

Lecture 8

Representation Learning: Intro & Reconstruction-based Methods

Speaker: Kaiming He



Overview

- Introduction to Representation Learning
- Unsupervised Learning and Self-supervised Learning
- Reconstruction-based Methods

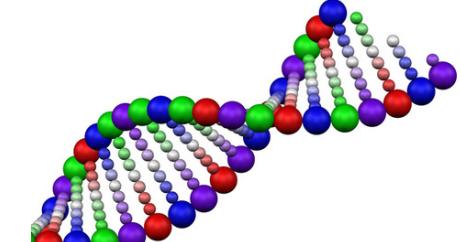
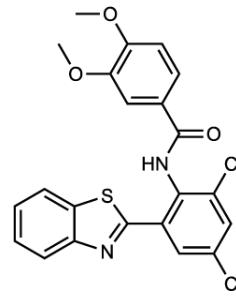
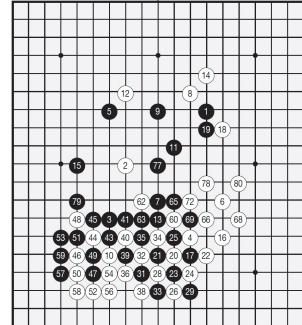
Deep Learning is Representation Learning

- Represent raw data to solve complex problems
 - compression, abstraction, conceptualization
- Raw data in different forms
 - pixels, words, waves, states, molecules, DNA, ...



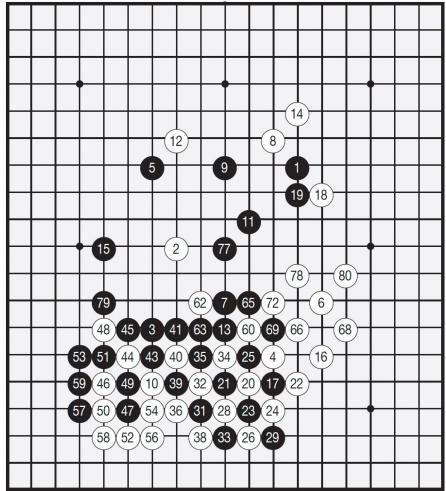
Current approaches to object recognition make essential use of machine learning methods. To improve their performance, we can collect larger datasets, learn more powerful models, and use better techniques for preventing overfitting. Until recently, datasets of labeled images were relatively small compared to the size of training sets used in NLP, such as the Penn Treebank [1] and the CIFAR-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 image classification model [2] was trained on a dataset of 15 million images [1]. 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 recognized for some time (e.g., [13]), so there has been a push to collect and to collect labeled 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 high-resolution images in over 22,000 categories.

To keep the number of parameters in a neural network reasonable, we need to learn a model with large learning capacity. However, the immense complexity of the object recognition task means that this problem cannot be solved even by a dataset as large as ImageNet, so our model should also have lots of power. The most common way to increase the power of a neural network is to add layers (CNNs) constitute one such class of models [16, 11, 13, 18, 15, 22, 26]. Their capacity can be controlled by varying their depth and breadth, and they also make strong and mostly correct assumptions about the nature of images (isometry, stationarity of statistics, and locality of dependencies). Thus, it is natural to compare feed-forward neural networks with shallowly-sized layers, CNNs have much fewer connections and parameters and so they are easier to train, while their theoretically-best performance is likely to be only slightly worse.



Defines a conference name: ICLR
(International Conference on Learning Representations)

Deep Learning is Representation Learning

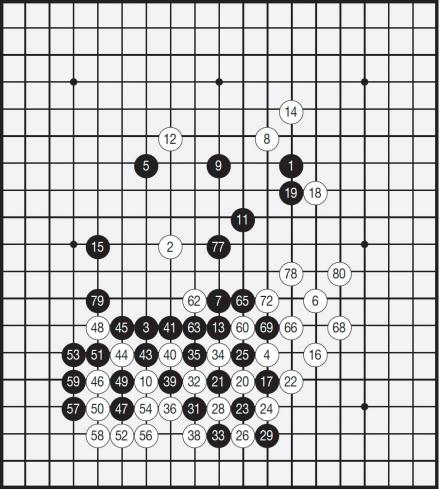


3^{361} states?

Game of Go

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

Deep Learning is Representation Learning



3^{361} states?

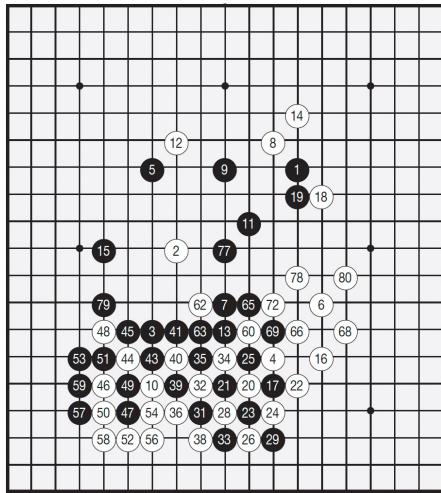
*Bad
representation*



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

- Image space has exponentially more states than Go.

Deep Learning is Representation Learning

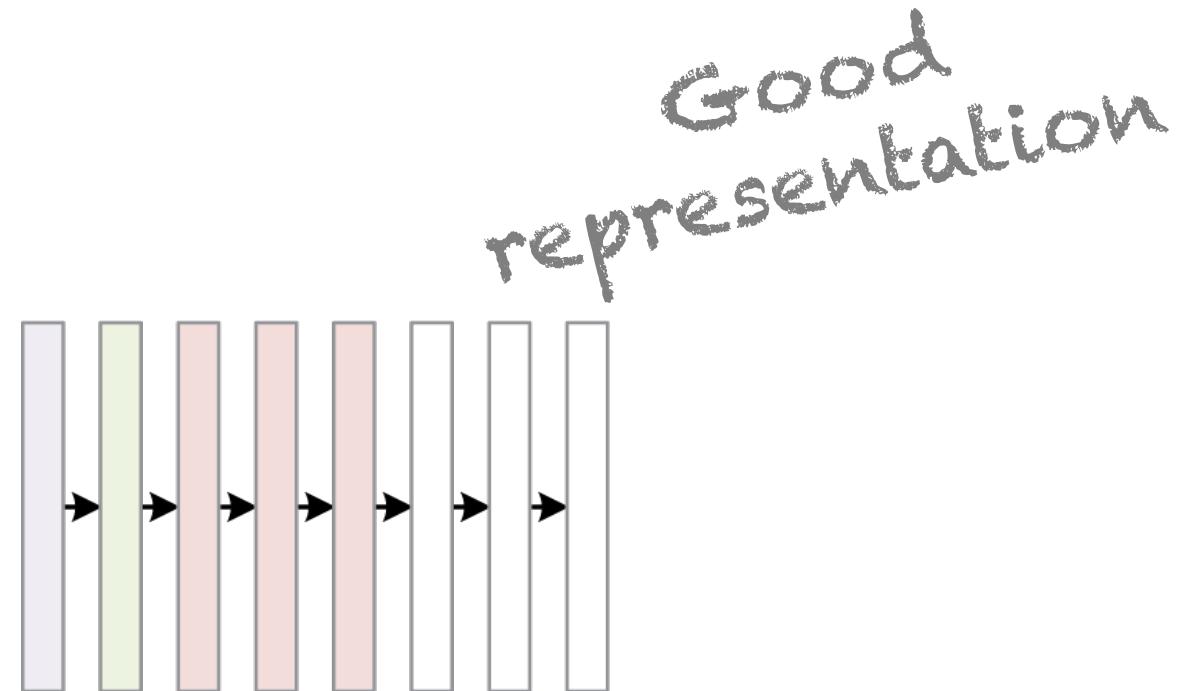


3^{361} states?

- Image recognition is solved in representation space.

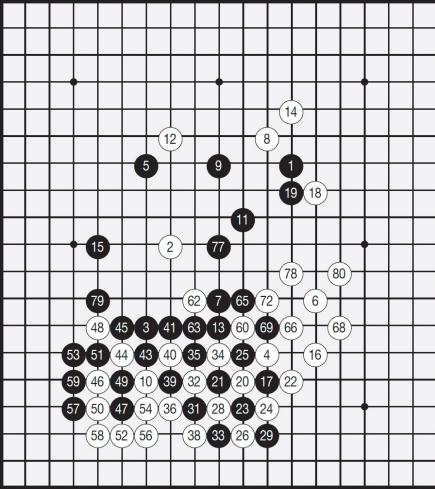


$256^{3 \times 500 \times 500}$?

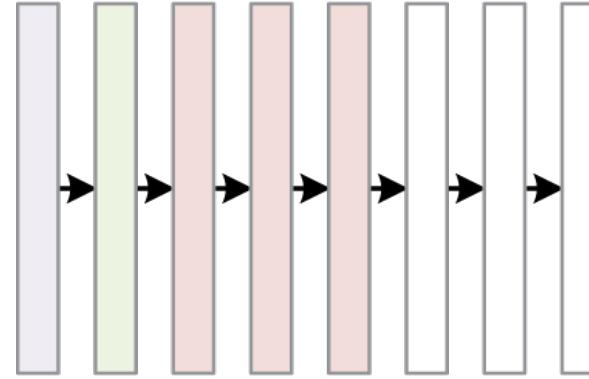


Deep Learning is Representation Learning

- Go playing can be solved in representation space.



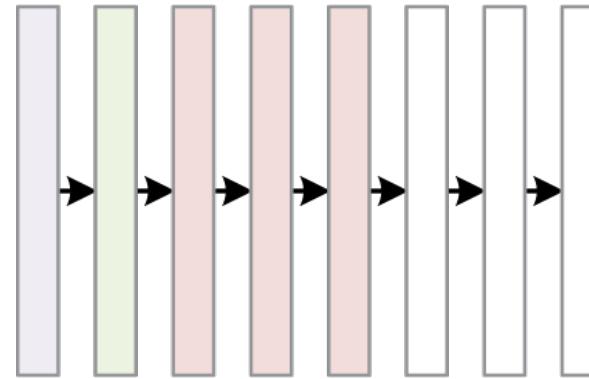
3^{361} states?



Good
representation



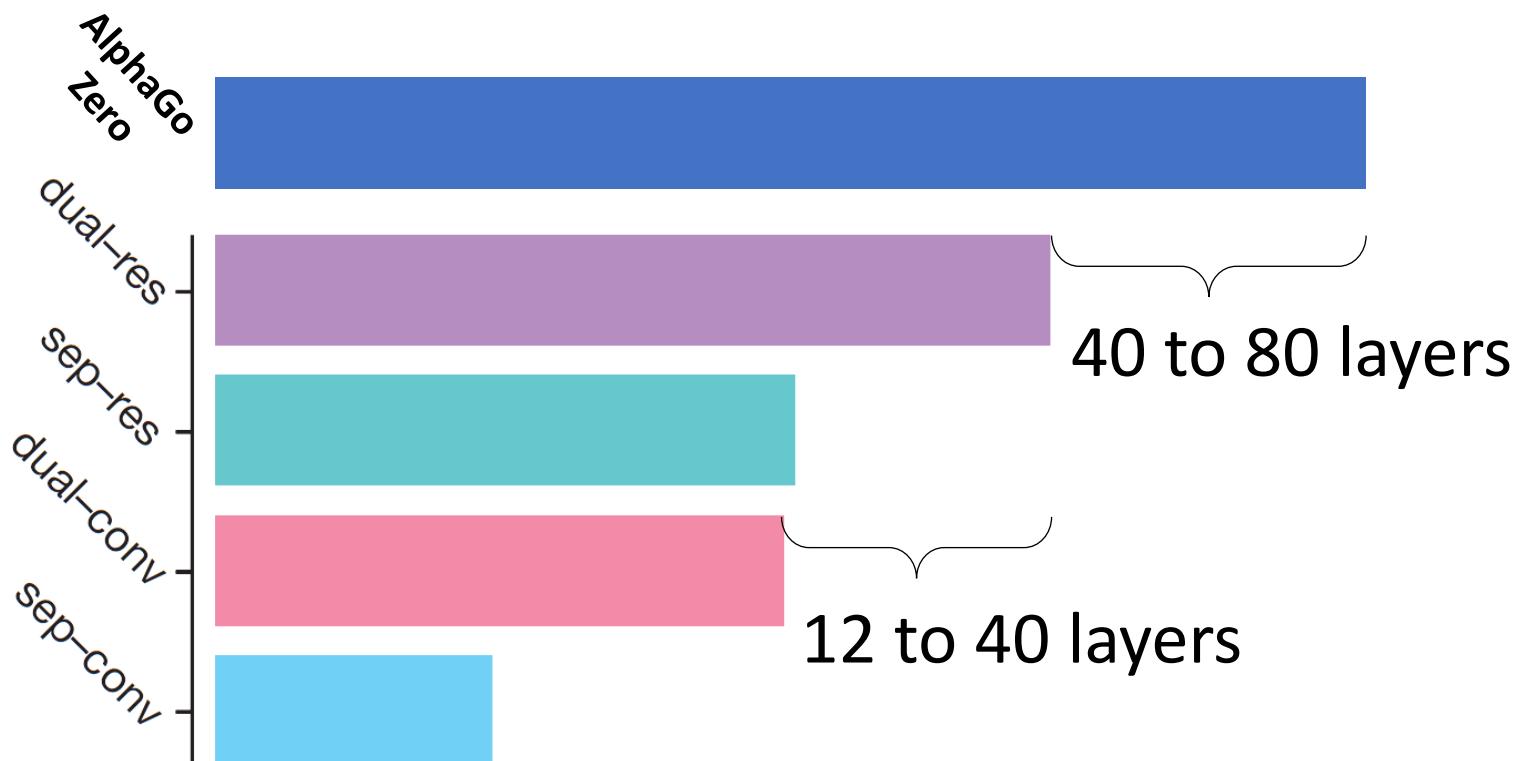
$256^{3 \times 500 \times 500}$?



Deep Learning is Representation Learning

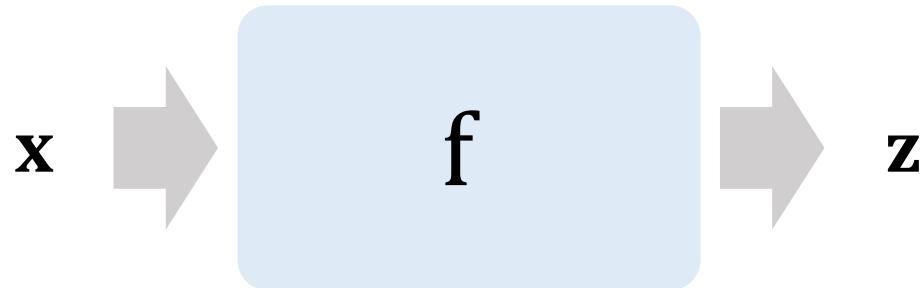
AlphaGo: better representation leads to higher rating

- outperform best human players
- without human knowledge



What is a representation?

- A representation of a data domain \mathcal{X} is a function $f: \mathcal{X} \rightarrow R^d$ that assigns a feature vector to each input in that domain.
- A representation of a datapoint x is a vector $z \in R^d$ with $z = f(x)$.



can be extended to:

- high-dim arrays: $d = d_1 \times d_2 \times \dots$
- other data structures: trees, list of arrays (pyramids), ...

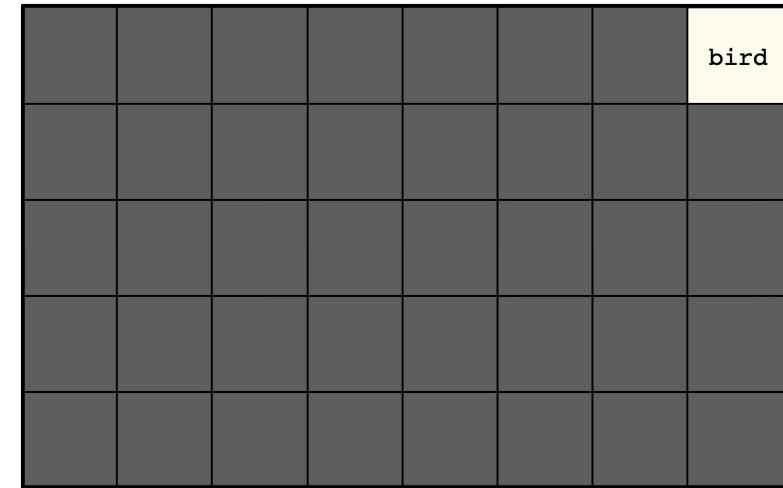
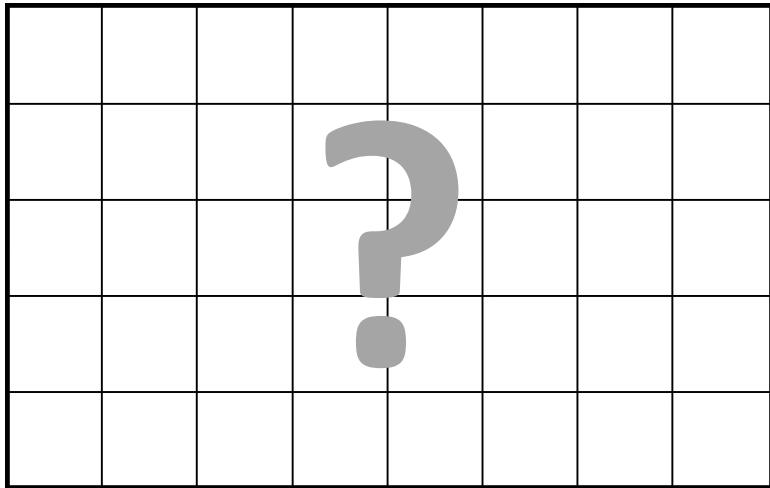
What does a representation represent?



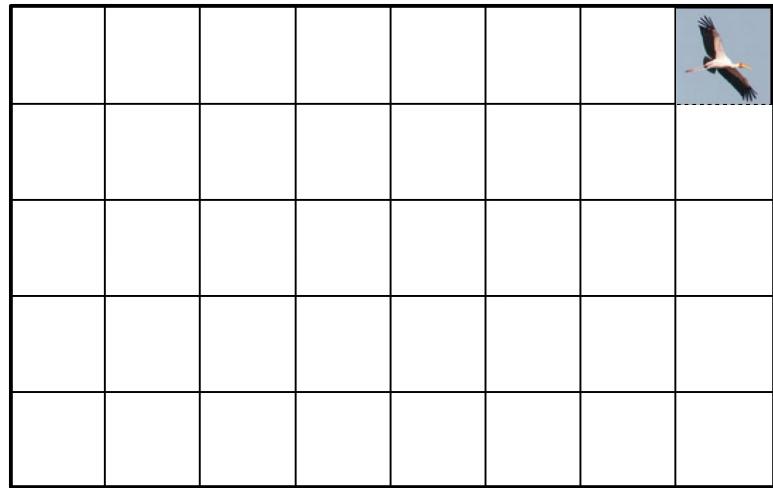
sky	sky	sky	sky	sky	sky	sky	bird
sky	sky	sky	sky	sky	sky	sky	sky
sky	sky	sky	sky	sky	sky	sky	sky
bird	bird	bird	bird	sky	bird	sky	sky
sky	sky	sky	bird	bird	sky	sky	sky

bird = [0, 0, 1, 0, 0, ..., 0]
sky = [0, 0, 0, 0, 1, ..., 0]

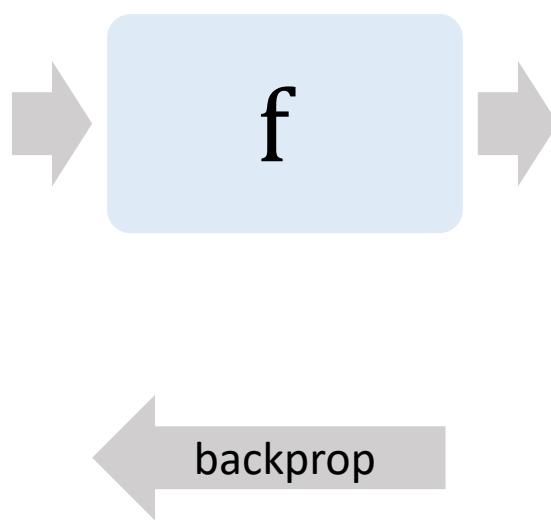
What does a representation represent?



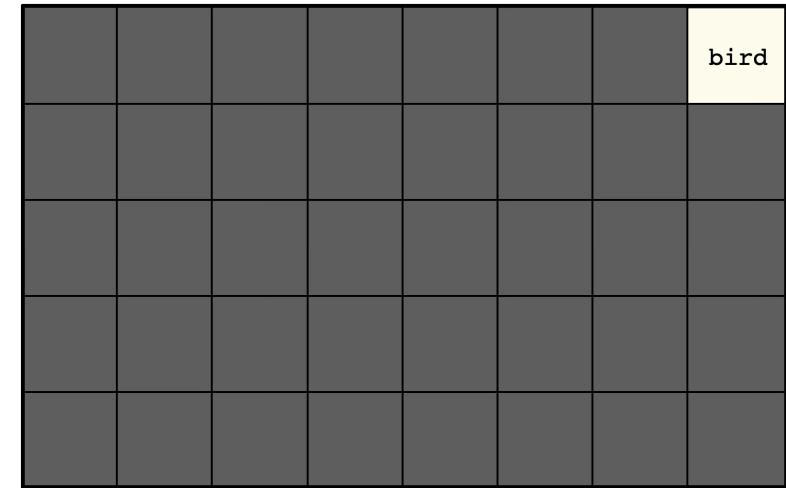
What does a representation represent?



x

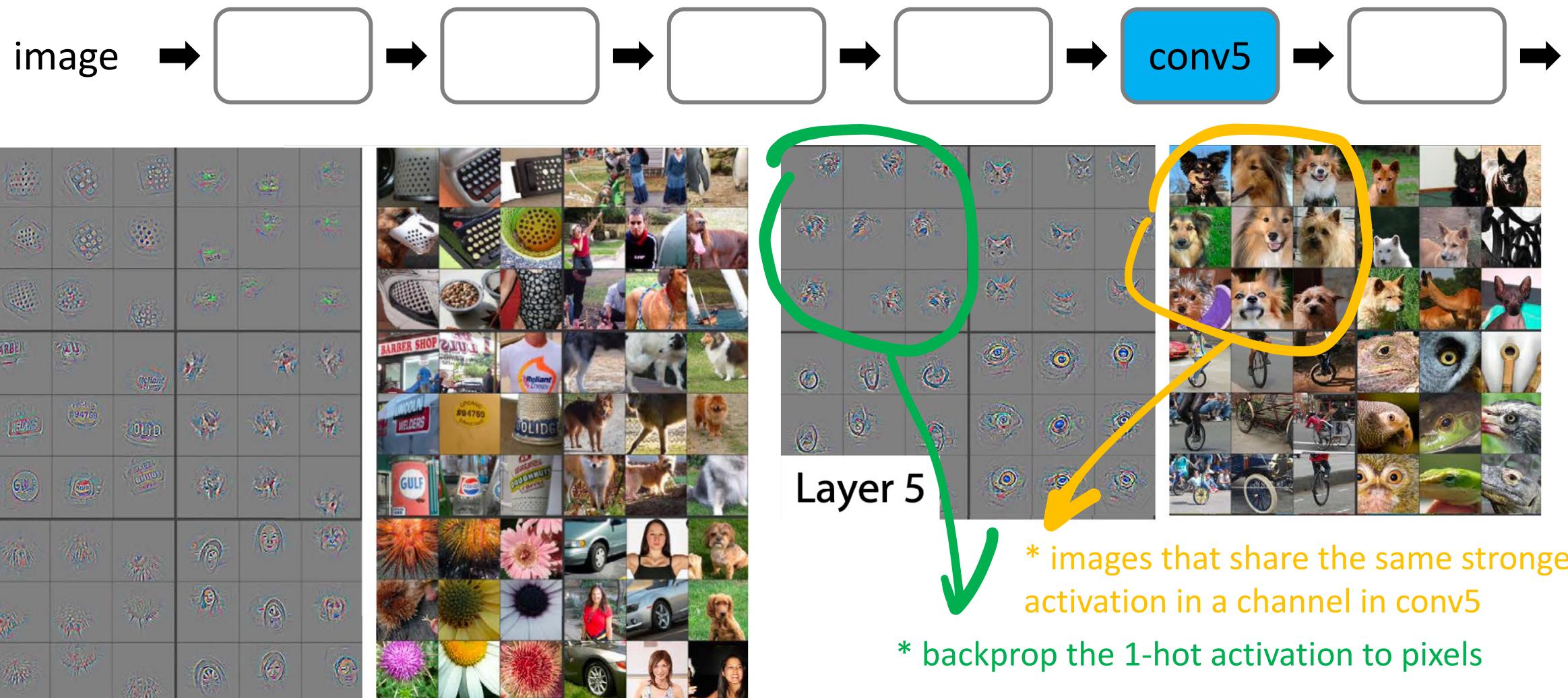


find x
s.t. $z = f(x)$

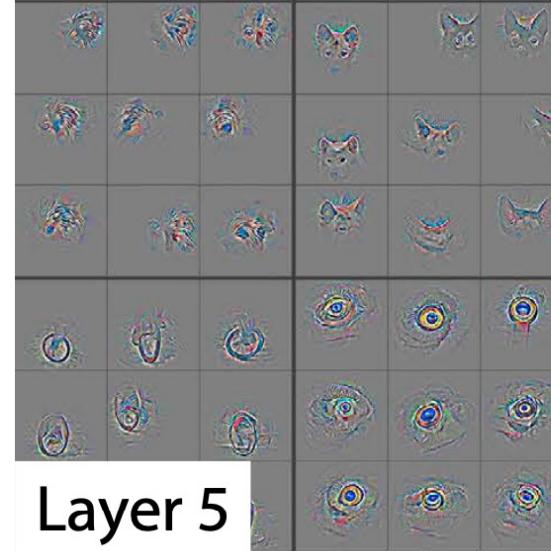
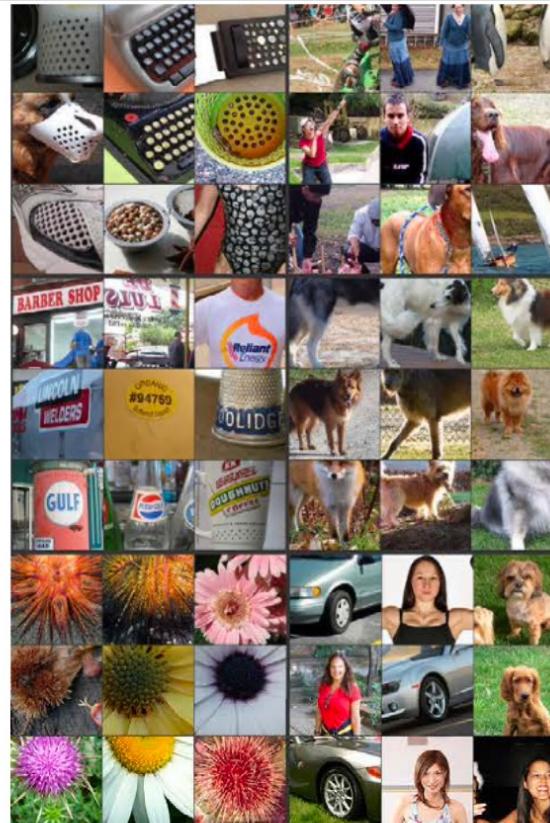
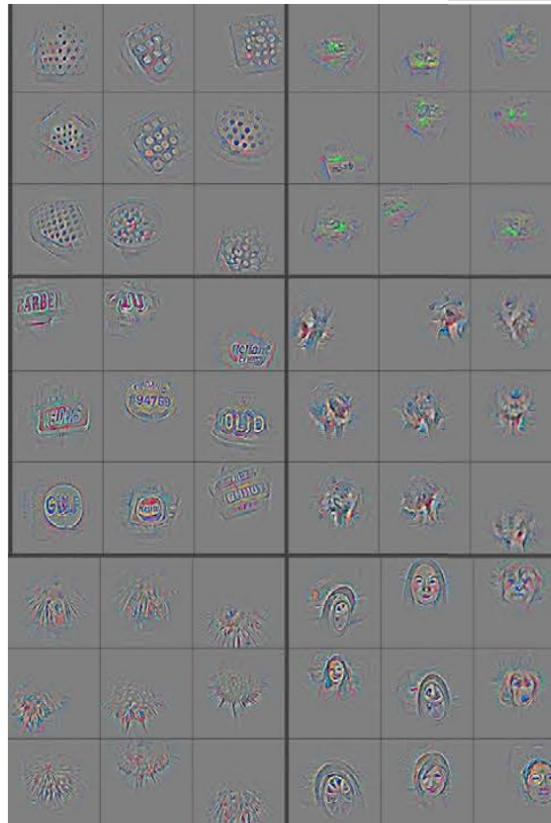


z

What does a representation represent?



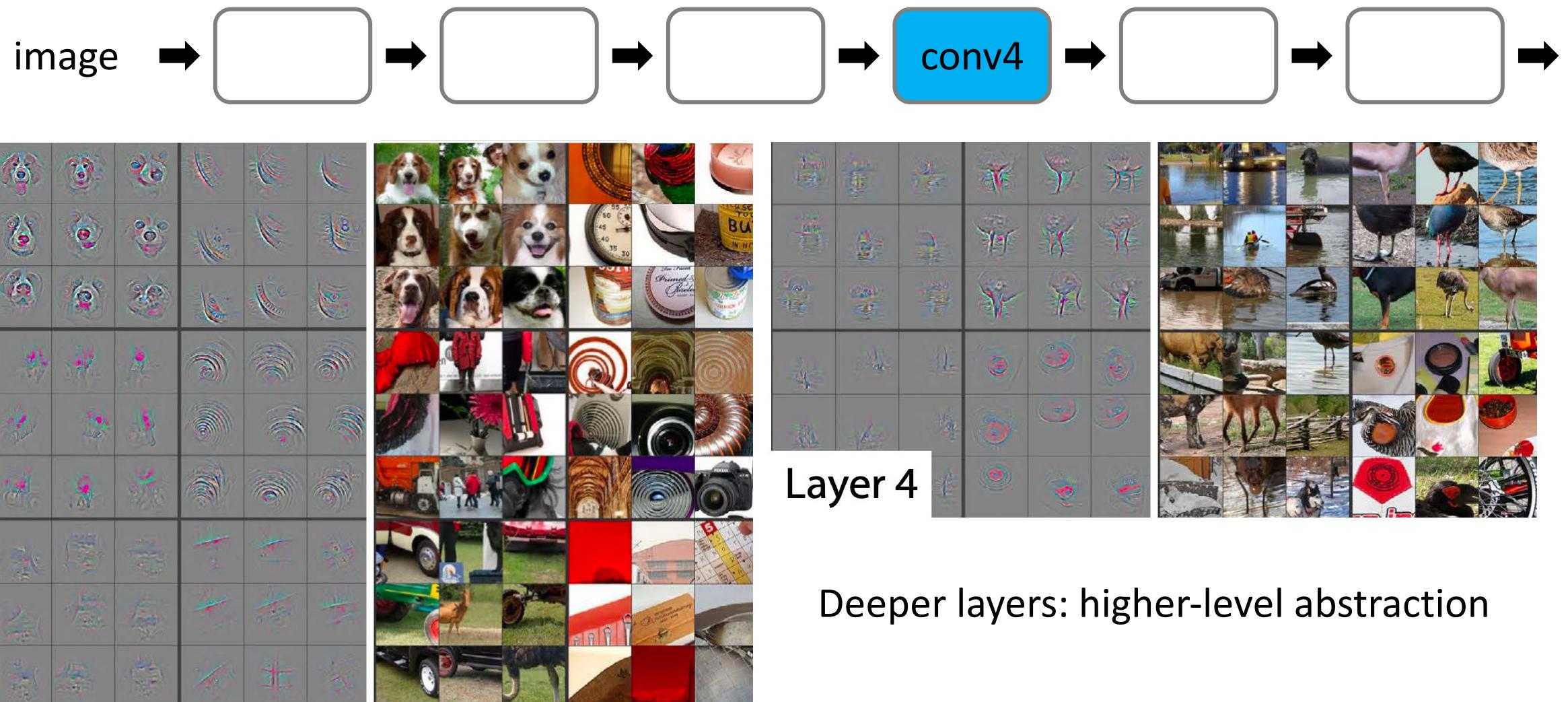
What does a representation represent?



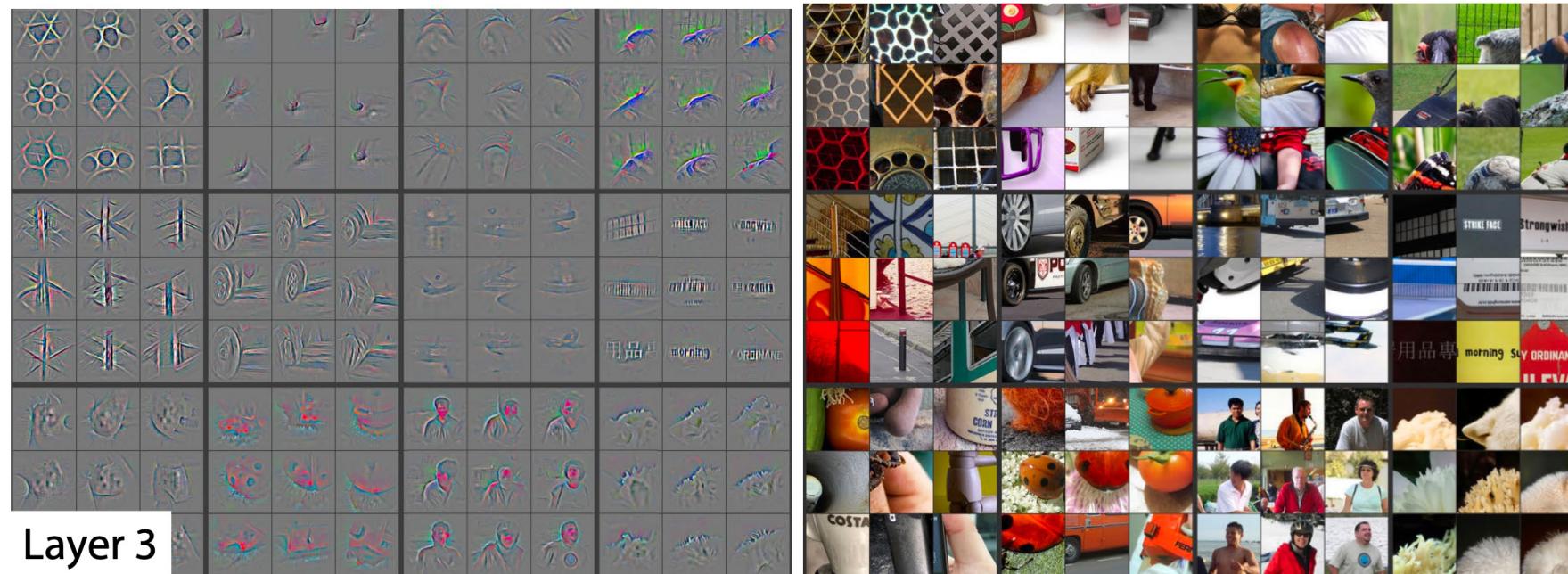
Layer 5

Deeper layers: higher-level abstraction

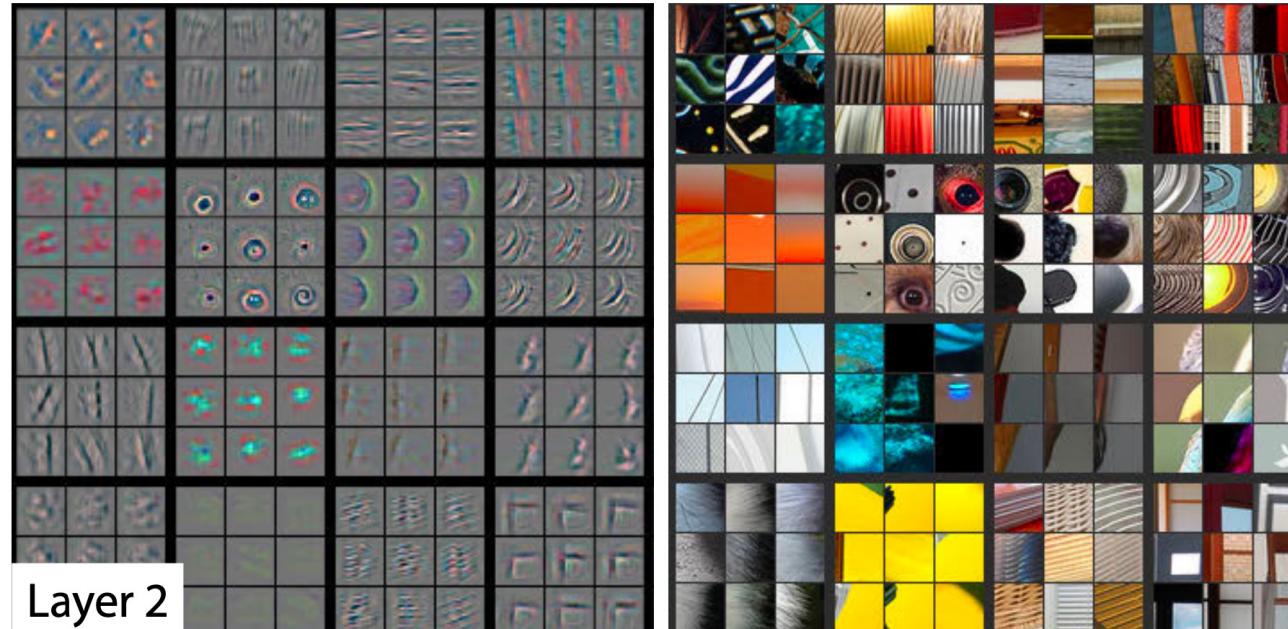
What does a representation represent?



What does a representation represent?

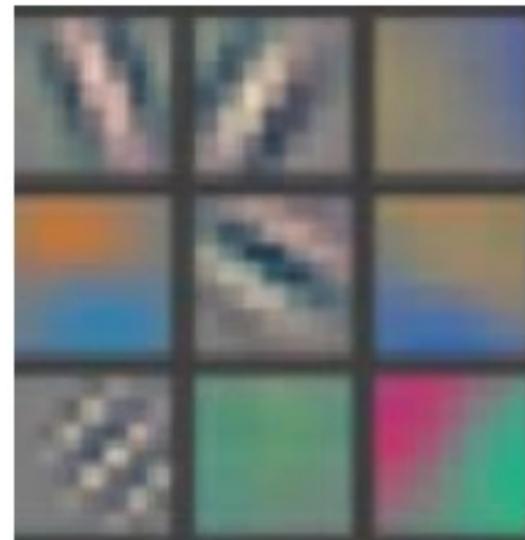


What does a representation represent?



low-level abstraction

What does a representation represent?



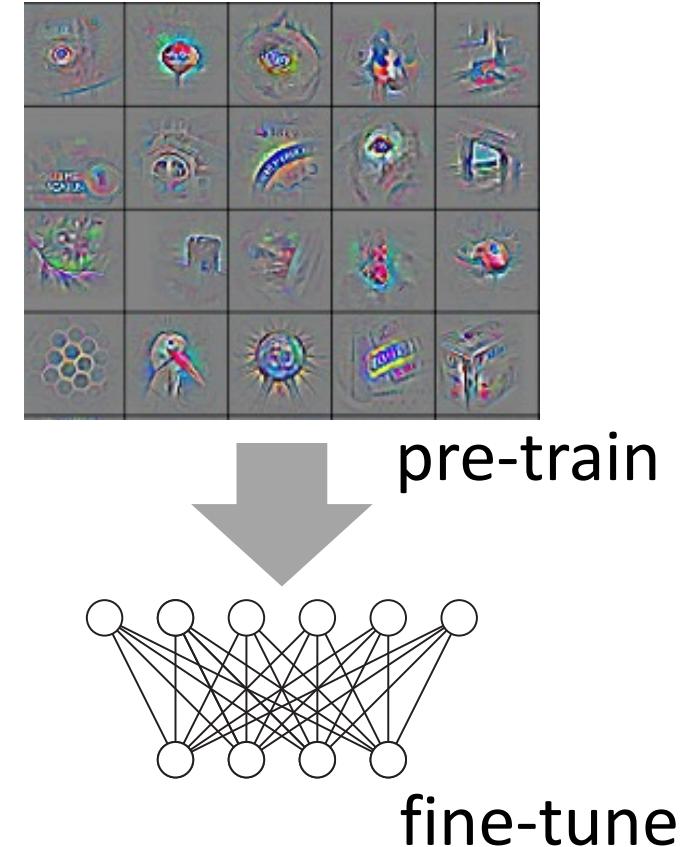
low-level abstraction

Deep Representations are Transferrable

Transfer learning:

- pre-train: large data
- fine-tune: small data
- enable DL for small data
- revolutionize CV and many areas (LLM)
- data: engine for learning representations

The single most important discovery
in DL revolution!



"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

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



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



Partial transfer

- pre-train and target domains may differ
- high-level representations too specialized
- randomly initialize new layers

Transfer learning

Pre-training



Fine-tuning

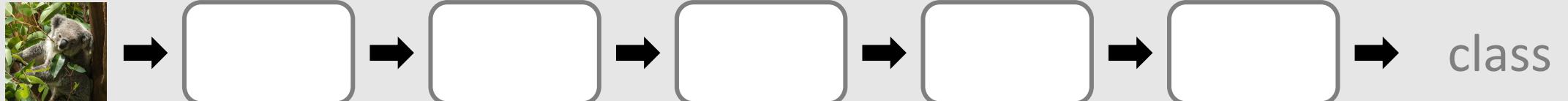


Frozen weights

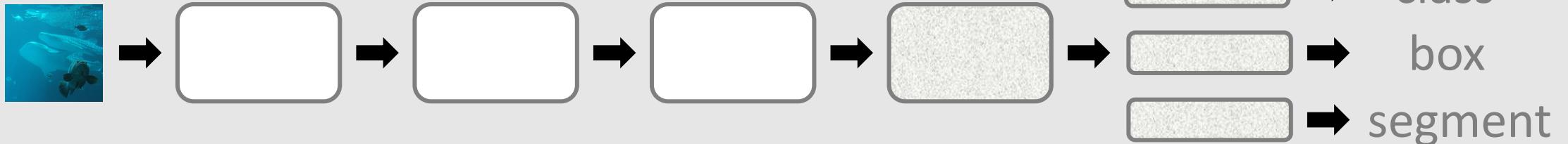
- freeze some/all pre-trained weights
- reduce overfitting if data is too little
- save memory, speed up training

Transfer learning

Pre-training



Fine-tuning



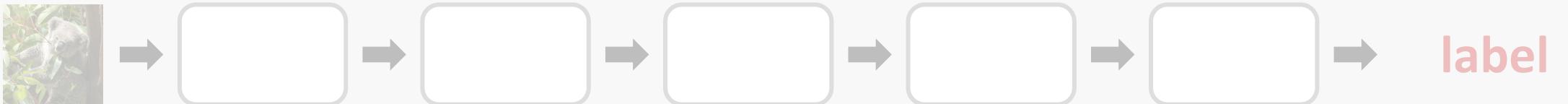
Network surgery

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

Pre-training Objectives

- **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

supervised



un-/self-supervised



transfer



Reconstruction-based Methods

for Unsupervised Representation Learning

Reconstruction: classical yet highly powerful

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



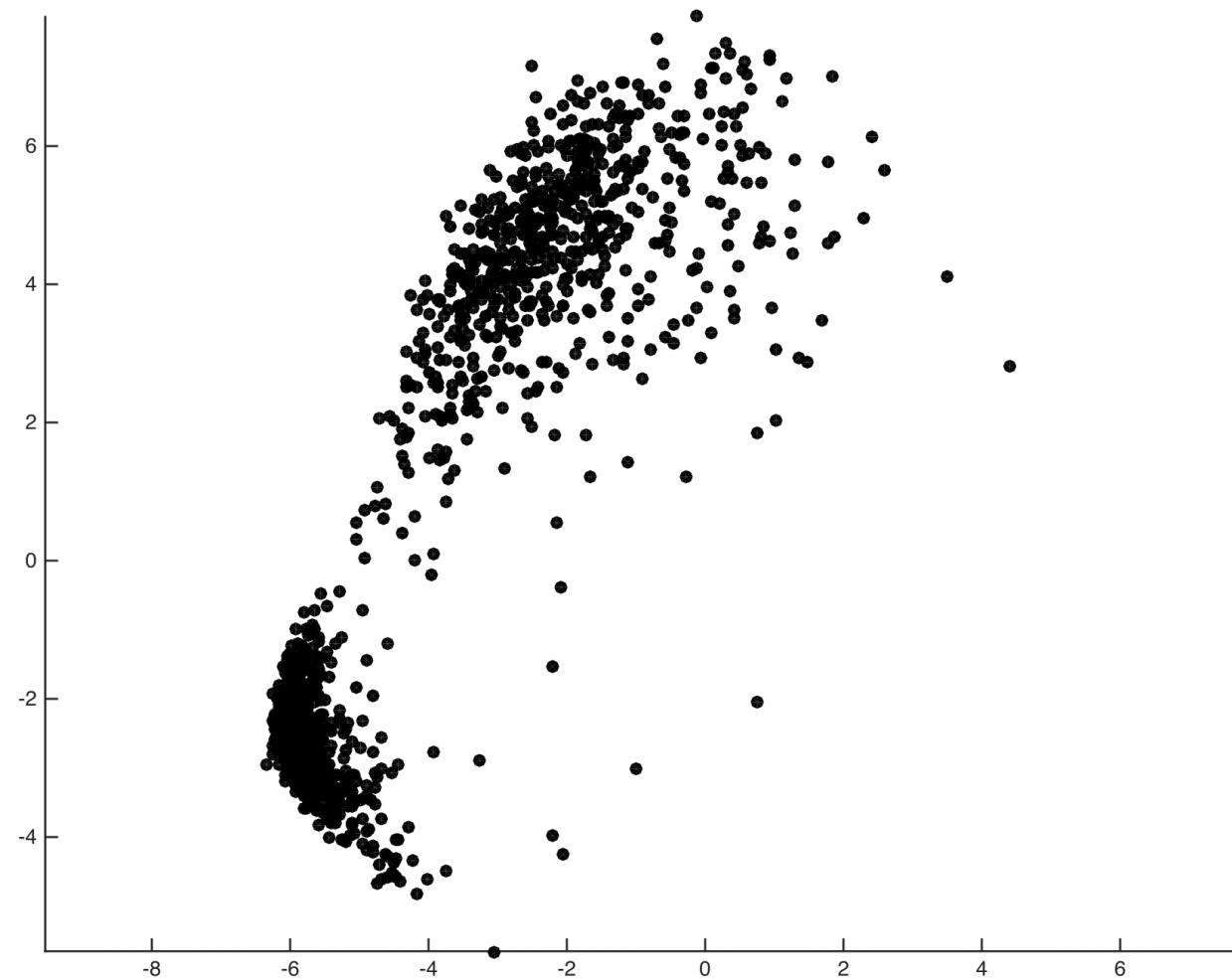
Reconstruction: classical yet highly powerful

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



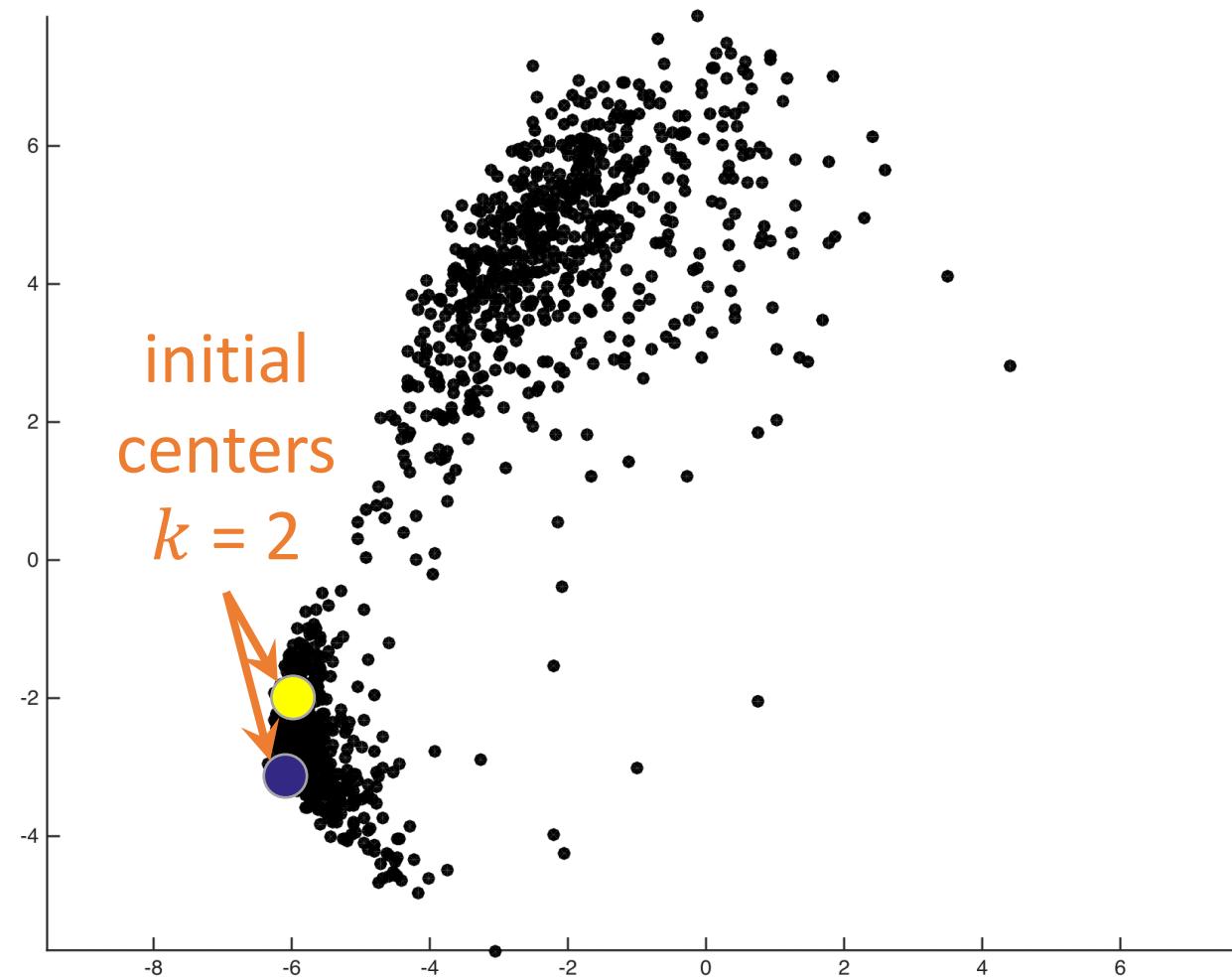
K-means clustering

Represent data by k centers



K-means clustering

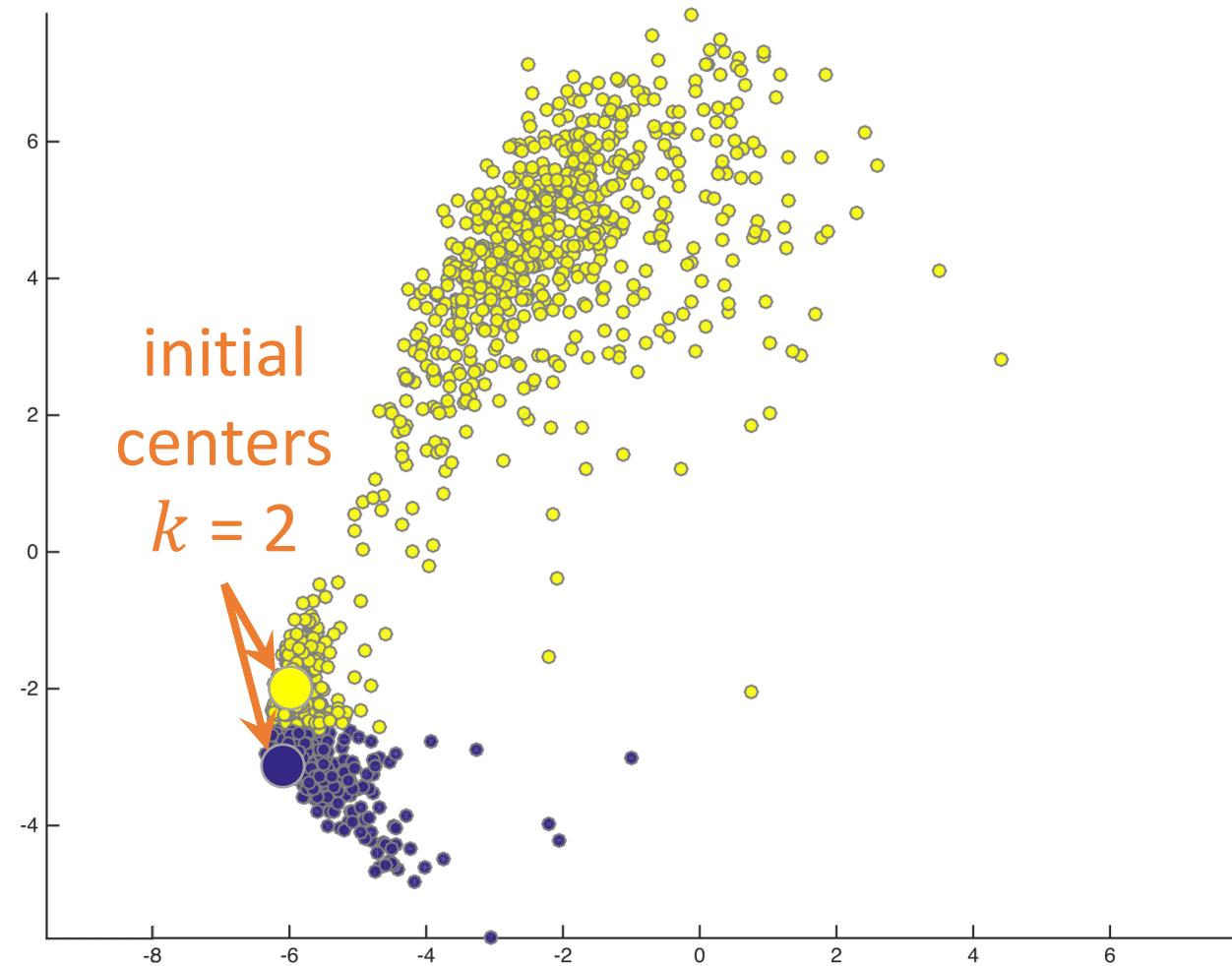
Represent data by k centers



K-means clustering

Represent data by k centers

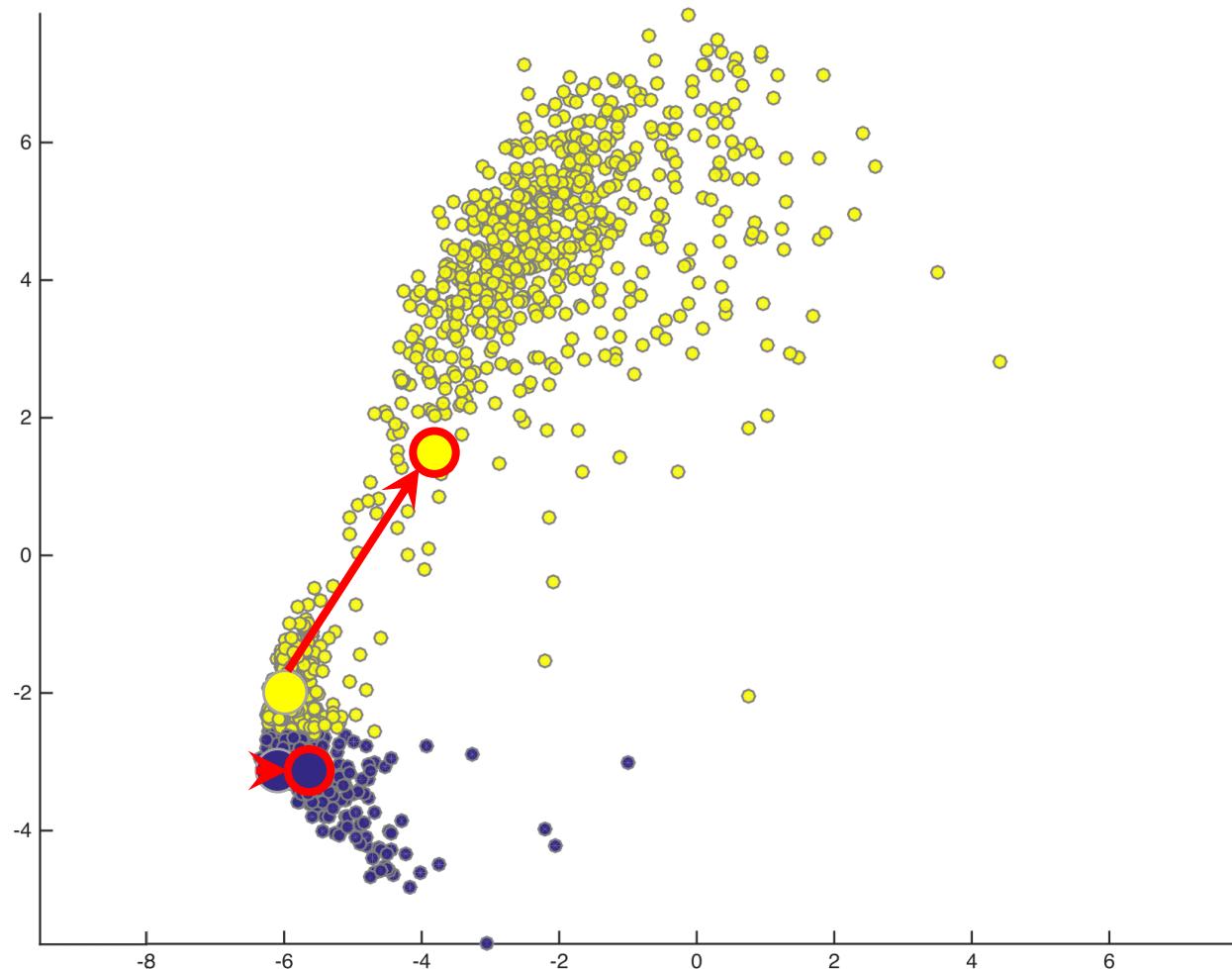
- Assignment: each sample is assigned to the nearest center



K-means clustering

Represent data by k centers

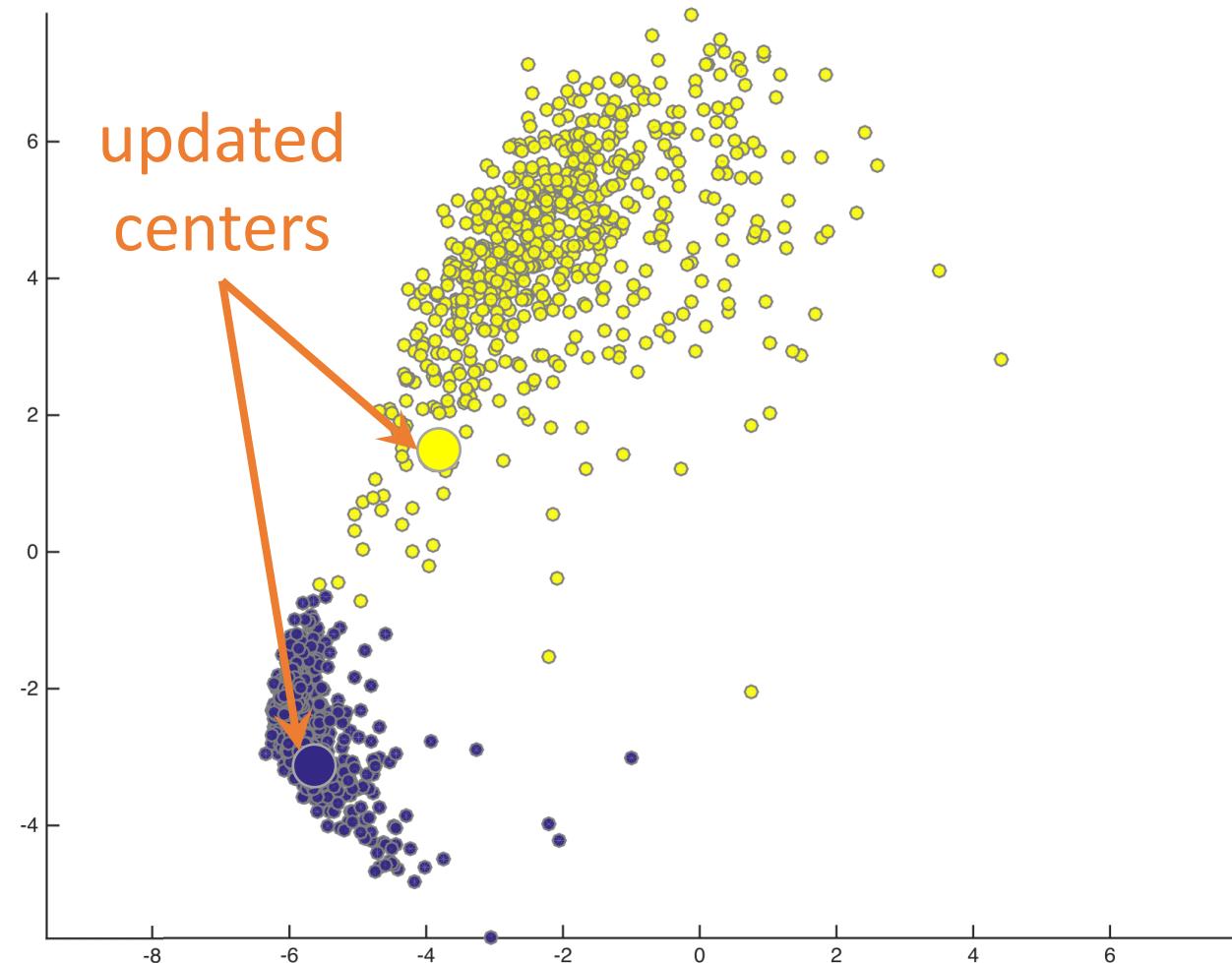
- **Assignment:** each sample is assigned to the nearest center
- **Update:** each center is updated by the mean of samples assigned to it



K-means clustering

Represent data by k centers

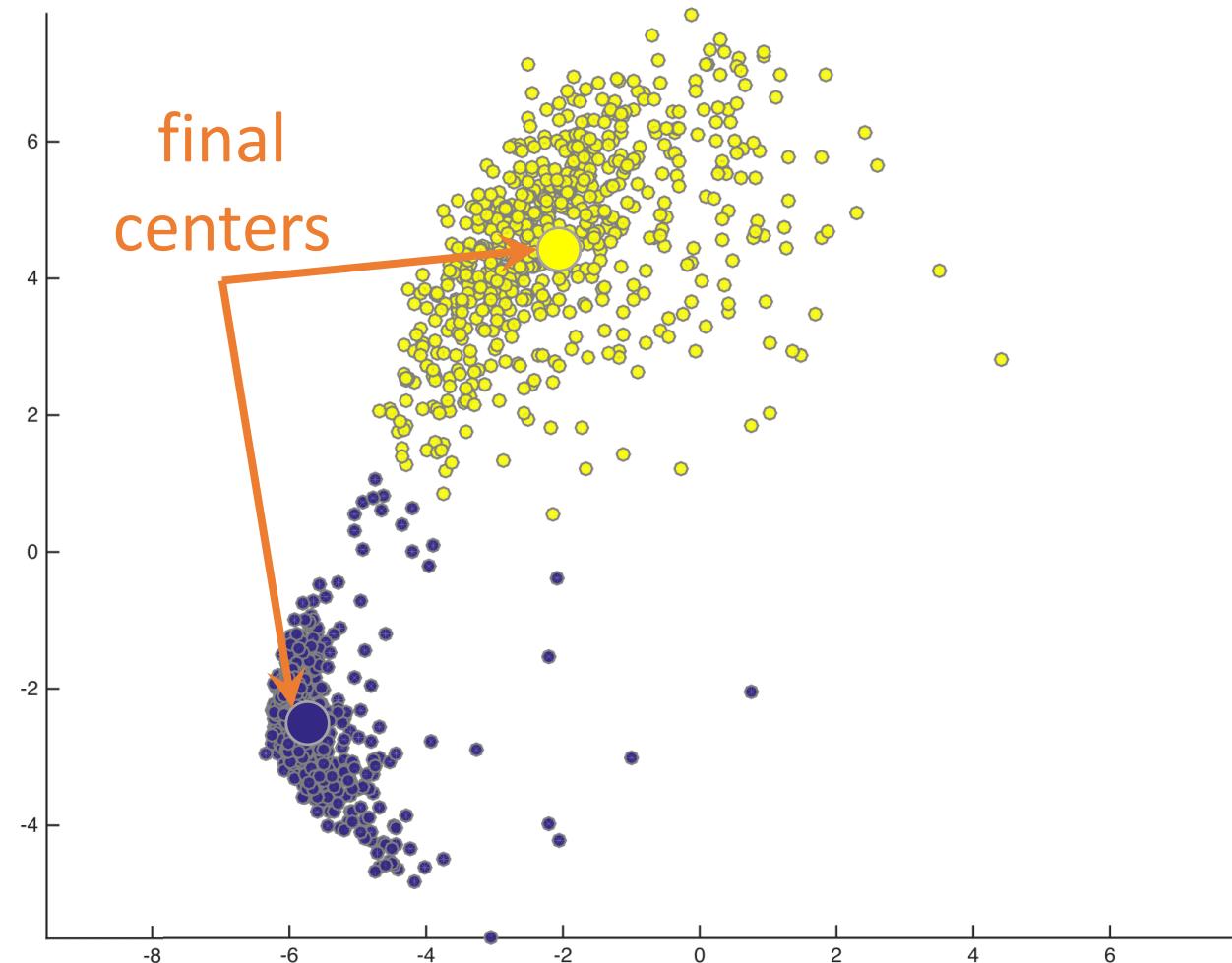
- **Assignment:** each sample is assigned to the nearest center
- **Update:** each center is updated by the mean of samples assigned to it



K-means clustering

Represent data by k centers

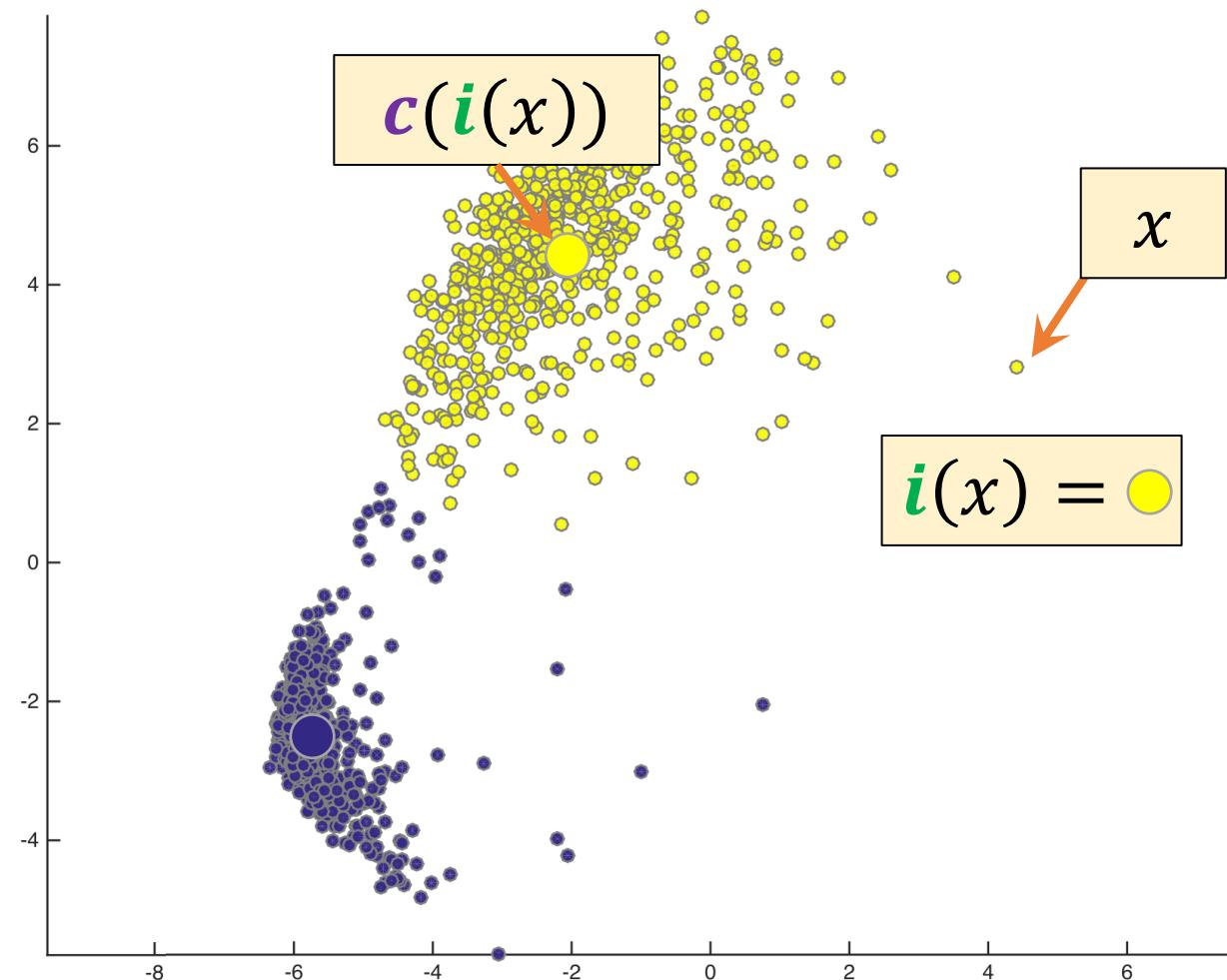
- **Assignment:** each sample is assigned to the nearest center
- **Update:** each center is updated by the mean of samples assigned to it
- Repeat until converge



K-means clustering

Loss function:

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

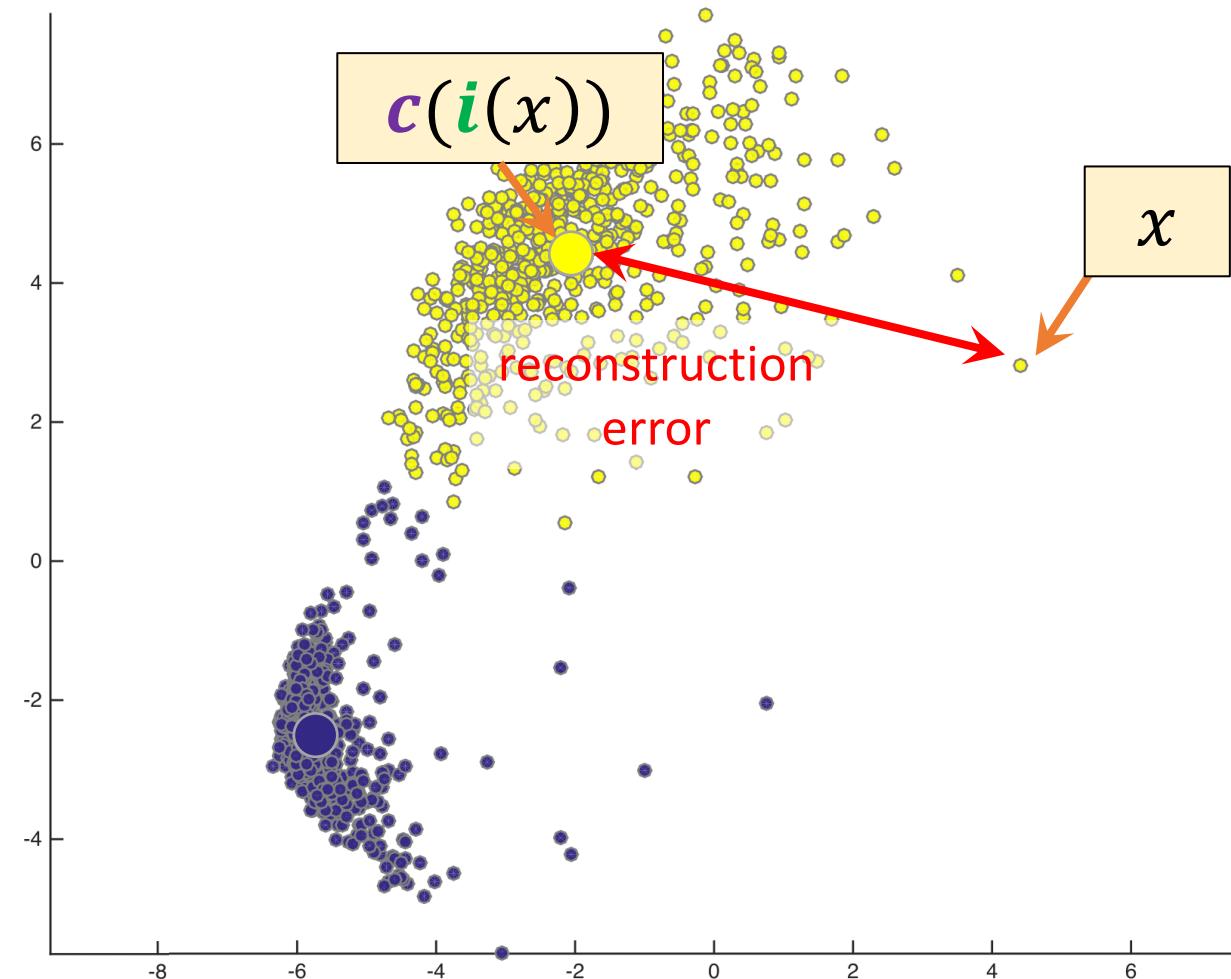


K-means clustering

Loss function:

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

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



K-means clustering

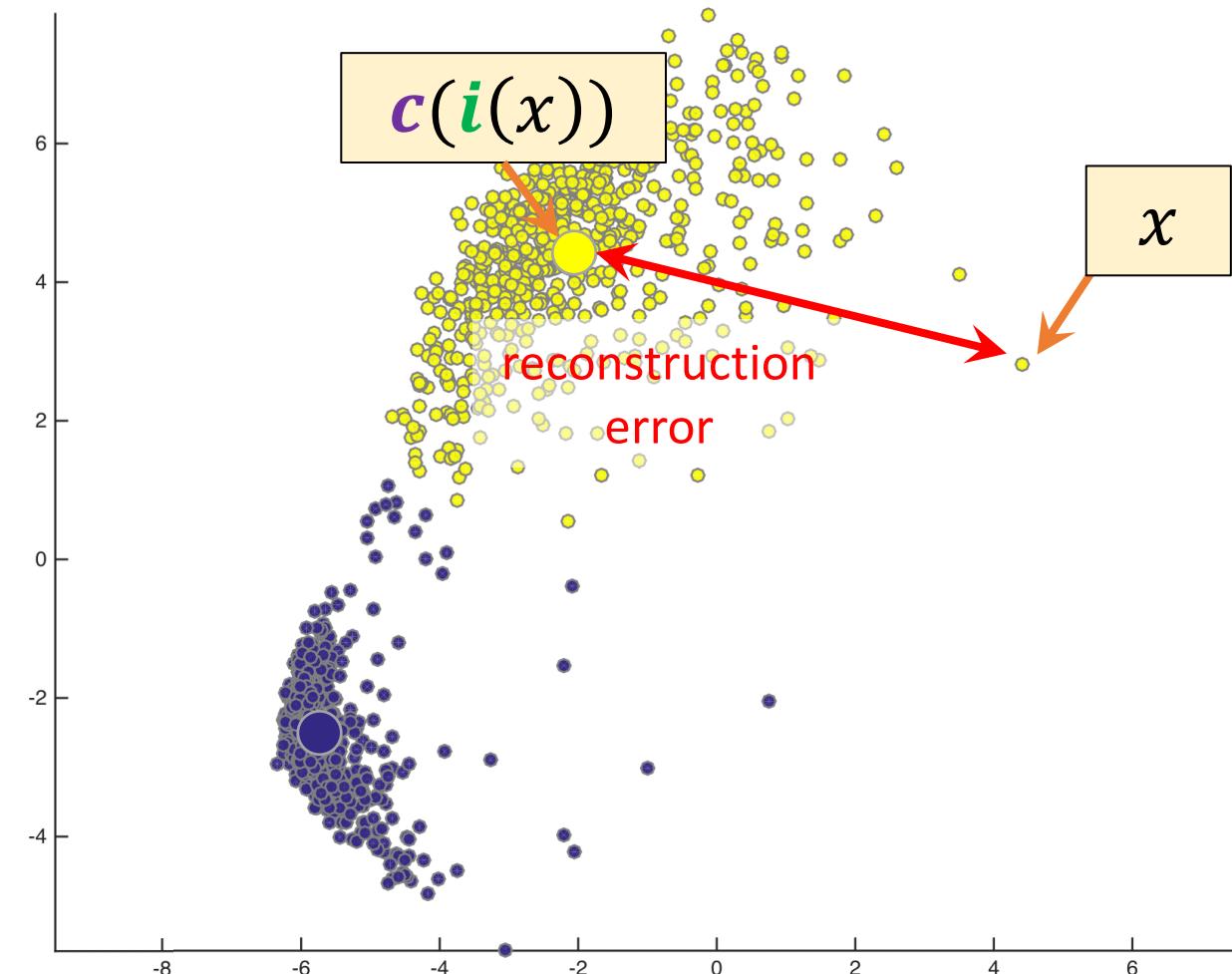
Loss function:

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

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:

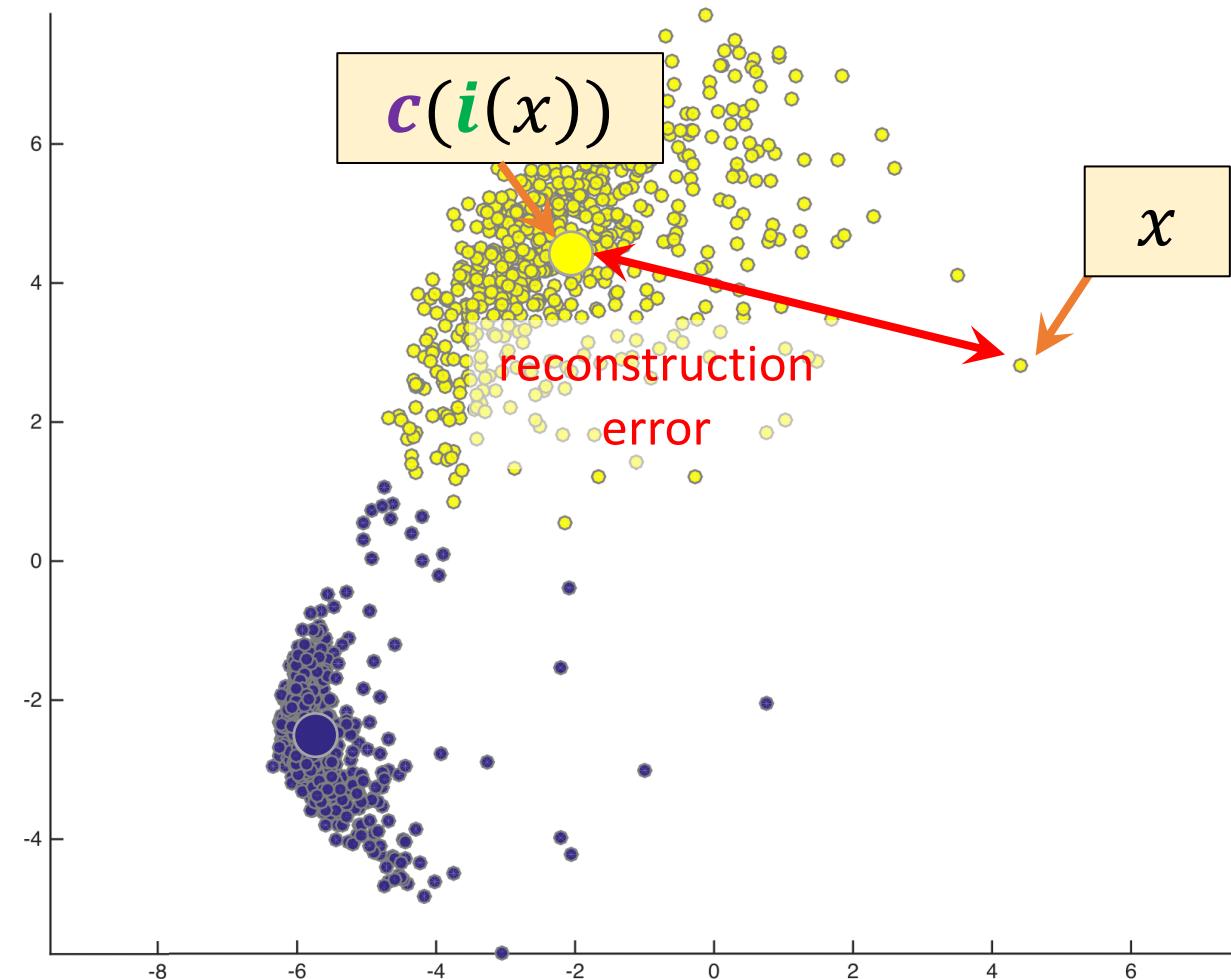
$$\min_{\mathbf{c}, \mathbf{i}} \sum_x \|x - \mathbf{c}(i(x))\|^2$$

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

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

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

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



K-means clustering

Loss function:

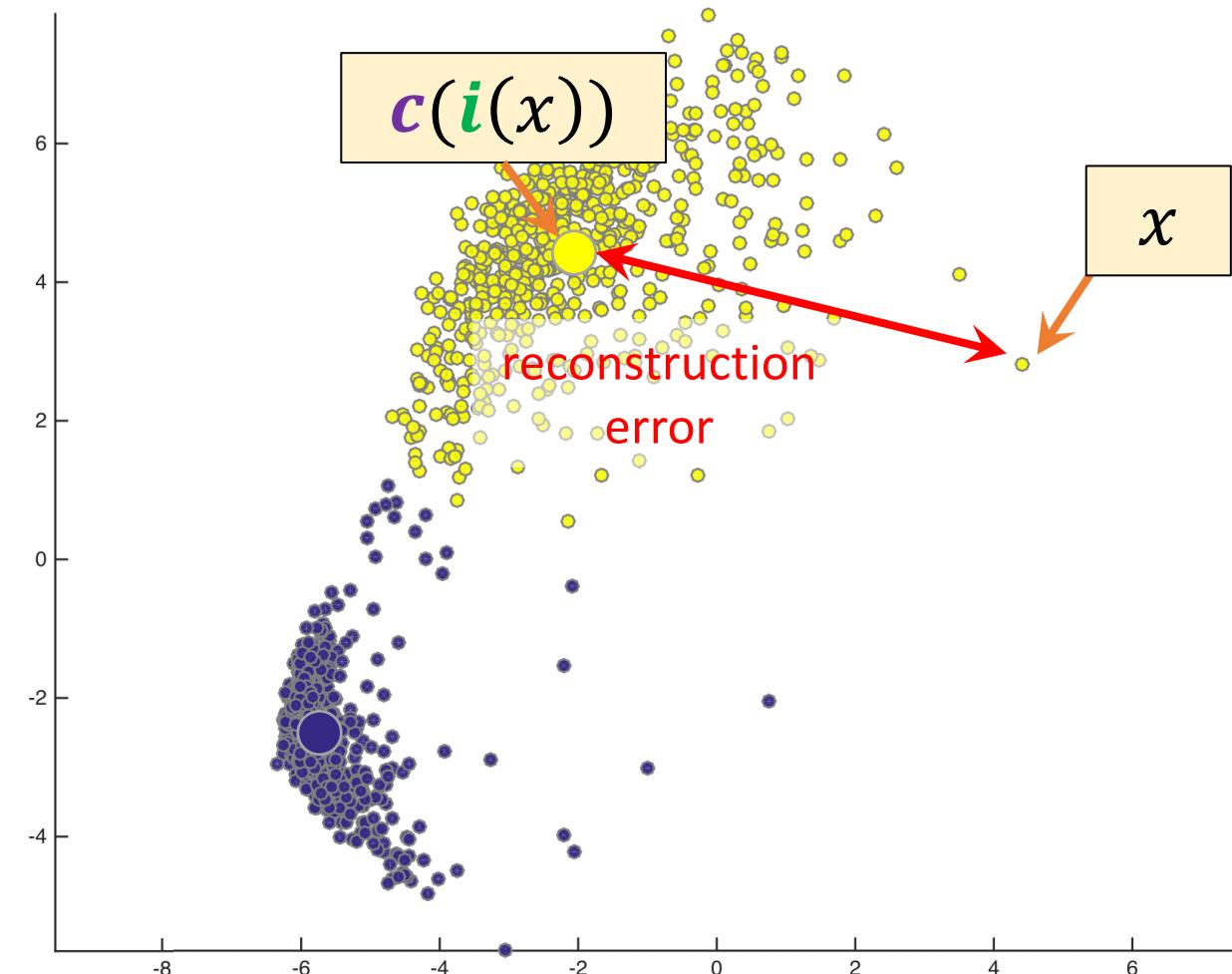
$$\min_{\mathbf{c}, \mathbf{i}} \sum_x \|x - \mathbf{c}(i(x))\|^2$$

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

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

Repeat until converge

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



K-means clustering

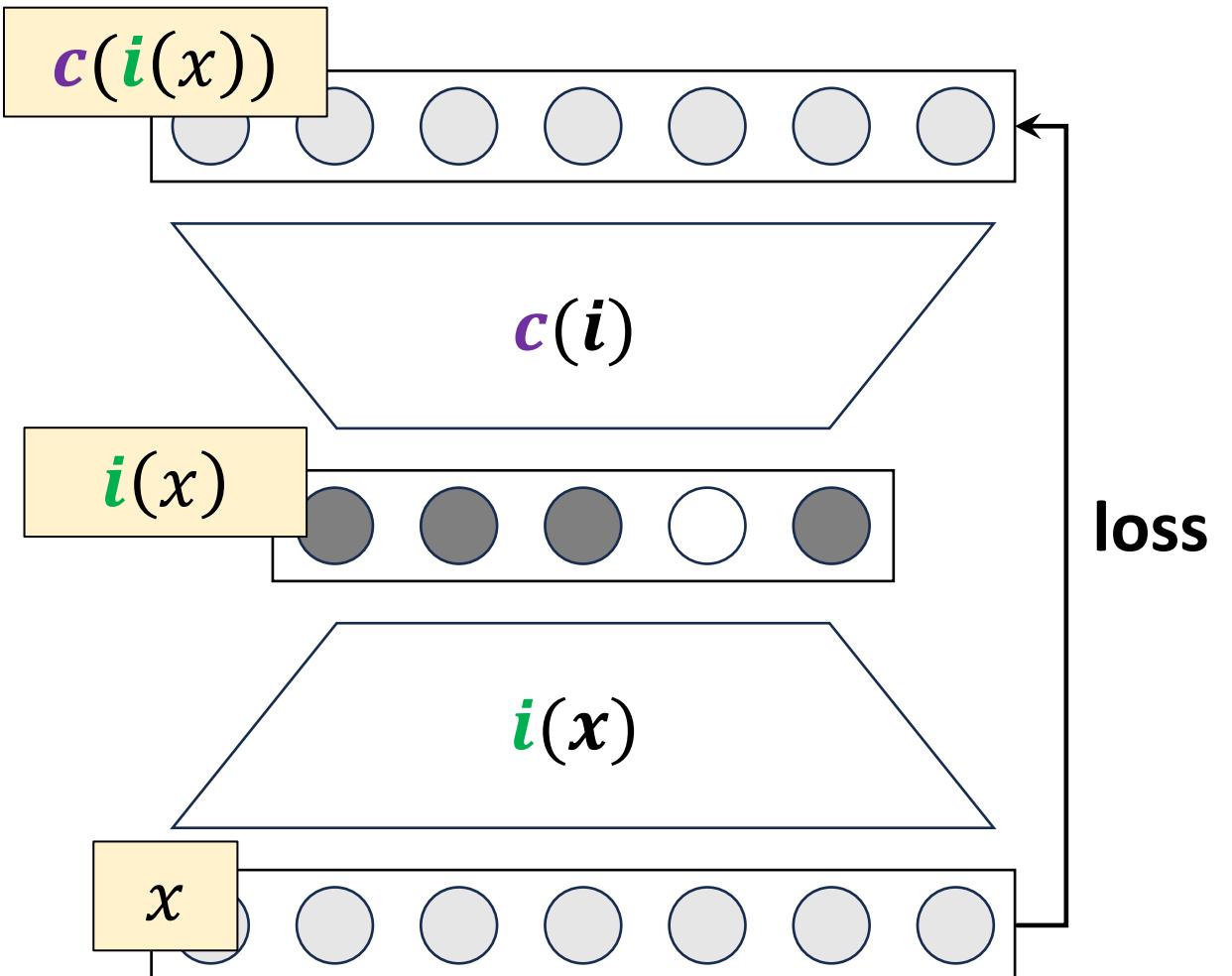
Loss function:

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

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

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

Repeat until converge



K-means clustering

Loss function:

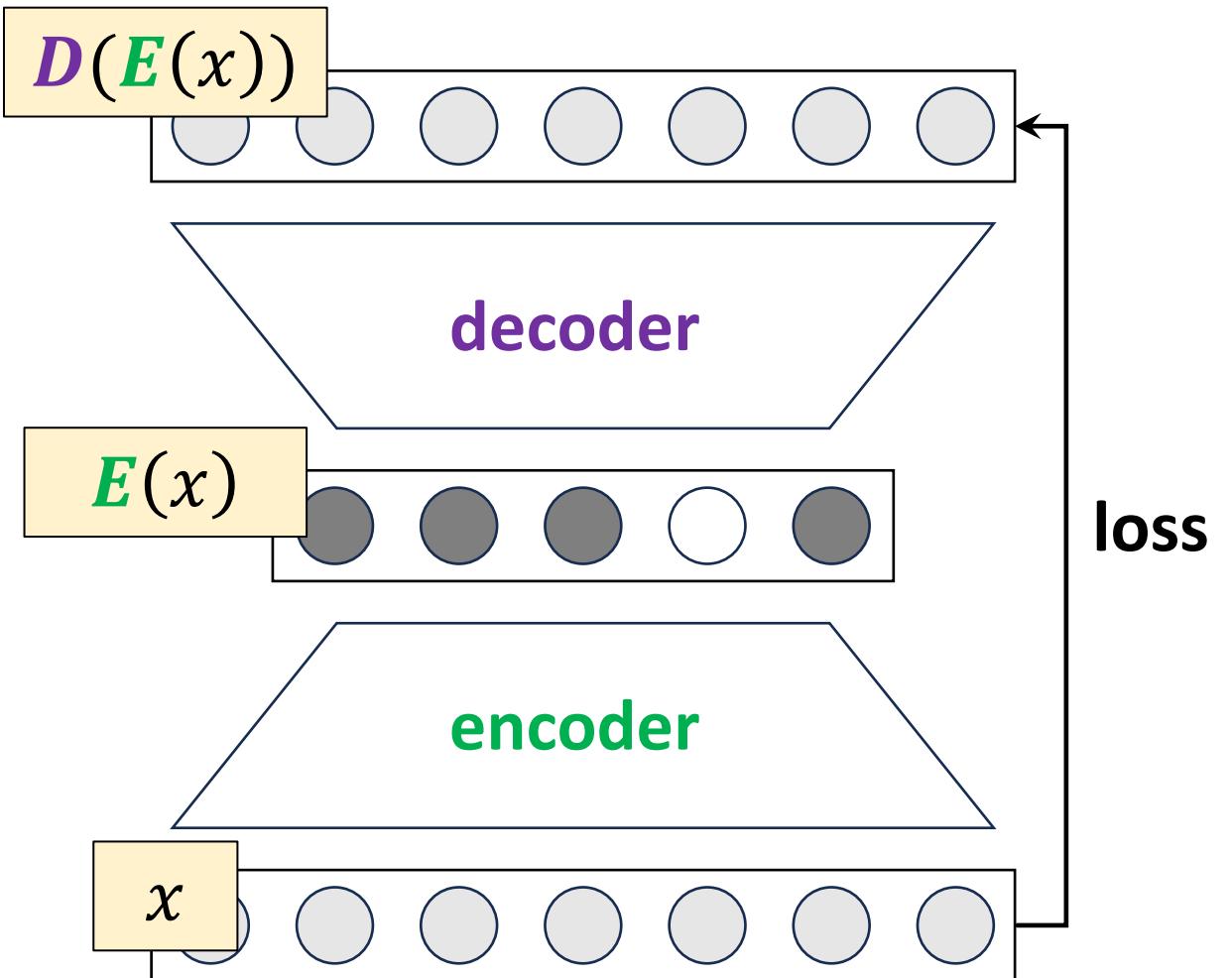
$$\min_{\mathbf{D}, \mathbf{E}} \sum_x \|x - \mathbf{D}(\mathbf{E}(x))\|^2$$

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

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

K-means performs Autoencoding

- w/ one-hot encoding
- w/ coupled encoder & decoder

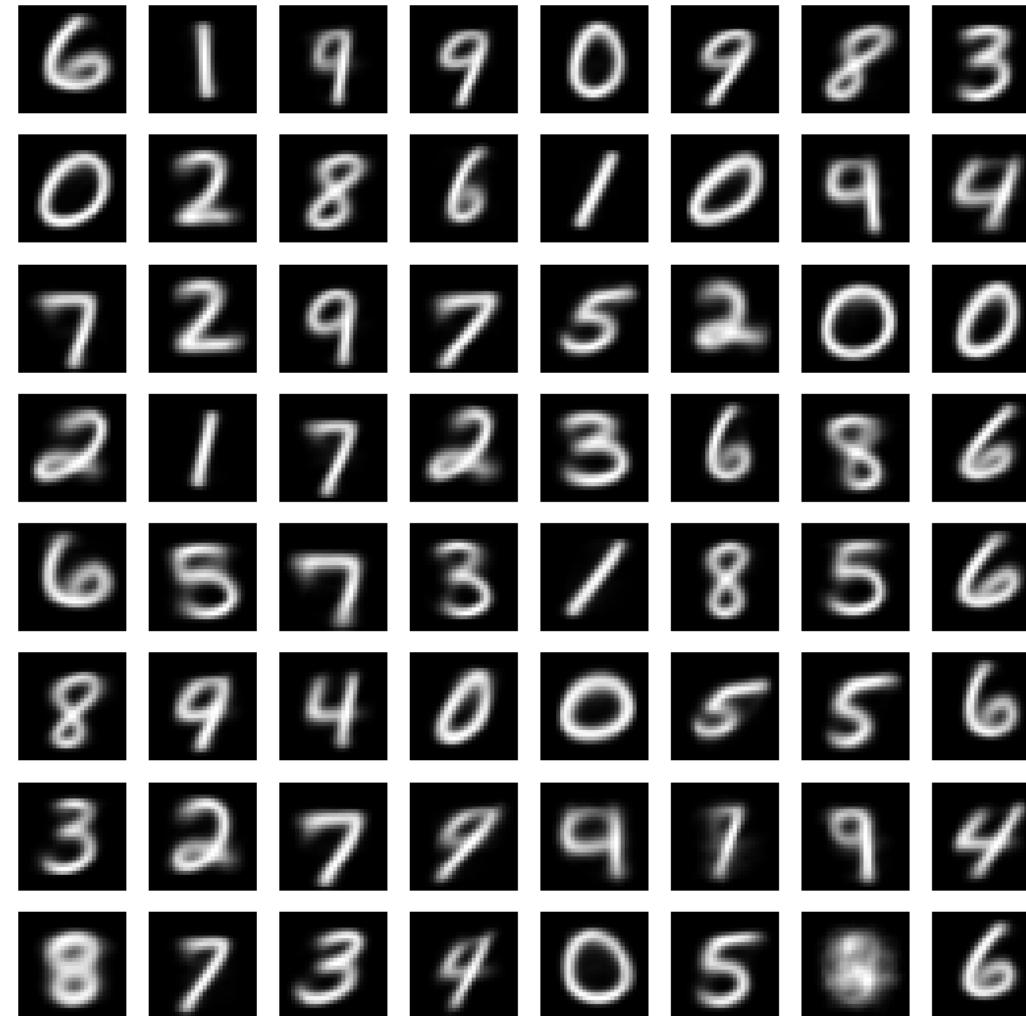


Example: K-means on MNIST digits

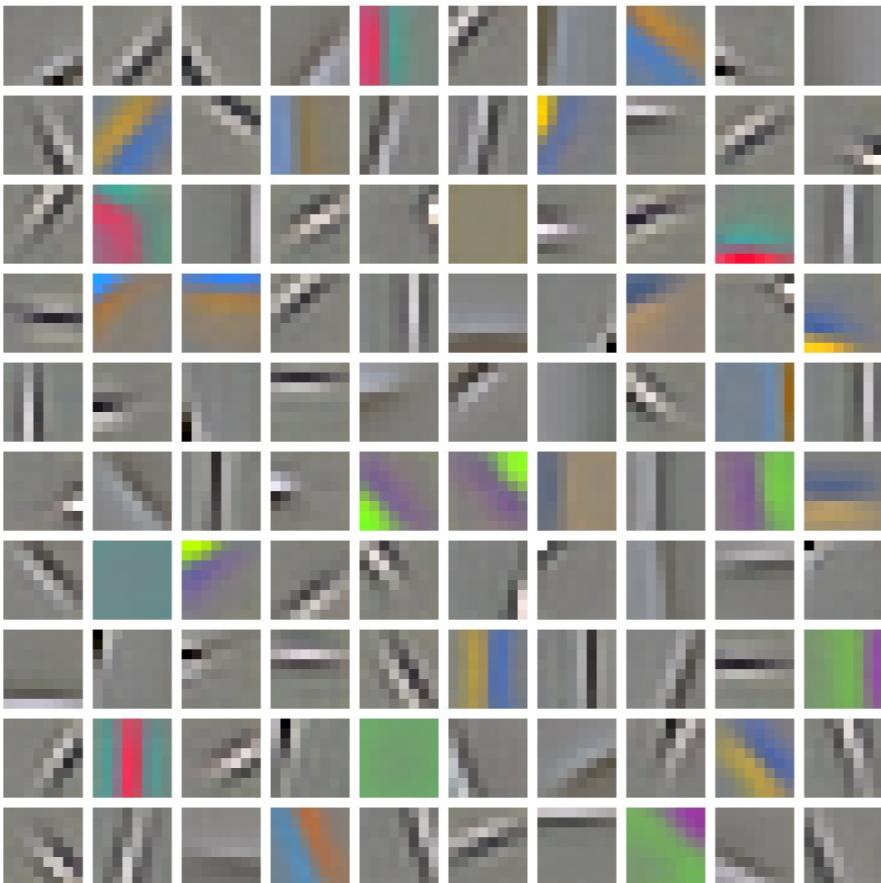
Visualized centers with $k = 64$

These are not real images.

These are the learned representations.



Example: K-means on 8×8 image patches



K-means
(w/ PCA whitening)



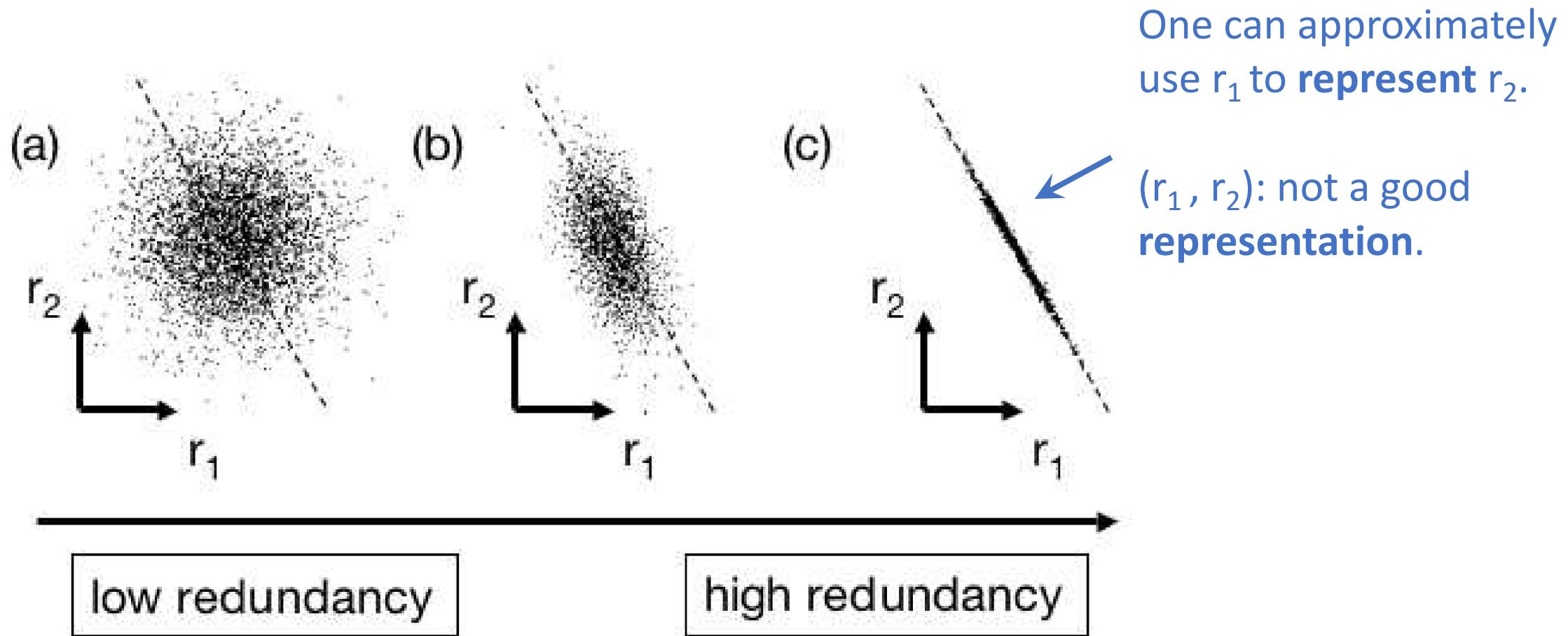
1st layer of AlexNet
(supervised learning)

Reconstruction: classical yet highly powerful

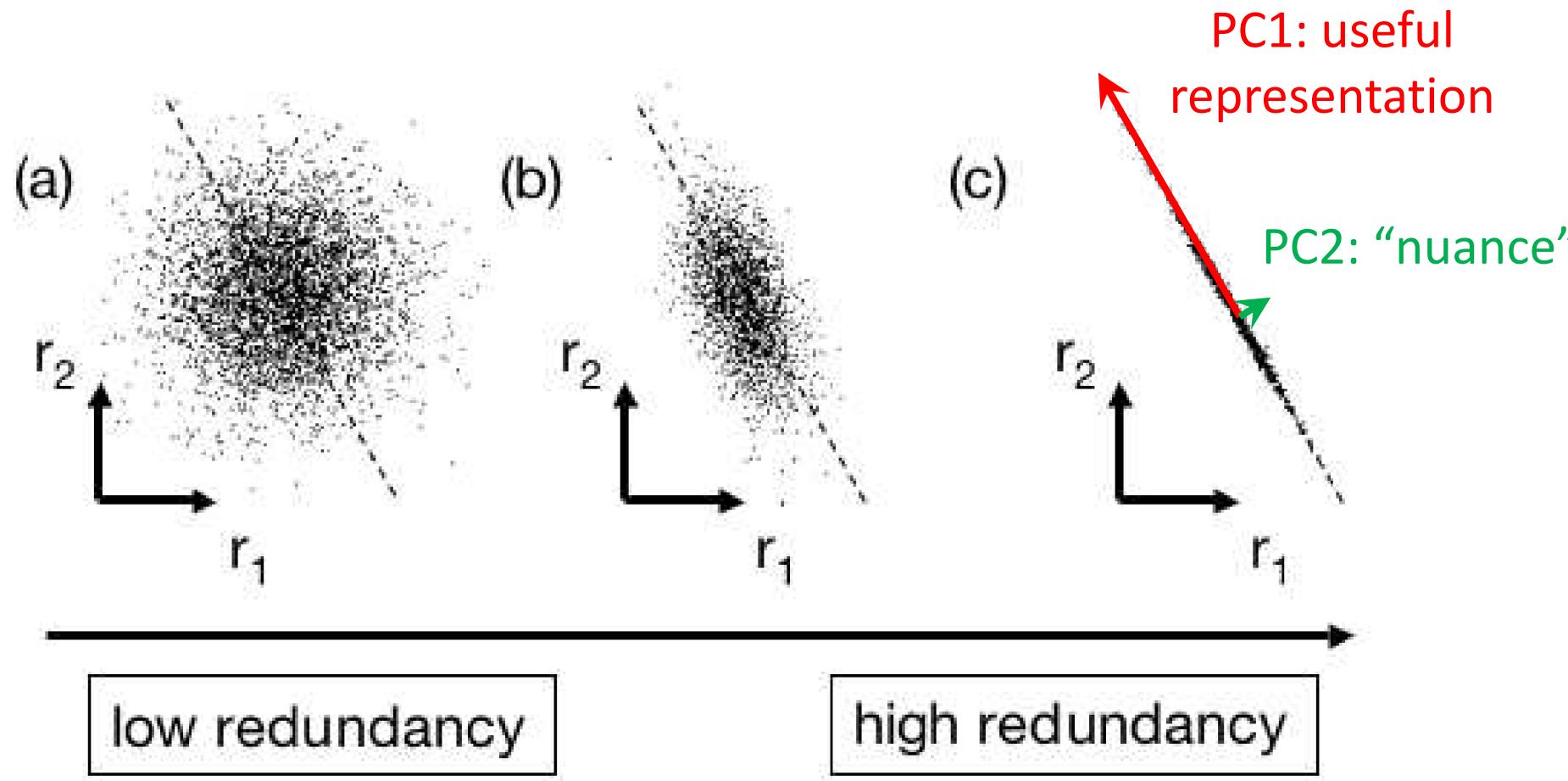
- K-means clustering
- Principal Component Analysis (PCA)
- Autoencoder



Principal Component Analysis (PCA)



Principal Component Analysis (PCA)



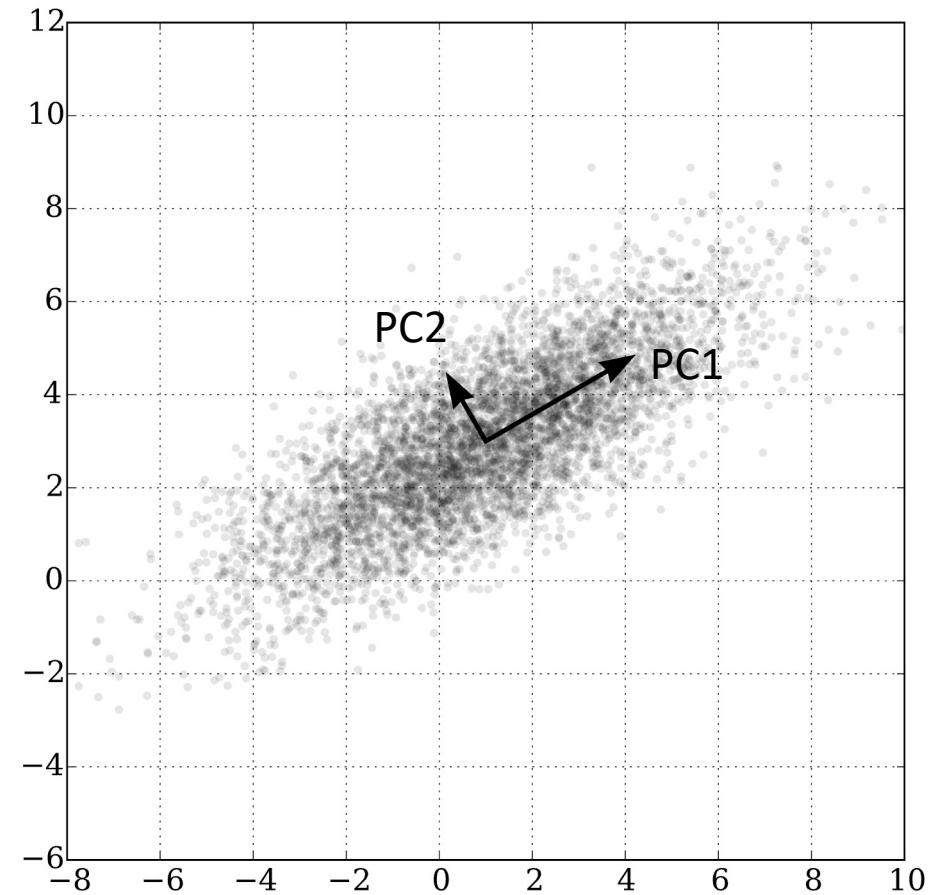
Principal Component Analysis (PCA)

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)

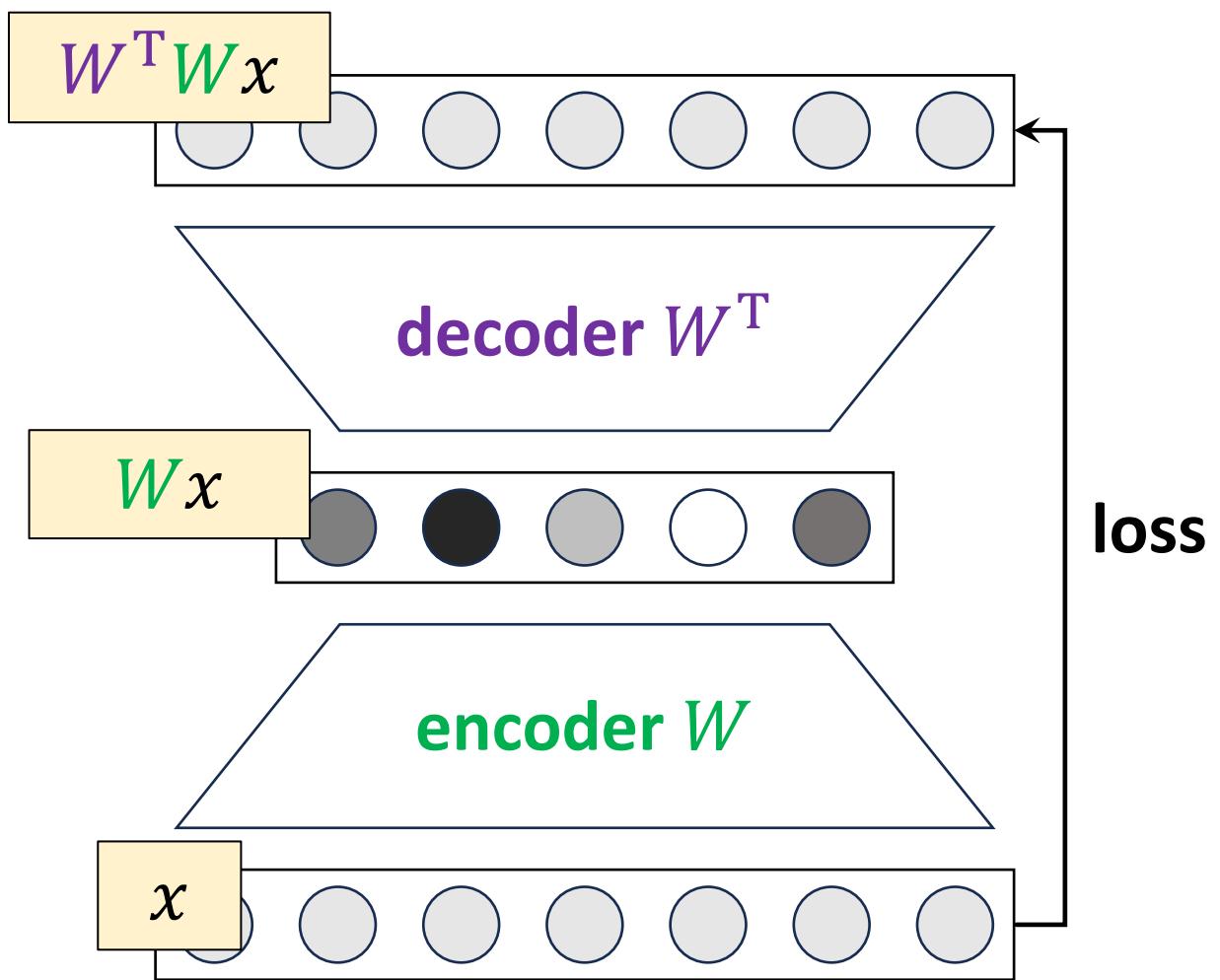
Represent data in the space of the largest k Principal Components (PCs)

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

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

PCA performs Autoencoding

- w/ linear encoder & decoder
- w/ orthogonal constraint



Example: PCA on face data (Eigenfaces)

These are not real images.

These are the learned representations.



Reconstruction: classical yet highly powerful

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

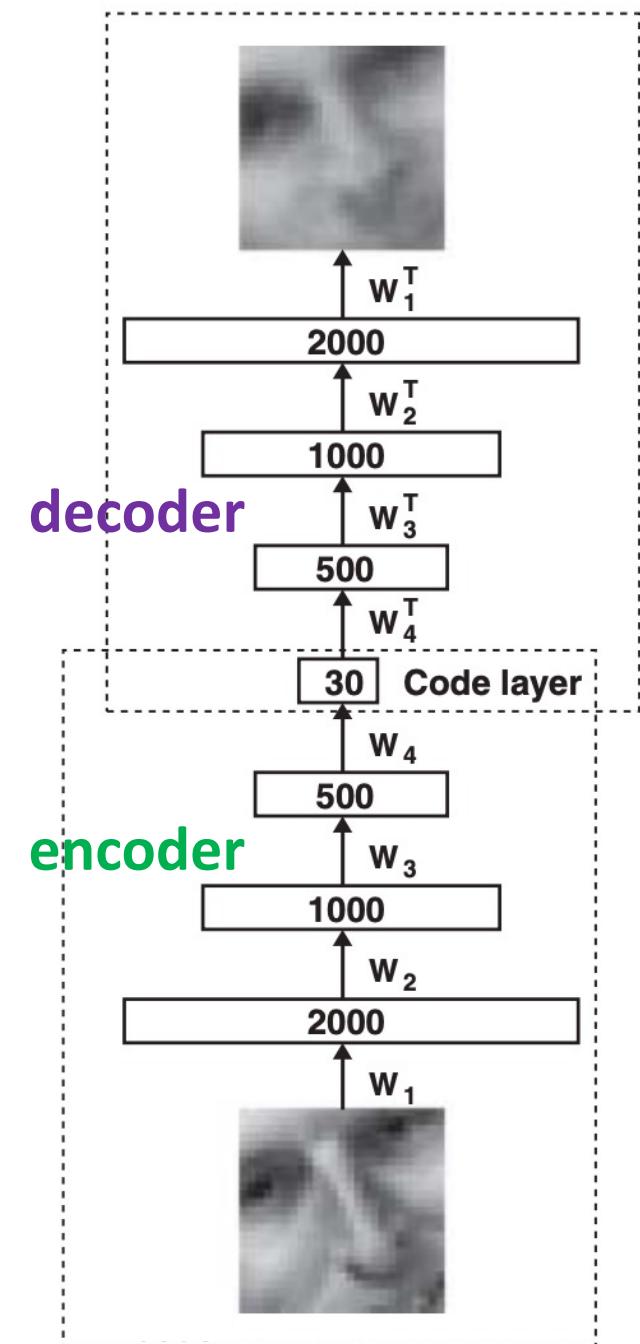


Autoencoder [2006]

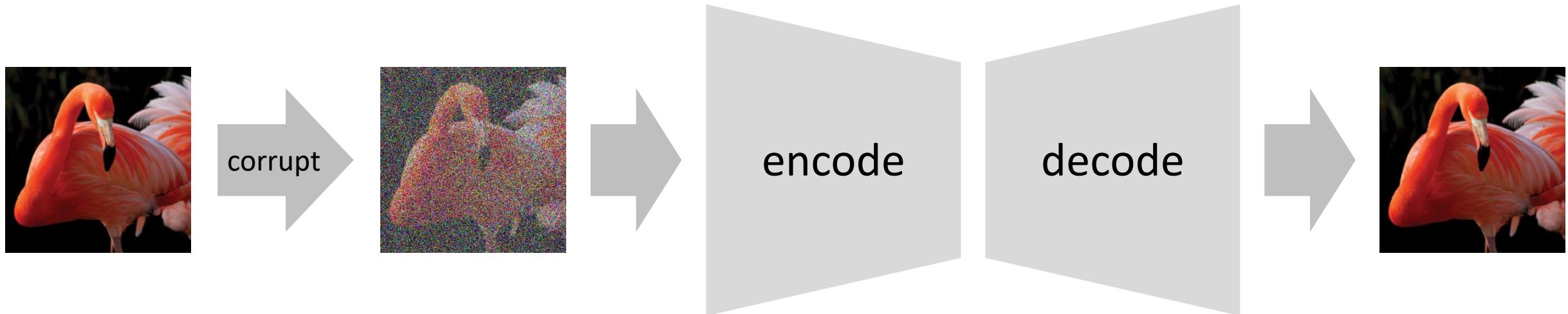
Represent data by compressing it into a low-dimensional latent space

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

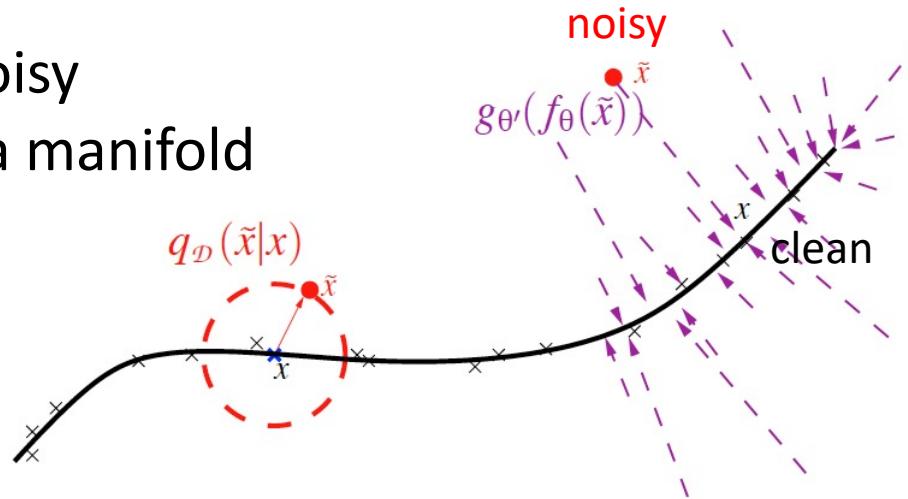
$$\min_{\mathbf{D}, \mathbf{E}} \sum_x \|x - \mathbf{D}(\mathbf{E}(x))\|^2$$



Denoising Autoencoder (DAE) [2010]



learning by projecting noisy
samples back to the data manifold



Masked Image Modeling [2016]

Context Encoders

- Autoencoding with masking
- in pixel space
- using ConvNet



Masked Language Modeling [2018]

BERT (Bidirectional Encoder Representations from Transformers)

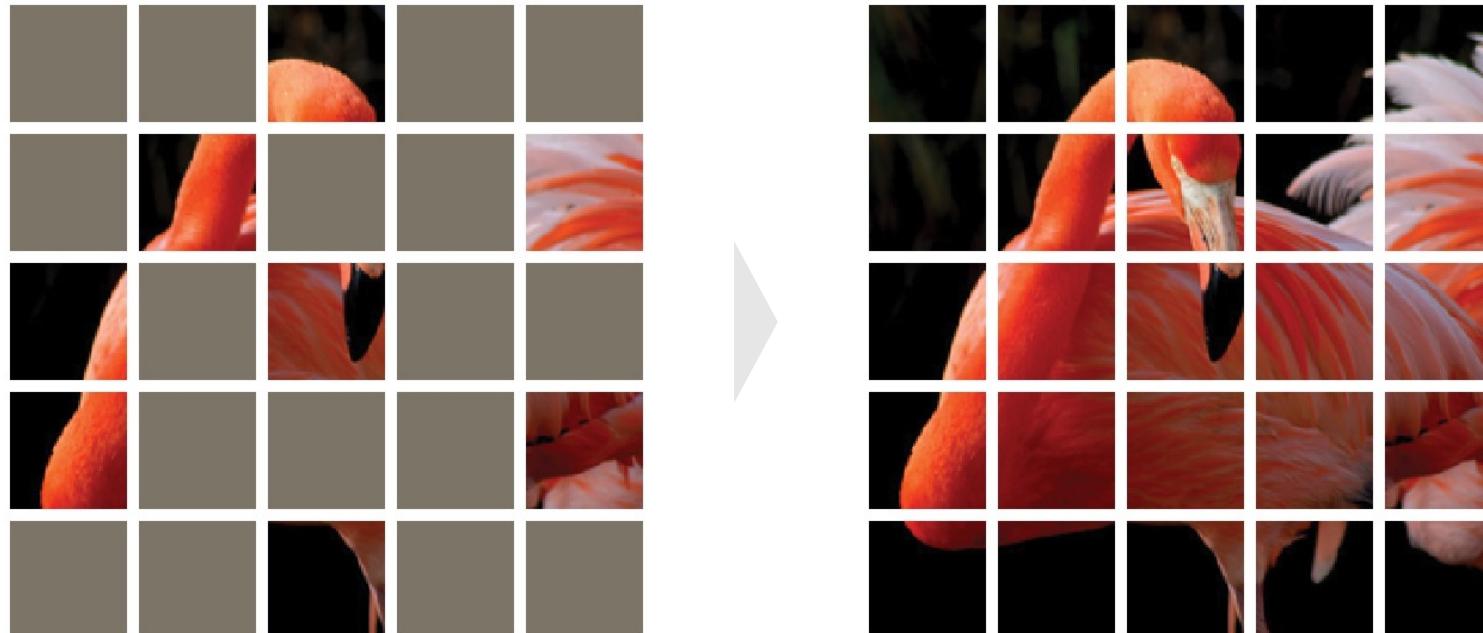
- Autoencoding with masking
- in token space
- using Transformers

The ___ opened their ___ and began to ___ → **net** → students .. books .. read

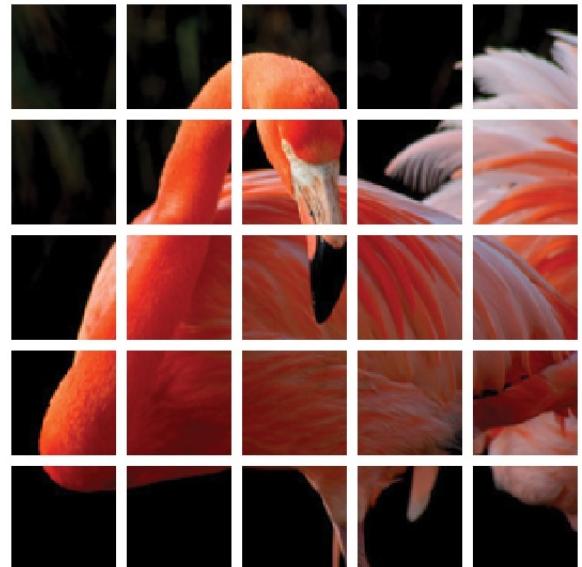
Masked Image Modeling [2022]

MAE (Masked Autoencoder)

- Autoencoding with masking
- in token/patch space
- using Transformers

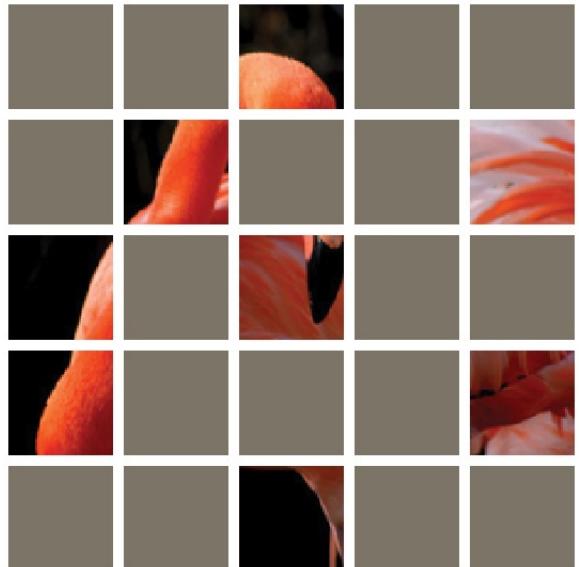


Masked Autoencoder (MAE)



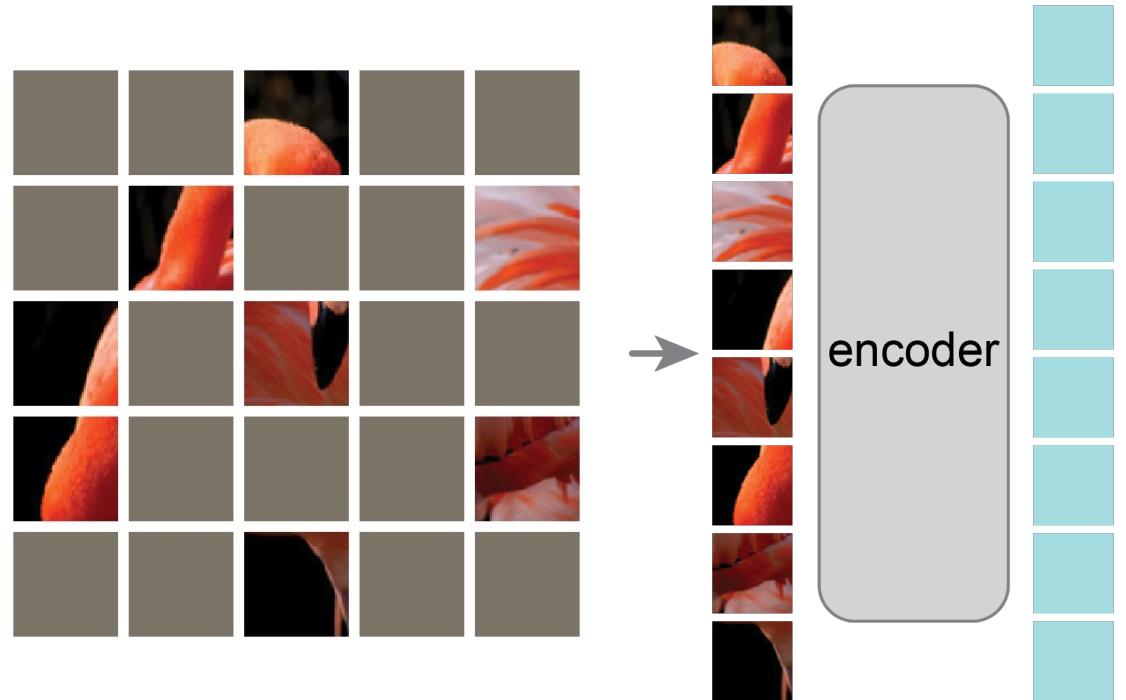
patches as visual tokens
(Vision Transformer)

Masked Autoencoder (MAE)



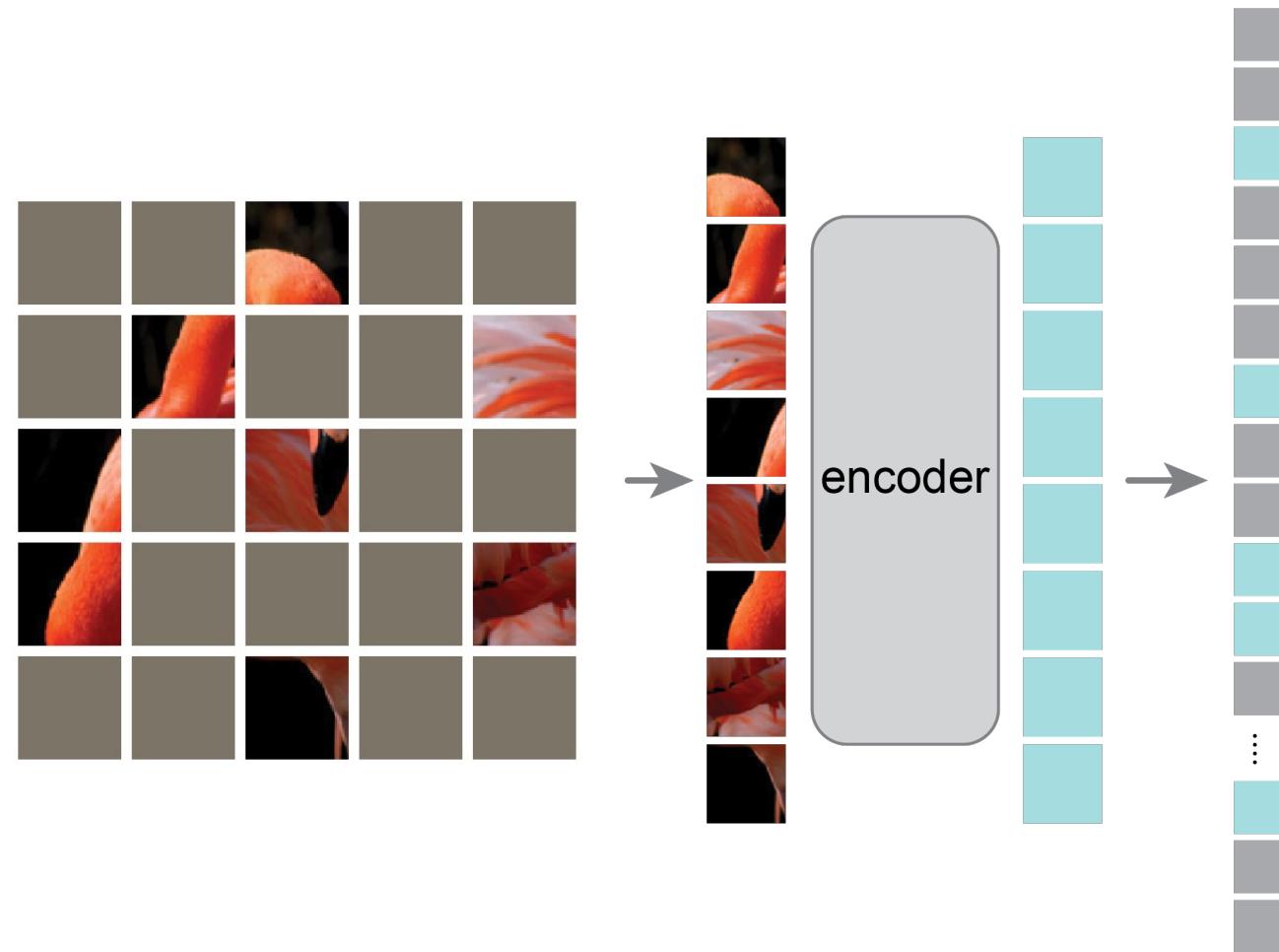
random masking

Masked Autoencoder (MAE)

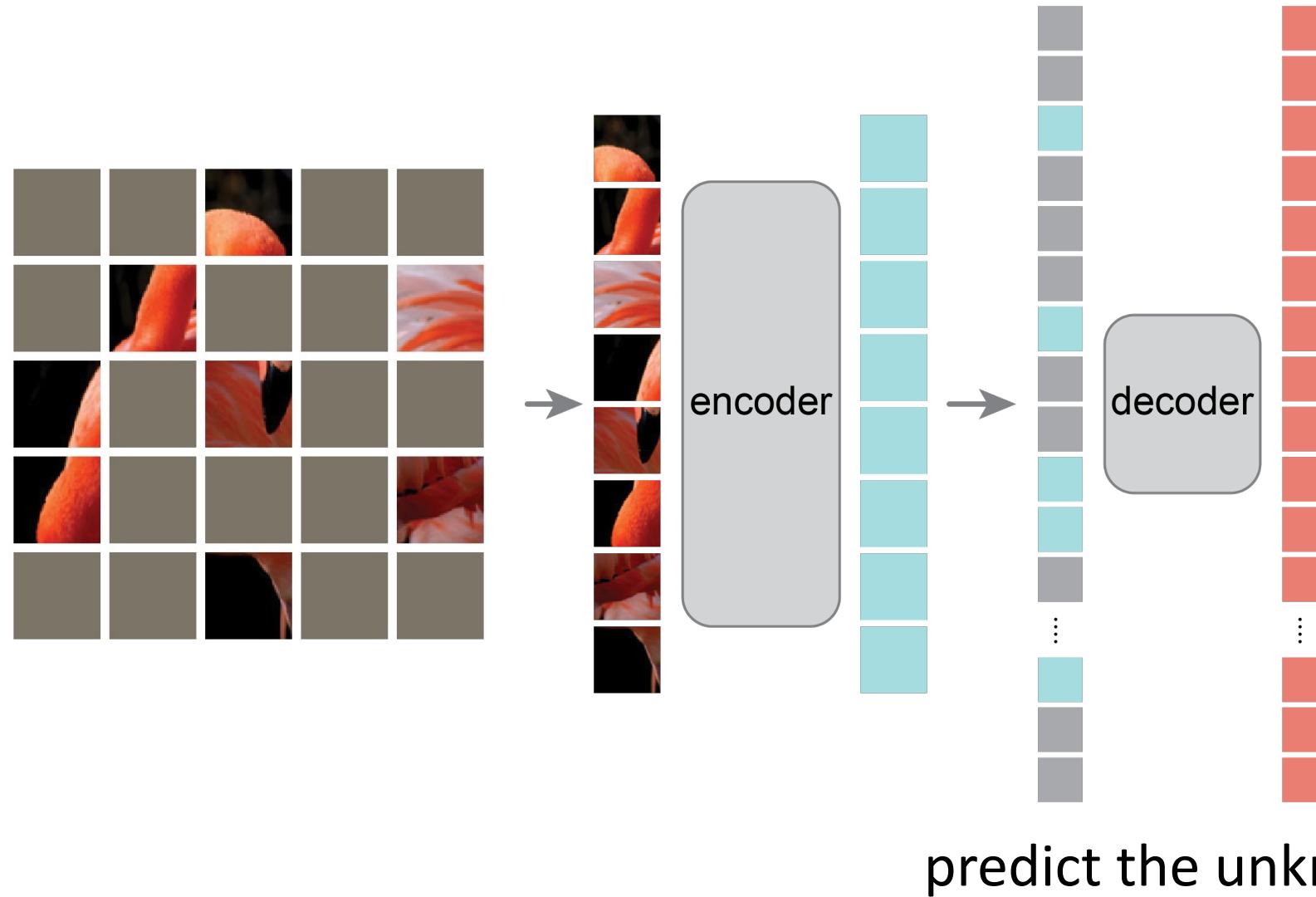


encode visible patches
w/ Transformer

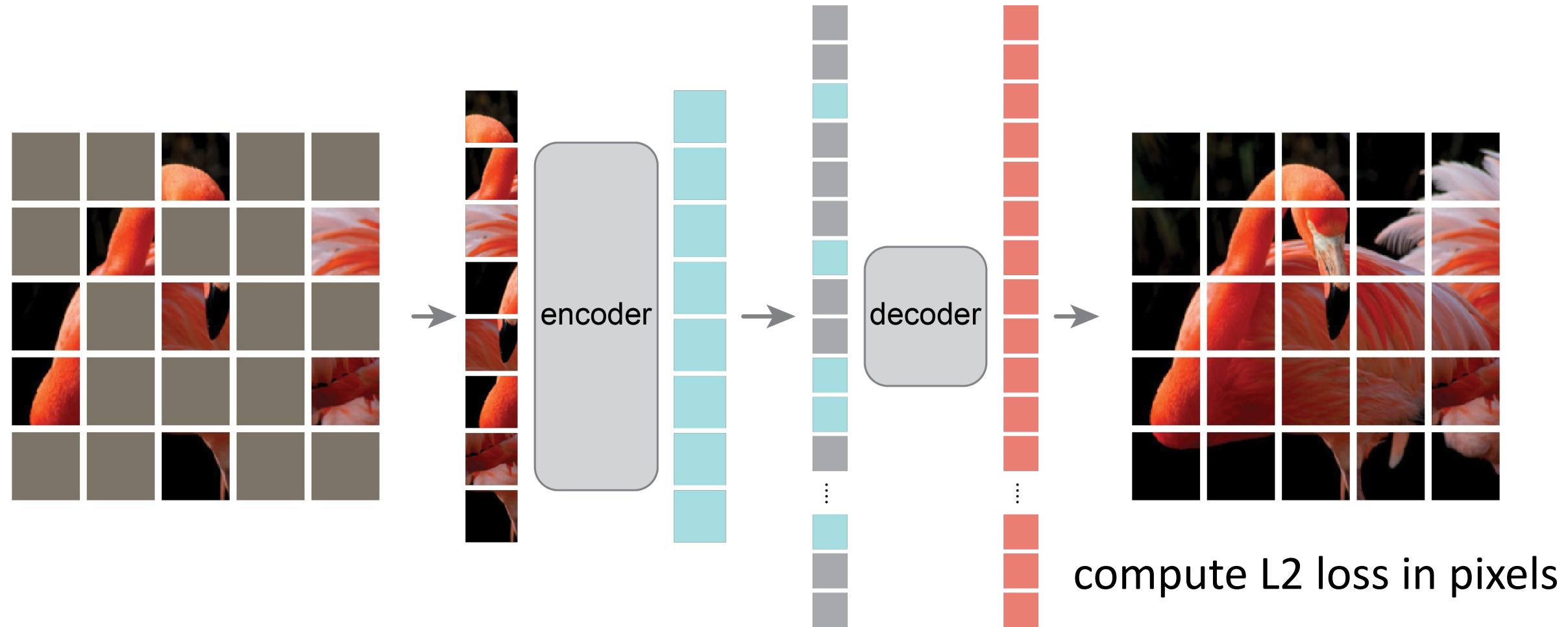
Masked Autoencoder (MAE)



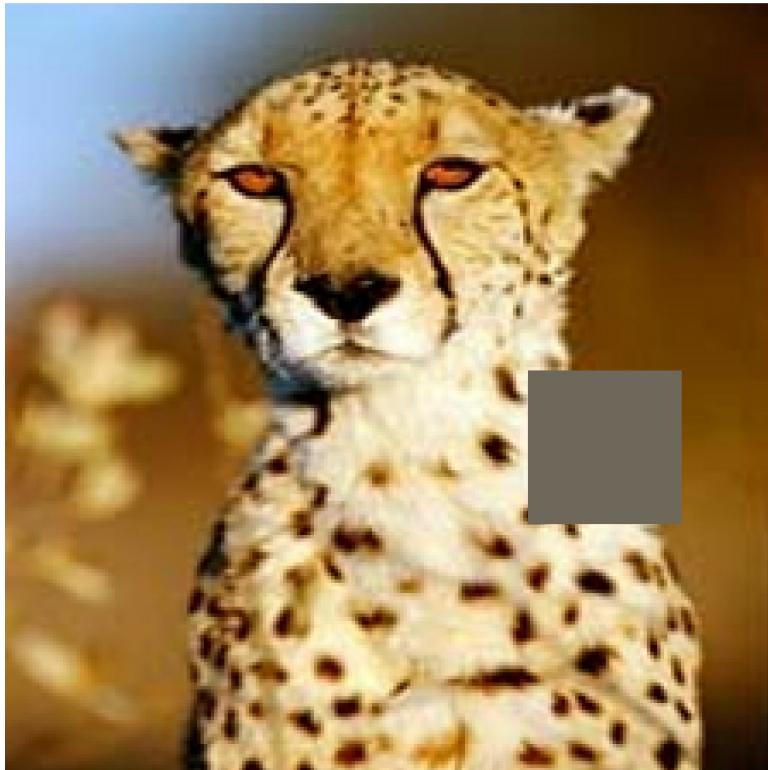
Masked Autoencoder (MAE)



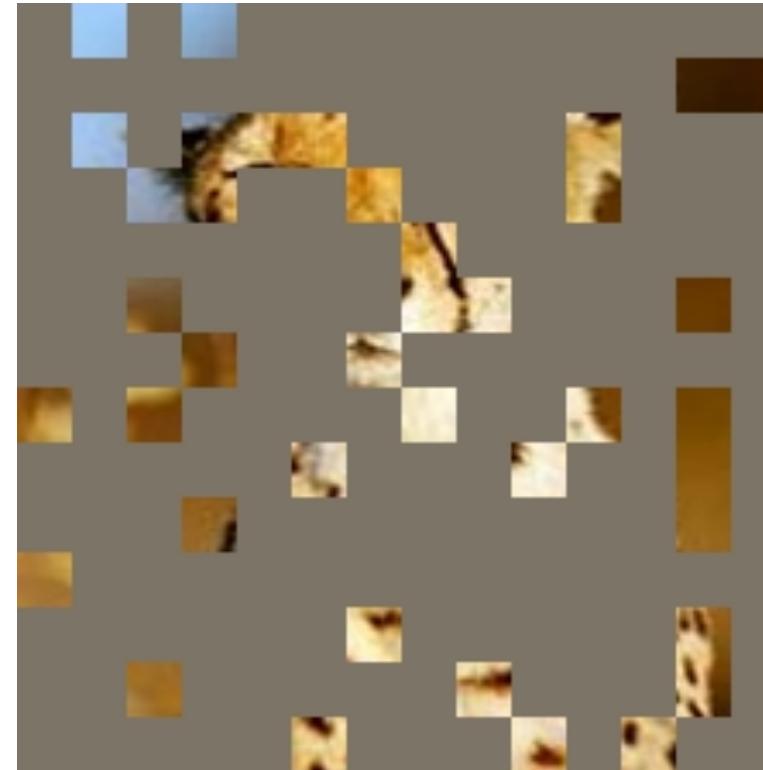
Masked Autoencoder (MAE)



How to encourage learning good representations?

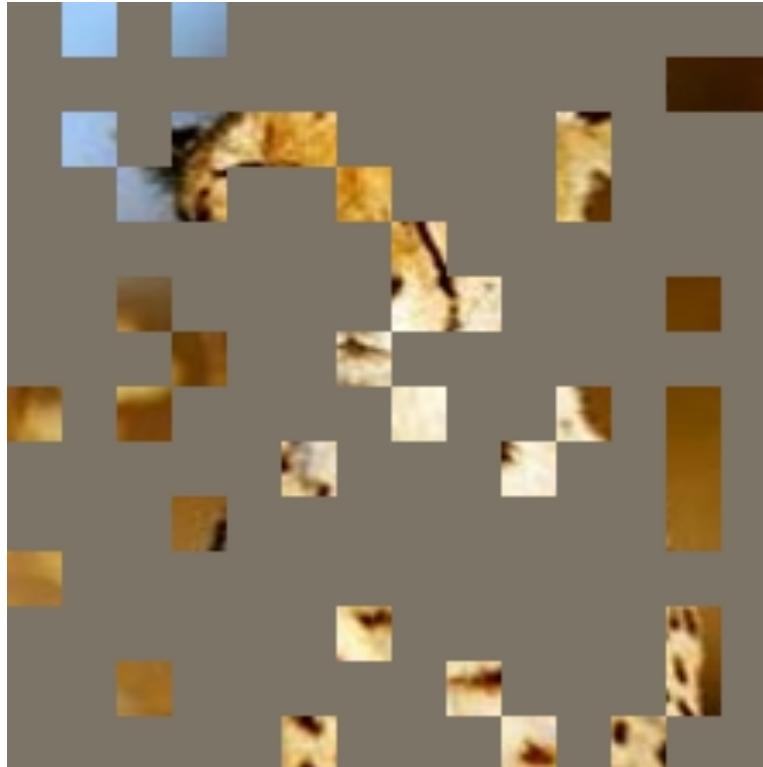


- predicting a small portion may not require high-level understanding

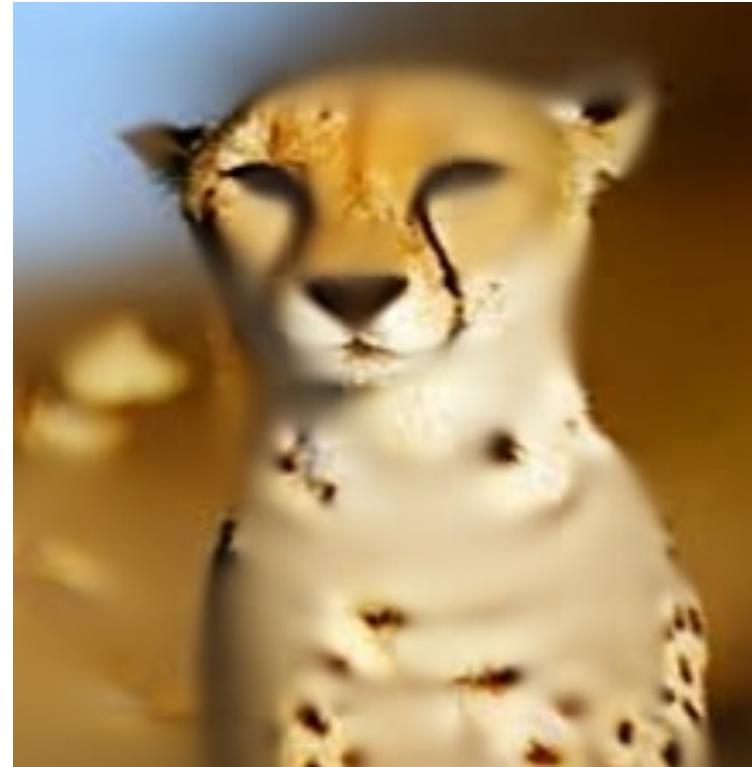


- predicting a large portion of unknown patches encourages to learn semantic features

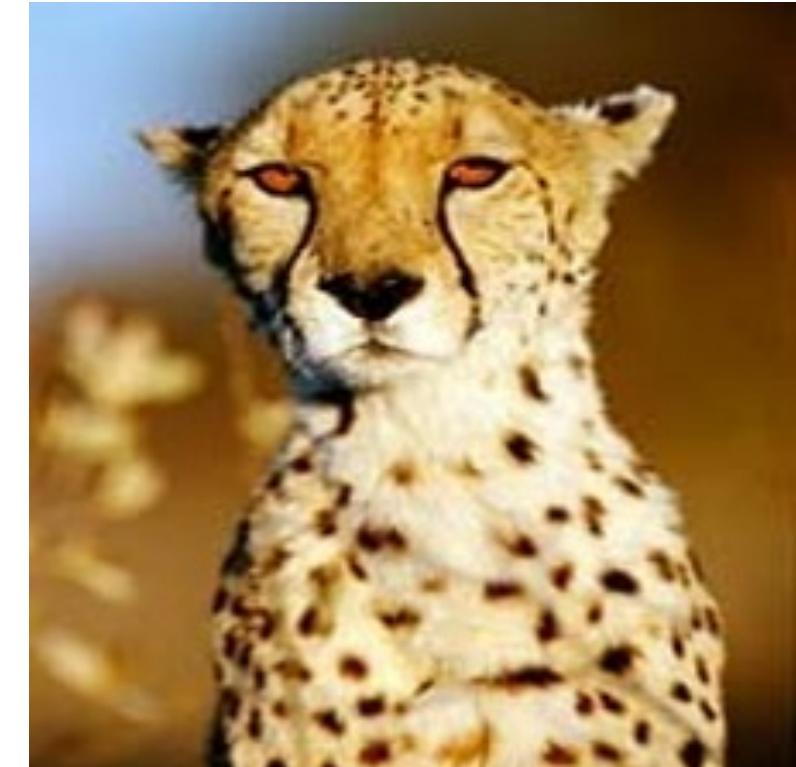
How to encourage learning good representations?



input

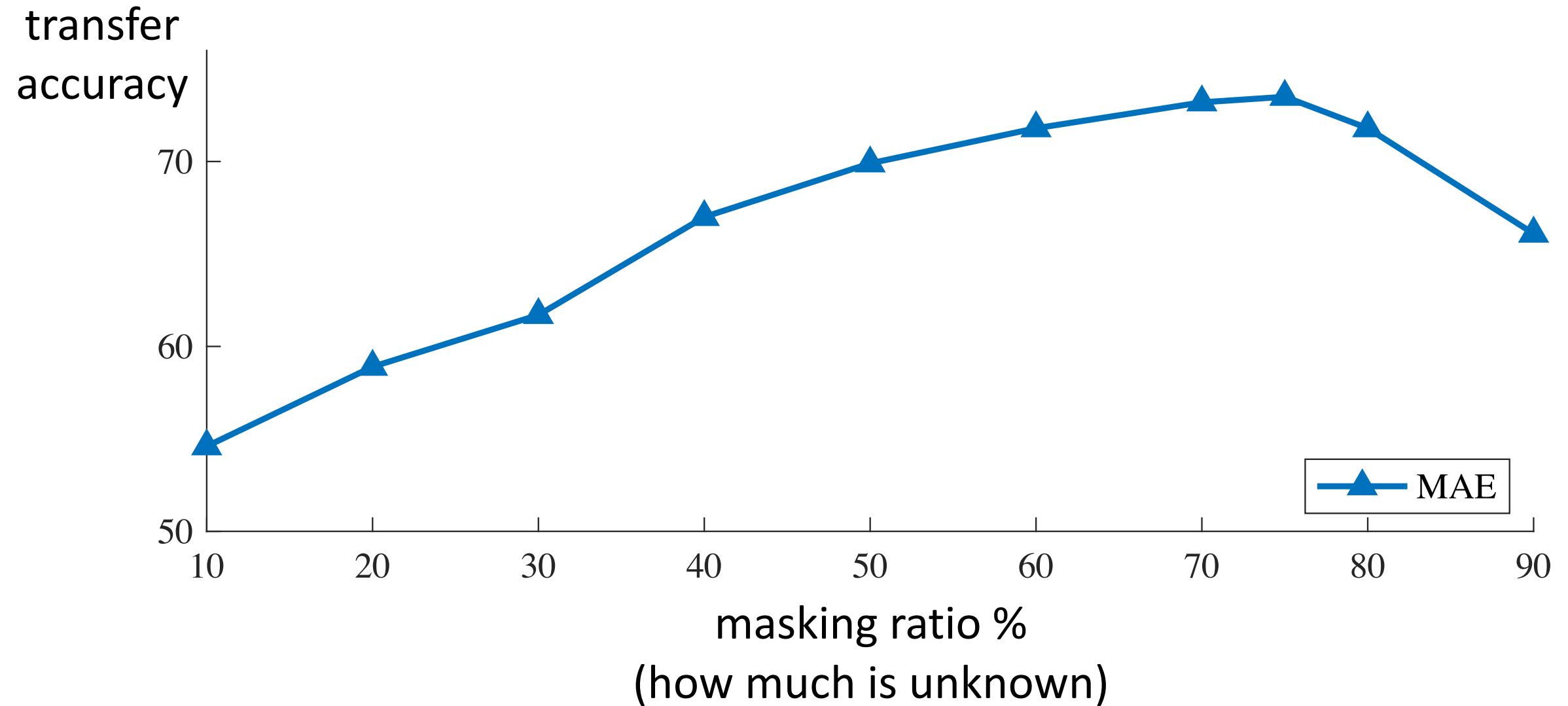


MAE prediction

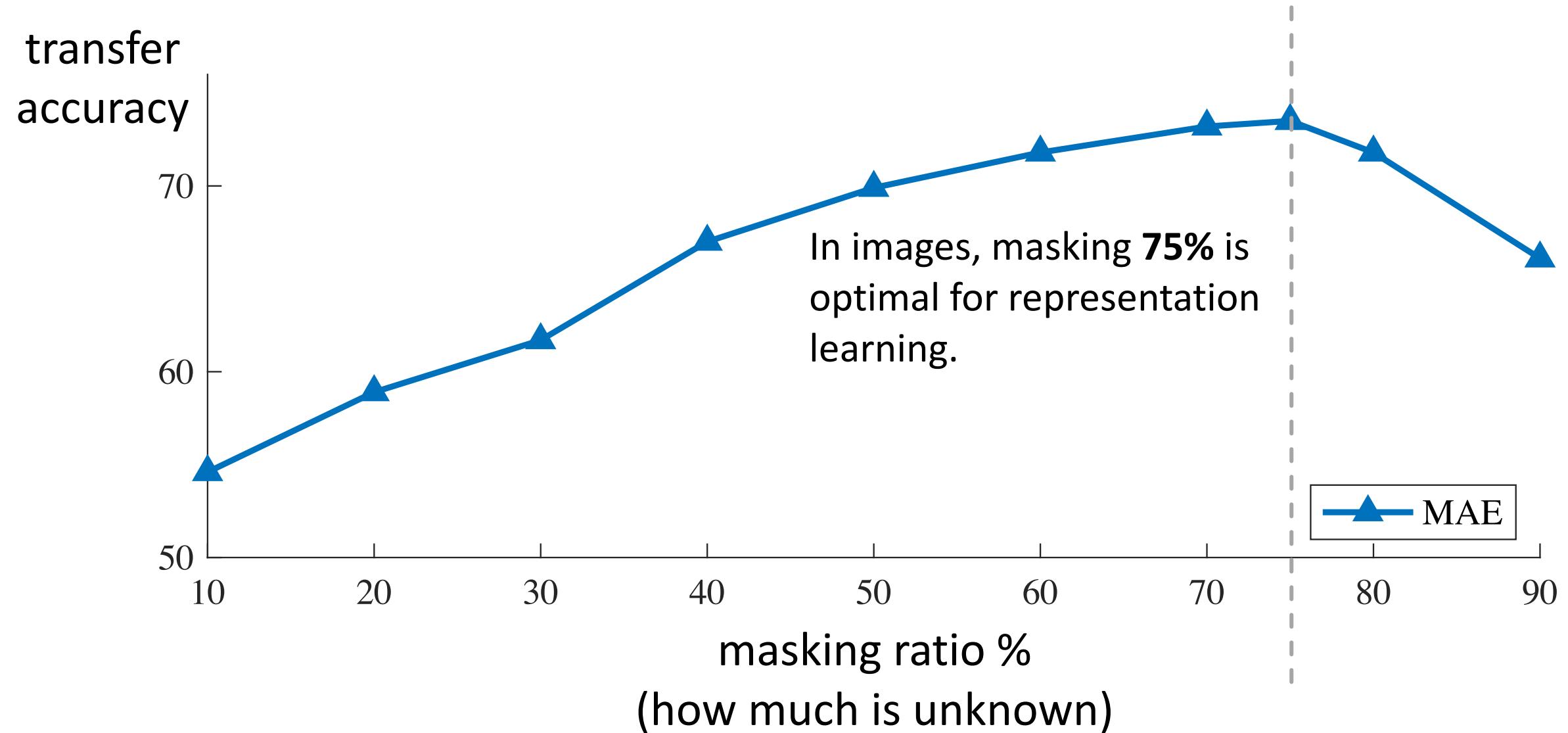


original

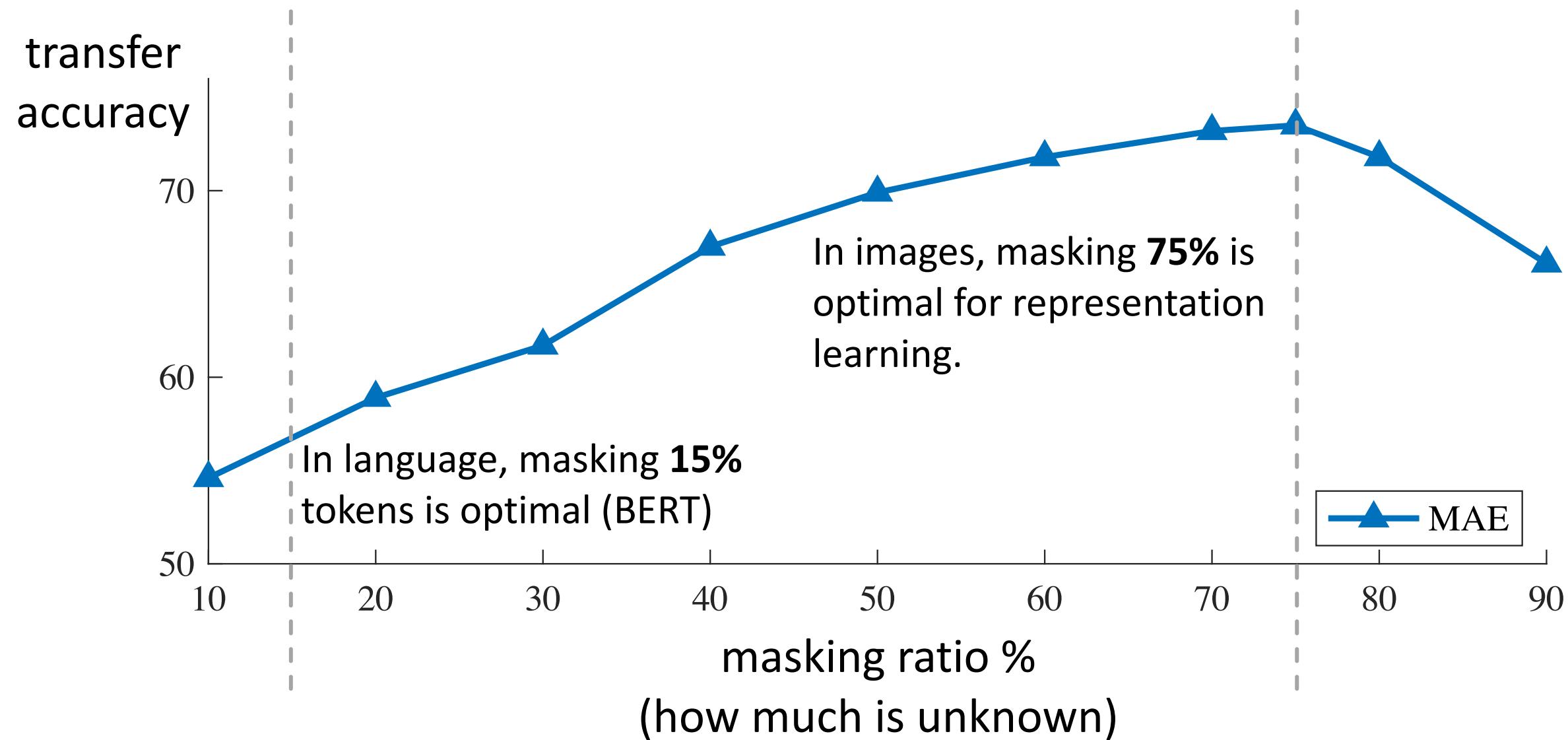
How to encourage learning good representations?



How to encourage learning good representations?



How to encourage learning good representations?



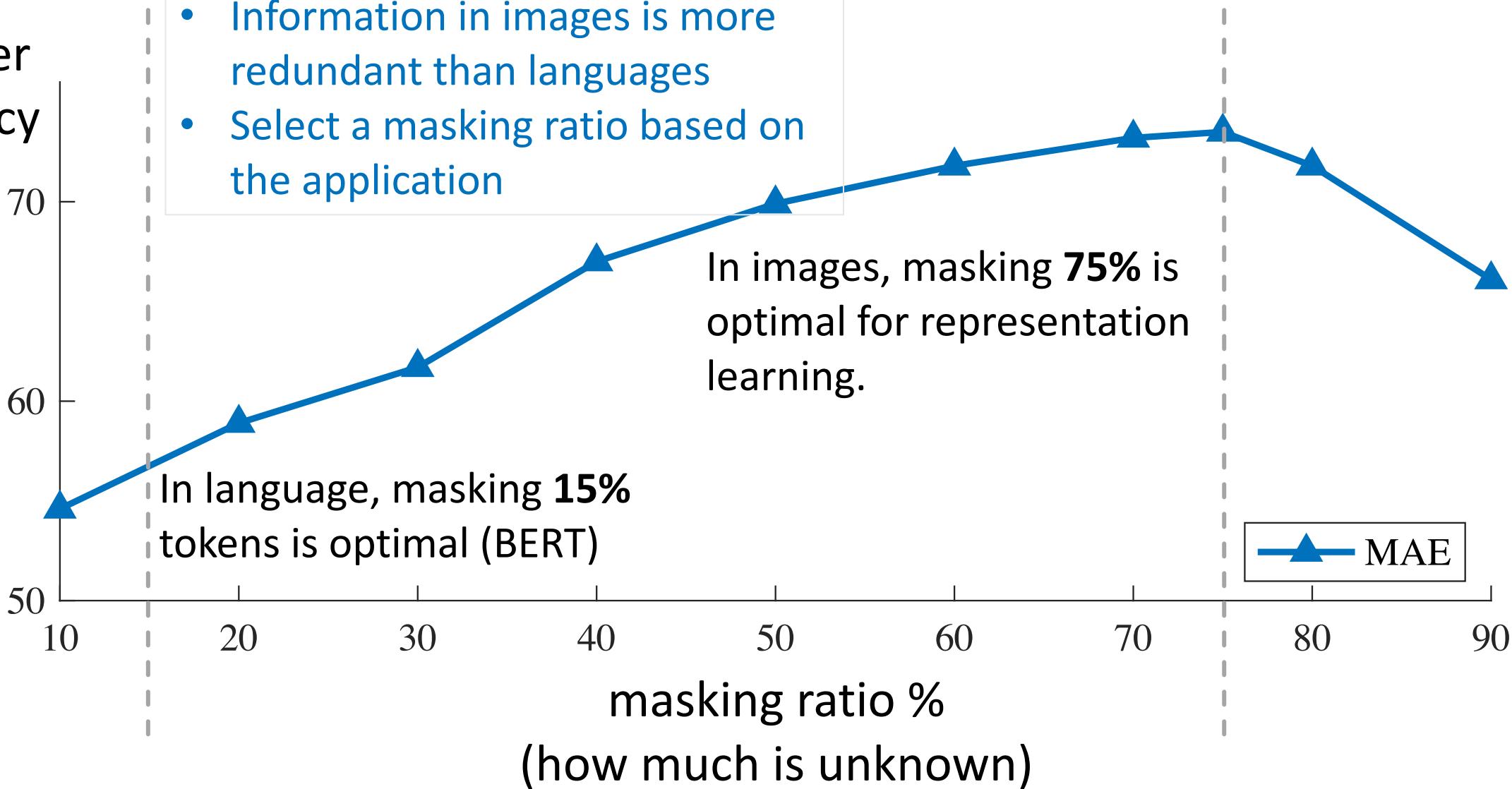
How to encourage learning good representations?

transfer accuracy

- Information in images is more redundant than languages
- Select a masking ratio based on the application

In images, masking **75%** is optimal for representation learning.

In language, masking **15%** tokens is optimal (BERT)

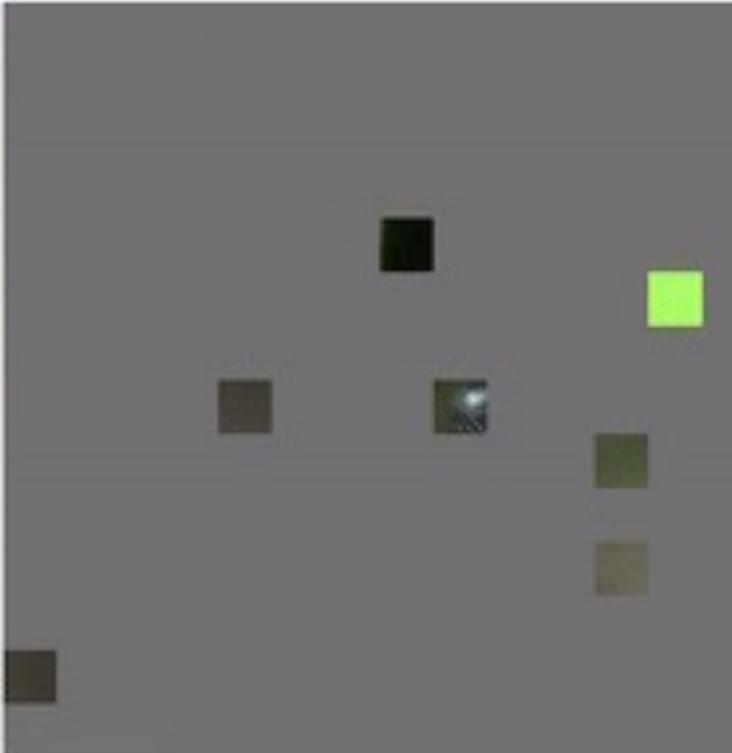


MAE on Videos

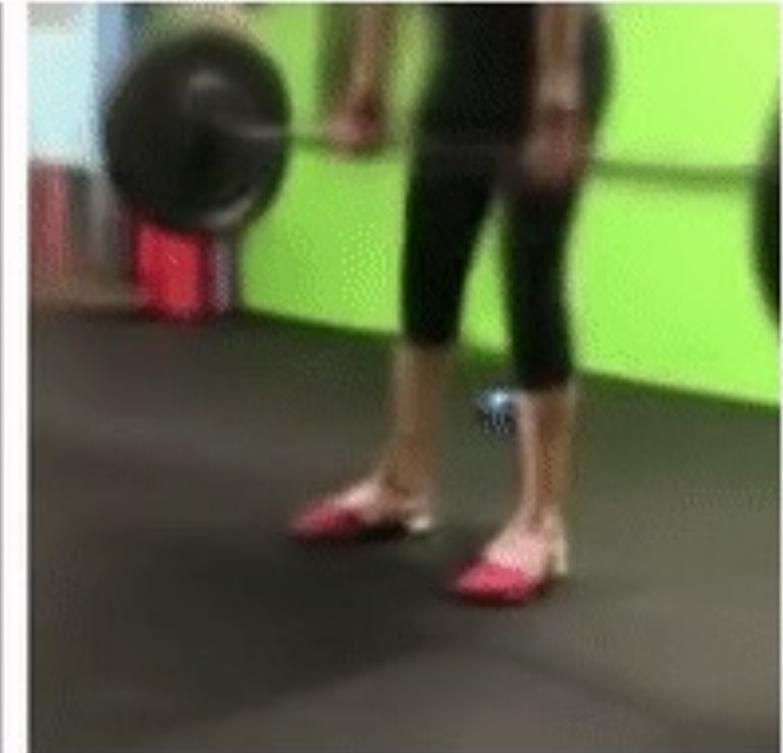
- 95% masking ratio
- videos are more redundant than images



original

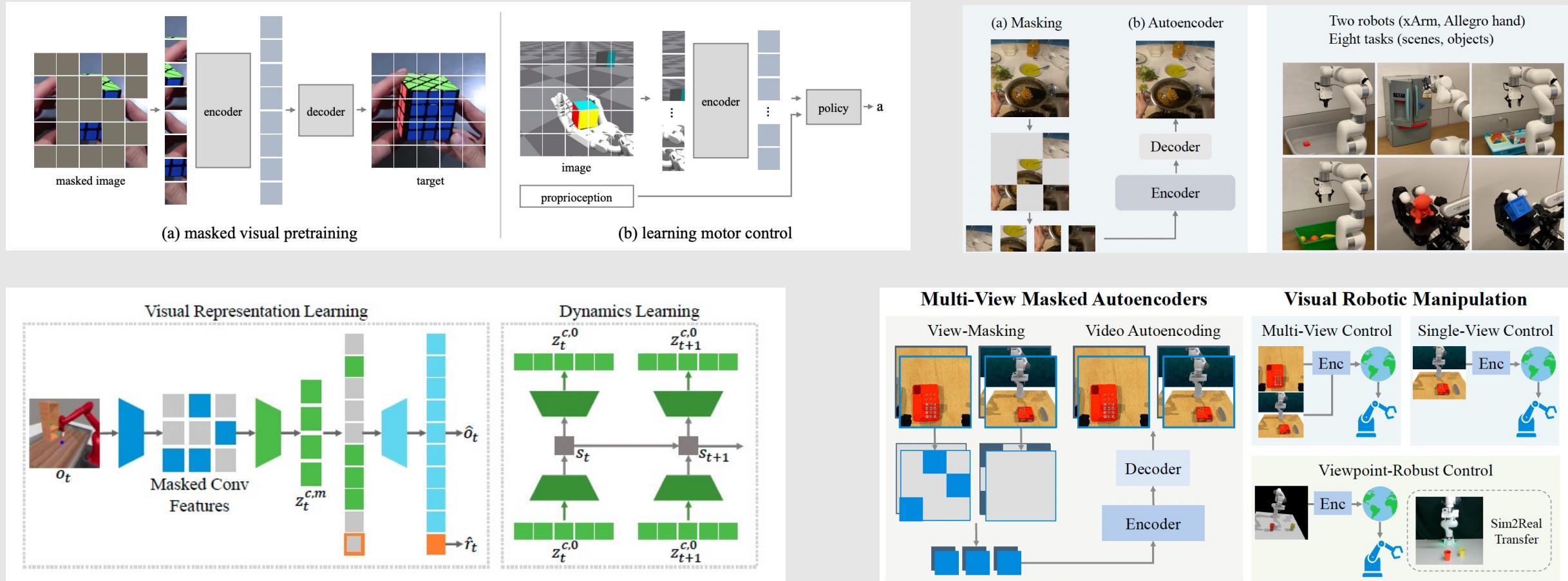


input



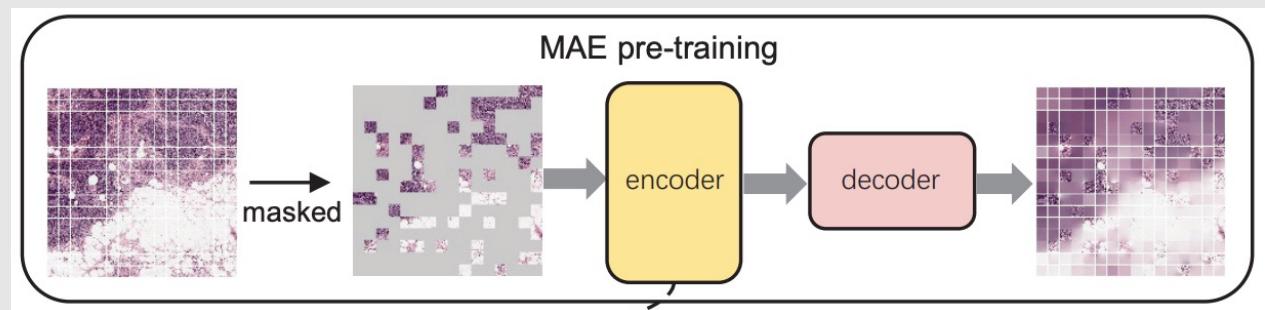
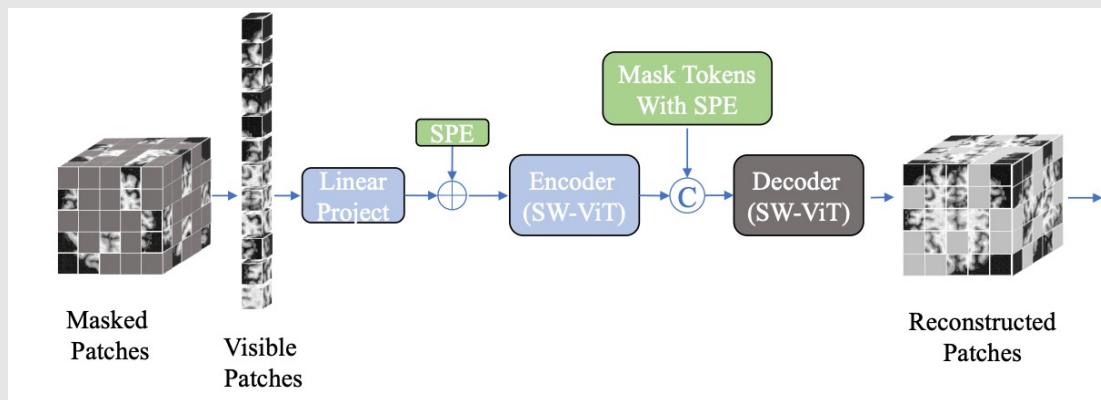
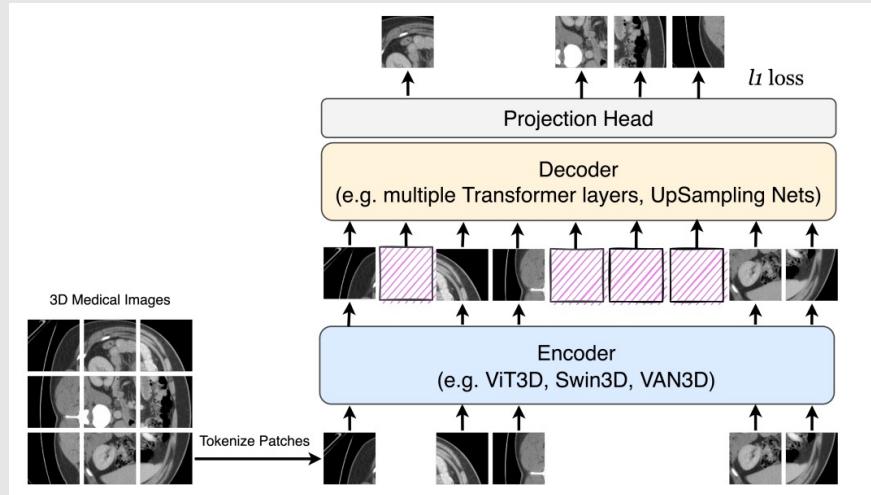
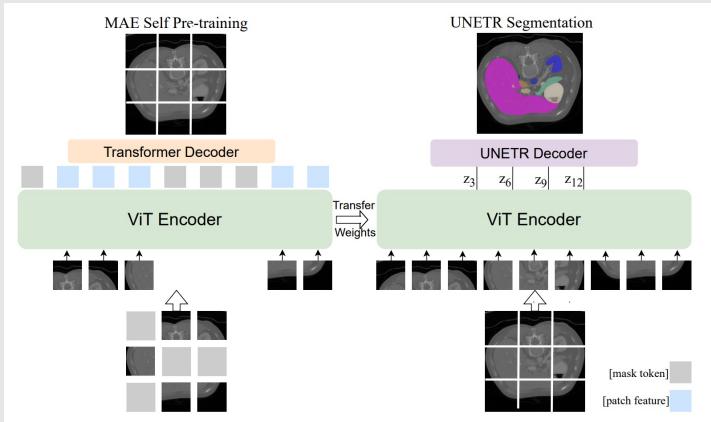
MAE prediction

Masked Autoencoder: Applications



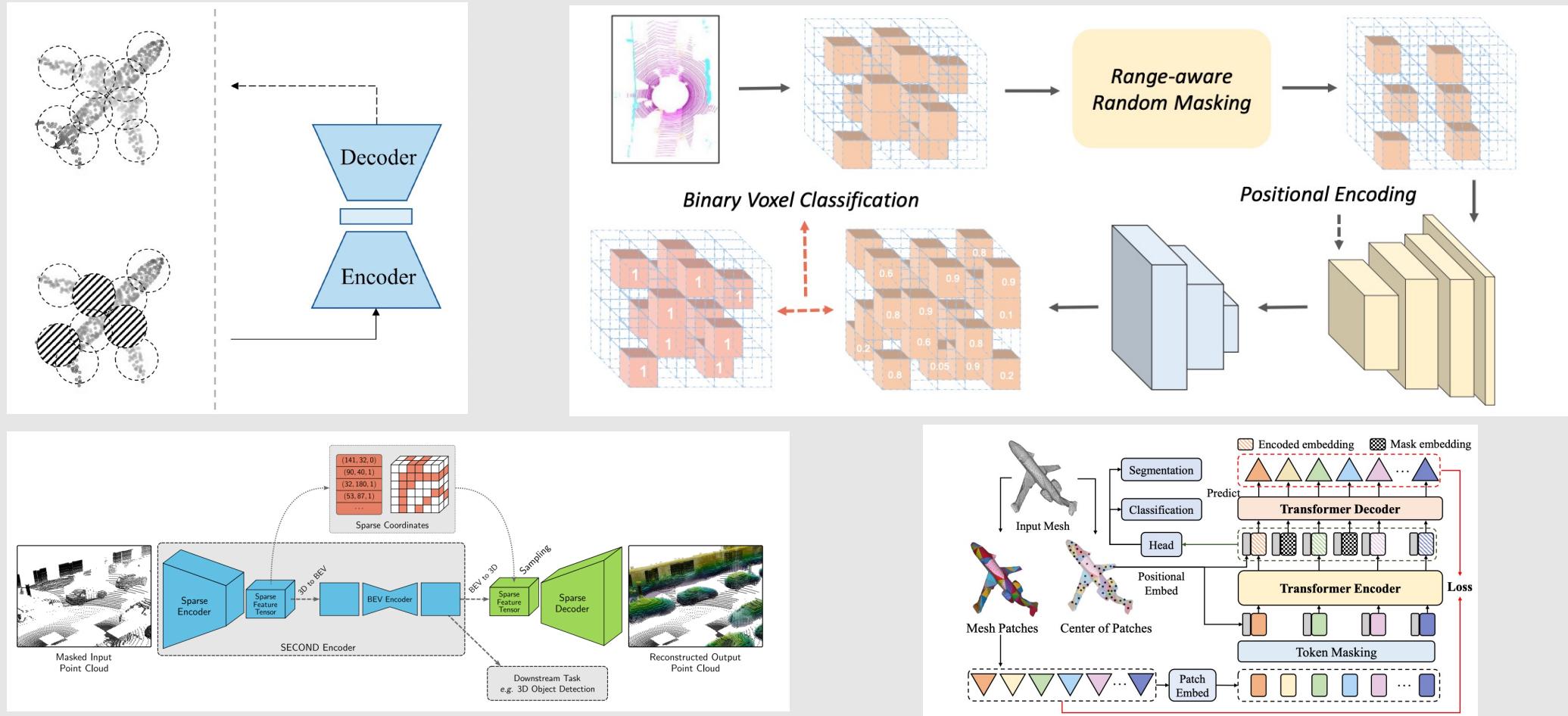
- **Robotics**

Masked Autoencoder: Applications



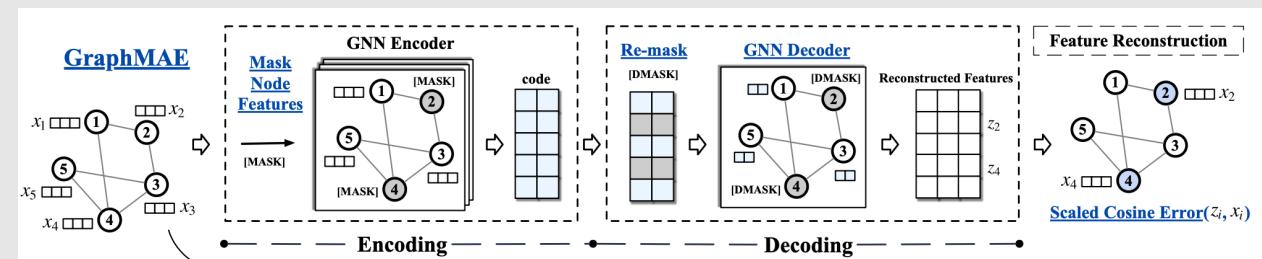
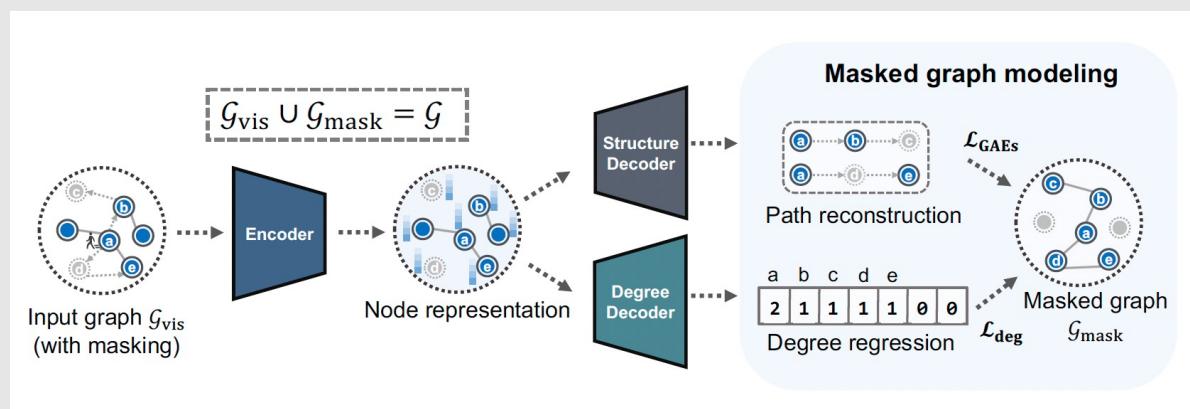
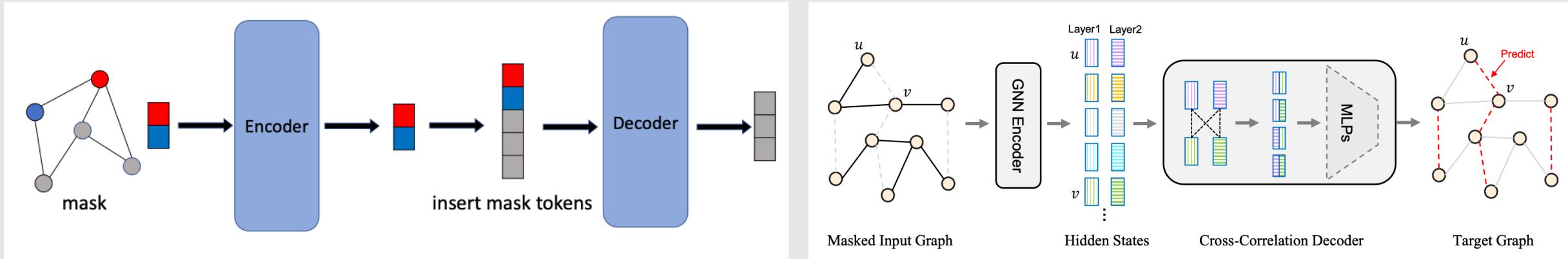
- Medical Imaging

Masked Autoencoder: Applications



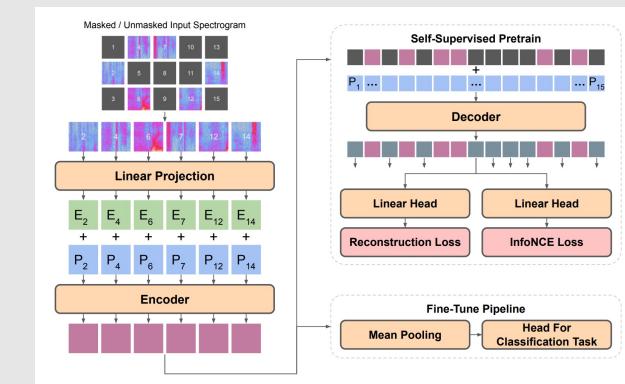
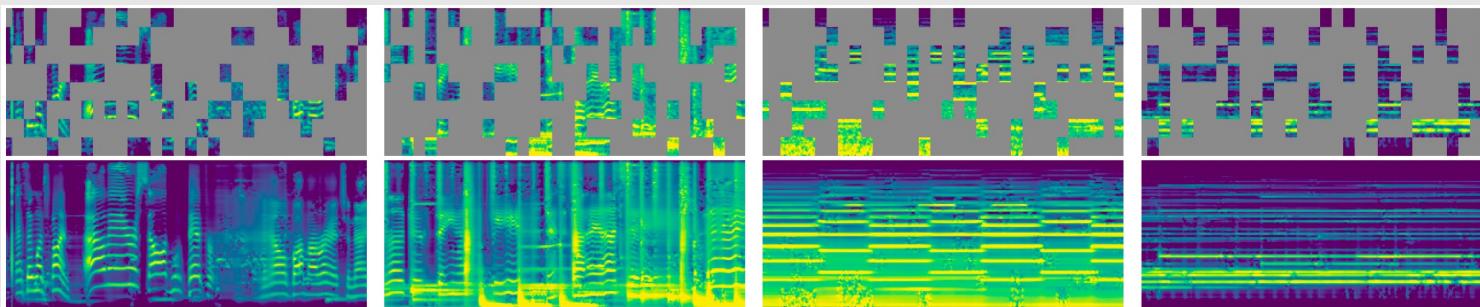
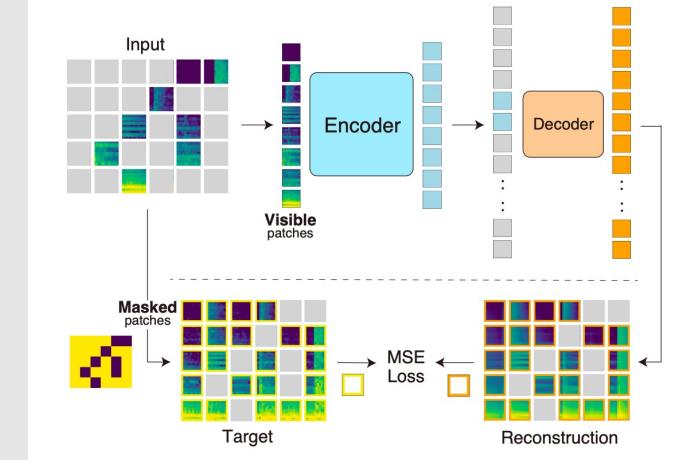
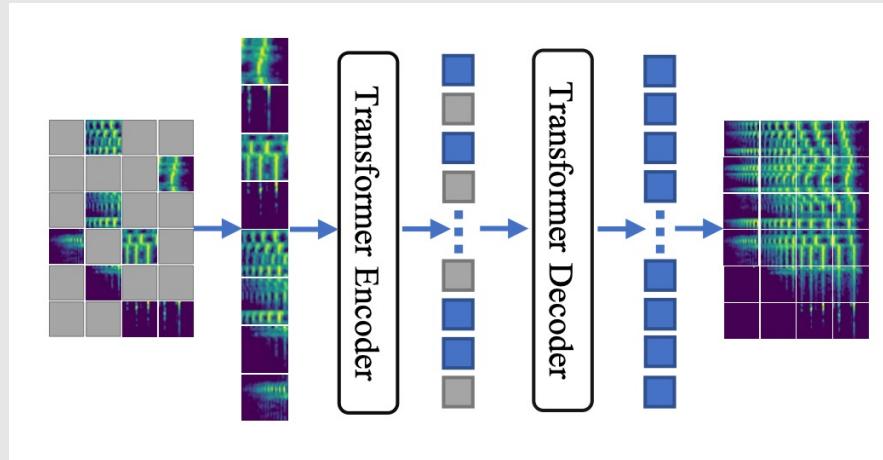
- **3D Geometry**

Masked Autoencoder: Applications



- **Graphs**

Masked Autoencoder: Applications



- **Audio**

Summary

Reconstruction-based Representation Learning

- Autoencoding (k-means, PCAs, AE, DAE, BERT, MAE, ...)
- (+) simple and easy to use
- (+) require little domain-knowledge
- (-/+) objective encourages preserving all nuances

