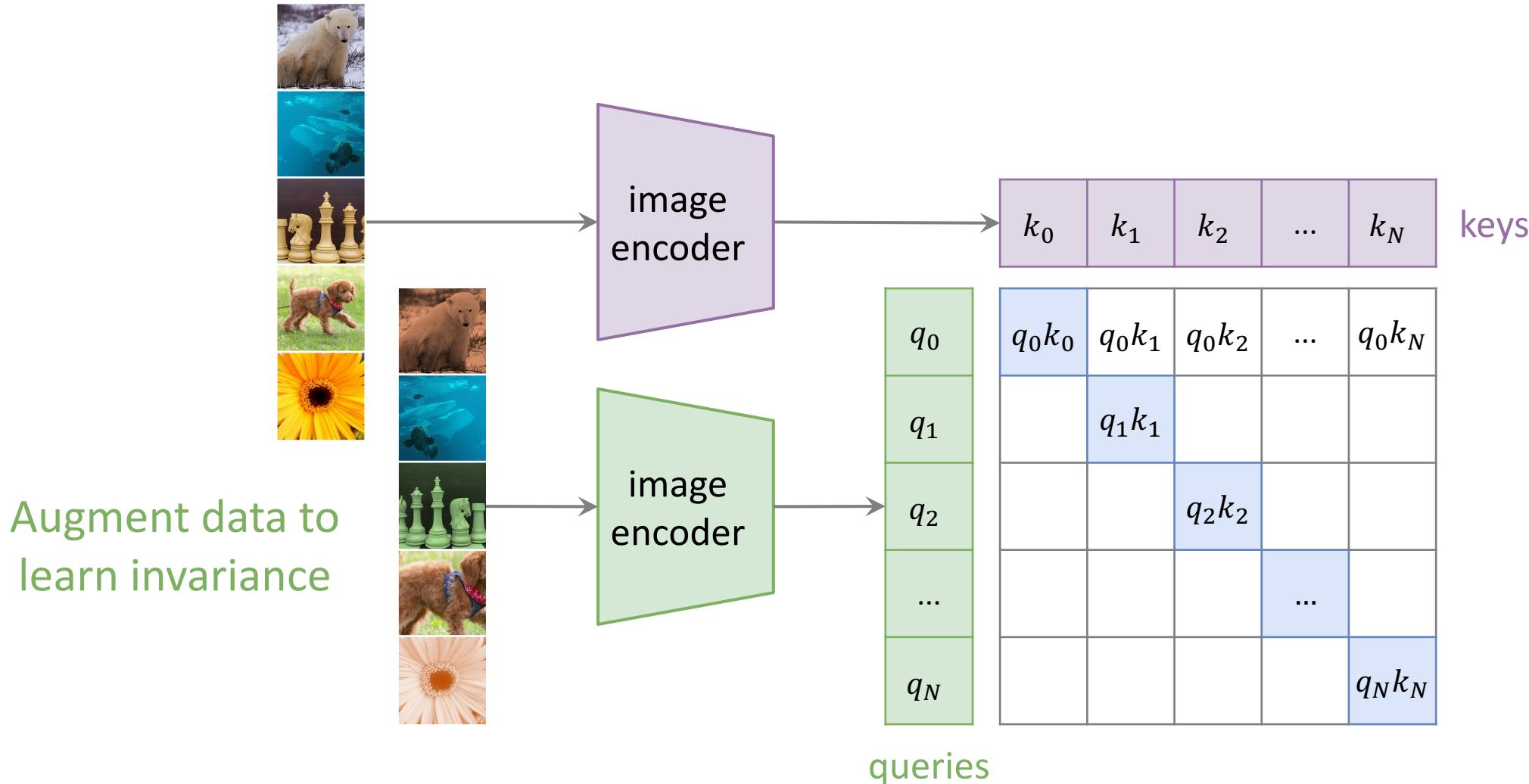
# Contrastive Learning

- What if we don't have the vocabulary? Form a vocabulary w/ data **itself**



# Contrastive Learning

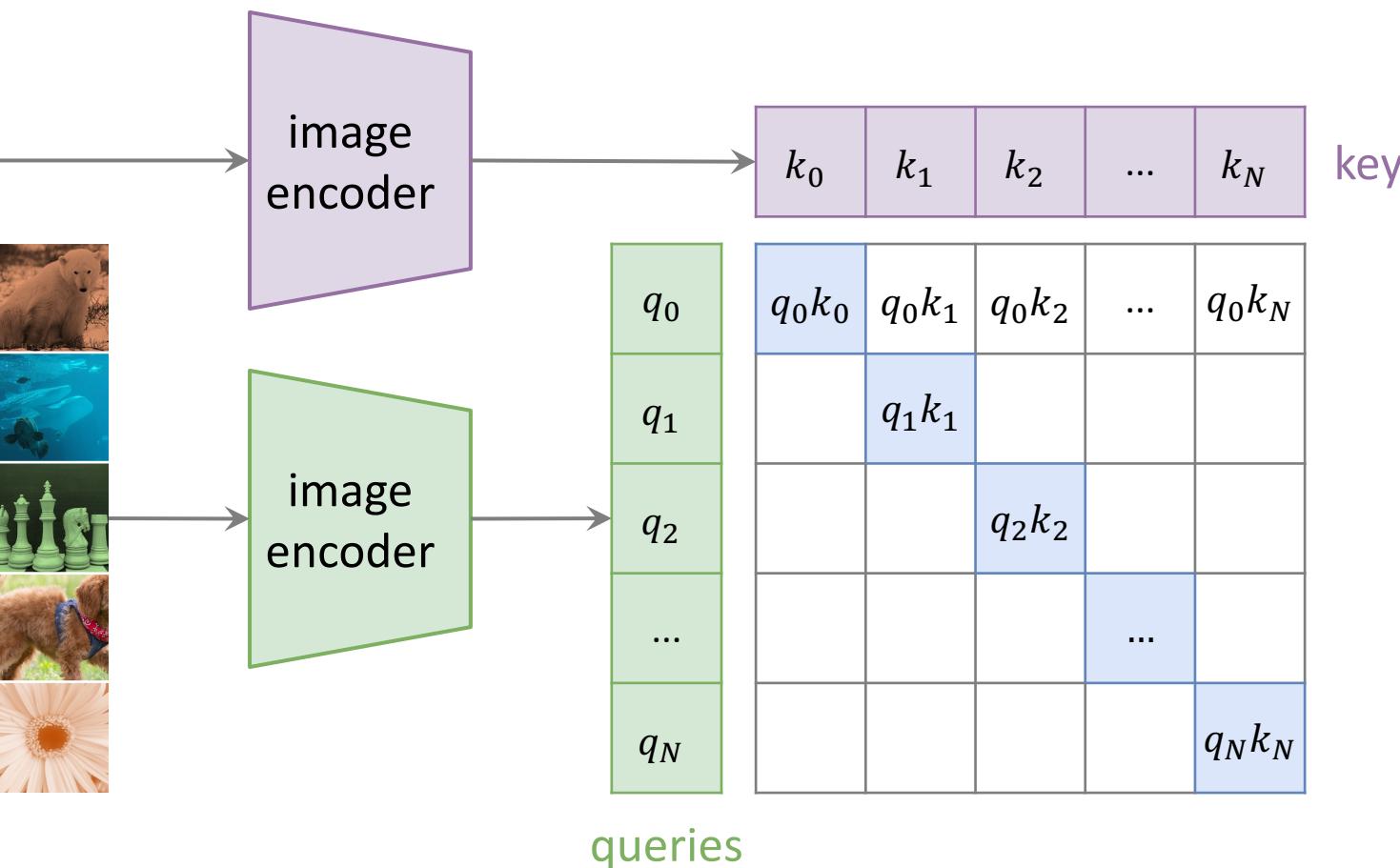
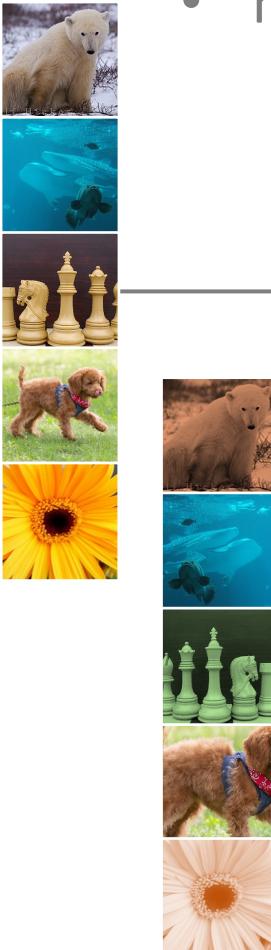
- What if we don't have the vocabulary? Form a vocabulary w/ data **itself**



# Contrastive Learning

contrastive loss

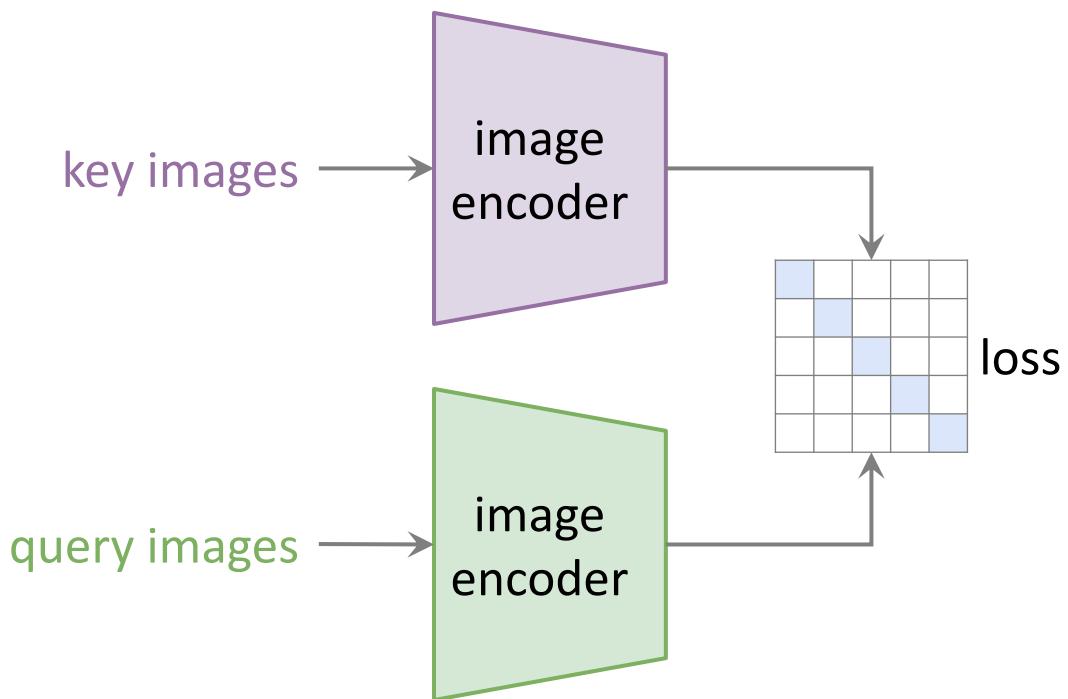
- maximize similarities of positive pairs
- minimize similarities of negative pairs



$$L = -\log \frac{\exp(\mathbf{q}\mathbf{k}_+/\tau)}{\sum_i \exp(\mathbf{q}\mathbf{k}_i/\tau)}$$

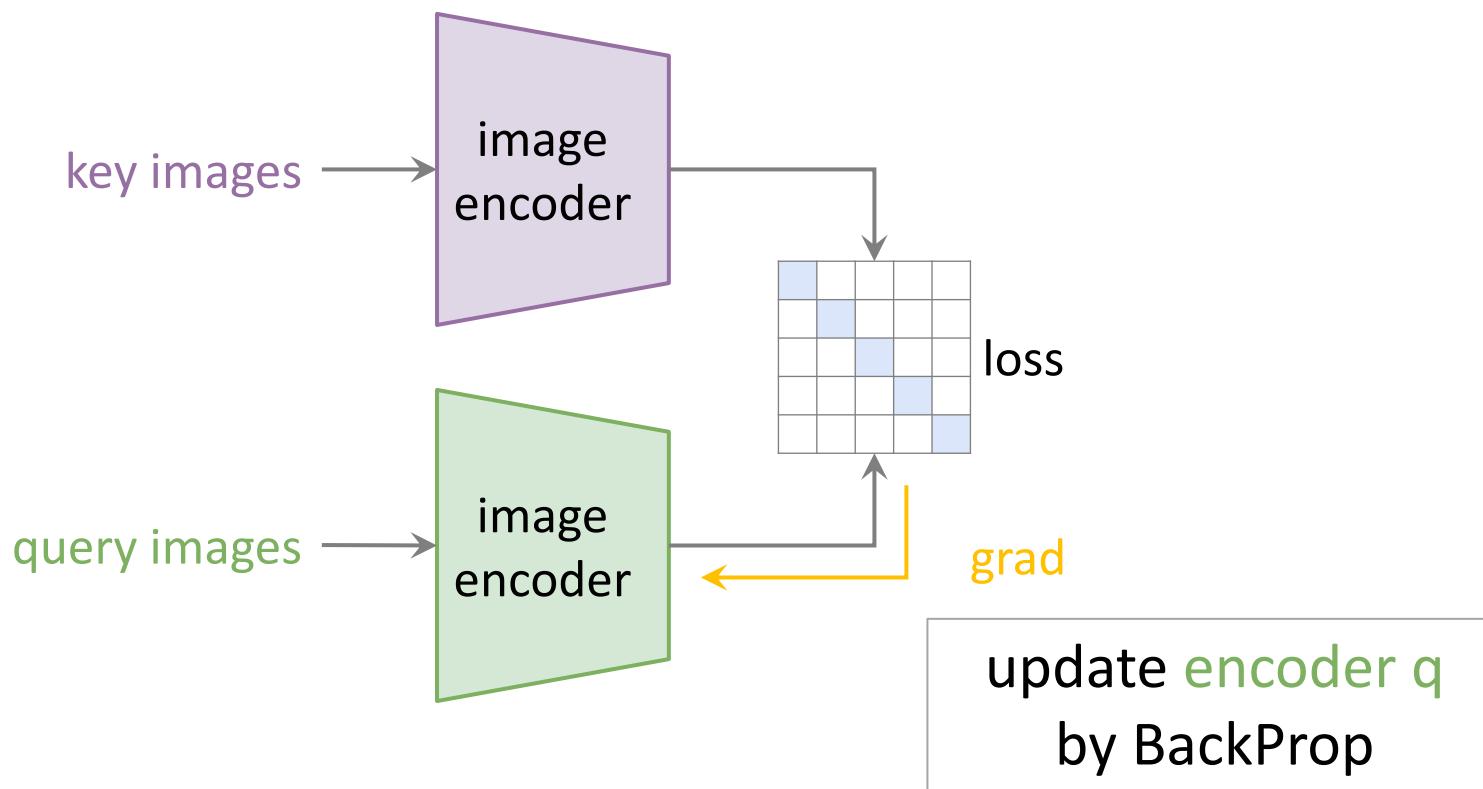
# Contrastive Learning

- maximize similarities of positive pairs
- minimize similarities of negative pairs



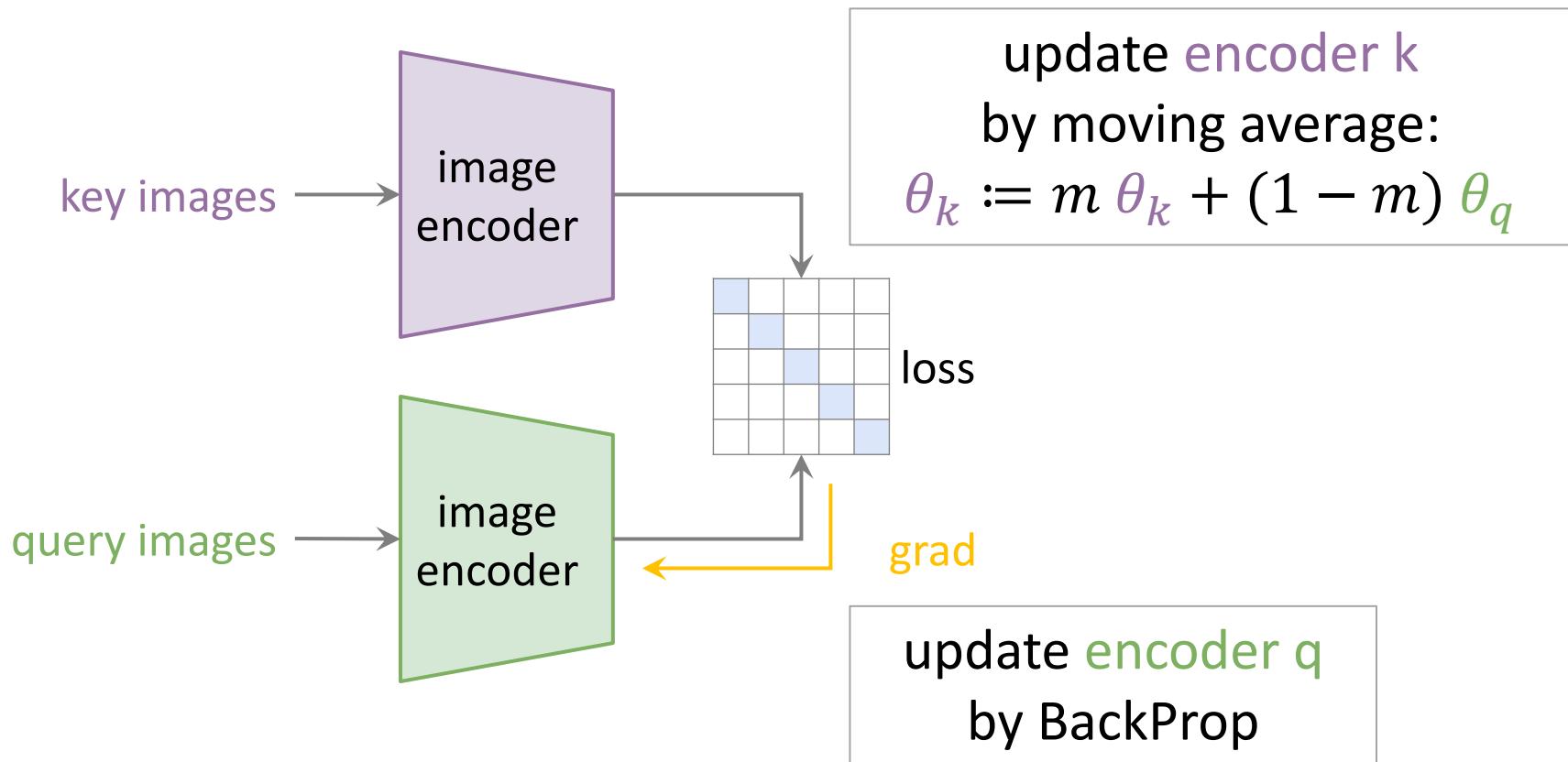
# Momentum Contrastive Learning (MoCo)

- maximize similarities of positive pairs
- minimize similarities of negative pairs



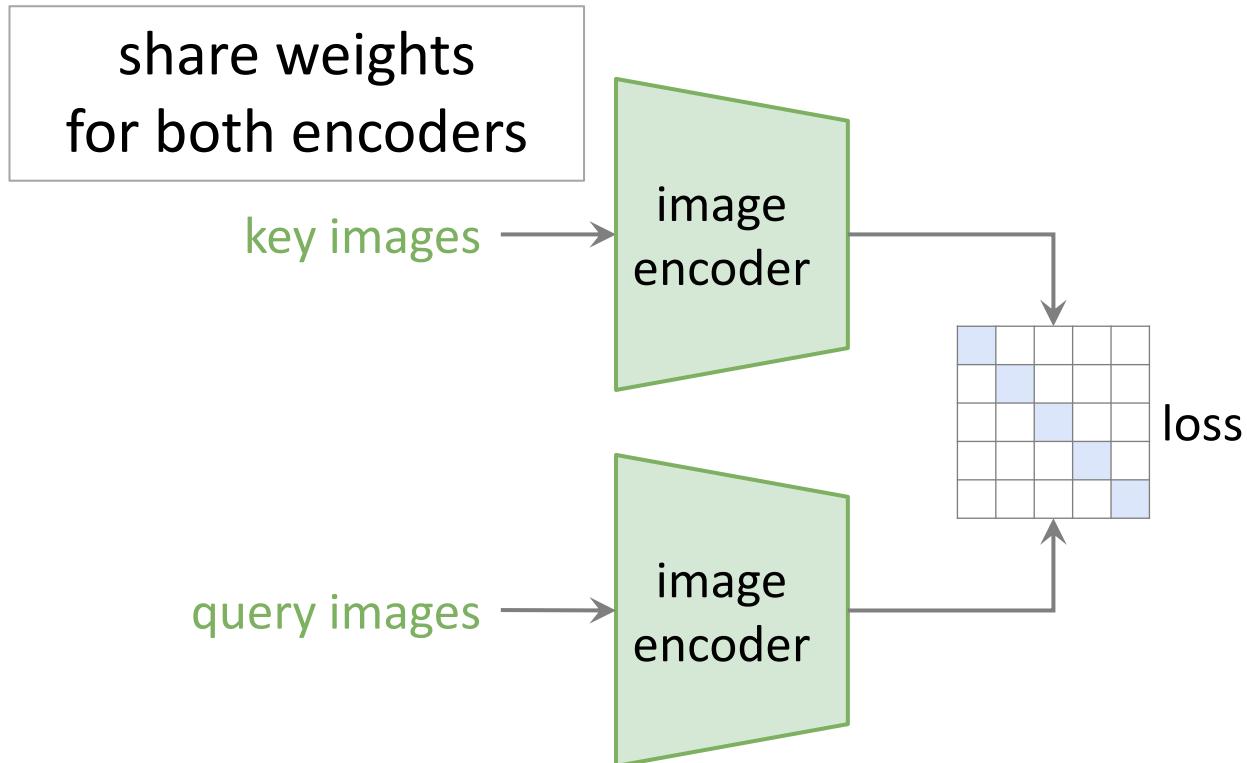
# Momentum Contrastive Learning (MoCo)

- maximize similarities of positive pairs
- minimize similarities of negative pairs



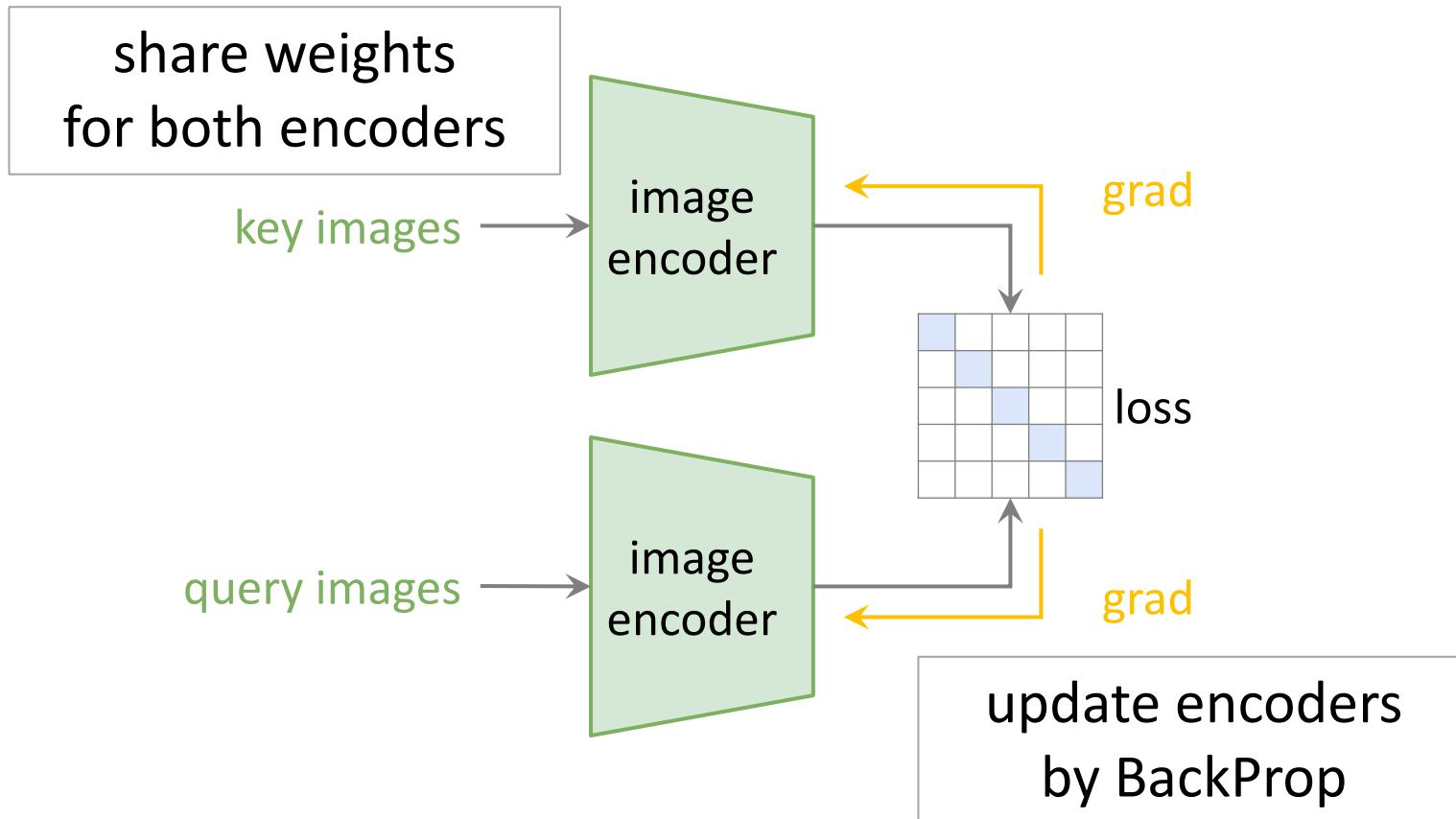
# Simple Contrastive Learning (SimCLR)

- maximize similarities of positive pairs
- minimize similarities of negative pairs



# Simple Contrastive Learning (SimCLR)

- maximize similarities of positive pairs
- minimize similarities of negative pairs

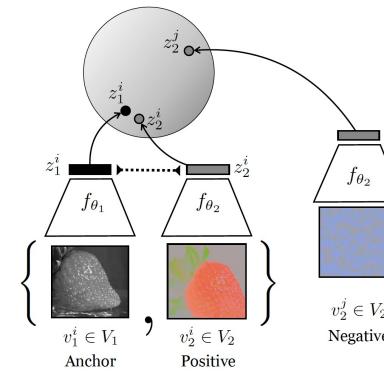


# Many More Variants

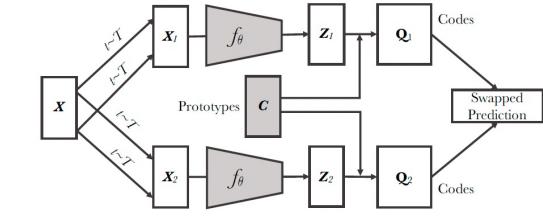
## Compare samples in the latent space

- Different data augmentation
- Different loss functions
- w/ or w/o negative samples
- w/ or w/o momentum encoder

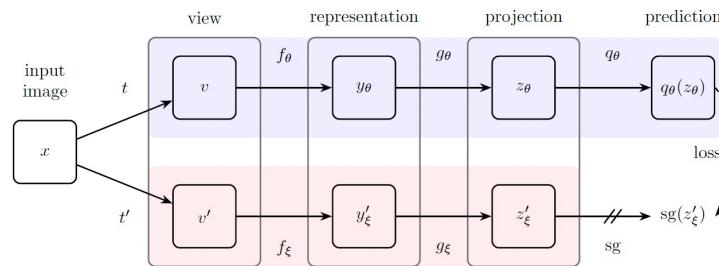
## A frontier research direction



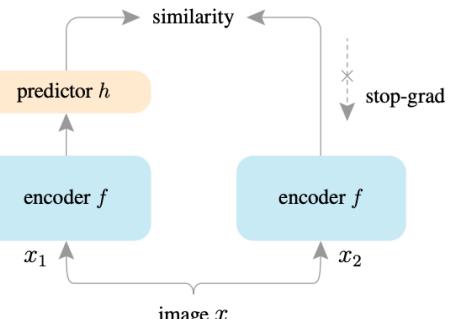
CMC



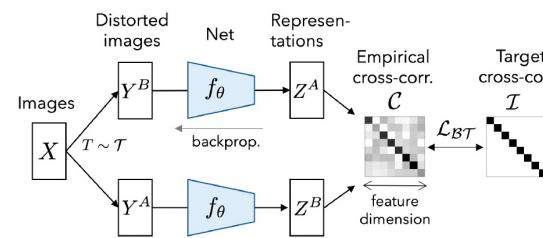
SwAV



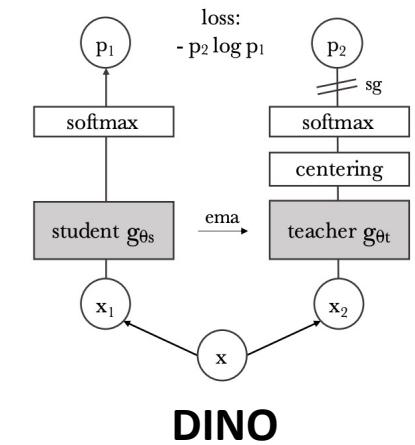
BYOL



SimSiam



Barlow Twins



DINO

# Self-supervised Learning

Modern self-supervised learning methods:

- Contrastive Learning
- **Predictive Learning**

# Predictive Learning: Language Models

## Next word prediction (GPT)

- Predict the next word (token) given a prefix

"The students opened their \_\_\_\_\_

model

books

# Predictive Learning: Language Models

## Masked language modeling (BERT)

- Predict the masked words (tokens) in a text



# Predictive Learning: Computer Vision

## Masked image modeling (Context Encoders)

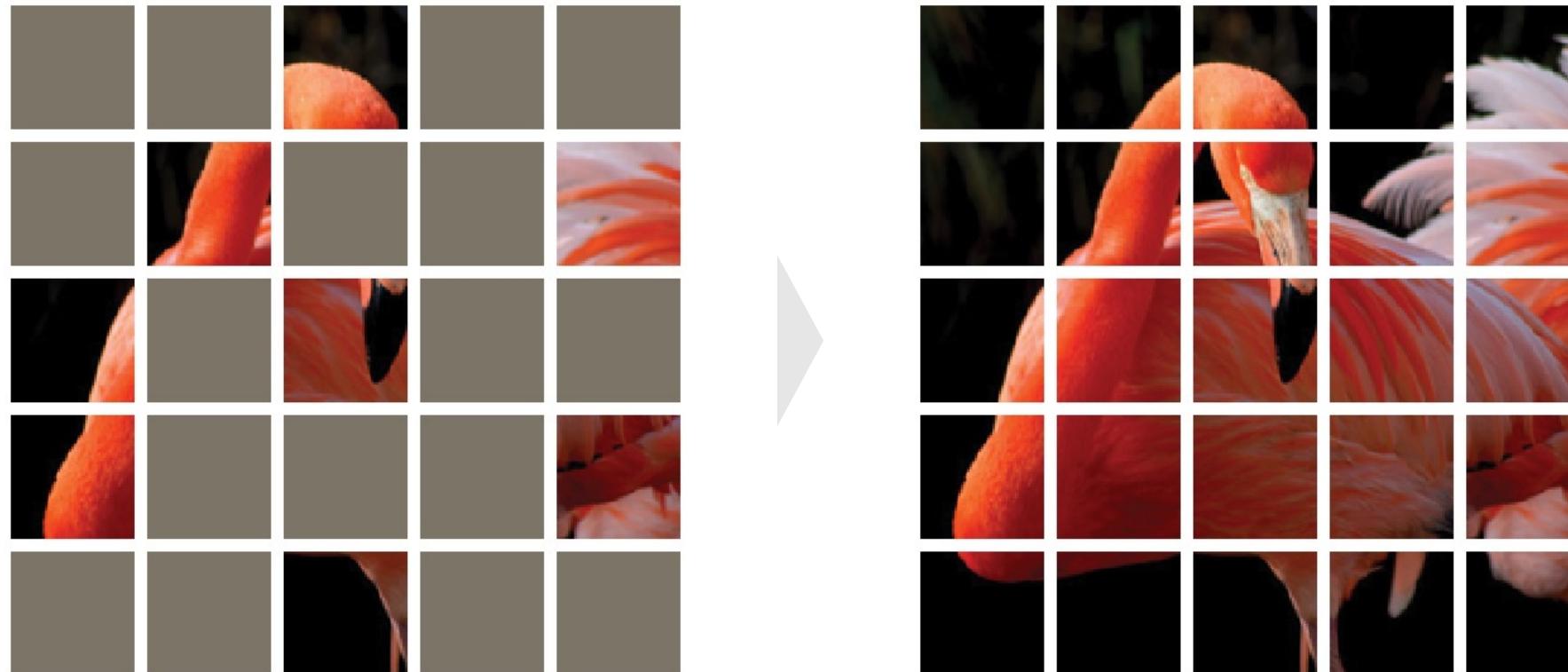
- Predict the masked regions using ConvNets



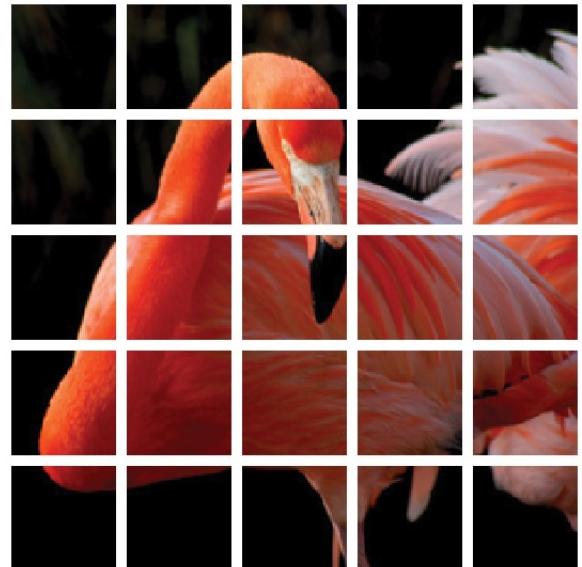
# Predictive Learning: Computer Vision

## Masked image modeling (Masked Autoencoder)

- Predict the masked patches using Transformers

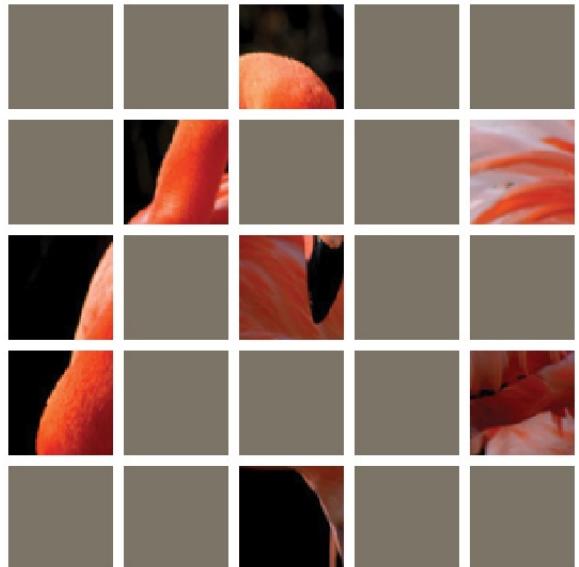


# Masked Autoencoder (MAE)



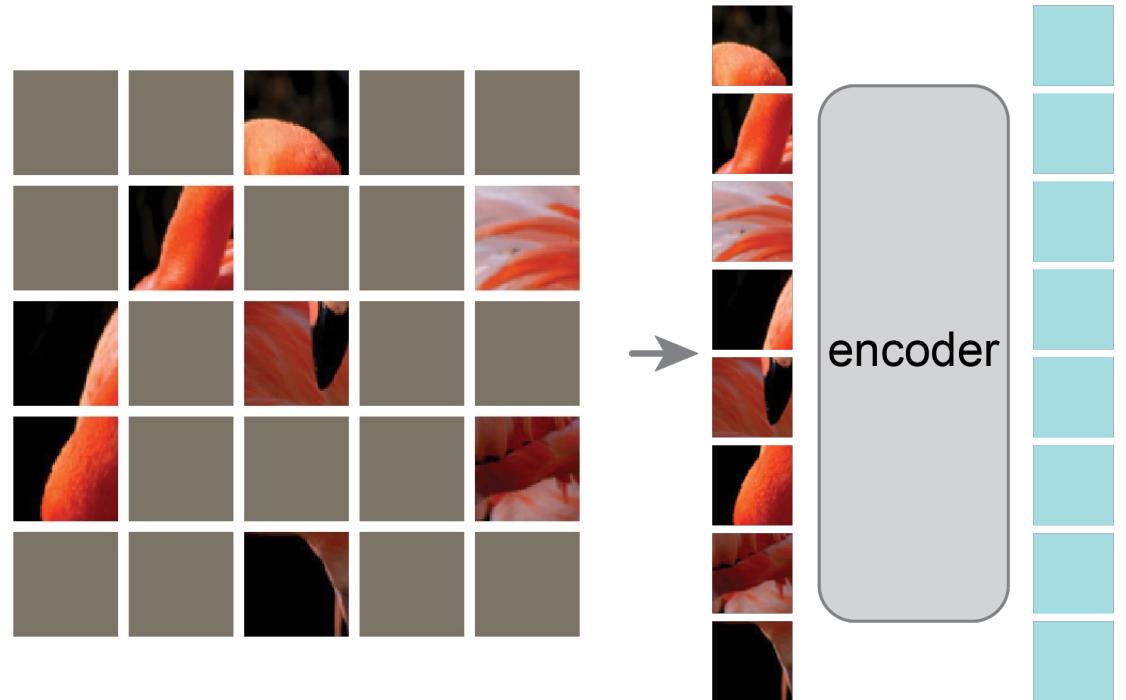
patches as visual tokens  
(Vision Transformer)

# Masked Autoencoder (MAE)



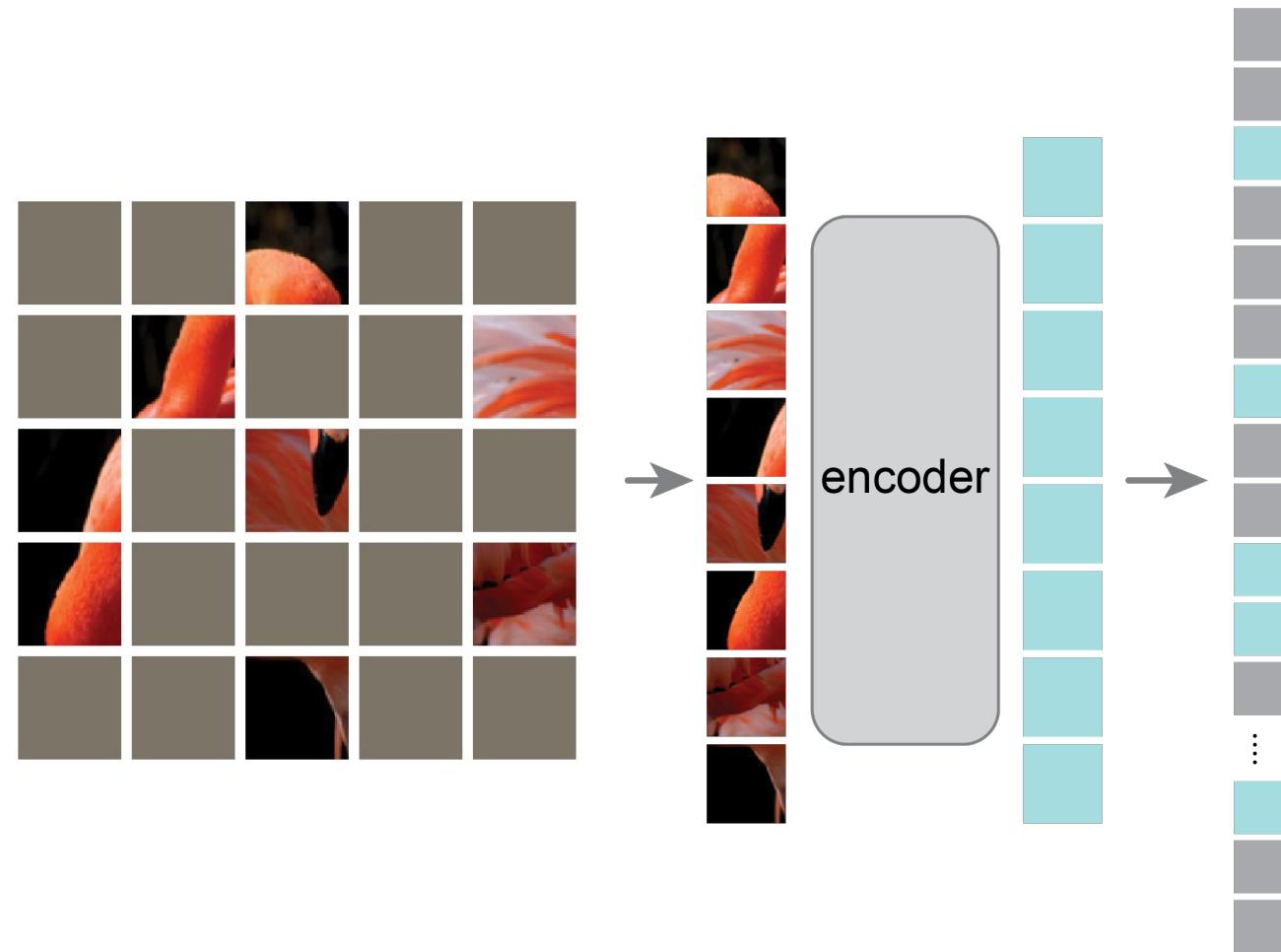
random masking

# Masked Autoencoder (MAE)

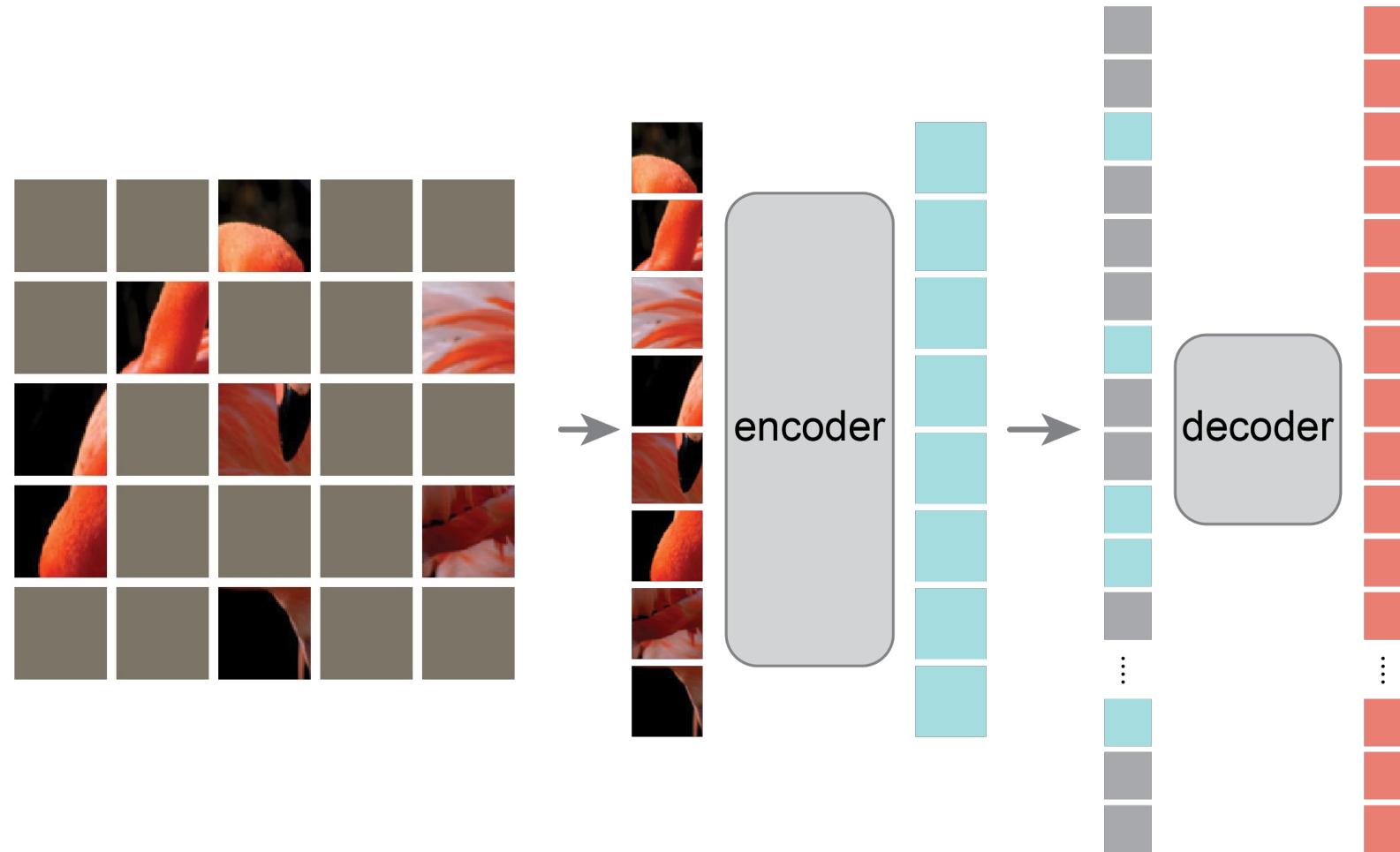


encode visible patches  
w/ Transformer

# Masked Autoencoder (MAE)

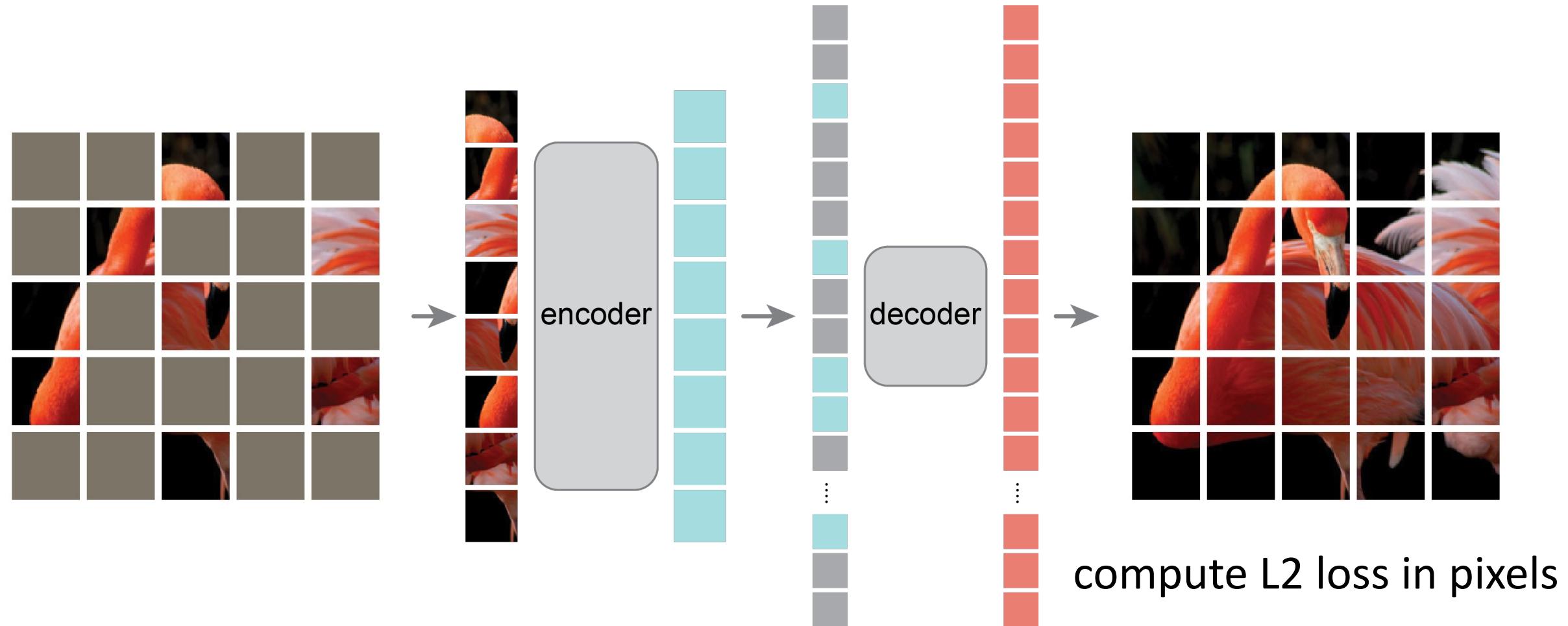


# Masked Autoencoder (MAE)

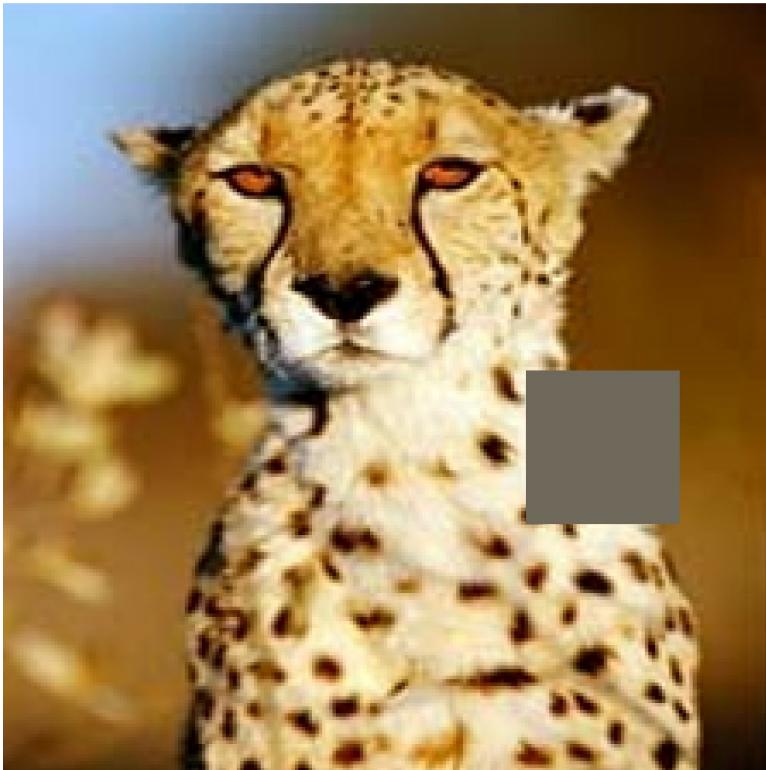


predict the unknown

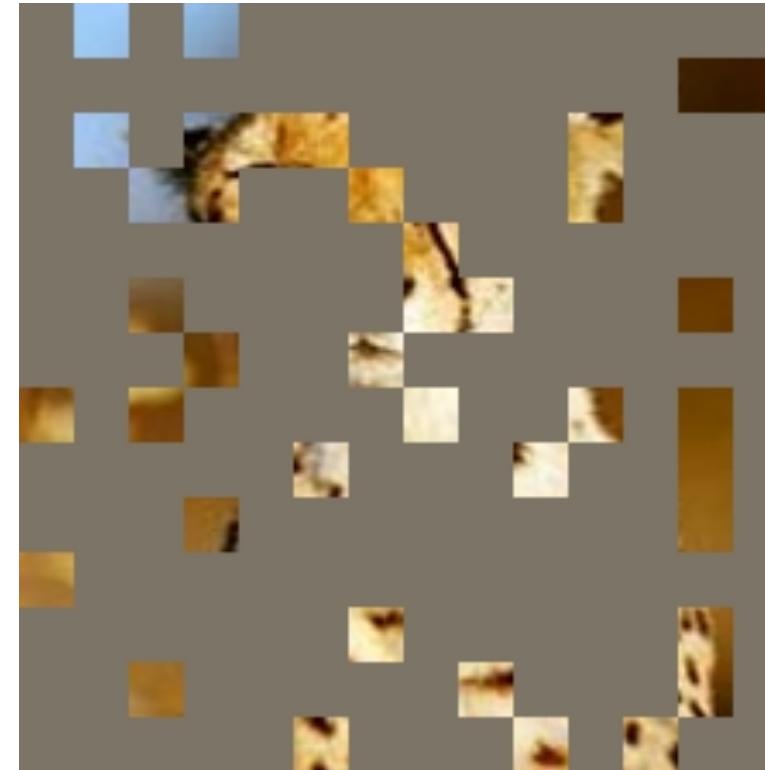
# Masked Autoencoder (MAE)



# How to learn good representations by predicting?

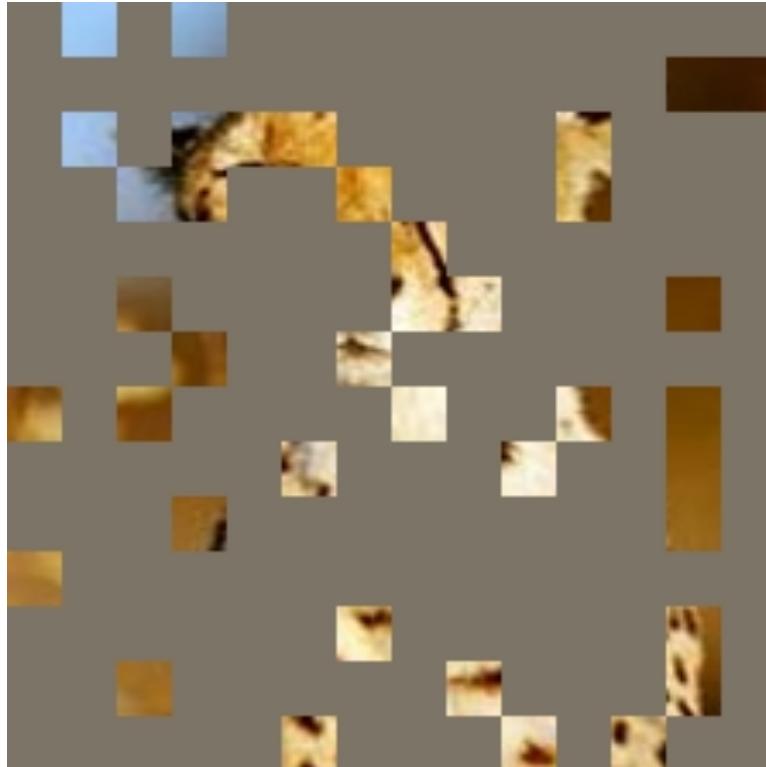


- predicting a small portion may not require high-level understanding

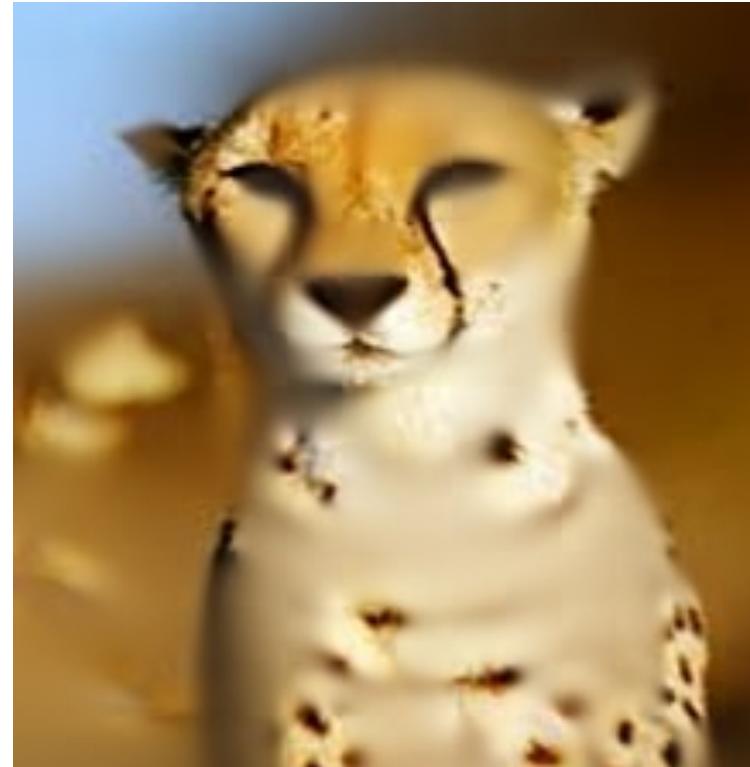


- predicting a large portion of unknown patches encourages to learn semantic features

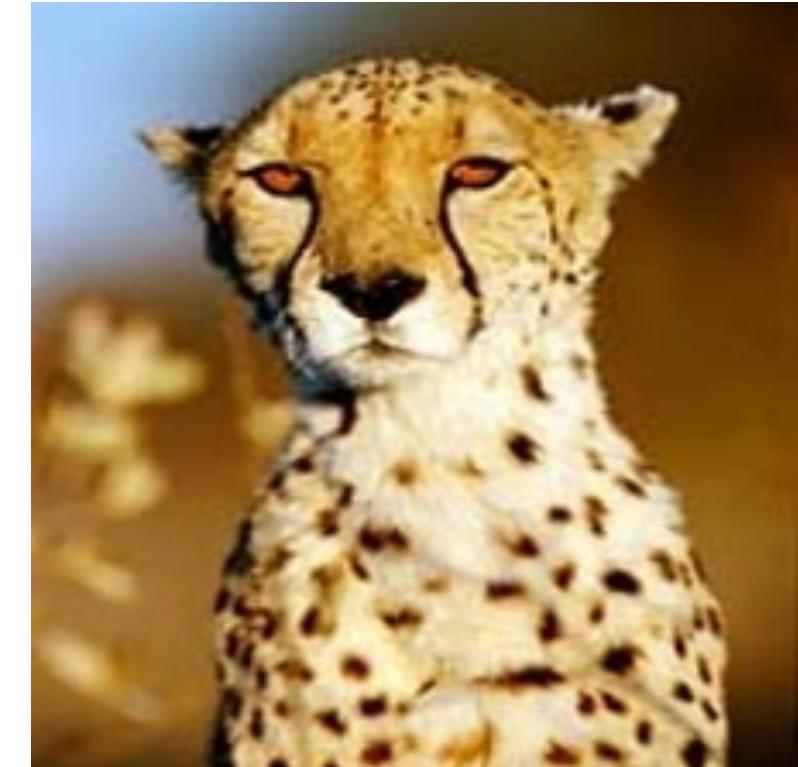
# How to learn good representations by predicting?



input



MAE prediction



original

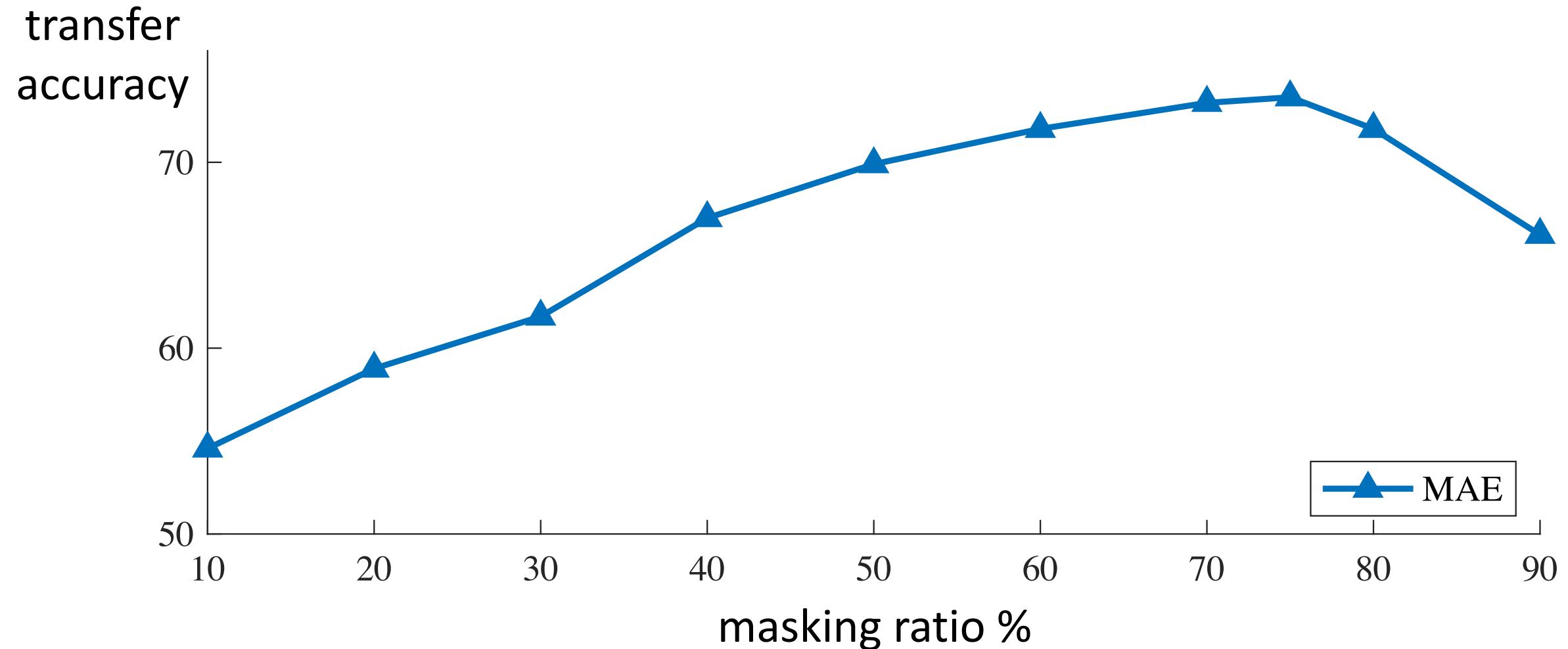


- **The learning process:** the network gradually makes sense of the semantic patterns

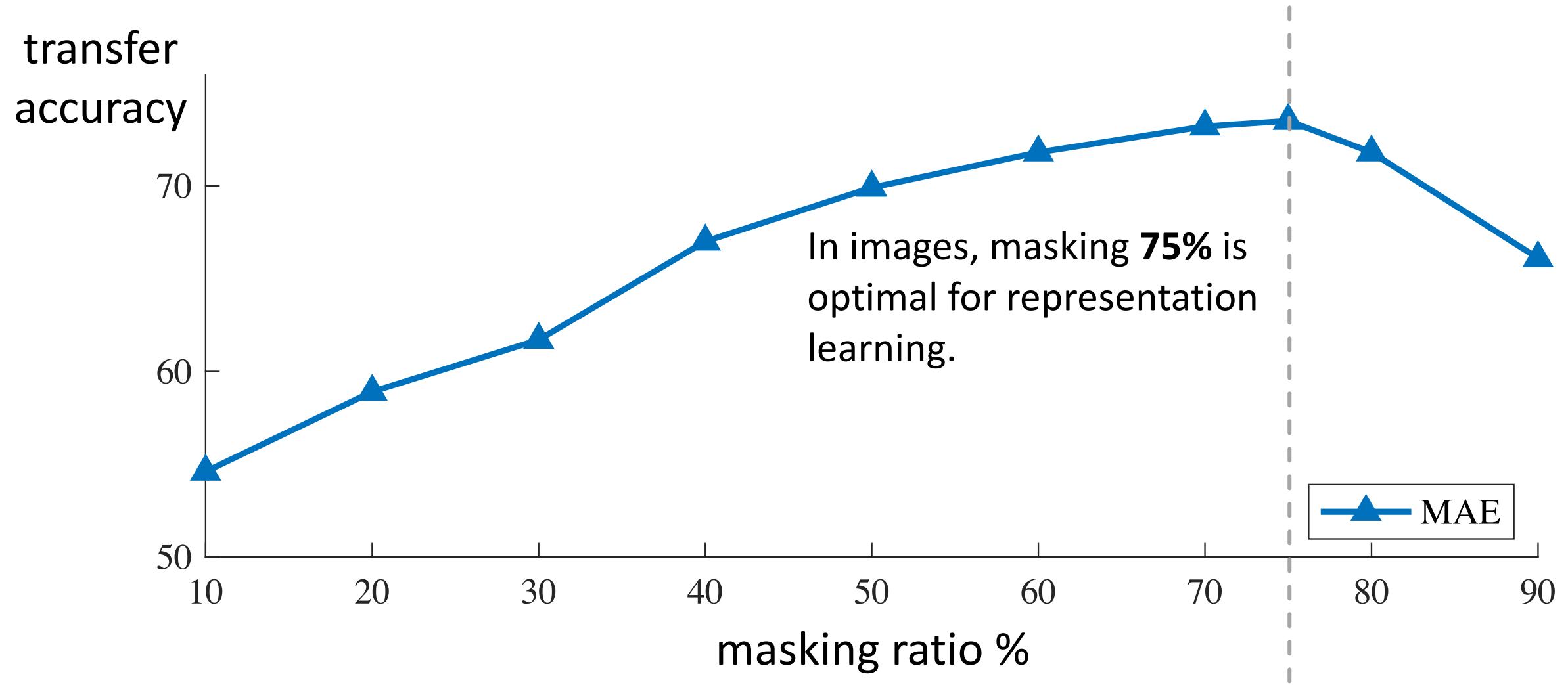


- **The learning process:** the network gradually makes sense of the semantic patterns

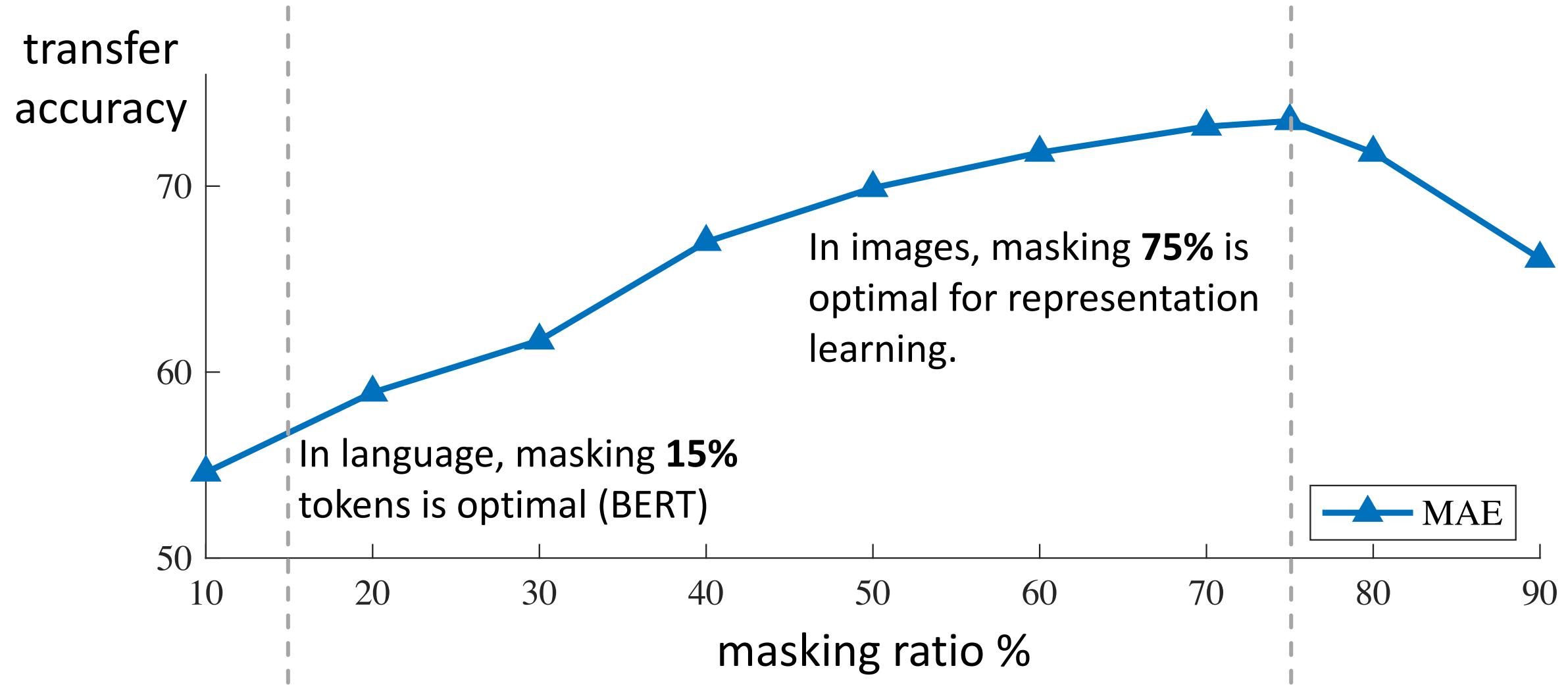
# How to learn good representations by predicting?



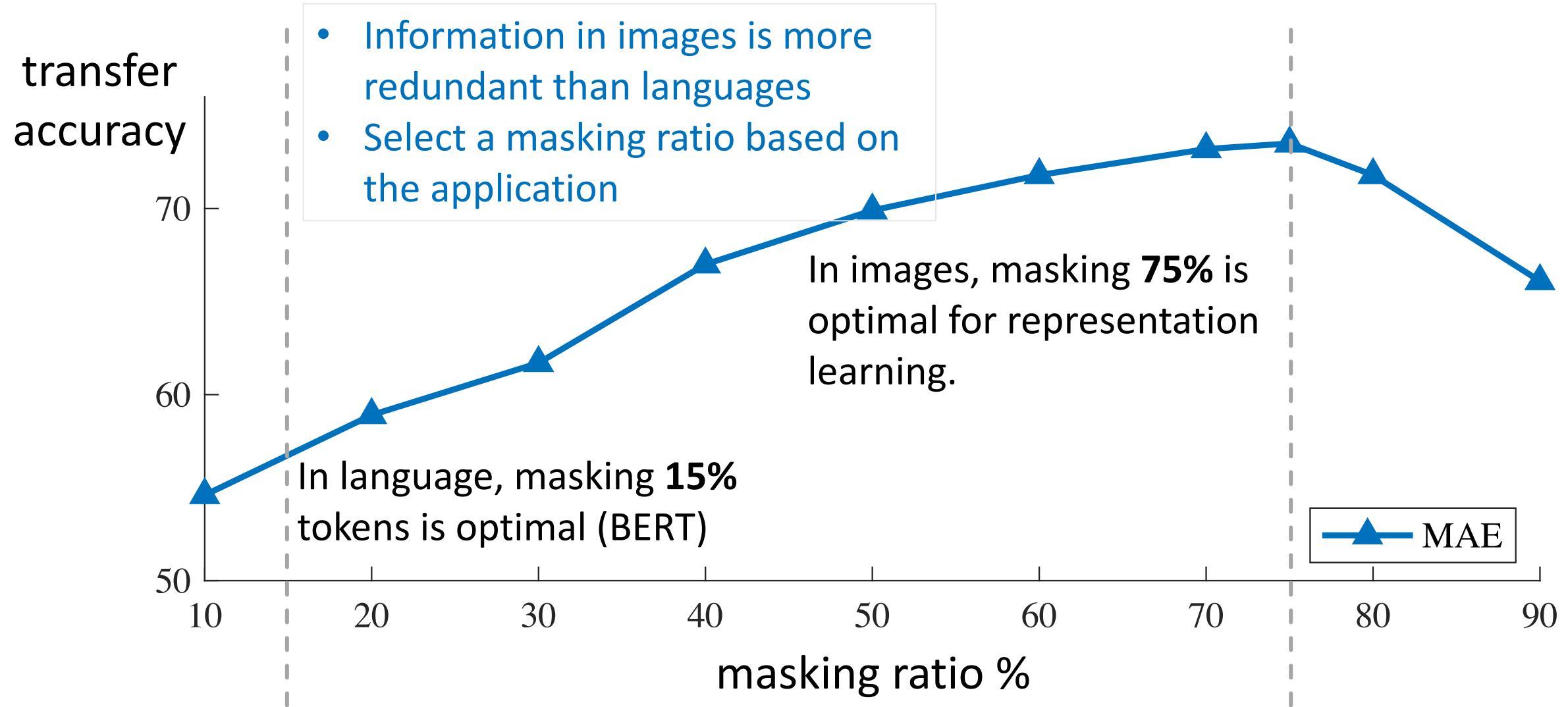
# How to learn good representations by predicting?



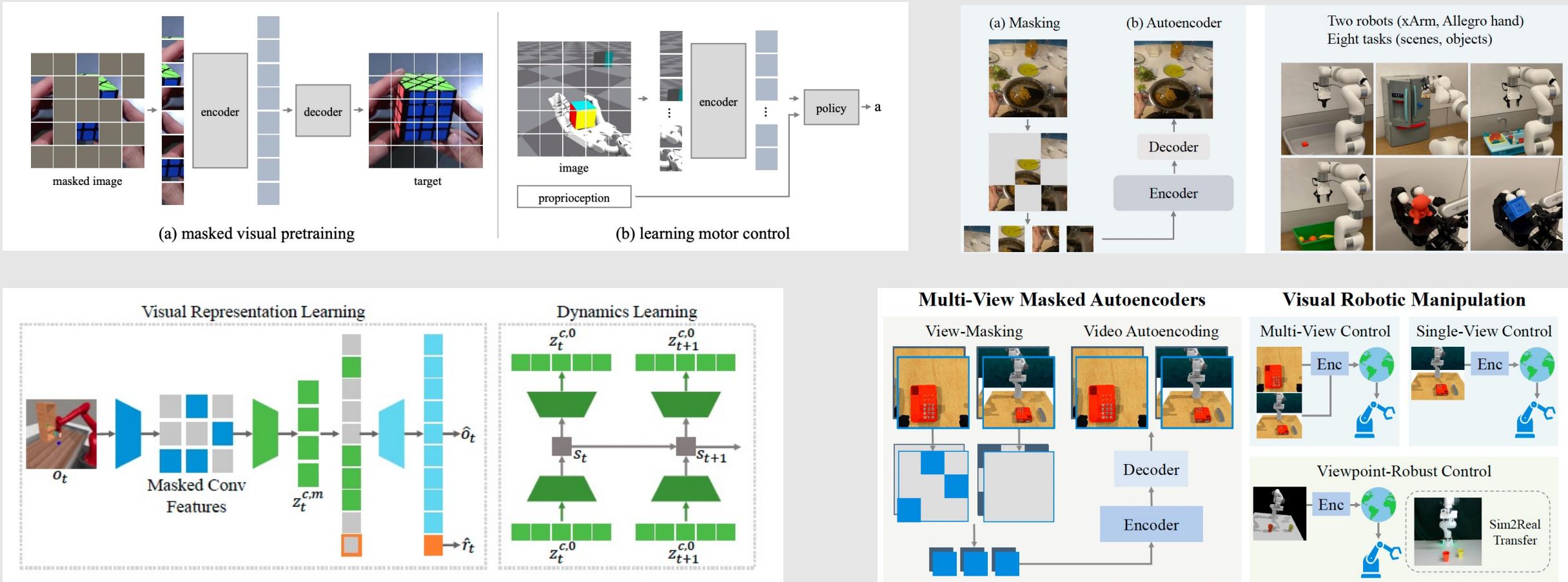
# How to learn good representations by predicting?



# How to learn good representations by predicting?

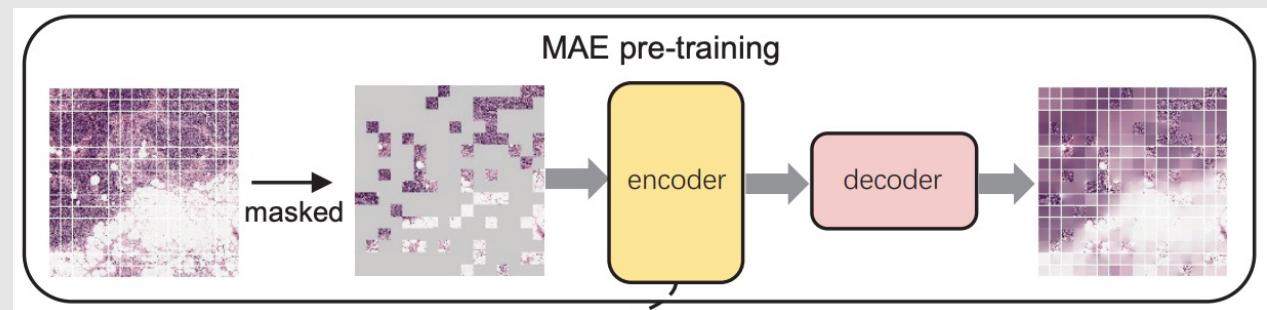
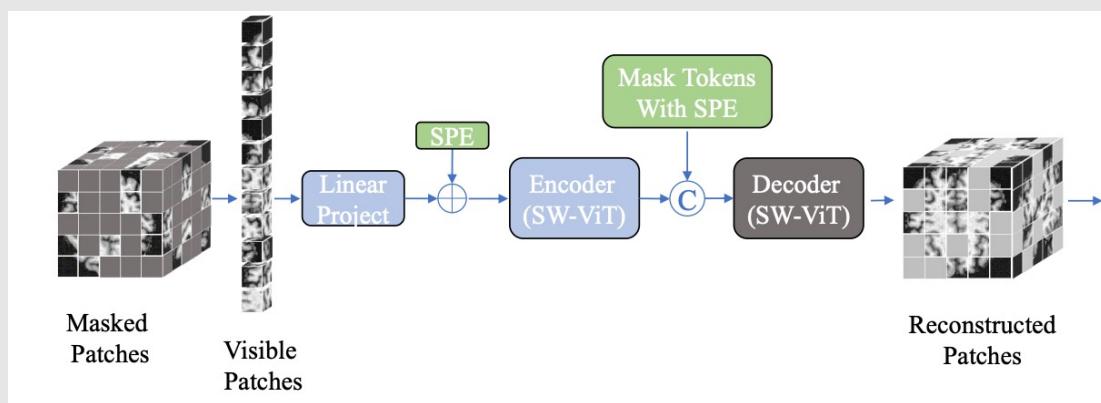
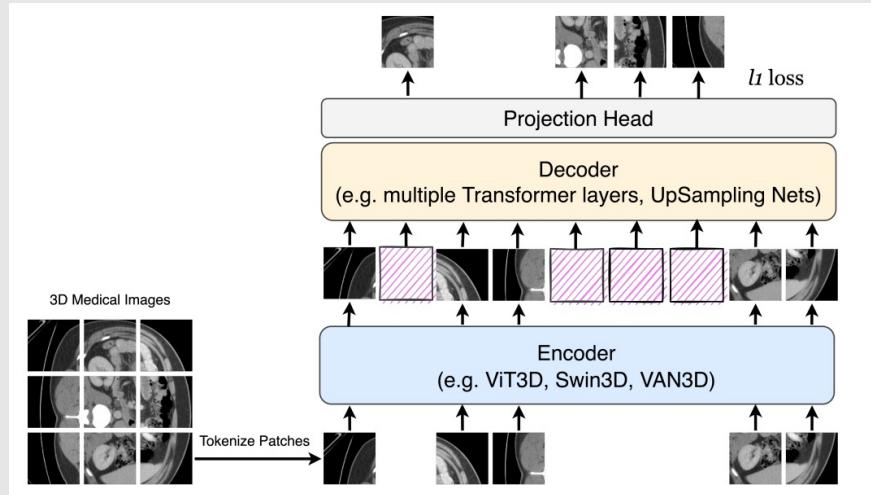
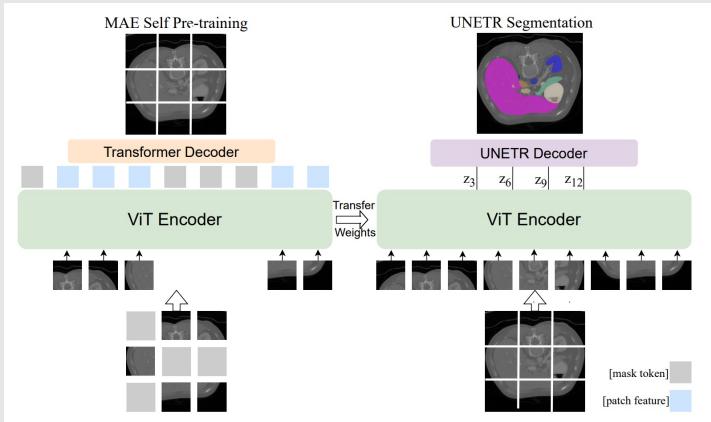


# Masked Autoencoder: Extensions



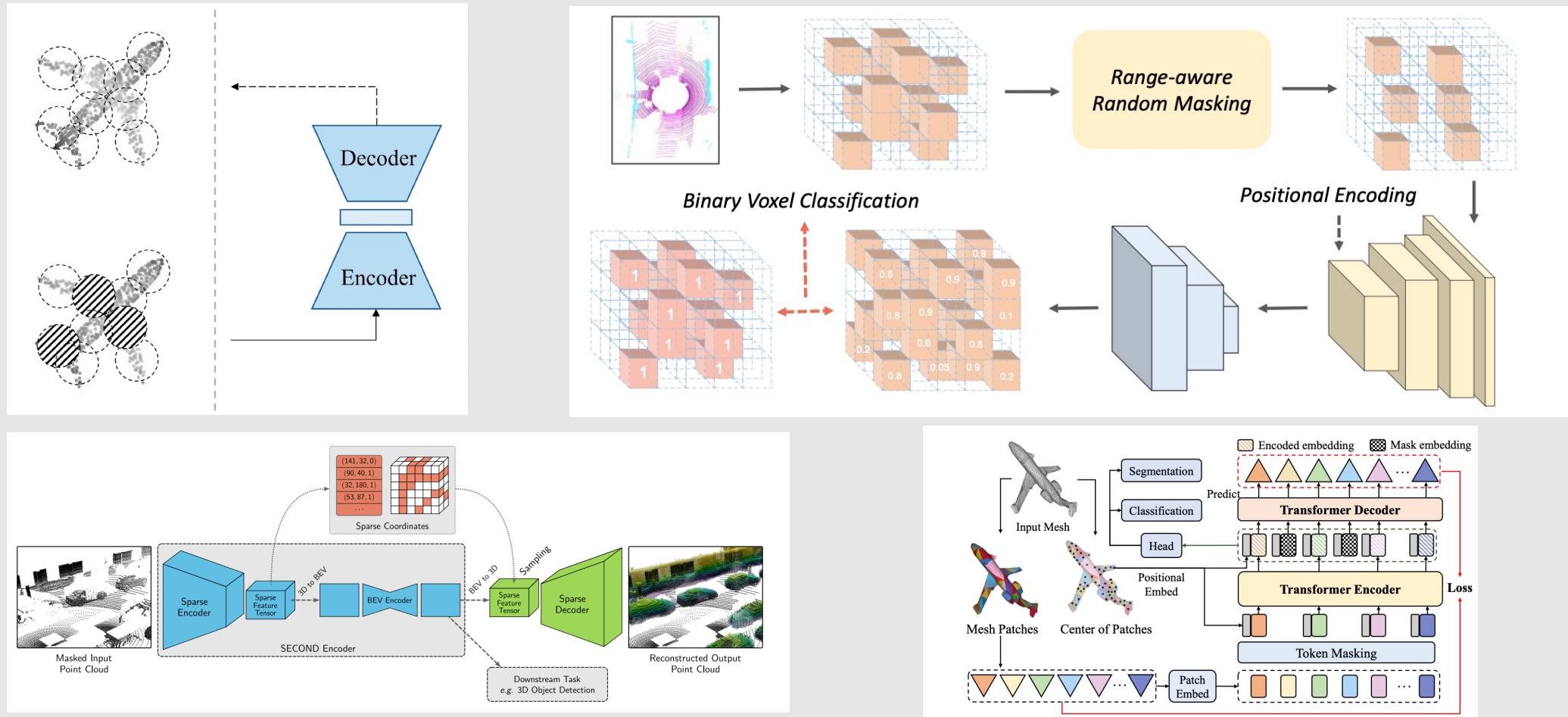
- Robotics

# Masked Autoencoder: Extensions



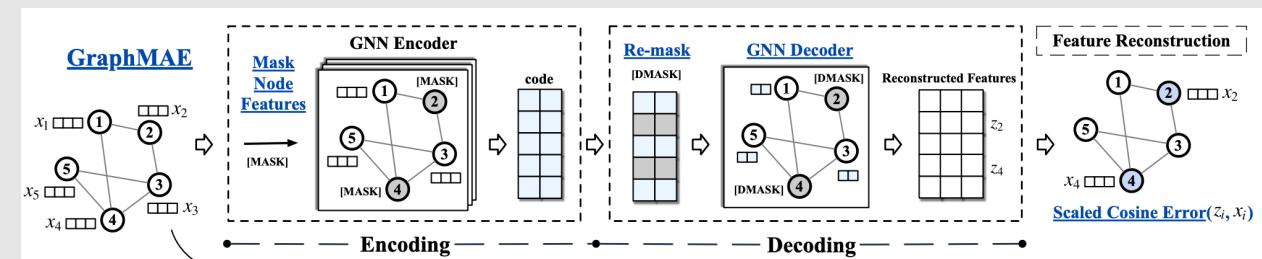
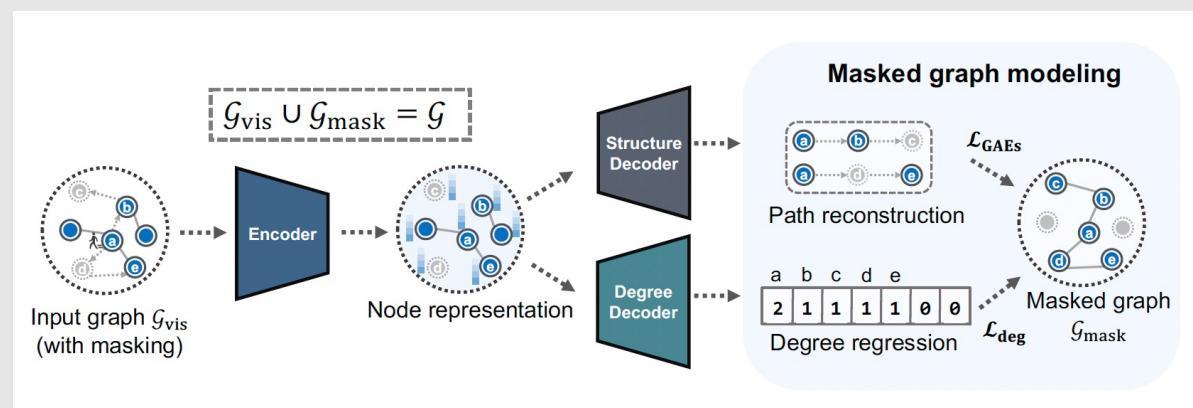
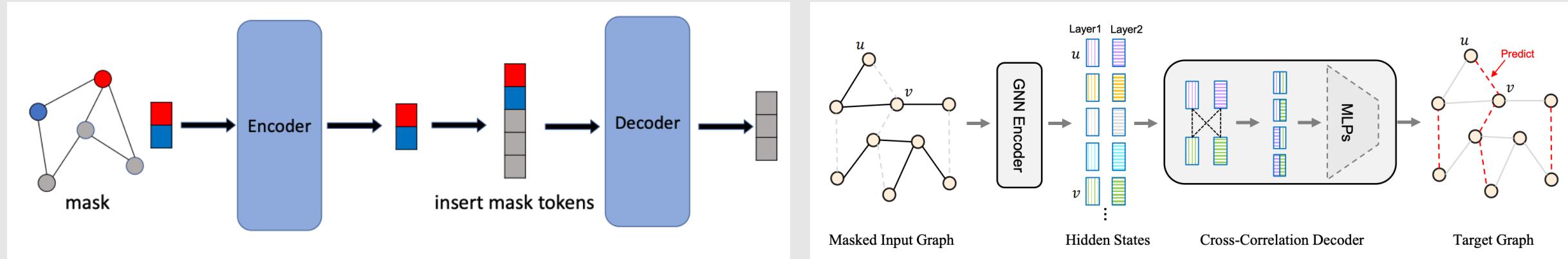
- **Medical Imaging**

# Masked Autoencoder: Extensions



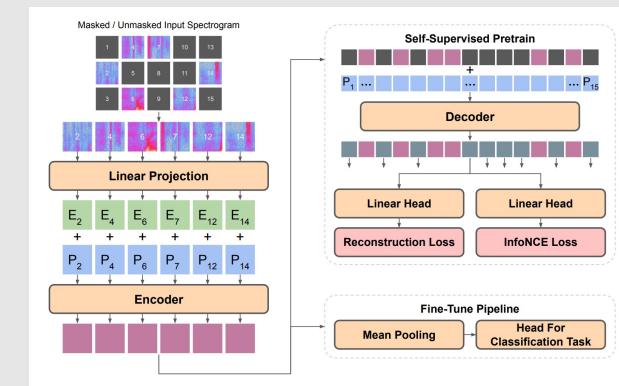
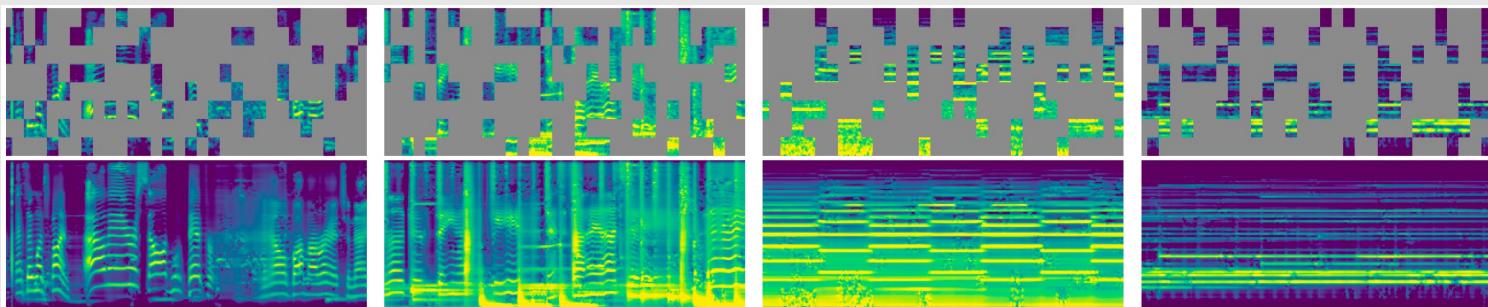
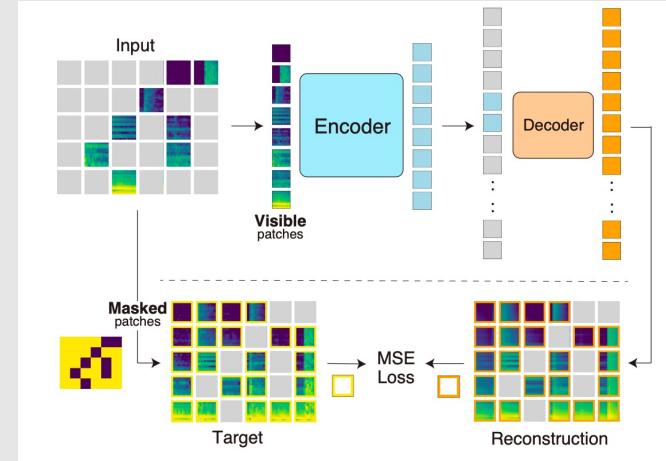
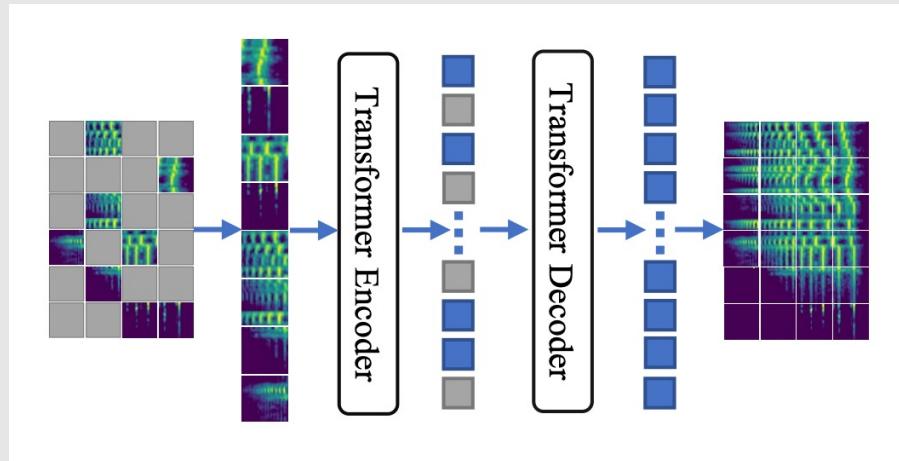
- **3D Geometry**

# Masked Autoencoder: Extensions

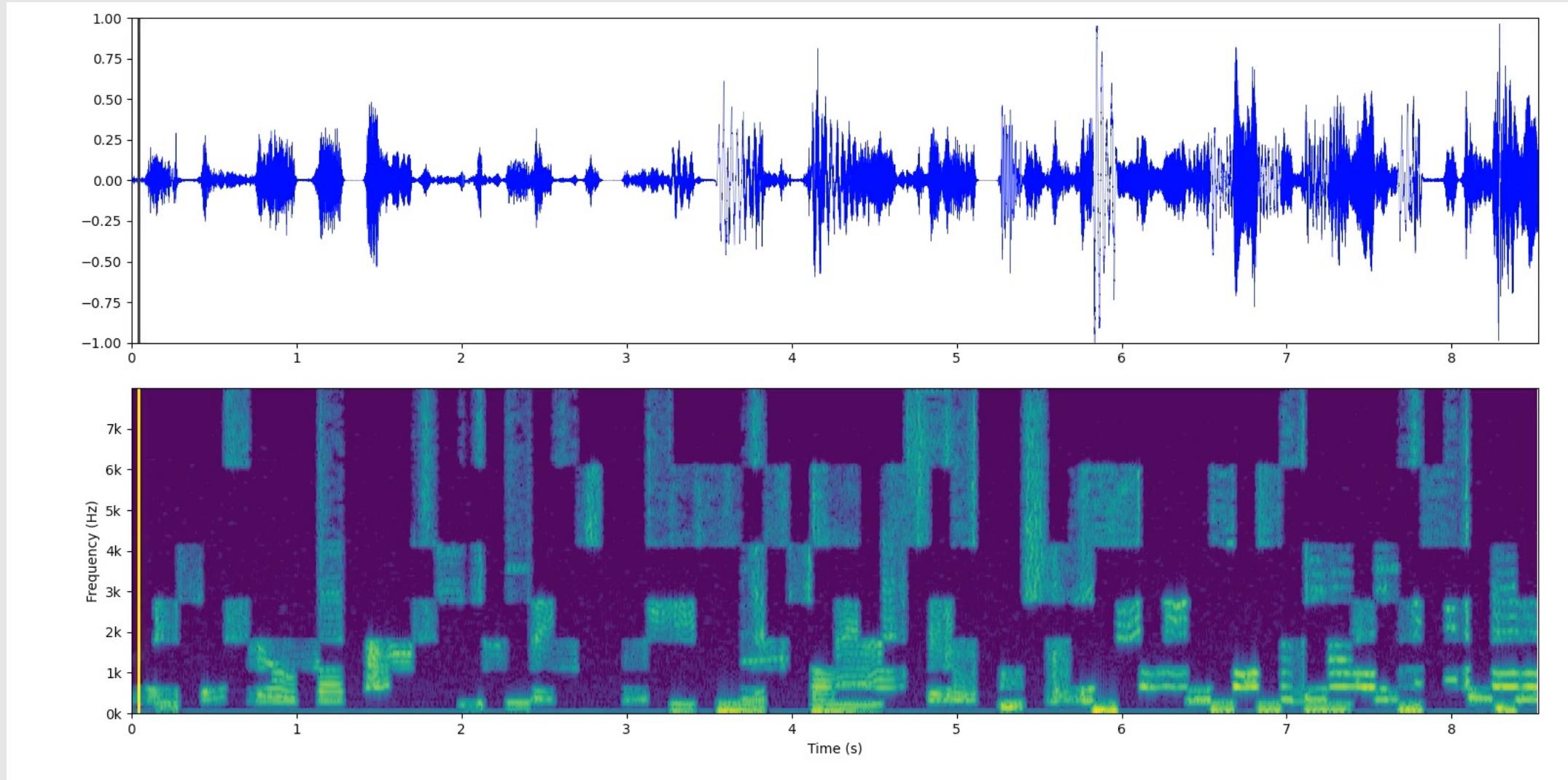


- **Graphs**

# Masked Autoencoder: Extensions



- **Audio**



# Audio MAE

# Summary

**Unsupervised learning** - classical (but very useful) methods:

- K-means clustering
- Principal Component Analysis (PCA)
- Autoencoder

**Self-supervised learning** - modern advances:

- Contrastive Learning
- Predictive Learning

# Lecture 12: Representation Learning

- Why should we care about Representation Learning?
- Transfer learning
- Unsupervised learning
- Self-supervised learning