

<https://introml.mit.edu/>

# 6.390 Intro to Machine Learning

## Lecture 8: Representation Learning

Shen Shen

Mar 30, 2026

3pm, Room 10-250

[Slides and Lecture Recording](#)

# Outline

# Outline

- Neural networks are *representation* learners

# Outline

- Neural networks are *representation* learners
- Self-supervised learning

# Outline

- Neural networks are *representation* learners
- Self-supervised learning
- Auto-encoders

# Outline

- Neural networks are *representation* learners
- Self-supervised learning
- Auto-encoders
- Word embeddings

# Outline

- Neural networks are *representation* learners
- Self-supervised learning
- Auto-encoders
- Word embeddings
- (Some recent representation learning ideas)

# Outline

- Neural networks are *representation* learners
- Self-supervised learning
- Auto-encoders
- Word embeddings
- (Some recent representation learning ideas)

Recap

# Recap

input

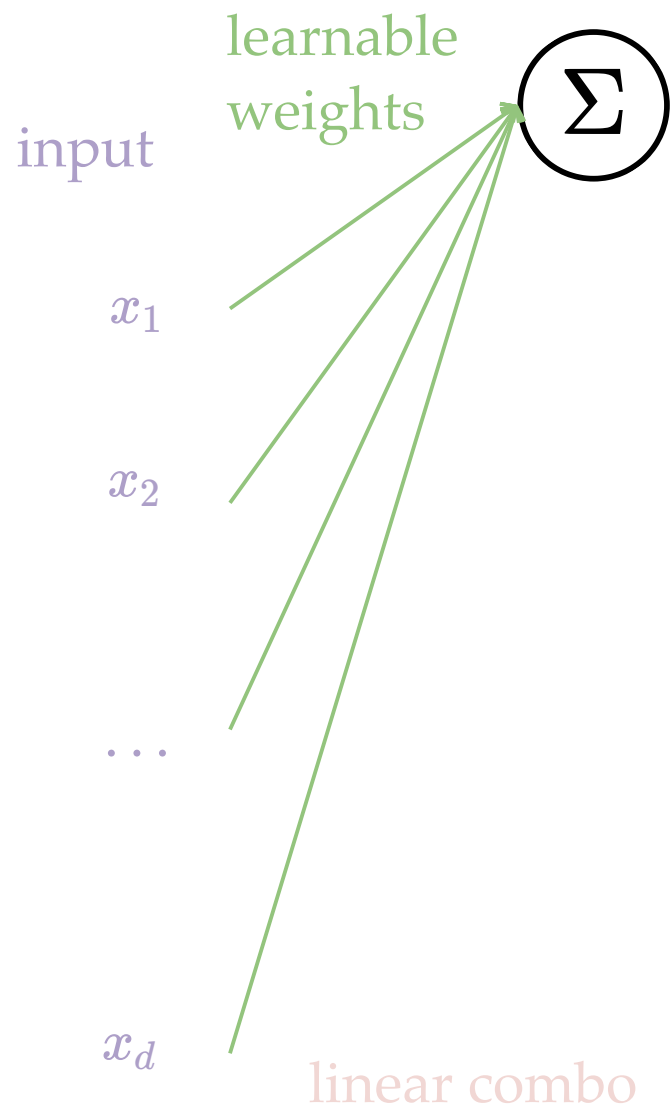
$x_1$

$x_2$

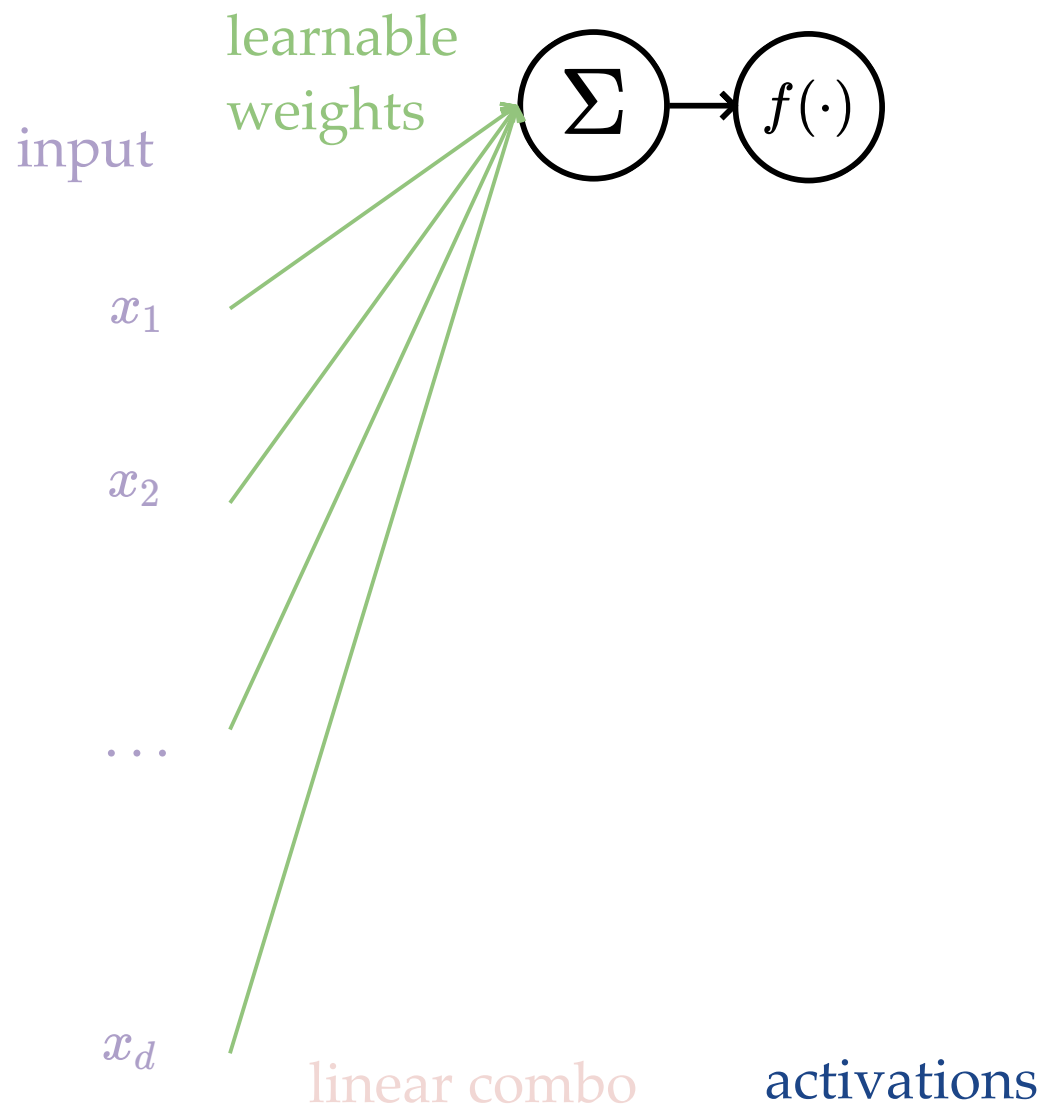
...

$x_d$

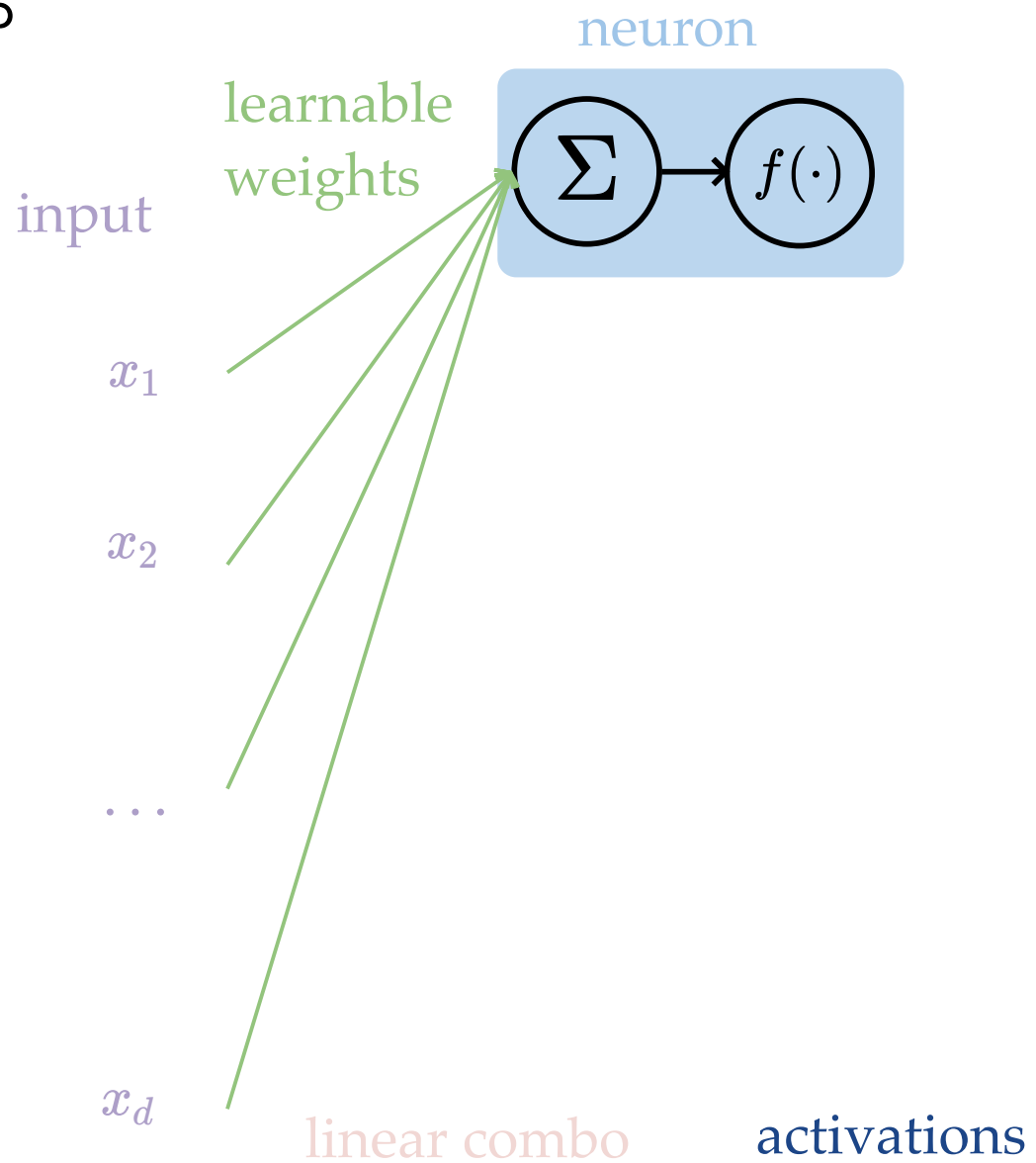
# Recap



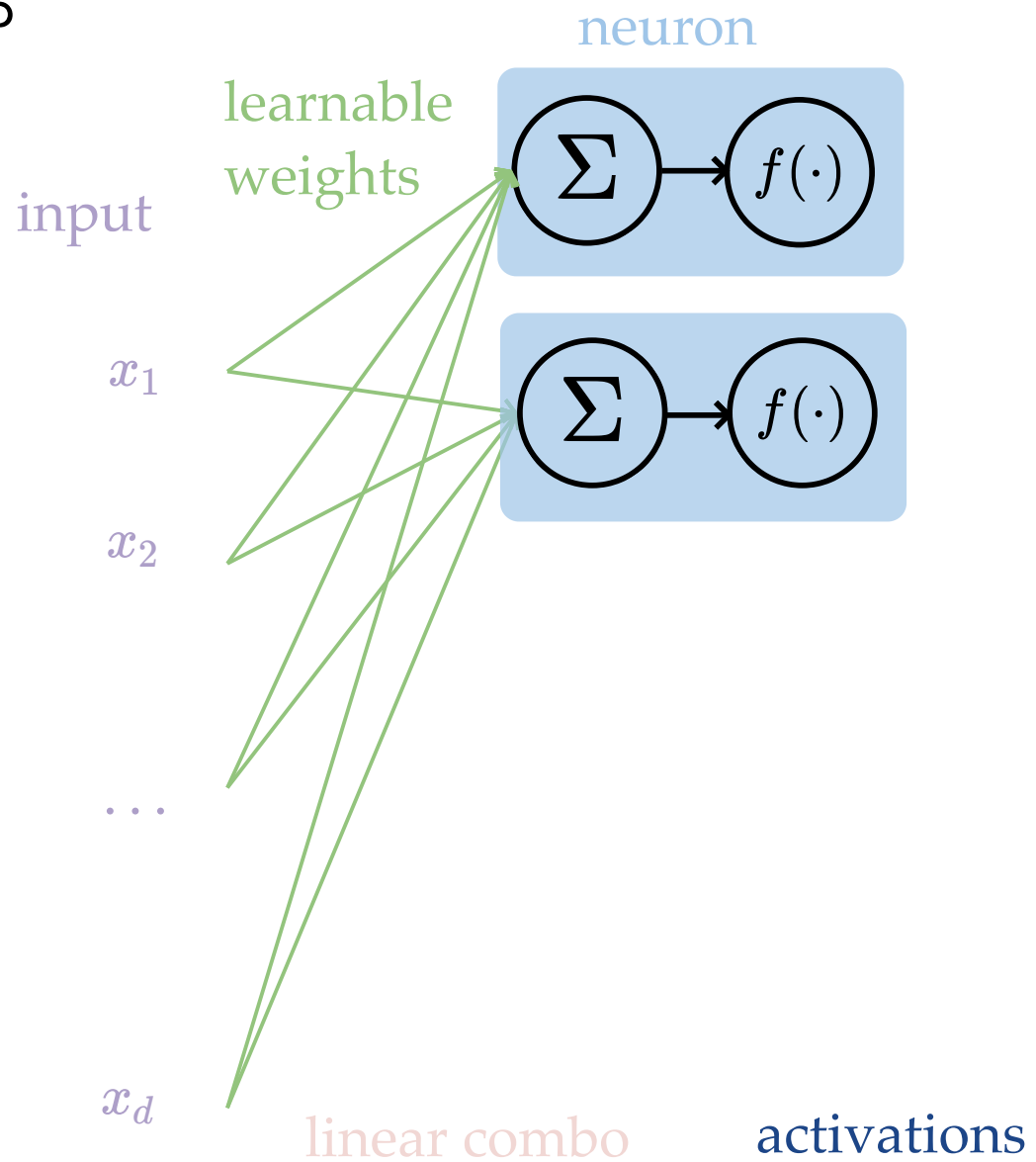
# Recap



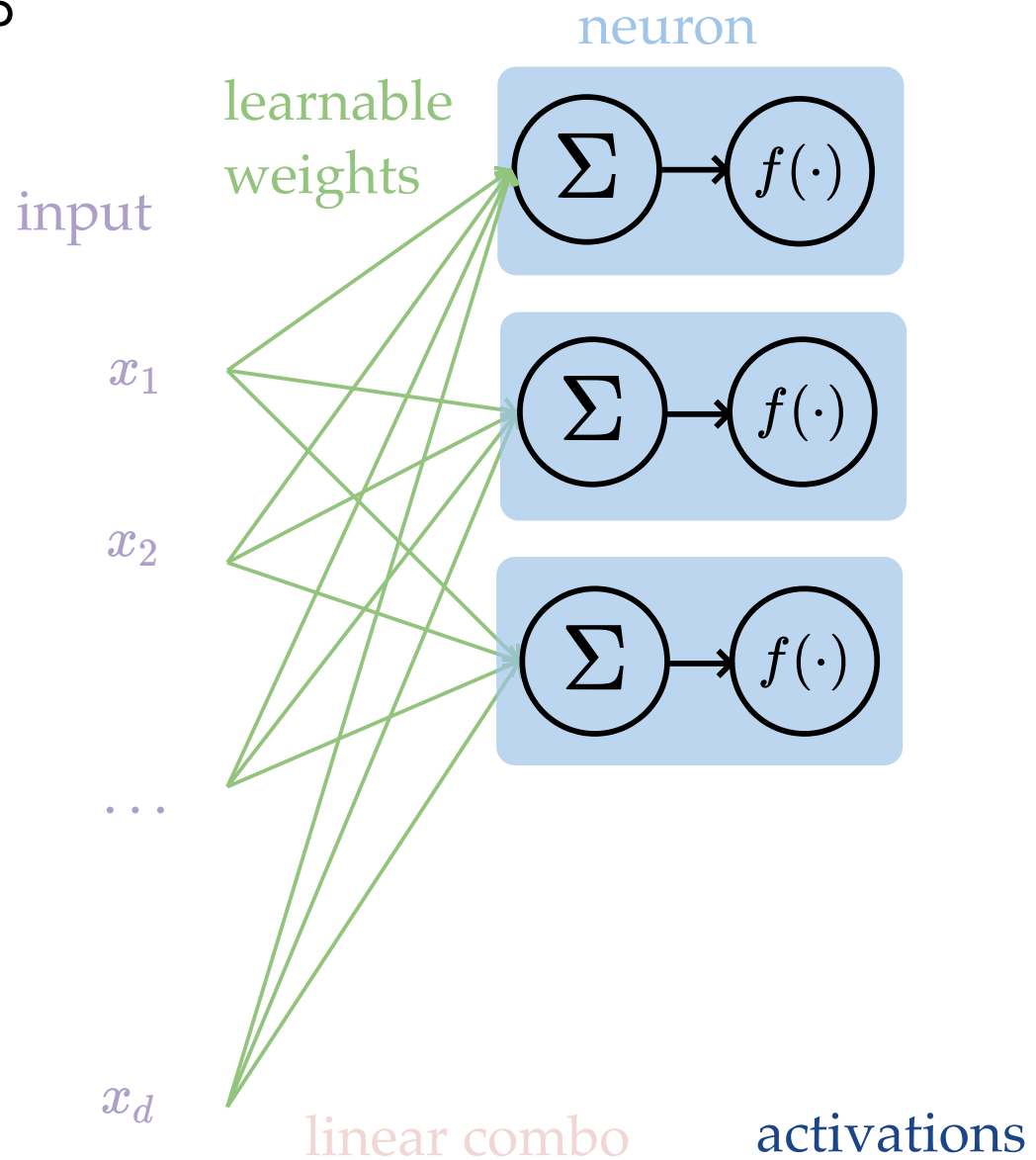
# Recap



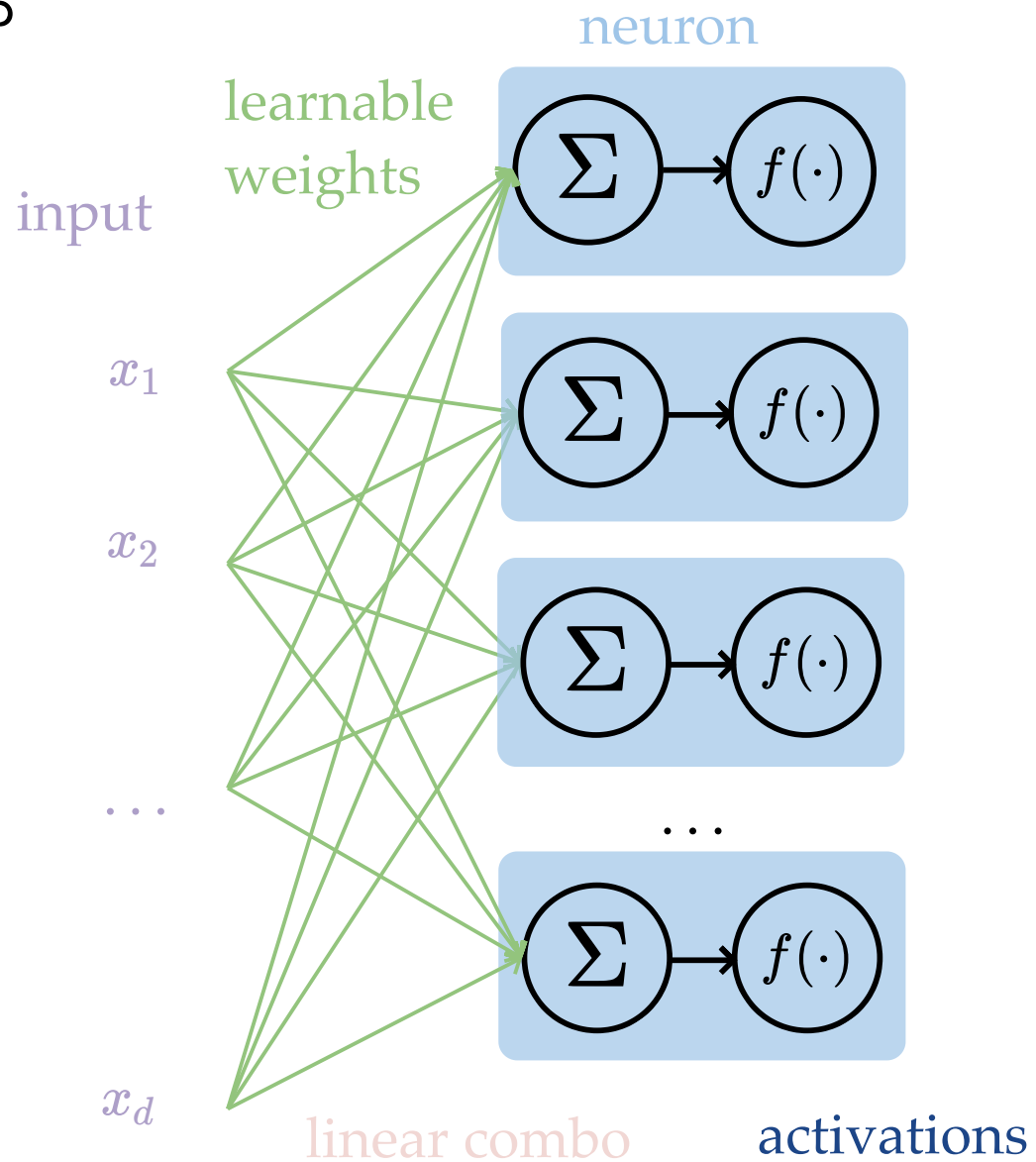
# Recap



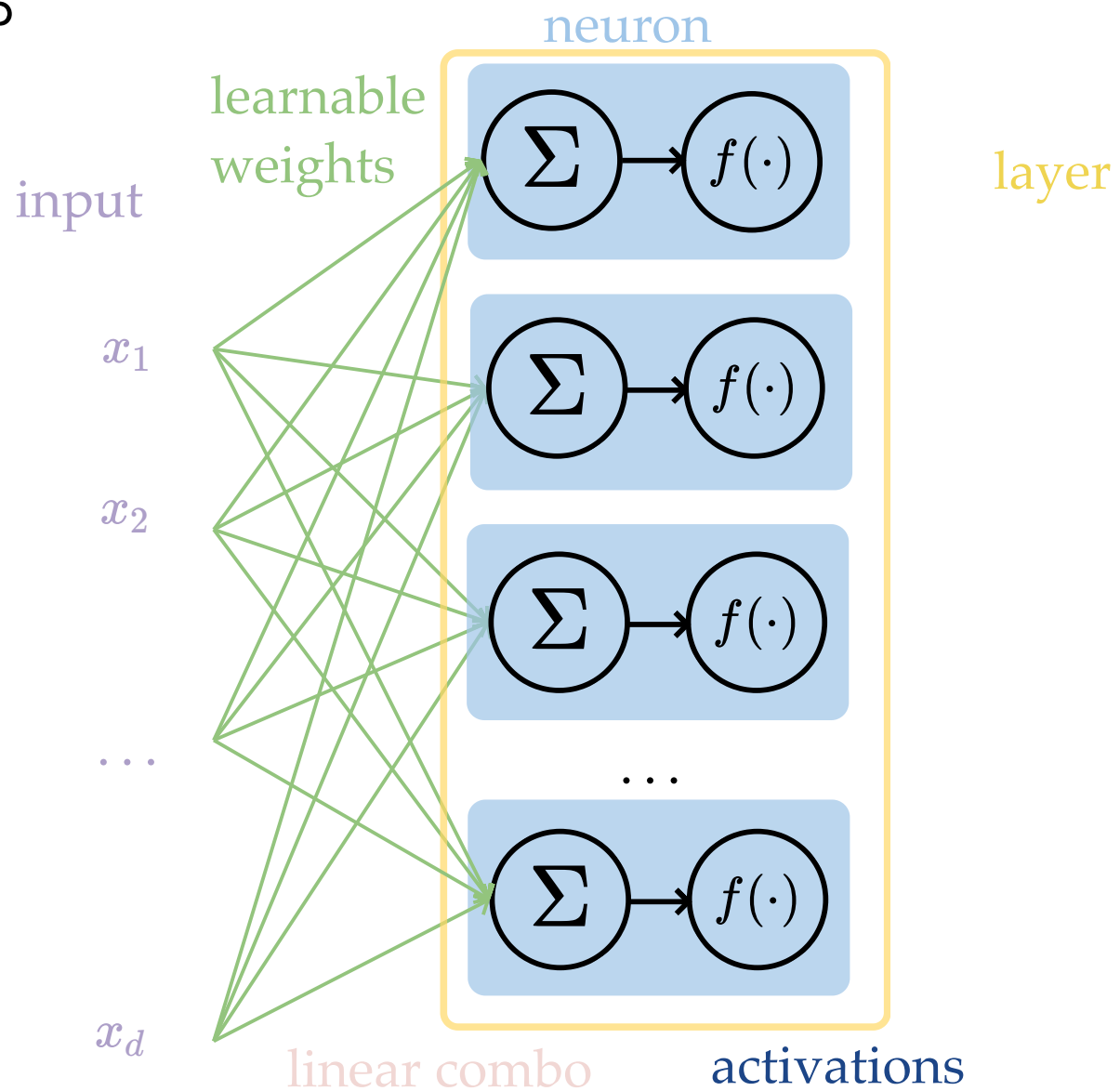
# Recap



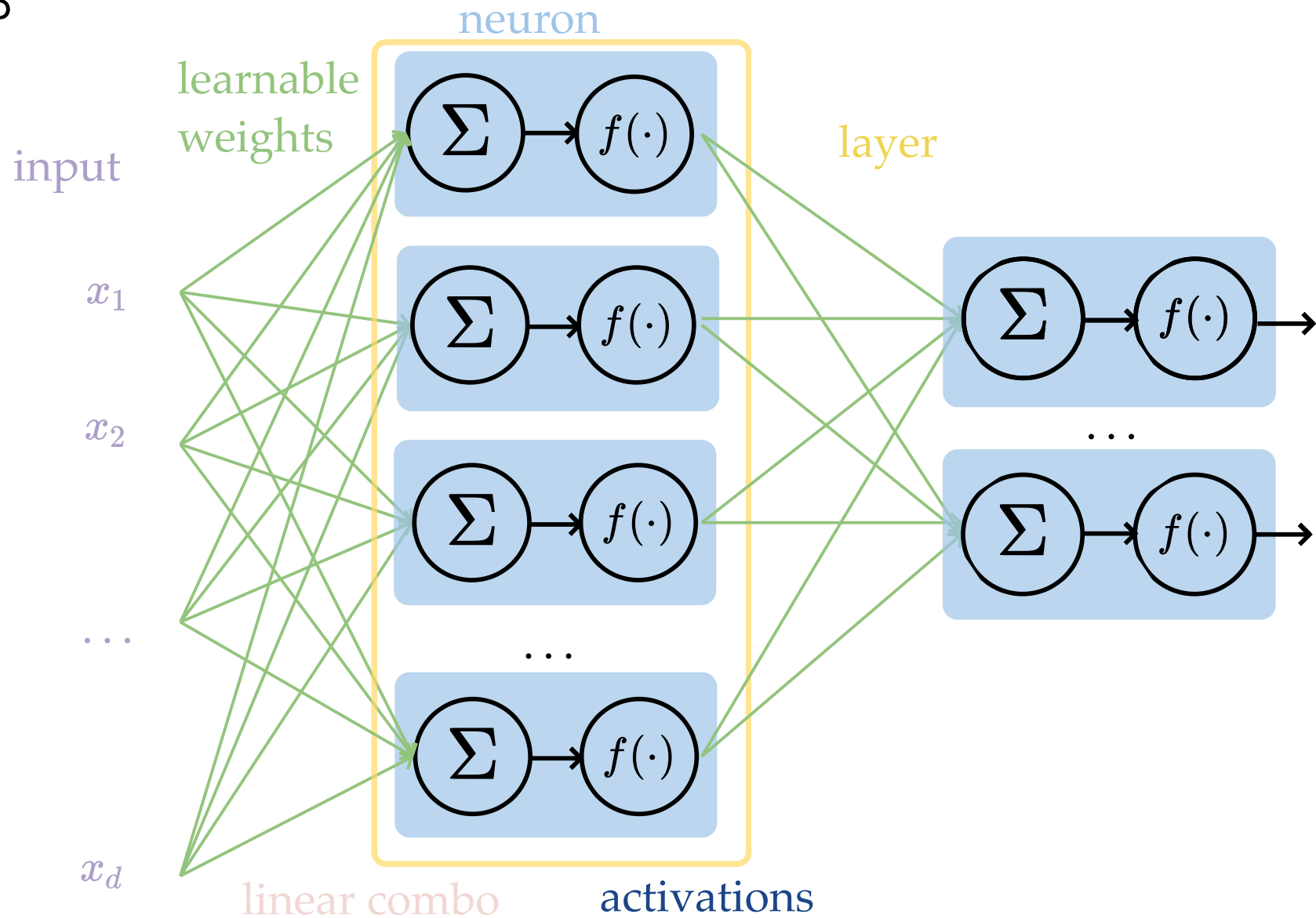
# Recap



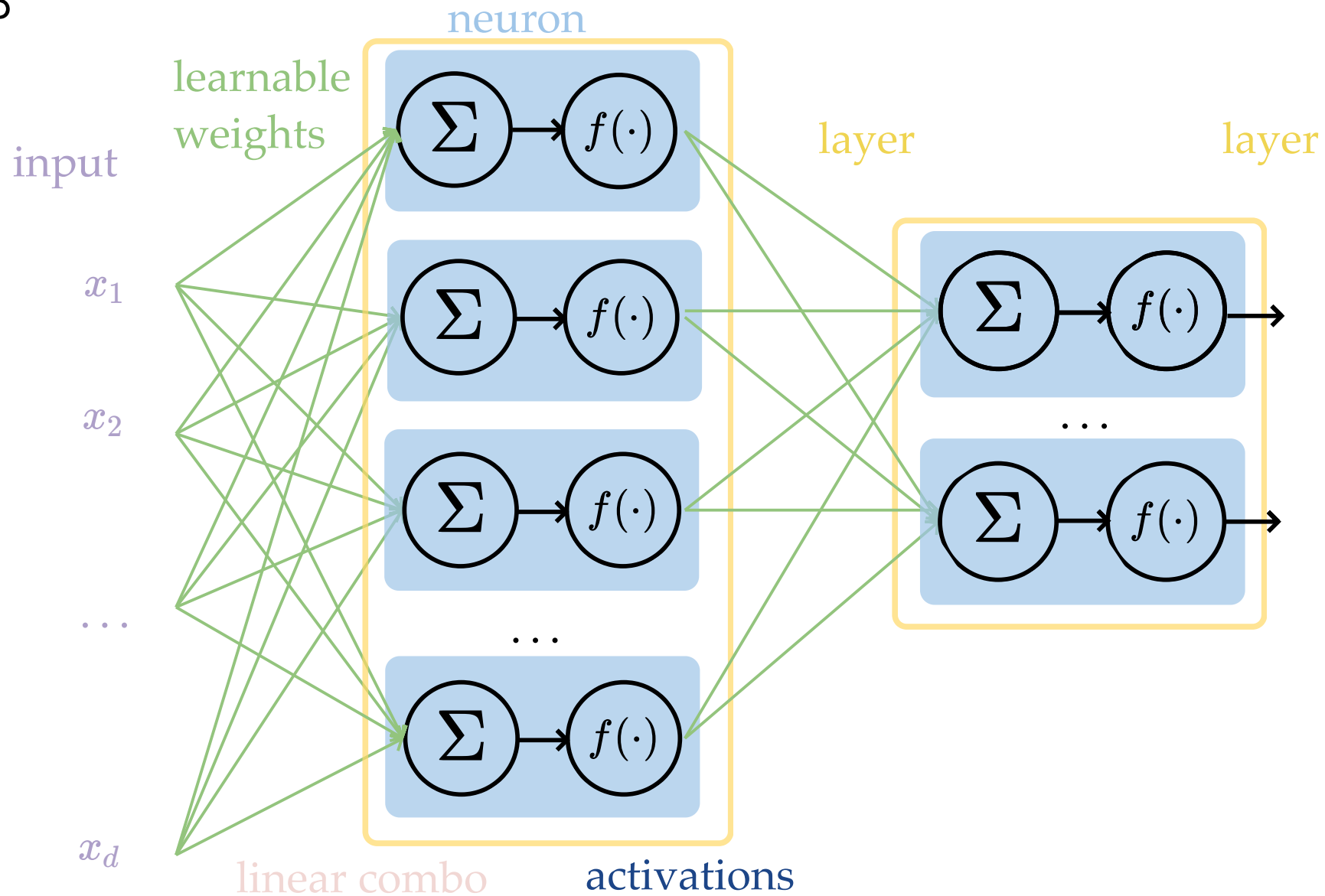
# Recap



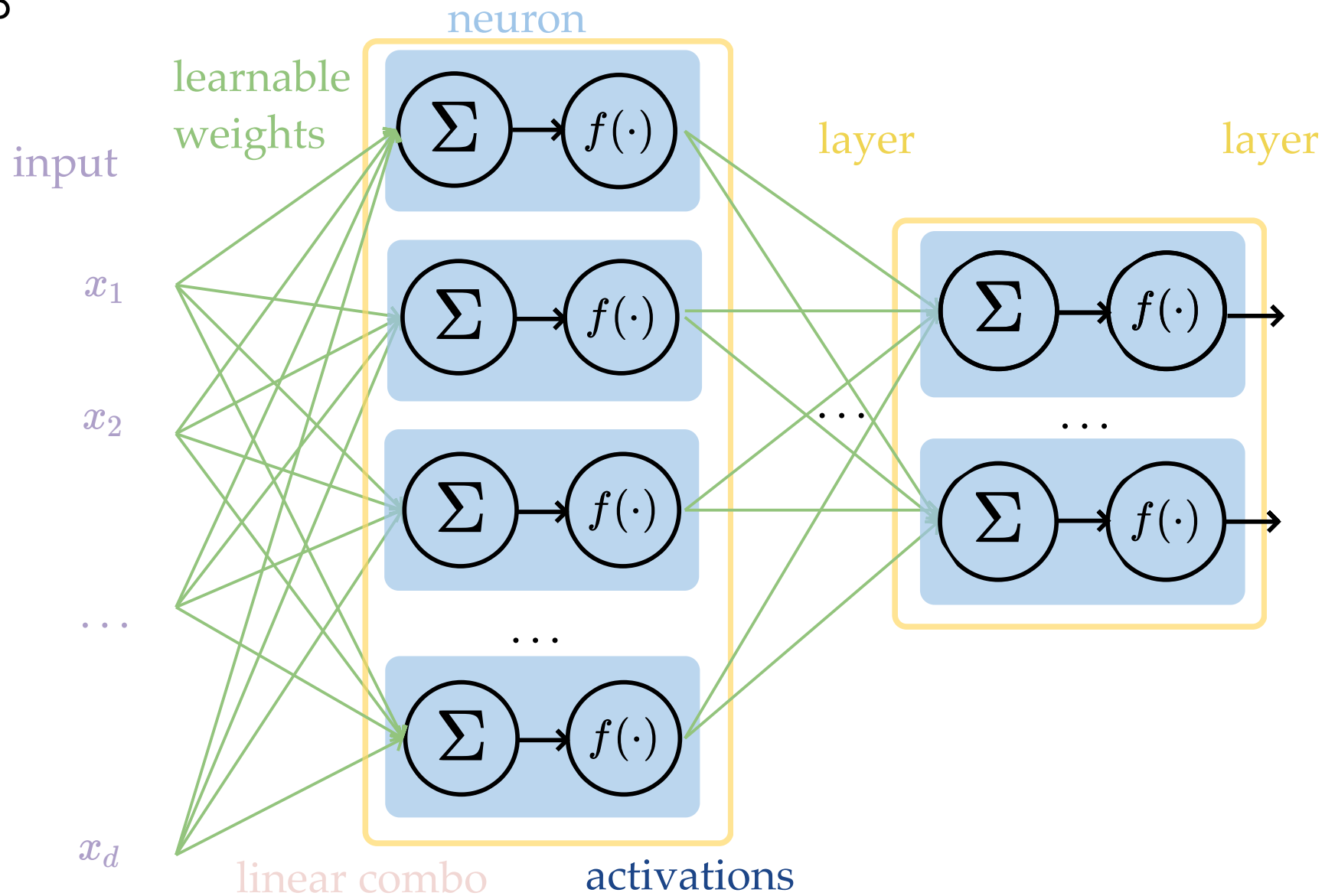
# Recap



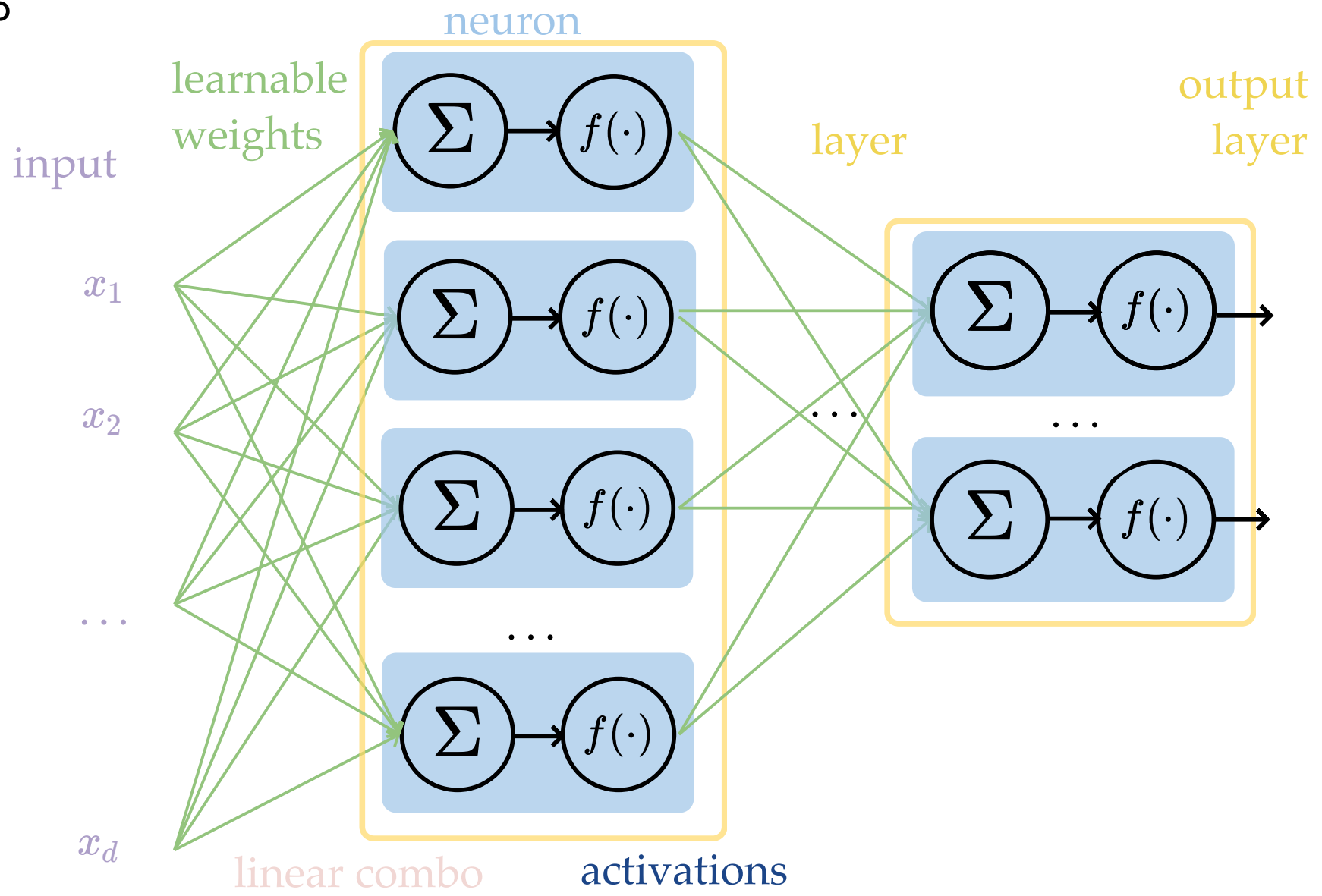
# Recap



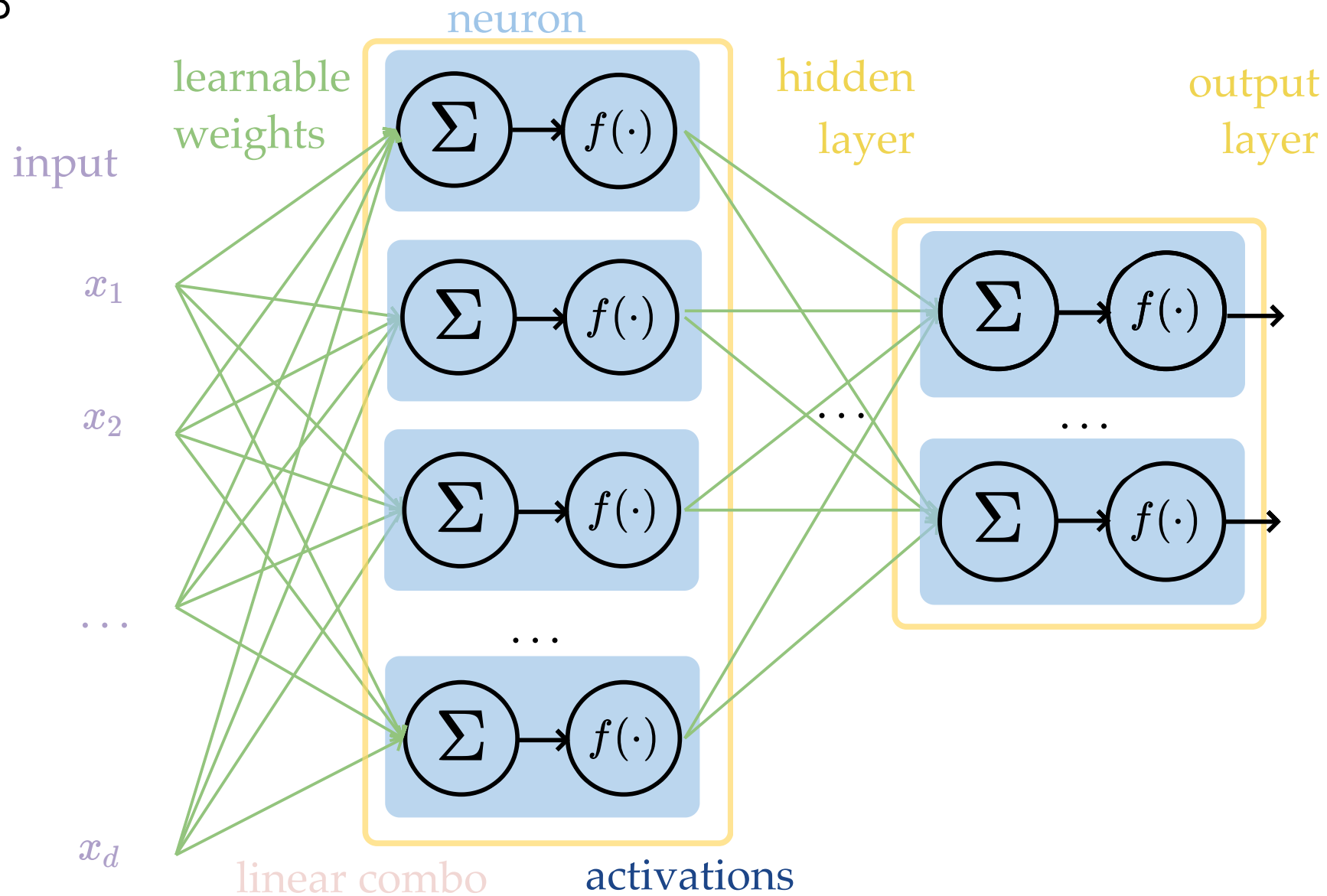
# Recap

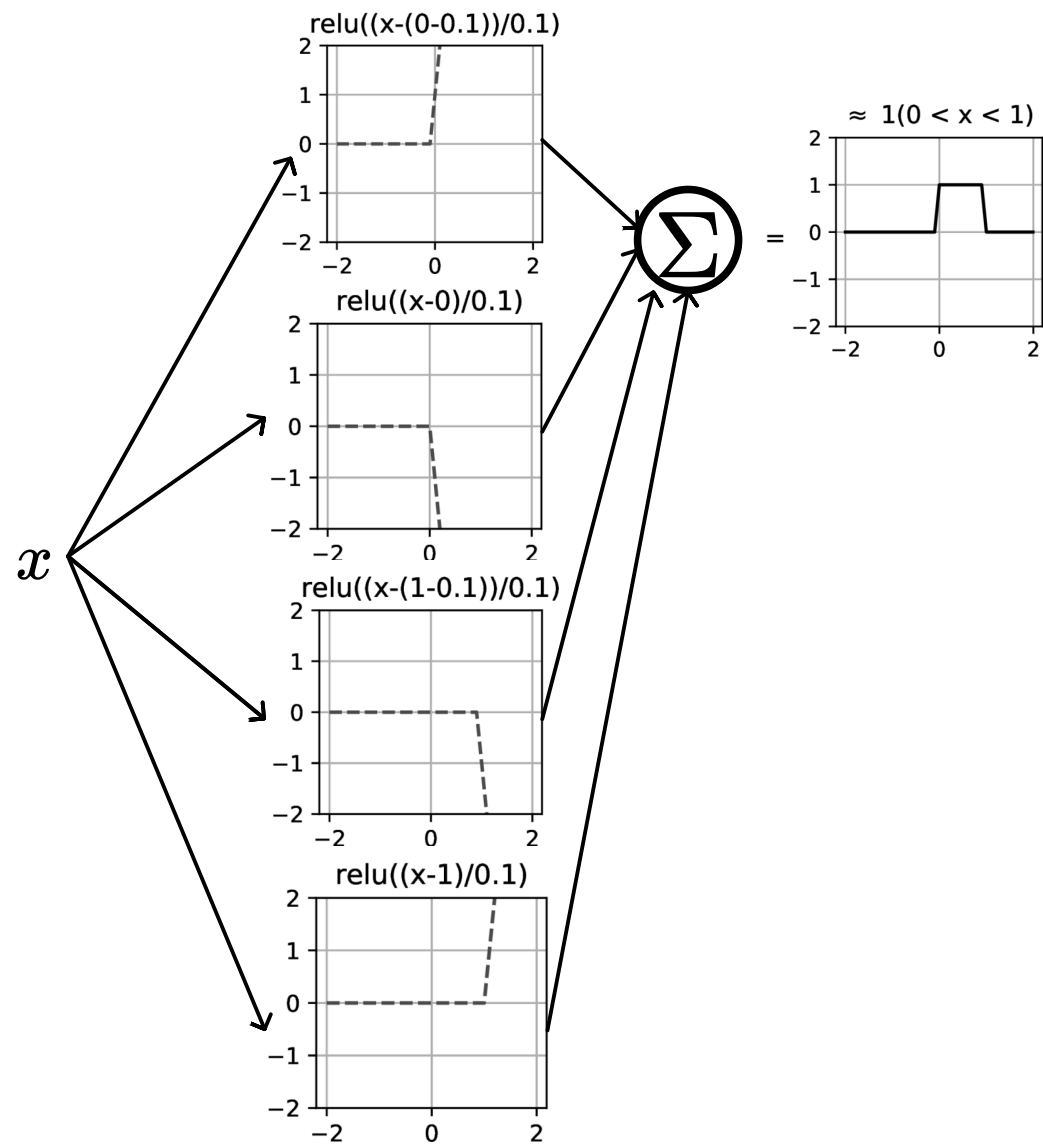


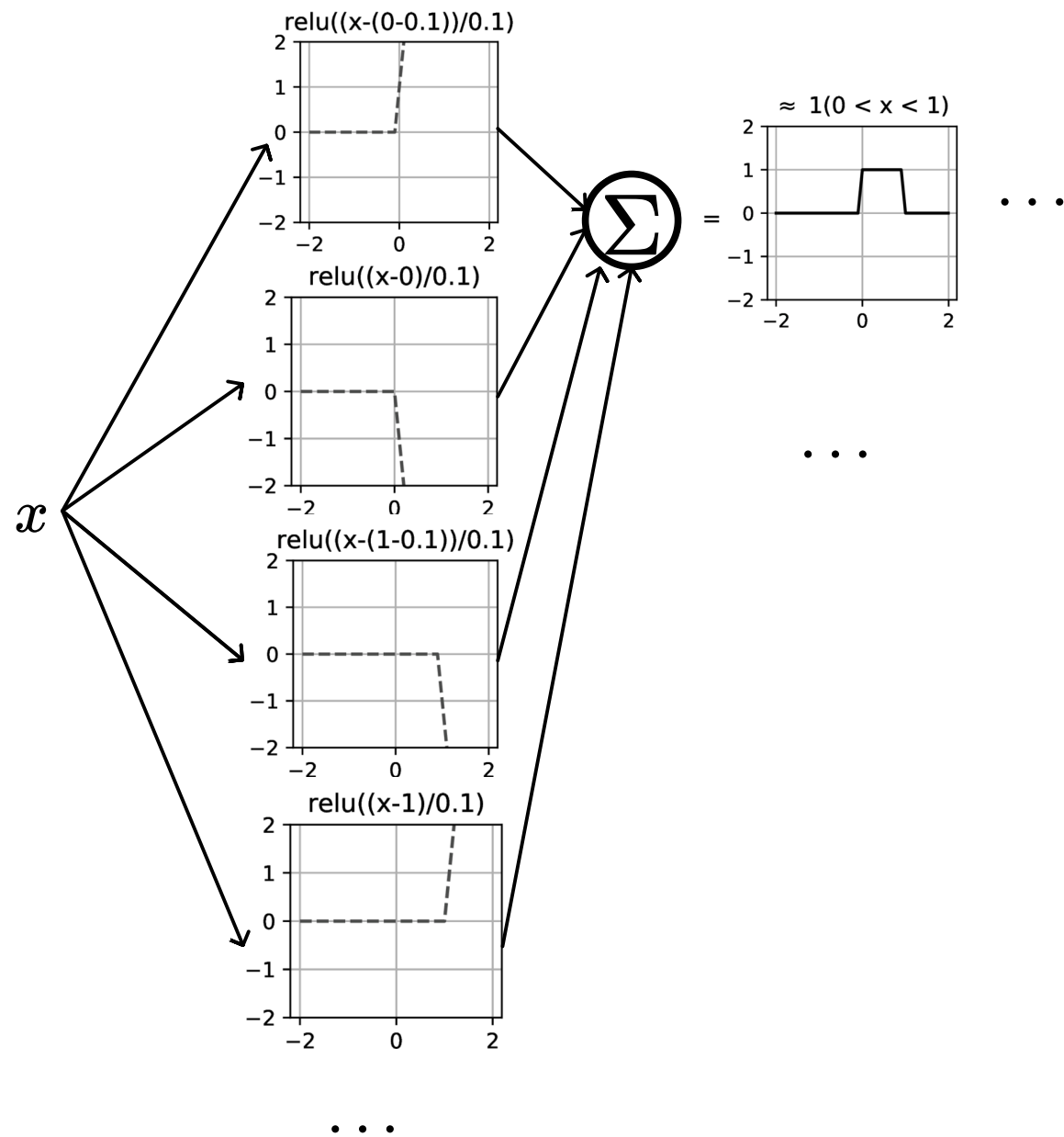
# Recap

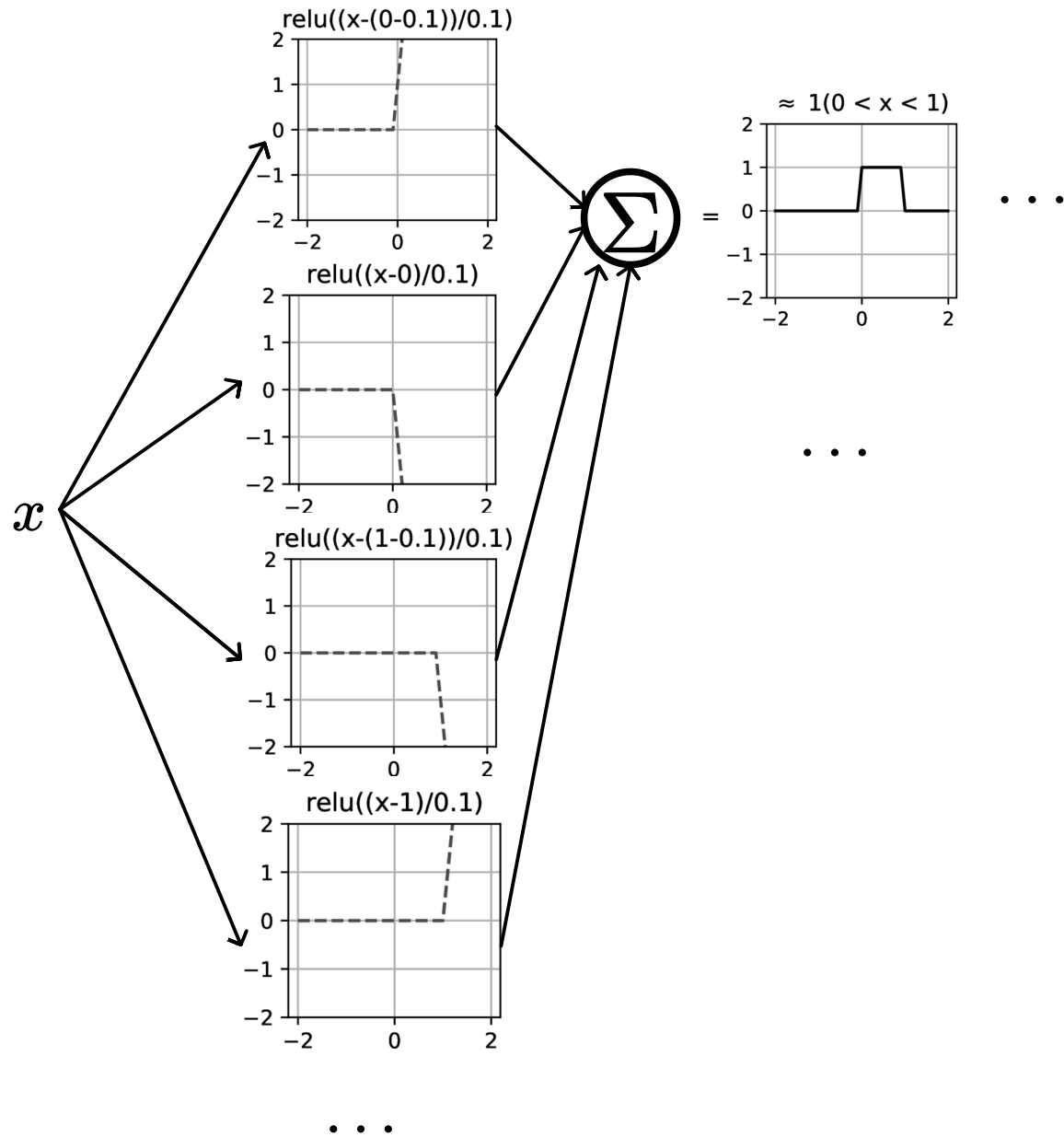


# Recap

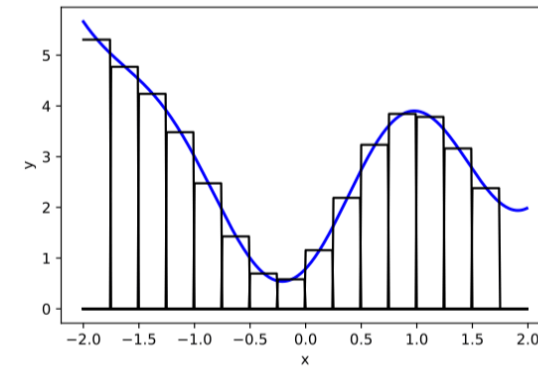








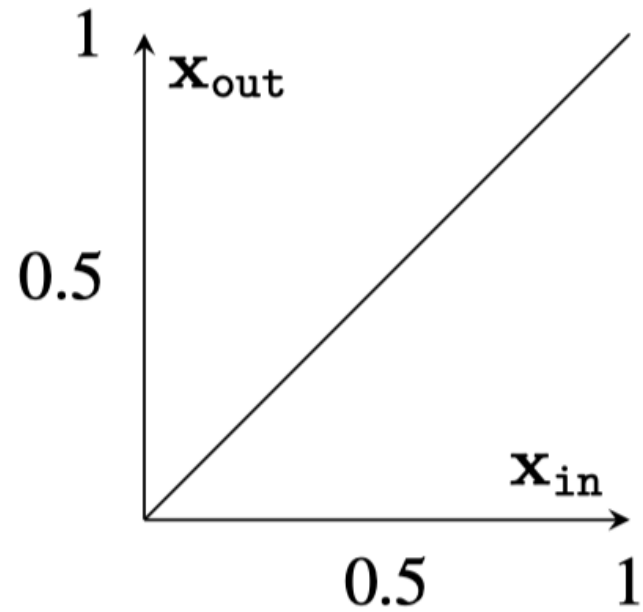
asymptotically, can approximate any function!



## Two different ways to visualize a function

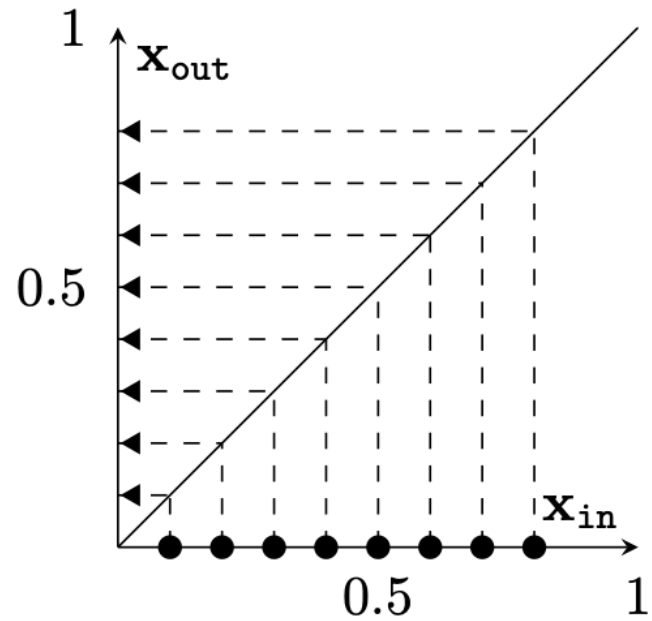
Two different ways to visualize a function

e.g. the identity function



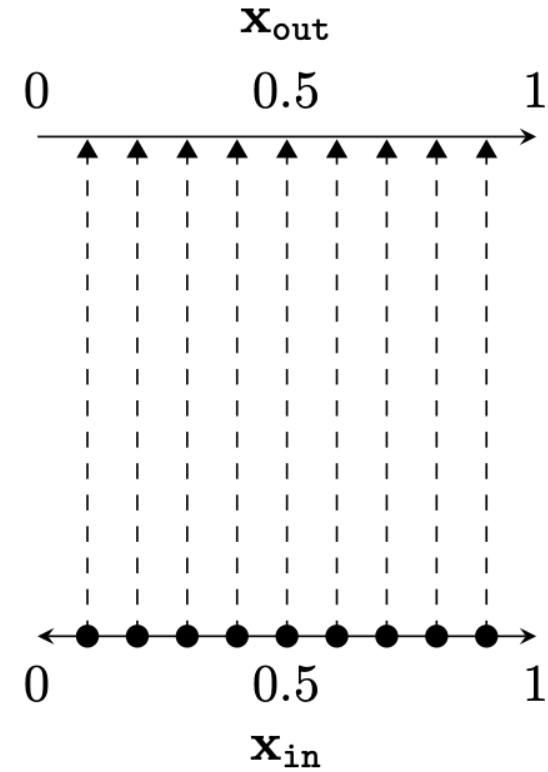
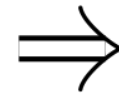
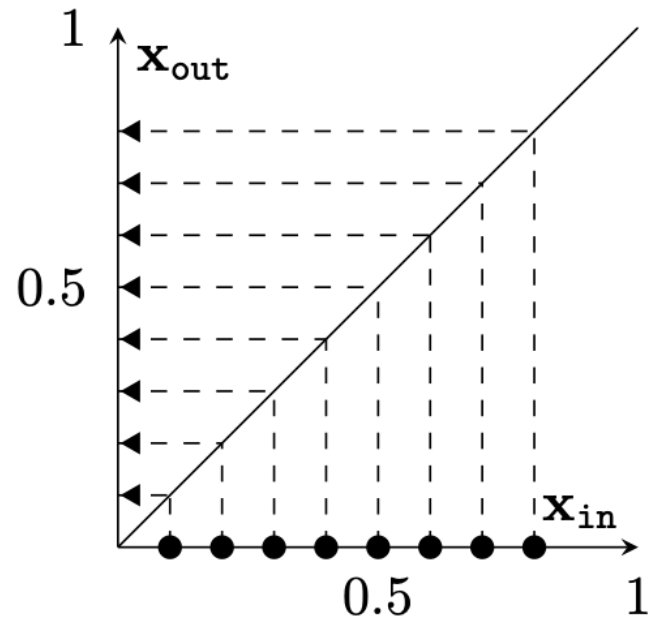
Two different ways to visualize a function

e.g. the identity function

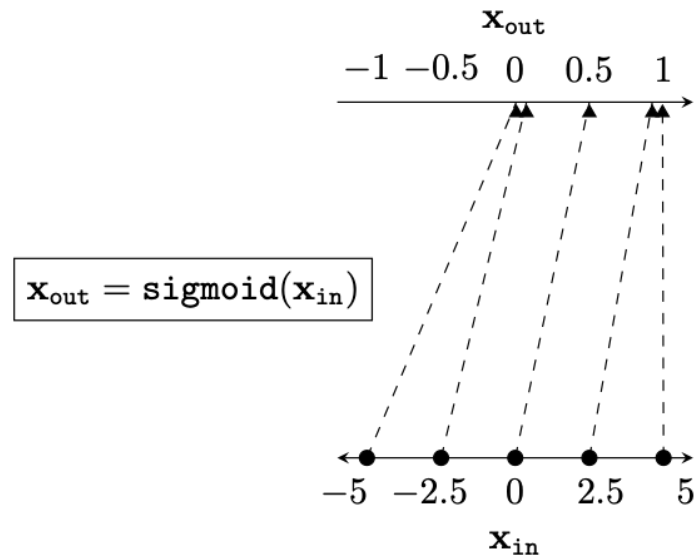
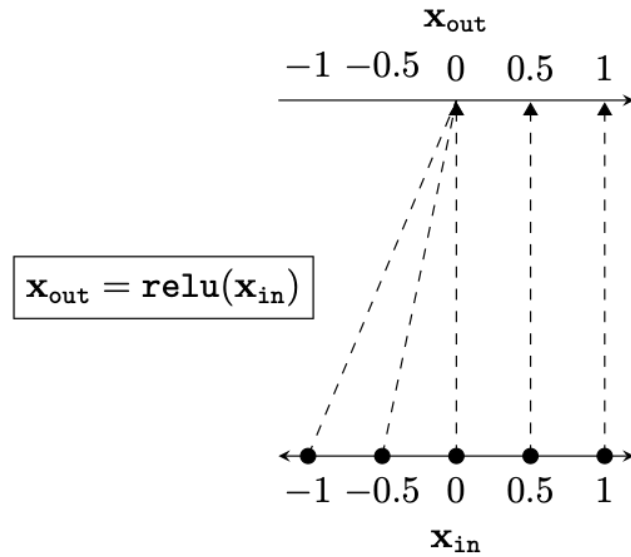
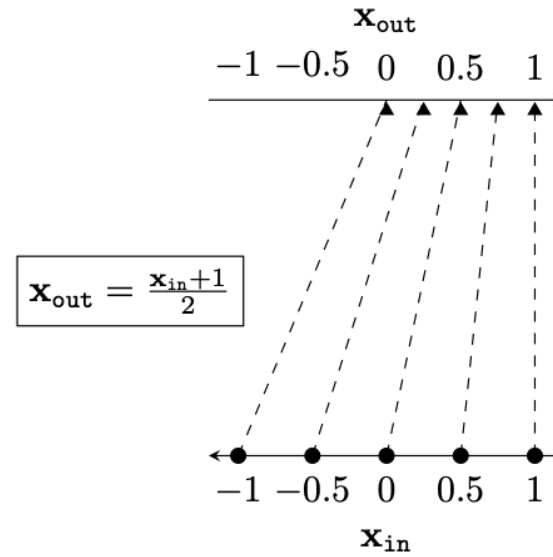
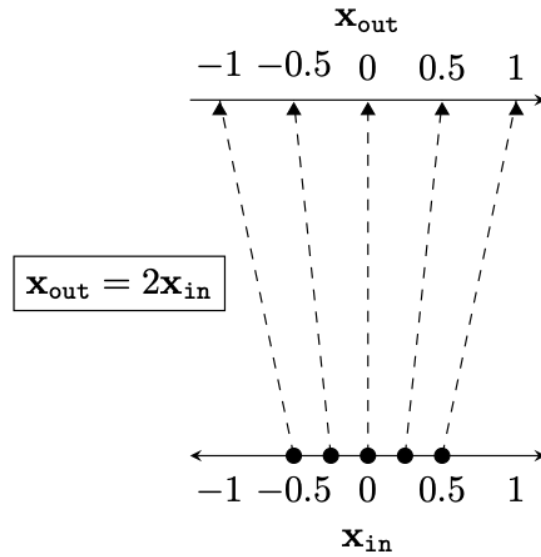


Two different ways to visualize a function

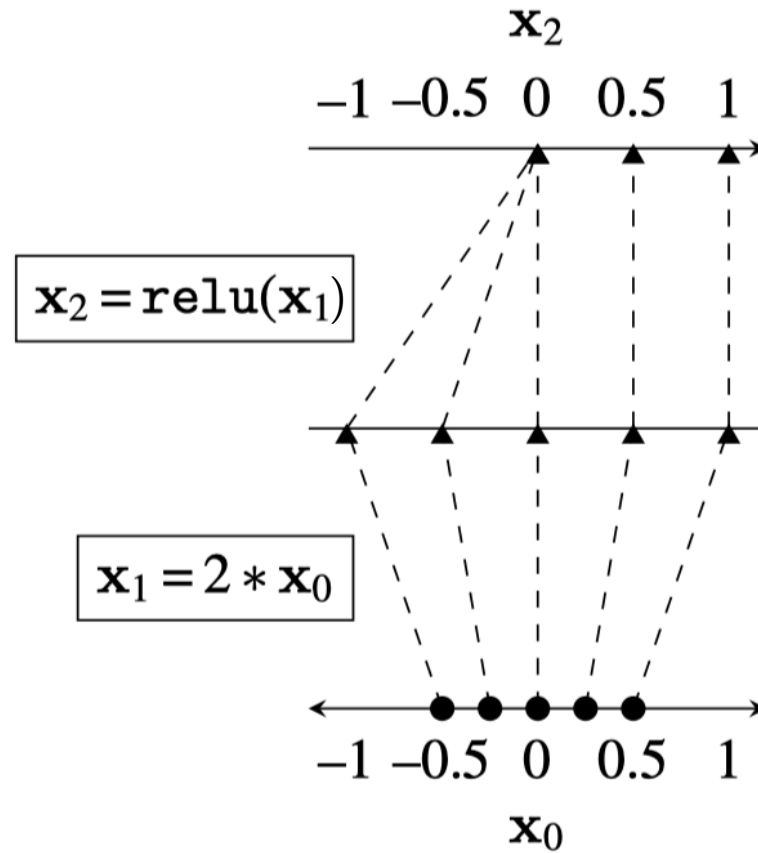
e.g. the identity function



# Representation transformations for a variety of neural net operations



and stack of neural net operations

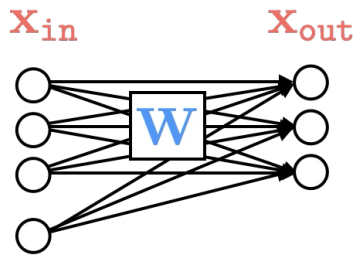


Parameters

Parameters

wiring graph

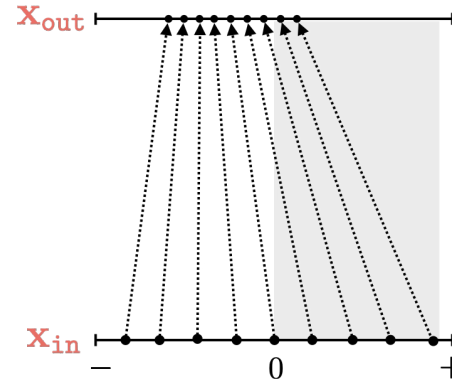
linear



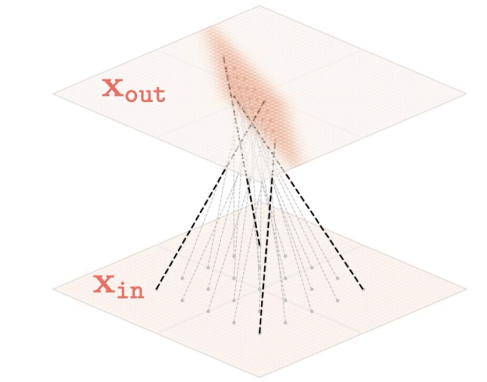
equation

$$X_{out} = W X_{in}$$

mapping 1D



mapping 2D



Parameters

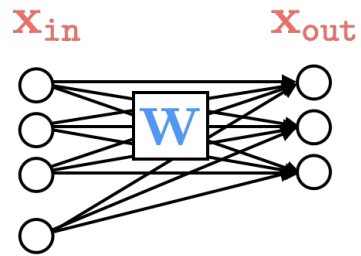
wiring graph

equation

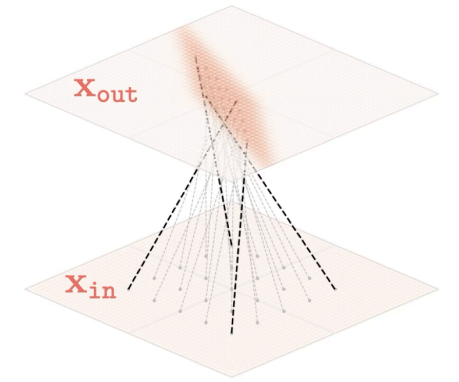
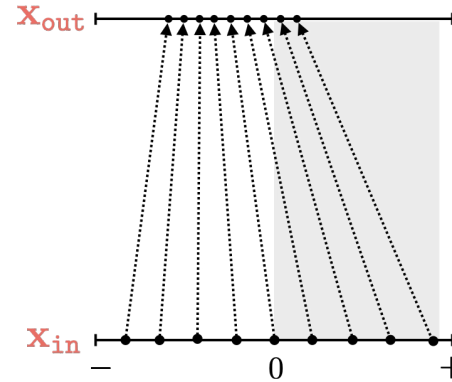
mapping 1D

mapping 2D

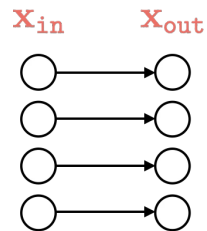
linear



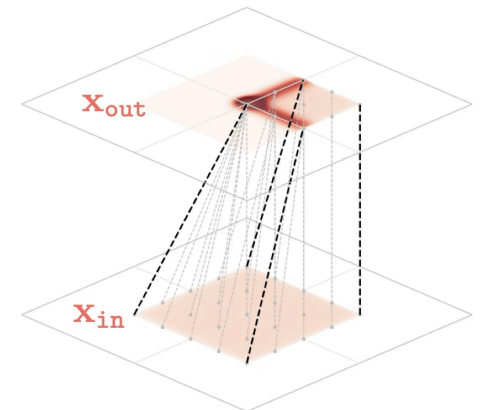
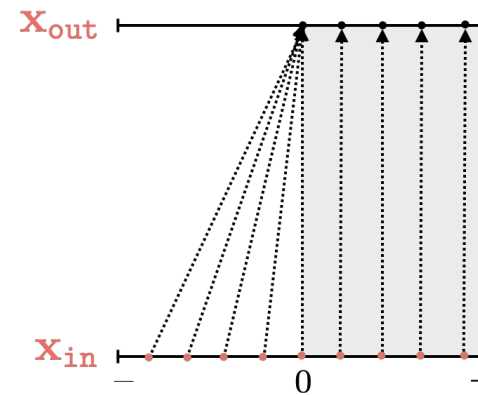
$$x_{out} = W x_{in}$$



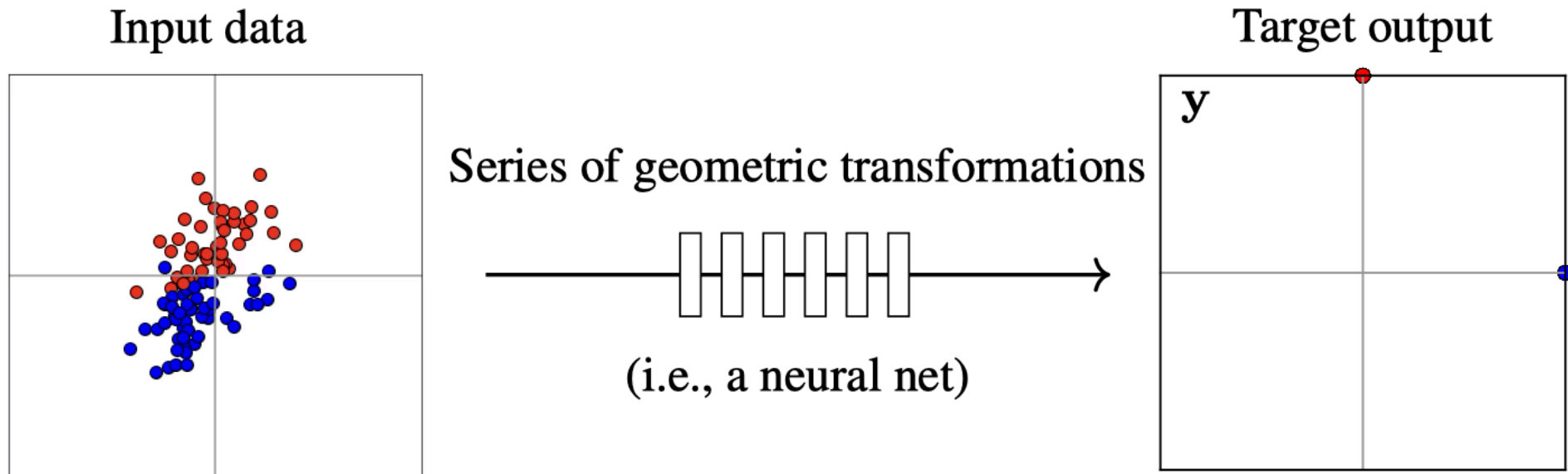
relu



$$x_{out_i} = \max(x_{in_i}, 0)$$



What does training a neural net classifier look like?



$$g = \text{softmax}(z_2)$$

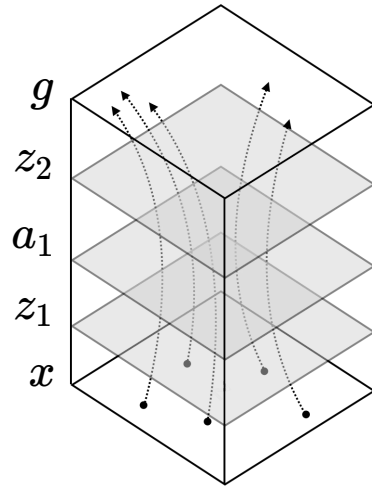
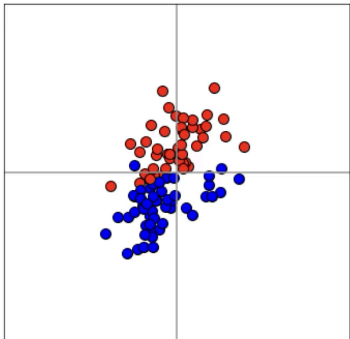
$$z_2 = \text{linear}(a_1) \in \mathbb{R}^2$$

$$a_1 = \text{ReLU}(z_1)$$

$$z_1 = \text{linear}(x) \in \mathbb{R}^2$$

$$x \in \mathbb{R}^2$$

Training data



$$g = \text{softmax}(z_2)$$

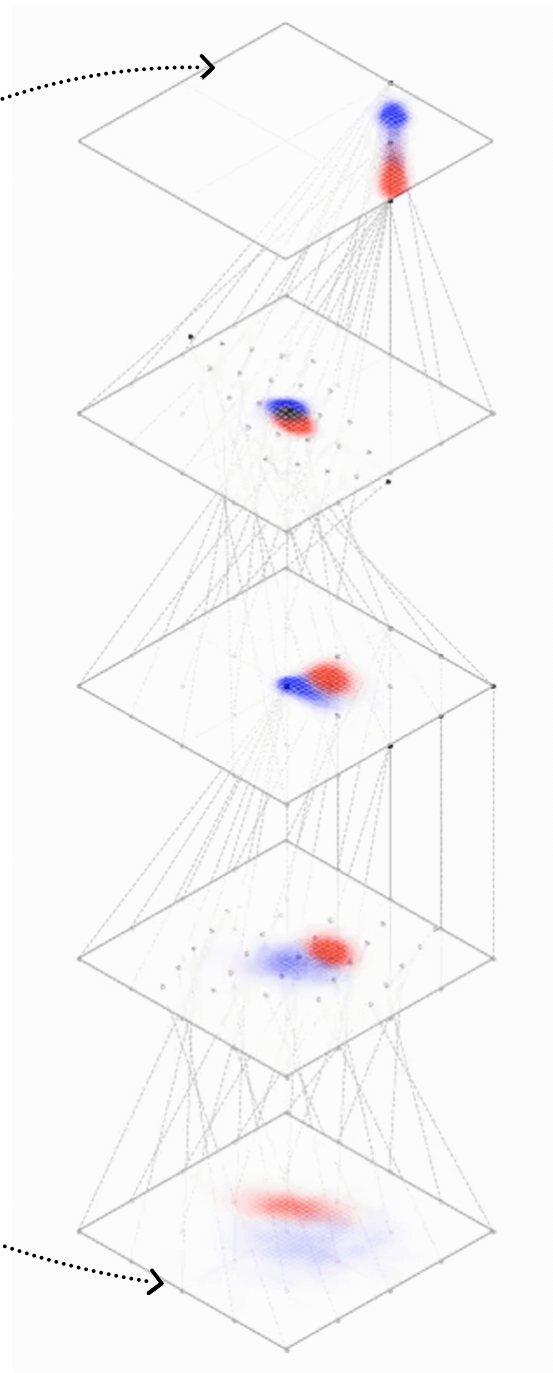
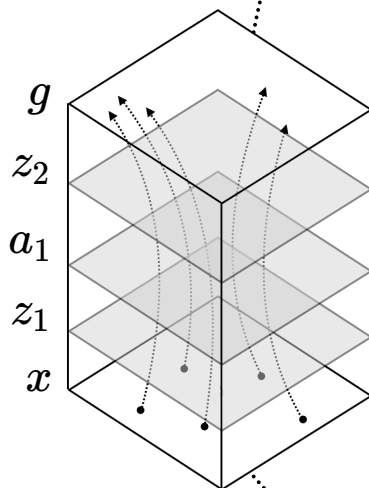
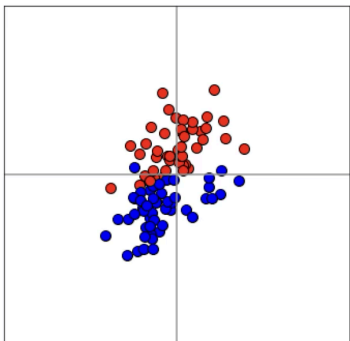
$$z_2 = \text{linear}(a_1) \in \mathbb{R}^2$$

$$a_1 = \text{ReLU}(z_1)$$

$$z_1 = \text{linear}(x) \in \mathbb{R}^2$$

$$x \in \mathbb{R}^2$$

Training data



$$g = \text{softmax}(z_2)$$

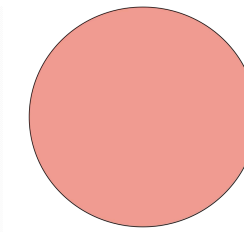
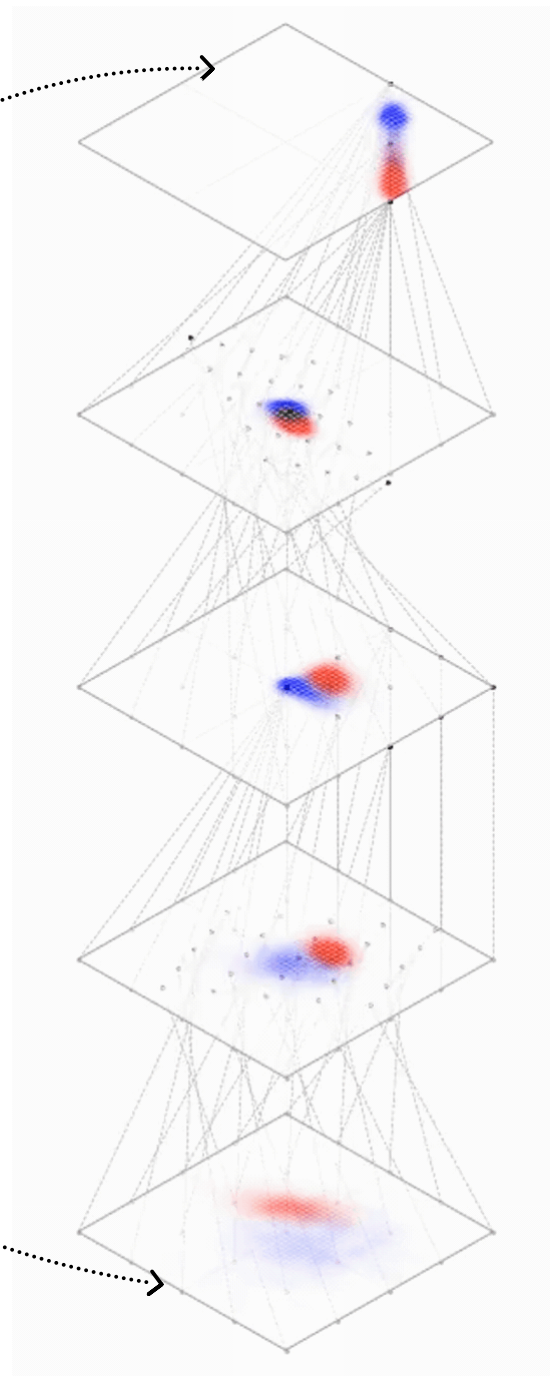
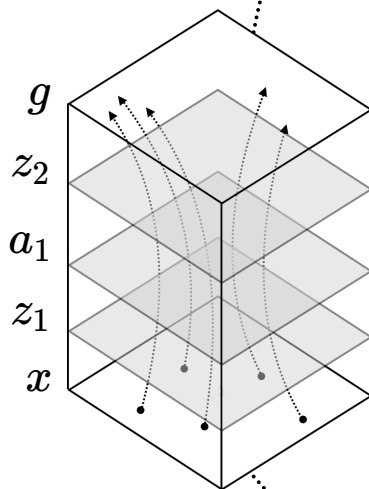
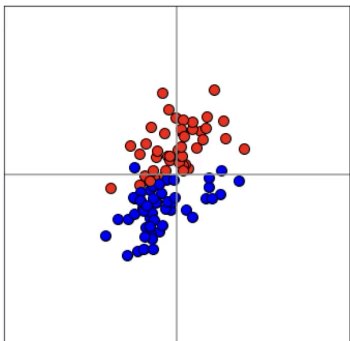
$$z_2 = \text{linear}(a_1) \in \mathbb{R}^2$$

$$a_1 = \text{ReLU}(z_1)$$

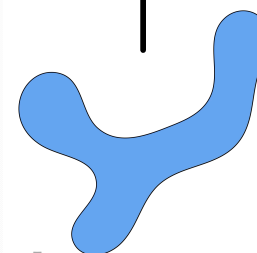
$$z_1 = \text{linear}(x) \in \mathbb{R}^2$$

$$x \in \mathbb{R}^2$$

Training data

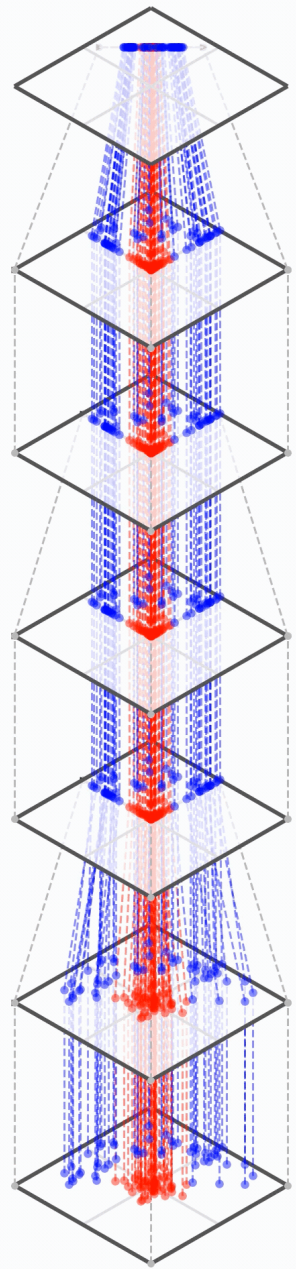


maps from  
complex data  
space to desired  
target space





Loss: 0.71



softmax

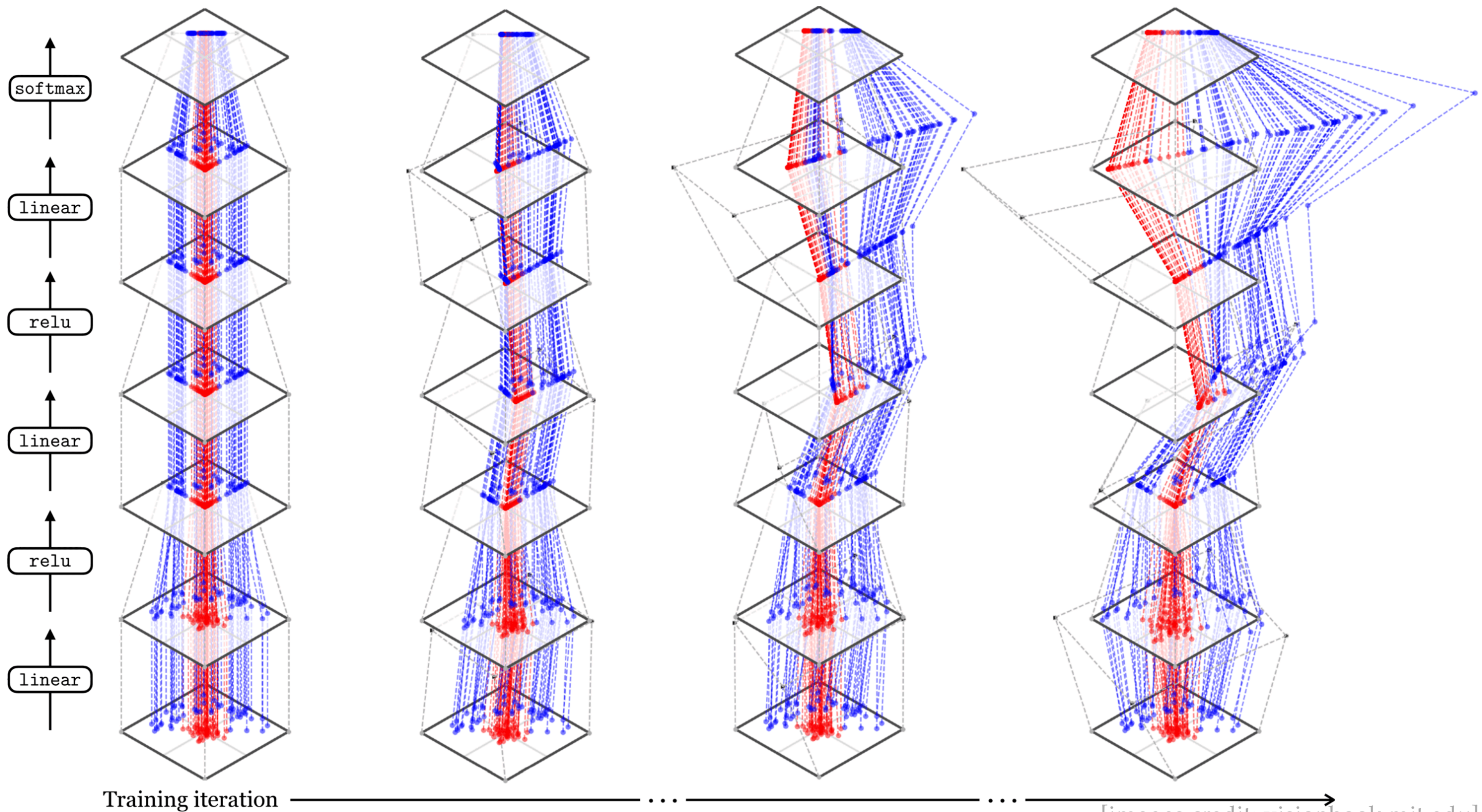
linear

relu

linear

relu

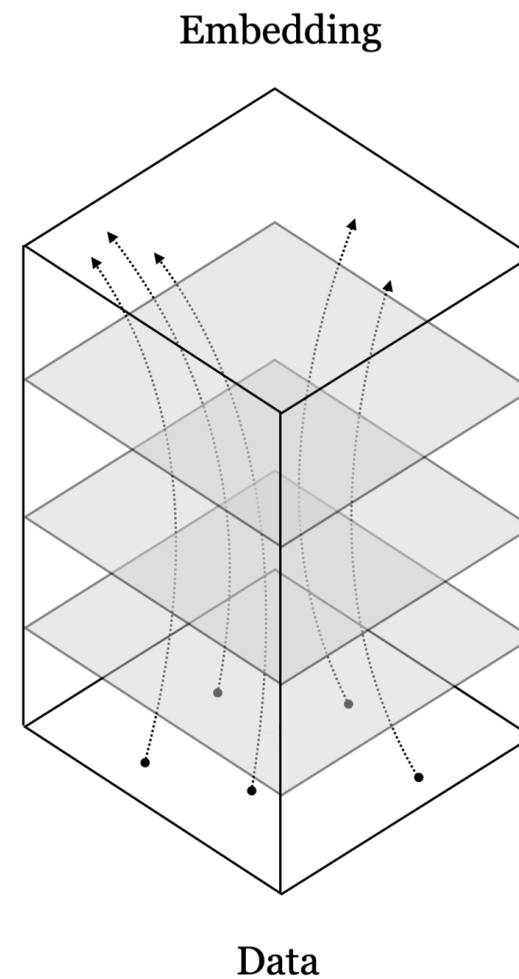
linear



Deep neural nets transform datapoints, layer by layer

Deep neural nets transform datapoints, layer by layer

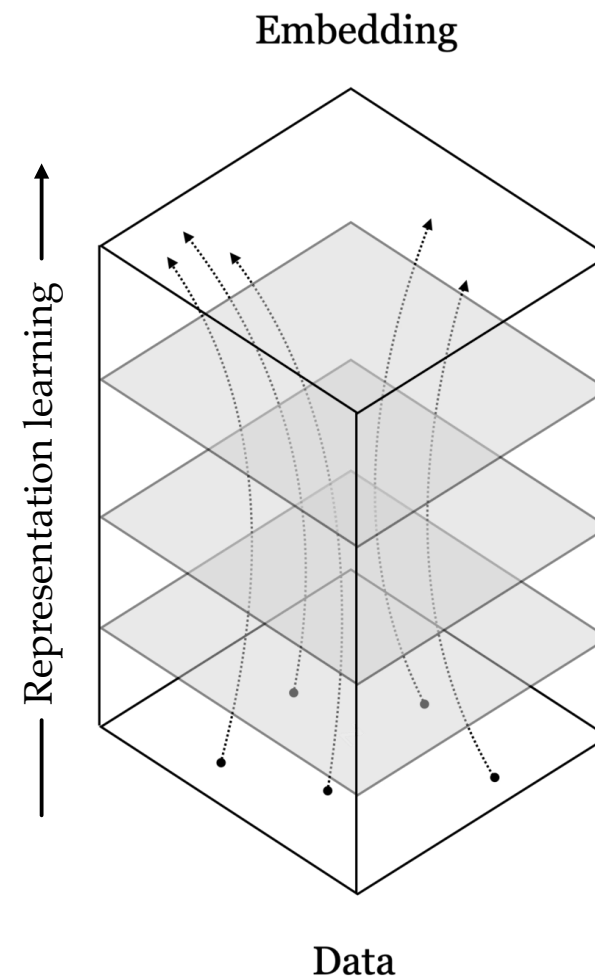
Each layer is a different *representation*, aka *embedding*, of the data



Deep neural nets transform datapoints, layer by layer

Each layer is a different *representation*, aka *embedding*, of the data

From data to latent embeddings: representation learning

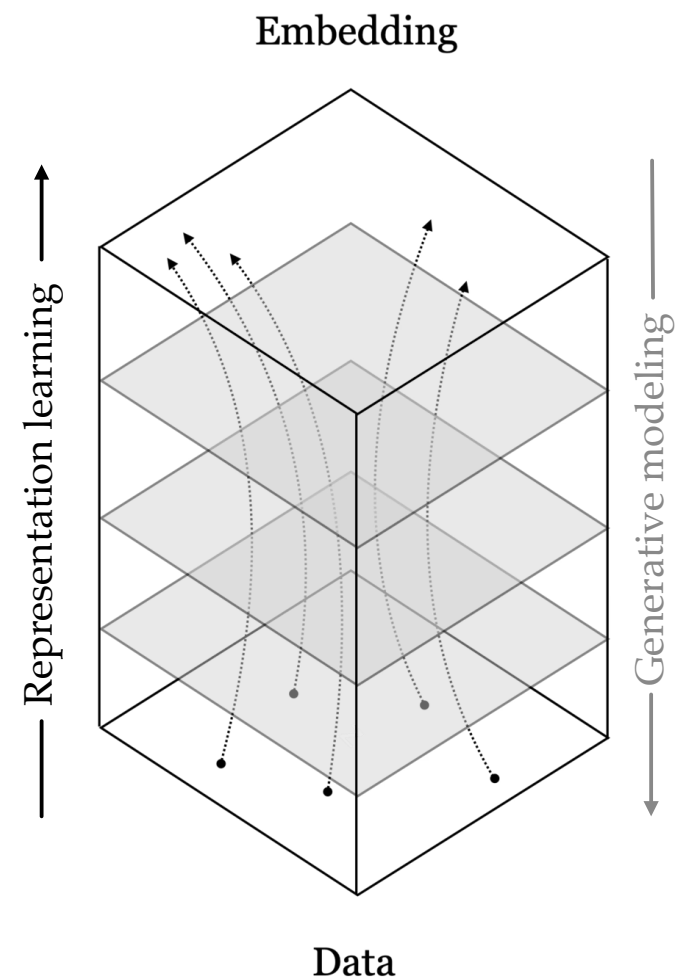


Deep neural nets transform datapoints, layer by layer

Each layer is a different *representation*, aka *embedding*, of the data

From data to latent embeddings: representation learning

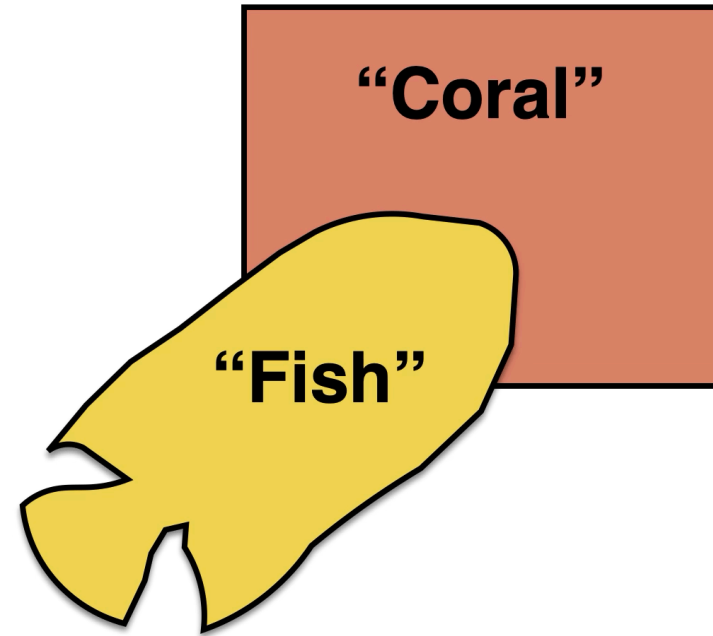
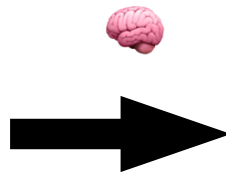
(From latent embeddings to data: generative modeling)



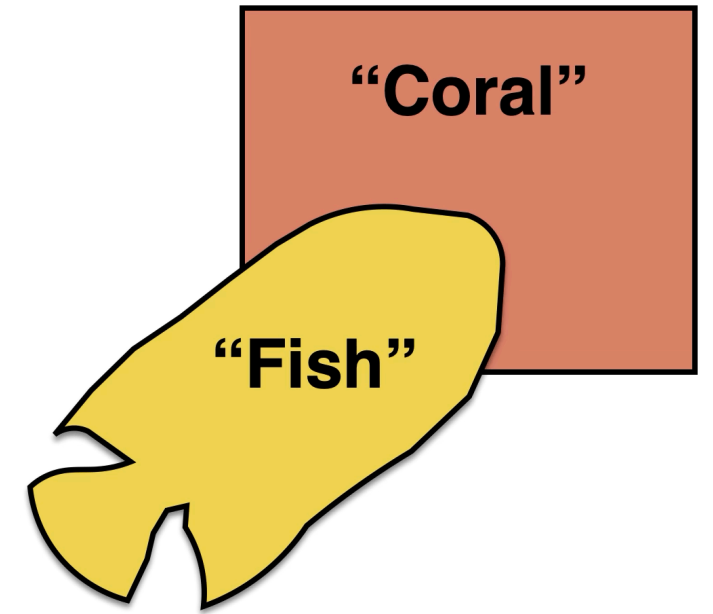
humans also learn representations



humans also learn representations



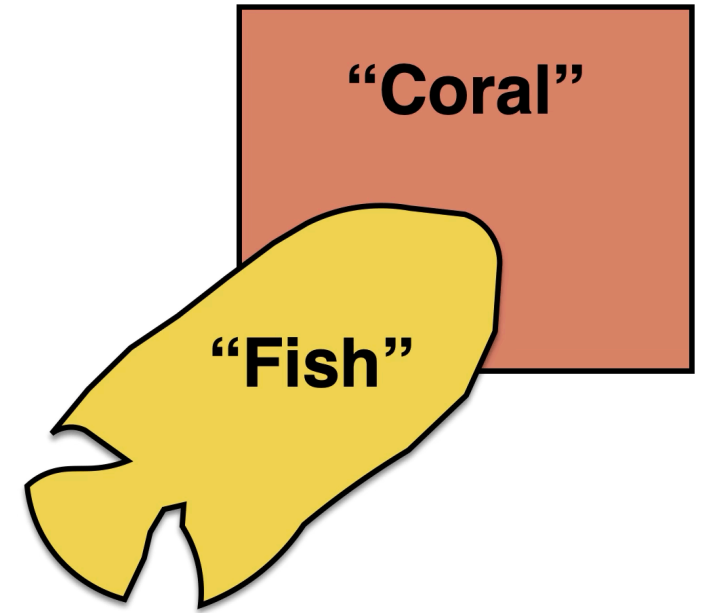
Good representations are:



[See "Representation Learning", Bengio 2013, for more commentary]

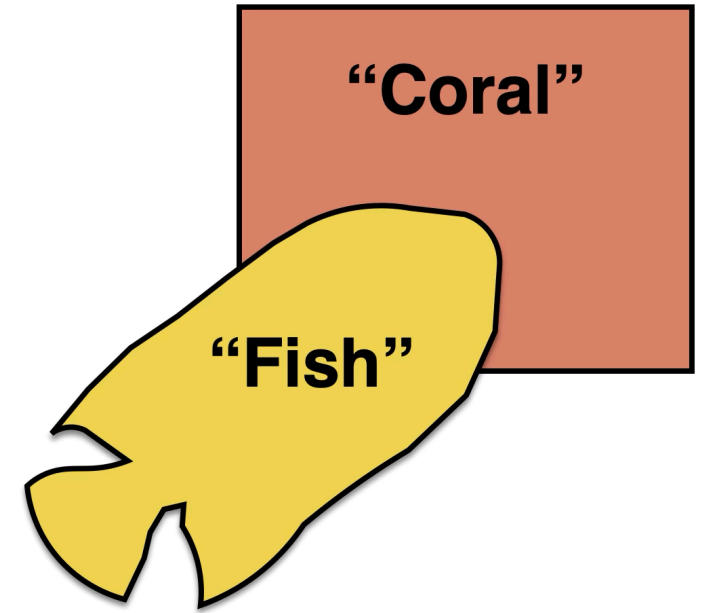
Good representations are:

- Compact (*minimal*)



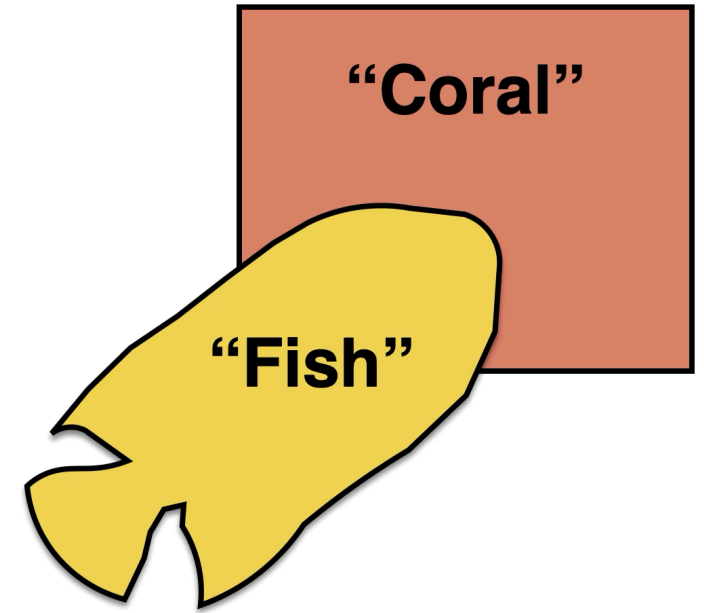
Good representations are:

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)



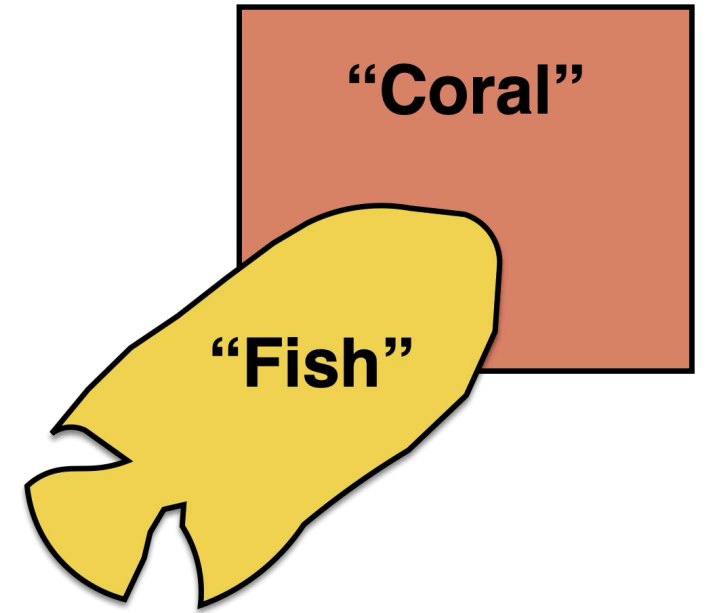
Good representations are:

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)
- Disentangled (*independent* factors)



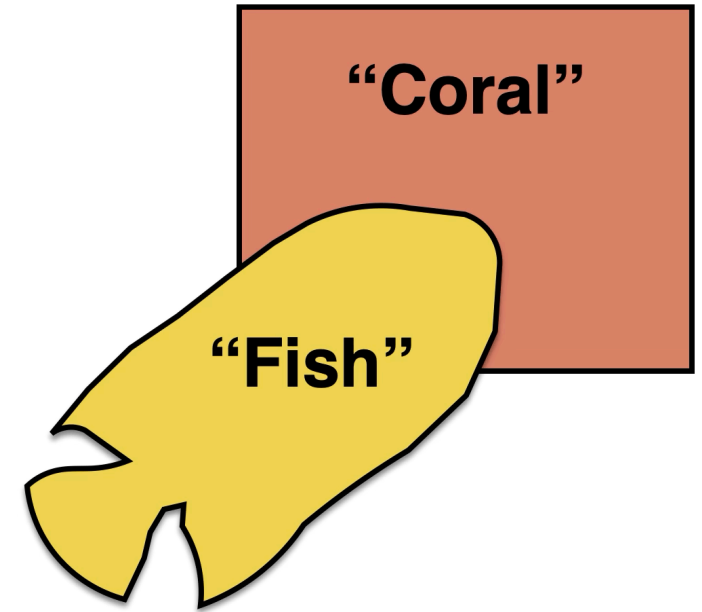
Good representations are:

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)
- Disentangled (*independent* factors)
- Interpretable (understandable by humans)



Good representations are:

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)
- Disentangled (*independent* factors)
- Interpretable (understandable by humans)
- Make *subsequent problem solving easy*







ImageNet also taught us that labeling 14 million images by hand is brutal.

ImageNet also taught us that labeling 14 million images by hand is brutal.

Undergrads were **time-consuming**, algorithms were flawed, and the team didn't have **money**—Li said the project failed to win any of the federal grants she applied for, receiving comments on proposals that it was shameful Princeton would research this topic, and that the only strength of proposal was that Li was a woman.

A solution finally surfaced in a chance hallway conversation with a graduate student who asked Li whether she had heard of Amazon Mechanical Turk, a service where hordes of humans sitting at computers around the world would complete small online tasks for pennies.

“He showed me the website, and I can tell you literally that day I knew the ImageNet project was going to happen,” she said.

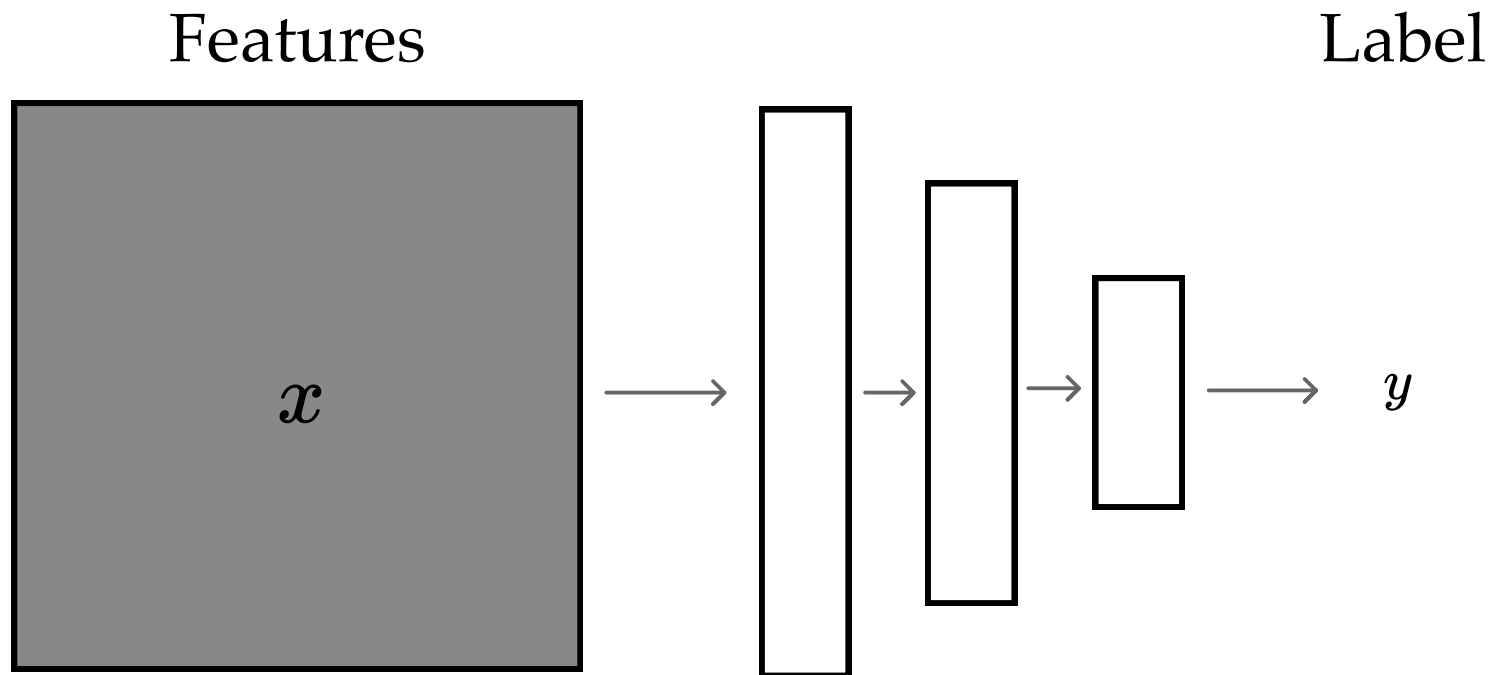
“Suddenly we found a tool that **could scale**, that we could not possibly dream of by hiring Princeton undergrads.”



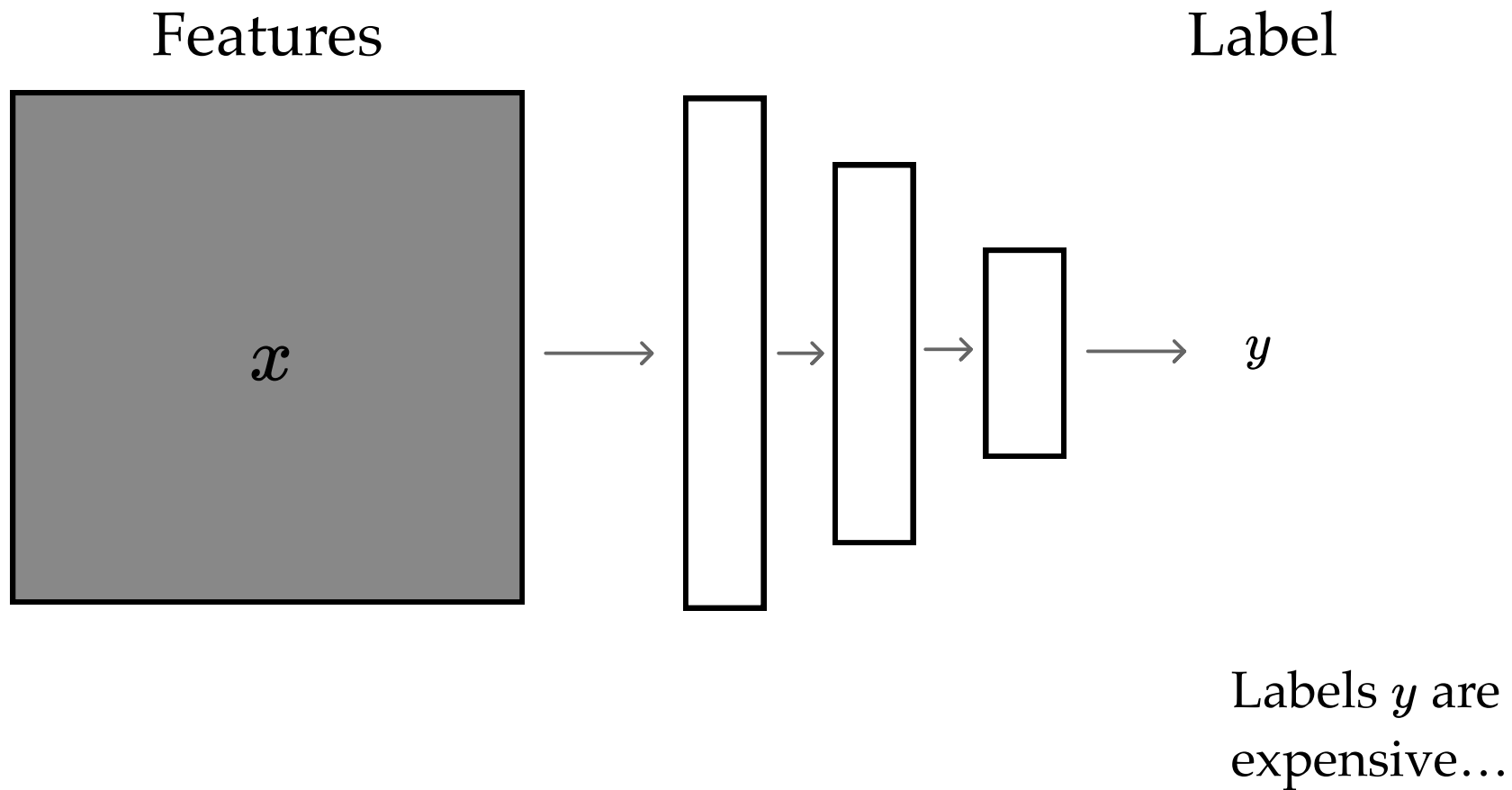
# Outline

- Neural networks are *representation* learners
- **Self-supervised learning**
- Auto-encoders
- Word embeddings
- (Some recent representation learning ideas)

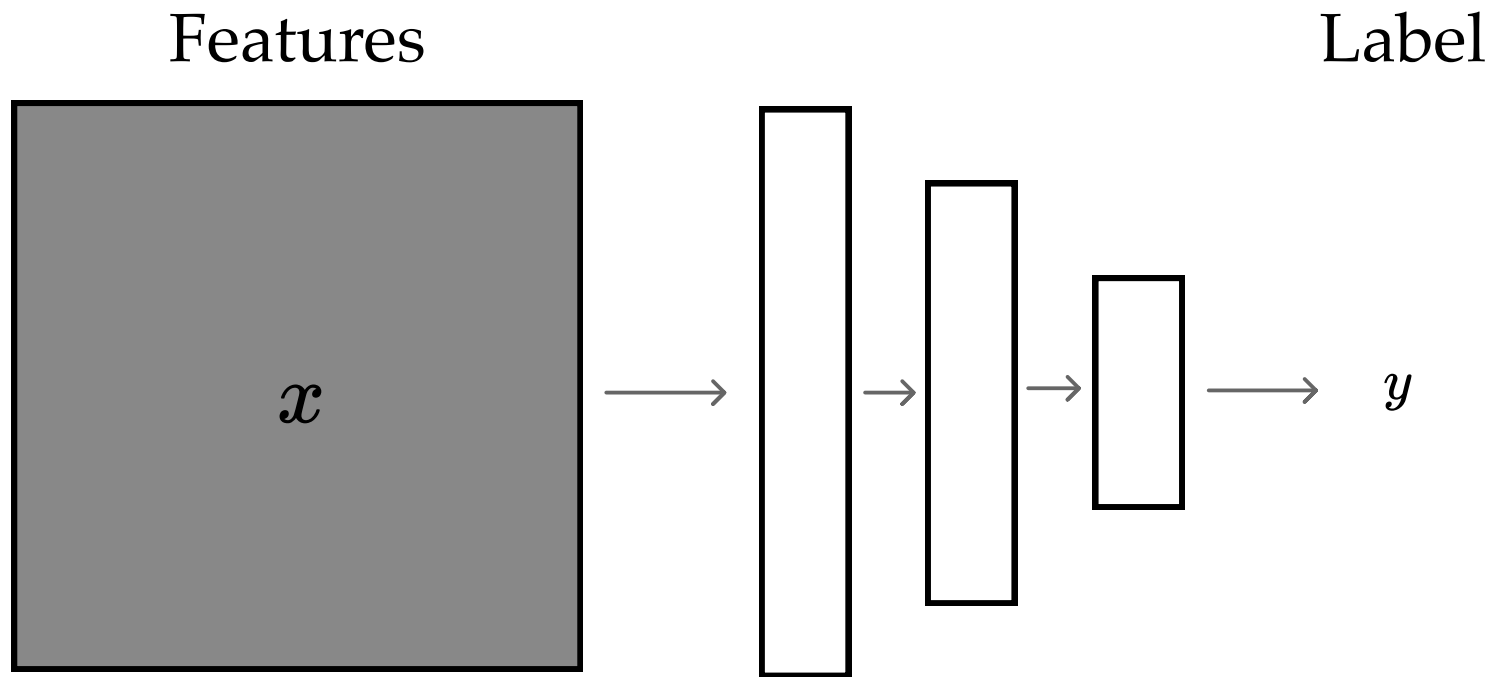
# Label prediction (supervised learning)



# Label prediction (supervised learning)



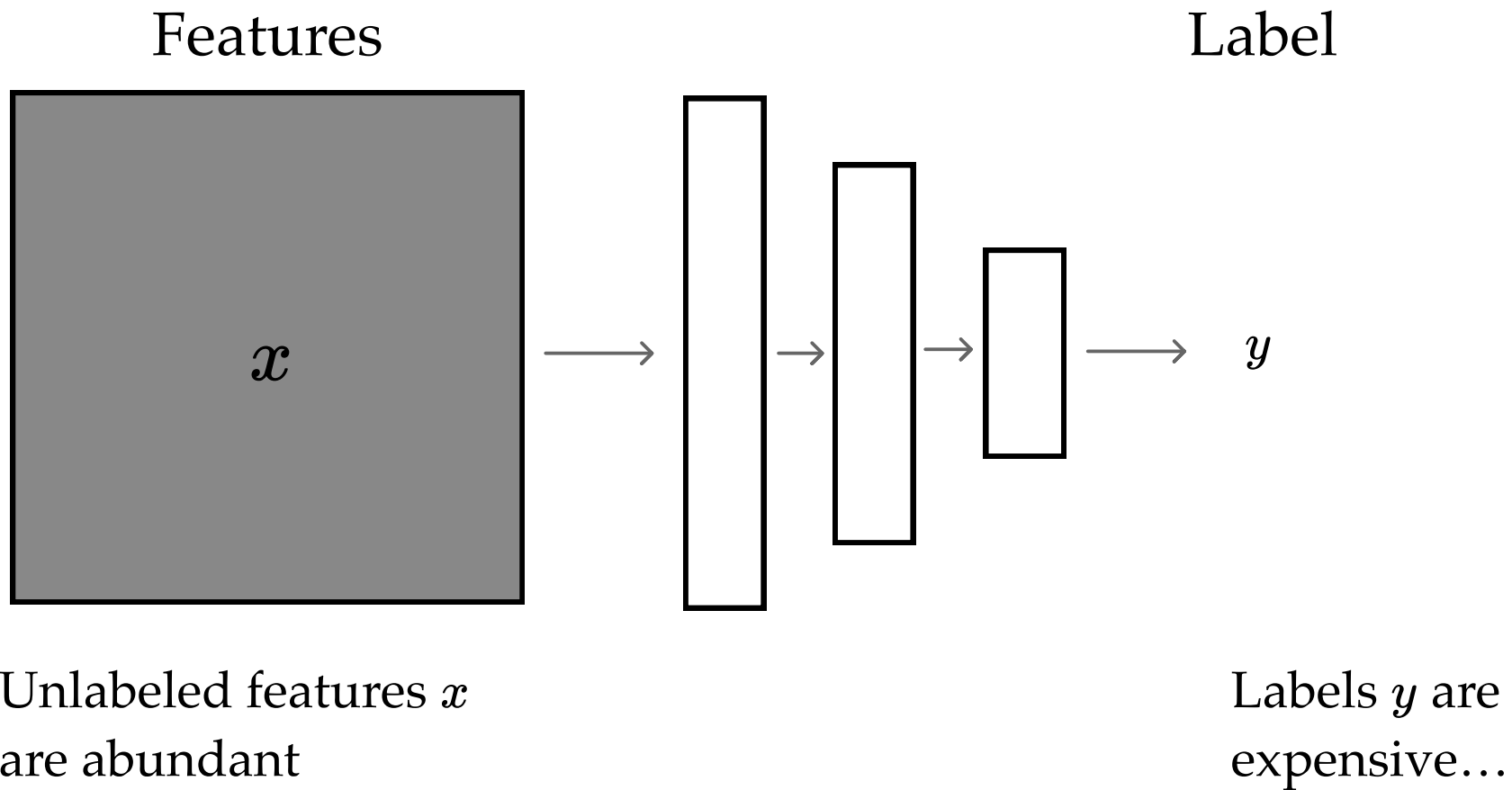
# Label prediction (supervised learning)



Unlabeled features  $x$   
are abundant

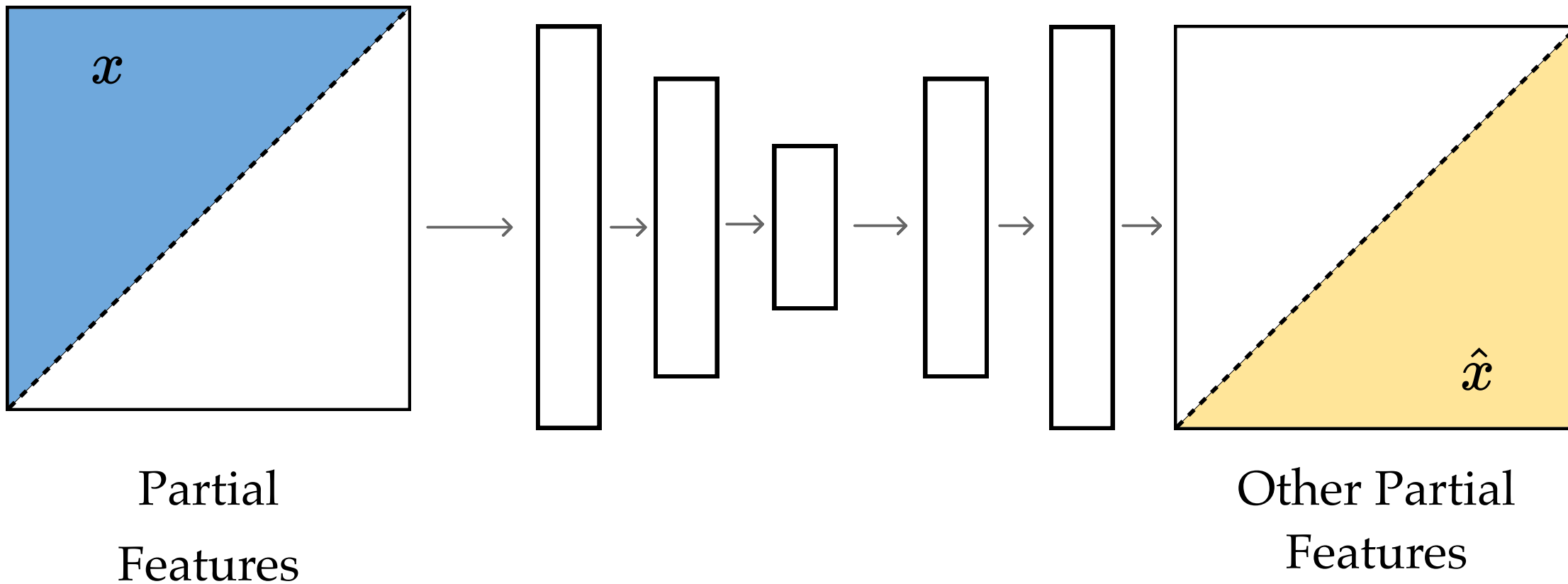
Labels  $y$  are  
expensive...

# Label prediction (supervised learning)



Can we learn useful things with just  $x$ ?

# Self-supervised Learning (partial feature reconstruction)



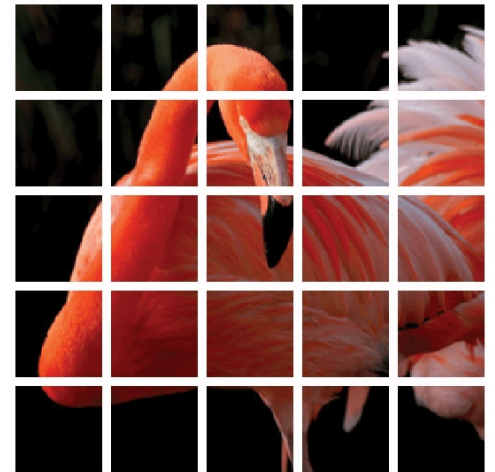
Hard reconstruction, forces understanding.

# Masked Auto-Encoder (Vision)

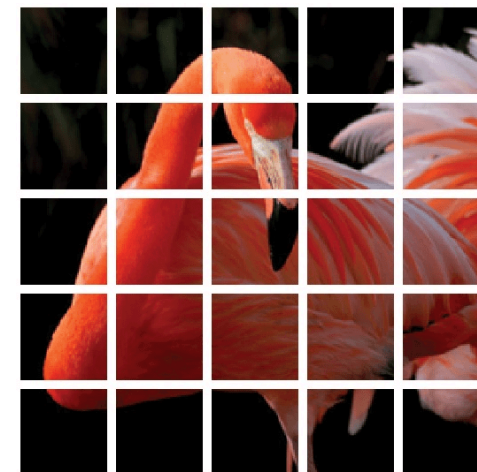
# Masked Auto-Encoder (Vision)



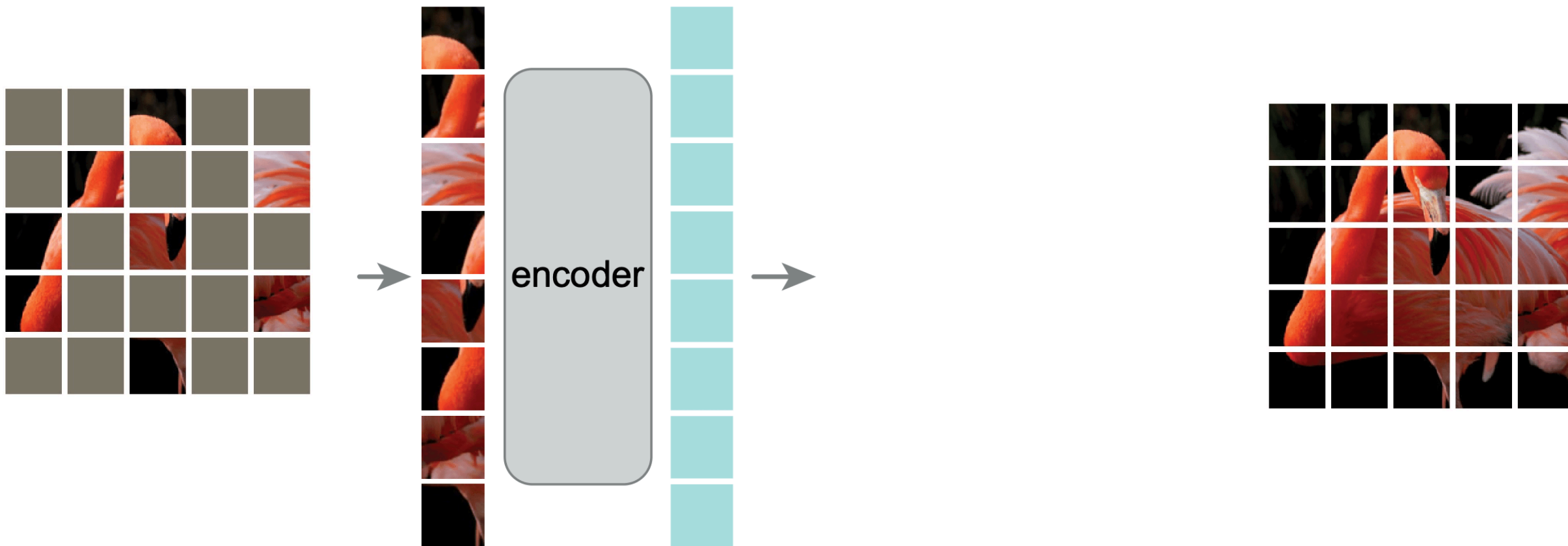
# Masked Auto-Encoder (Vision)



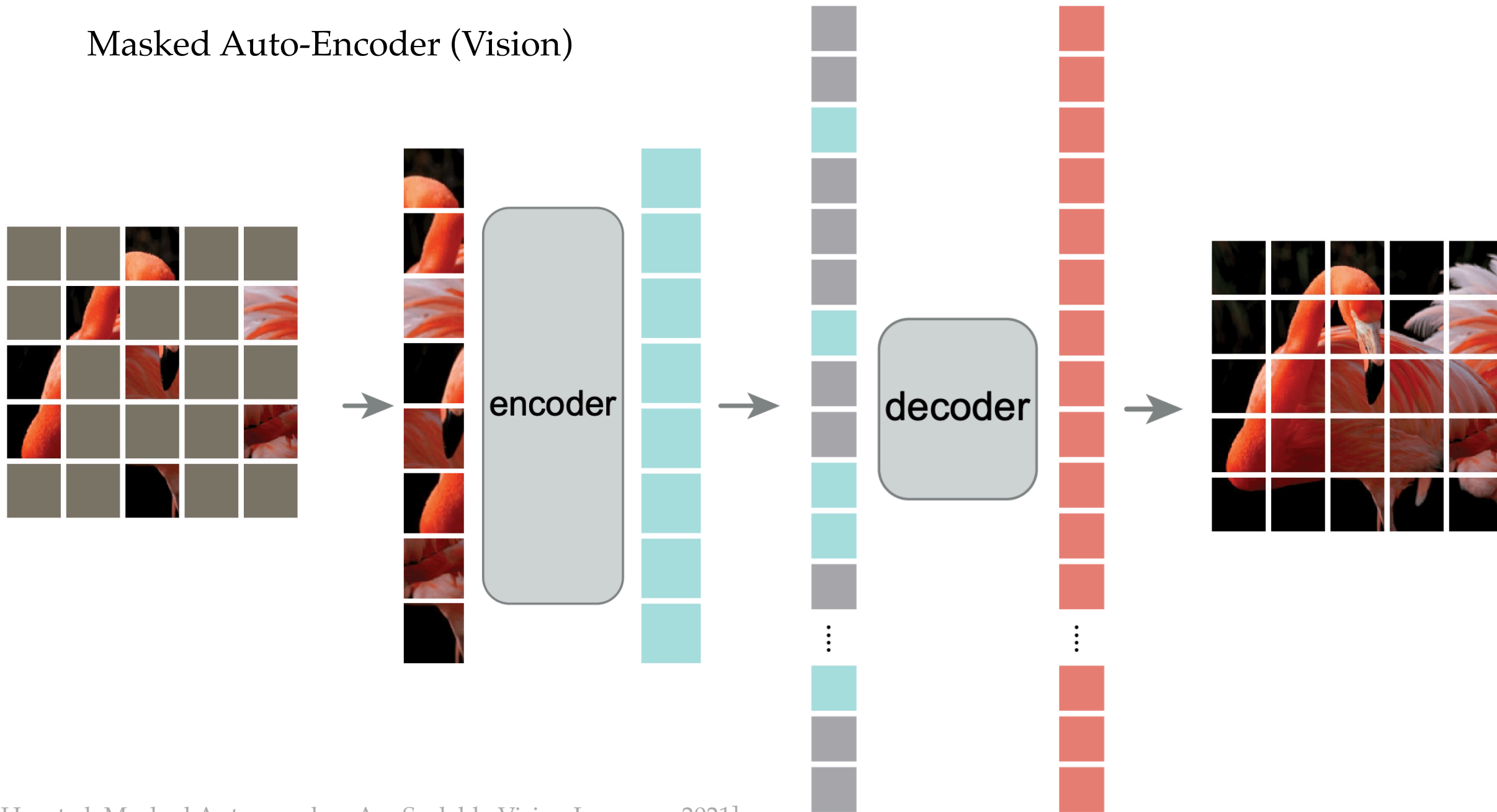
# Masked Auto-Encoder (Vision)



# Masked Auto-Encoder (Vision)



# Masked Auto-Encoder (Vision)



## Masking channels

e.g. masking the color

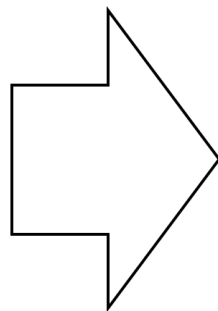
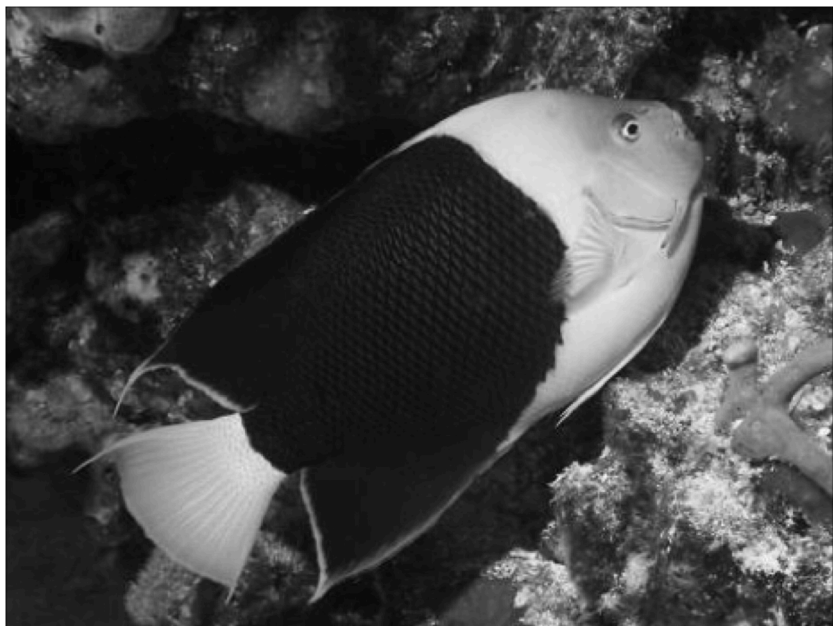


## Masking channels



# Masking channels

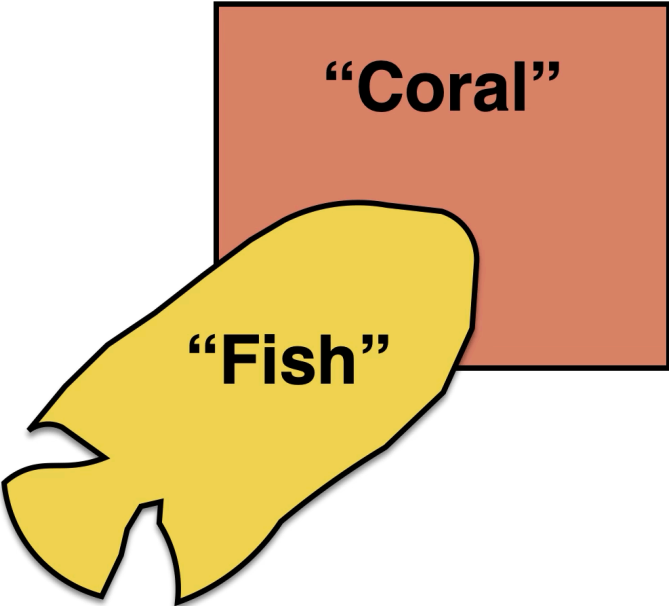
predict color from gray-scale



# Outline

- Neural networks are *representation* learners
- Self-supervised learning
- **Auto-encoders**
- Word embeddings
- (Some recent representation learning ideas)

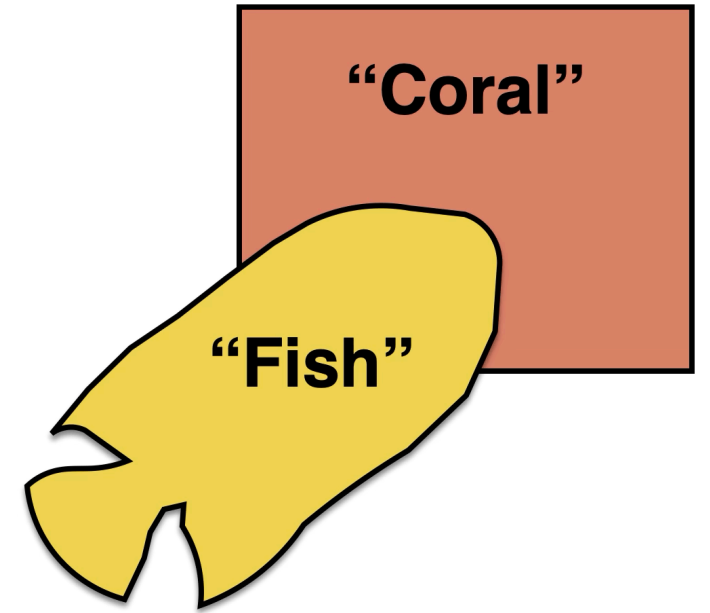
Good representations are:



[See "Representation Learning", Bengio 2013, for more commentary]

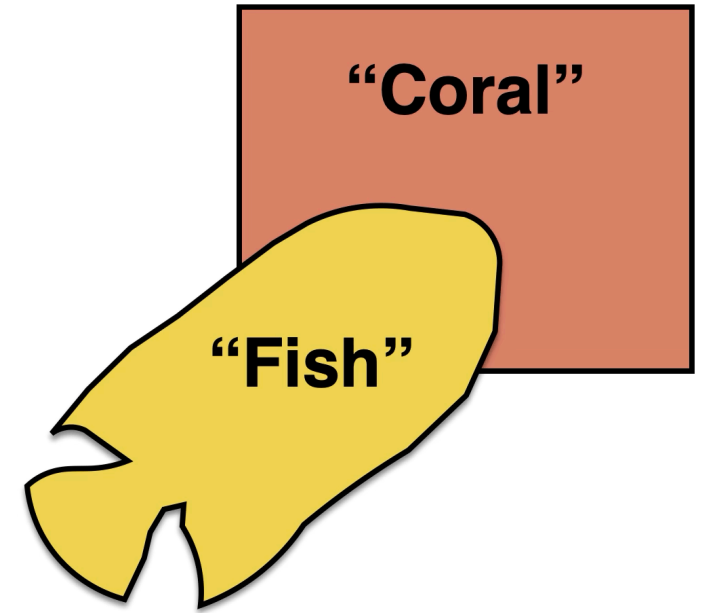
Good representations are:

- Compact (*minimal*)



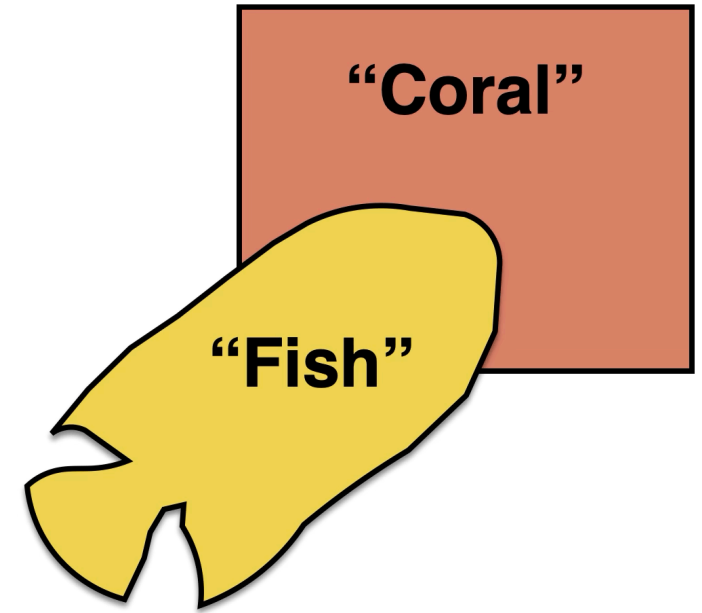
Good representations are:

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)



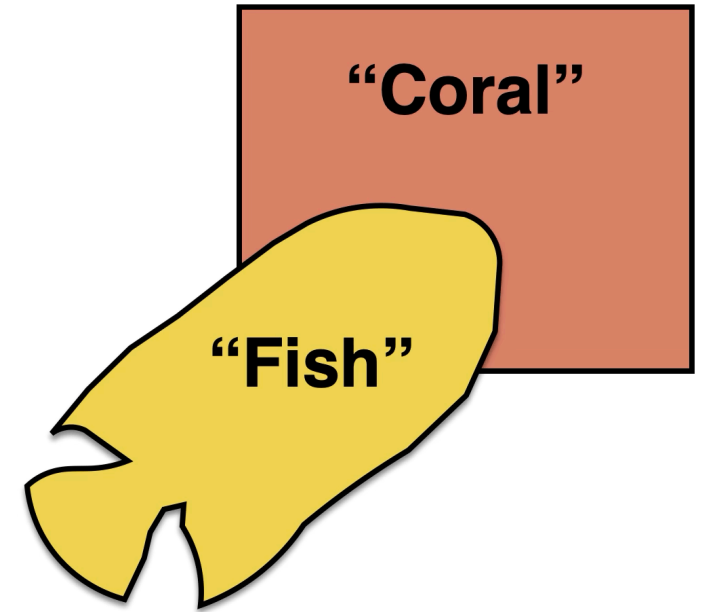
Good representations are:

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)
- Disentangled (*independent* factors)



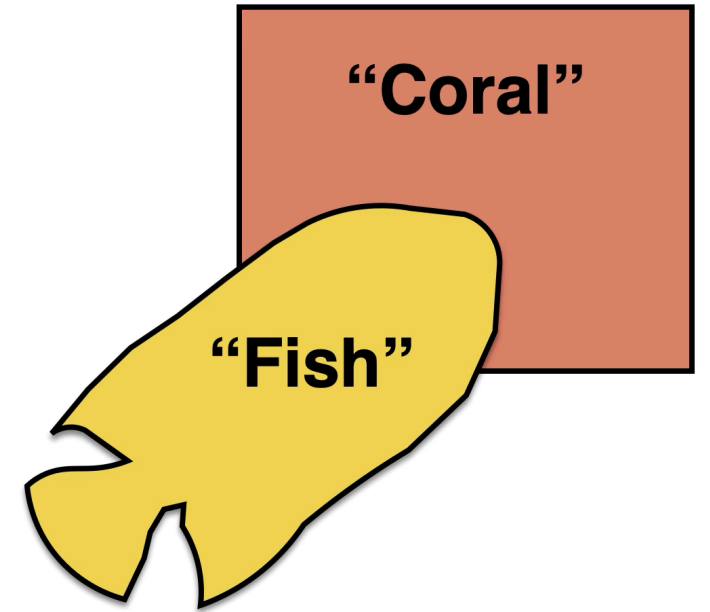
Good representations are:

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)
- Disentangled (*independent* factors)
- Interpretable (understandable by humans)



Good representations are:

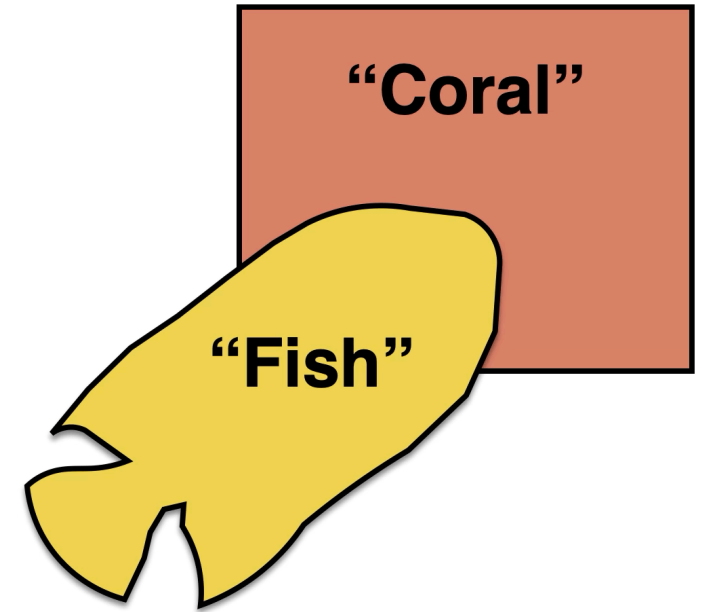
- Compact (*minimal*)
- Explanatory (*roughly sufficient*)
- Disentangled (*independent* factors)
- Interpretable (understandable by humans)
- Make *subsequent problem solving easy*



Good representations are:

Auto-encoders  
explicitly aims

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)
- Disentangled (*independent* factors)
- Interpretable (understandable by humans)
- Make *subsequent problem solving easy*

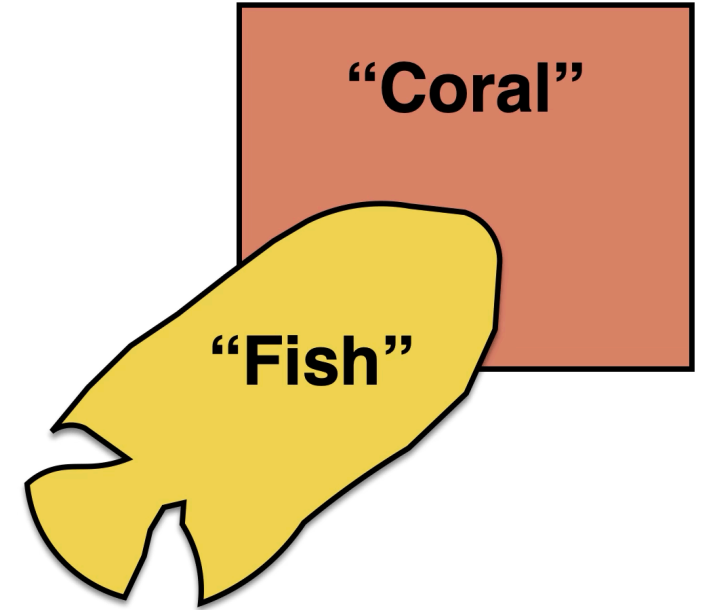


Good representations are:

Auto-encoders  
explicitly aims

these may just  
emerge as well

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)
- Disentangled (*independent* factors)
- Interpretable (understandable by humans)
- Make *subsequent problem solving easy*



Observed image



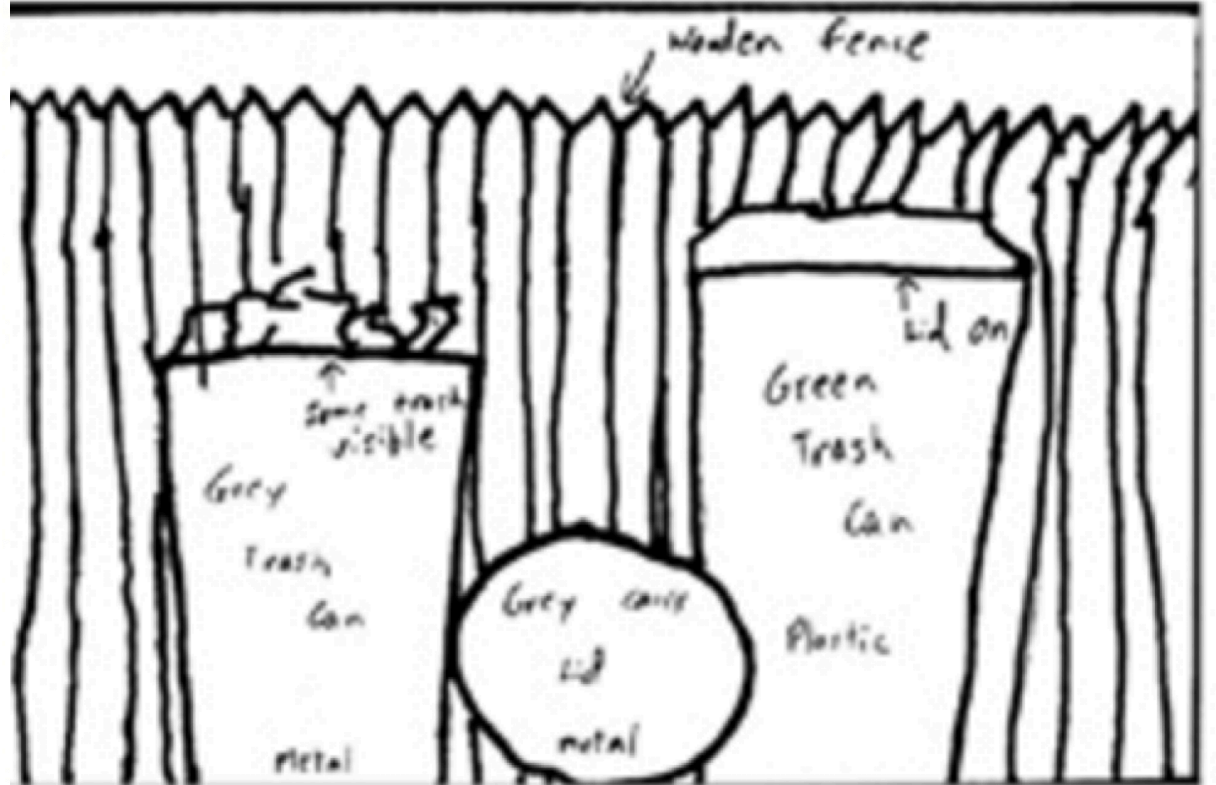
[Bartlett, 1932]

[Intraub & Richardson, 1989]

Observed image



Drawn from memory



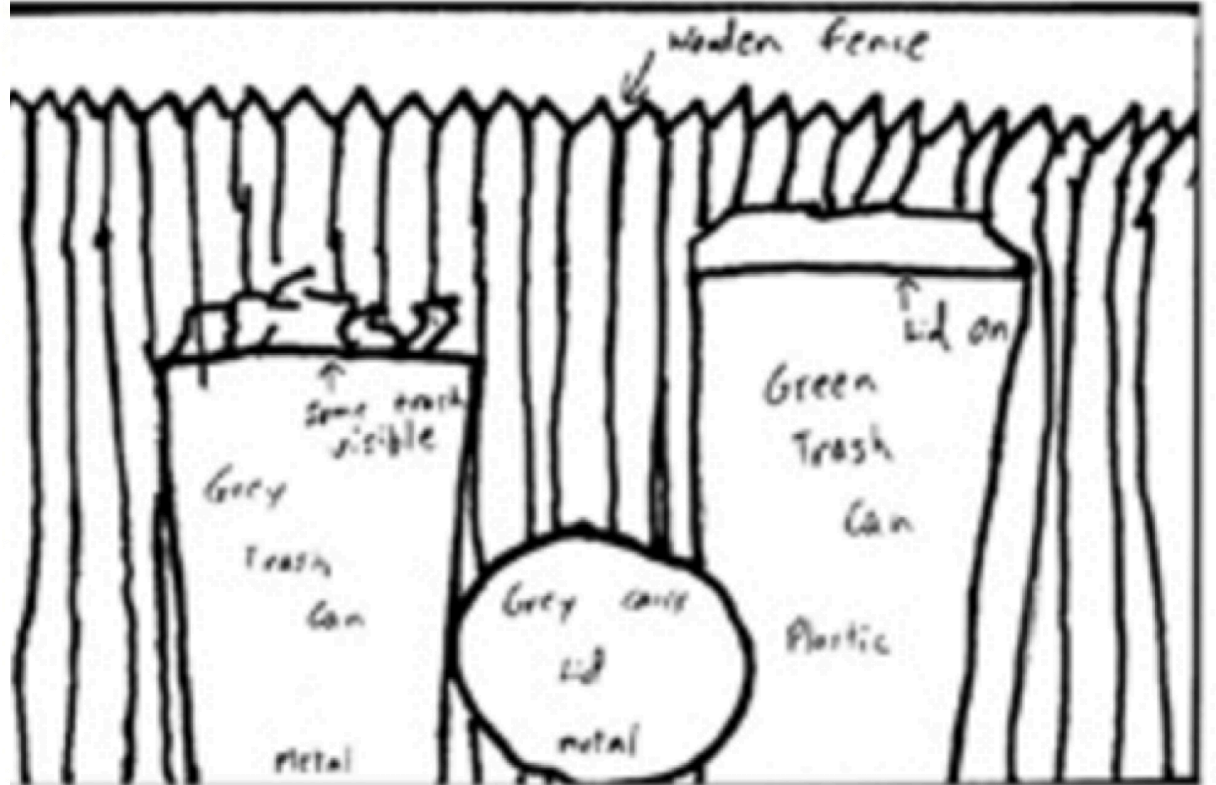
[Bartlett, 1932]

[Intraub & Richardson, 1989]

Observed image



Drawn from memory

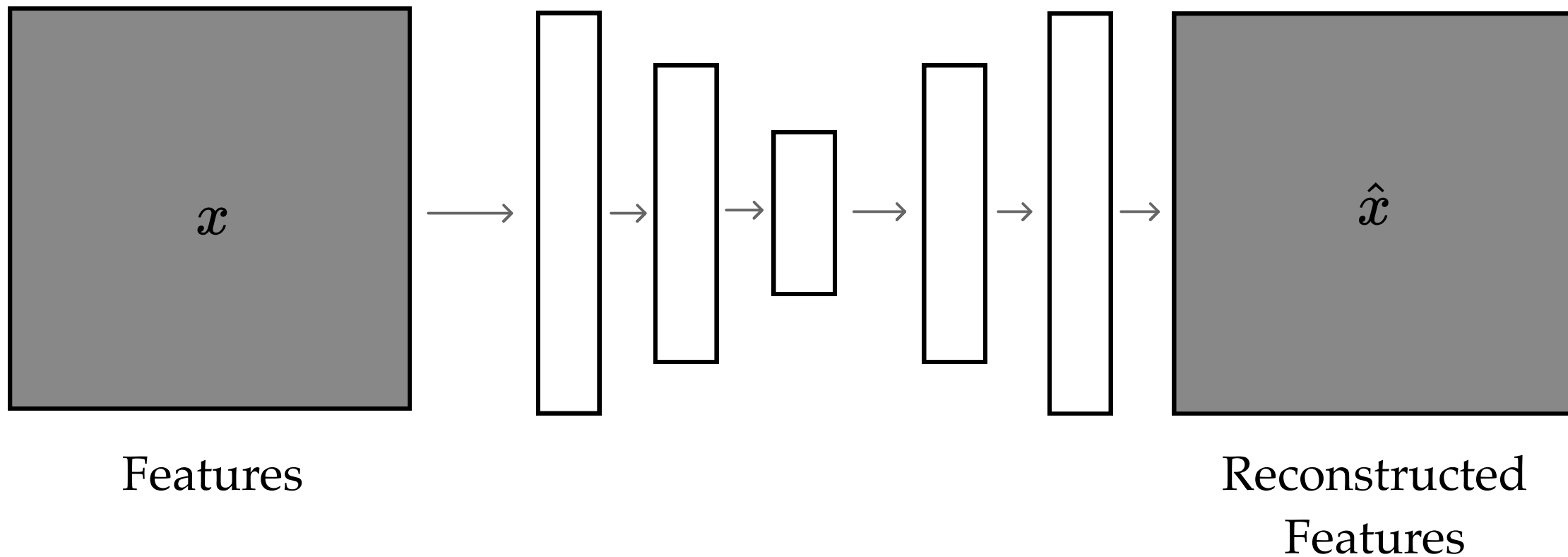


[Bartlett, 1932]

[Intraub & Richardson, 1989]



# Unsupervised Learning (feature reconstruction)



# Auto-encoder

Auto-encoder



Image

# Auto-encoder



Image

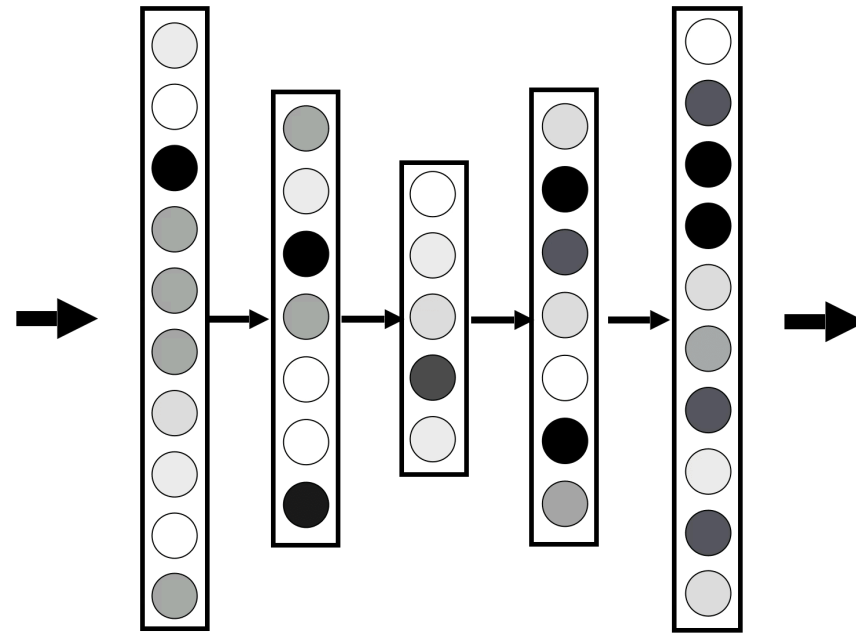


Reconstructed  
image

# Auto-encoder



Image

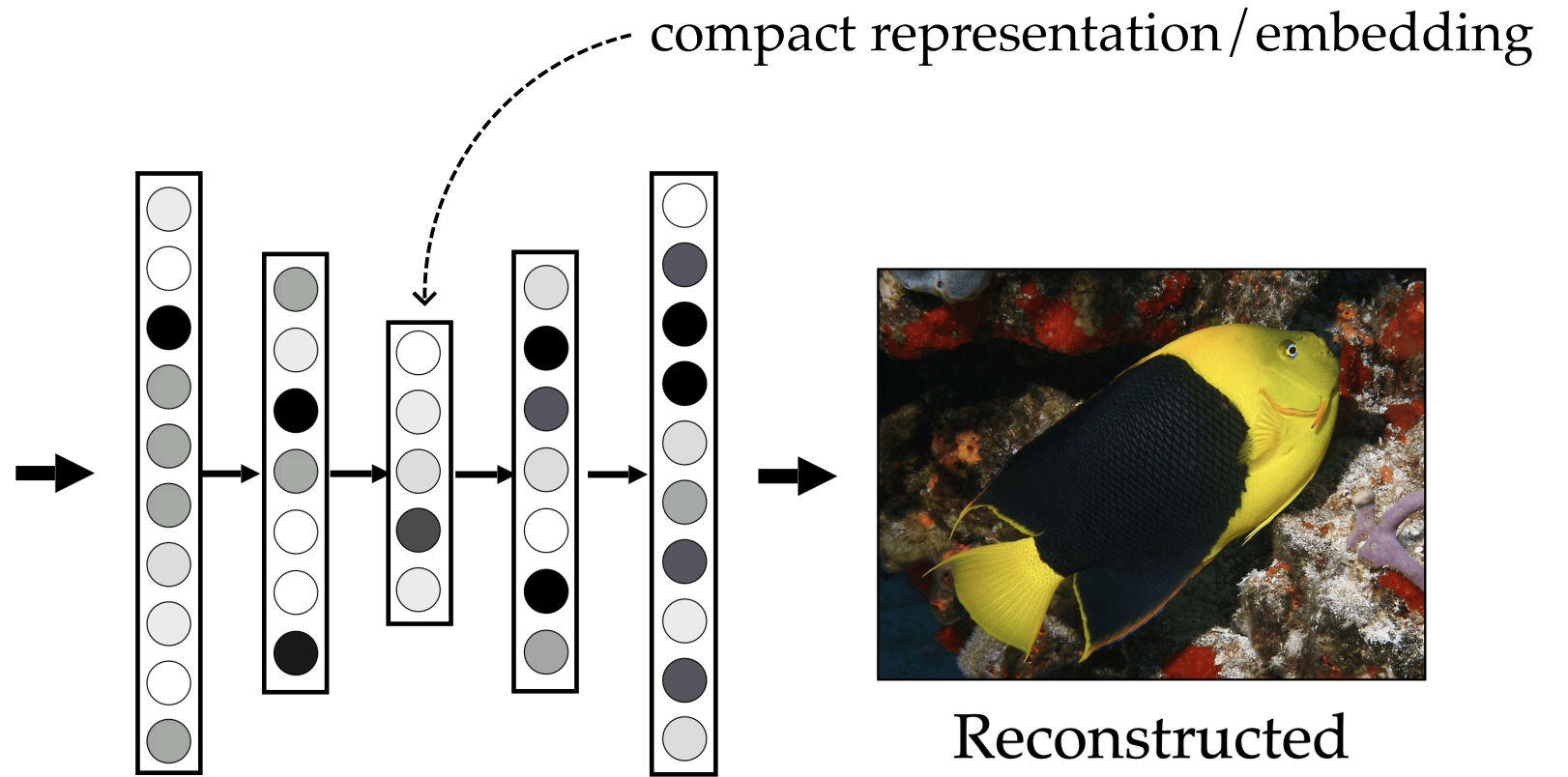


Reconstructed  
image

# Auto-encoder



Image

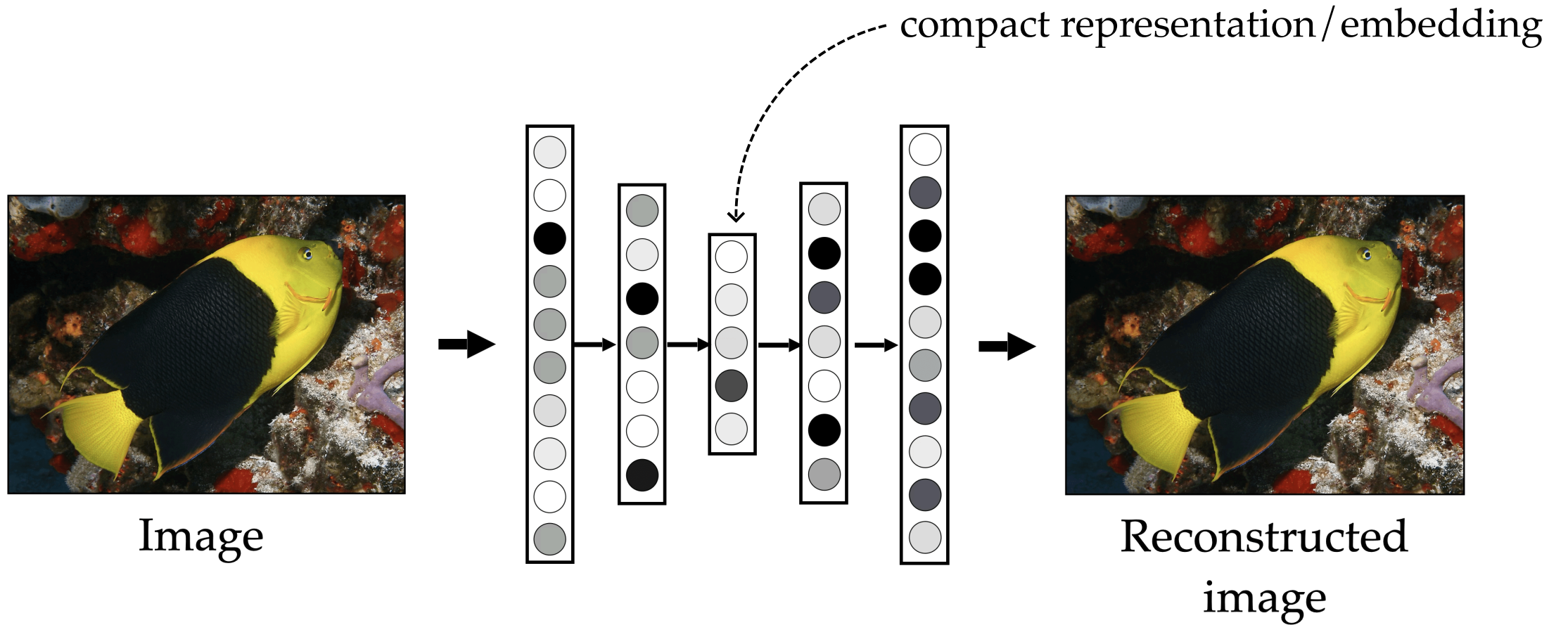


compact representation / embedding



Reconstructed  
image

# Auto-encoder



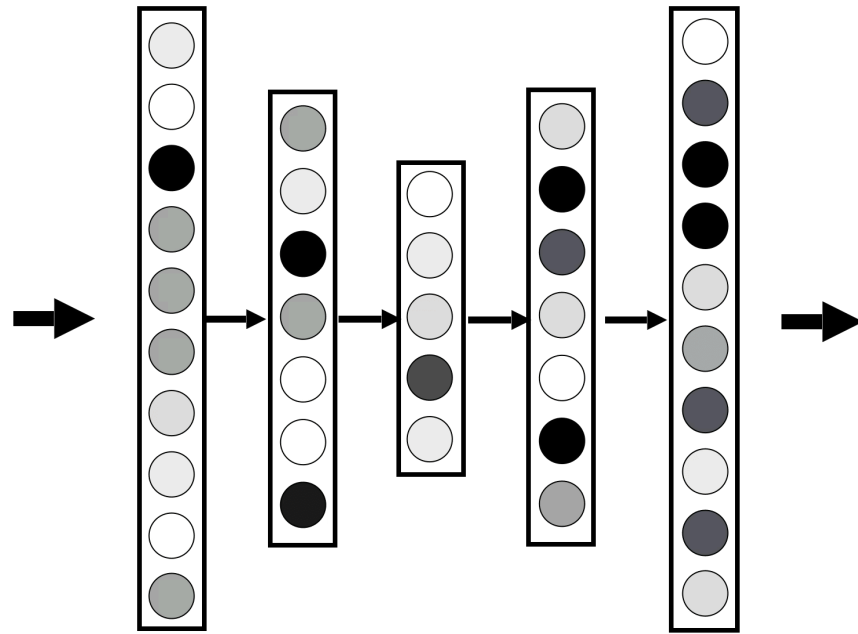
*"What I cannot create, I do not understand."* Feynman

# Auto-encoder

$x$



Image



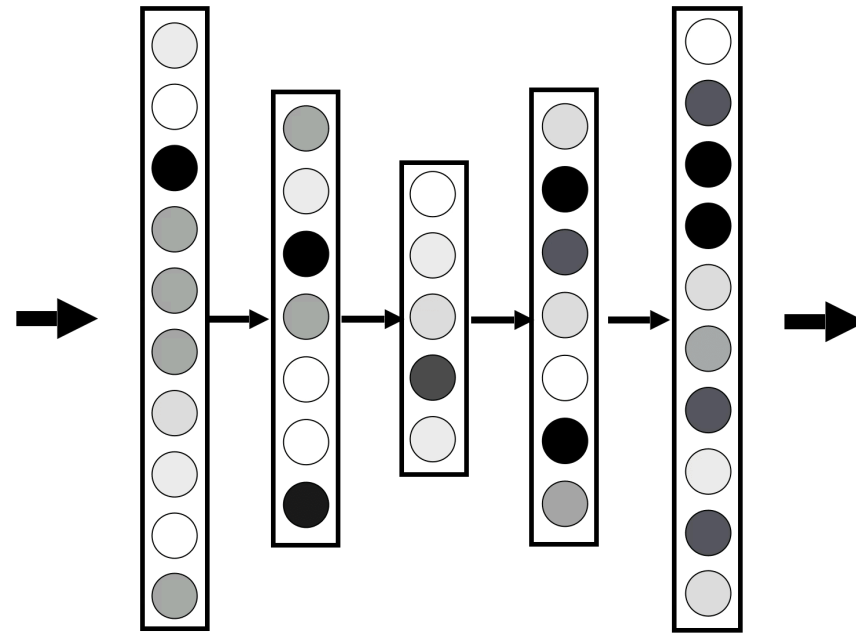
Reconstructed  
image

# Auto-encoder

$x$



Image

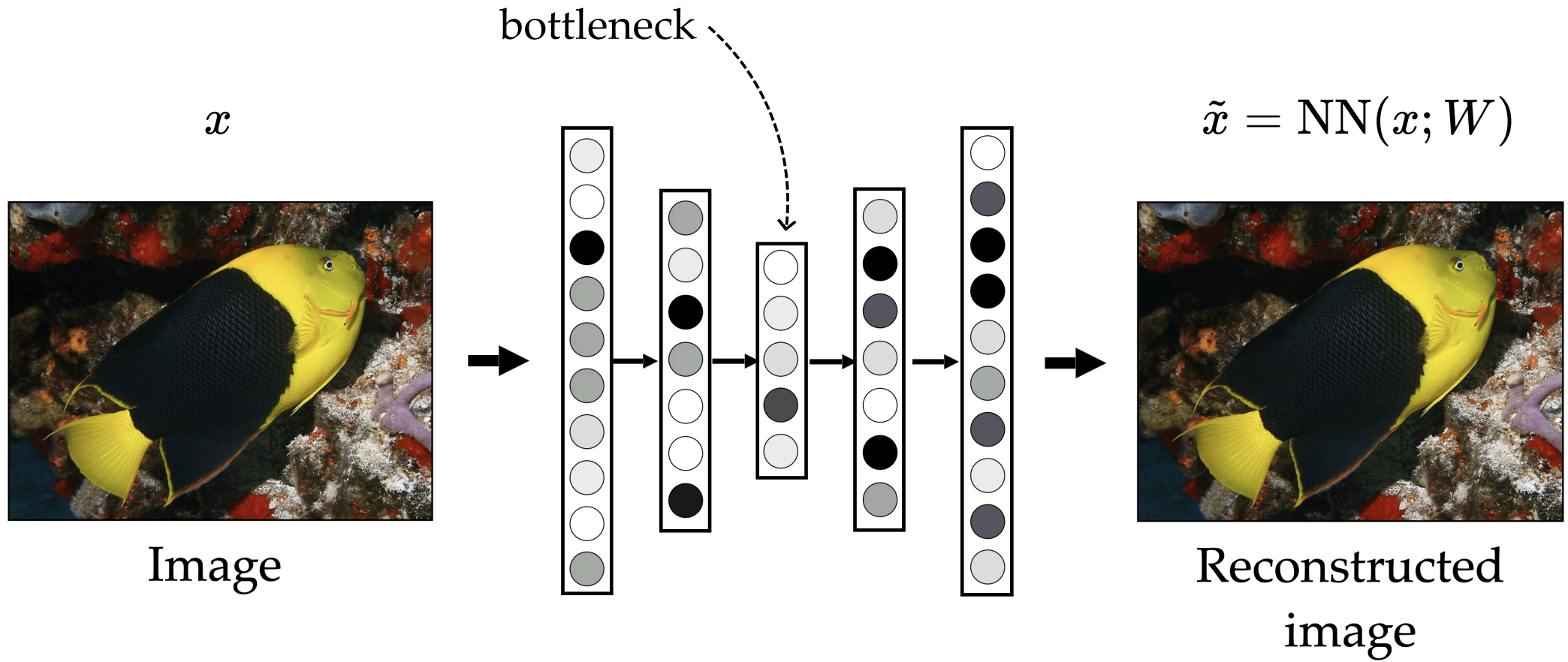


$$\tilde{x} = \text{NN}(x; W)$$



Reconstructed  
image

# Auto-encoder



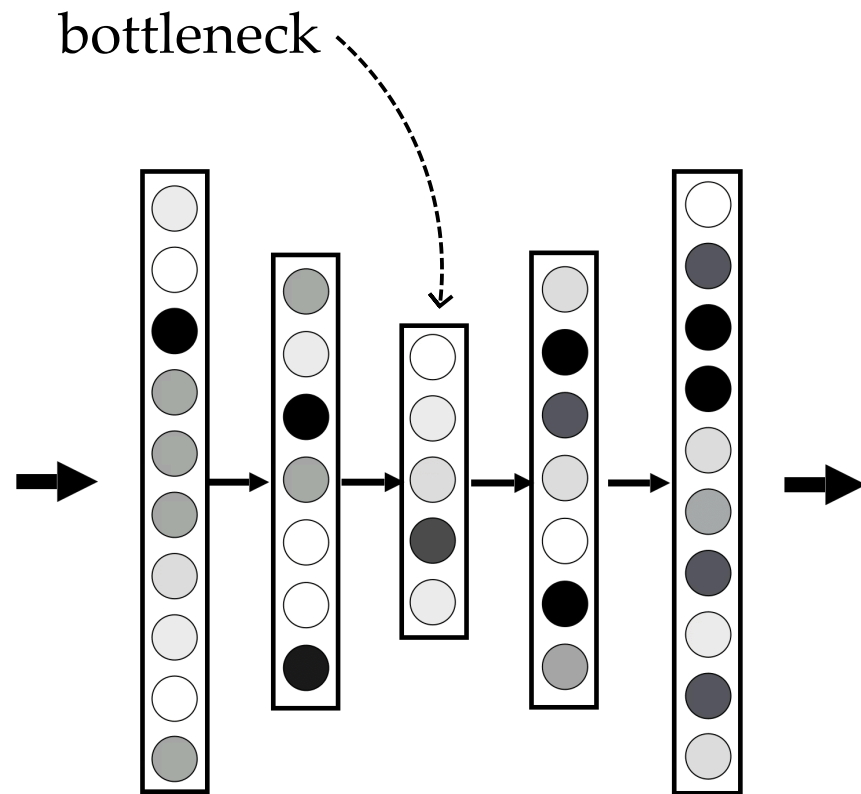
Auto-encoder

$$\min_W ||x - \tilde{x}||^2$$

$x$



Image



$$\tilde{x} = \text{NN}(x; W)$$



Reconstructed  
image

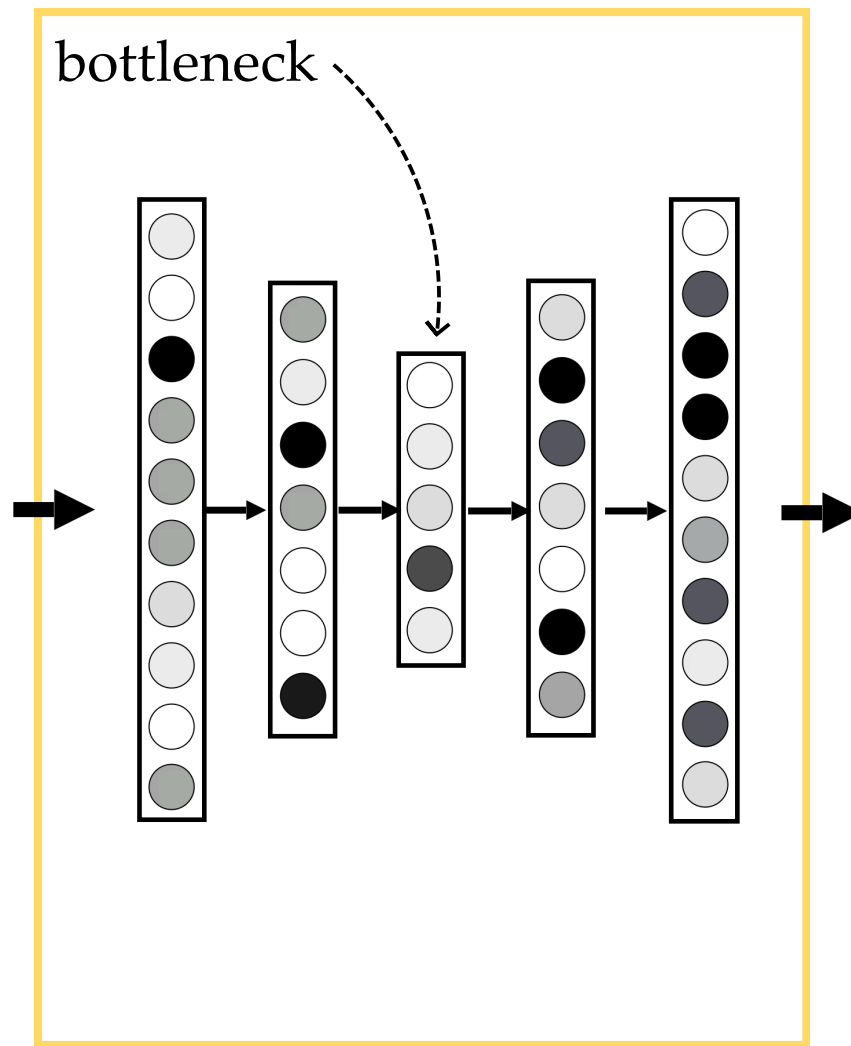
# Auto-encoder

$$\min_W ||x - \tilde{x}||^2$$

$x$



Image



$$\tilde{x} = \text{NN}(x; W)$$



Reconstructed  
image

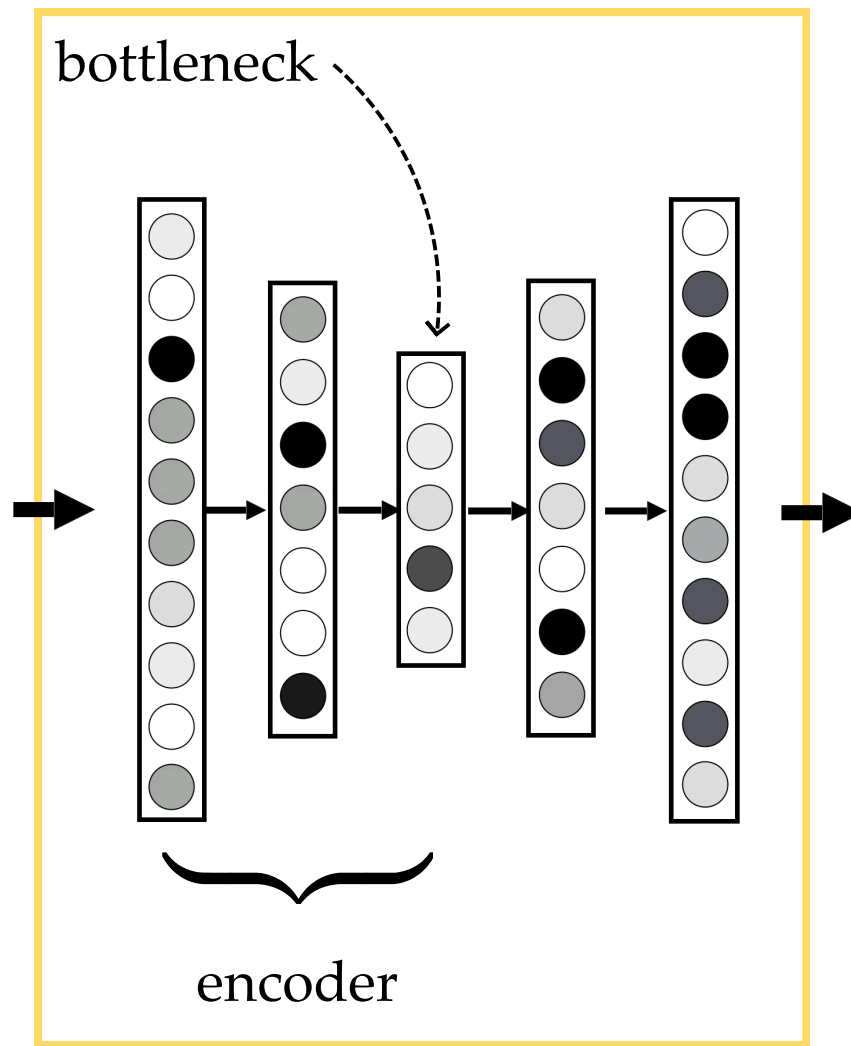
# Auto-encoder

$$\min_W ||x - \tilde{x}||^2$$



Image

$x$



$$\tilde{x} = \text{NN}(x; W)$$



Reconstructed image

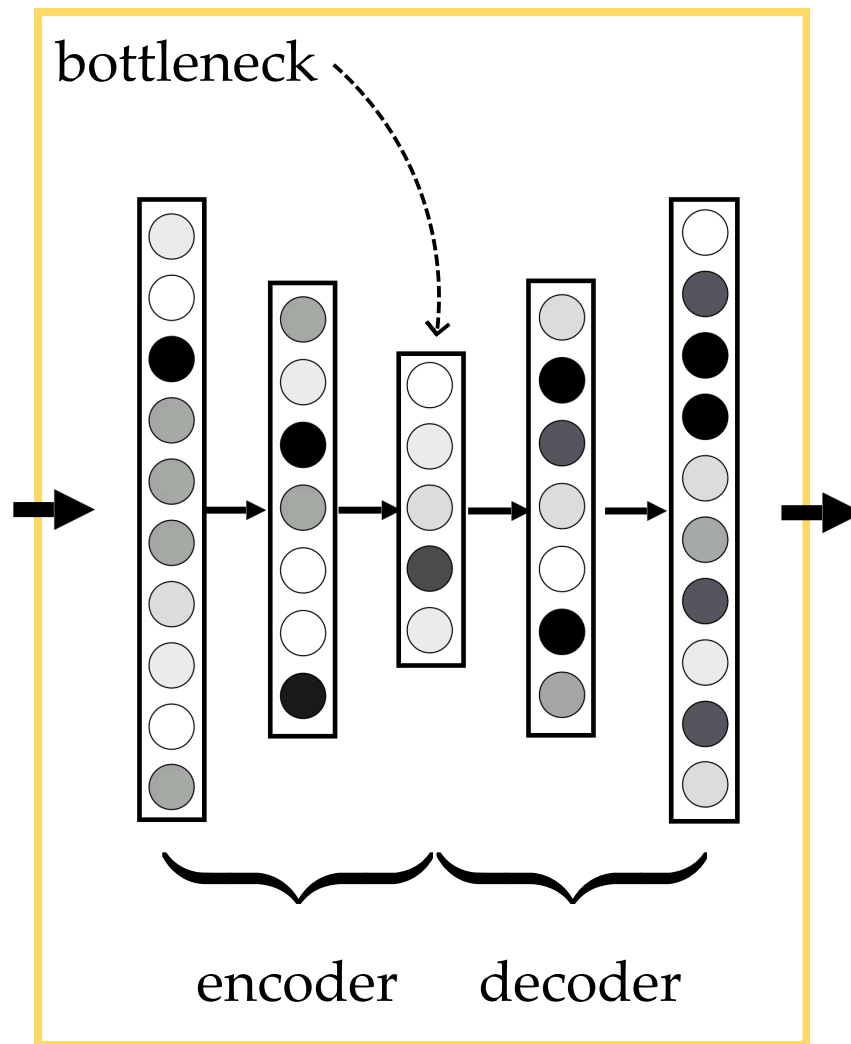
# Auto-encoder

$$\min_W ||x - \tilde{x}||^2$$

$x$



Image



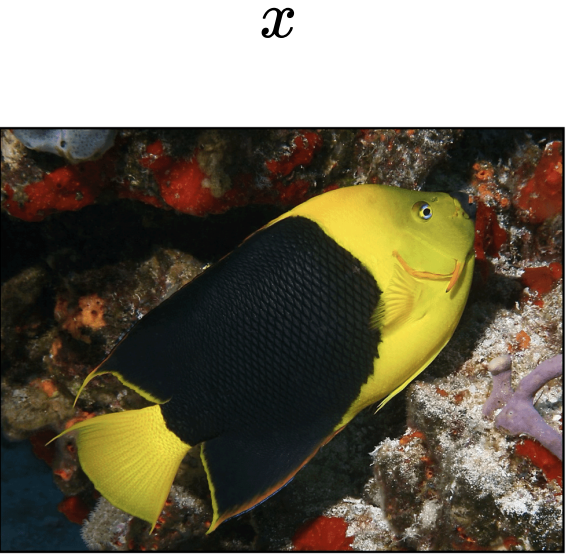
$$\tilde{x} = \text{NN}(x; W)$$



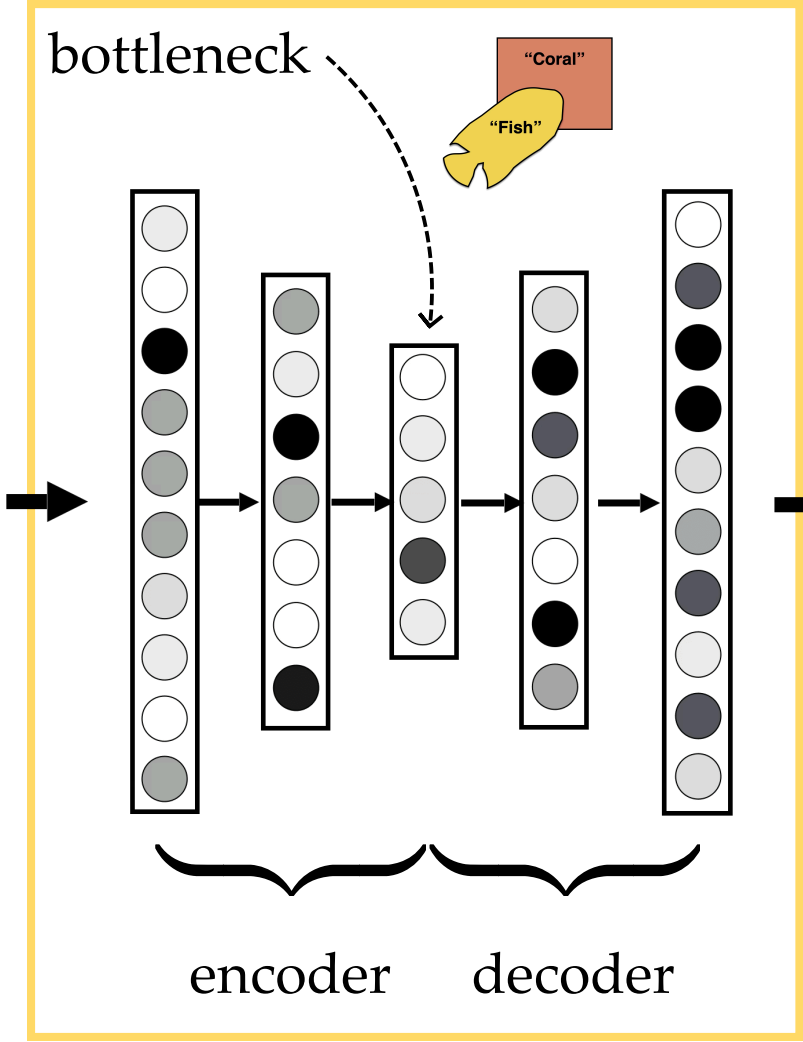
Reconstructed image

Auto-encoder

$$\min_W ||x - \tilde{x}||^2$$



Image



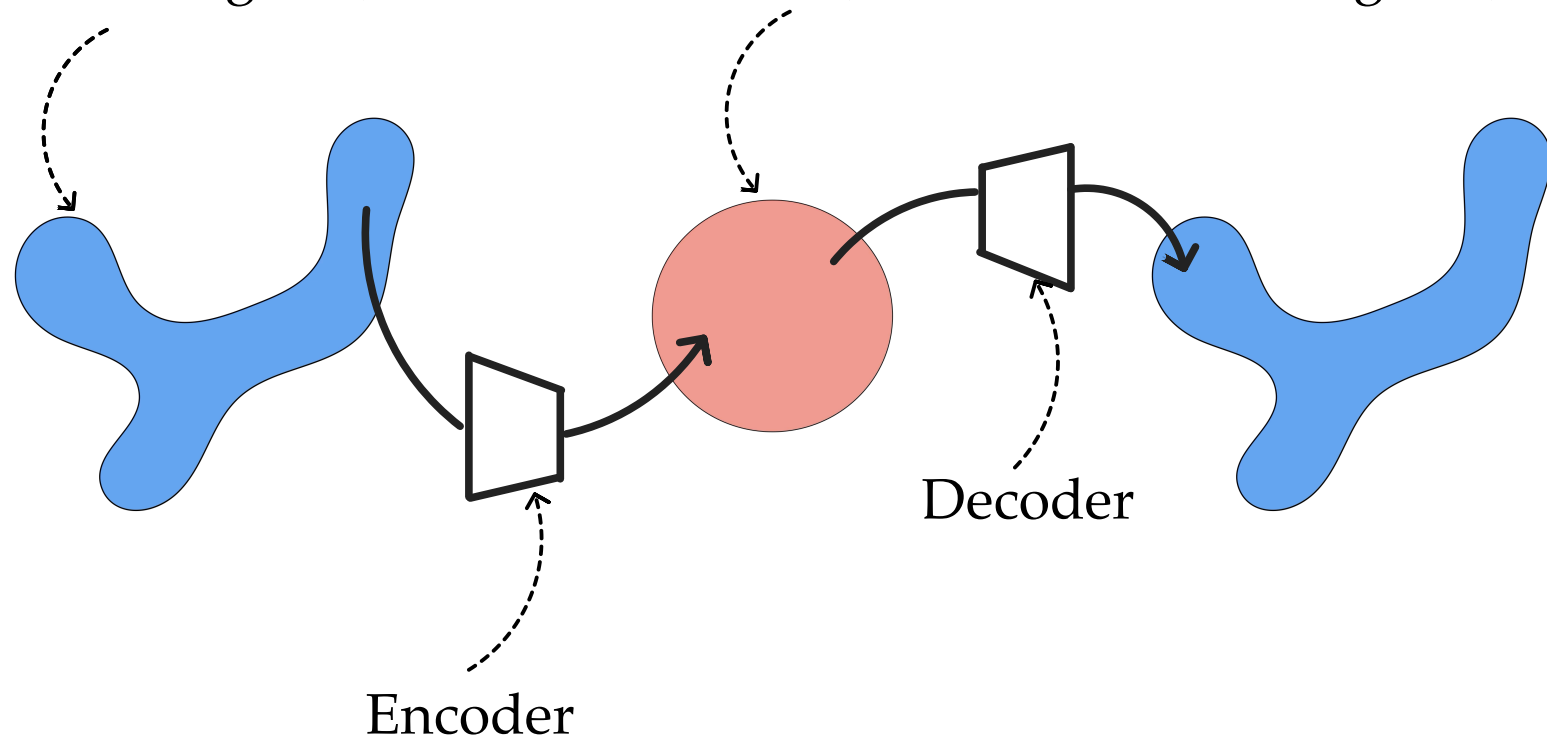
$$\tilde{x} = \text{NN}(x; W)$$



Reconstructed image

Data space:  
(high dimensional, irregular)

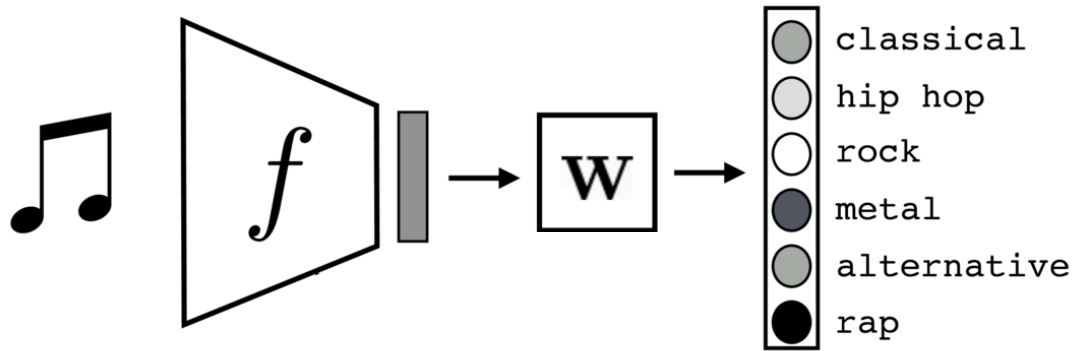
Representation space  
(low dimensional, regular)



[Typically, encoders can serve as a translator to get "good representations",  
whereas decoders can serve as "generative models"]

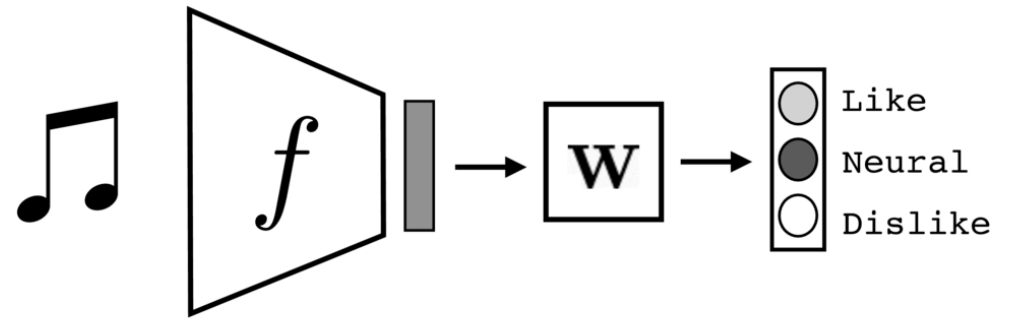
# Training

Genre recognition



# Testing

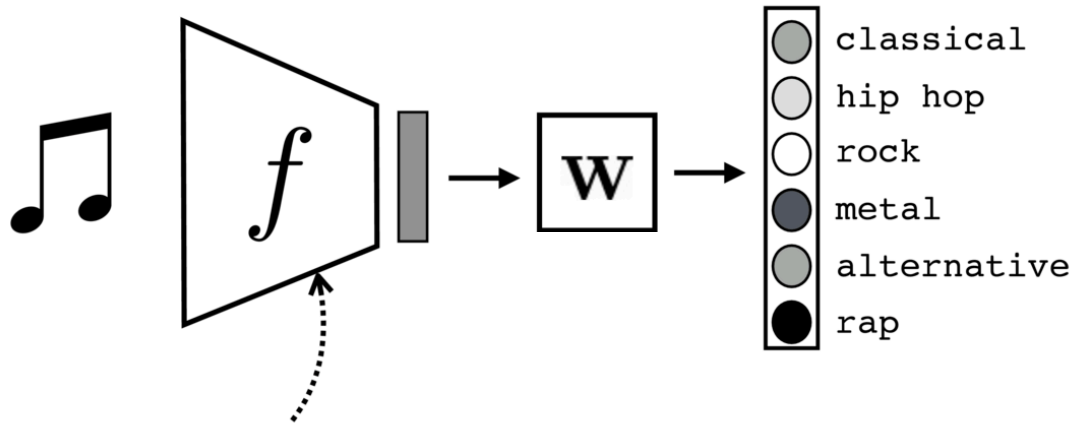
Preference prediction



Often, what we will be “tested” on is not what we were trained on.

# Training

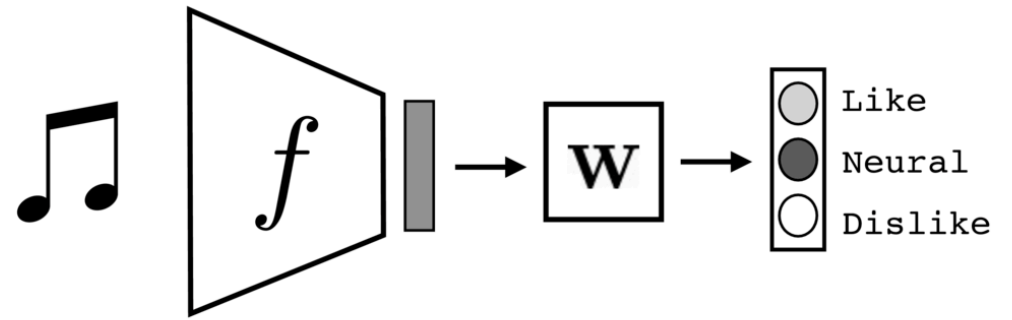
Genre recognition



"common"-sense representation

# Testing

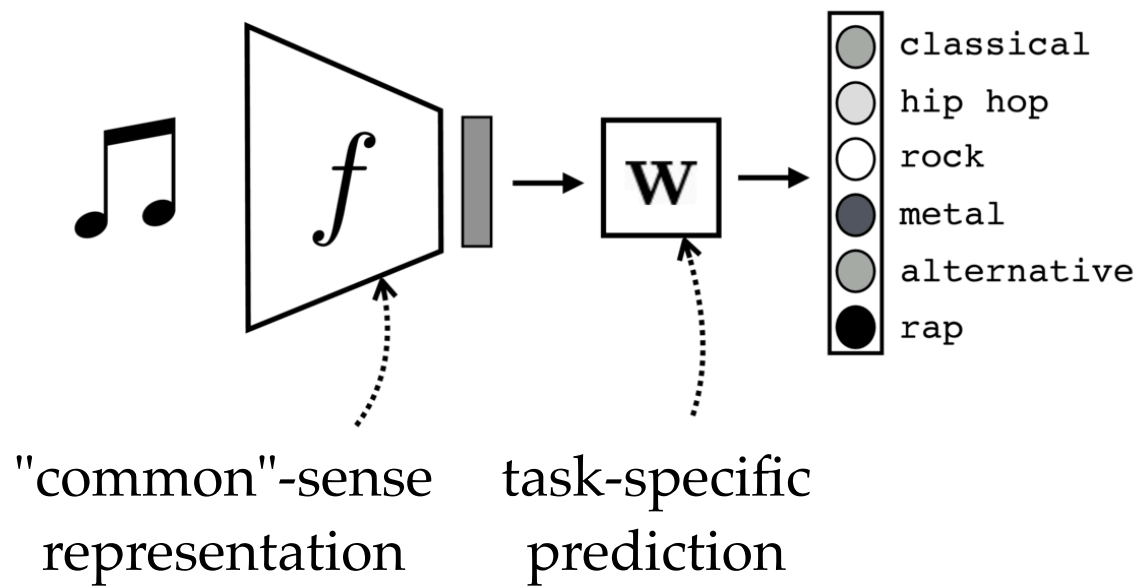
Preference prediction



Often, what we will be “tested” on is not what we were trained on.

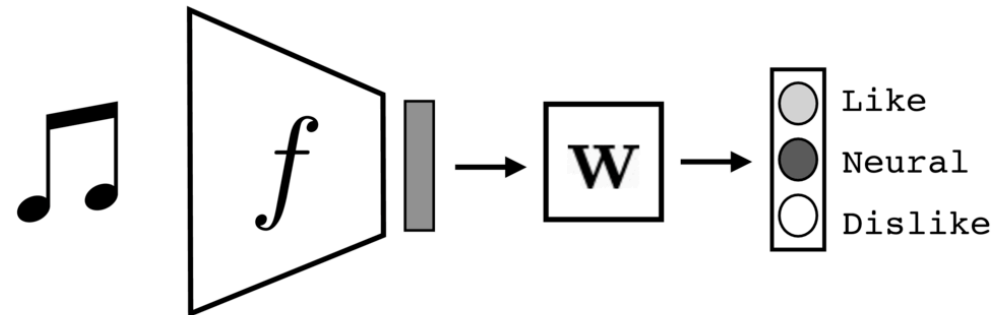
# Training

## Genre recognition



# Testing

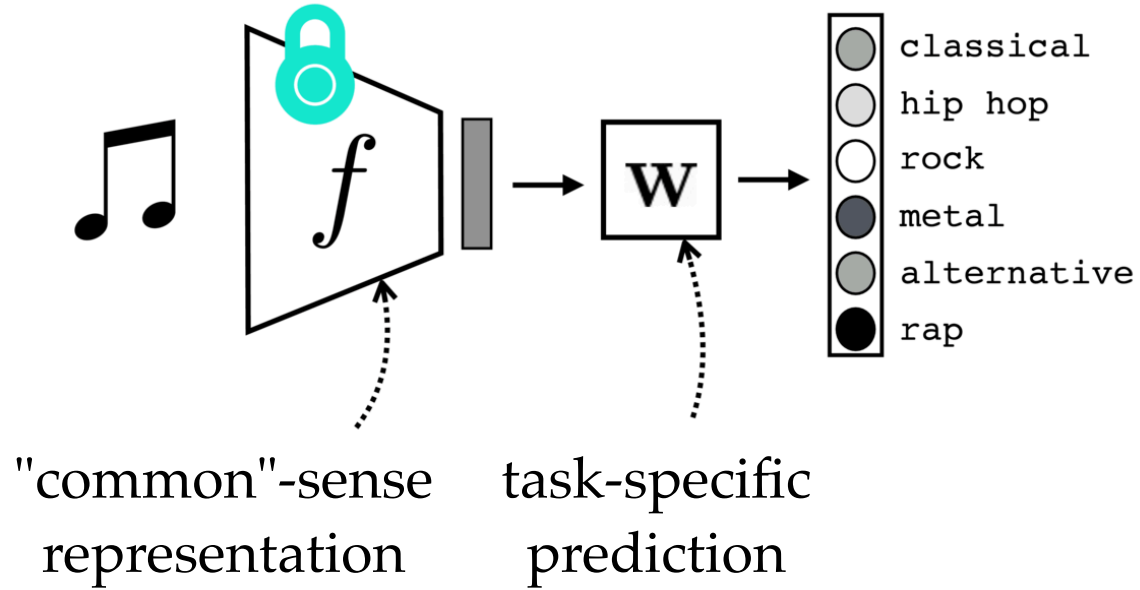
## Preference prediction



Often, what we will be “tested” on is not what we were trained on.

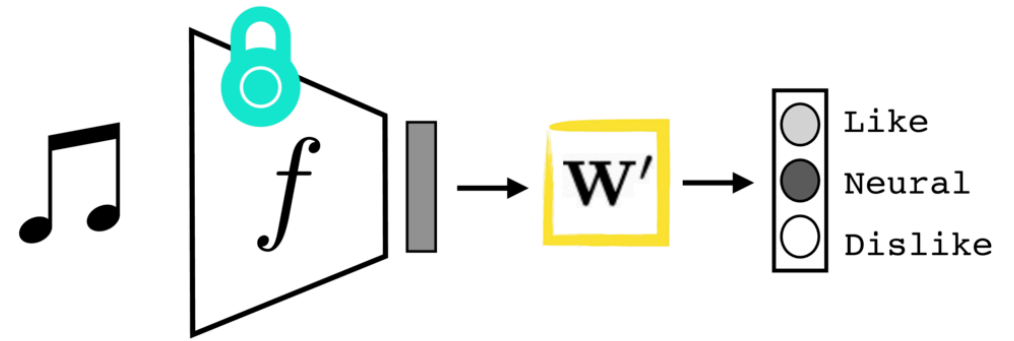
# Training

## Genre recognition



# Adapting

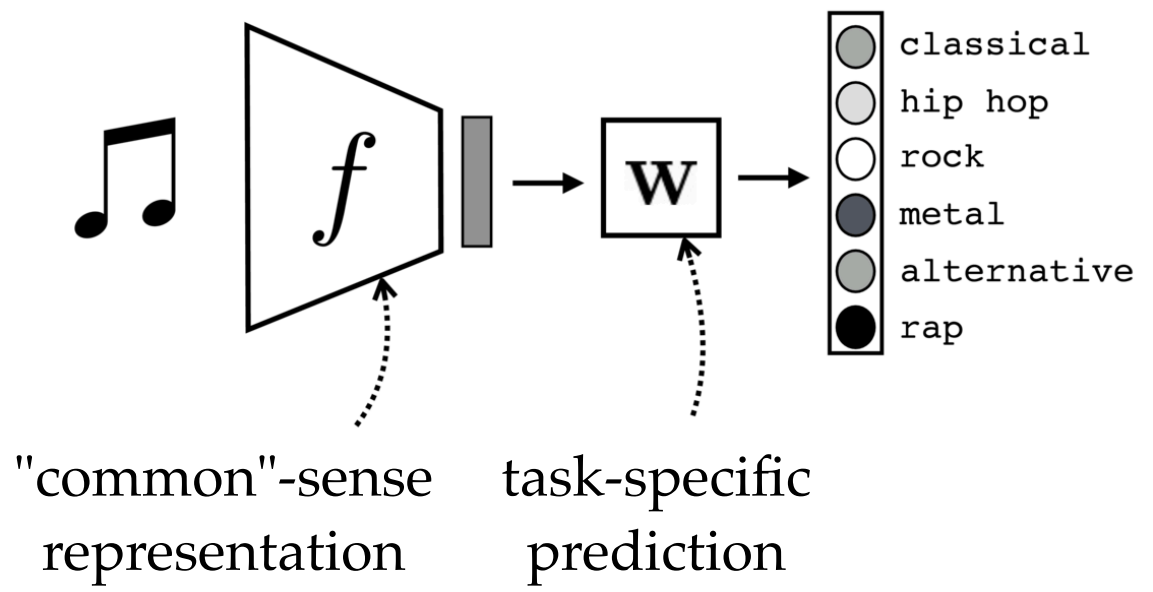
## Preference prediction



Final-layer adaptation: freeze  $f$ , train a new final layer to new target data

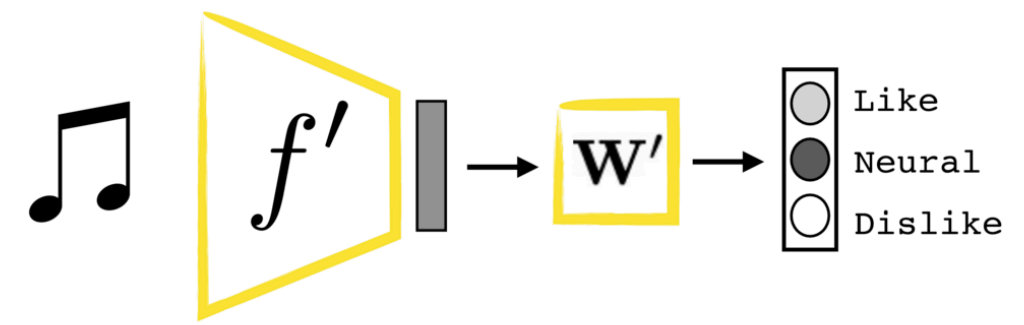
# Training

## Genre recognition



# Adapting

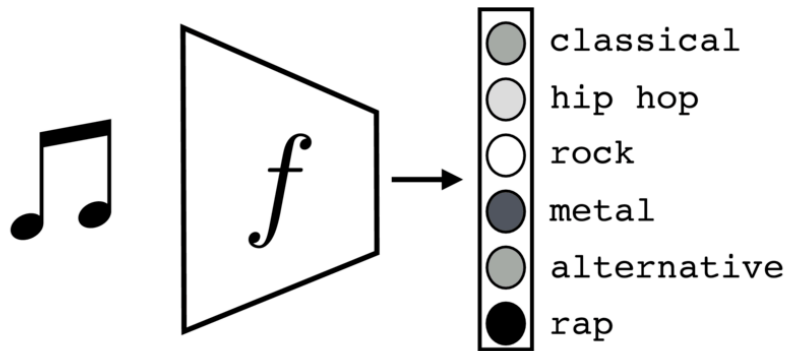
## Preference prediction



Finetuning: initialize  $f'$  as  $f$ , then continue training for  $f'$  as well, on new target data

# Pretraining

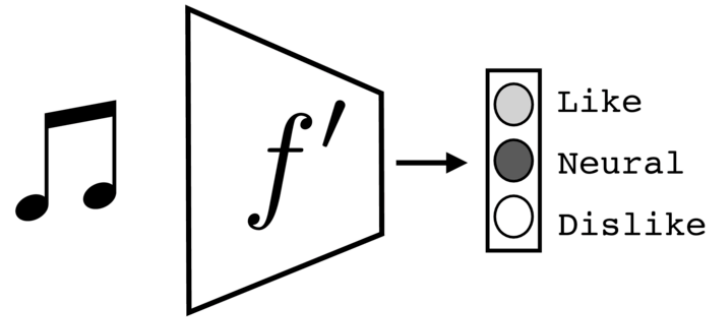
Genre recognition



*A lot of data*

# Adapting

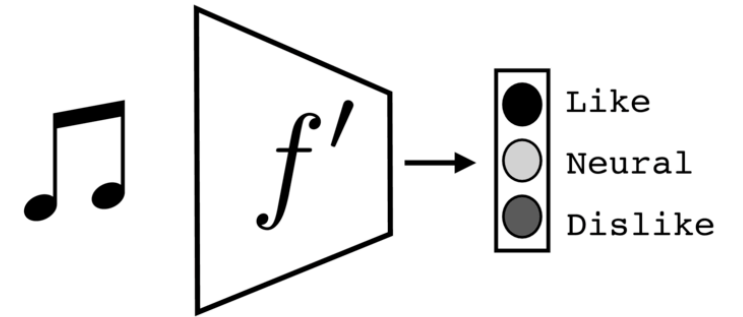
Preference prediction



*A little data*

# Testing

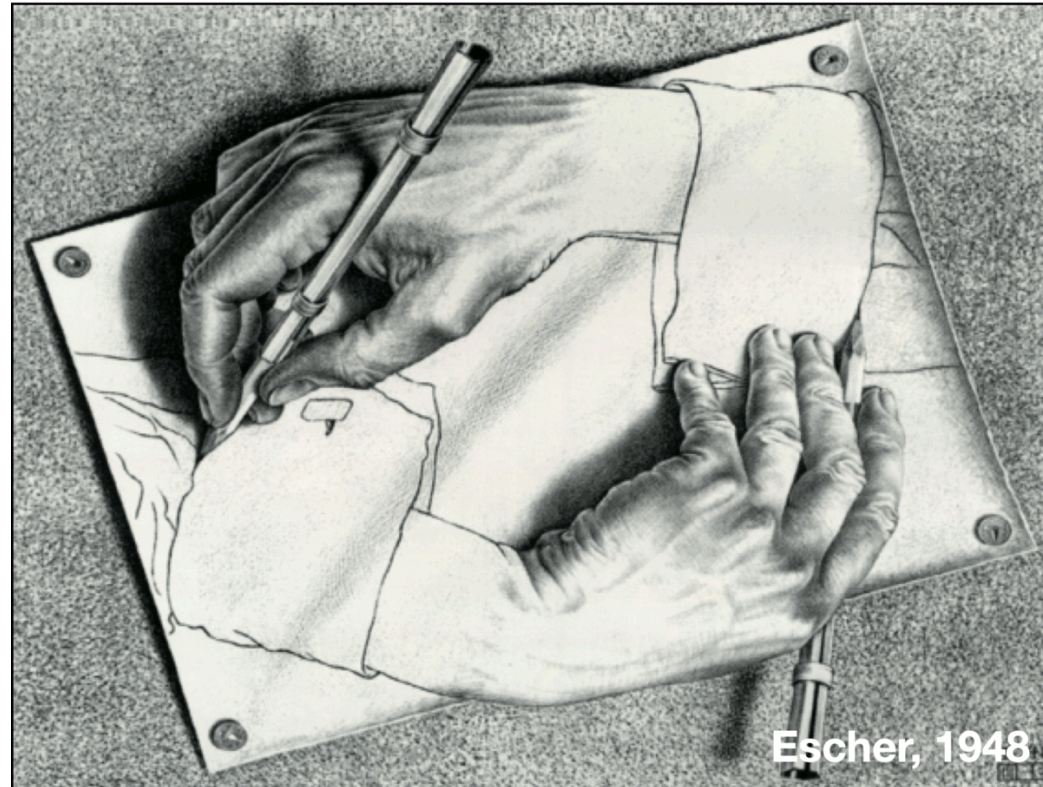
Preference prediction



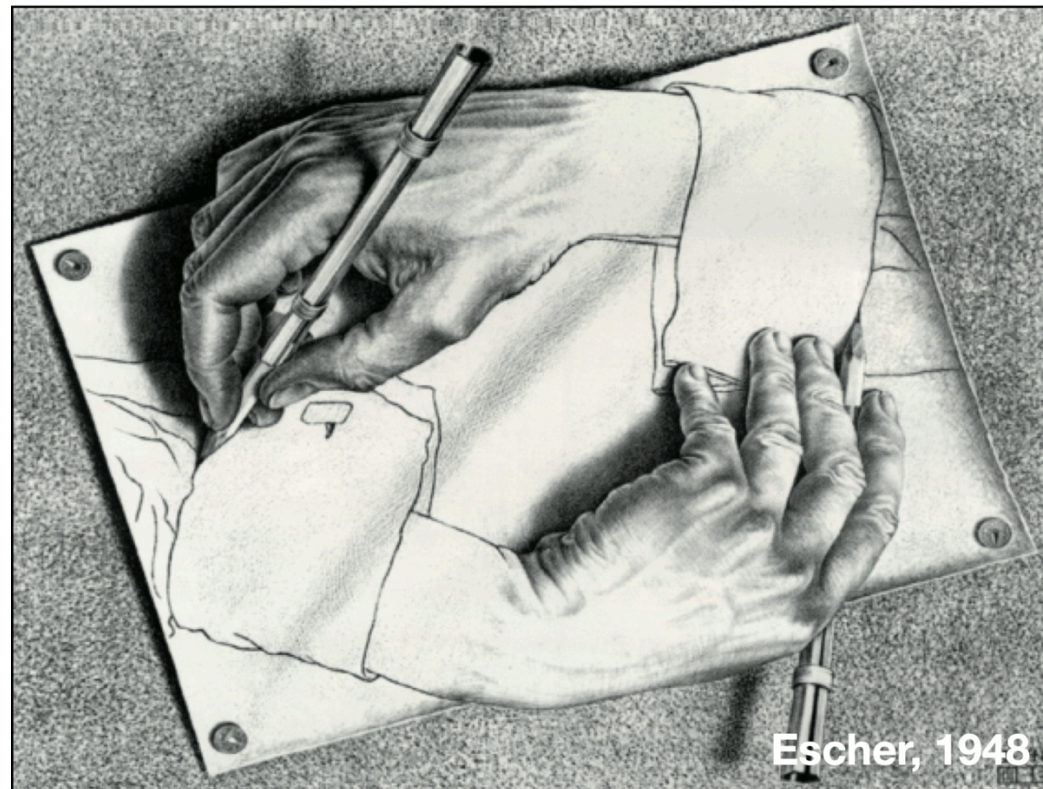
# Outline

- Neural networks are *representation* learners
- Self-supervised learning
- Auto-encoders
- Word embeddings
- (Some recent representation learning ideas)

Large Language Models (LLMs) are trained in this self-supervised way

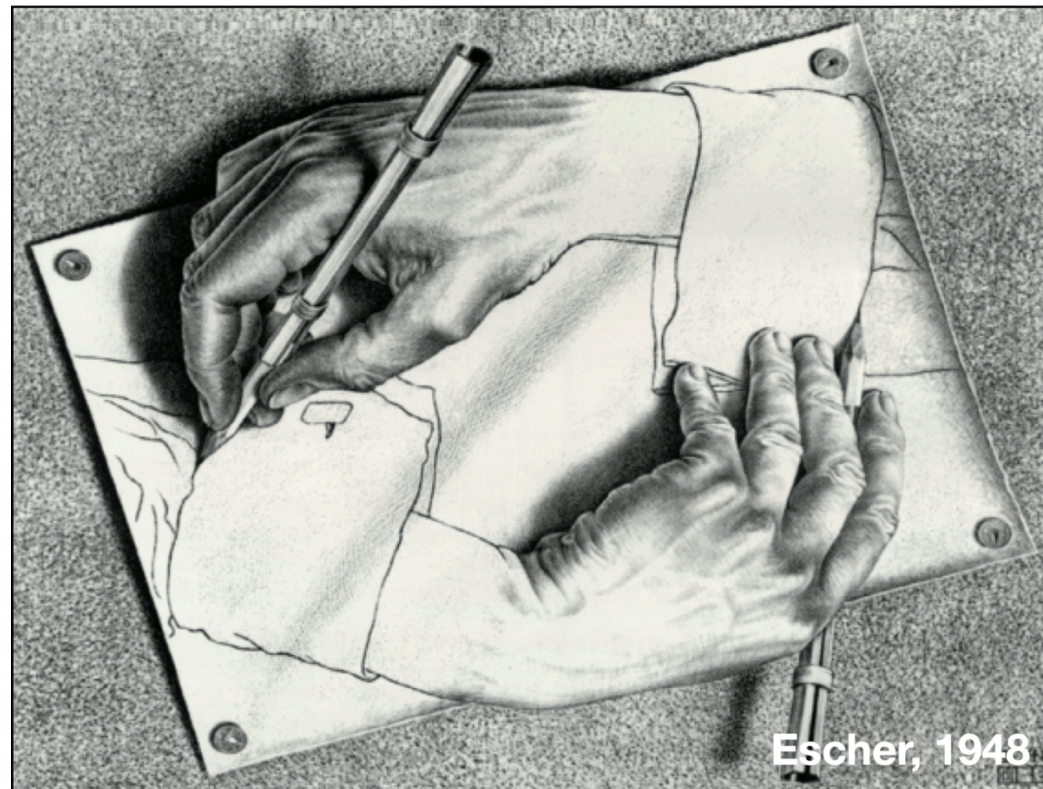


Large Language Models (LLMs) are trained in this self-supervised way



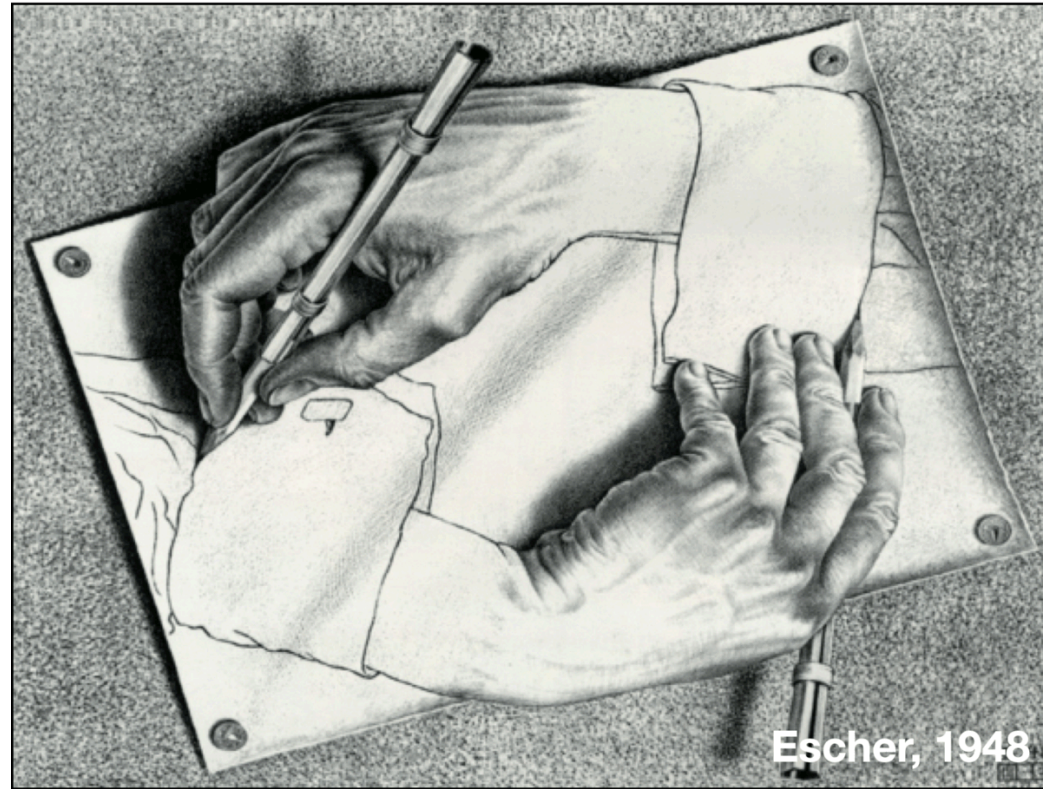
- Scrape the internet for plain texts.

Large Language Models (LLMs) are trained in this self-supervised way



- Scrape the internet for plain texts.
- Cook up “labels” (prediction targets) from these texts.

Large Language Models (LLMs) are trained in this self-supervised way



- Scrape the internet for plain texts.
- Cook up “labels” (prediction targets) from these texts.
- Convert “unsupervised” problem into “supervised” setup.

"To date, the cleverest thinker of all time was Issac. "

"To date, the cleverest thinker of all time was Issac. "



"To date, the cleverest thinker of all time was Issac. "



feature

⋮

To date, the

label

⋮

cleverest

"To date, the cleverest thinker of all time was Issac. "



feature

⋮

To date, the

To date, the cleverest

label

⋮

cleverest

thinker

"To date, the cleverest thinker of all time was Issac. "



feature

⋮

To date, the

To date, the cleverest

To date, the cleverest thinker

label

⋮

cleverest

thinker

was

"To date, the cleverest thinker of all time was Issac. "



feature

⋮

To date, the

To date, the cleverest

To date, the cleverest thinker

⋮

label

⋮

cleverest

thinker

was

⋮

"To date, the cleverest thinker of all time was Issac. "



feature

⋮

To date, the

To date, the cleverest

To date, the cleverest thinker

⋮

To date, the cleverest thinker of all time was

label

⋮

cleverest

thinker

was

⋮

Issac

"To date, the cleverest thinker of all time was Issac. "



feature

⋮

To date, the

To date, the cleverest

To date, the cleverest thinker

⋮

To date, the cleverest thinker of all time was

label

⋮

cleverest

thinker

was

⋮

Issac

To date, the \_\_\_\_\_



model

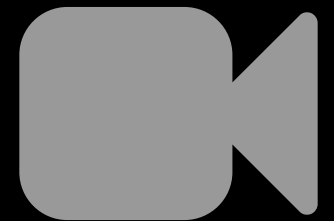
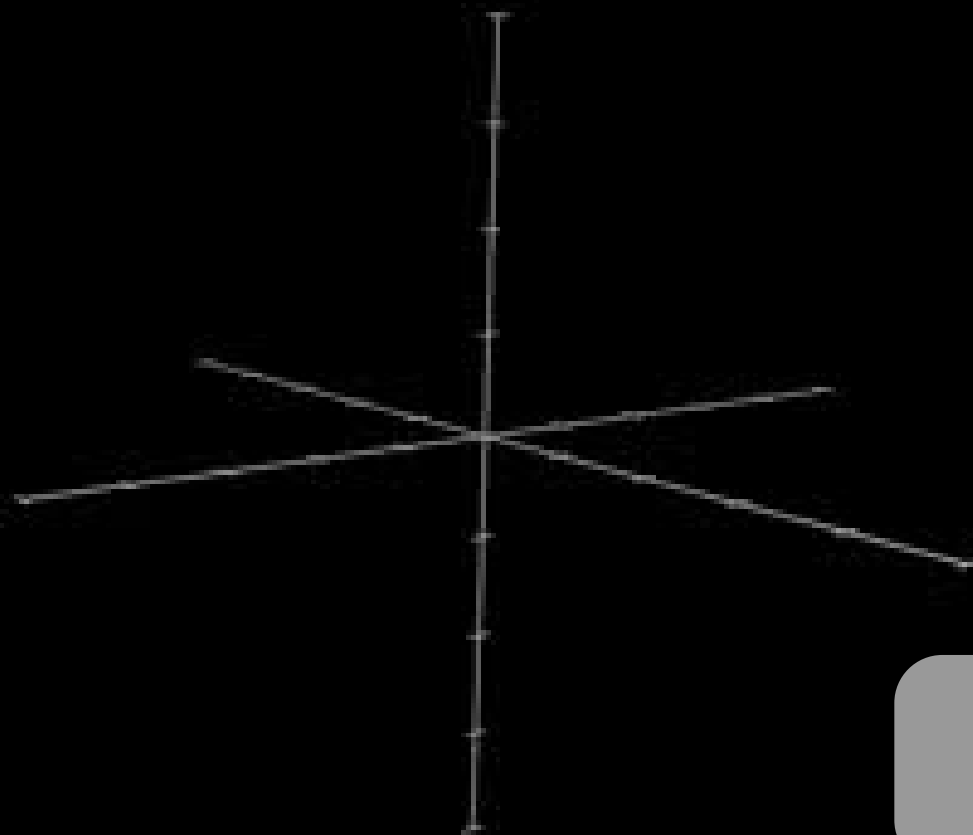


# Word embedding

Words



Vectors



# dot-product similarity



$$\begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \vdots \\ v_n \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_n \end{bmatrix}$$

Dot product

$$= v_1 w_1 + v_2 w_2 + v_3 w_3 + \vdots + v_n w_n$$



Good word embeddings space is equipped with sensible dot-product similarity

Good word embeddings space is equipped with sensible dot-product similarity

For now, let's look at how good embeddings enable "soft" dictionary look-up:

Good word embeddings space is equipped with sensible dot-product similarity

For now, let's look at how good embeddings enable "soft" dictionary look-up:

```
1 dict_en2fr = {  
2     "apple" : "pomme",  
3     "banana" : "banane",  
4     "lemon" : "citron"}
```

Good word embeddings space is equipped with sensible dot-product similarity

For now, let's look at how good embeddings enable "soft" dictionary look-up:

```
1 dict_en2fr = {  
2     "apple" : "pomme",  
3     "banana" : "banane",  
4     "lemon" : "citron"}
```

Key

Value

Good word embeddings space is equipped with sensible dot-product similarity

For now, let's look at how good embeddings enable "soft" dictionary look-up:

```
1 dict_en2fr = {  
2   "apple" : "pomme",  
3   "banana" : "banane",  
4   "lemon" : "citron"}
```

Key            Value

apple : pomme

Good word embeddings space is equipped with sensible dot-product similarity

For now, let's look at how good embeddings enable "soft" dictionary look-up:

```
1 dict_en2fr = {  
2   "apple" : "pomme",  
3   "banana" : "banane",  
4   "lemon" : "citron"}
```

Key            Value

apple : pomme

banana : banane

Good word embeddings space is equipped with sensible dot-product similarity

For now, let's look at how good embeddings enable "soft" dictionary look-up:

```
1 dict_en2fr = {  
2   "apple" : "pomme",  
3   "banana" : "banane",  
4   "lemon" : "citron"}
```

Key            Value

apple : pomme

banana : banane

lemon : citron

```
1 dict_en2fr = {
2     "apple" : "pomme",
3     "banana" : "banane",
4     "lemon" : "citron"}
5
6 query = "lemon"
7 output = dict_en2fr[query]
```

Key            Value

apple : pomme

banana : banane

lemon : citron

```
1 dict_en2fr = {
2     "apple" : "pomme",
3     "banana" : "banane",
4     "lemon" : "citron"}
5
6 query = "lemon"
7 output = dict_en2fr[query]
```

Query

Key

Value

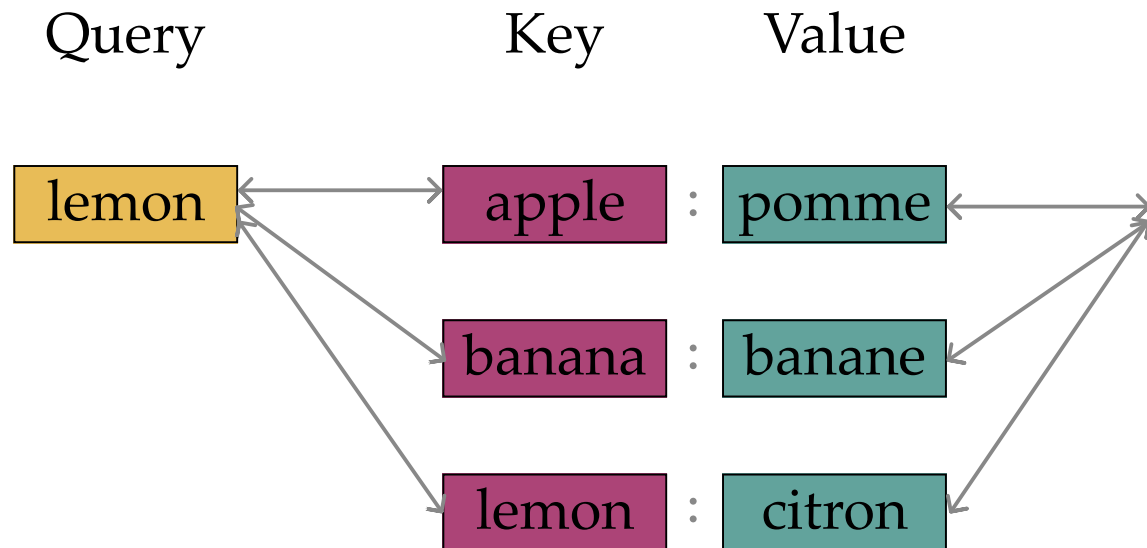
lemon

apple : pomme

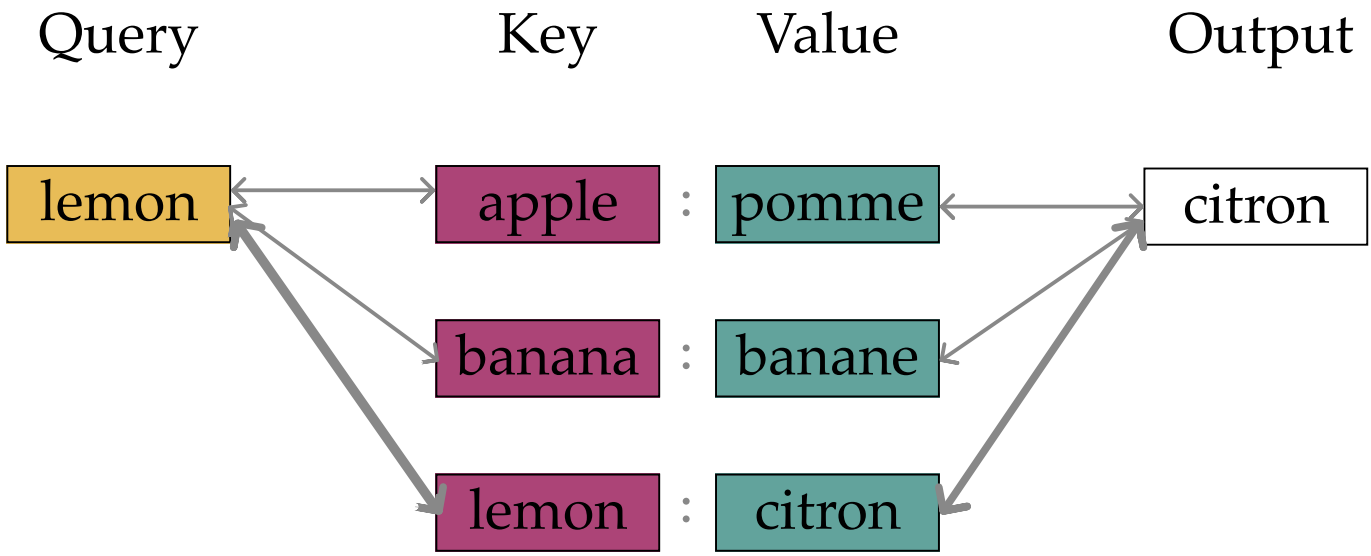
banana : banane

lemon : citron

```
1 dict_en2fr = {
2     "apple" : "pomme",
3     "banana" : "banane",
4     "lemon" : "citron"}
5
6 query = "lemon"
7 output = dict_en2fr[query]
```

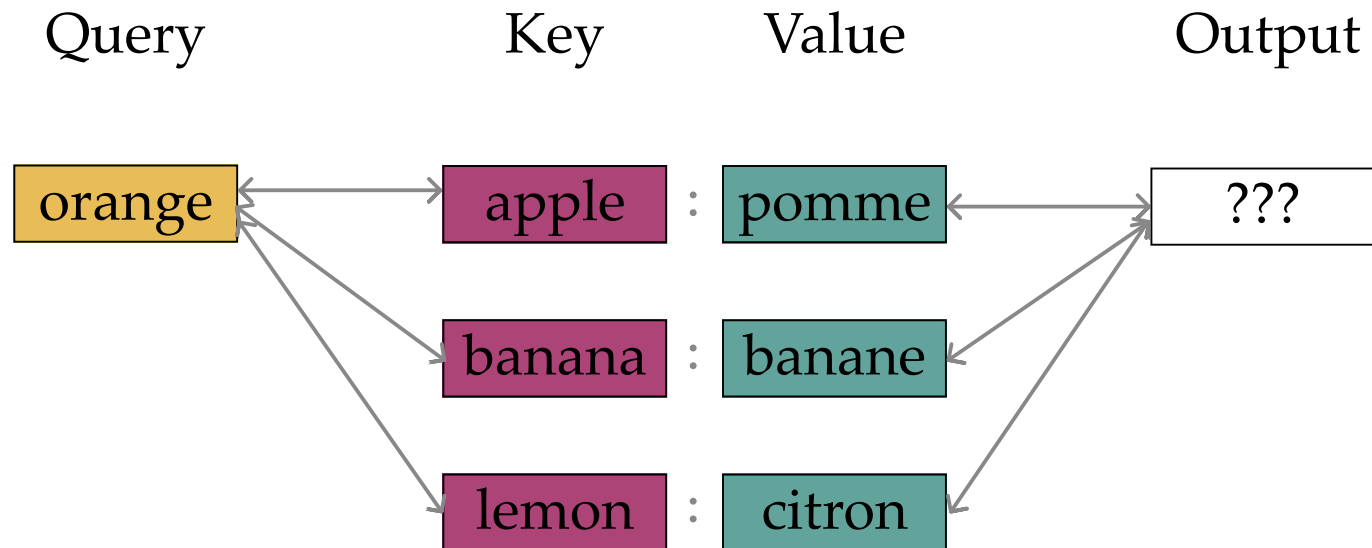


```
1 dict_en2fr = {
2     "apple" : "pomme",
3     "banana" : "banane",
4     "lemon" : "citron"}
5
6 query = "lemon"
7 output = dict_en2fr[query]
```



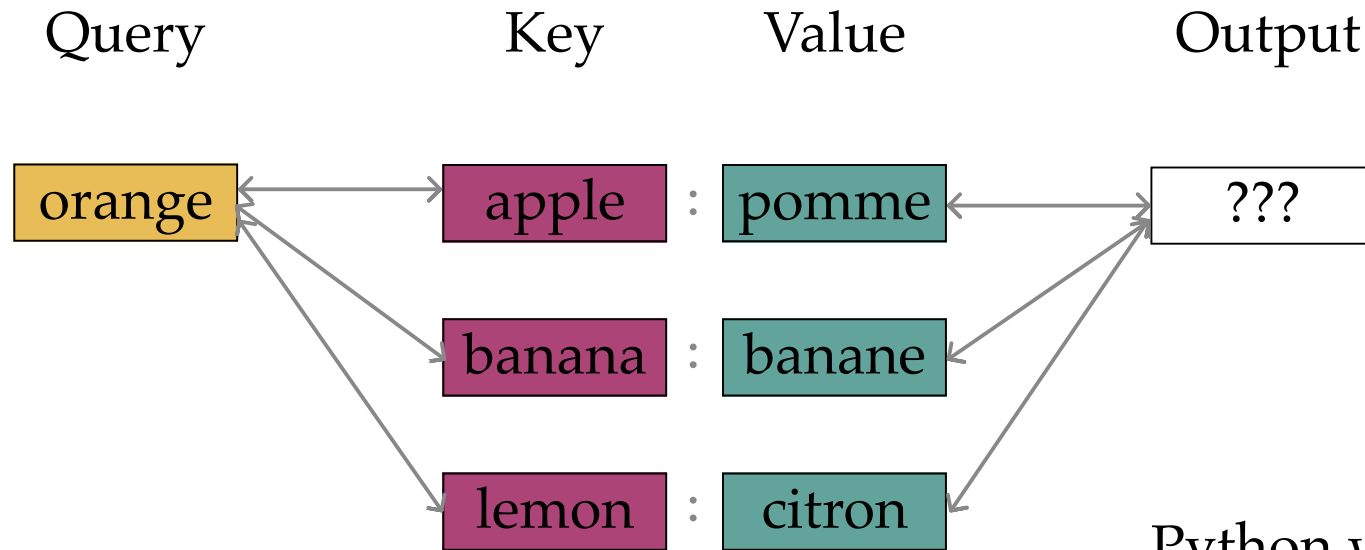
What if we run

```
1 dict_en2fr = {  
2     "apple" : "pomme",  
3     "banana" : "banane",  
4     "lemon" : "citron"}  
5  
6 query = "orange"  
7 output = dict_en2fr[query]
```



# What if we run

```
1 dict_en2fr = {  
2     "apple" : "pomme",  
3     "banana" : "banane",  
4     "lemon" : "citron"}  
5  
6 query = "orange"  
7 output = dict_en2fr[query]
```

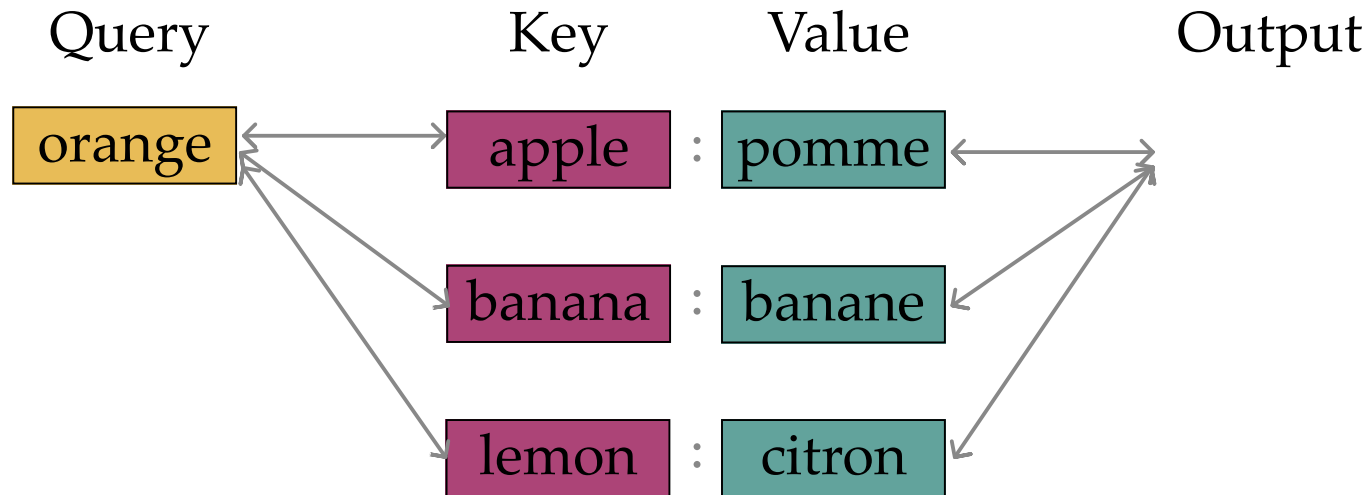


Python would complain. 🤪

What if we run

```
1 dict_en2fr = {  
2     "apple" : "pomme",  
3     "banana" : "banane",  
4     "lemon" : "citron"}  
5  
6 query = "orange"  
7 output = dict_en2fr[query]
```

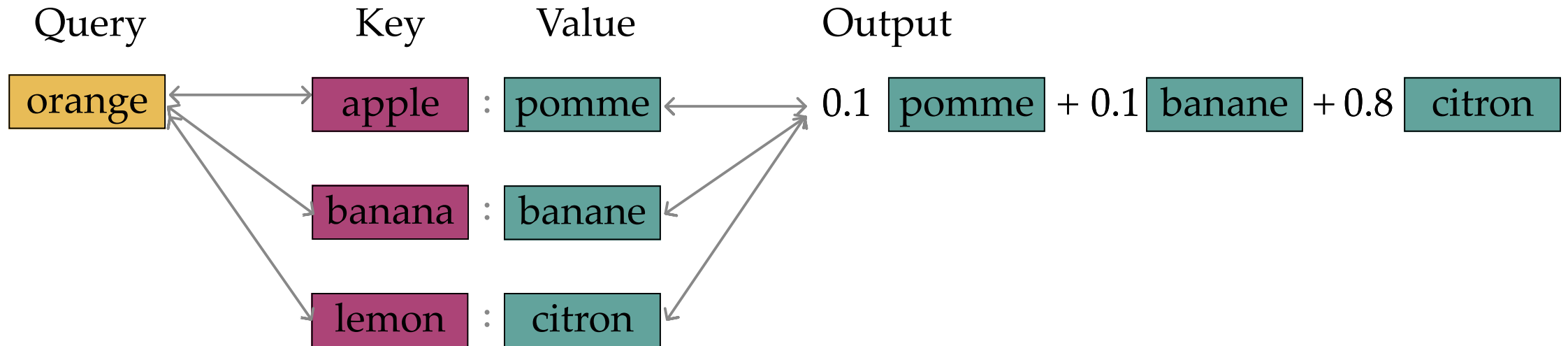
But we can probably see the rationale behind this:



What if we run

```
1 dict_en2fr = {  
2     "apple" : "pomme",  
3     "banana" : "banane",  
4     "lemon" : "citron"}  
5  
6 query = "orange"  
7 output = dict_en2fr[query]
```

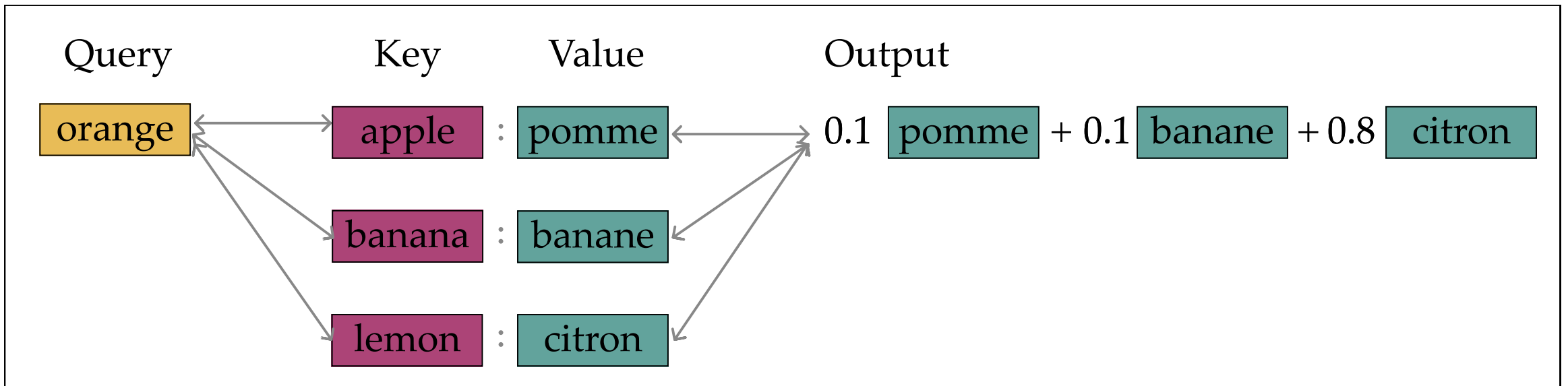
But we can probably see the rationale behind this:

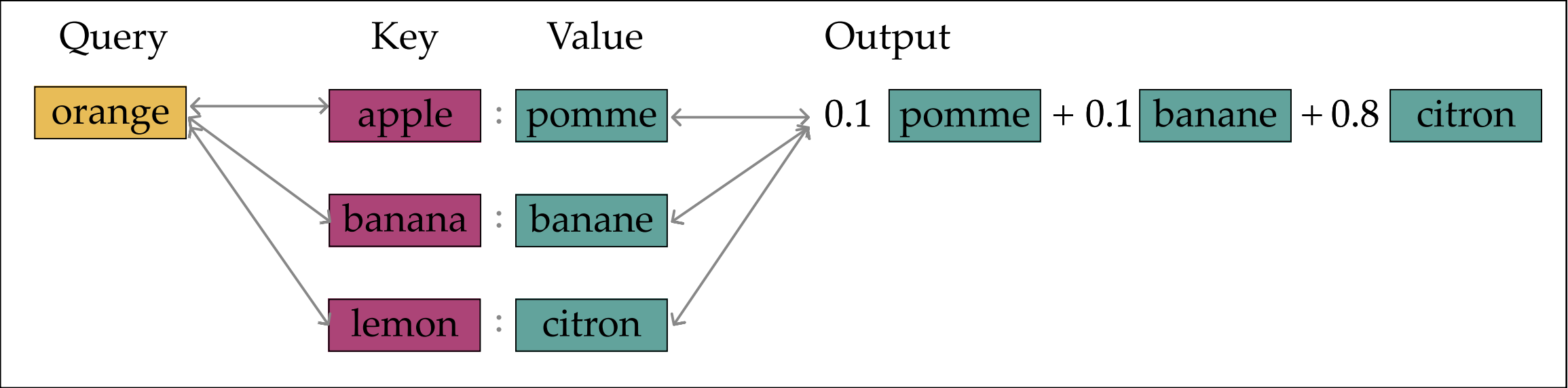


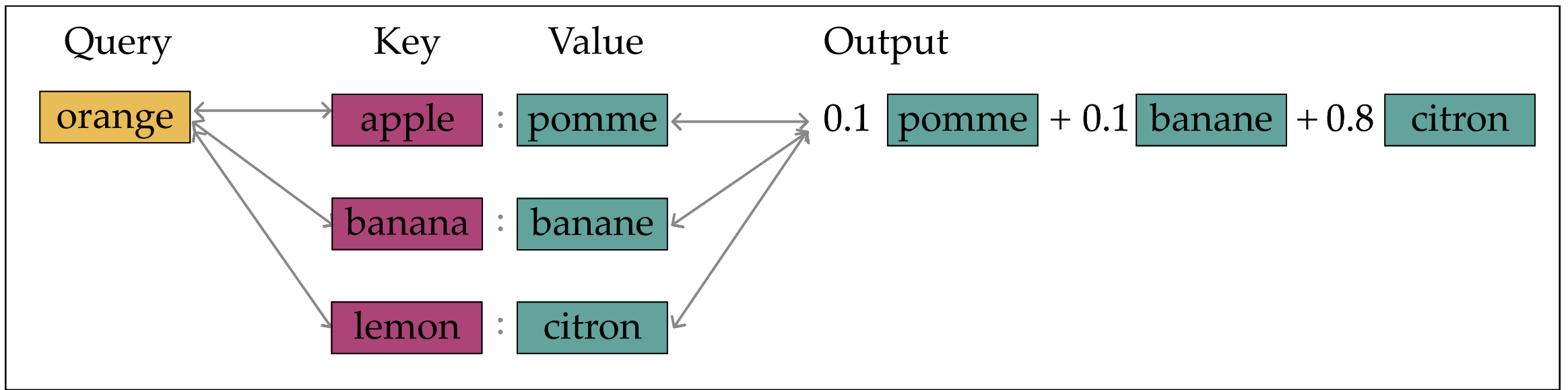
What if we run

```
1 dict_en2fr = {  
2     "apple" : "pomme",  
3     "banana" : "banane",  
4     "lemon" : "citron"}  
5  
6 query = "orange"  
7 output = dict_en2fr[query]
```

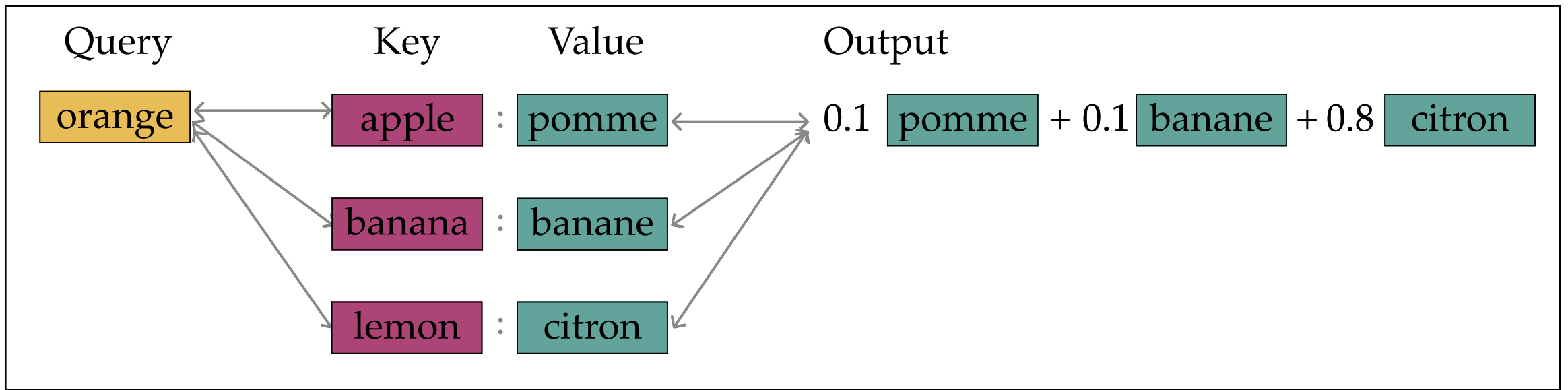
But we can probably see the rationale behind this:





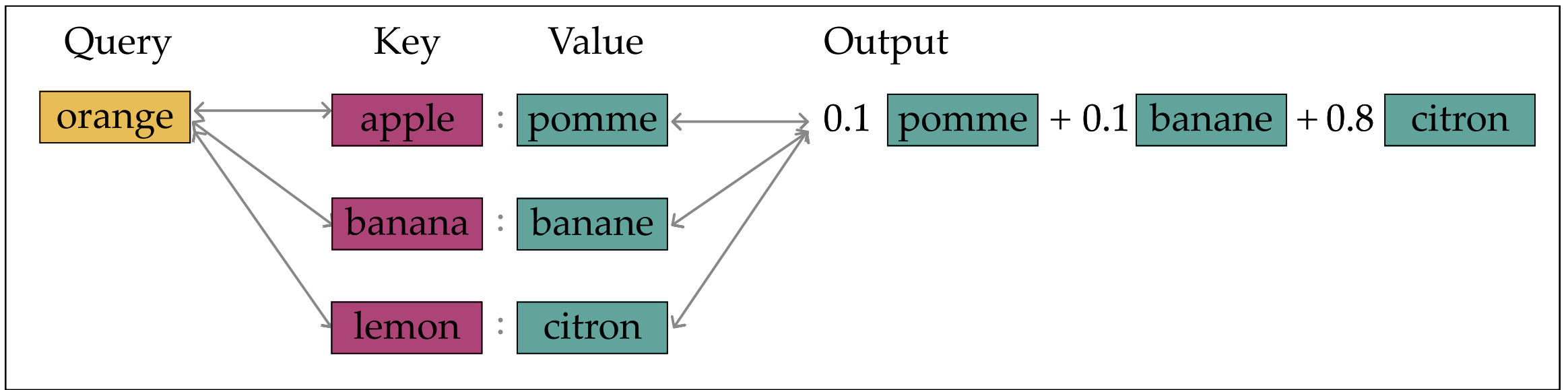


We put (query, key, value) in "good" embeddings in our human brain



We put (query, key, value) in "good" embeddings in our human brain

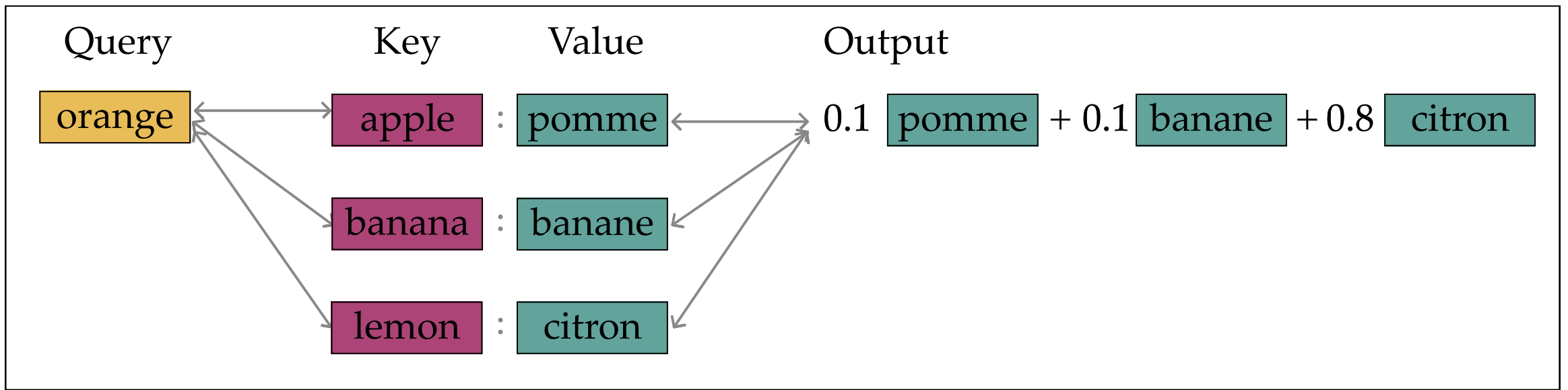
such that "merging" the values  $0.1 \text{ pomme} + 0.1 \text{ banane} + 0.8 \text{ citron}$



We put (query, key, value) in "good" embeddings in our human brain

such that "merging" the values  $0.1 \text{ pomme} + 0.1 \text{ banane} + 0.8 \text{ citron}$

via these merging percentages  $[0.1 \ 0.1 \ 0.8]$  made sense



We put (query, key, value) in "good" embeddings in our human brain

such that "merging" the values  $0.1 \text{ pomme} + 0.1 \text{ banane} + 0.8 \text{ citron}$

via these merging percentages  $[0.1 \ 0.1 \ 0.8]$  made sense

**very roughly**, the attention mechanism in transformers automates this process. We'll see the details next week.

# Outline

- Neural networks are *representation* learners
- Self-supervised learning
- Auto-encoders
- Word embeddings
- (Some recent representation learning ideas)

# Outline

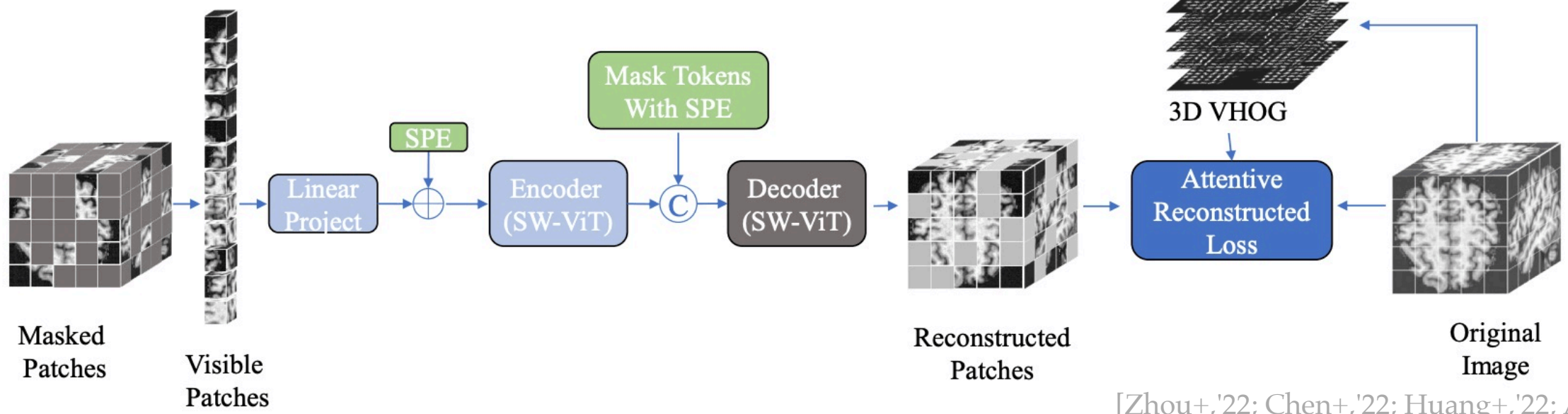
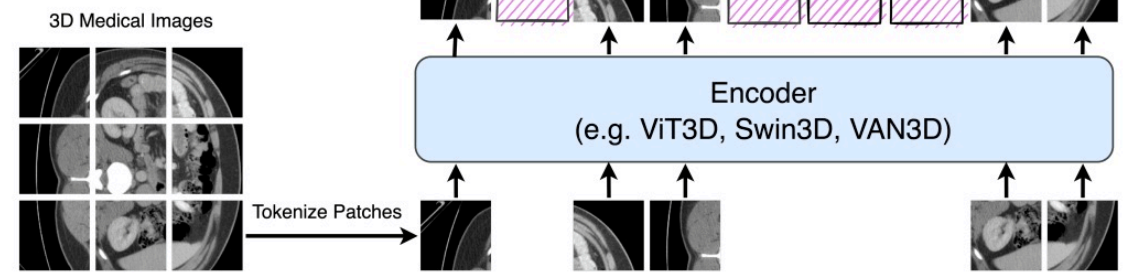
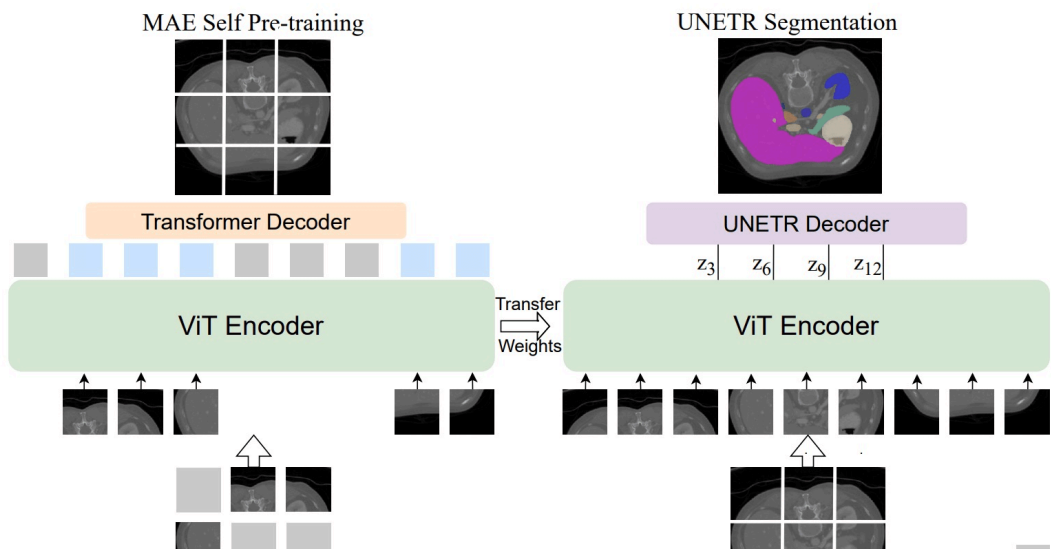
- Neural networks are *representation* learners
- Self-supervised learning
- Auto-encoders
- Word embeddings
- (Some recent representation learning ideas)
  - 1. masking: other applications
  - 2. contrastive: similarity

# Outline

- Neural networks are *representation* learners
- Self-supervised learning
- Auto-encoders
- Word embeddings
- (Some recent representation learning ideas)
  - 1. masking: other applications
  - 2. contrastive: similarity
  - 3. multi-modality: alignment

# 1. Masking

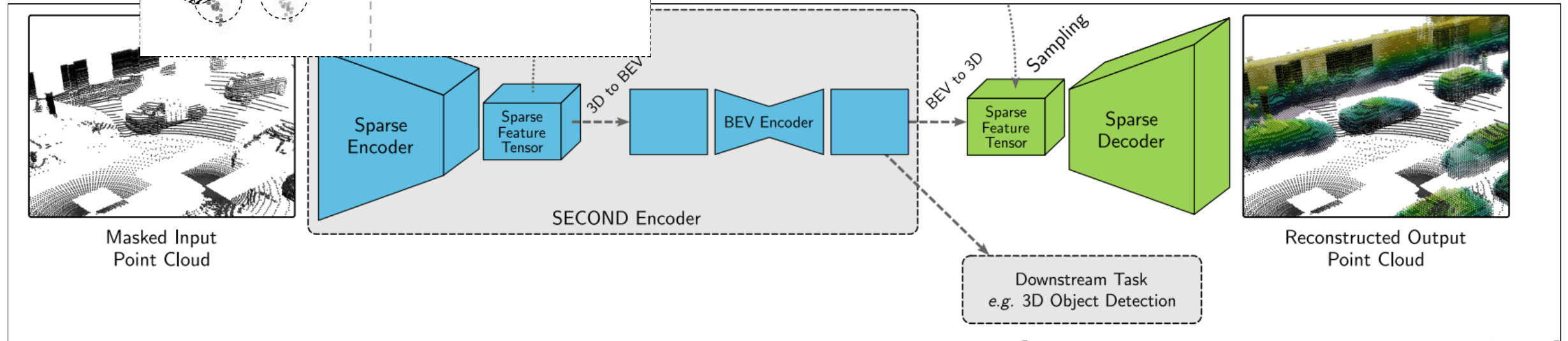
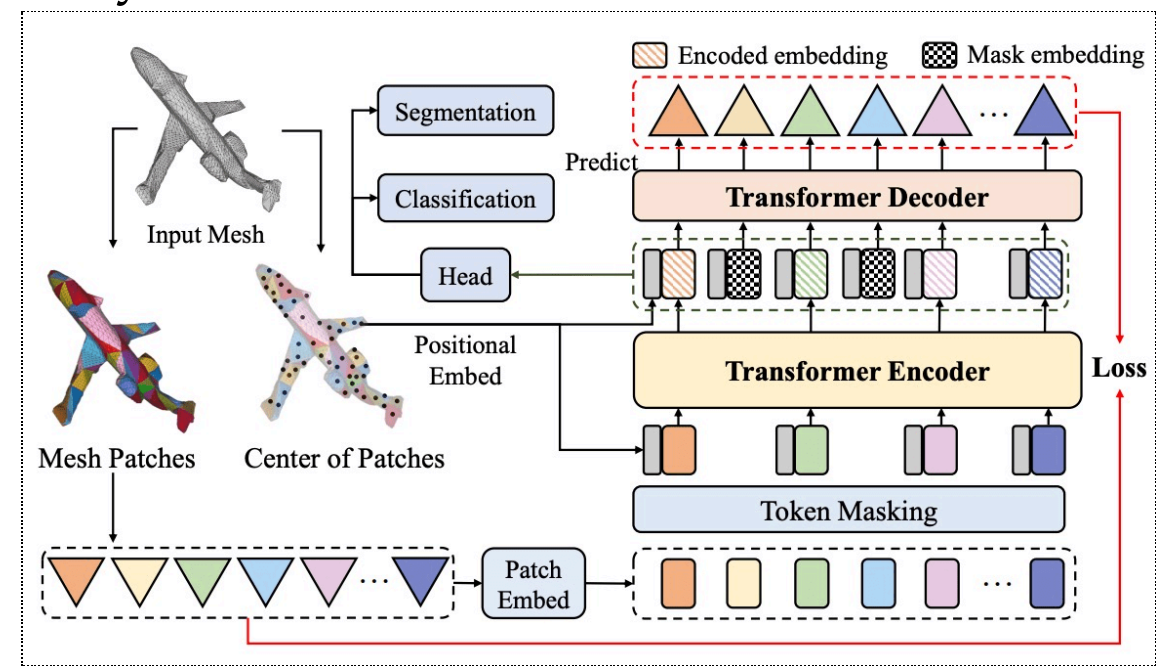
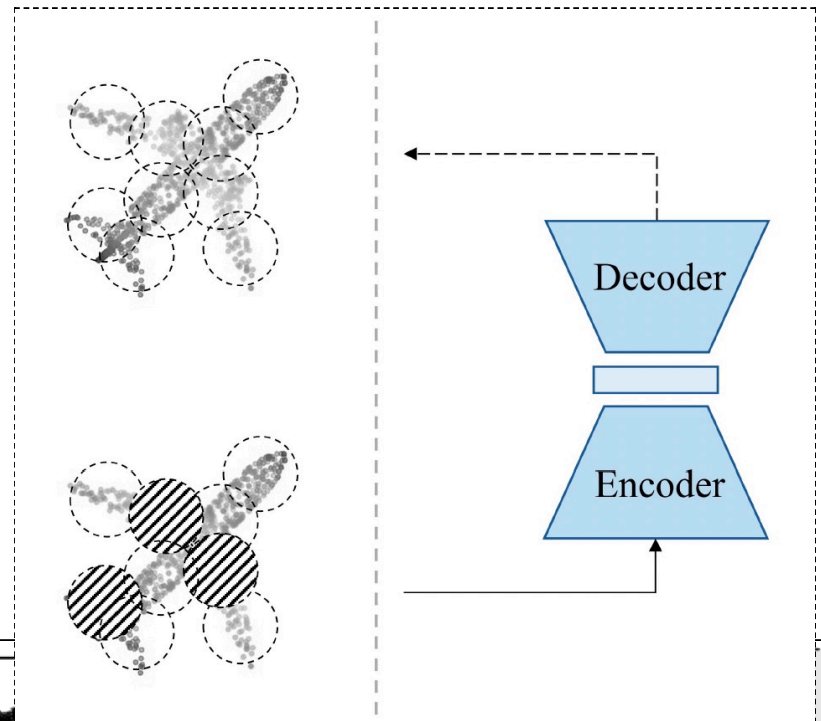
e.g. medical imaging



[Zhou+, '22; Chen+, '22; Huang+, '22; An+, '22]

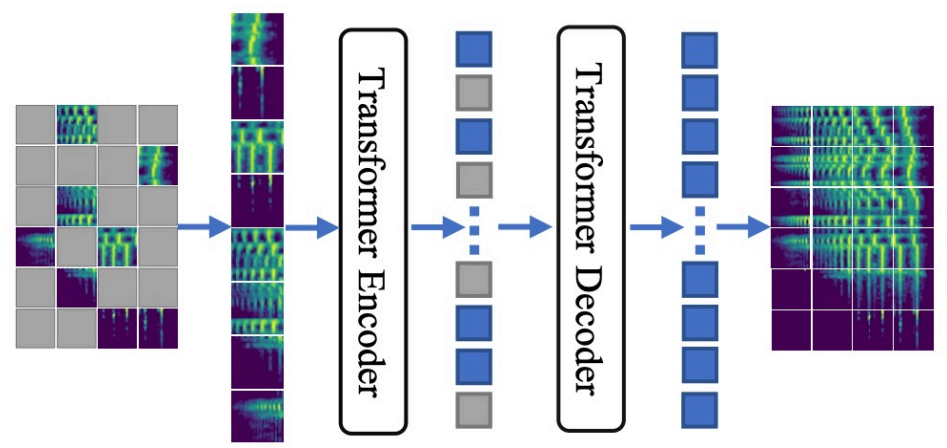
# 1. Masking

e.g. 3d geometry

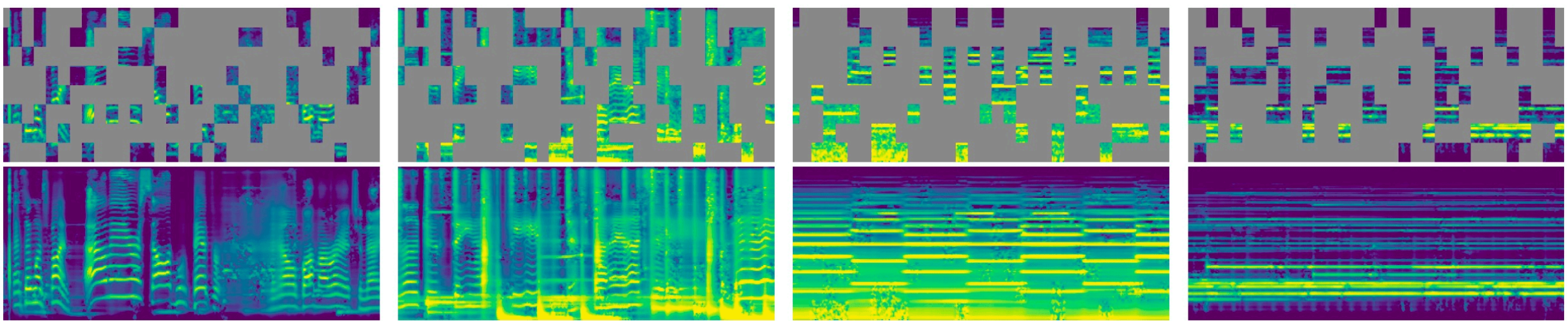
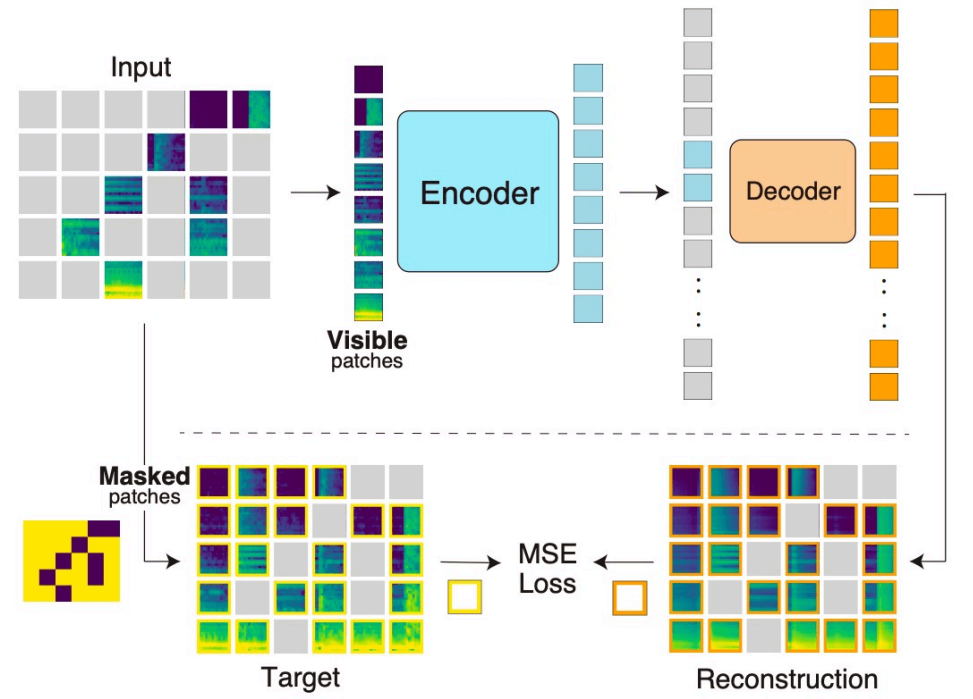


[Pang+, '22; Liang+, '22; Min+, '22; Krispel+, '22]

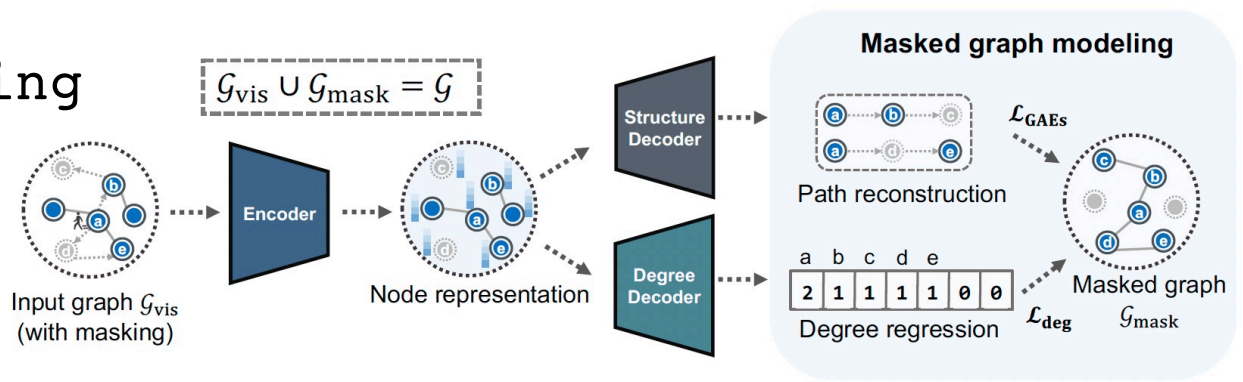
# 1. Masking



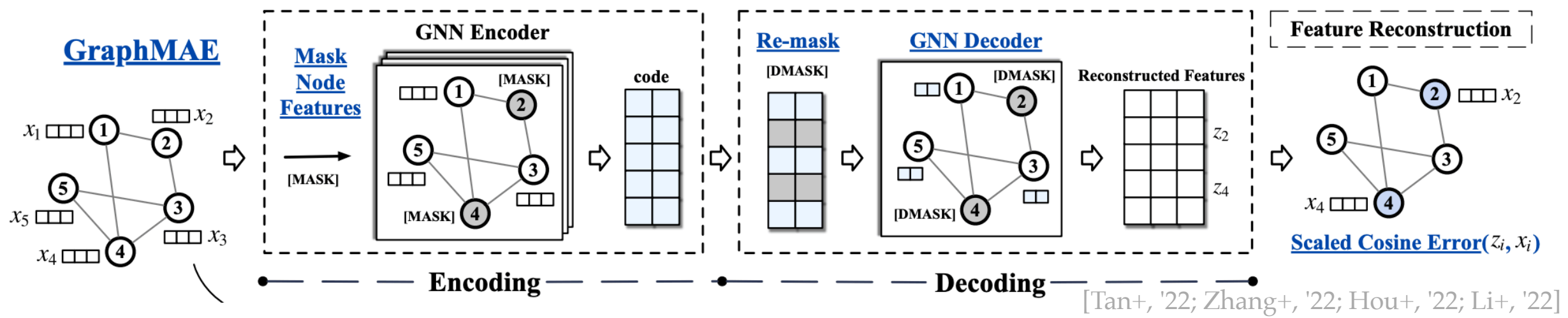
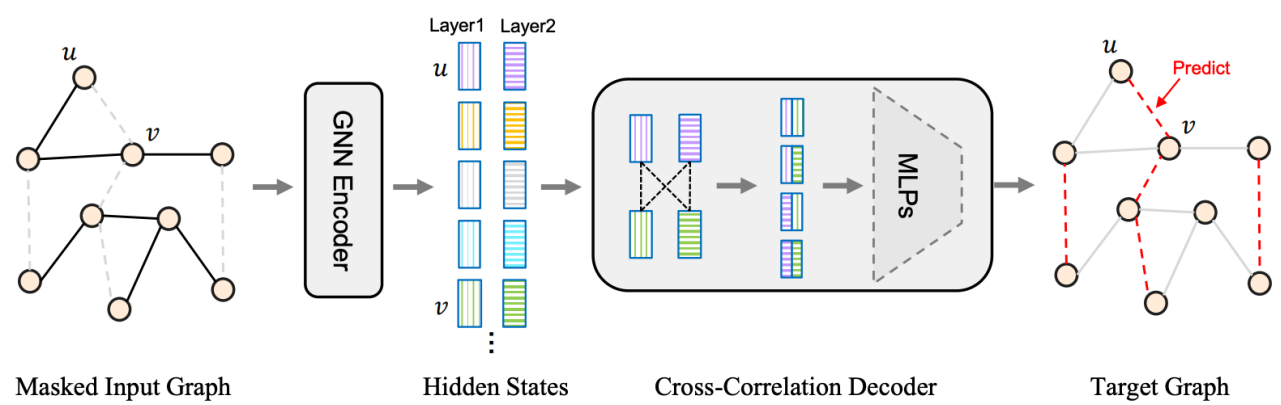
e.g. audio



# 1. Masking

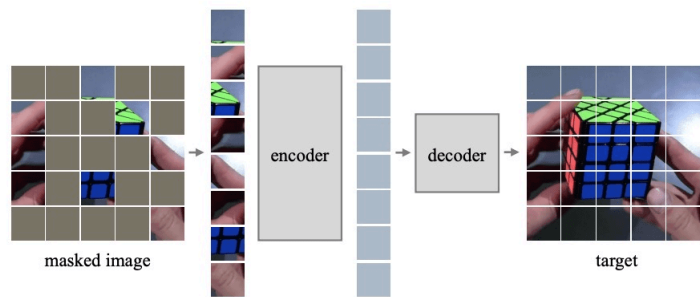


e.g. graphs

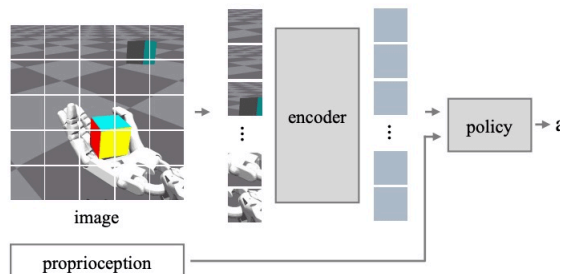


# 1. Masking

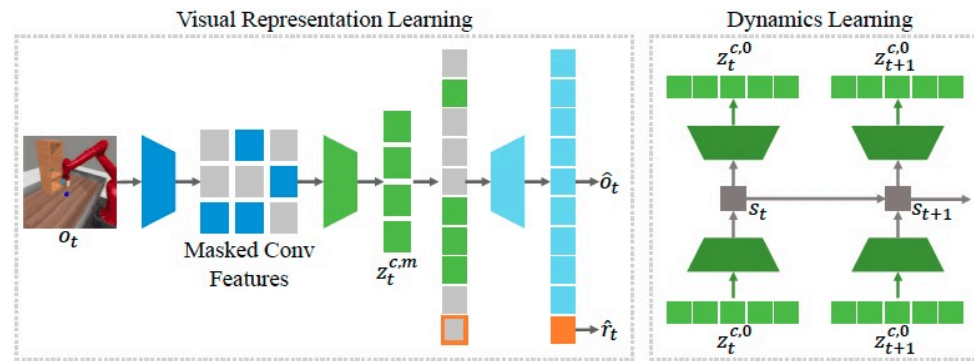
e.g. robotics



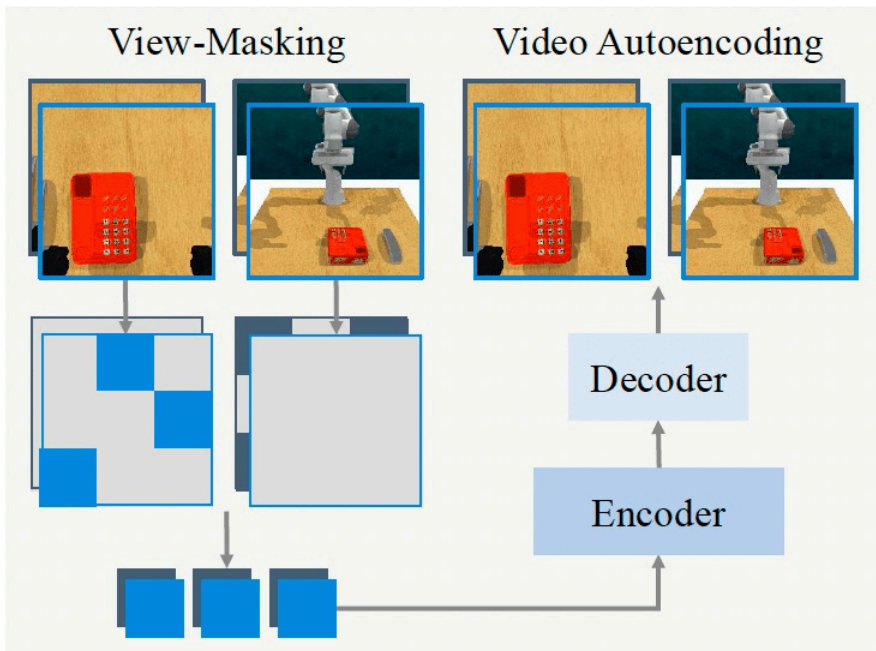
(a) masked visual pretraining



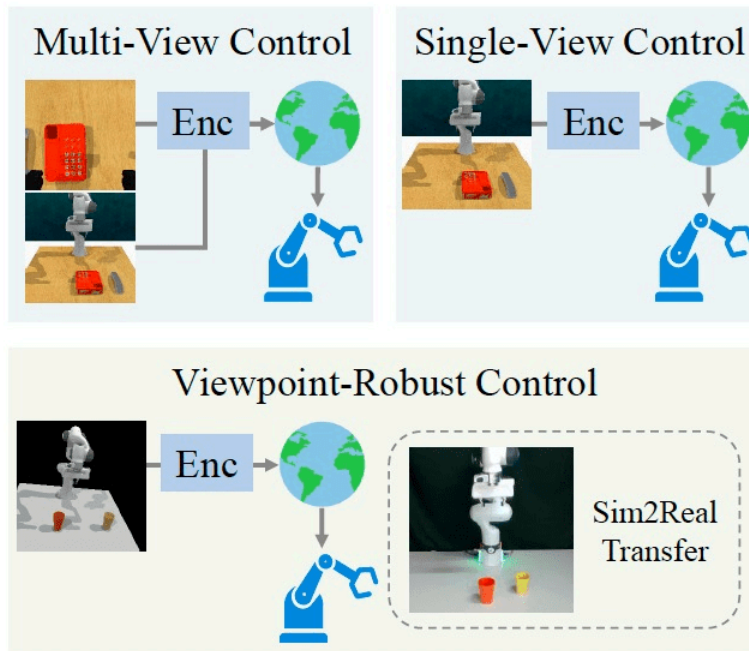
(b) learning motor control



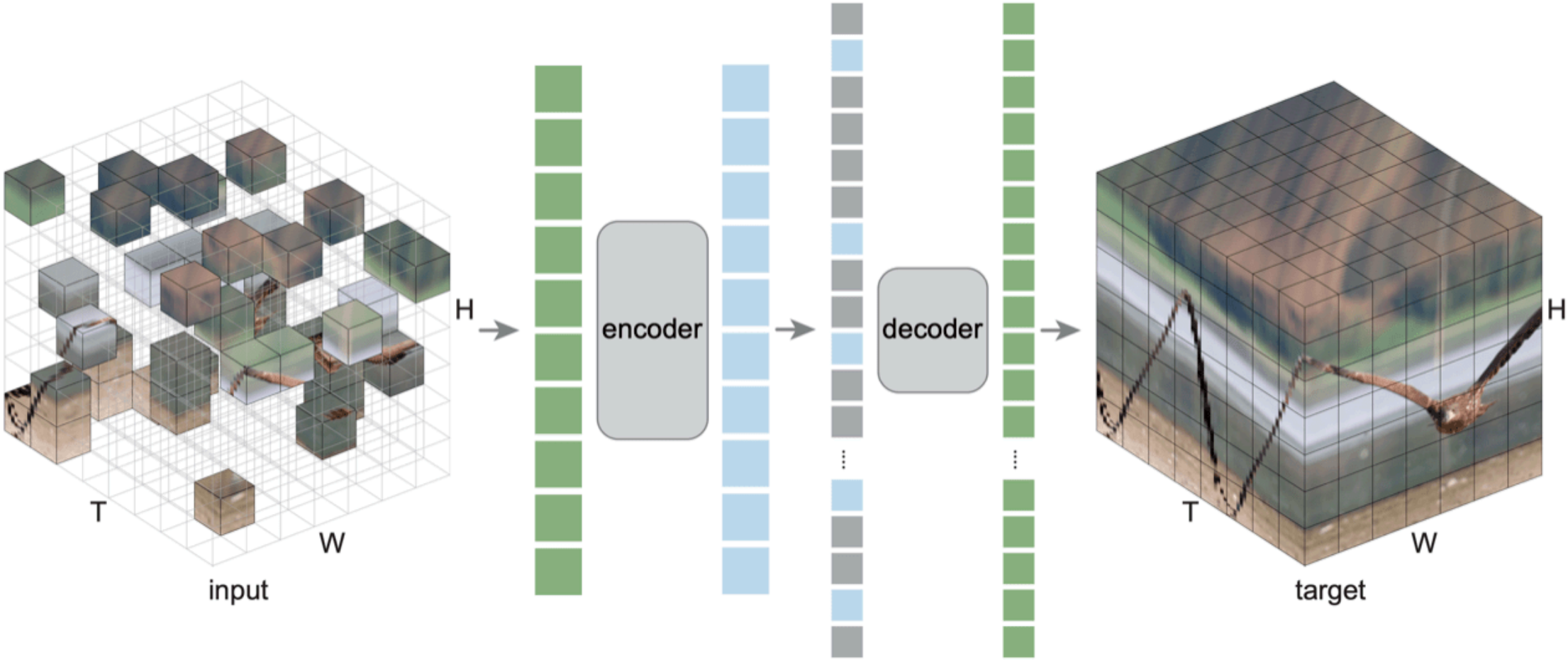
## Multi-View Masked Autoencoders



## Visual Robotic Manipulation

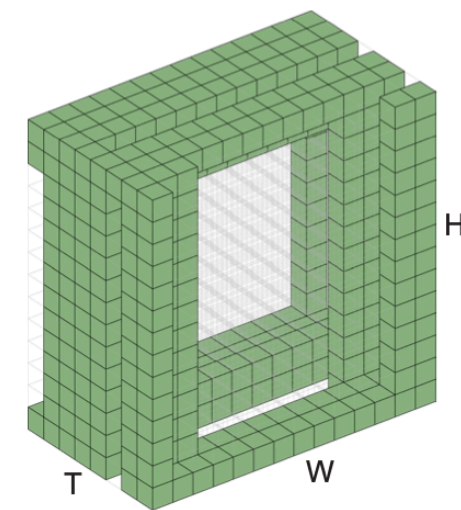
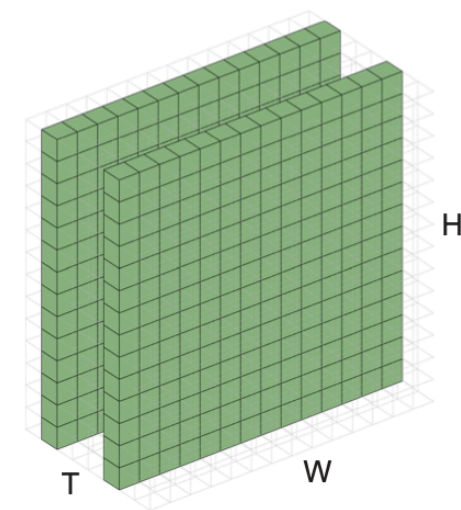
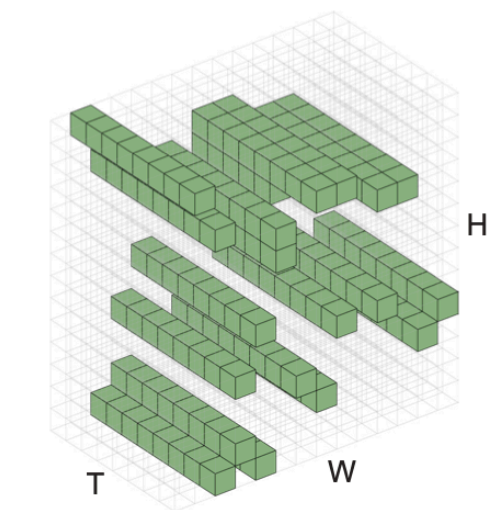
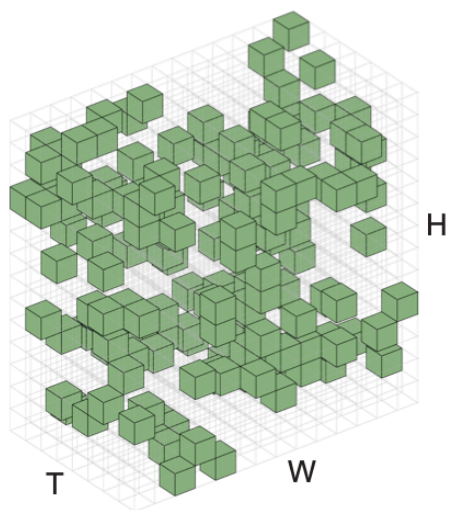
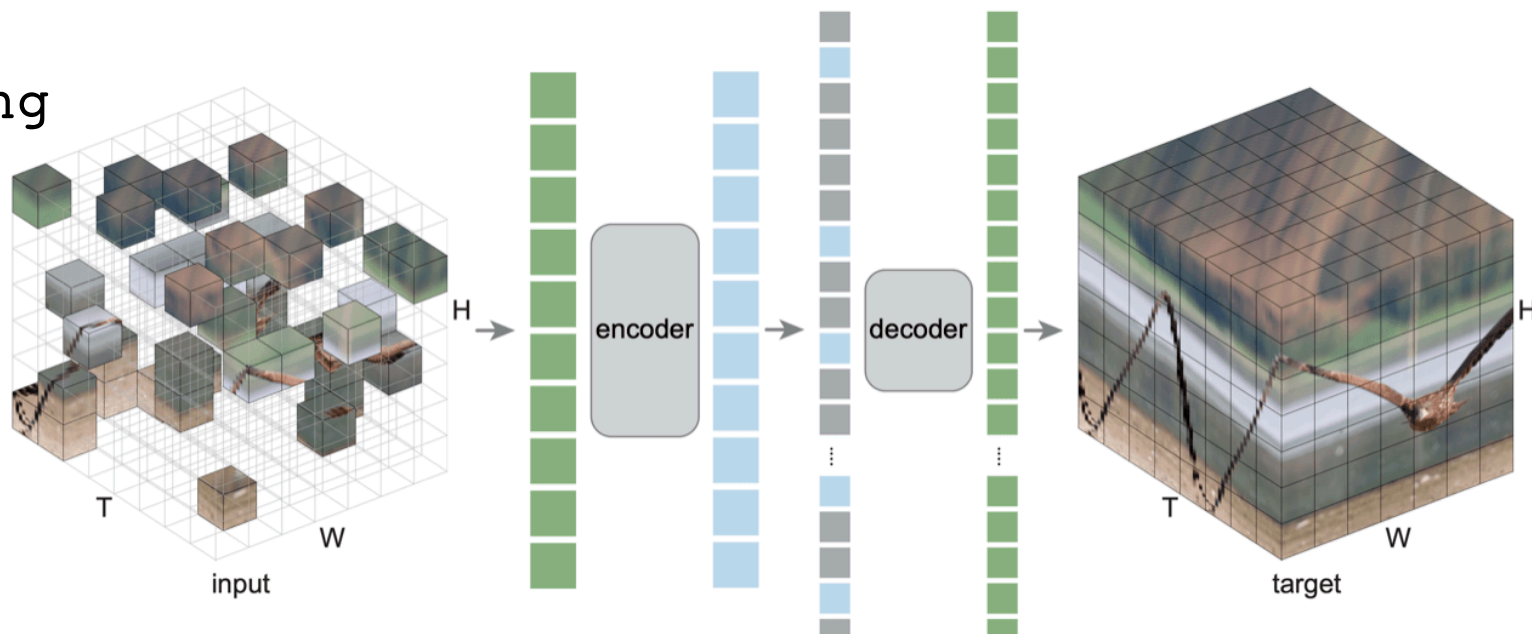


# 1. Masking



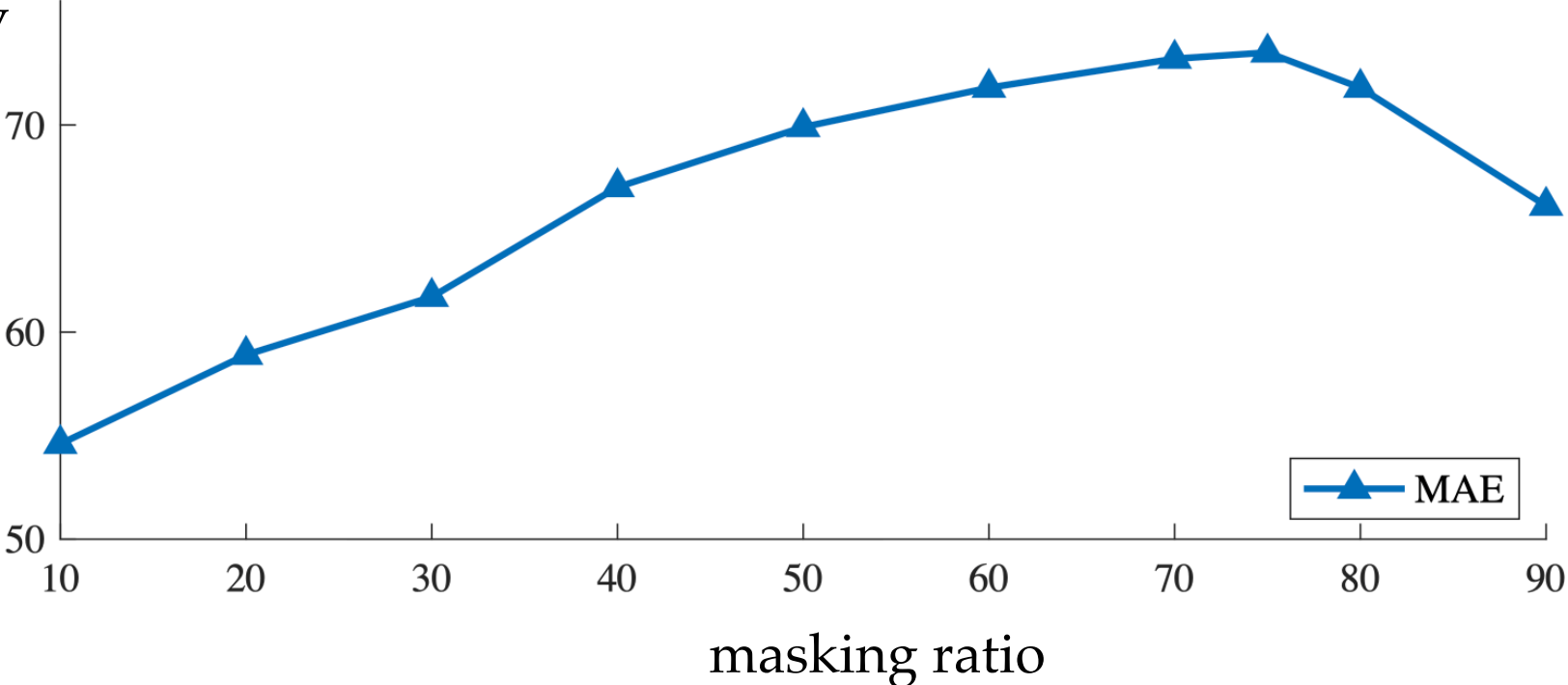
[Feichtenhofer, et al., "Masked Autoencoders As Spatiotemporal Learners", NeurIPS 2022]

# 1. Masking



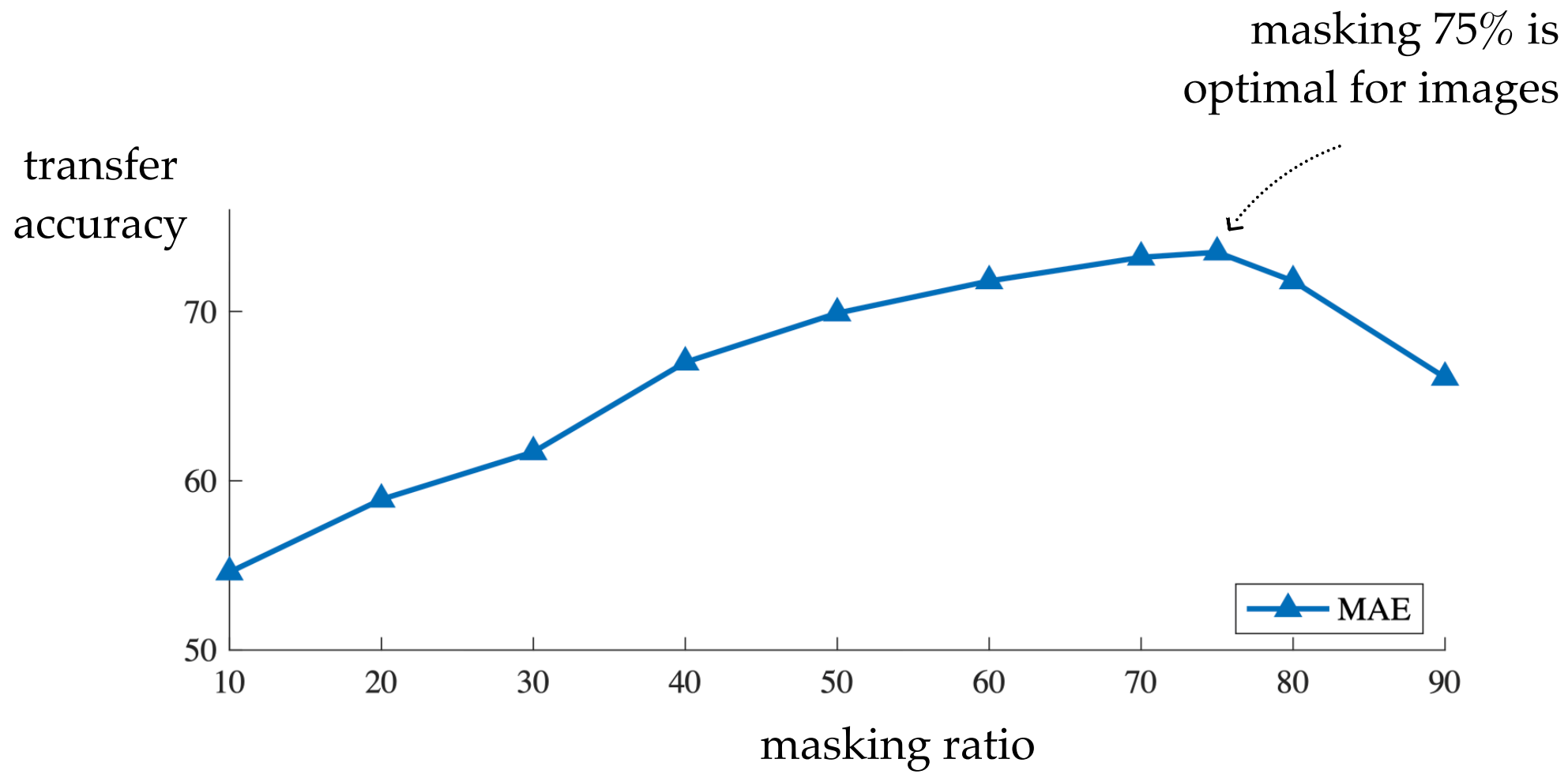
# 1. Masking

transfer  
accuracy



[He, et al. Masked Autoencoders Are Scalable Vision Learners, 2021]

# 1. Masking

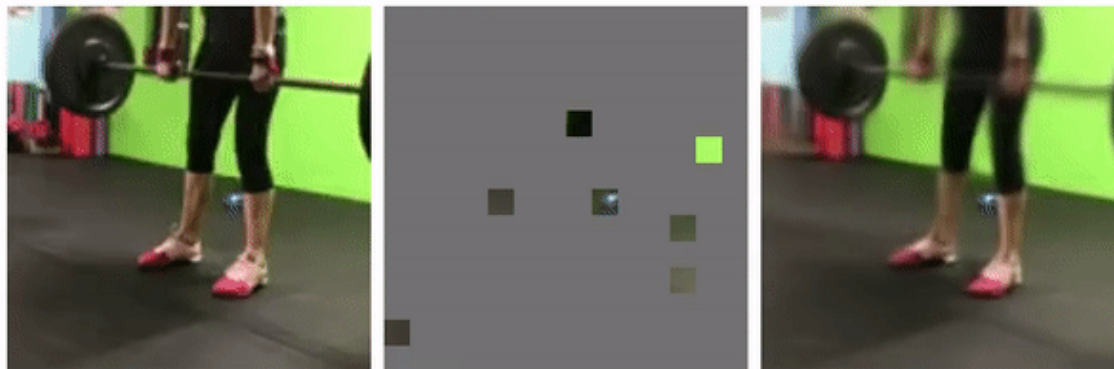


[He, et al. Masked Autoencoders Are Scalable Vision Learners, 2021]

# 1. Masking

# 1. Masking

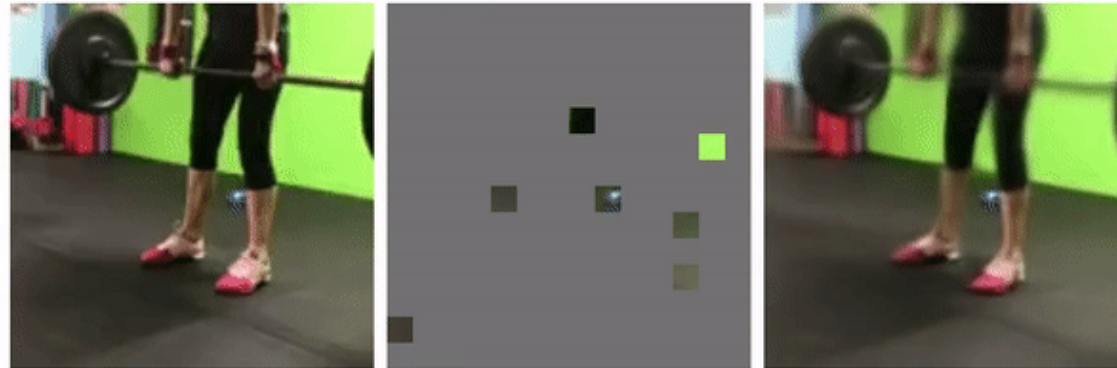
95% masked



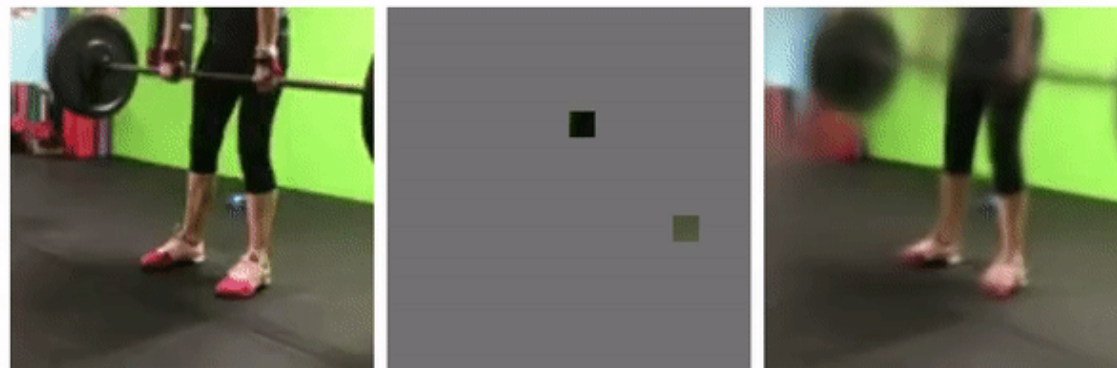
# 1. Masking

Similar empirical studies shows 15% as optimal for languages, and 95% for videos

95% masked



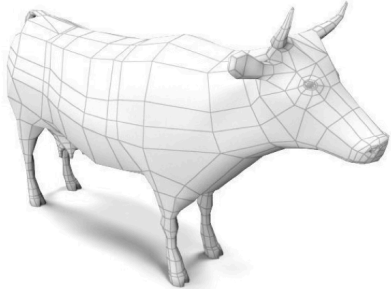
98% masked



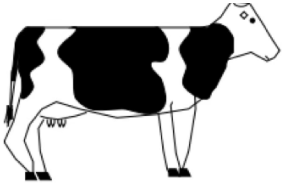
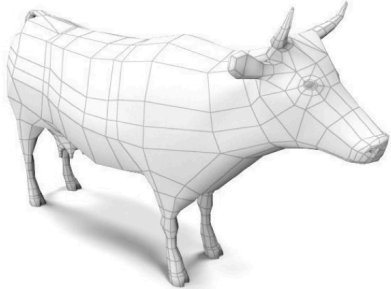
# 2. Contrastive learning



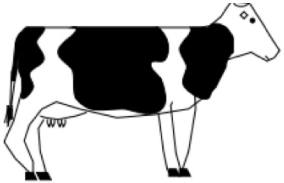
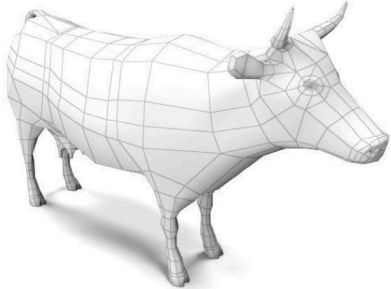
# 2. Contrastive learning



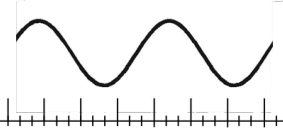
# 2. Contrastive learning



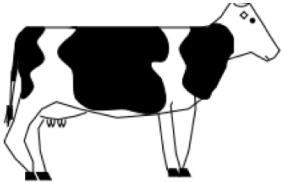
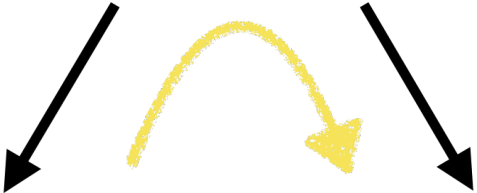
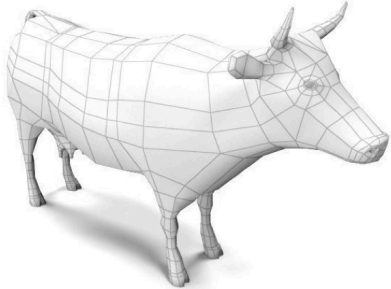
# 2. Contrastive learning



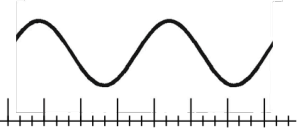
moo



# 2. Contrastive learning



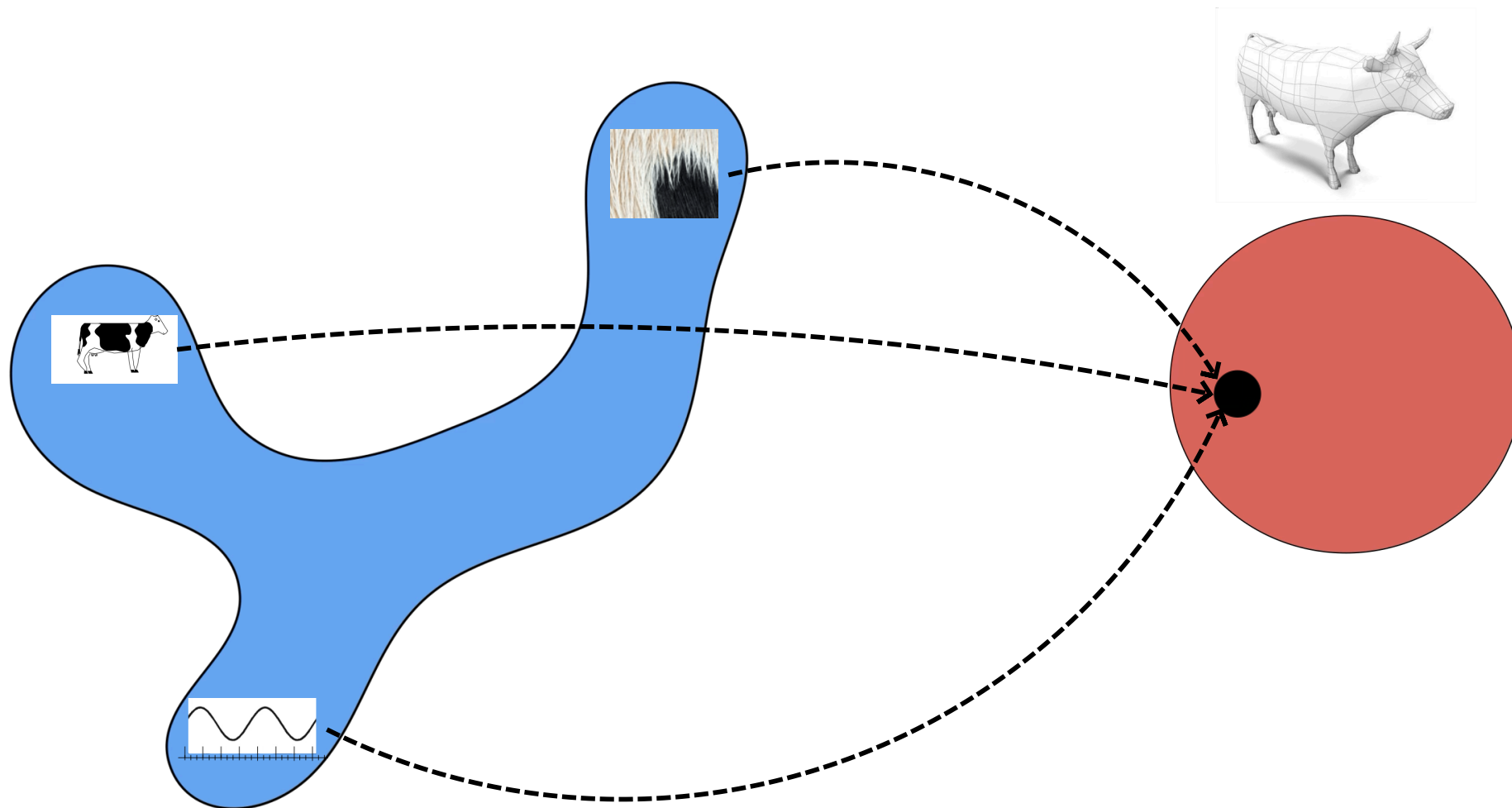
moo



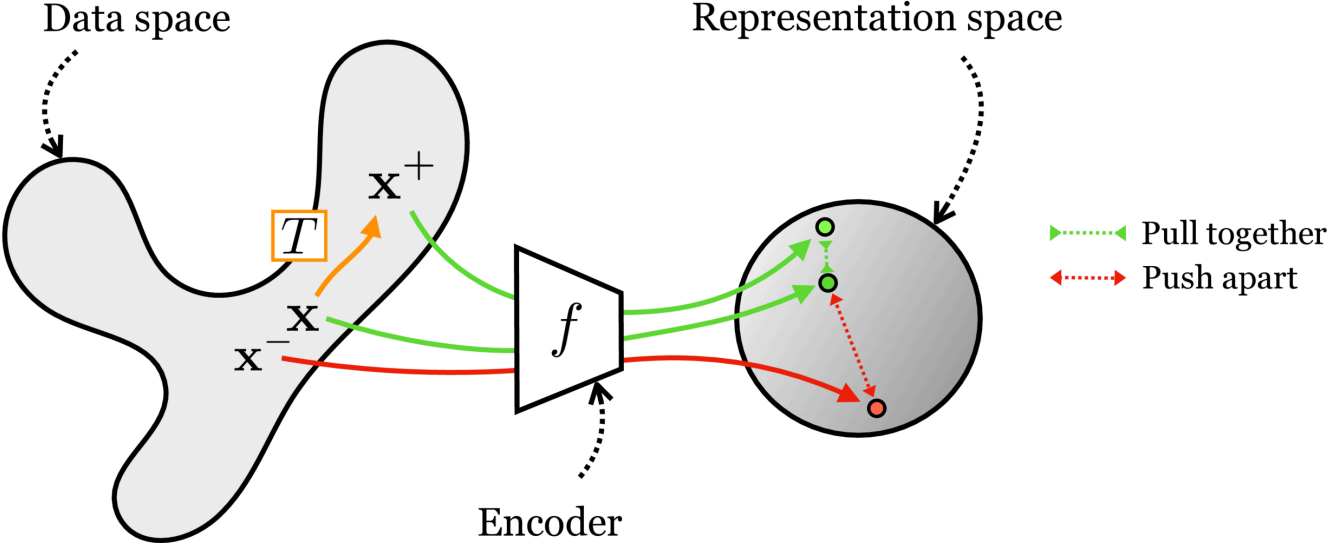
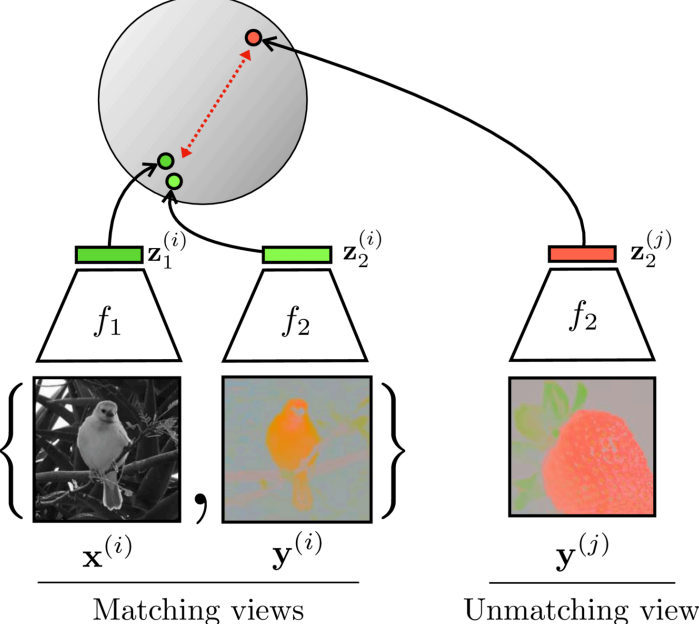
## 2. Contrastive learning

Observations

State



# 2. Contrastive learning



[images credit: visionbook.mit.edu]

## 2. Contrastive learning

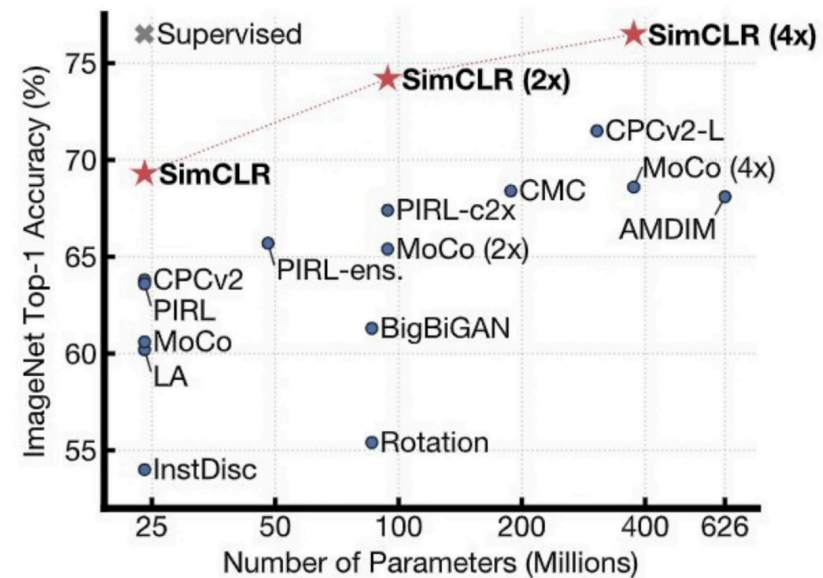
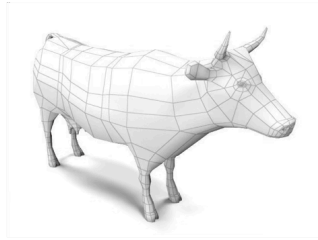


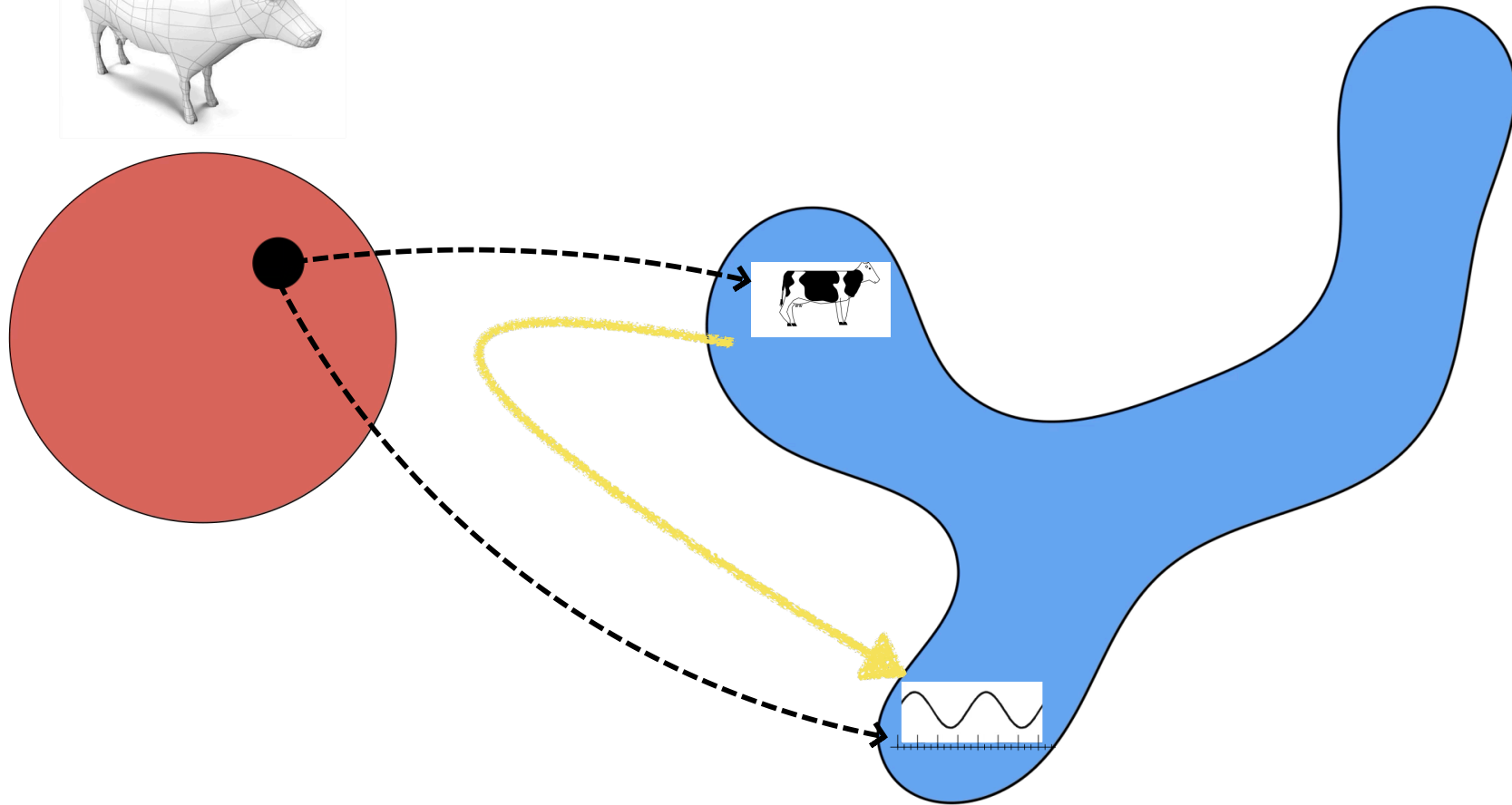
Figure 1. ImageNet top-1 accuracy of linear classifiers trained on representations learned with different self-supervised methods (pretrained on ImageNet). Our method, SimCLR, is shown in bold.

### 3. Multi-modality

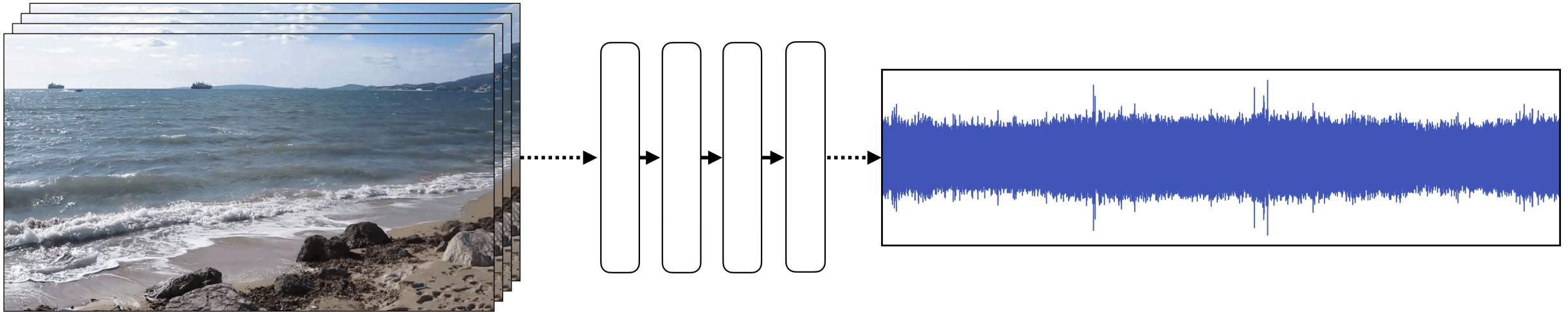
State



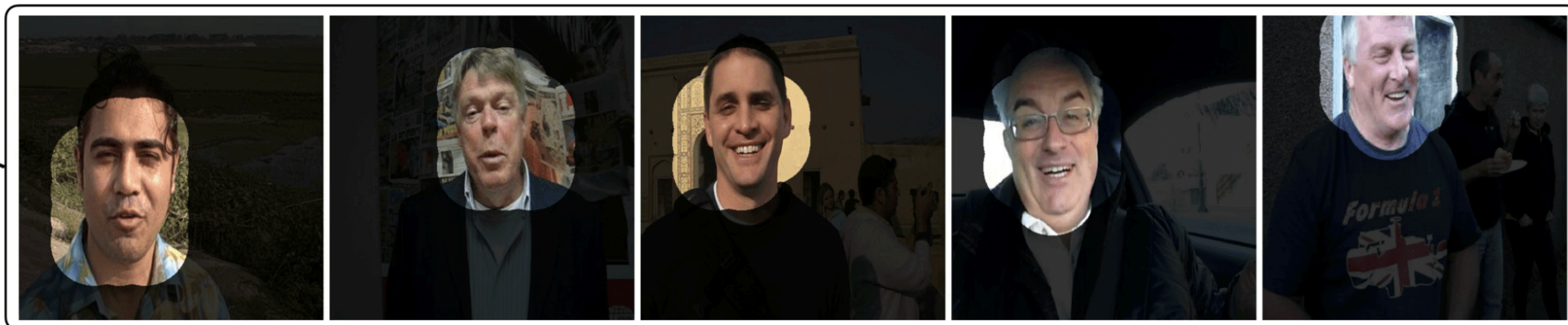
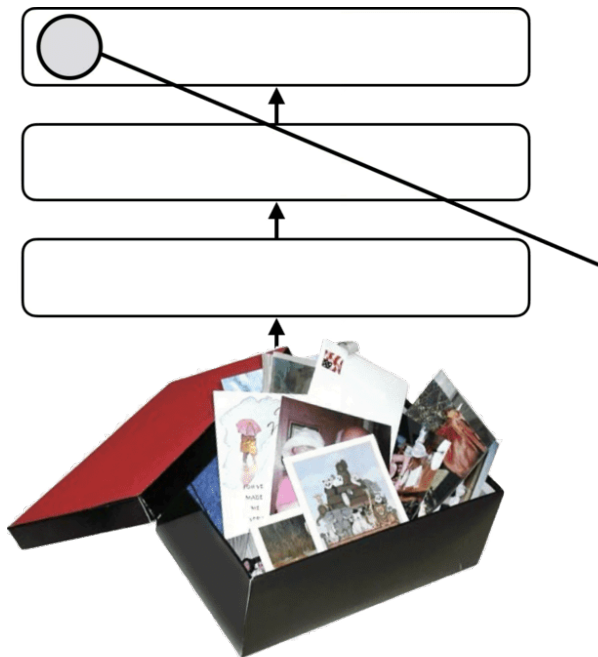
Observations



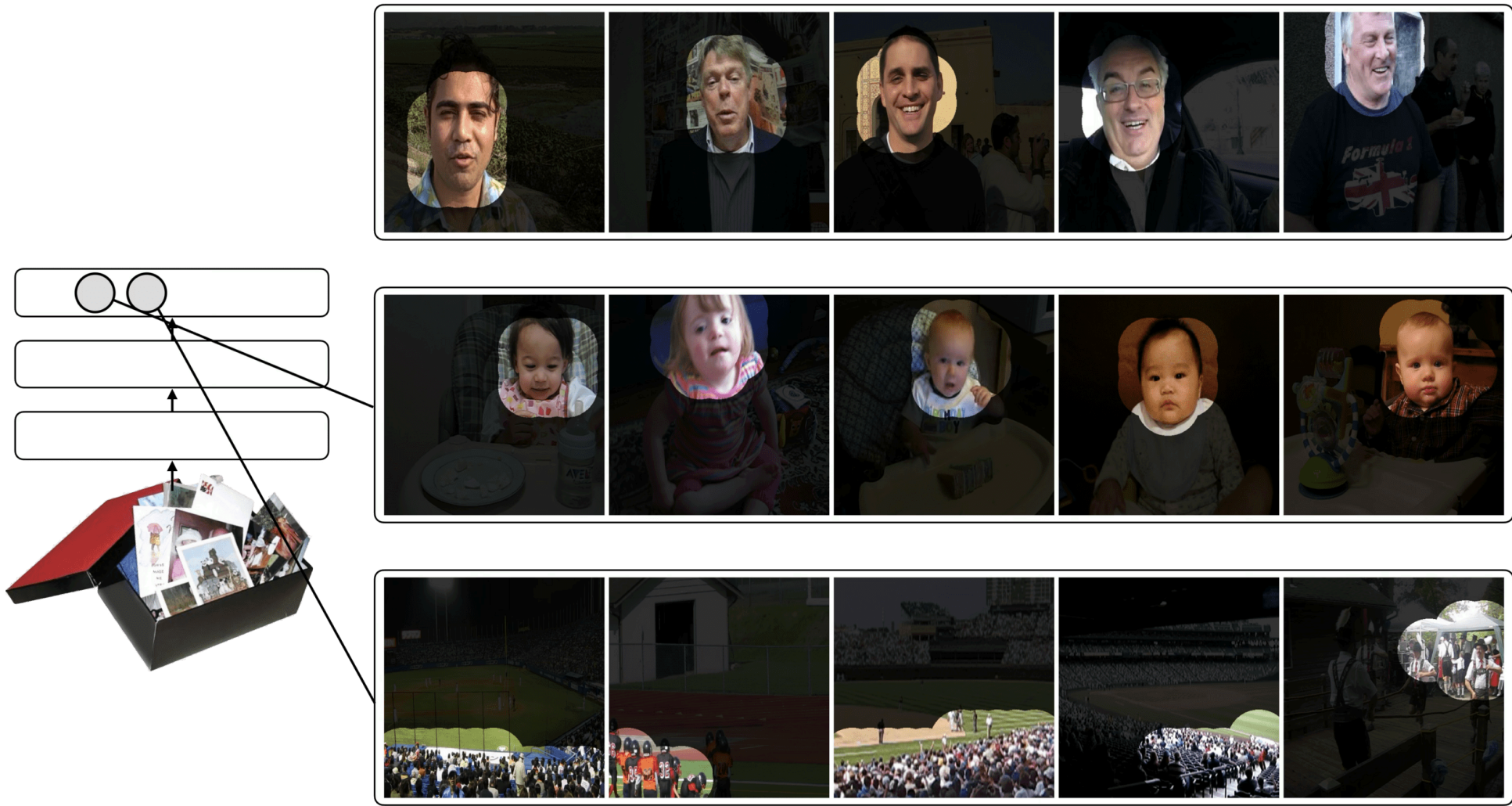
e.g. video, audio, images



What did the model learn?

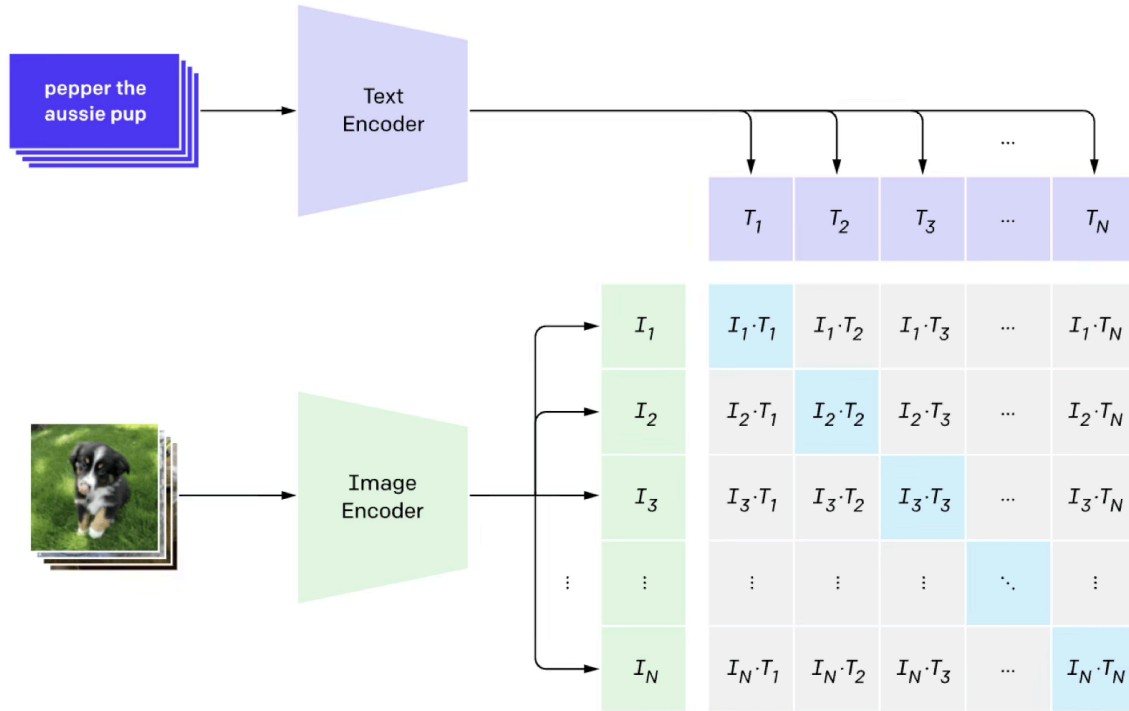


Strongest responses in dataset



e.g. image classification (done in the contrastive way)

### 1. Contrastive pre-training



```
# extract feature representations of each modality
```

```
I_f = image_encoder(I) #[n, d_i]
```

```
T_f = text_encoder(T)  #[n, d_t]
```

```
# joint multimodal embedding [n, d_e]
```

```
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
```

```
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)
```

```
# scaled pairwise cosine similarities [n, n]
```

```
logits = np.dot(I_e, T_e.T) * np.exp(t)
```

```
# symmetric loss function
```

```
labels = np.arange(n)
```

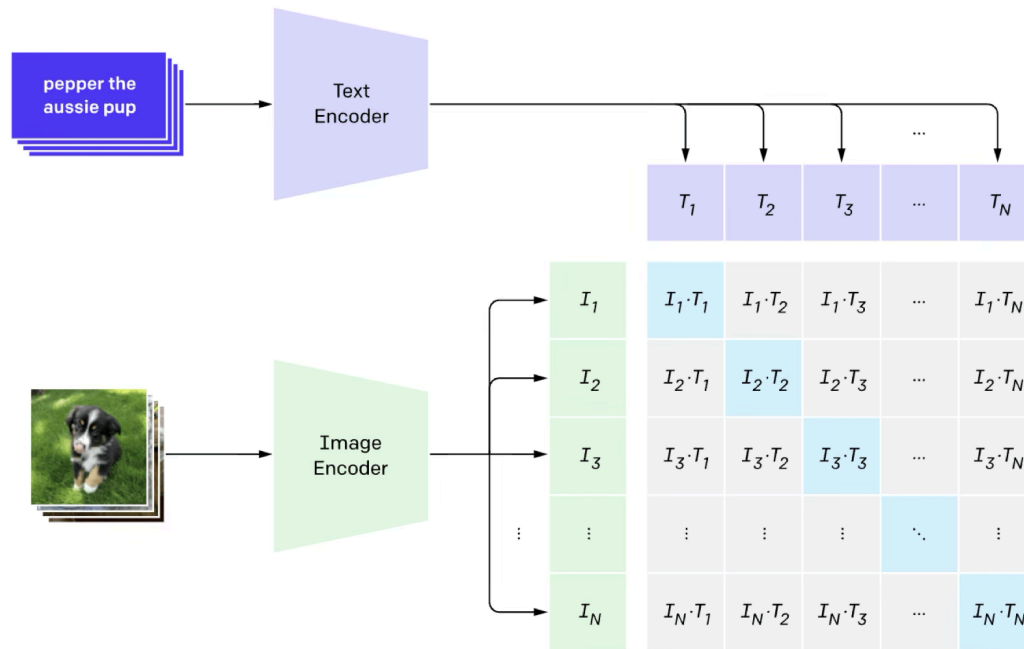
```
loss_i = cross_entropy_loss(logits, labels, axis=0)
```

```
loss_t = cross_entropy_loss(logits, labels, axis=1)
```

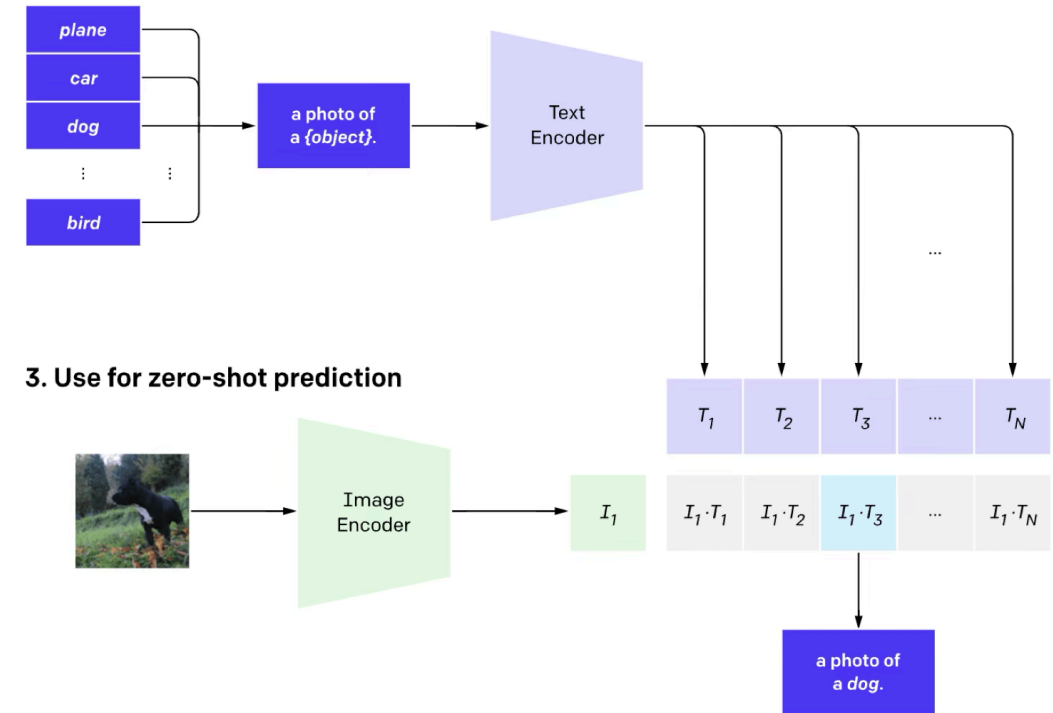
```
loss   = (loss_i + loss_t)/2
```

e.g. image classification (done in the contrastive way)

### 1. Contrastive pre-training



### 2. Create dataset classifier from label text



e.g Dall-E: text-image generation

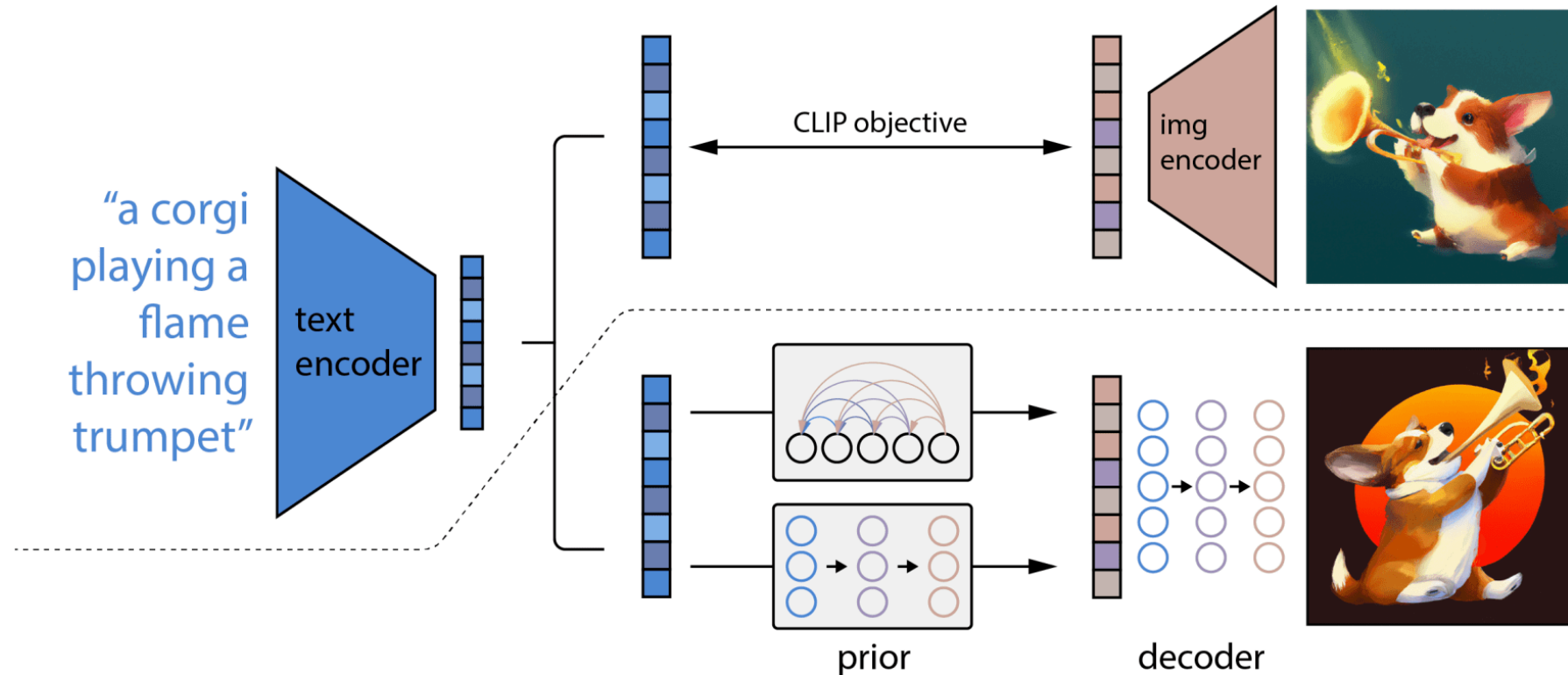


Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or diffusion prior to produce an image embedding, and then this embedding is used to condition a diffusion decoder which produces a final image. Note that the CLIP model is frozen during training of the prior and decoder.

# 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**



# Summary

- We looked at the mechanics of NN. Today we see they learn representations, just like our brains do.
- This is useful because representations transfer — they act as prior knowledge that enables quick learning on new tasks.
- Representations can also be learned without labels, e.g. as we do in unsupervised, or self-supervised learning. This is great since labels are expensive and limiting.
- Without labels there are many ways to learn representations:
  - representations as compressed codes, auto-encoder with bottleneck
  - (representations that are predictive of their context)
  - (representations that are shared across sensory modalities)