

Lecture 11: Representation Learning

Today

- What representations do neural nets learn?
- Transfer learning
- Unsupervised learning

Observed image



Drawn from memory



[Bartlett, 1932]

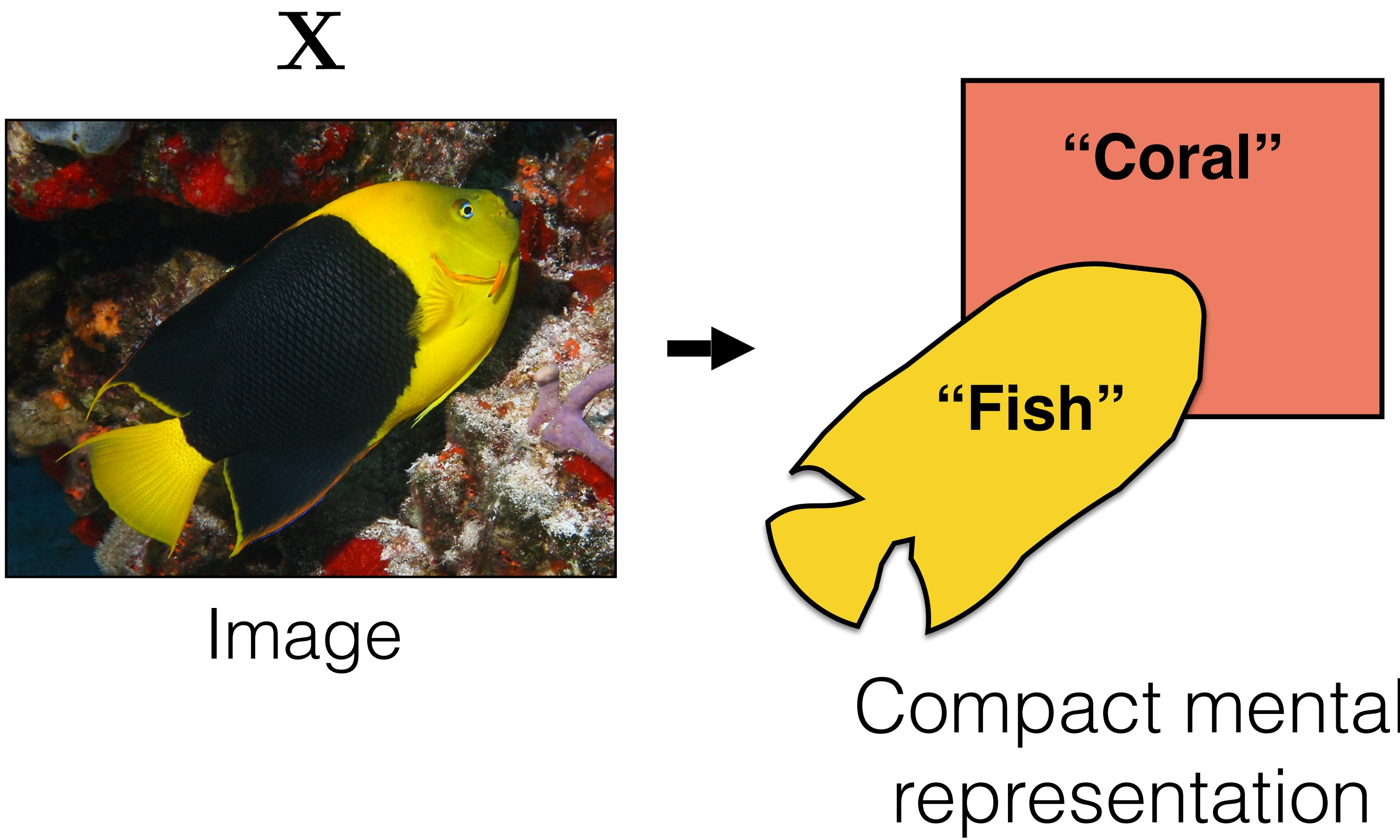
[Intraub & Richardson, 1989]



"I stand at the window and see a house, trees, sky. Theoretically I might say there were 327 brightnesses and nuances of colour. Do I have "327"? No. I have sky, house, and trees."

— Max Wertheimer, 1923

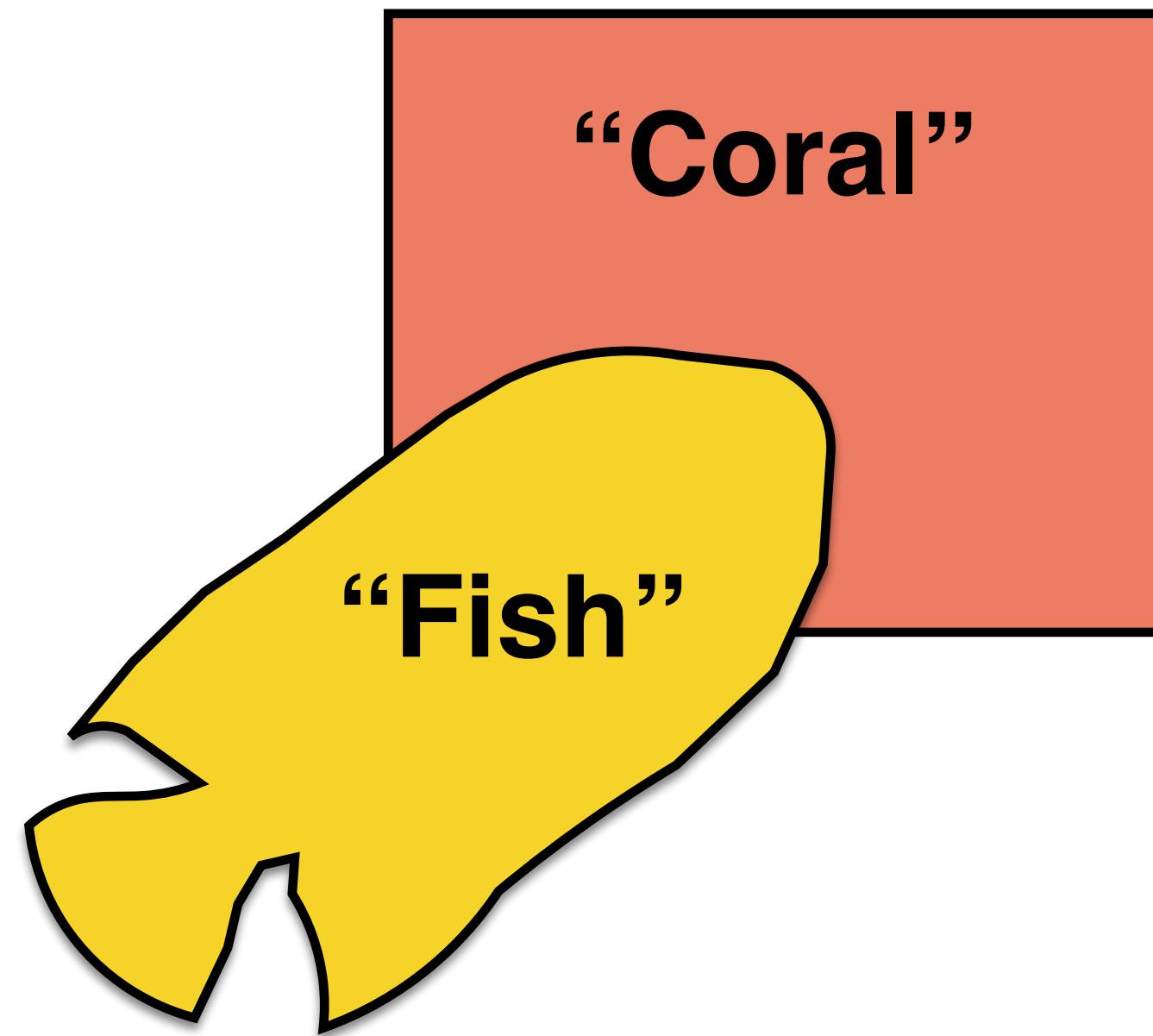
Representation learning



Representation learning

Good representations are:

1. Compact (*minimal*)
2. Explanatory (*sufficient*)
3. Disentangled (*independent factors*)
4. Interpretable
5. *Make subsequent problem solving easy*



8

[See “Representation Learning”, Bengio 2013, for more commentary]

Transfer learning

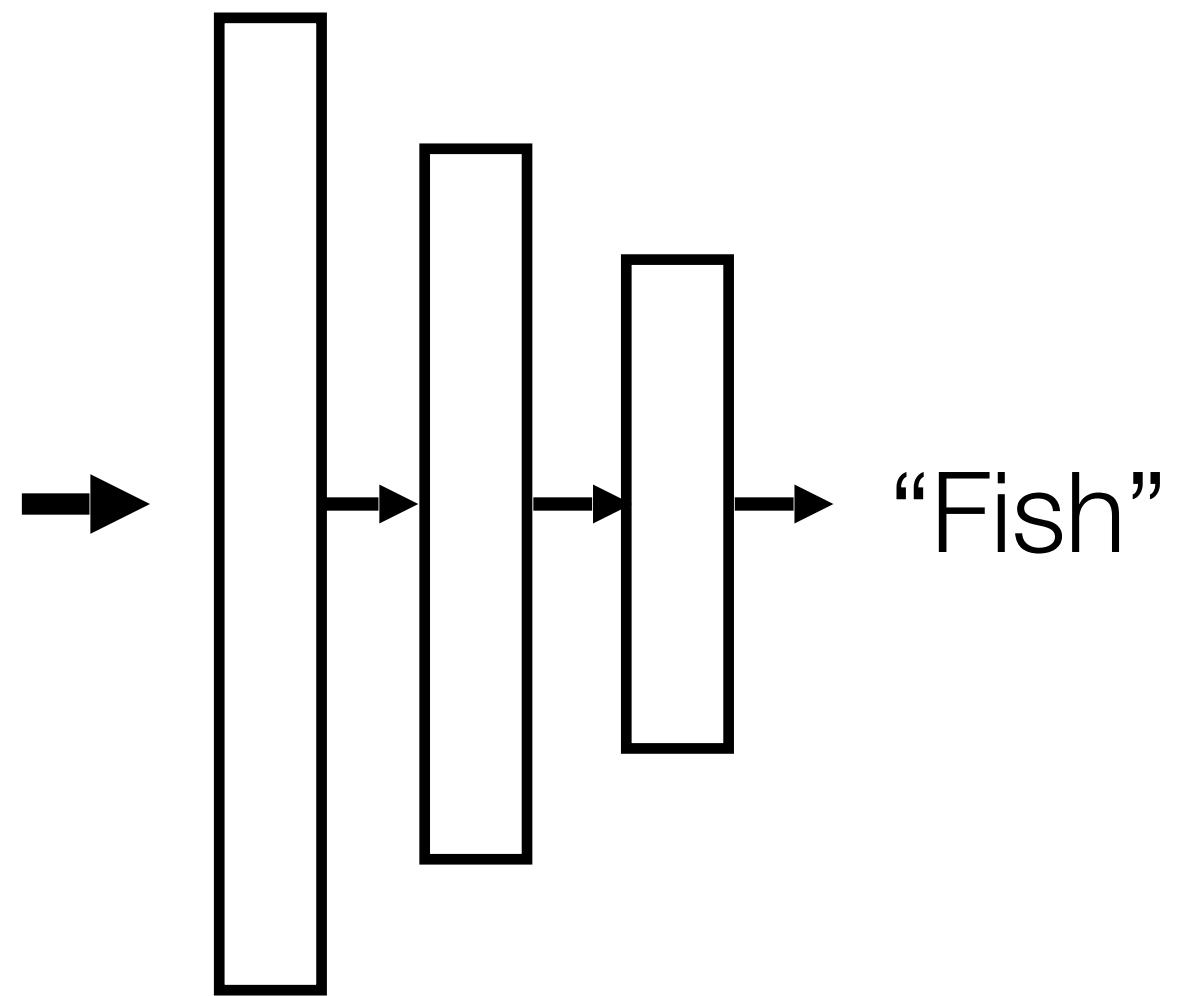
“Generally speaking, a good representation is one that makes a subsequent learning task easier.” – *Deep Learning*, Goodfellow et al. 2016



?

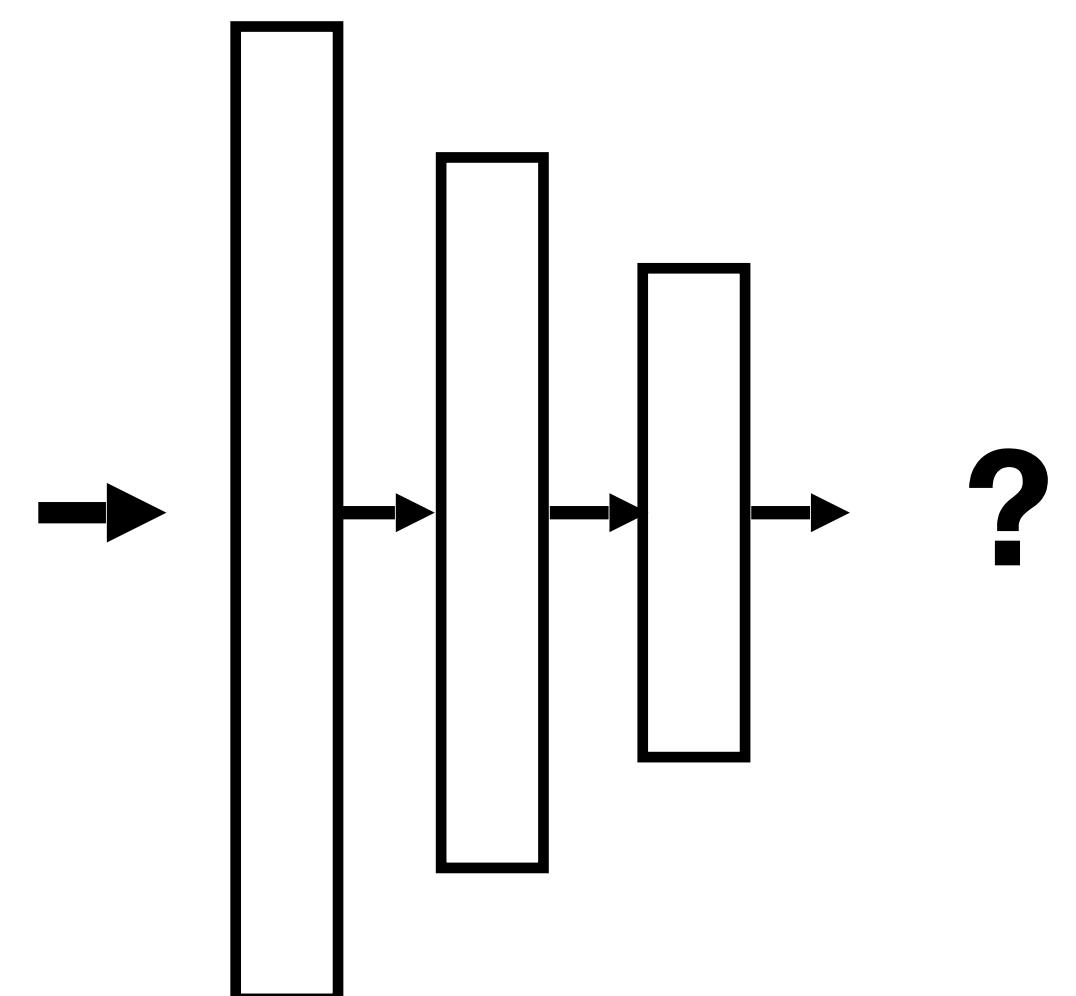
Training

Object recognition



Testing

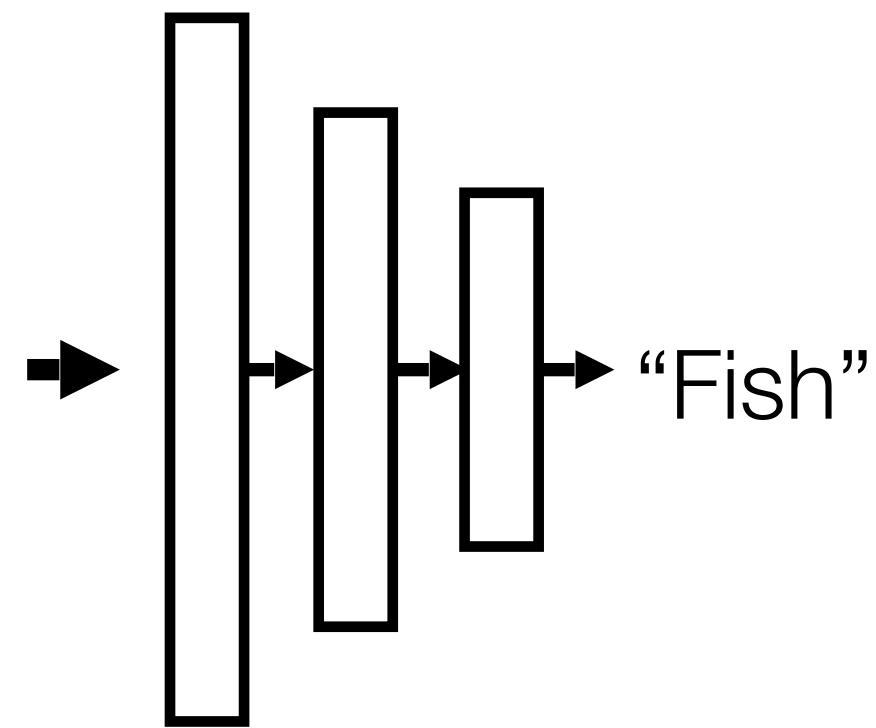
Place recognition



Often, what we will be “tested” on is to learn to do a new thing.

Pretraining

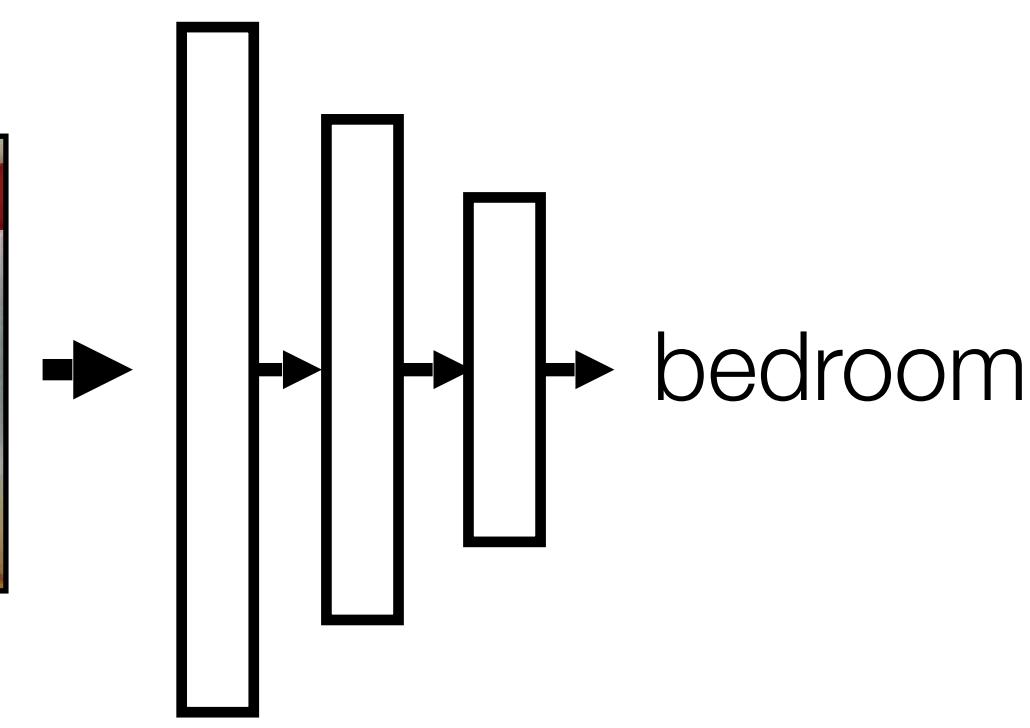
Object recognition



A lot of data

Finetuning

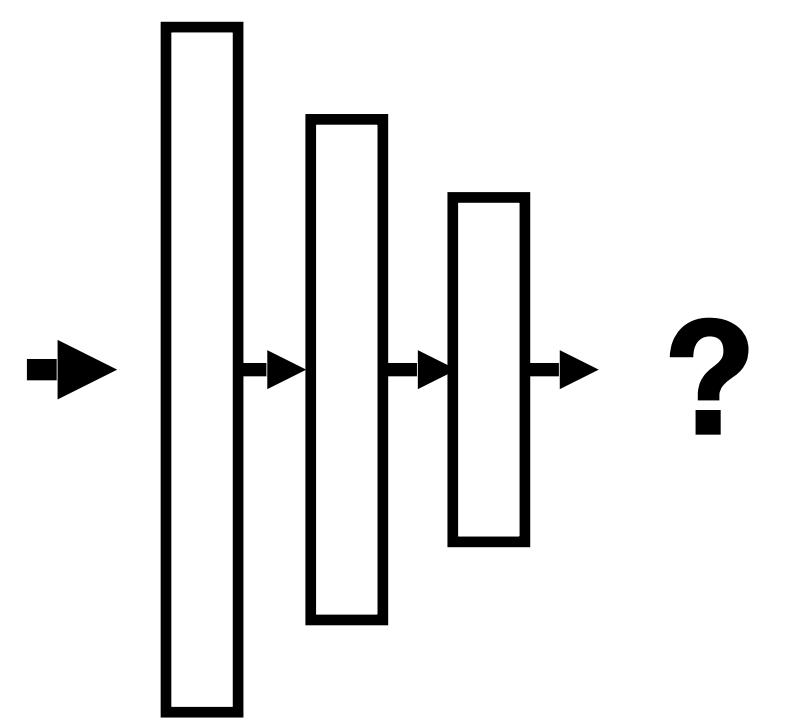
Place recognition



A little data

Testing

Place recognition

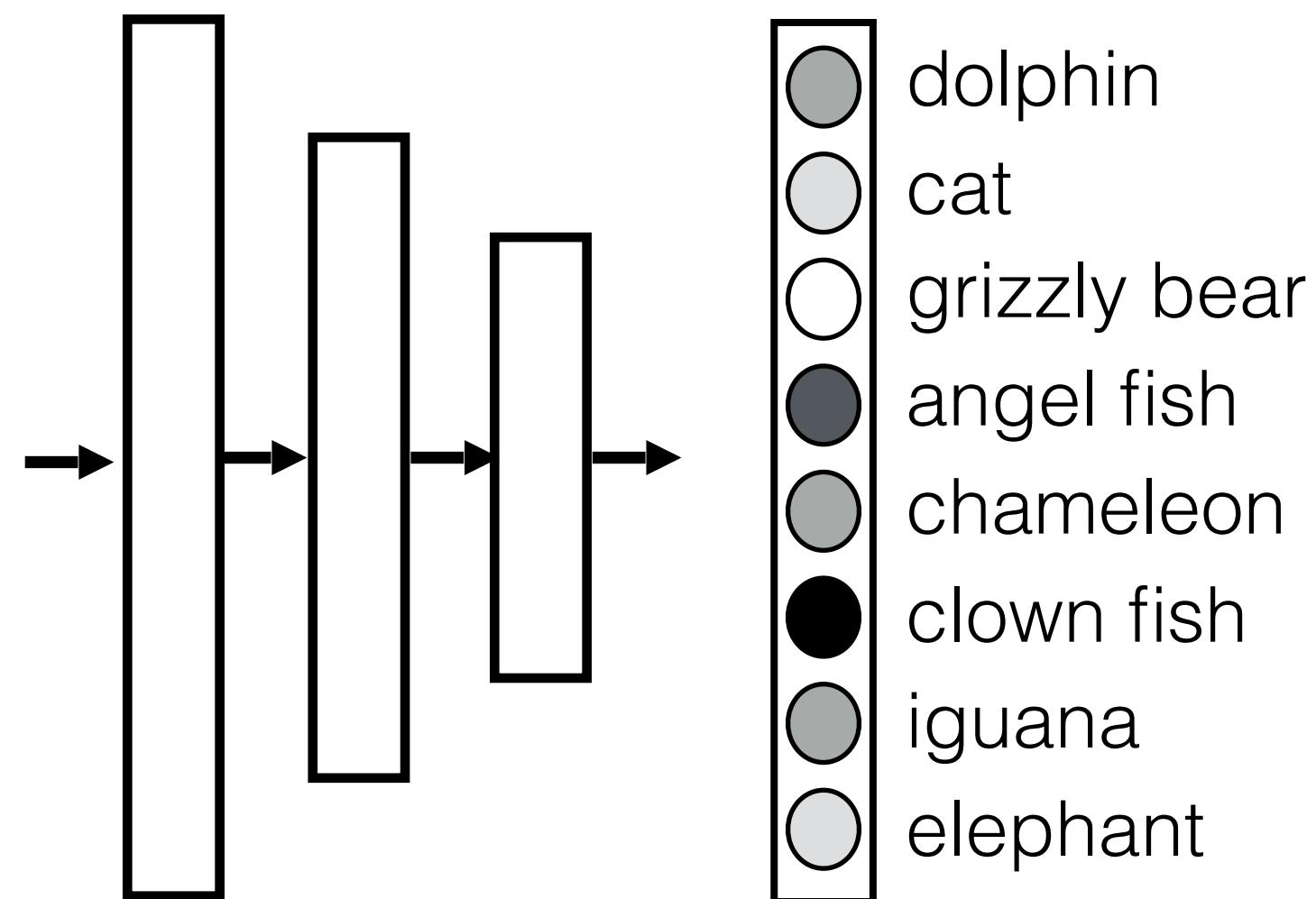


Finetuning starts with the representation learned on a previous task, and adapts it to perform well on a new task.

Finetuning

Pretraining

Object recognition



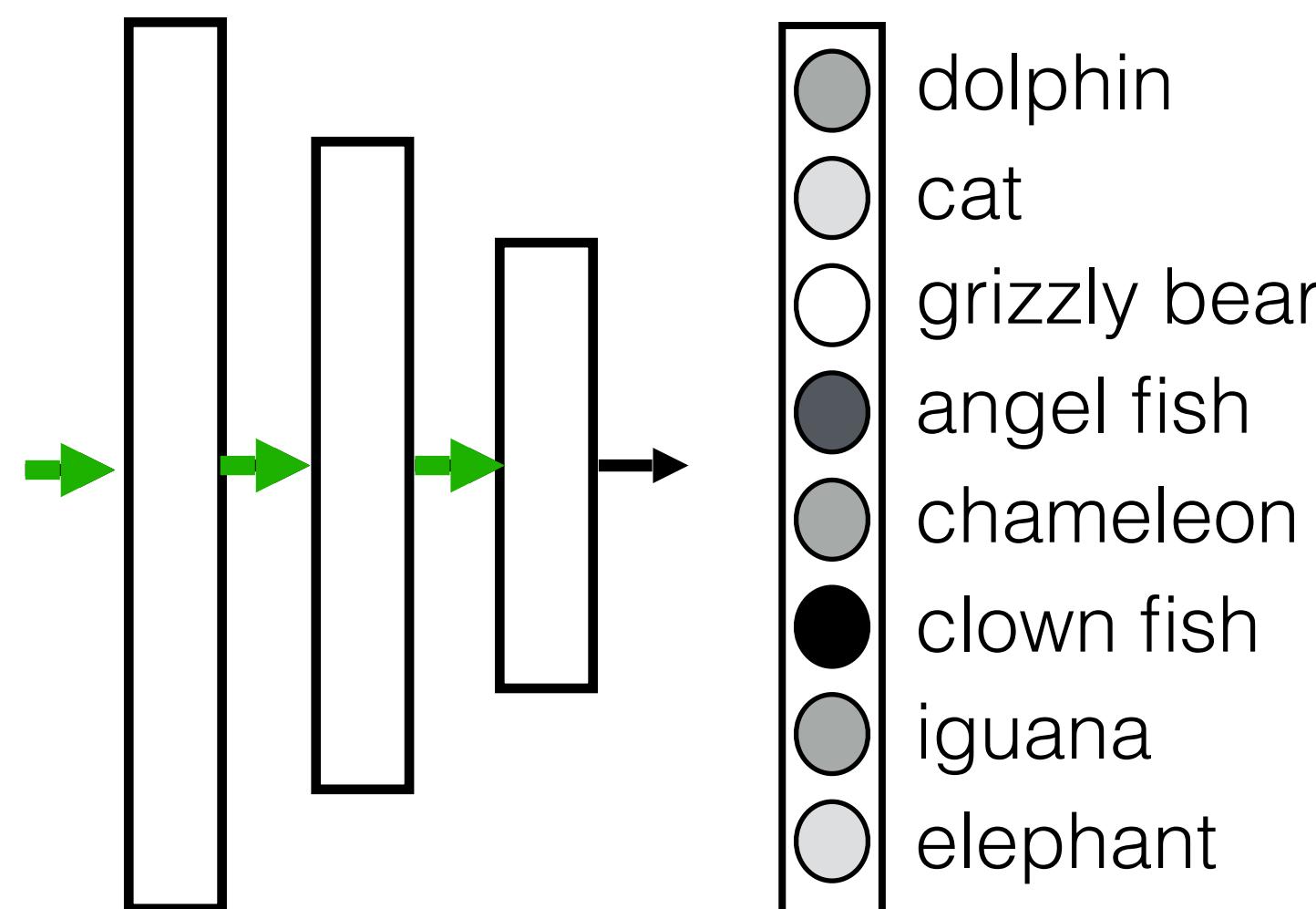
Finetuning

Place recognition

Finetuning

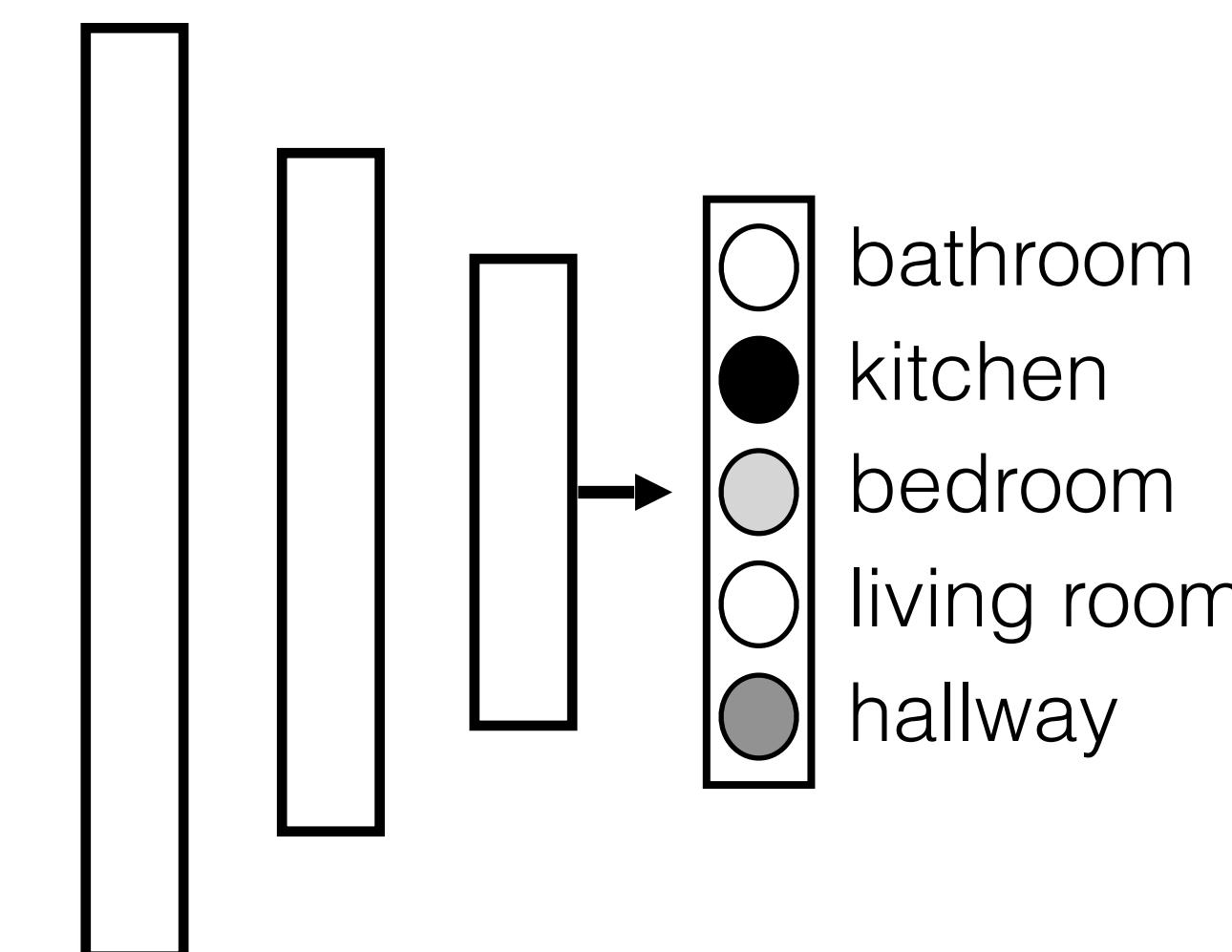
Pretraining

Object recognition



Finetuning

Place recognition



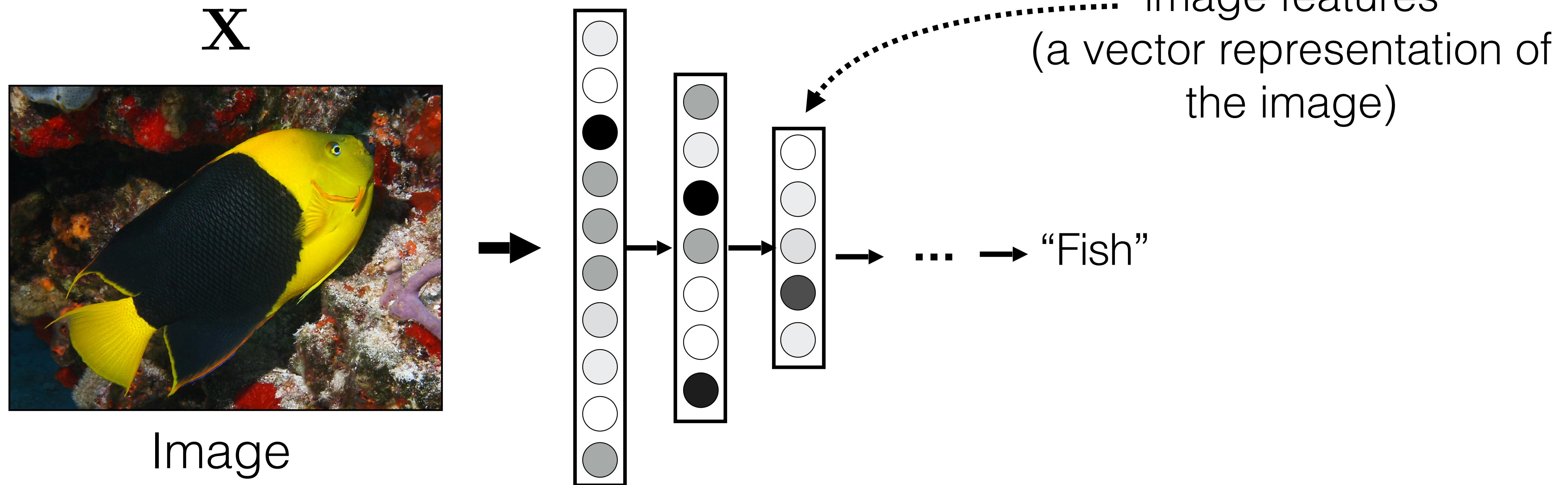
The “learned representation” is just the weights and biases, so that’s what we transfer

Finetuning

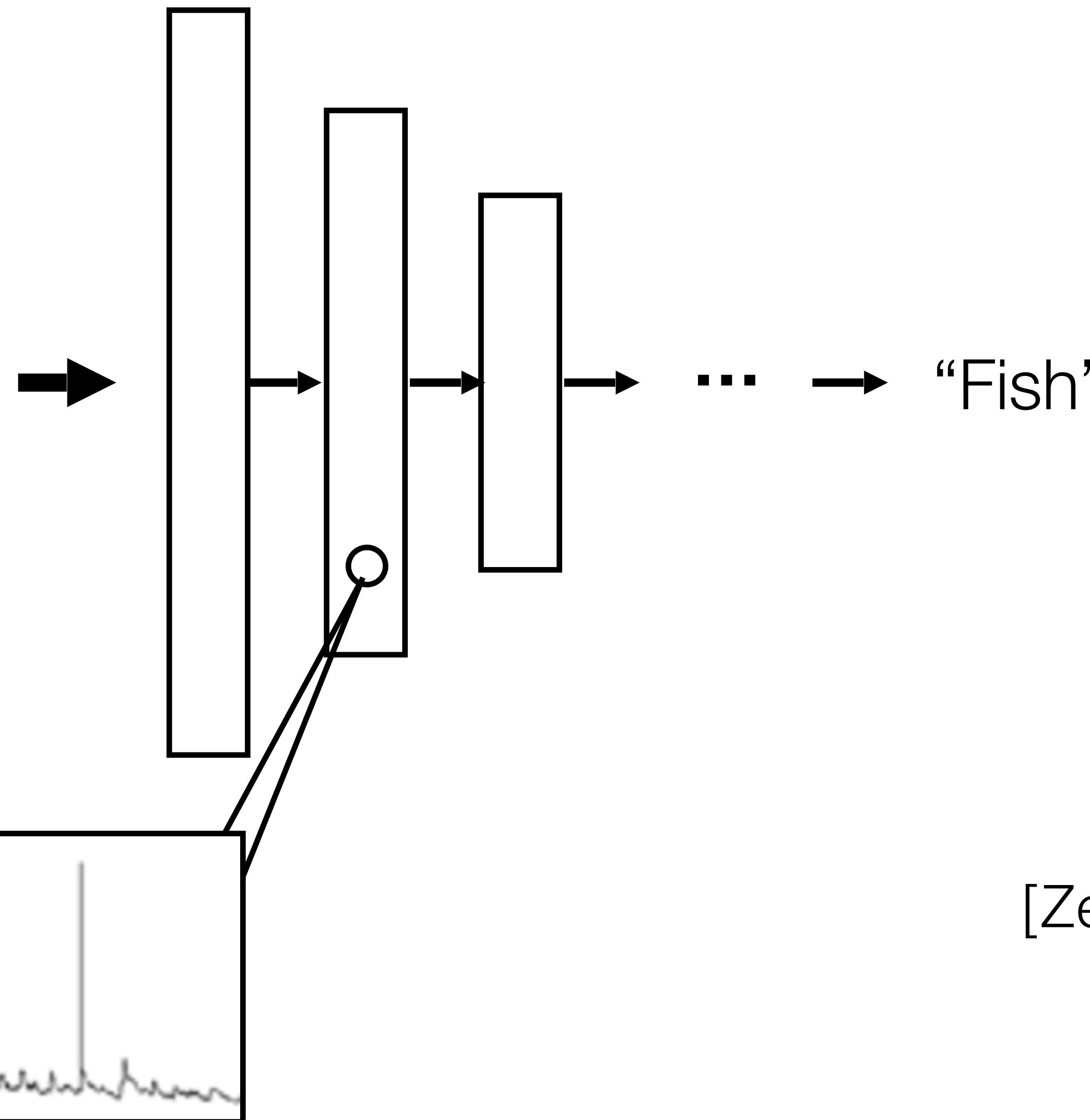
- Pretrain a network on task A (often object recognition), resulting in parameters **W**
- Initialize a second network with some or all of **W**
- Train the second network on task B, resulting in parameters **W'**
- Why would we expect this to work?

Visualizing representations

What do deep nets internally learn?



Deep net “electrophysiology”



[Zeiler & Fergus, ECCV 2014]

[Zhou et al., ICLR 2015]

Visualizing and Understanding CNNs

[Zeiler and Fergus, 2014]

Gabor-like filters learned by **layer 1**

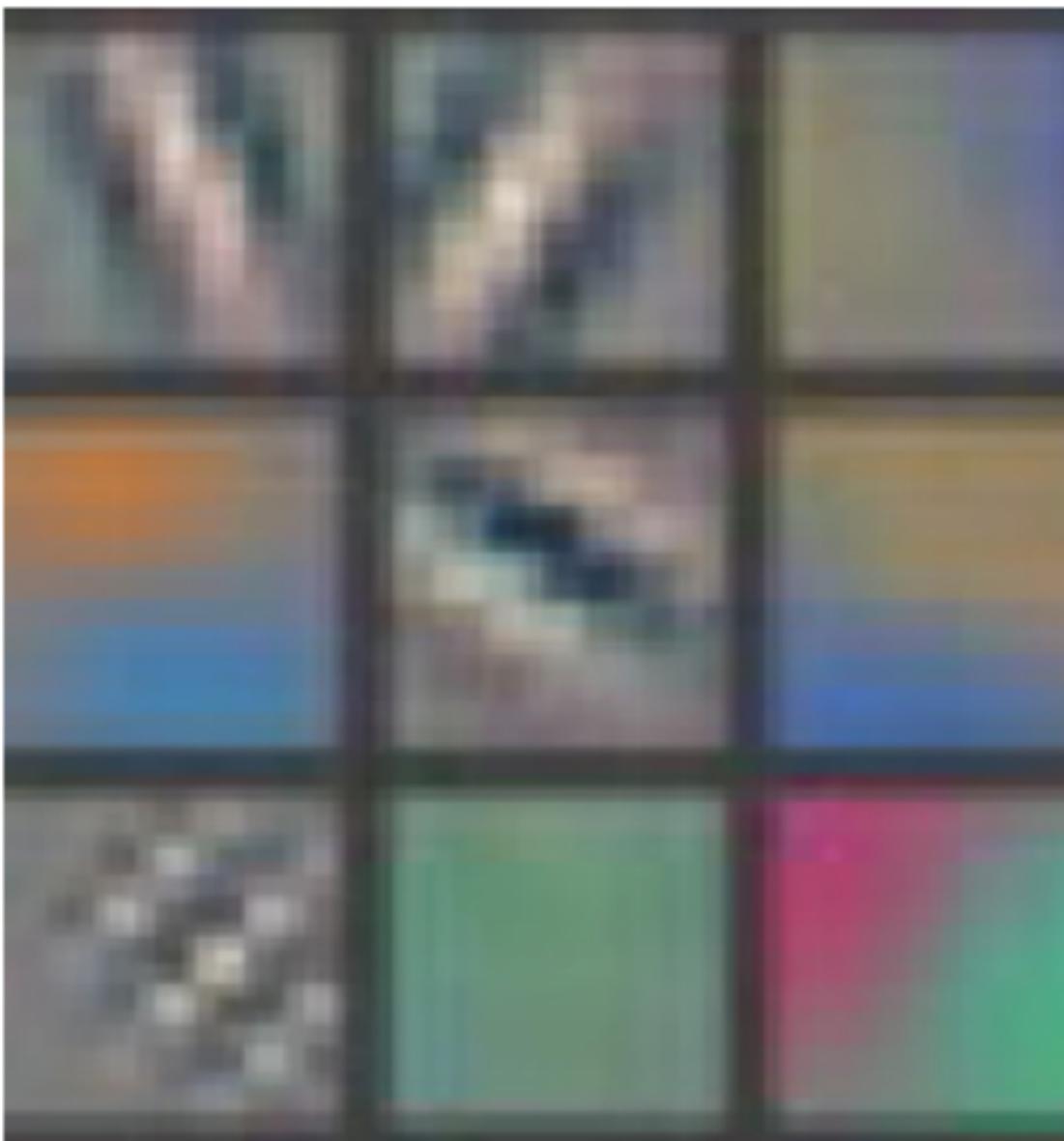
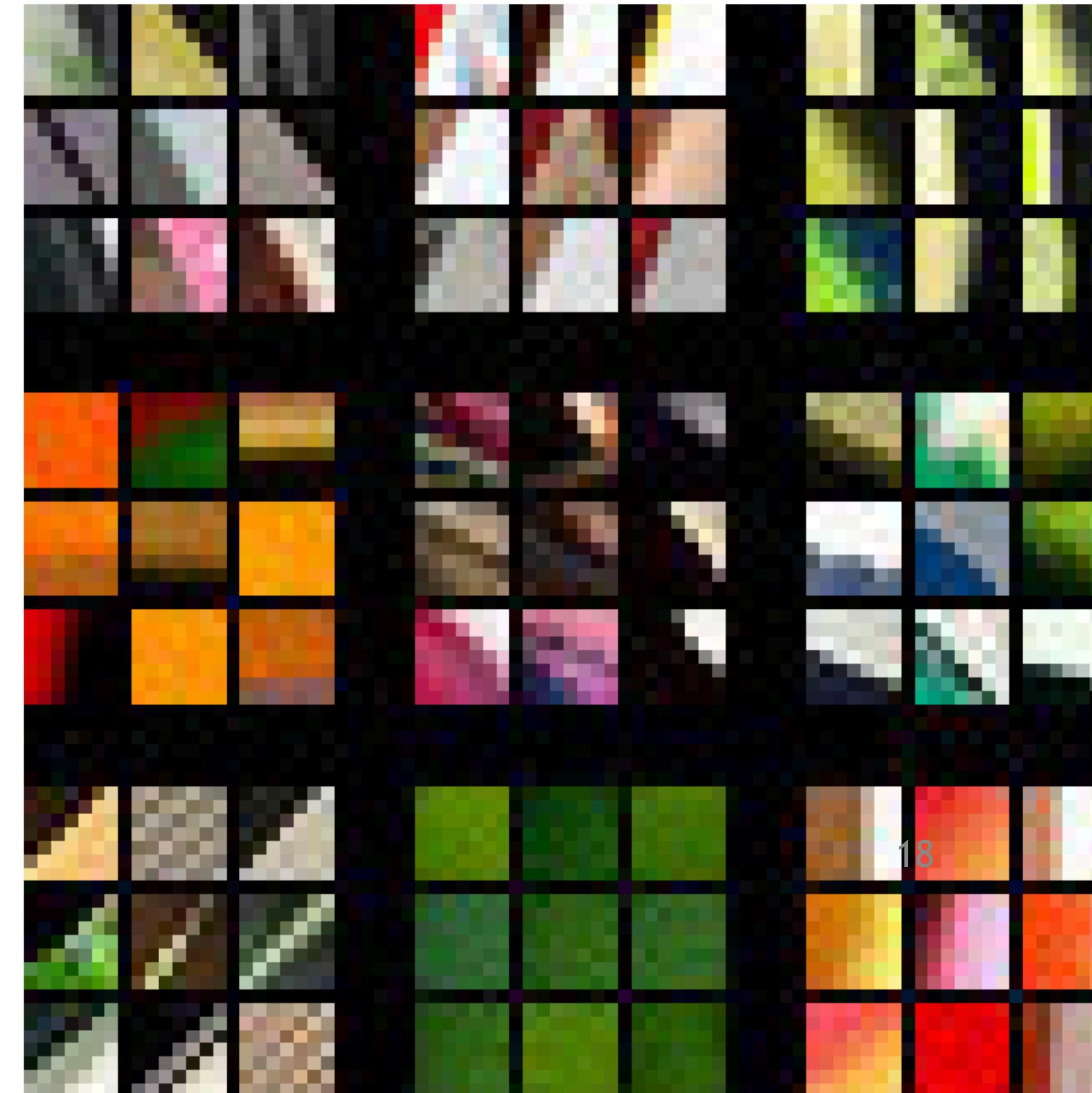
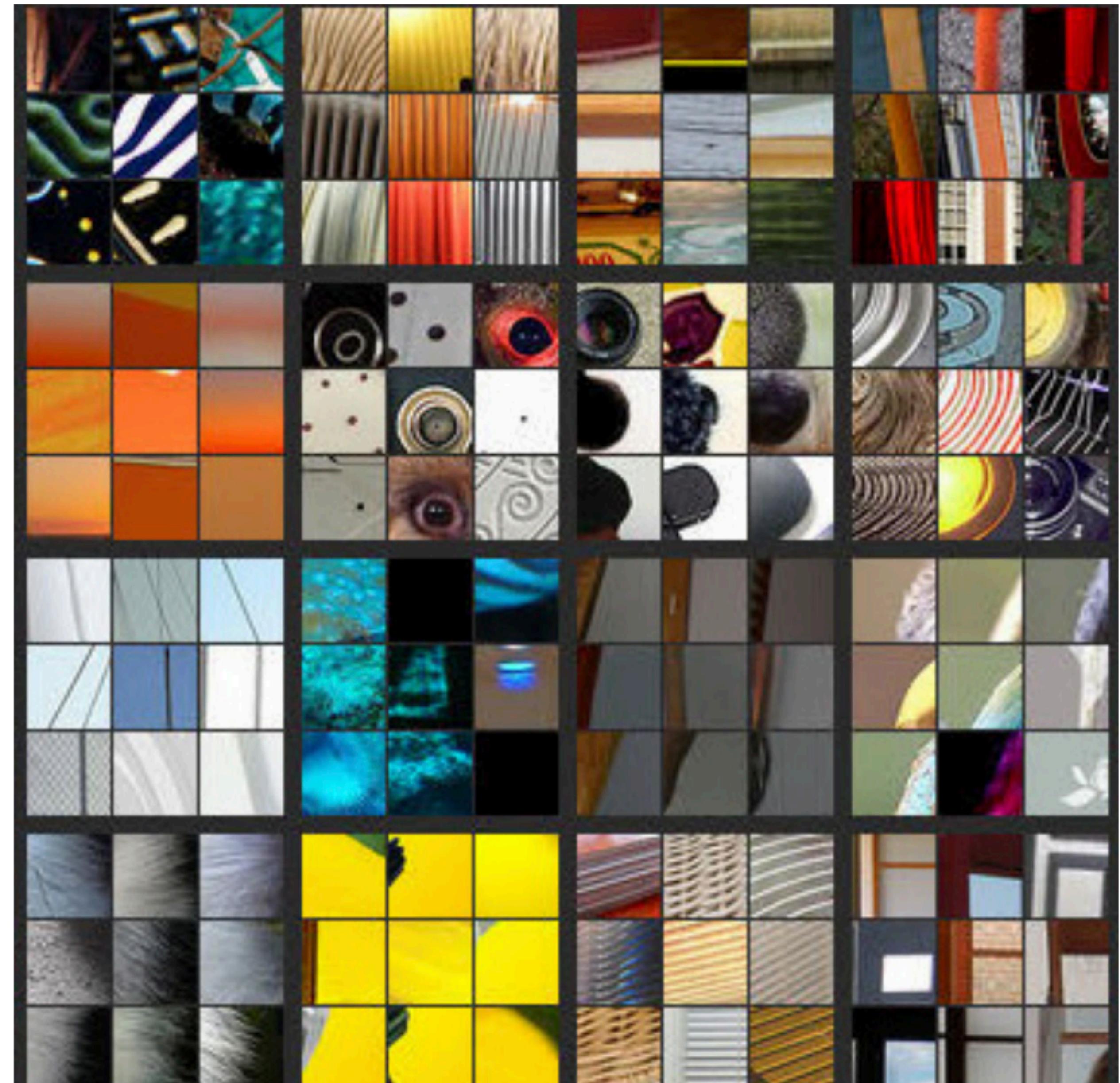


Image patches that activate each of the
layer 1 filters most strongly



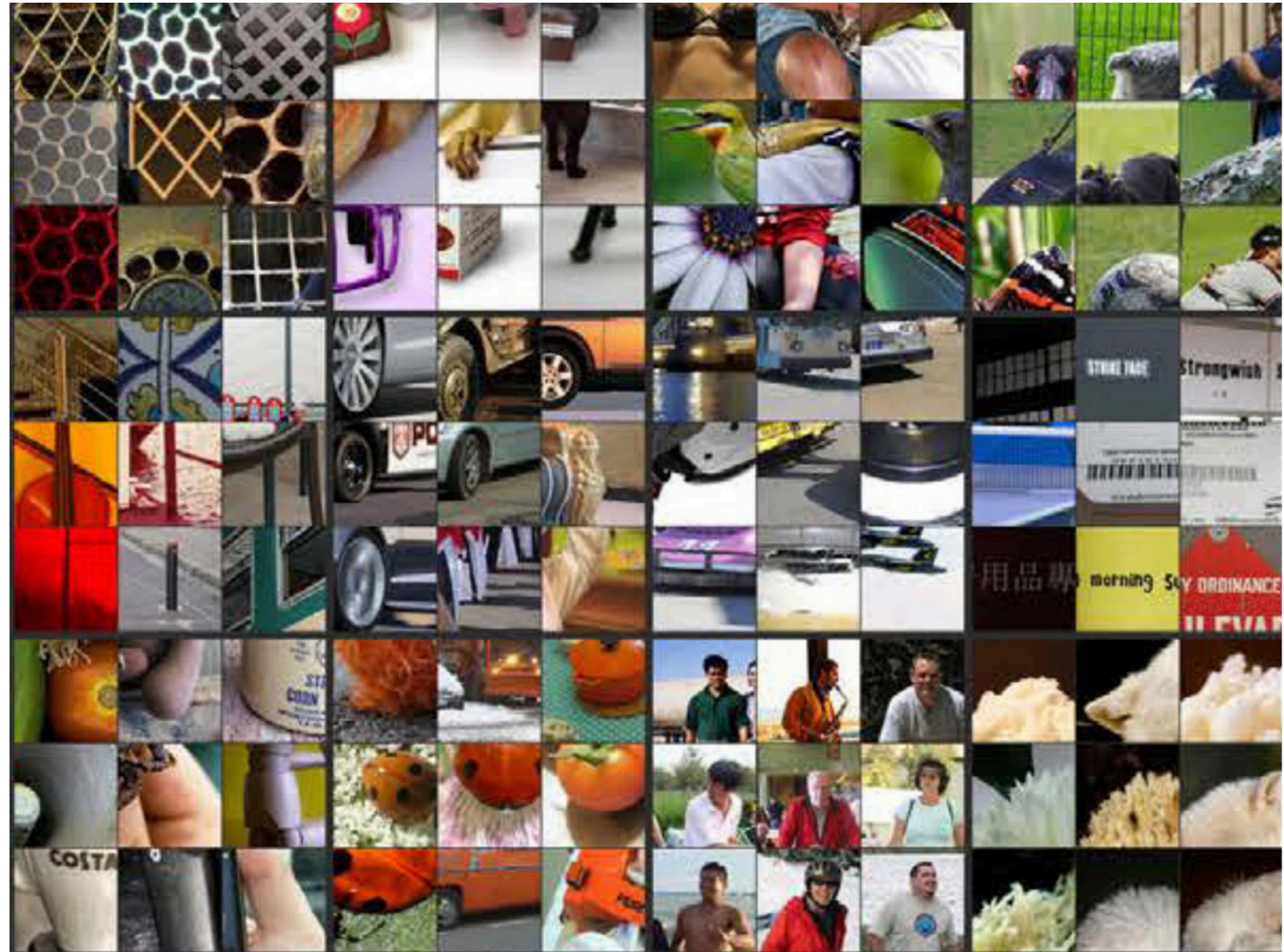
[Zeiler and Fergus, 2014]

Image patches that activate each of the **layer 2** neurons most strongly



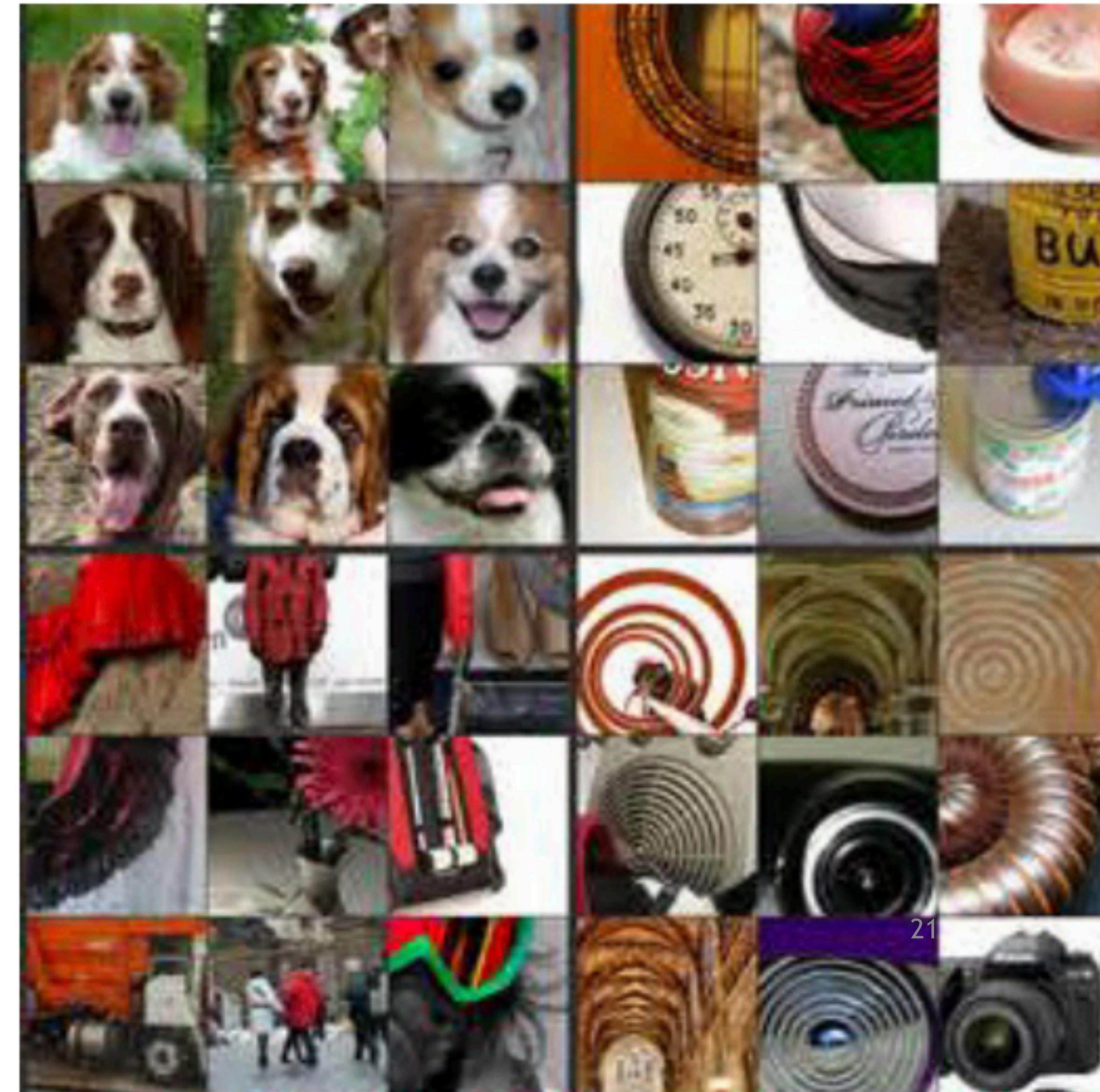
[Zeiler and Fergus, 2014]

Image patches that activate each of the **layer 3** neurons most strongly



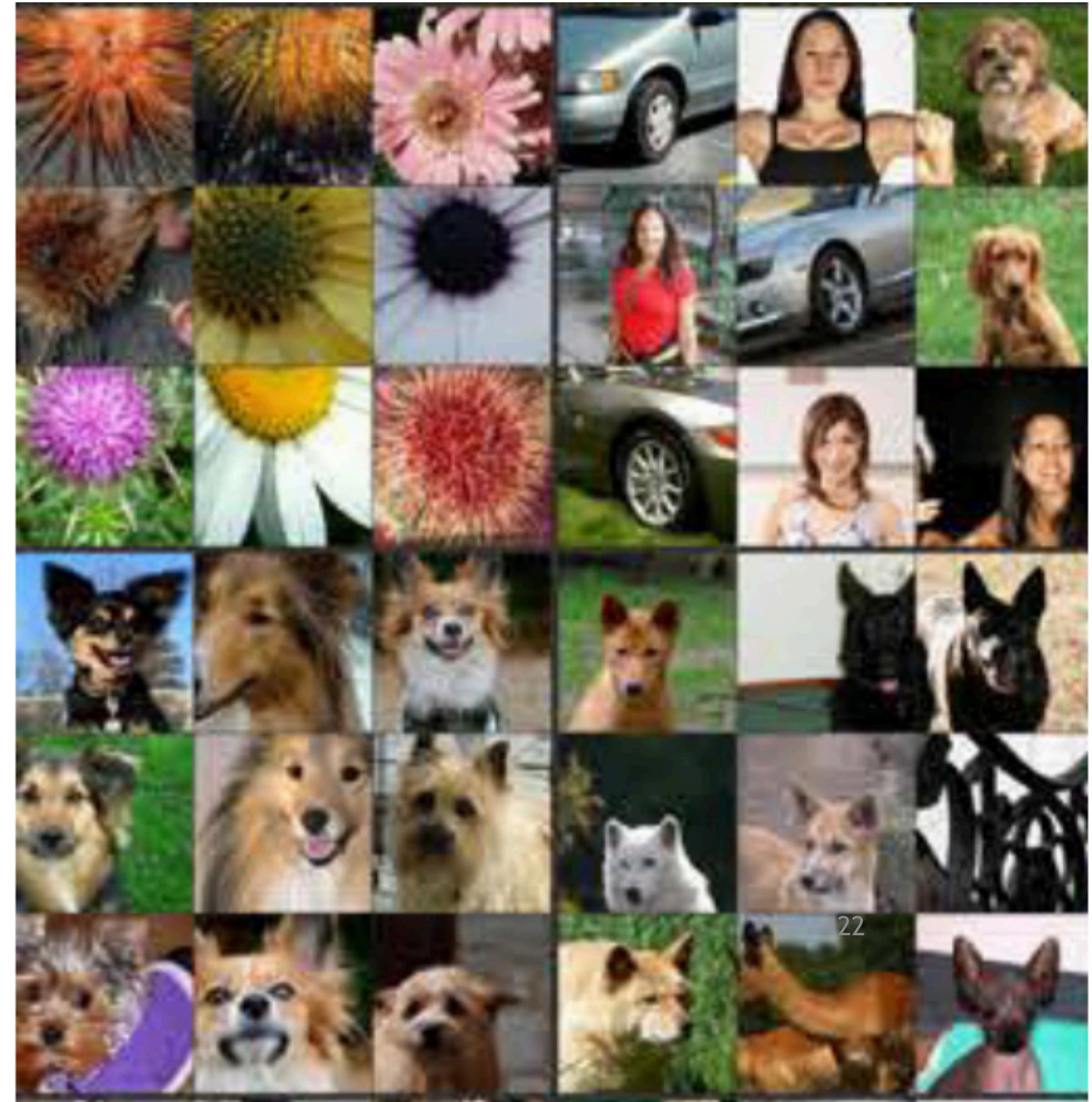
[Zeiler and Fergus, 2014]

Image patches that activate each of the **layer 4** neurons most strongly

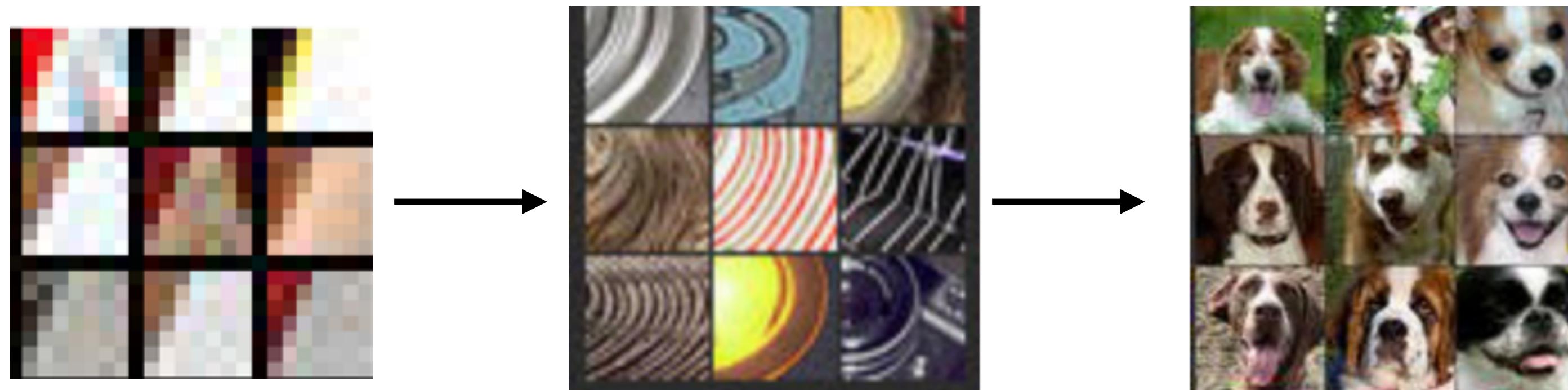
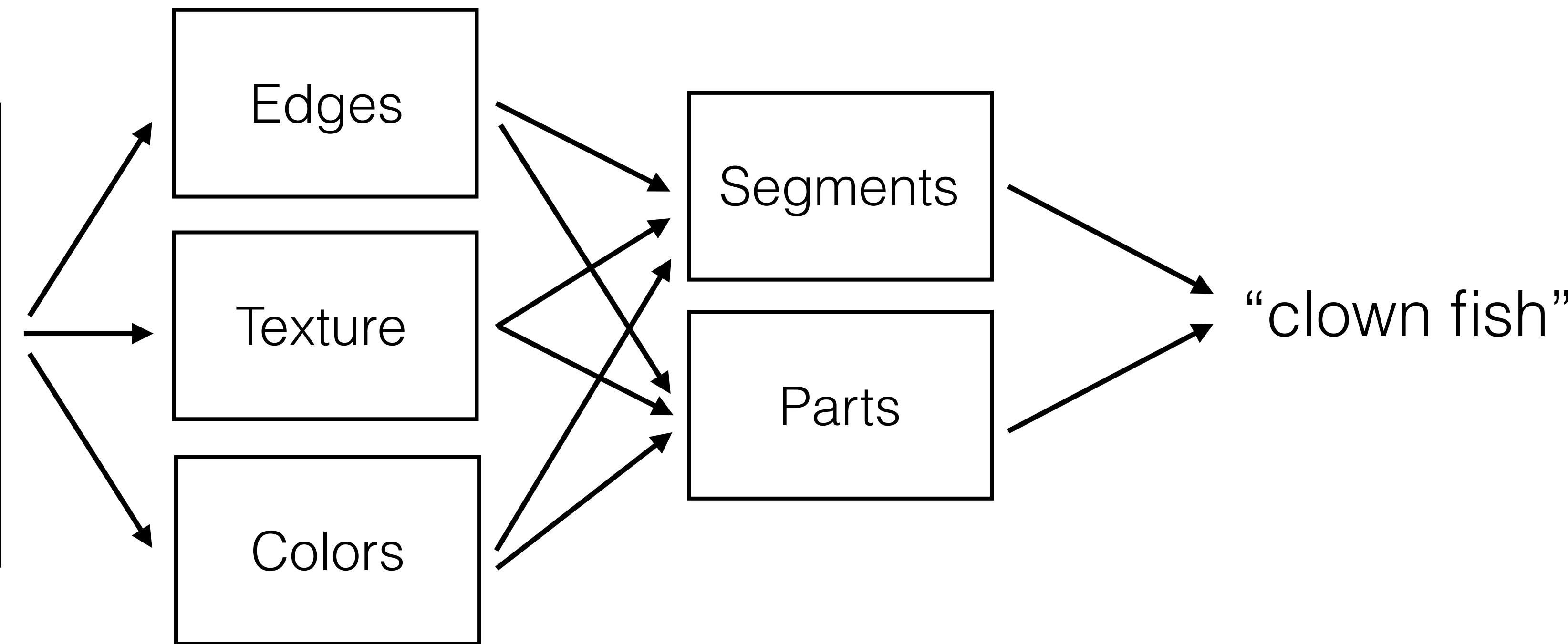


[Zeiler and Fergus, 2014]

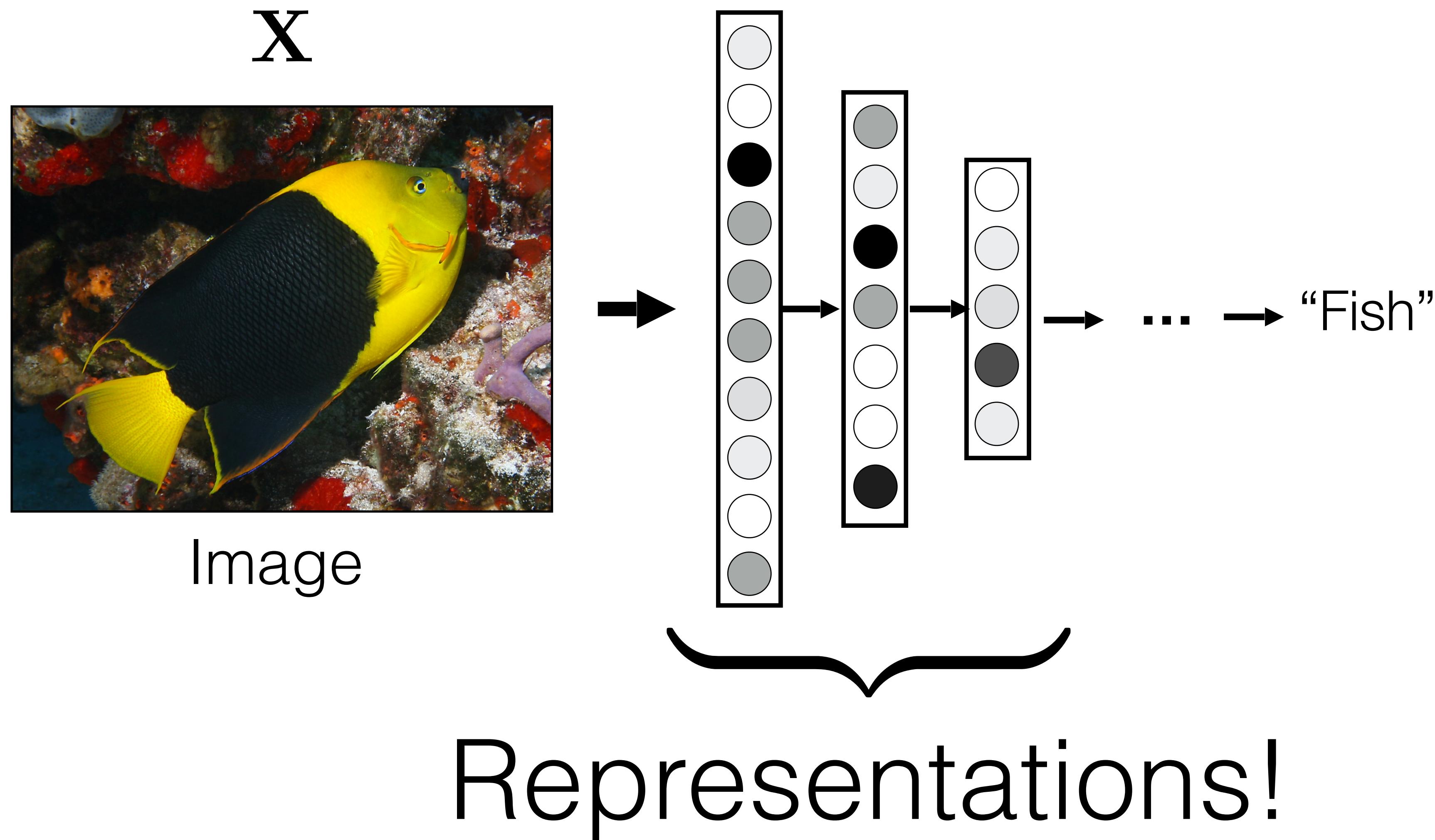
Image patches that activate each of the **layer 5** neurons most strongly

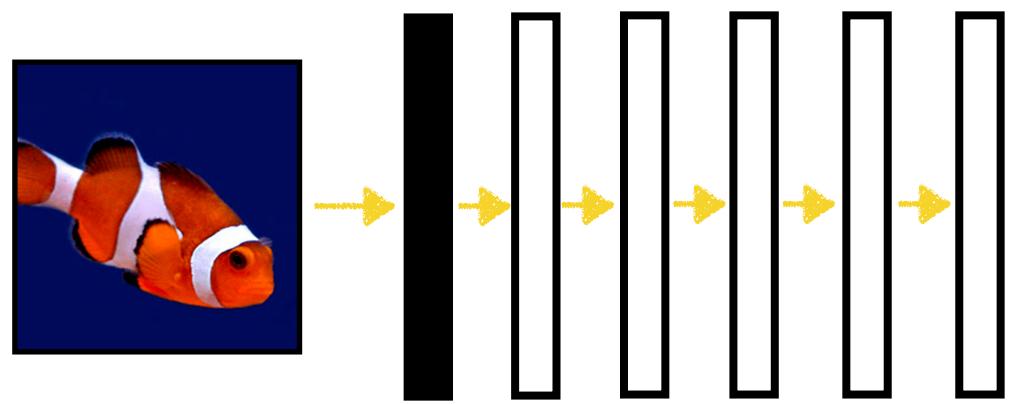


CNNs learned the classical visual recognition pipeline

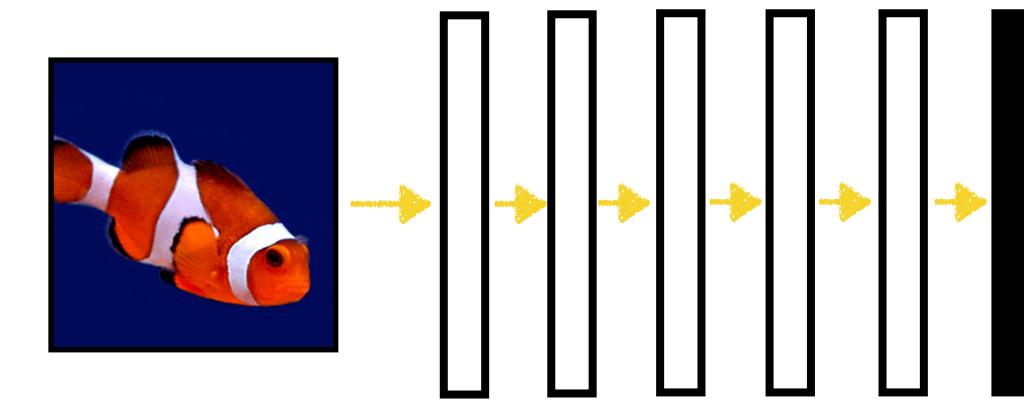


What do deep nets internally learn?

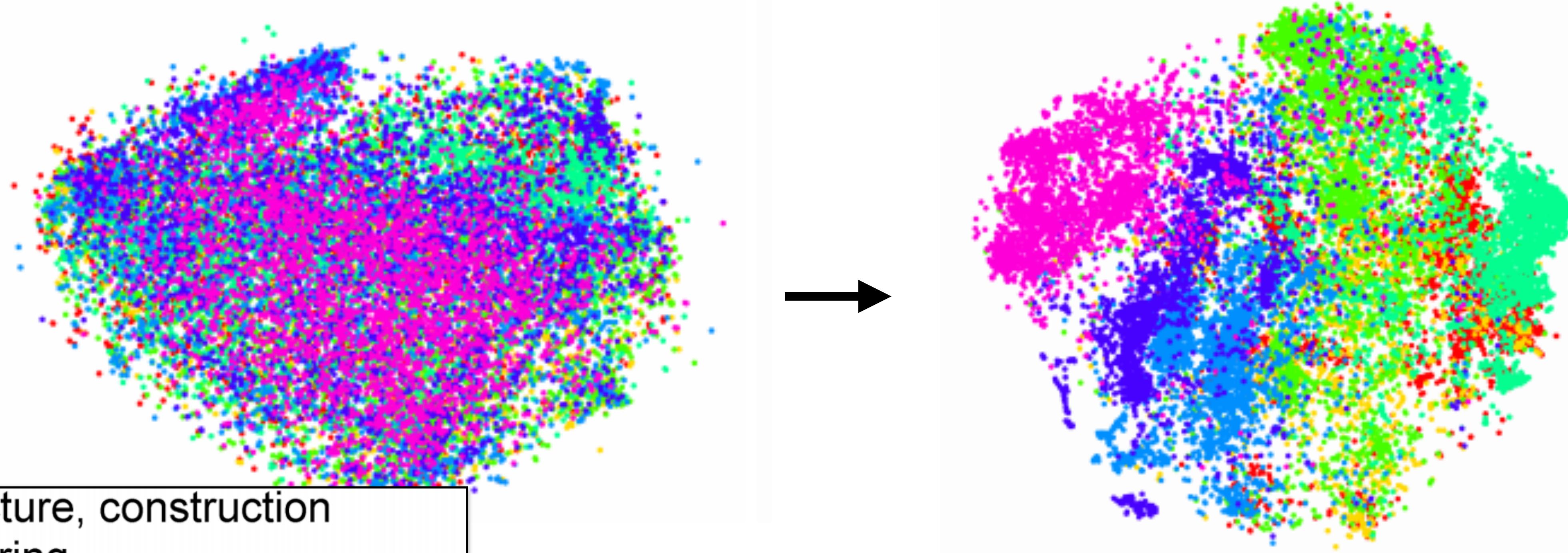




Layer 1 representation



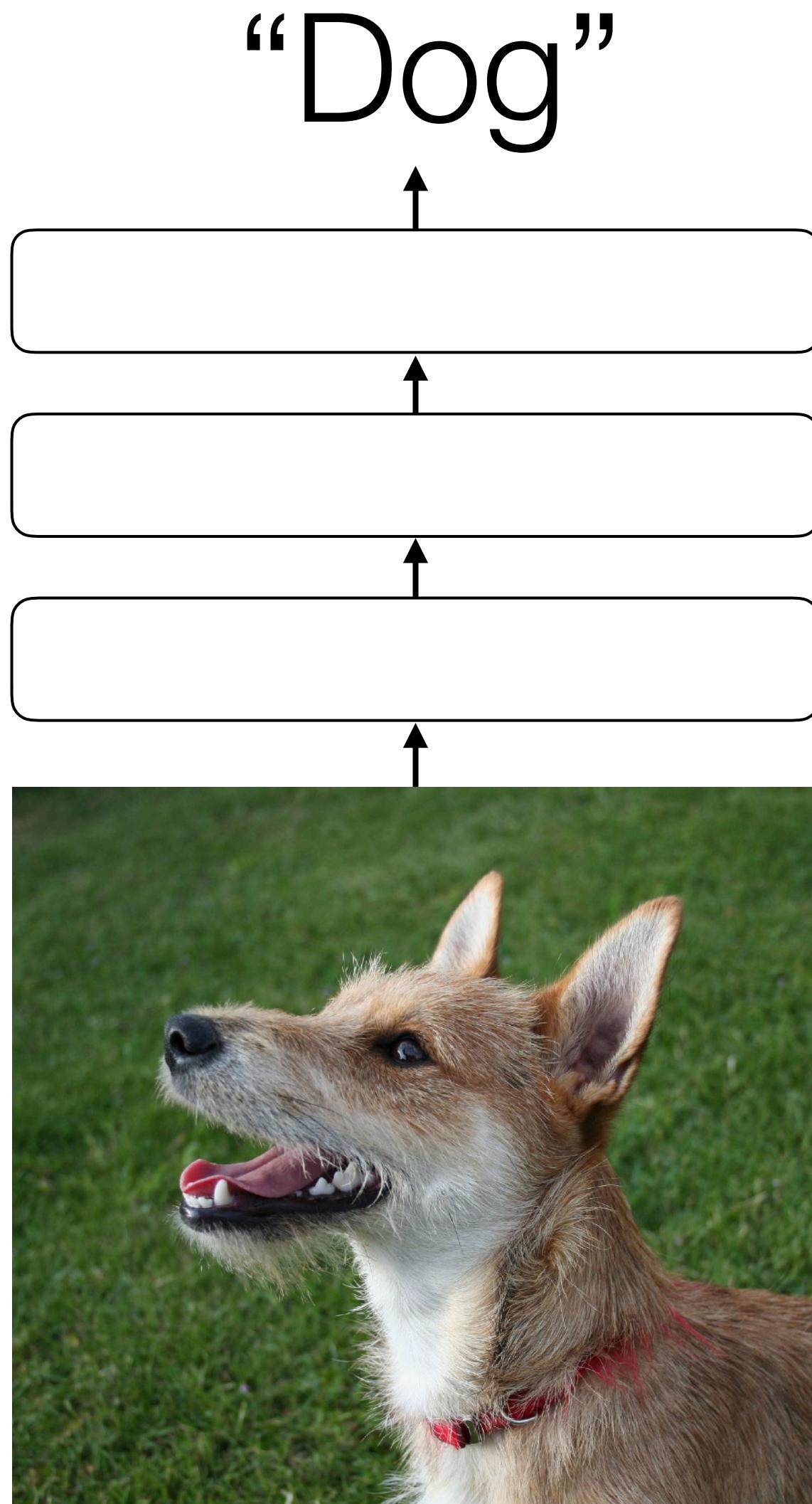
Layer 6 representation



- structure, construction
- covering
- commodity, trade good, good
- conveyance, transport
- invertebrate
- bird
- hunting dog

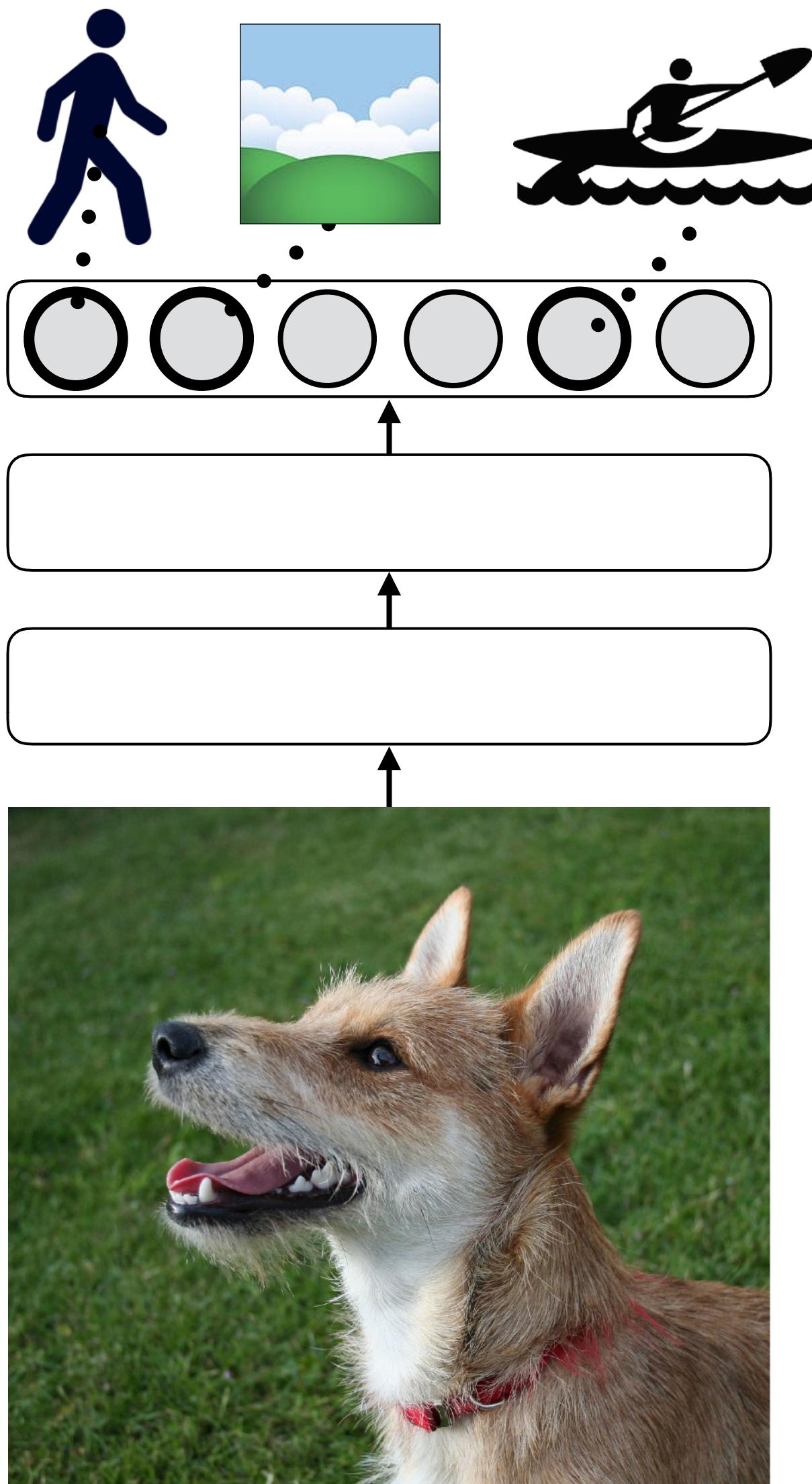
[Visualization technique : t-sne, van der Maaten & Hinton, 2008]
[DeCAF, Donahue, Jia, et al. 2013]

Transferring CNN features



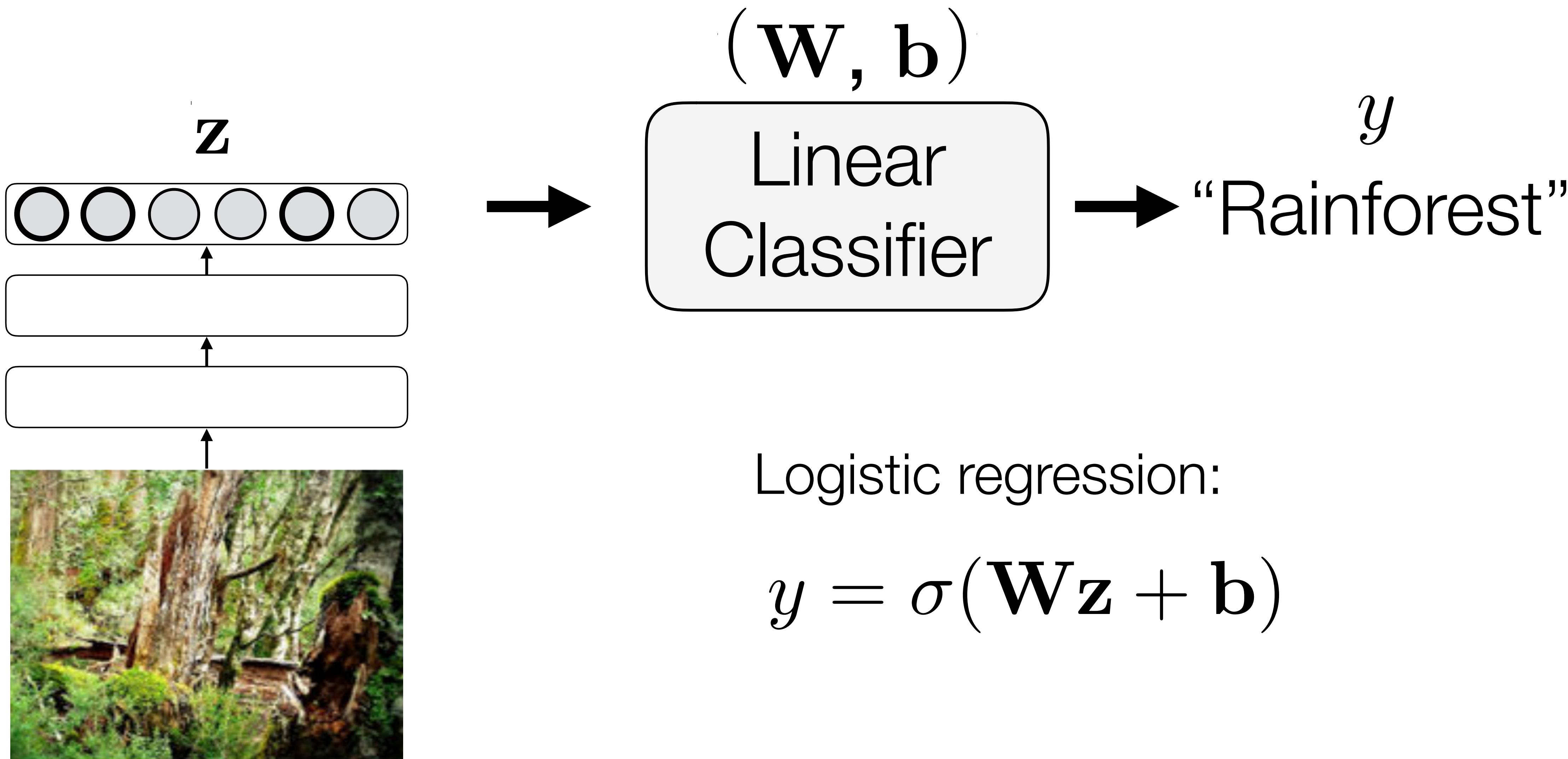
Object recognition net

Transferring CNN features



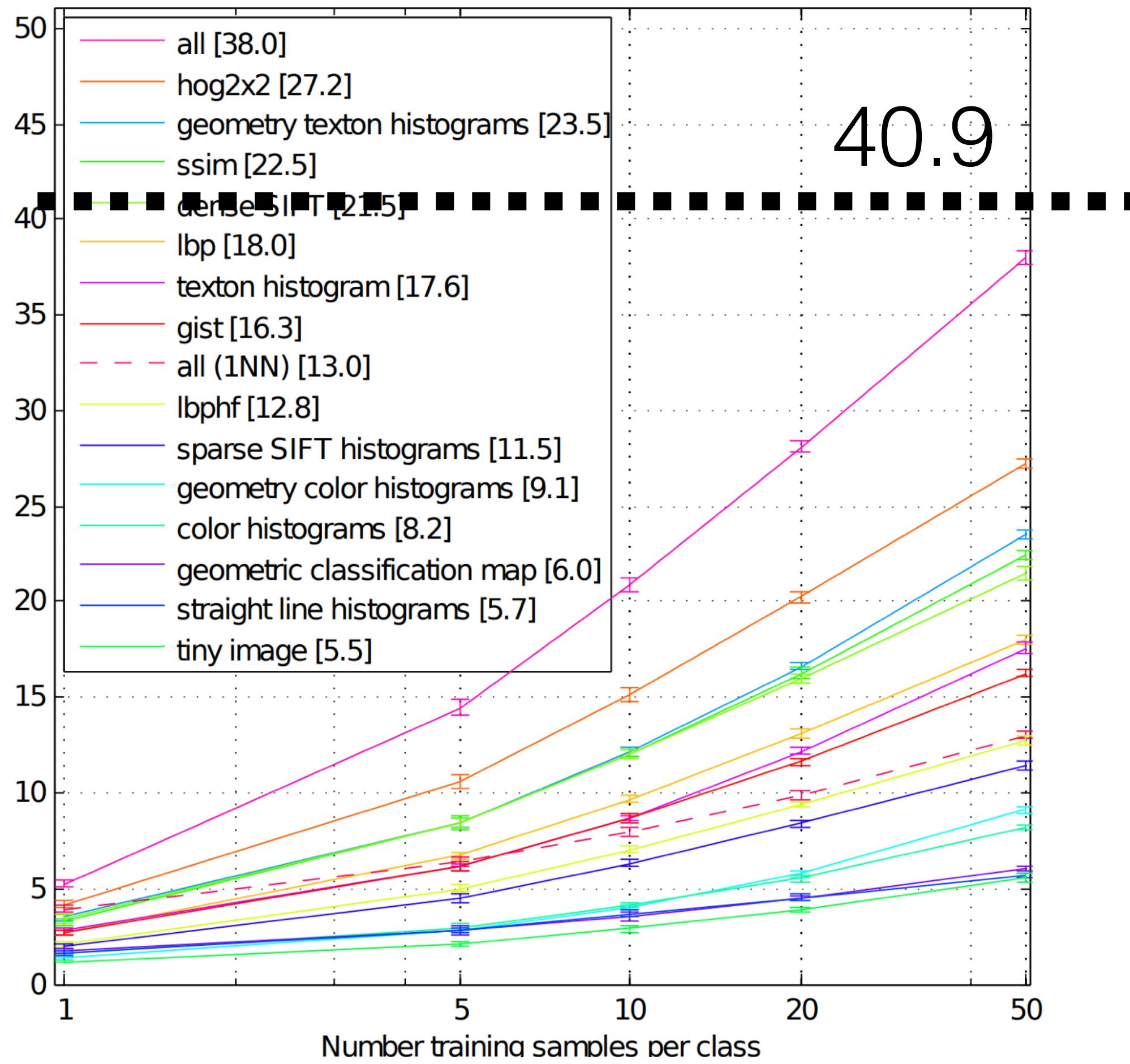
Object recognition net

Simple feature transfer



Transferring CNN features

Hand-crafted features



CNN features pretrained on ImageNet
+ linear classifier [Donahue et al. 2013]

How do we learn good representations?

Supervised object recognition

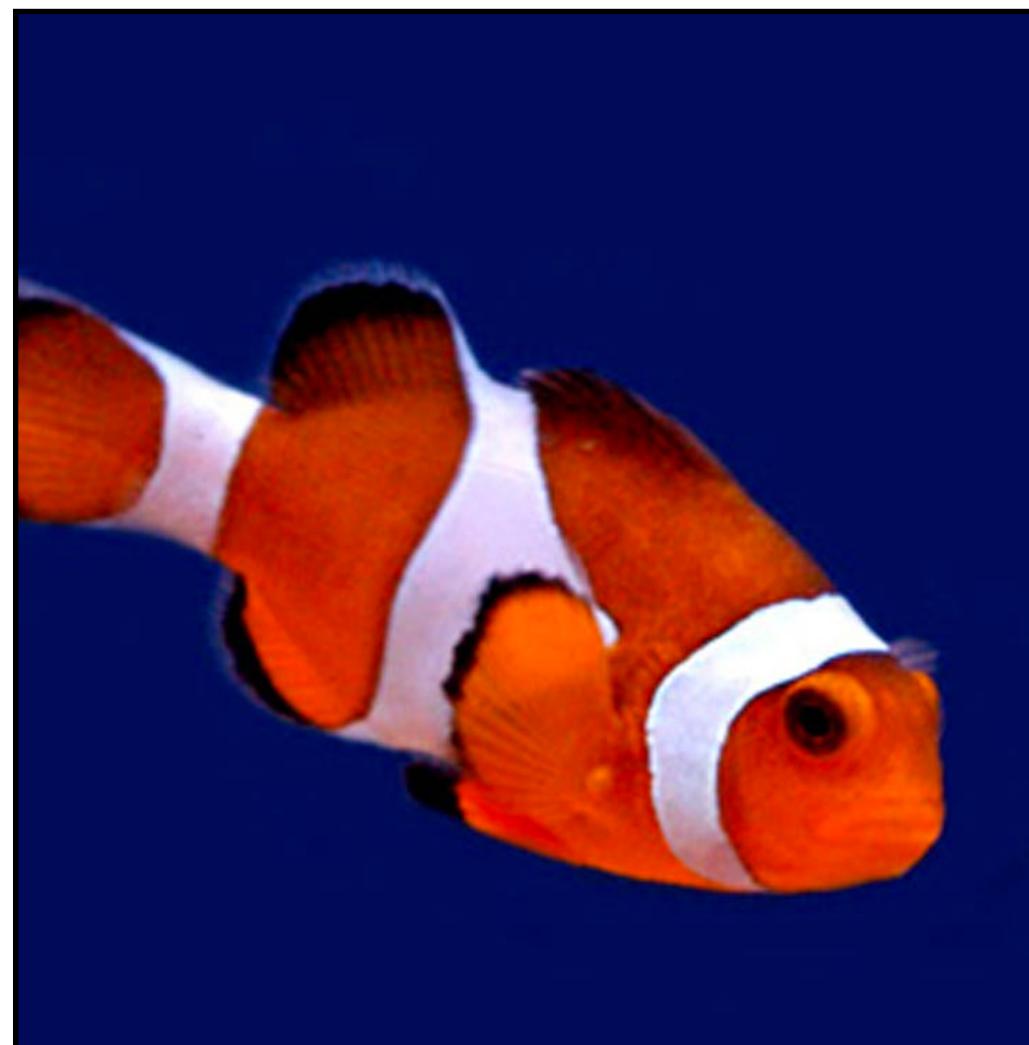


image X



→ “Fish”

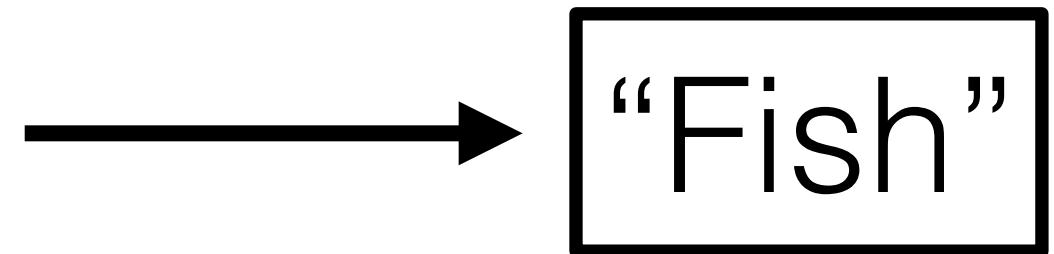
A horizontal black arrow pointing from the learner to the output label.

label Y

Supervised object recognition



image X

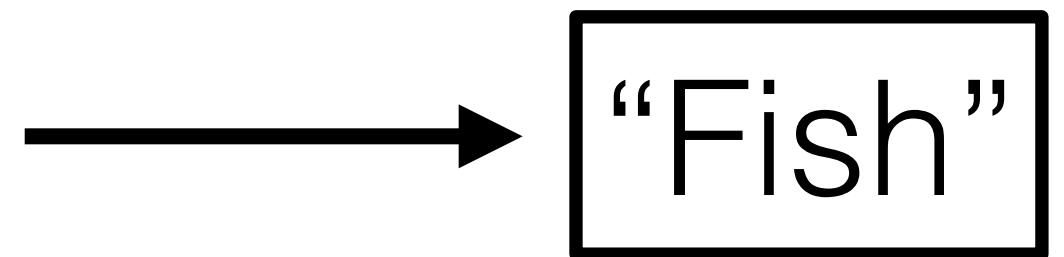


label Y

Supervised object recognition



image X



label Y

Supervised object recognition

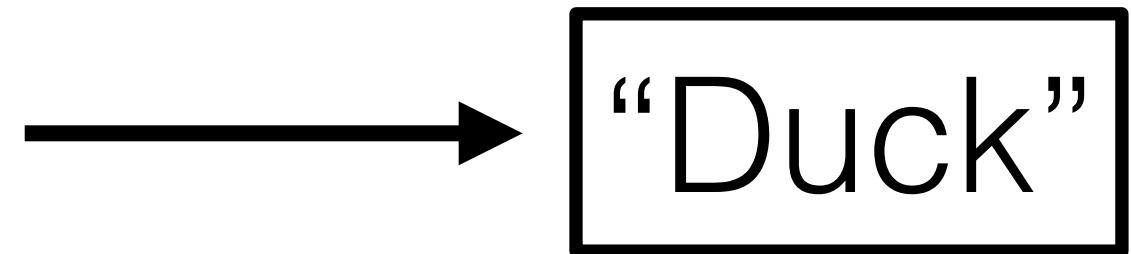
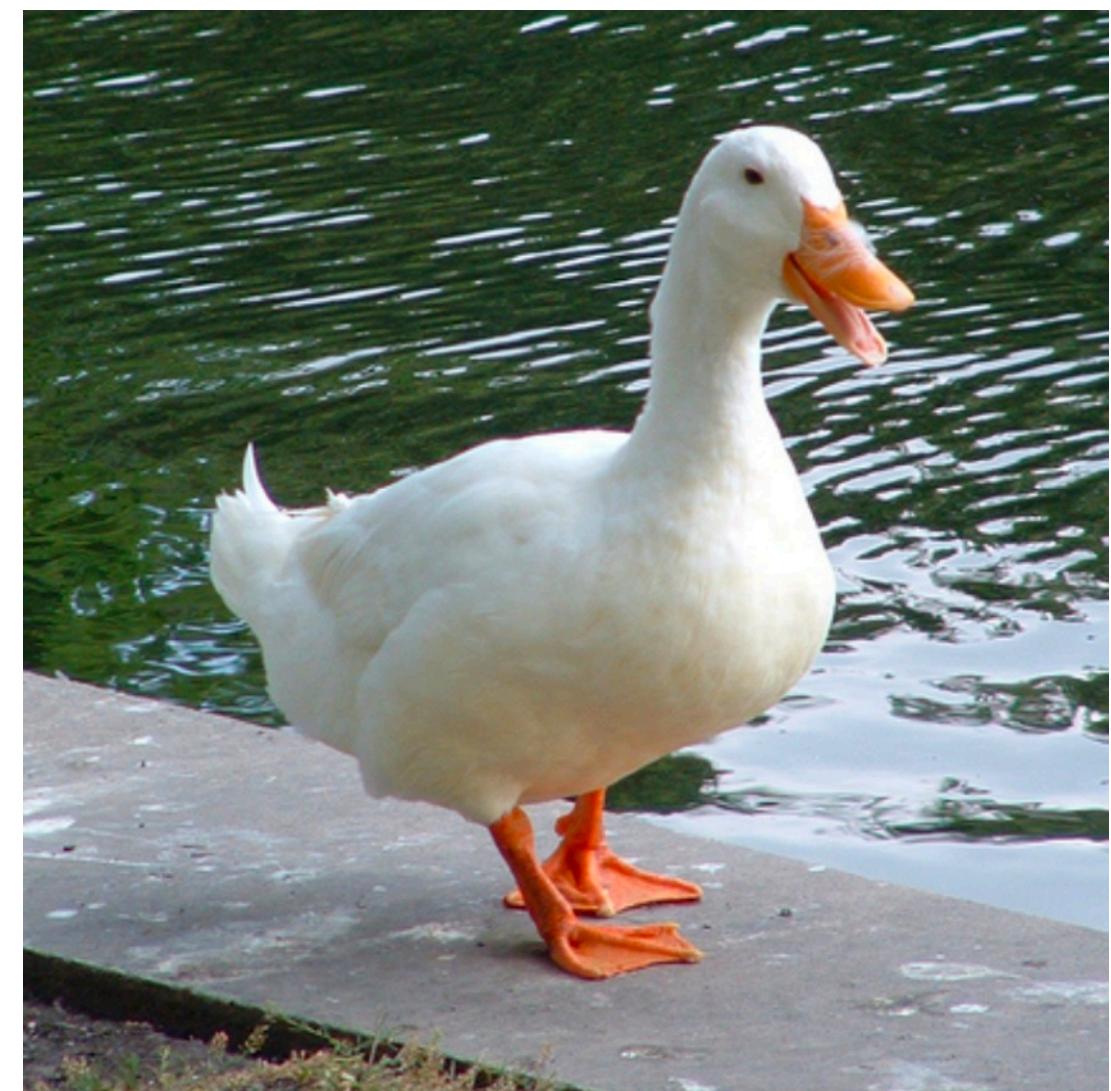


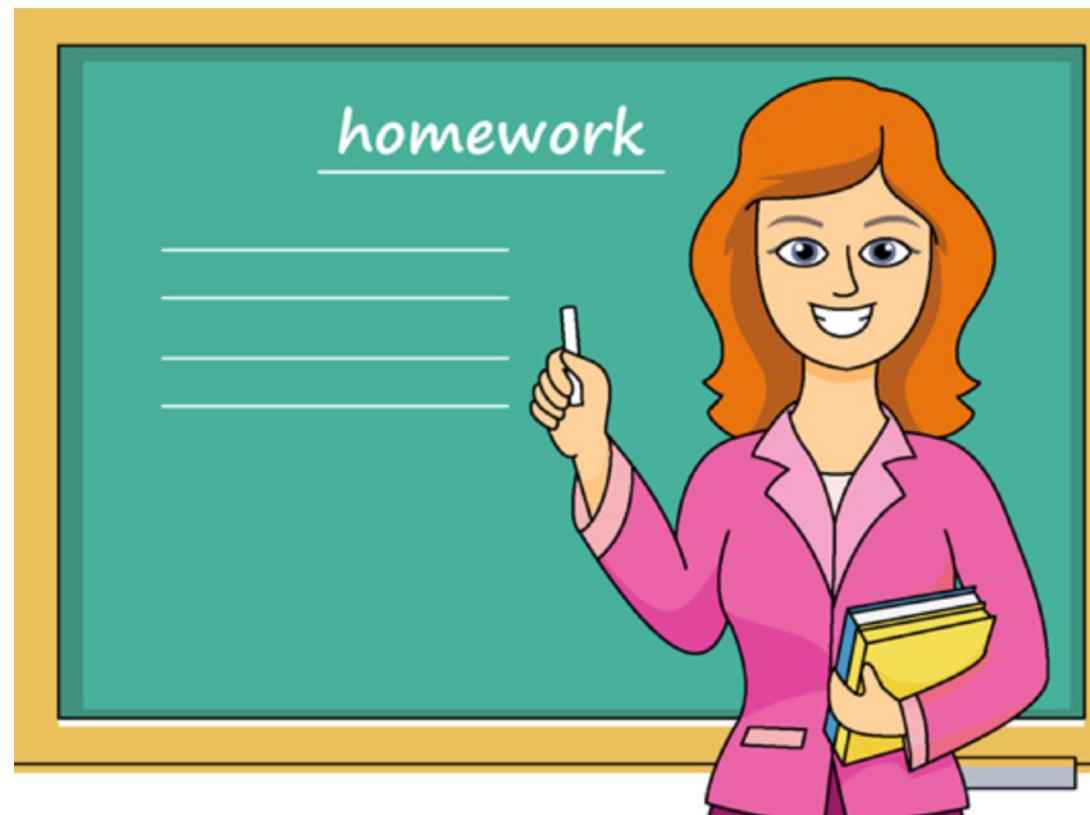
image X

label Y

Supervised computer vision

Hand-curated training data

- + Informative
- Expensive
- Limited to teacher's knowledge



Vision in nature

Raw unlabeled training data

- + Cheap
- Noisy
- Harder to interpret



Learning from examples

(aka **supervised learning**)

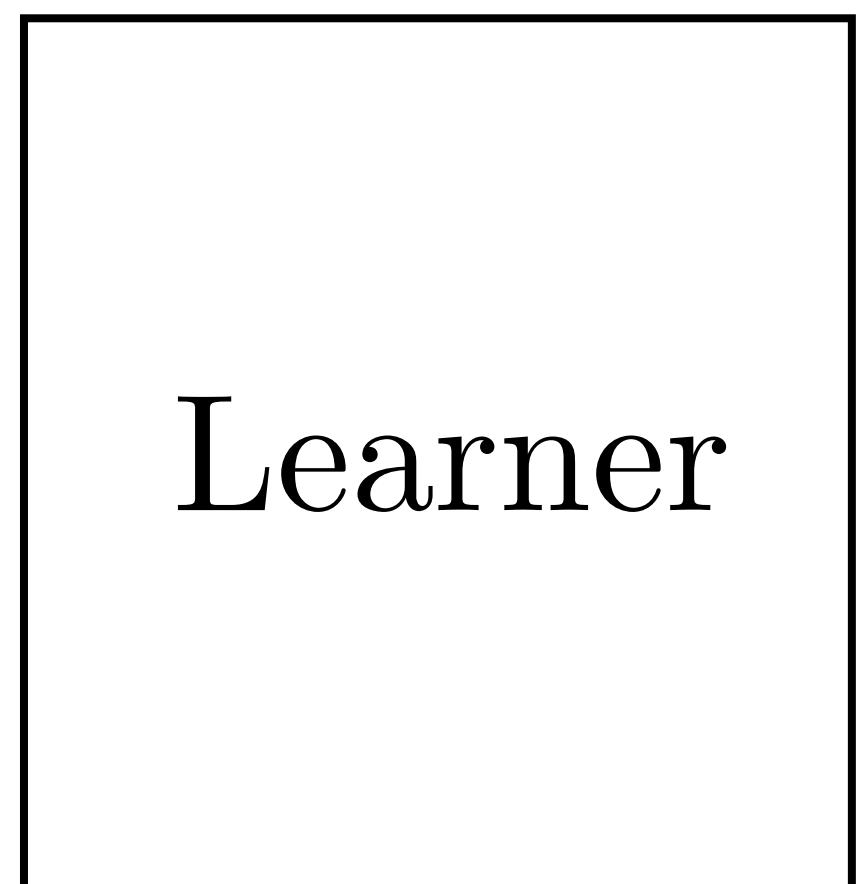
Training data

$$\{x_1, y_1\}$$

$$\{x_2, y_2\} \rightarrow$$

$$\{x_3, y_3\}$$

...



$$\rightarrow f : X \rightarrow Y$$

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(x_i), y_i)$$

Representation Learning

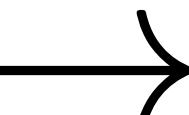
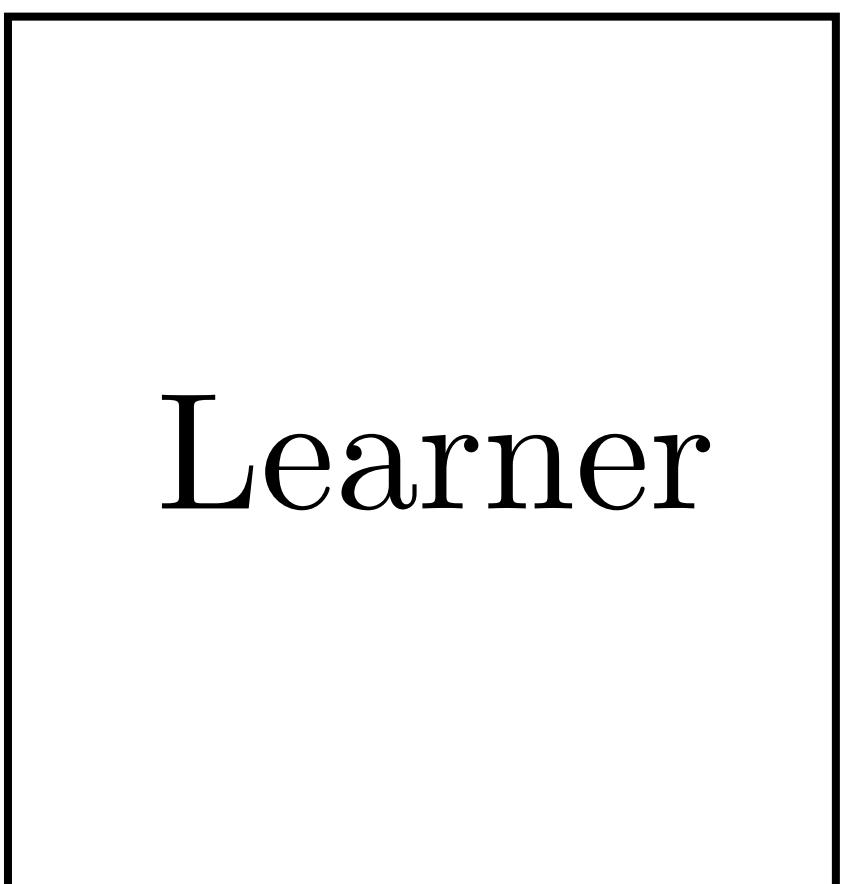
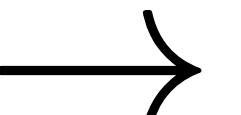
Data

$\{x_1\}$

$\{x_2\}$

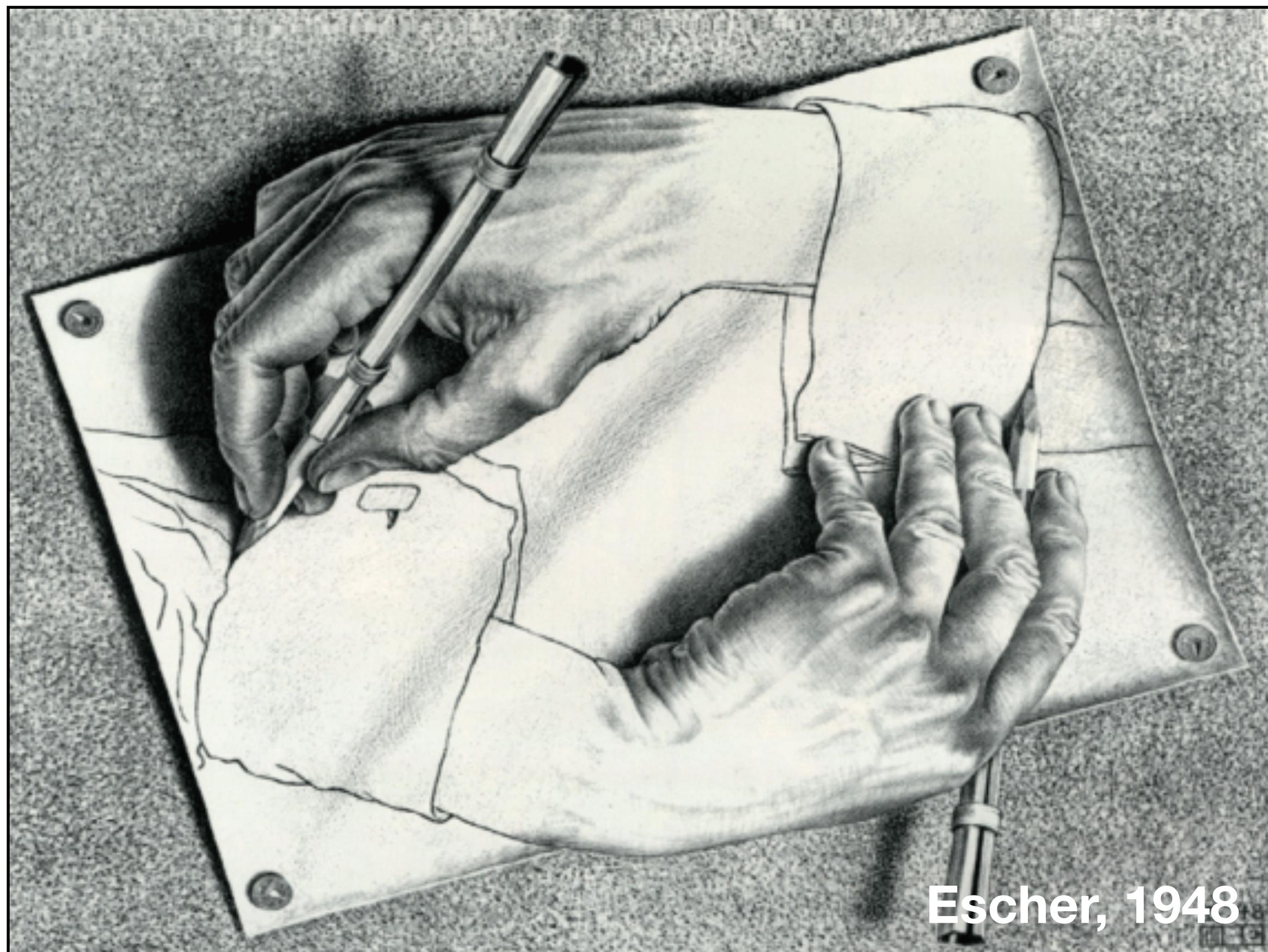
$\{x_3\}$

...



Representations

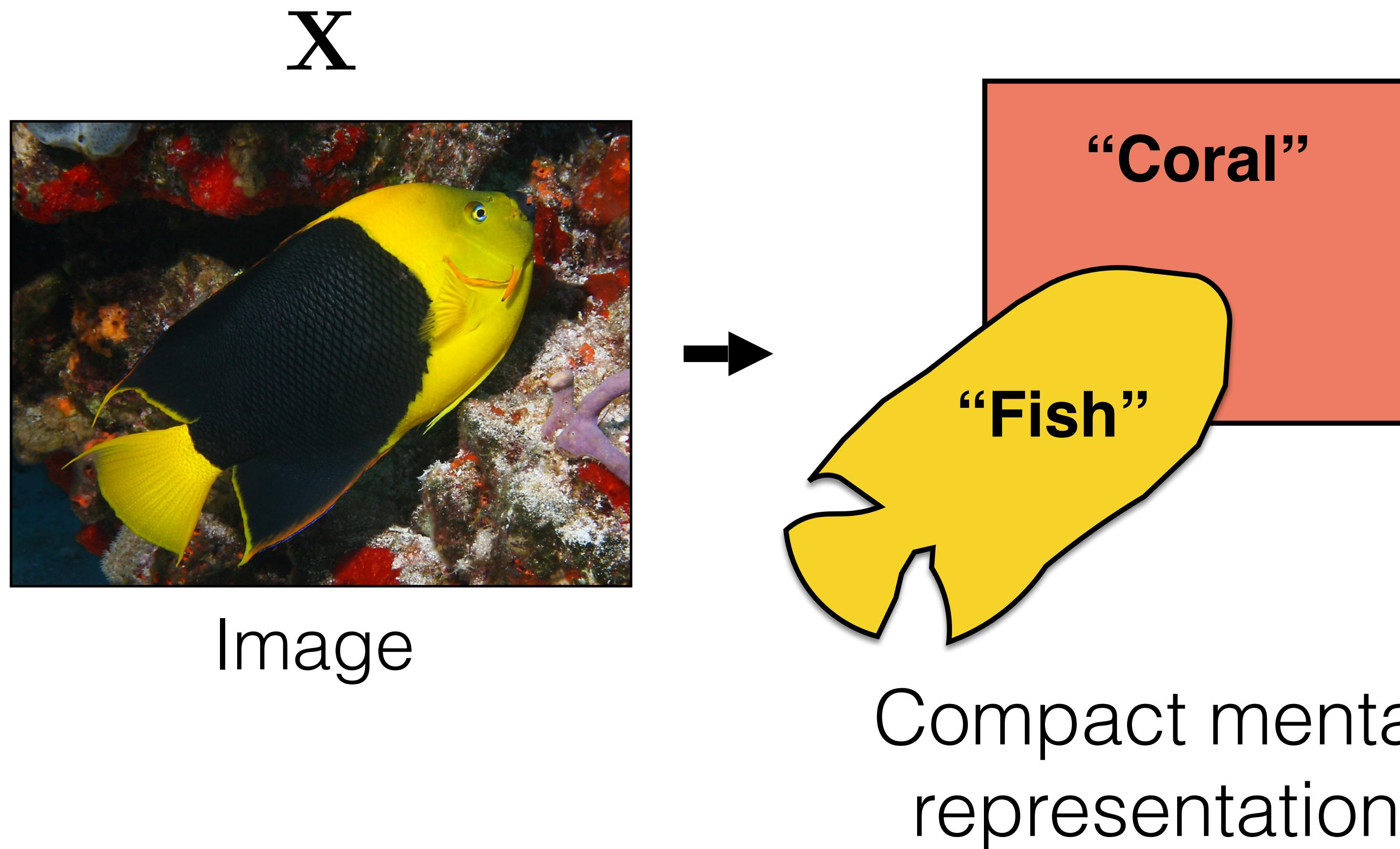
Self-supervised learning



Common trick:

- Convert “unsupervised” problem into “supervised” empirical risk minimization
- Do so by cooking up “labels” (prediction targets) from the raw data itself
- Designing new algorithms still takes a lot of trial and error.

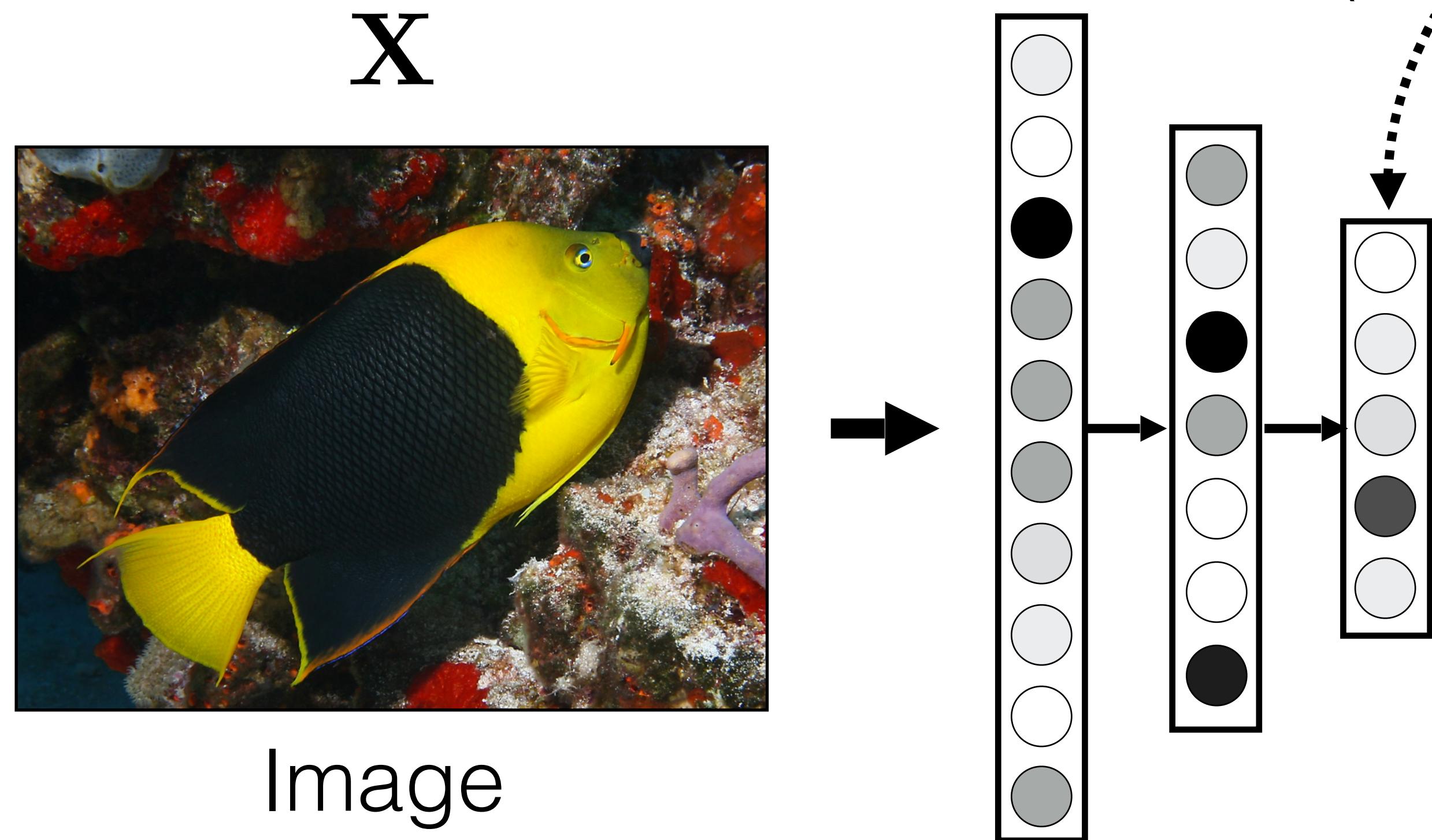
Unsupervised Representation Learning



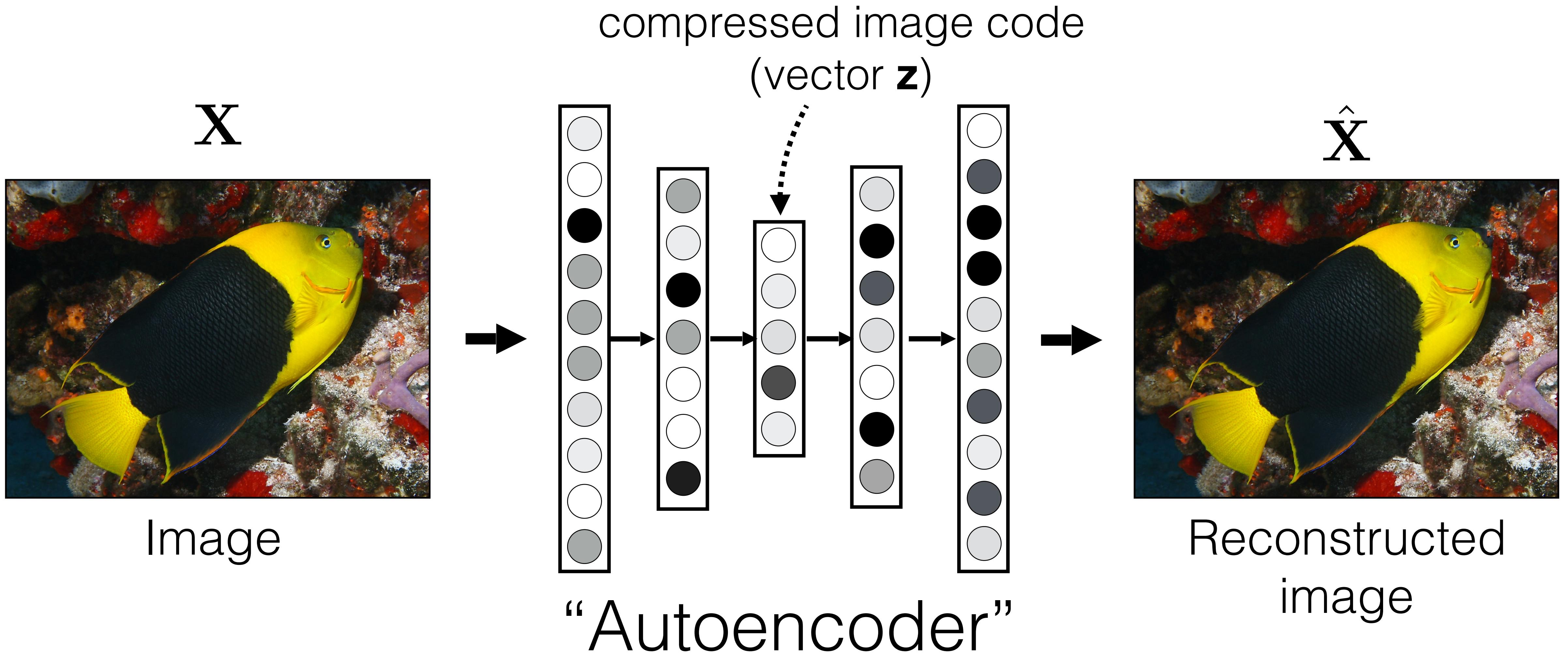
47

Unsupervised Representation Learning

compressed image code
(vector \mathbf{z})



Unsupervised Representation Learning

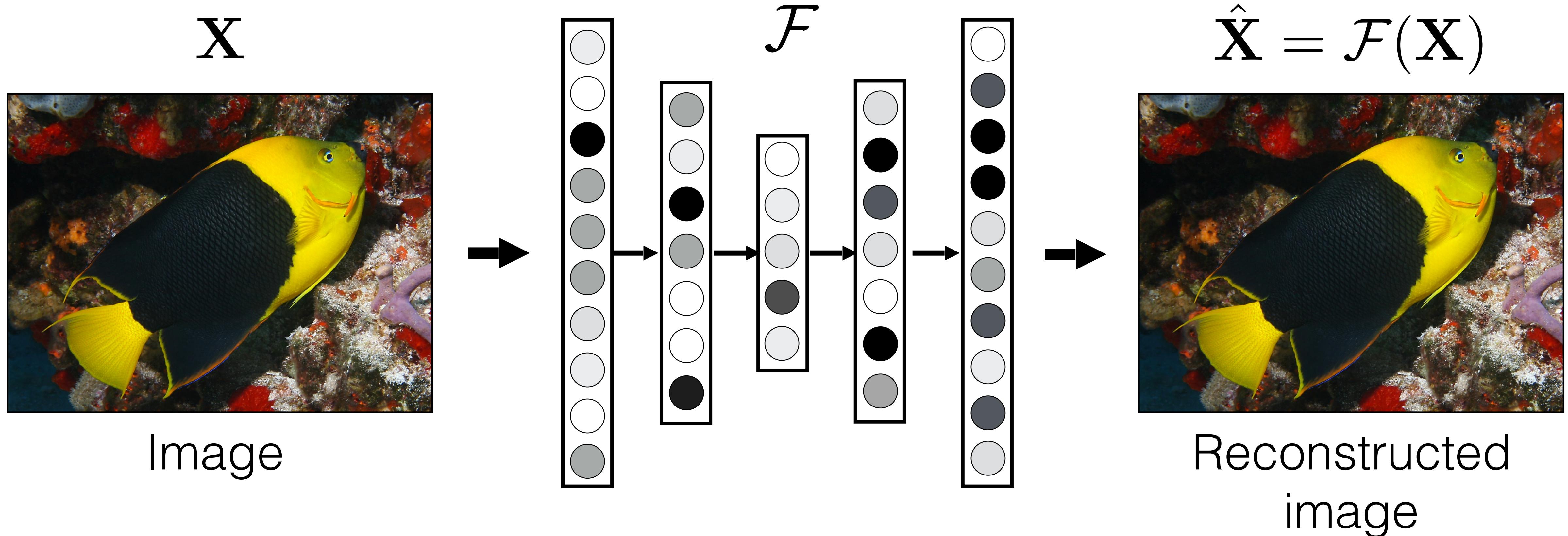


49

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

Source: Isola, Freeman, Torralba

Autoencoder



$$\arg \min_{\mathcal{F}} \mathbb{E}_{\mathbf{X}} [||\mathcal{F}(\mathbf{X}) - \mathbf{X}||]$$

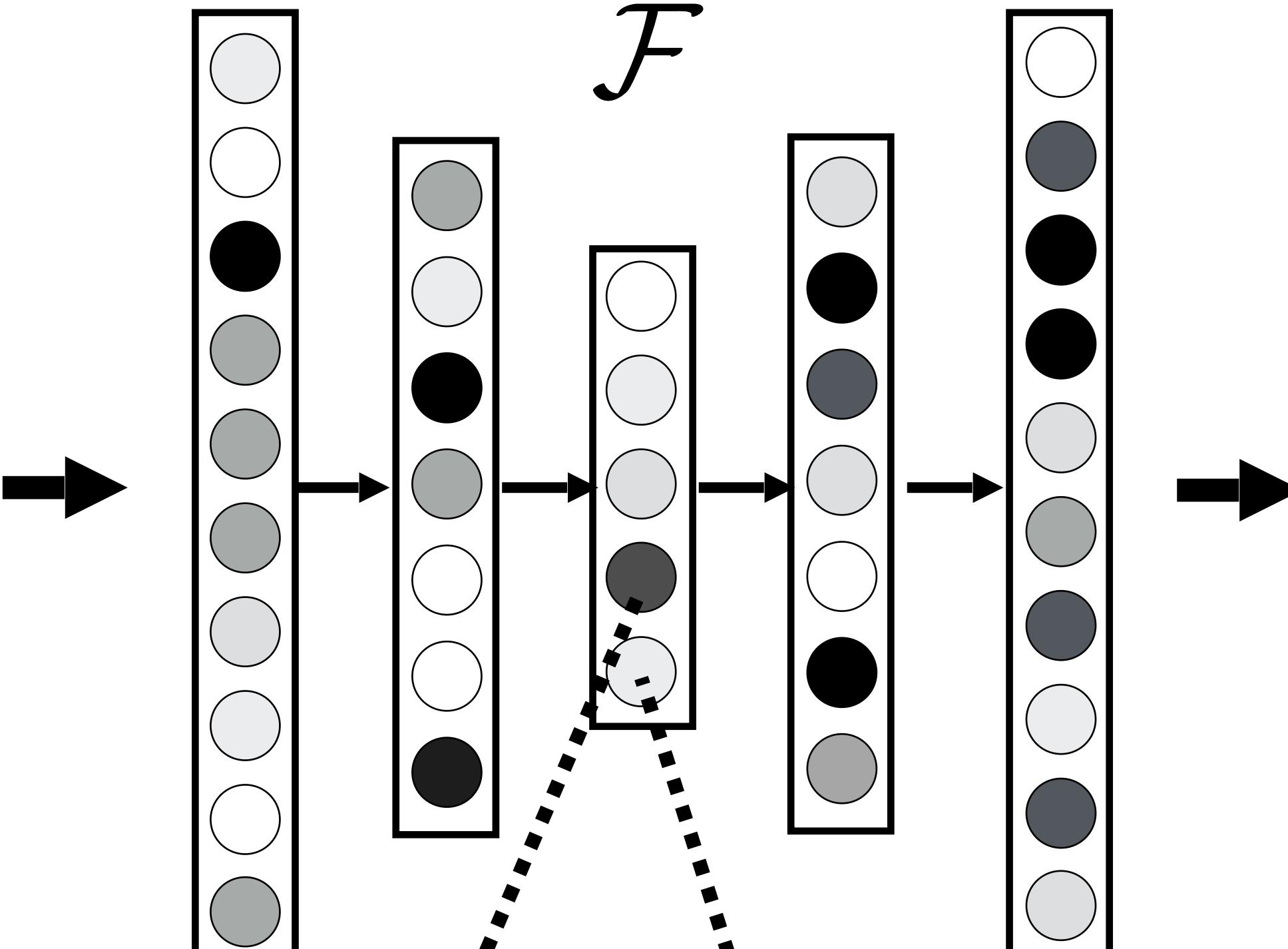
50

\mathbf{X}



Image

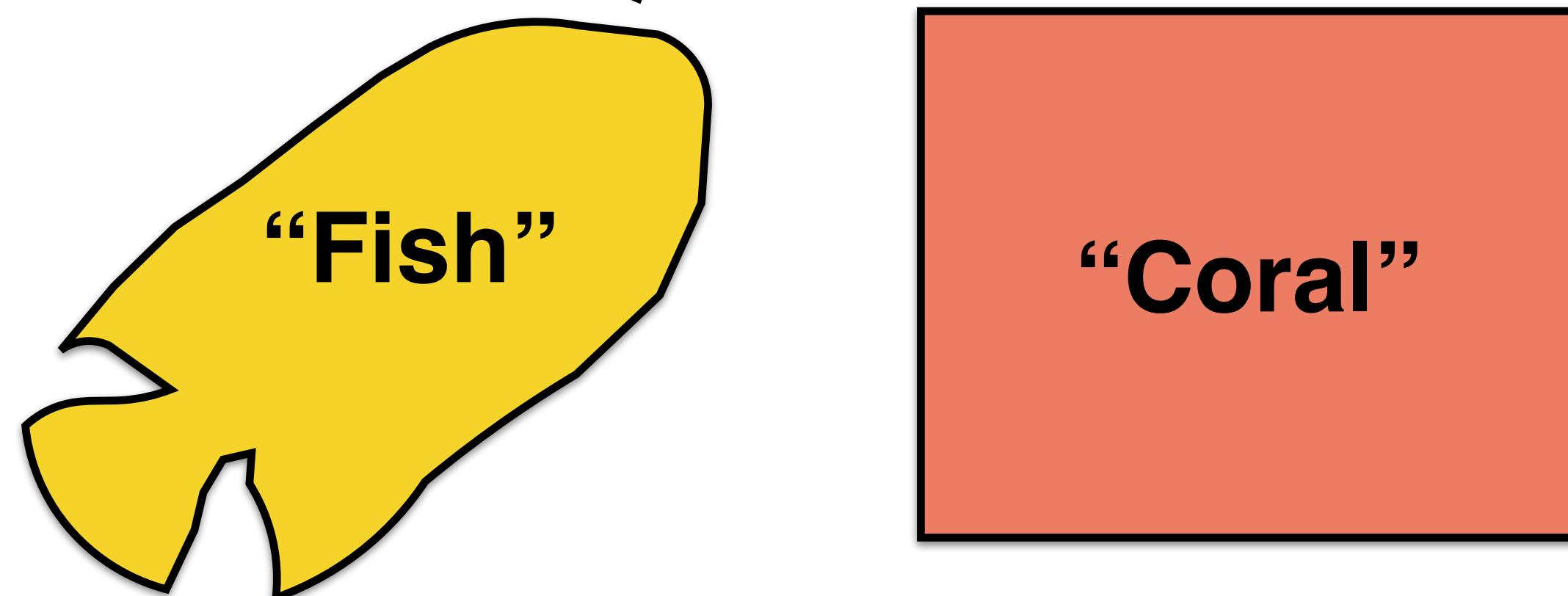
\mathcal{F}



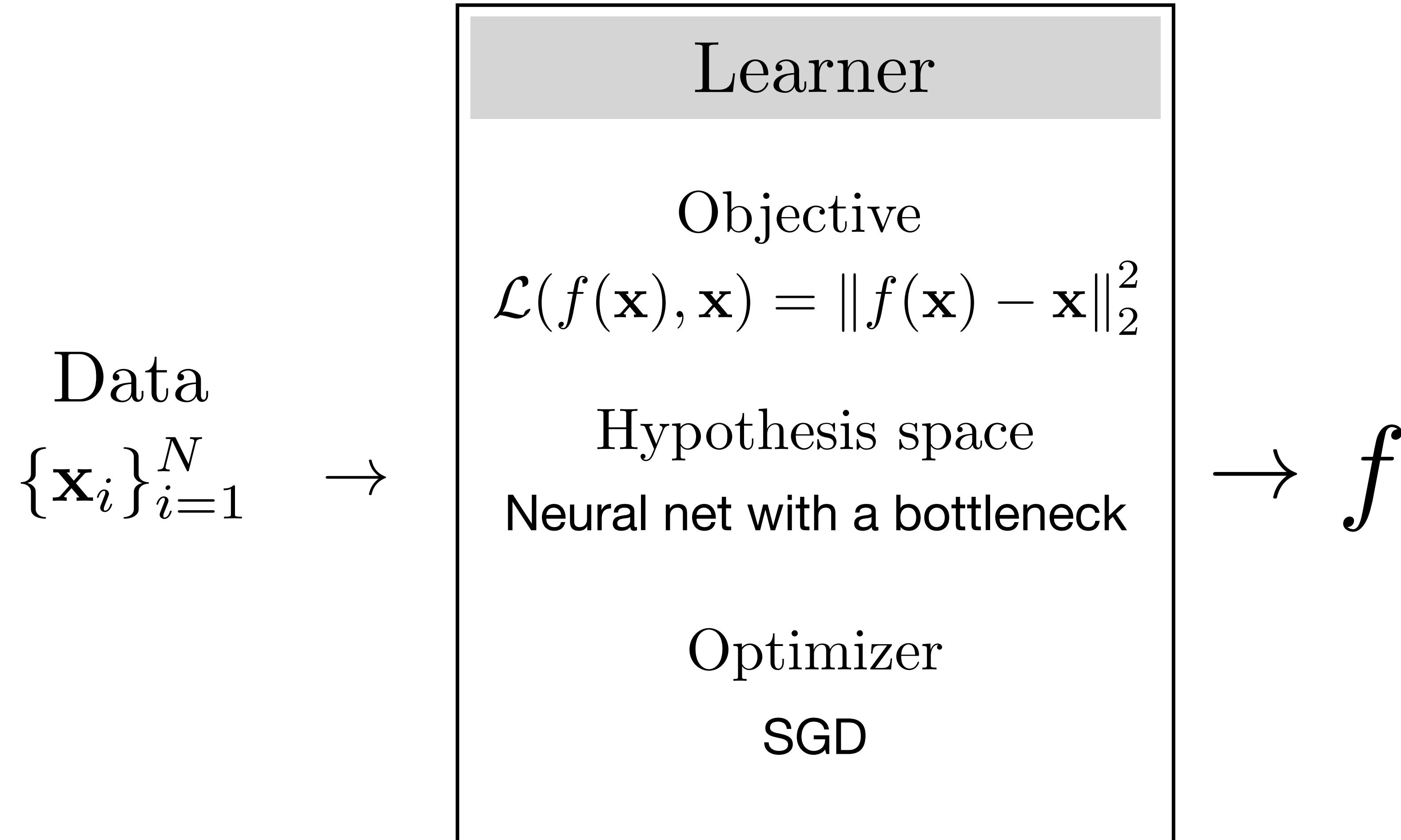
$\hat{\mathbf{X}} = \mathcal{F}(\mathbf{X})$



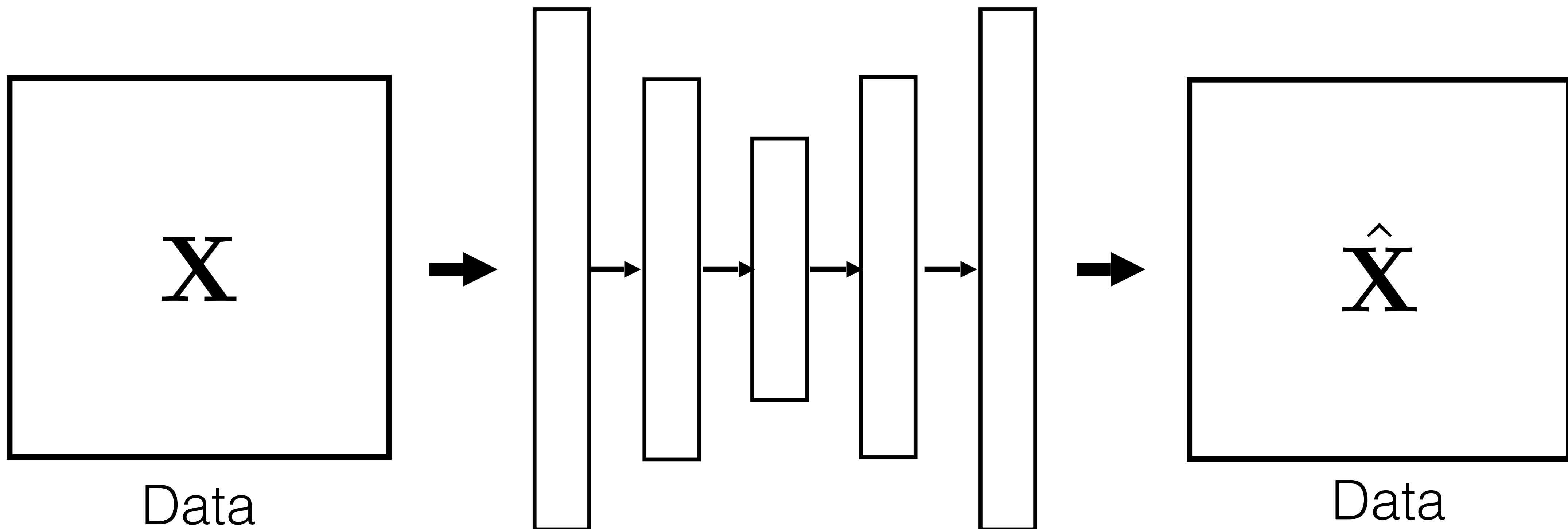
Reconstructed
image



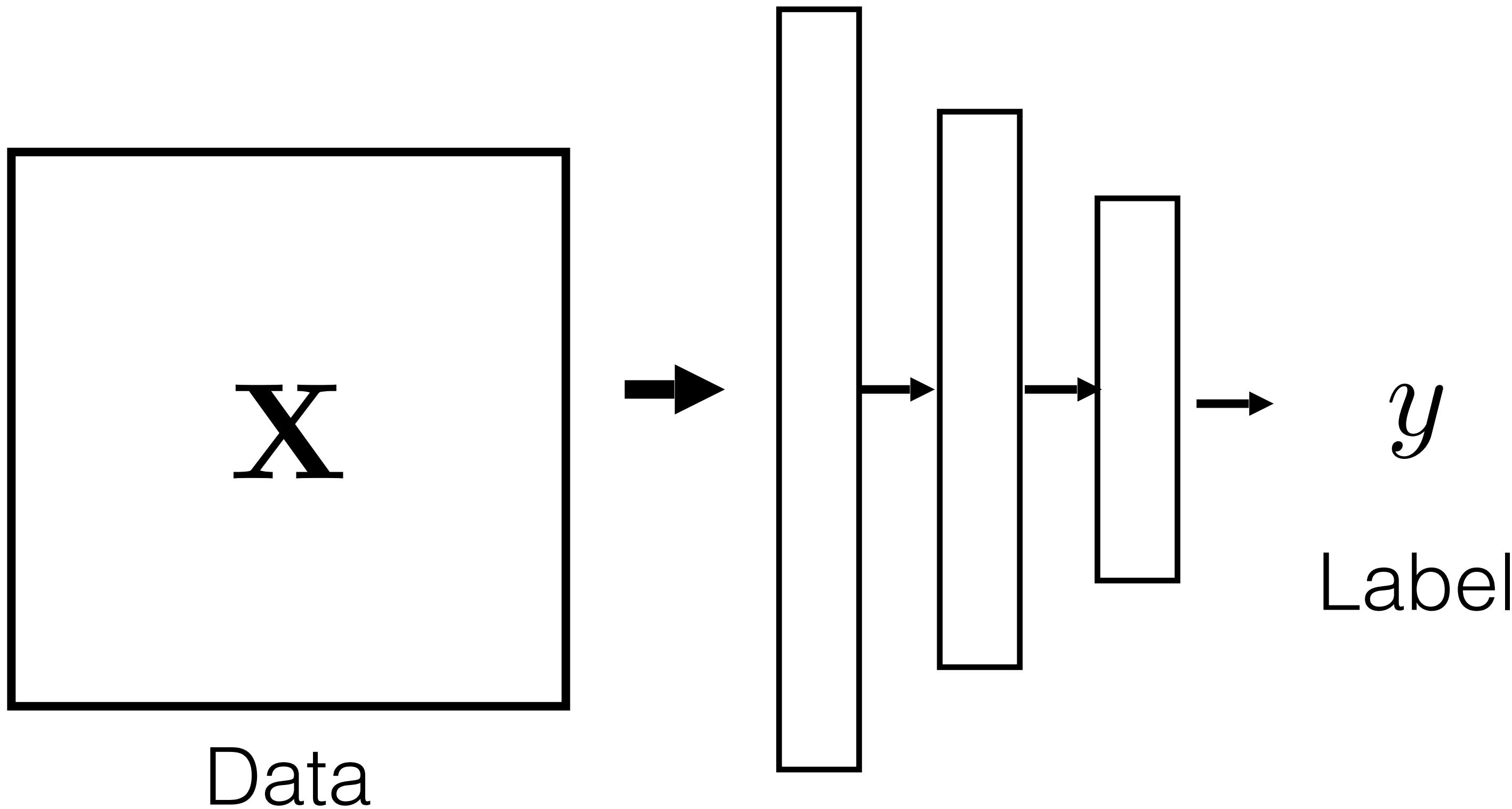
Autoencoder



Data compression



Label prediction

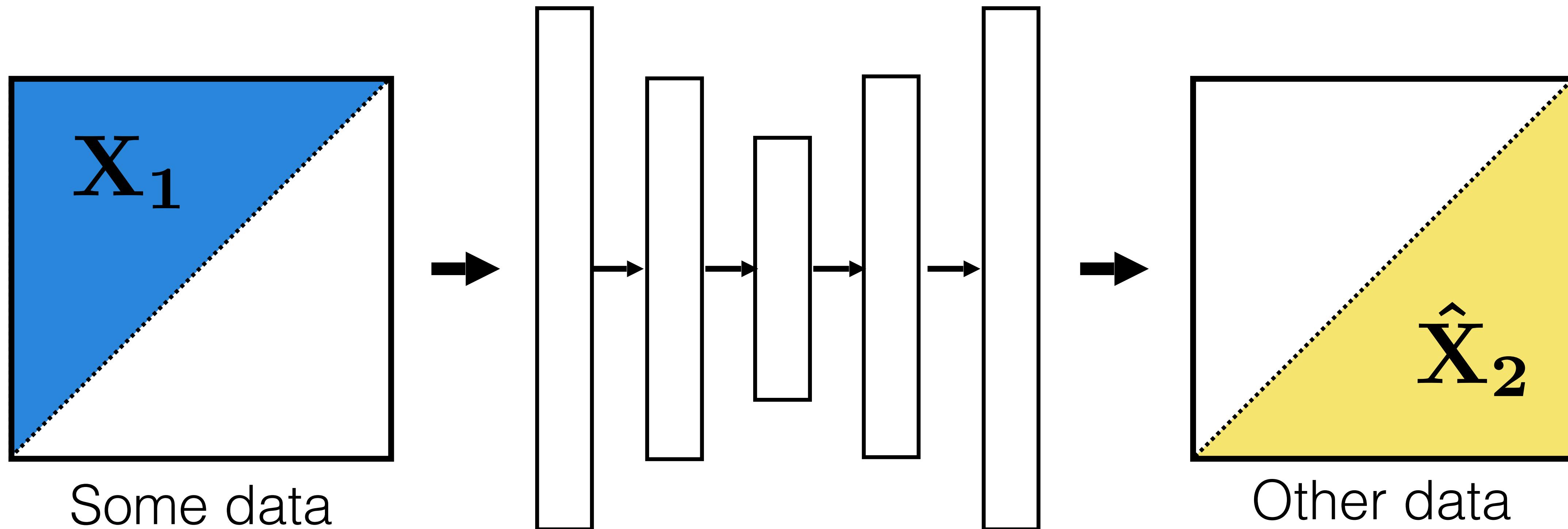


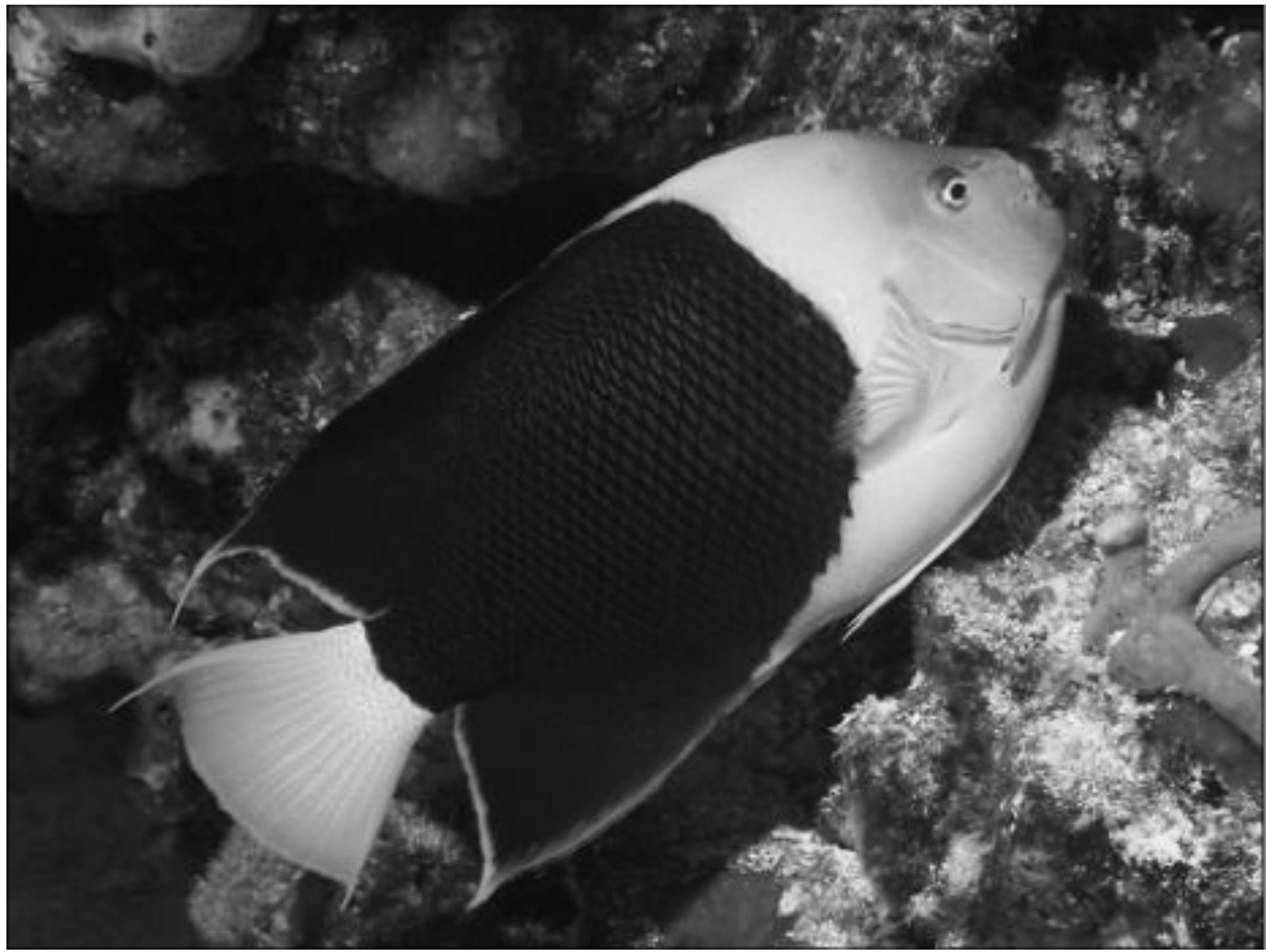
e.g., image classification

54

Data prediction

aka “self-supervised learning”





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

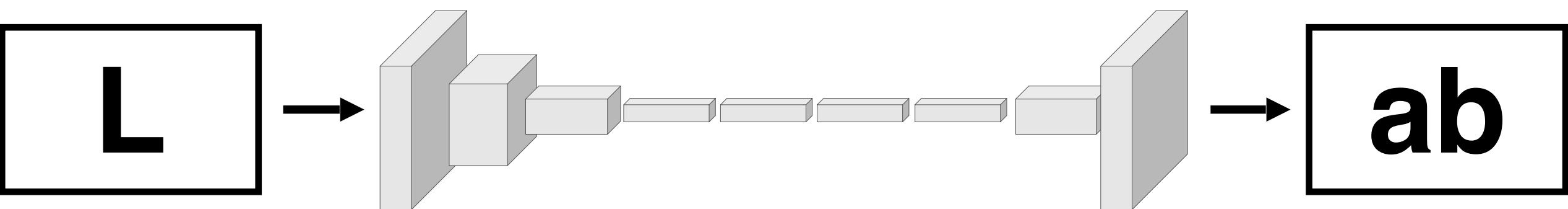


Grayscale image: L channel

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Color information: ab channels

$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$

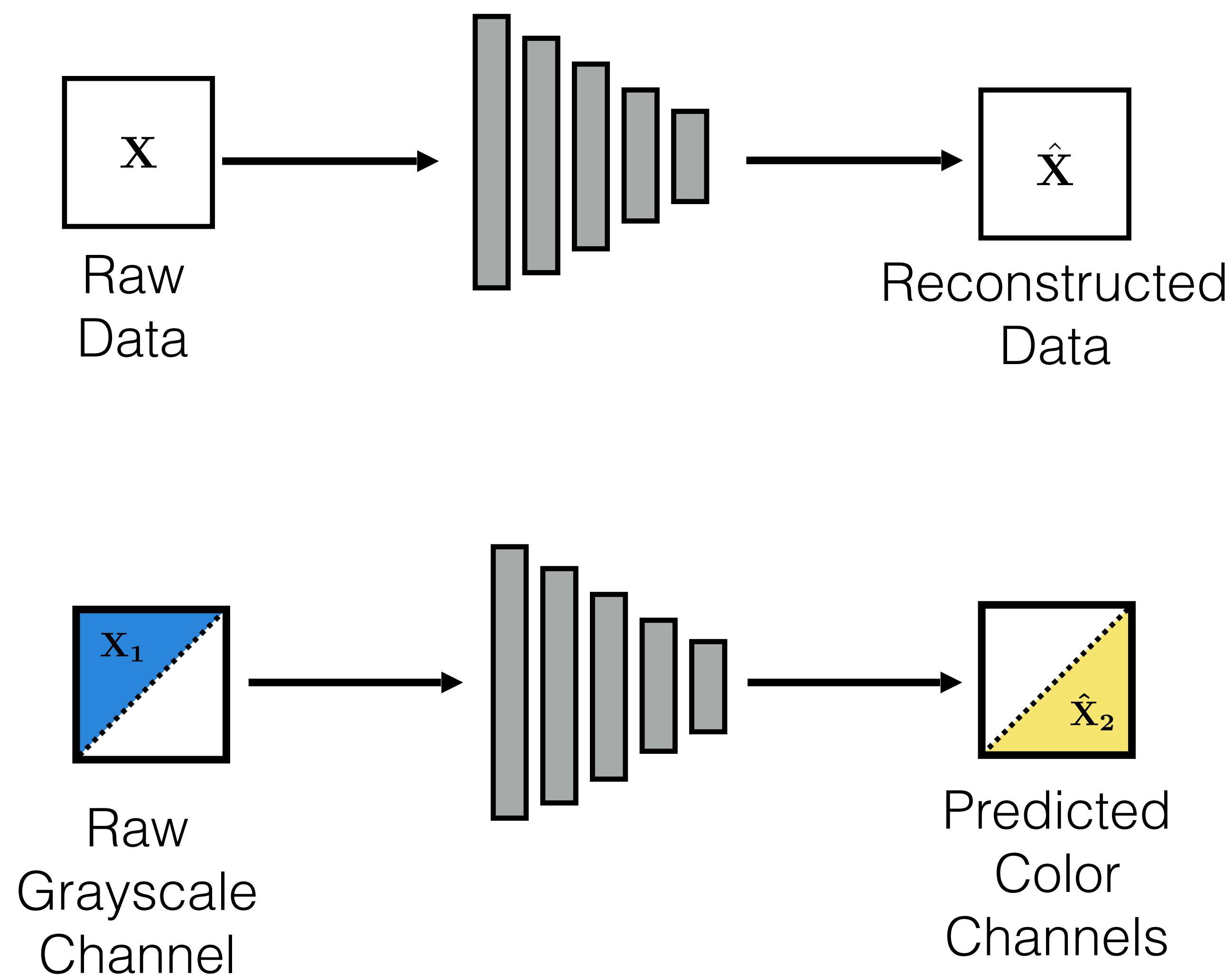
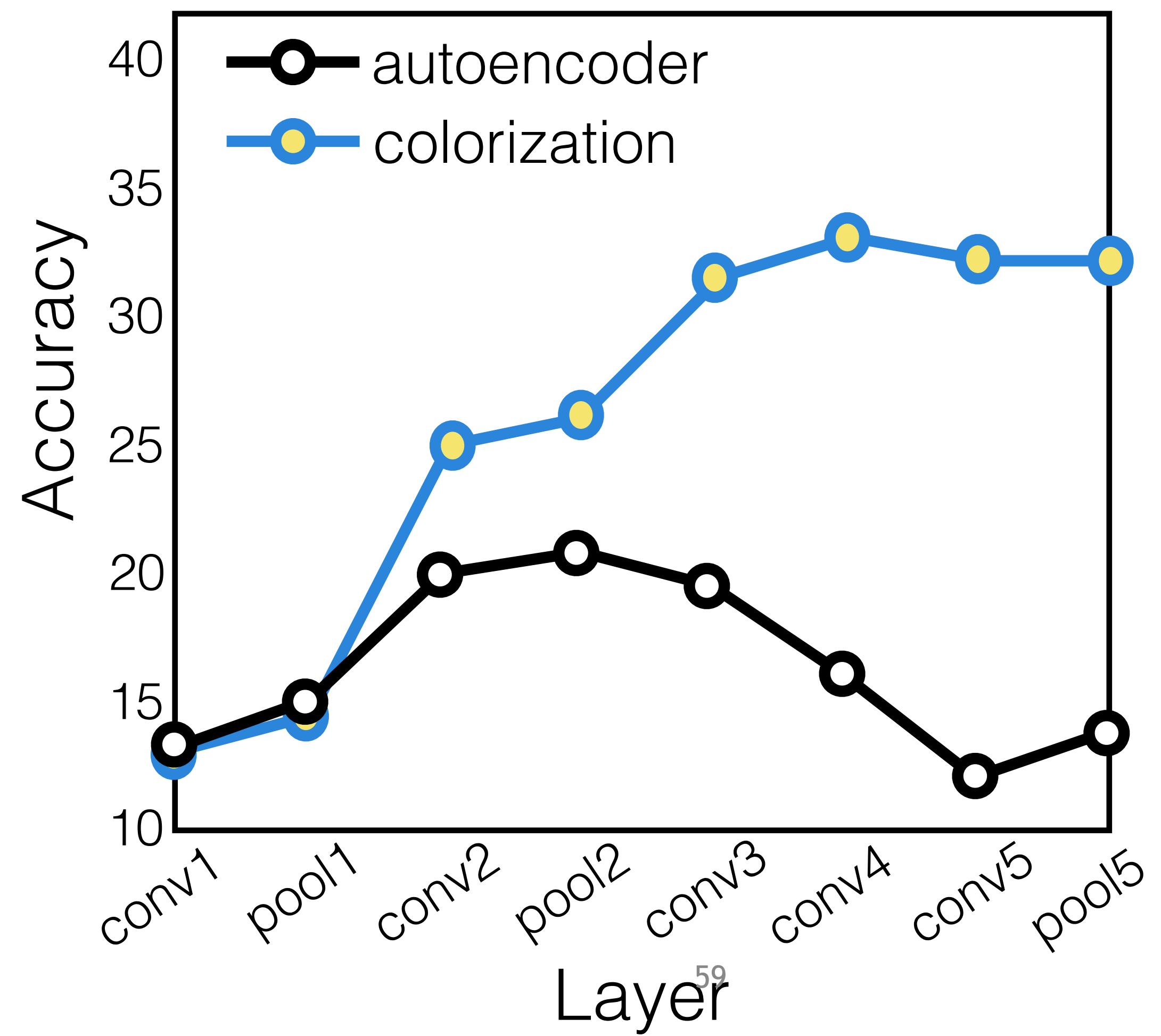


56

[Zhang, Isola, Efros, ECCV 2016]

Classification performance

ImageNet Task [Russakovsky et al. 2015]

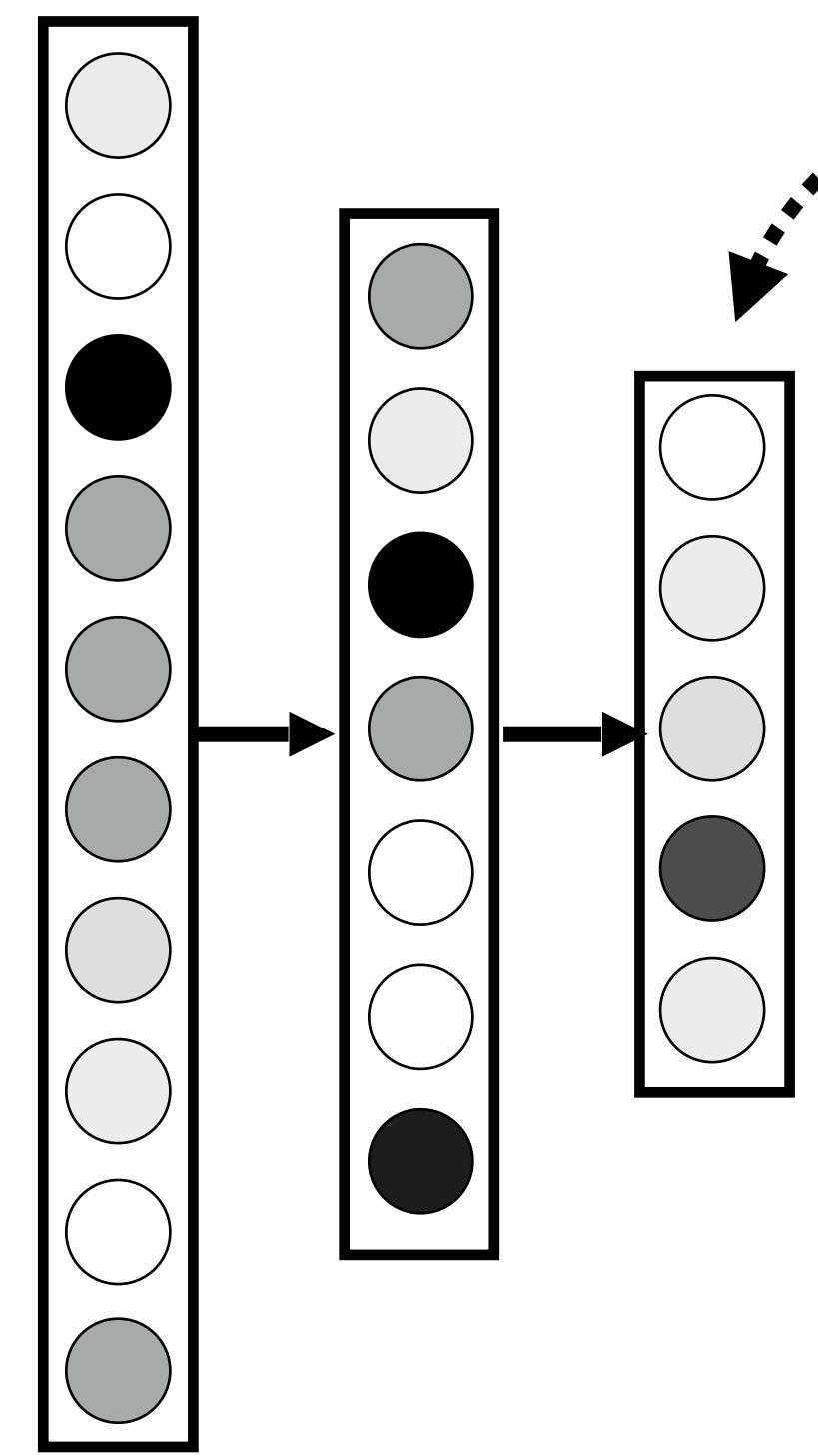
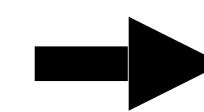


im2vec

X

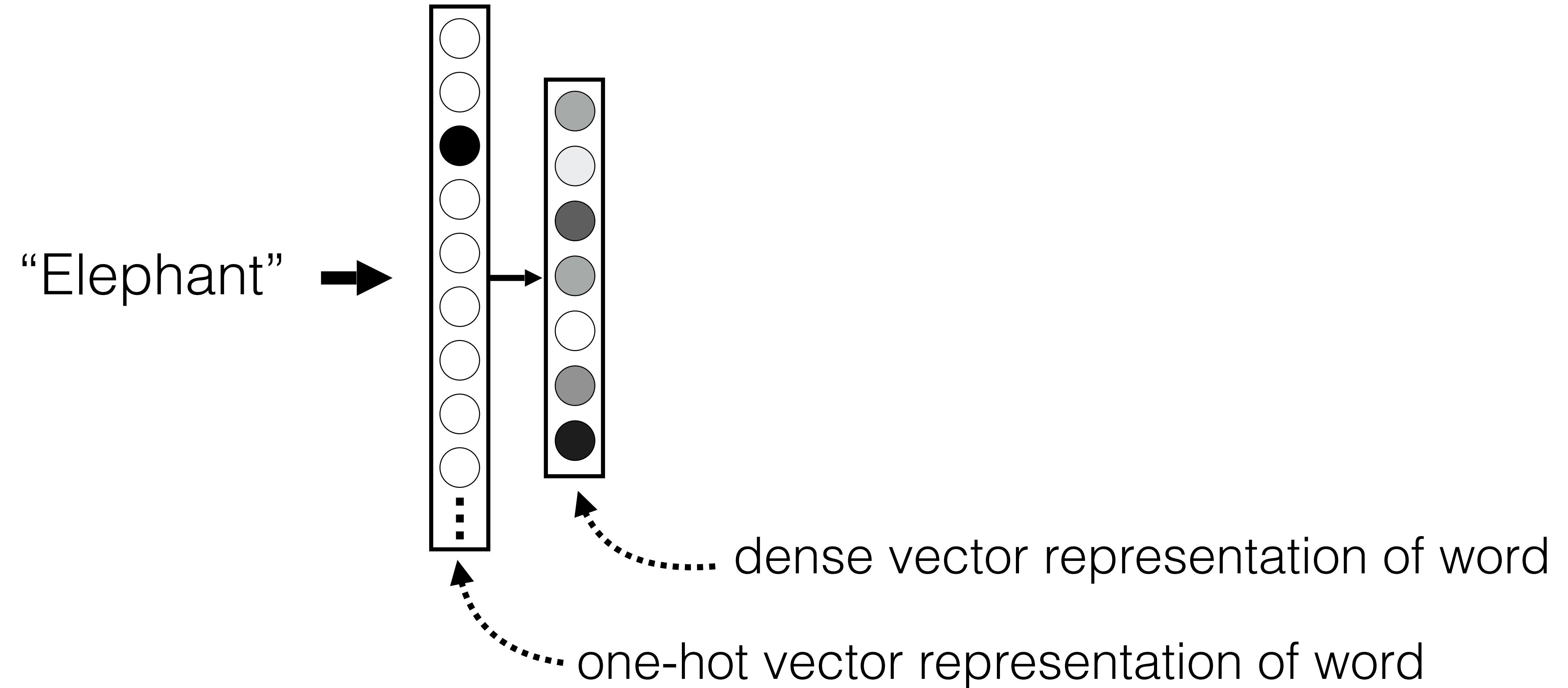


Image

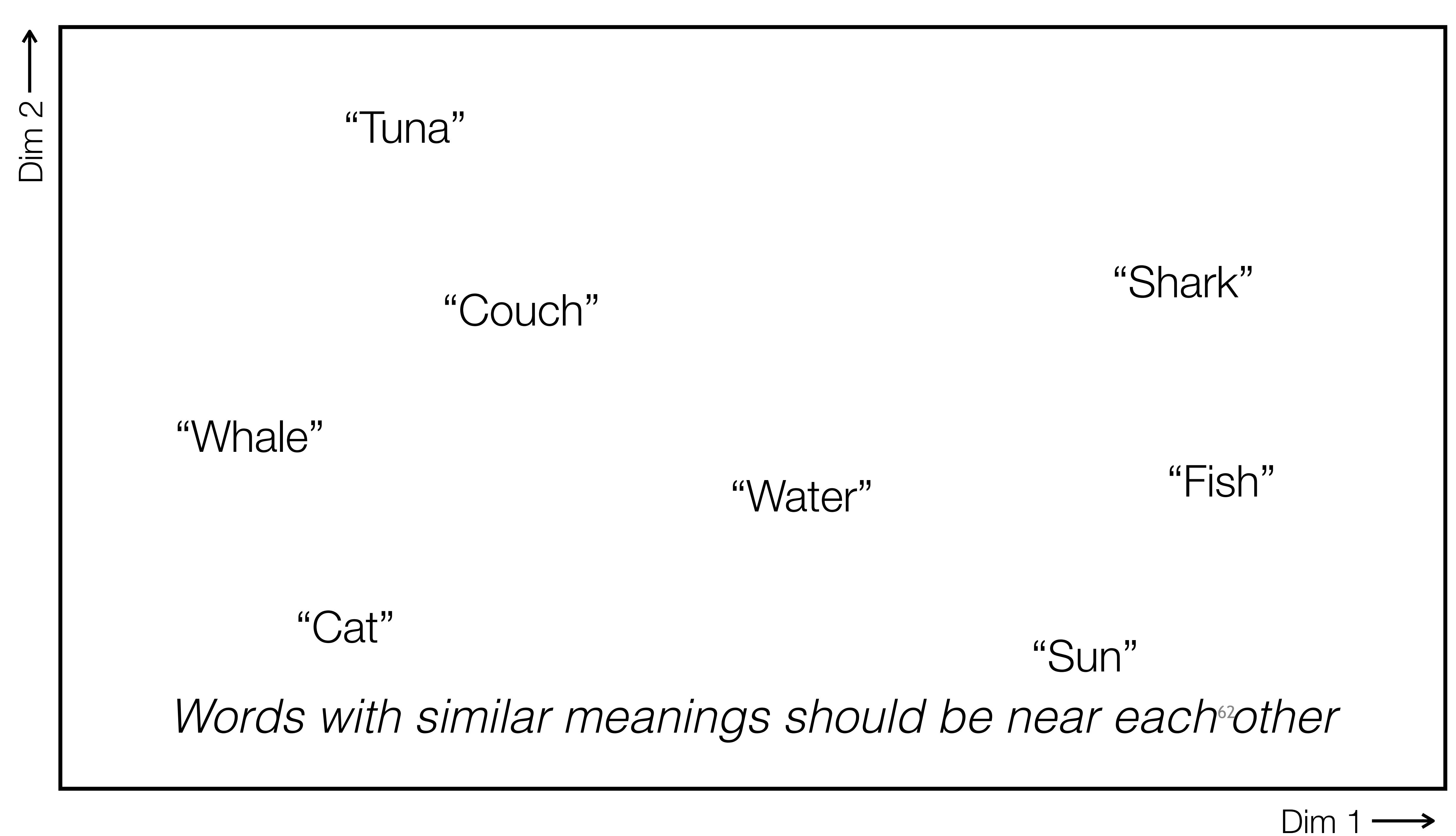


Represent image as a vector of neural activations
(perhaps representing a vector of detected texture patterns or object parts)

word2vec



X2vec methods are also called embeddings of X, e.g., a **word embedding**⁶¹



word2vec

Words with similar meanings should be near each other

Proxy: words that are used in the same context tend to have similar meanings

words with similar contexts should be near each other

Next to the 'sofa' is a desk, and a 'person' is sitting behind it.

'armchair'

'bench'

'chair'

'deck chair'

'ottoman'

'seat'

'stool'

'swivel chair'

'loveseat'

...

'man'

'woman'

'child'

'teenager'

'girl'

'boy'

'baby'

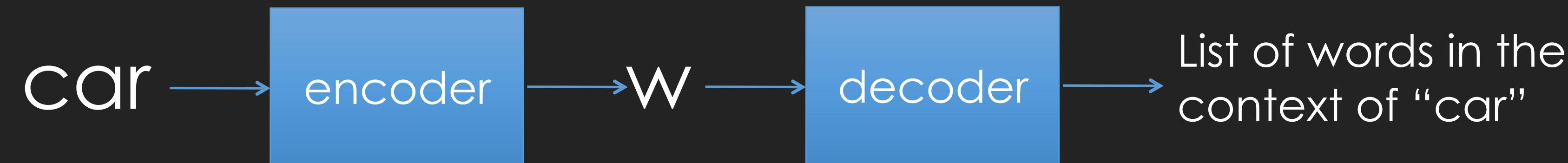
'daughter'

'son'

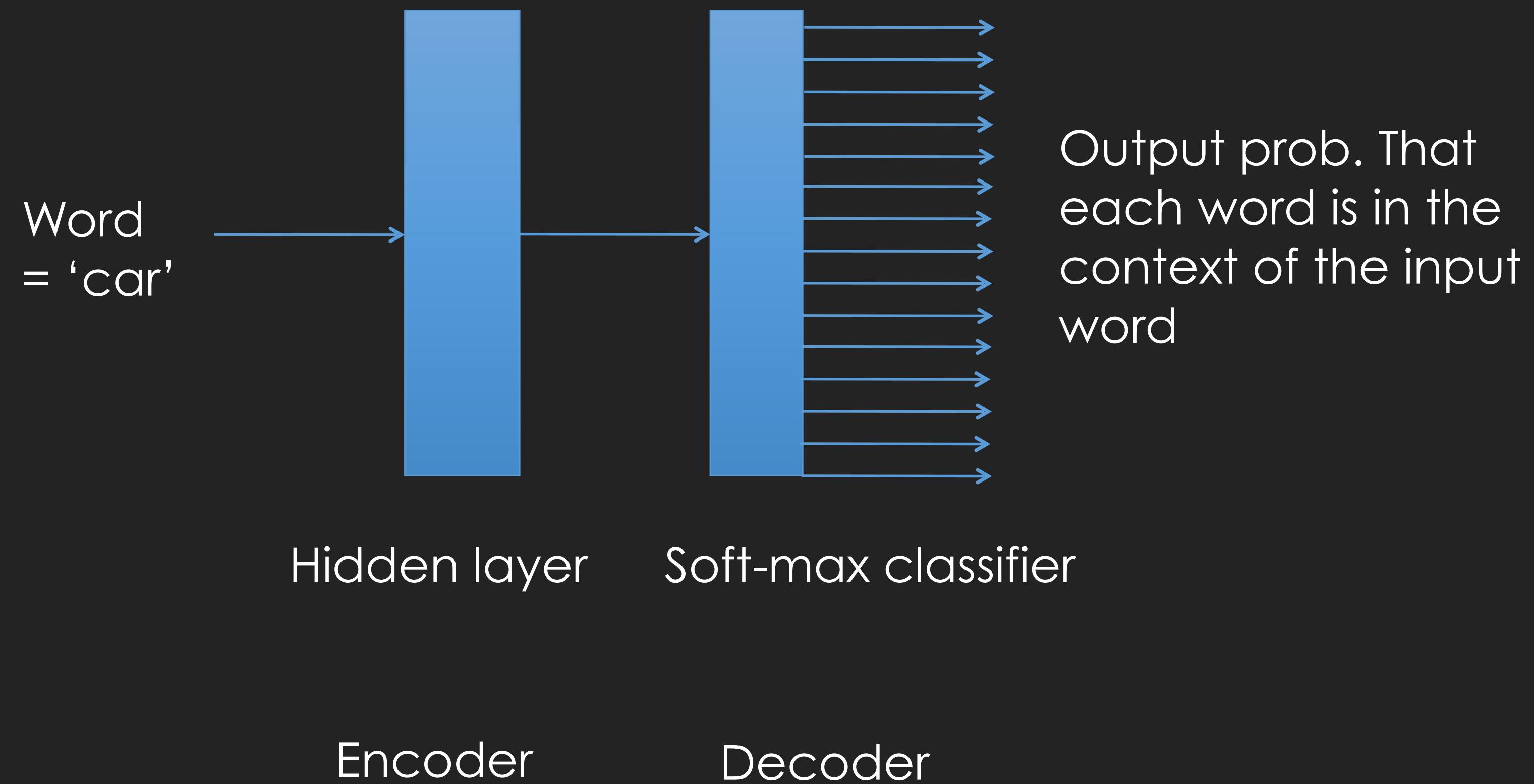
...

word2vec

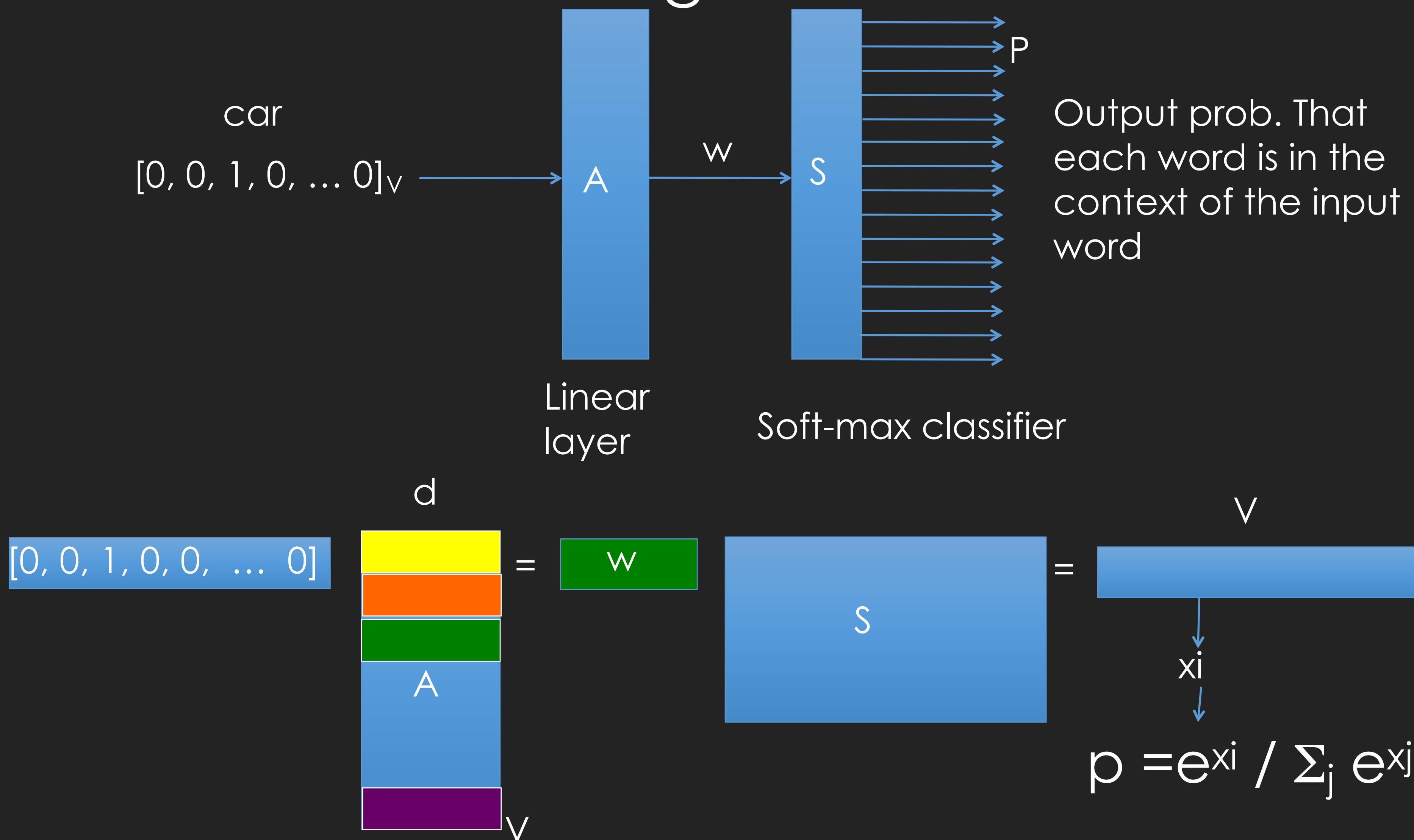
I parked the **car** in a nearby street. It is a red **car** with two doors, ...



word2vec



word2vec, training



Algebraic operations with the vector representation of words

$$X = \text{Vector}(\text{"Paris"}) - \text{vector}(\text{"France"}) + \text{vector}(\text{"Italy"})$$

Closest nearest neighbor to X is $\text{vector}(\text{"Rome"})$

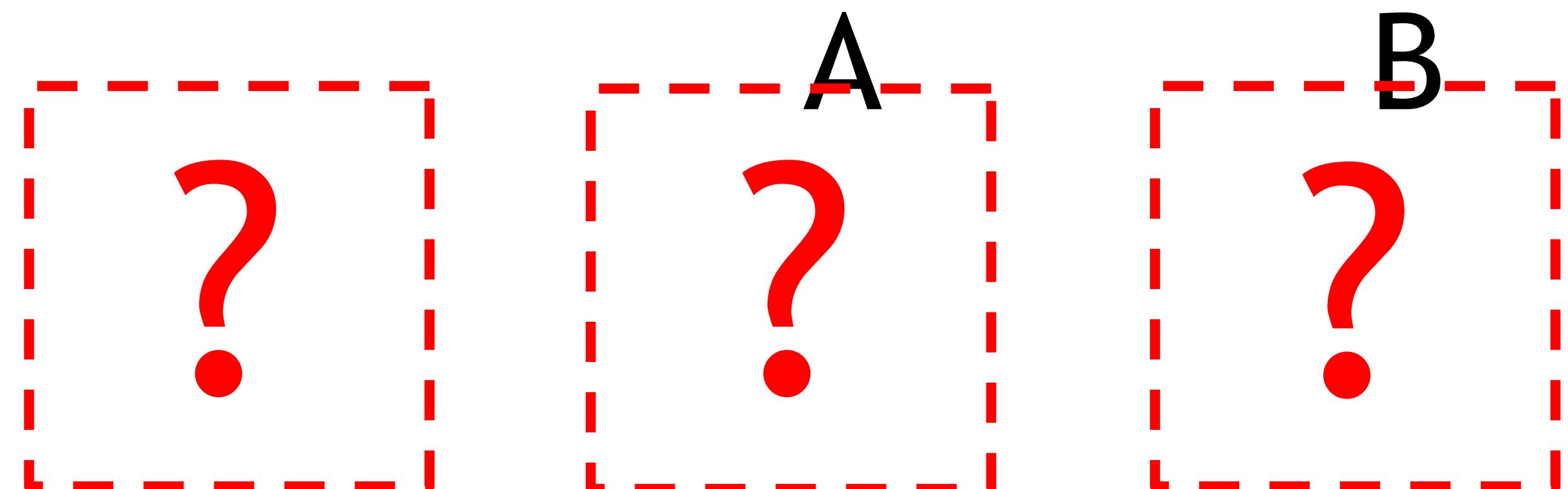
Context as Supervision

[Collobert & Weston 2008; Mikolov et al. 2013]

house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal milk; but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and appliances, the toasters and carpet sweepers of Lilliput, but she knew that most adult vis

Deep
Net

Context Prediction as Supervision

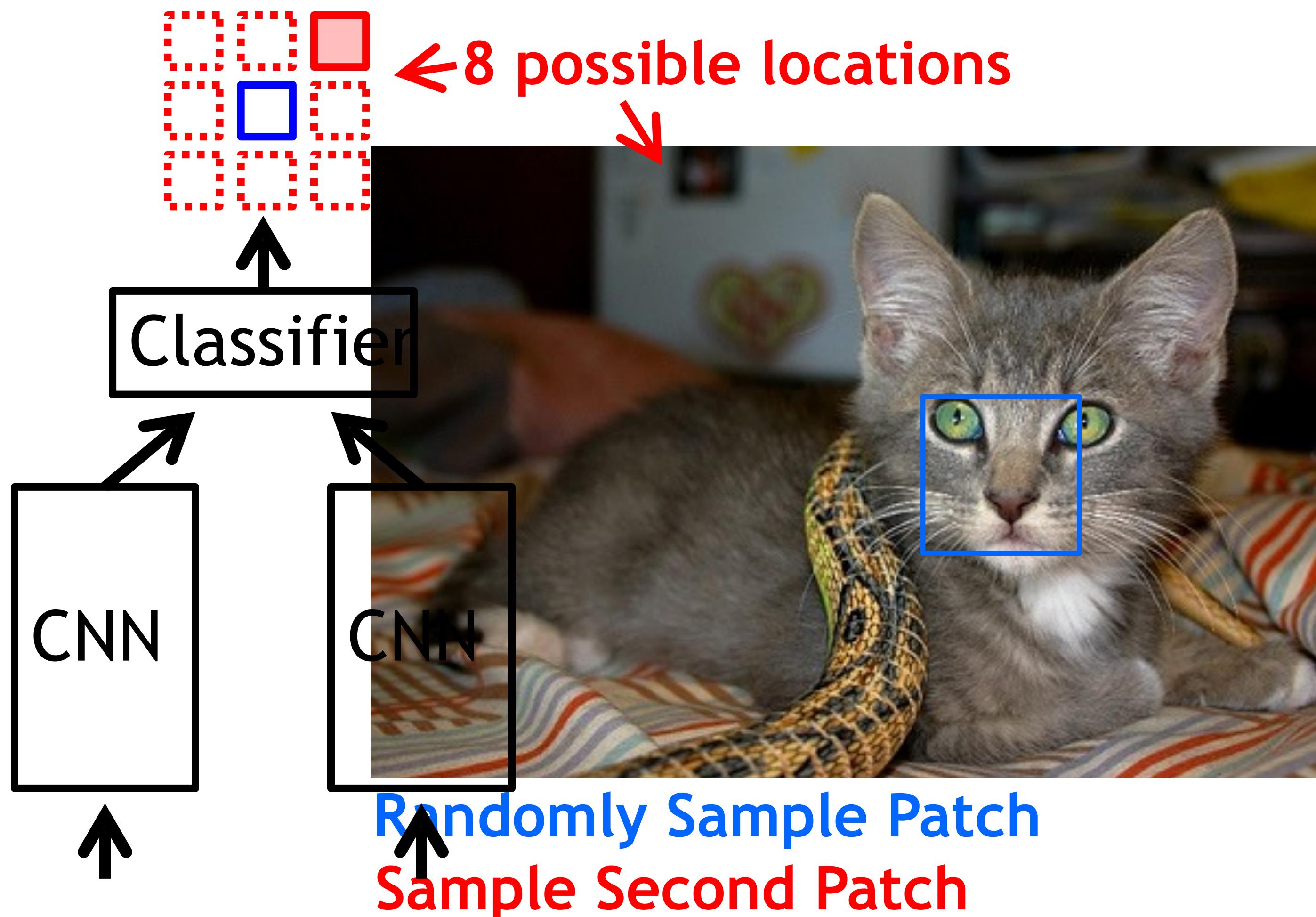


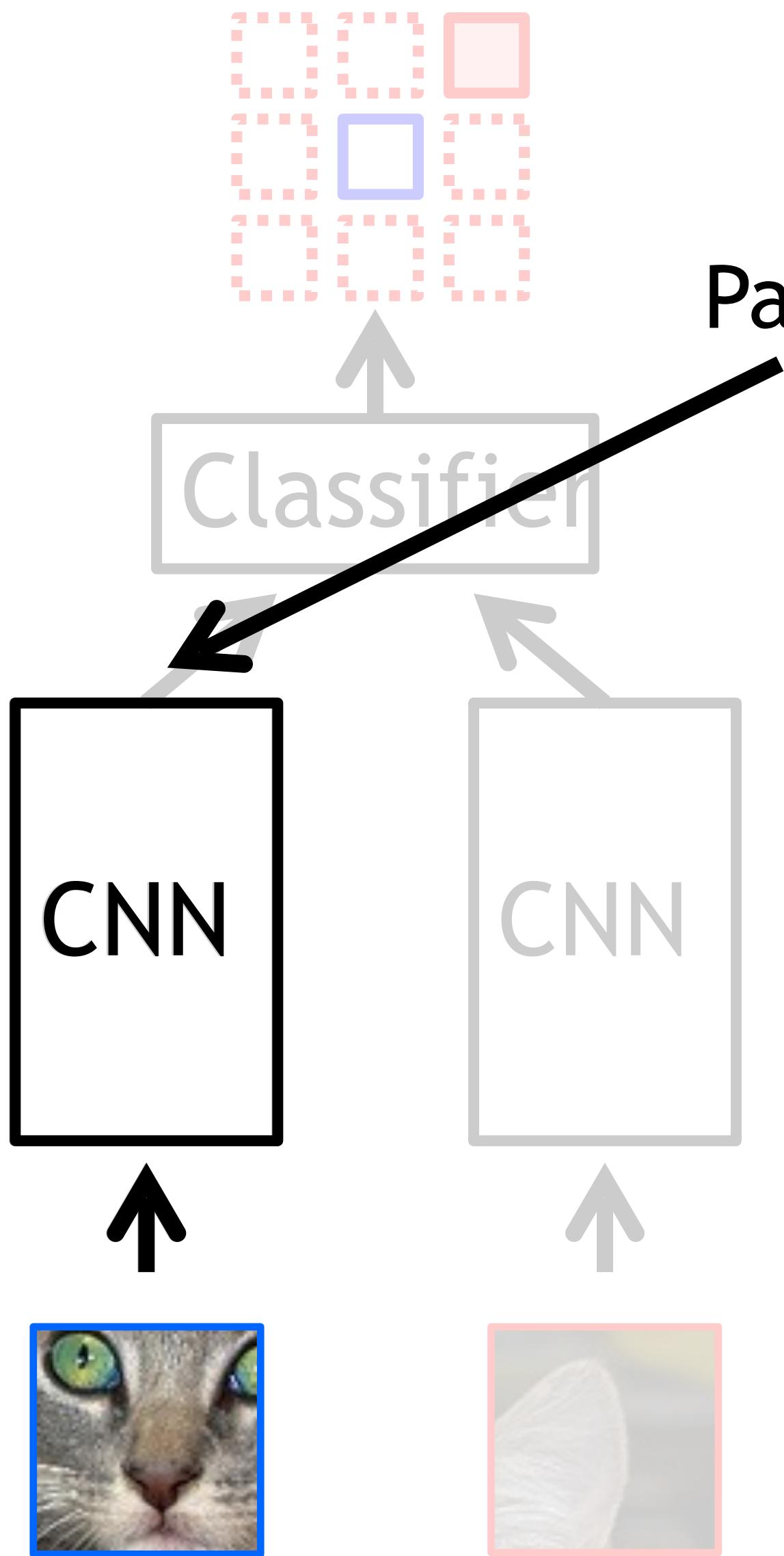
Semantics from a non-semantic task



[Slide credit: Carl Doersch] ⁷²

Relative Position Task





Patch Embedding (representation)

Input



Nearest Neighbors



Note: connects *across* instances!

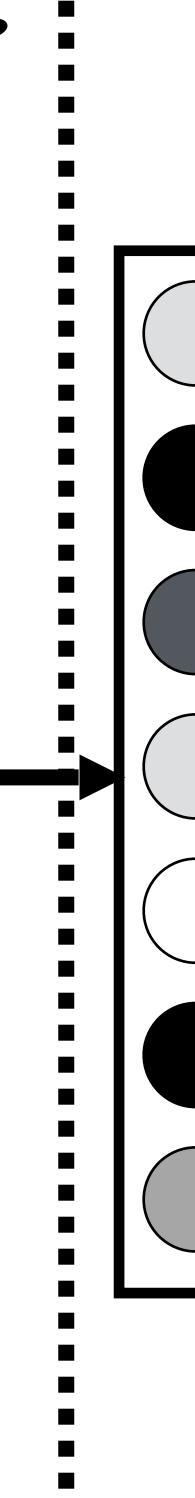
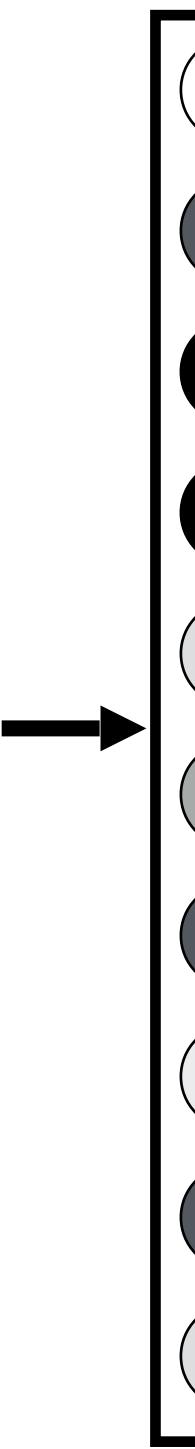
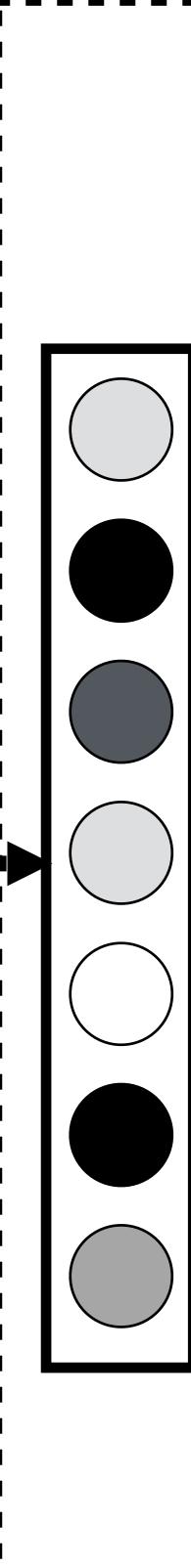
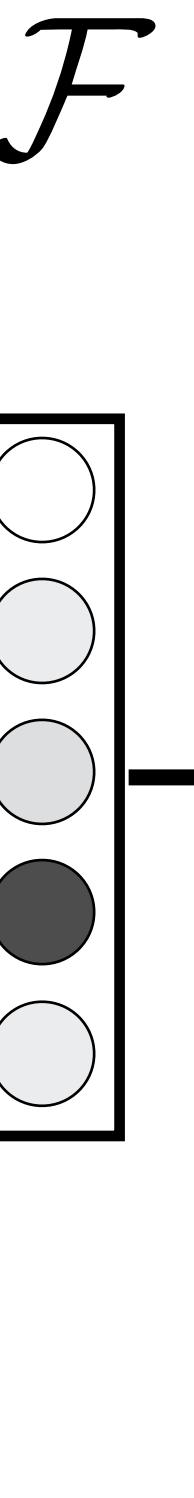
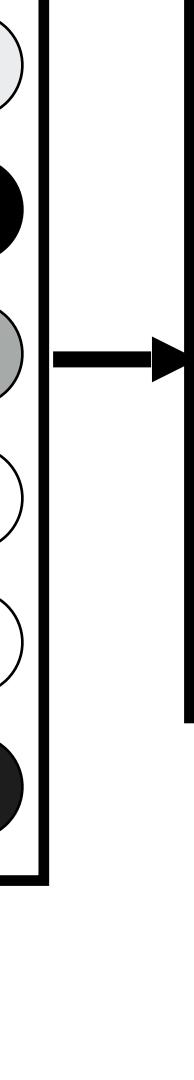
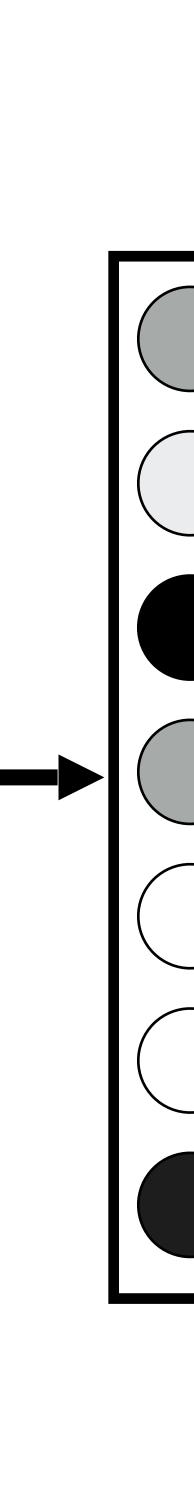
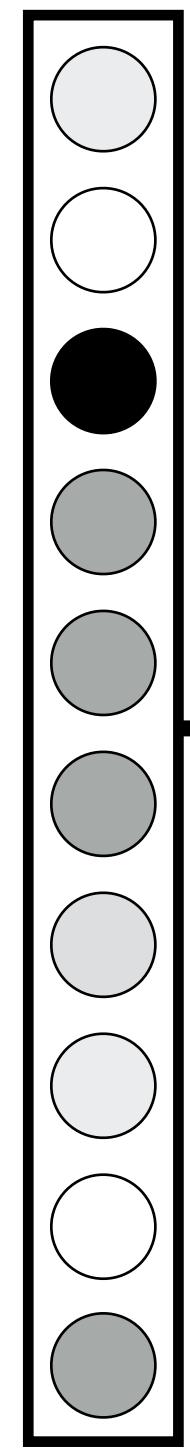
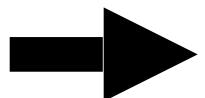
Revisiting autoencoders

Is reconstruction necessary?

\mathbf{X}



Image

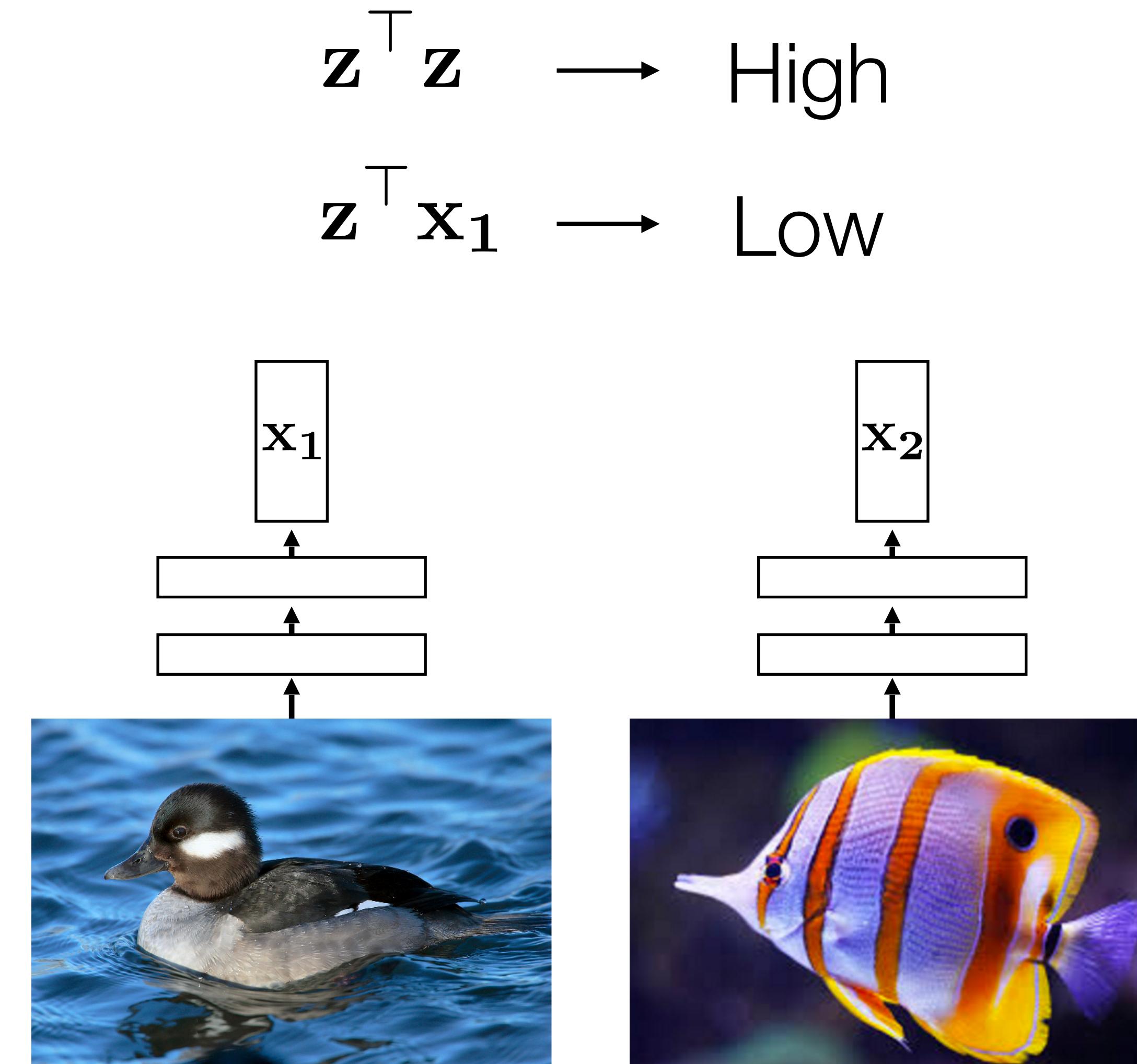
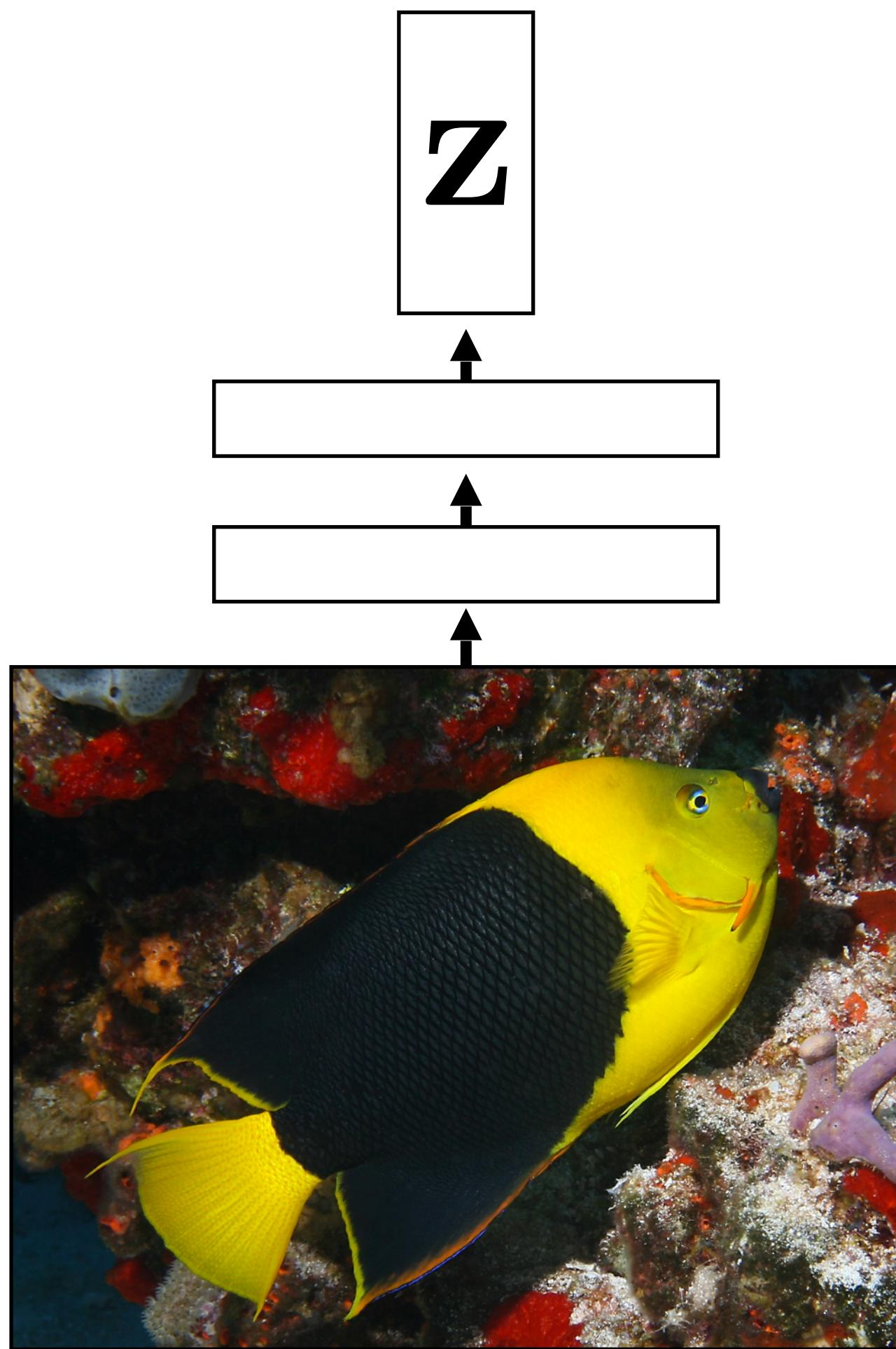


$\hat{\mathbf{X}} = \mathcal{F}(\mathbf{X})$

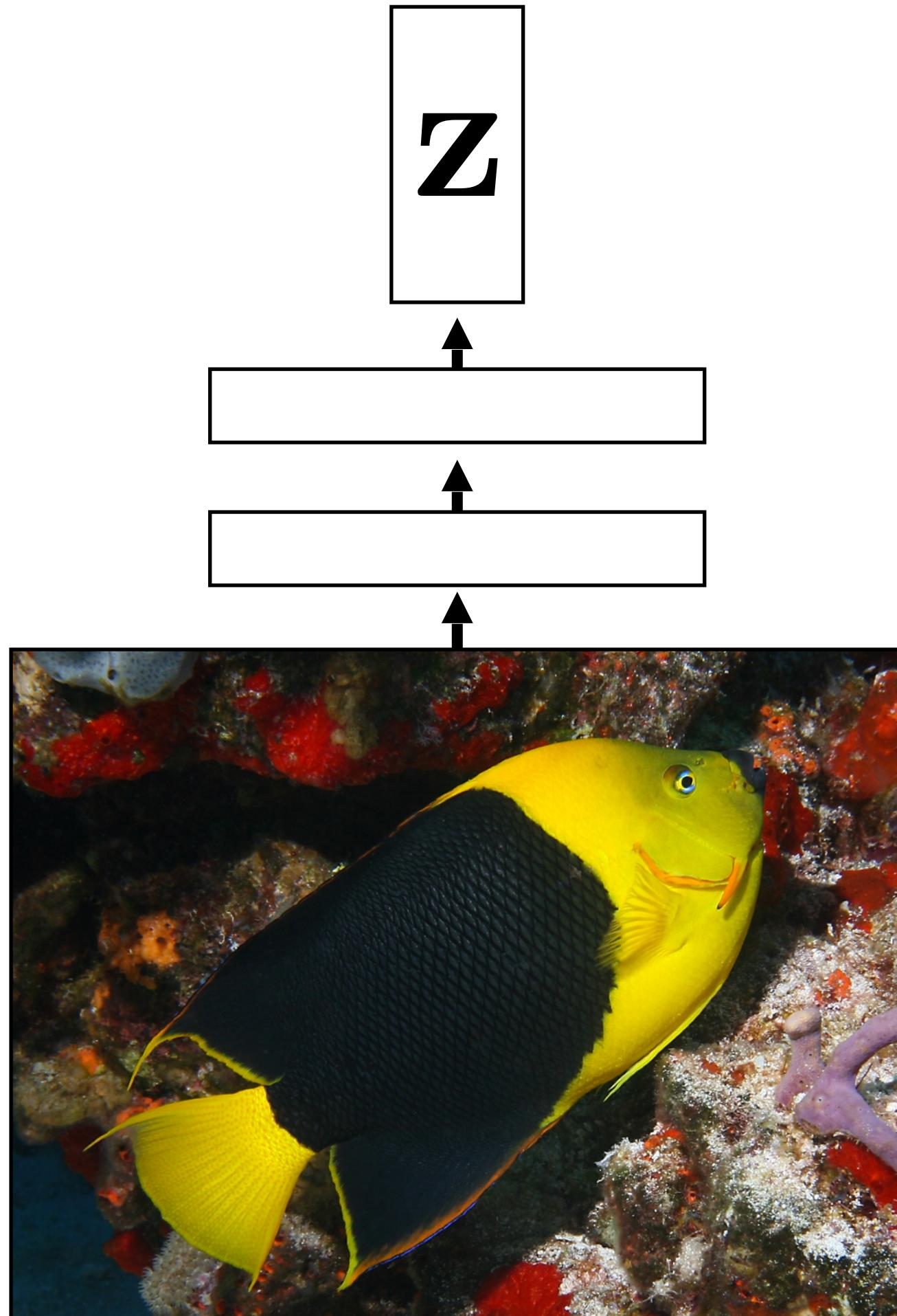


Reconstructed
image

Contrastive learning



Contrastive learning

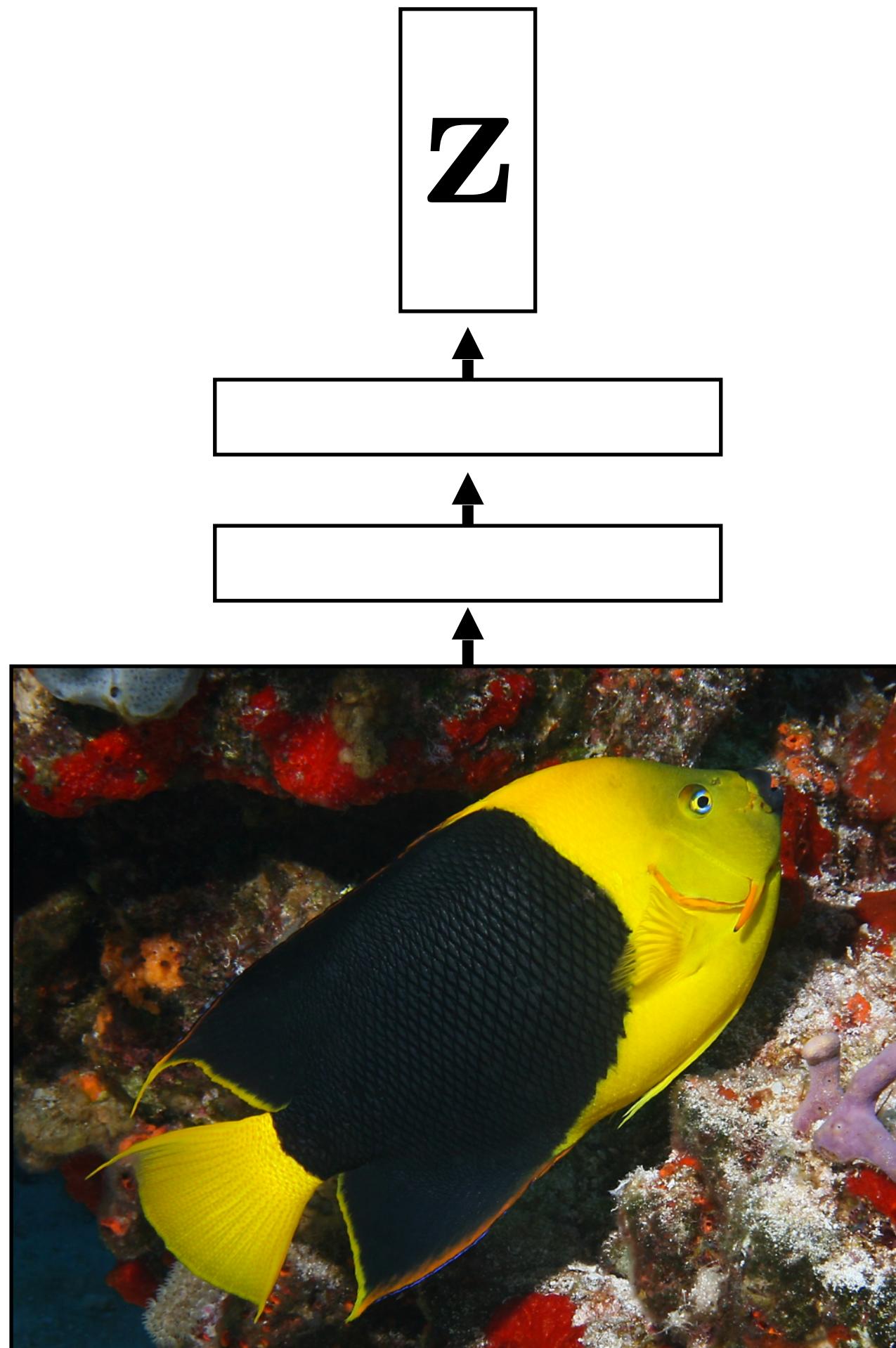


Maximize:

$$\frac{\exp\{z^T z\}}{\exp\{z^T z + \sum_i z^T x_i\}}$$

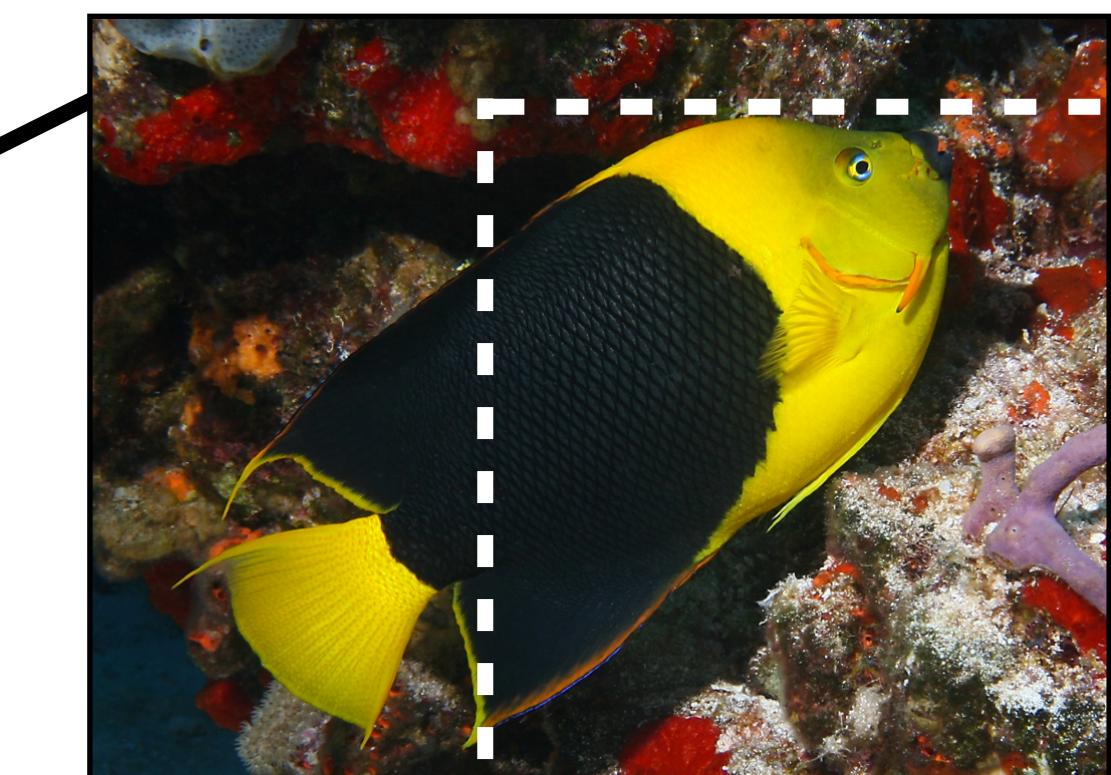
Equivalent to softmax loss with
each image as a category

Contrastive learning

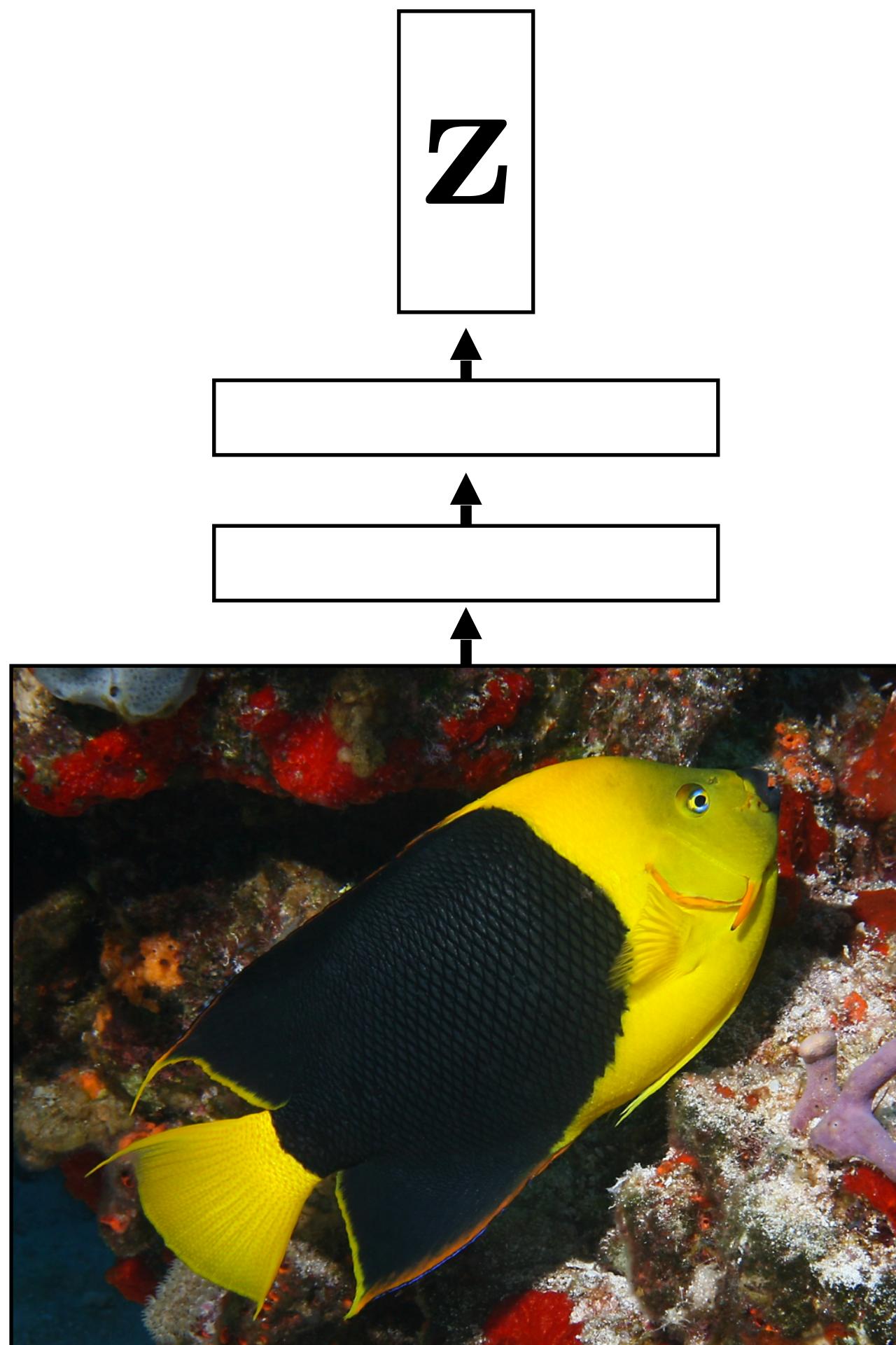


Can build invariance.
Compare to warped images.

$$\frac{\exp\{z^\top \tilde{z}\}}{\exp\{z^\top \tilde{z} + \sum_i z^\top x_i\}}$$

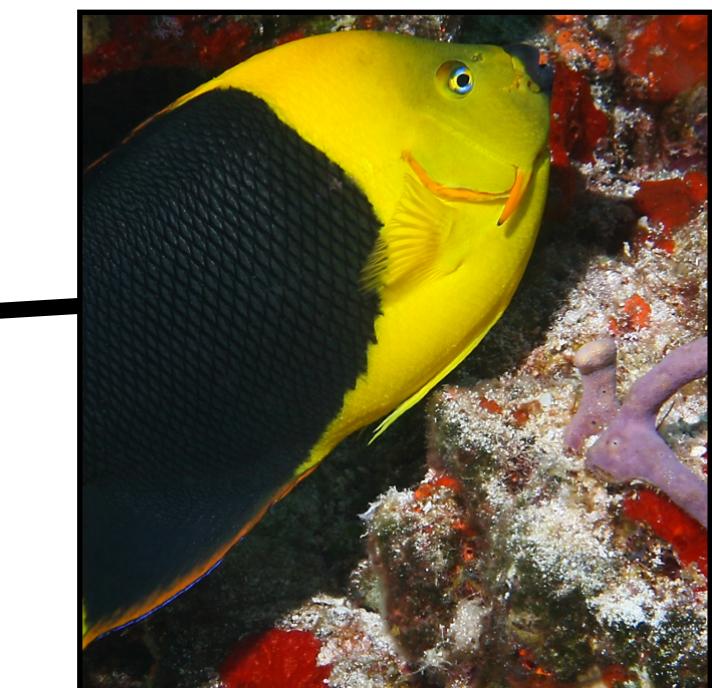


Contrastive learning



Can build invariance.
Compare to warped images.

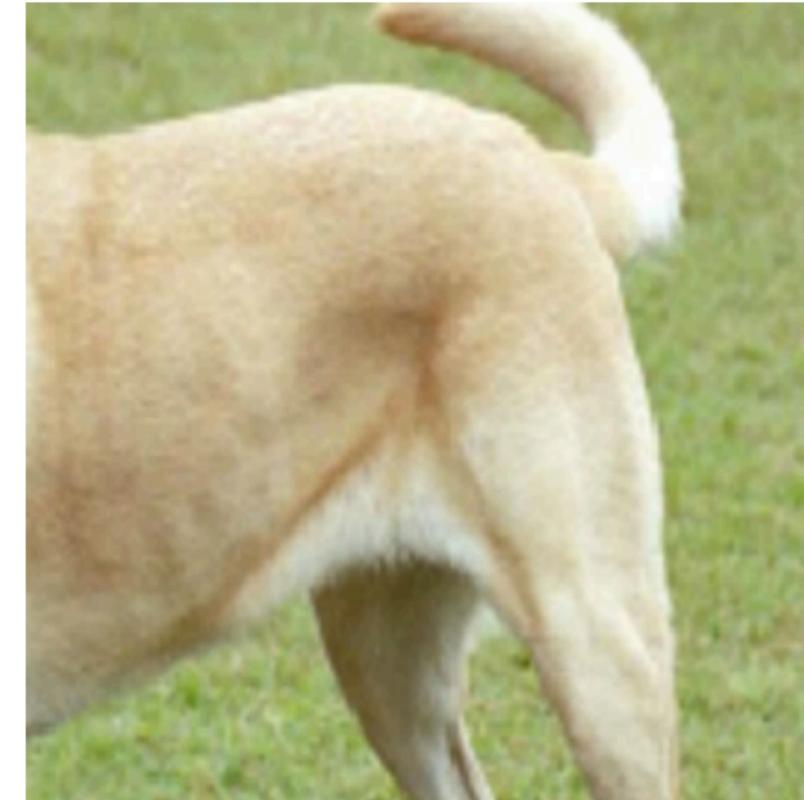
$$\frac{\exp\{z^\top \tilde{z}\}}{\exp\{z^\top \tilde{z} + \sum_i z^\top x_i\}}$$



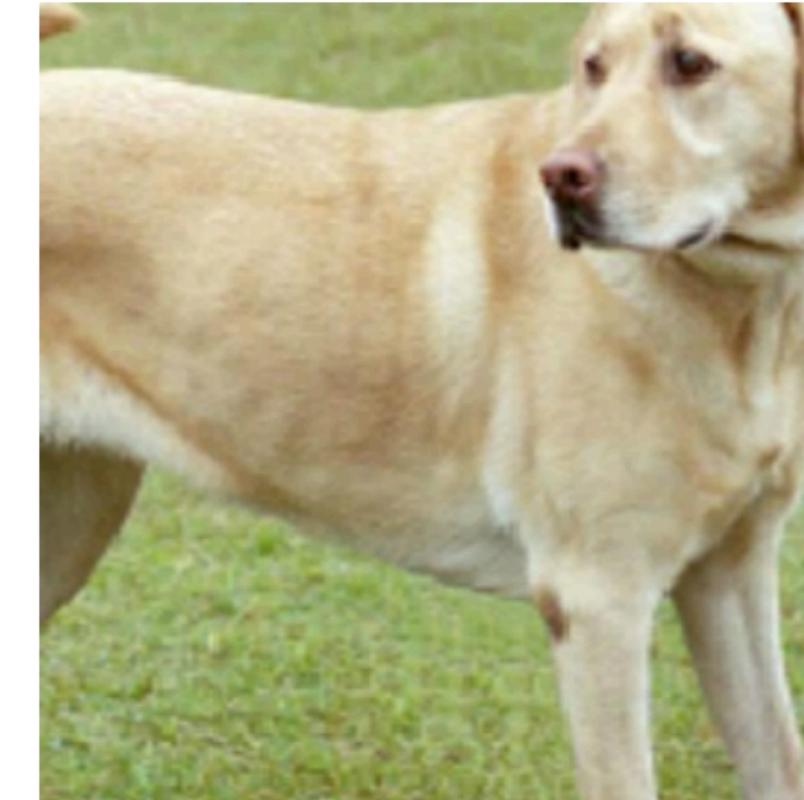
Contrastive learning



(a) Original



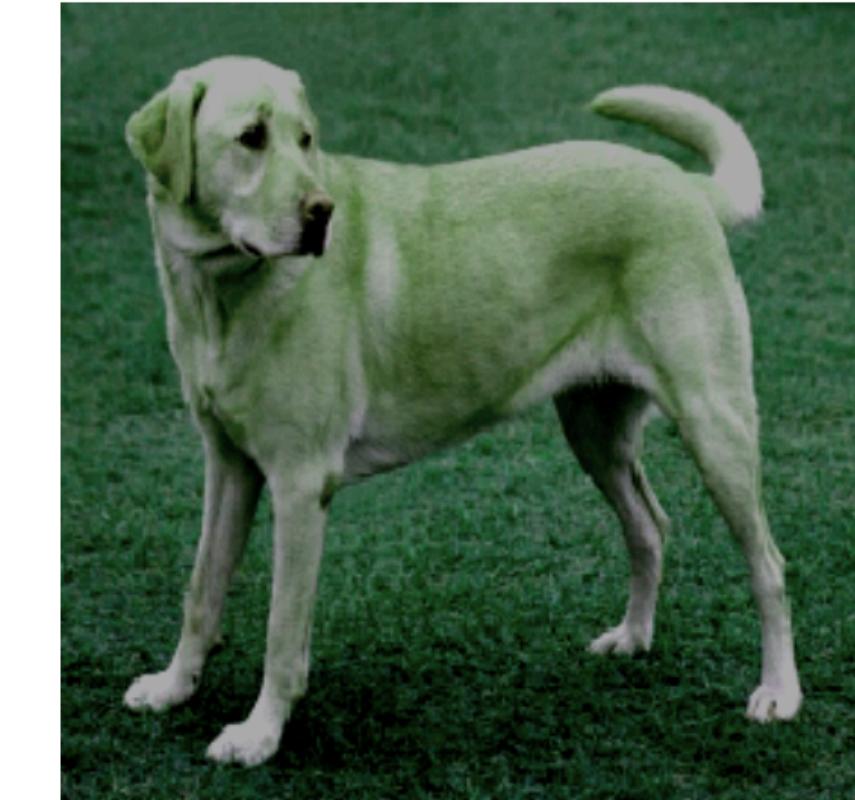
(b) Crop and resize



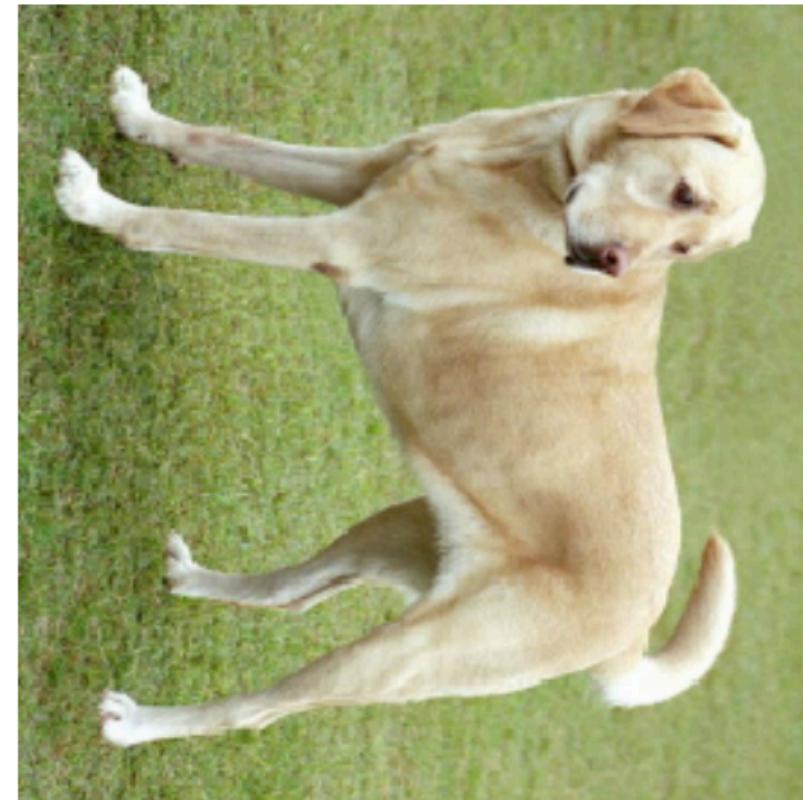
(c) Crop, resize (and flip)



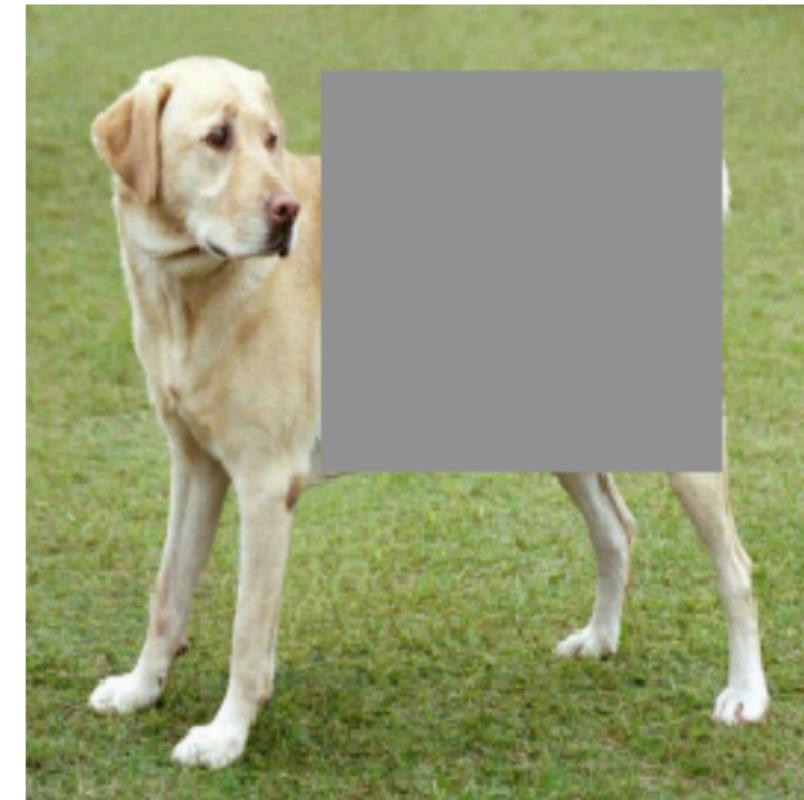
(d) Color distort. (drop)



(e) Color distort. (jitter)



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



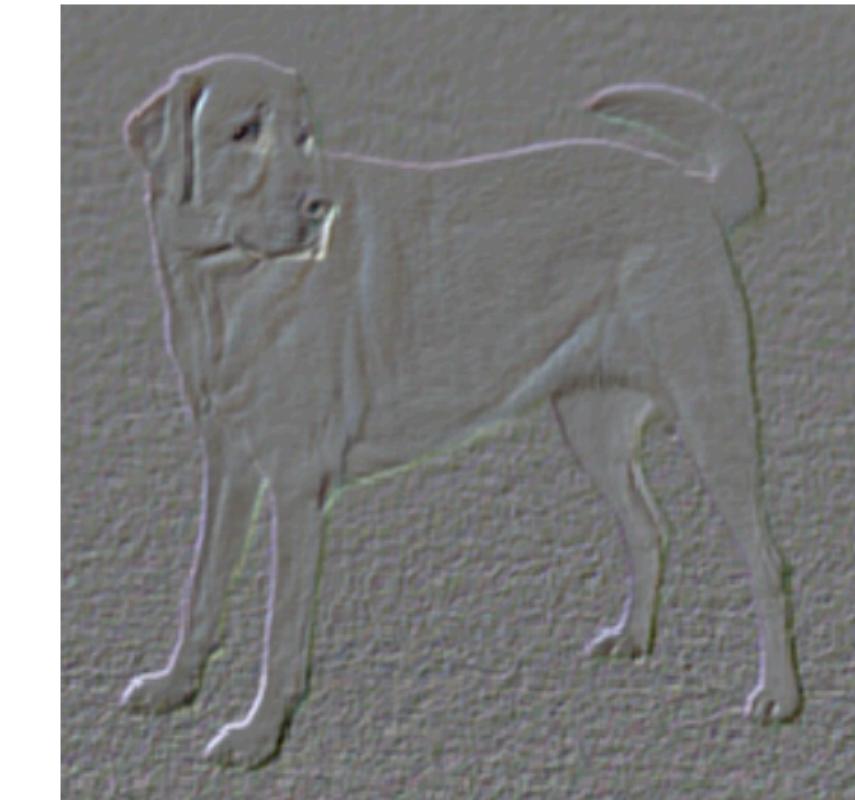
(g) Cutout



(h) Gaussian noise



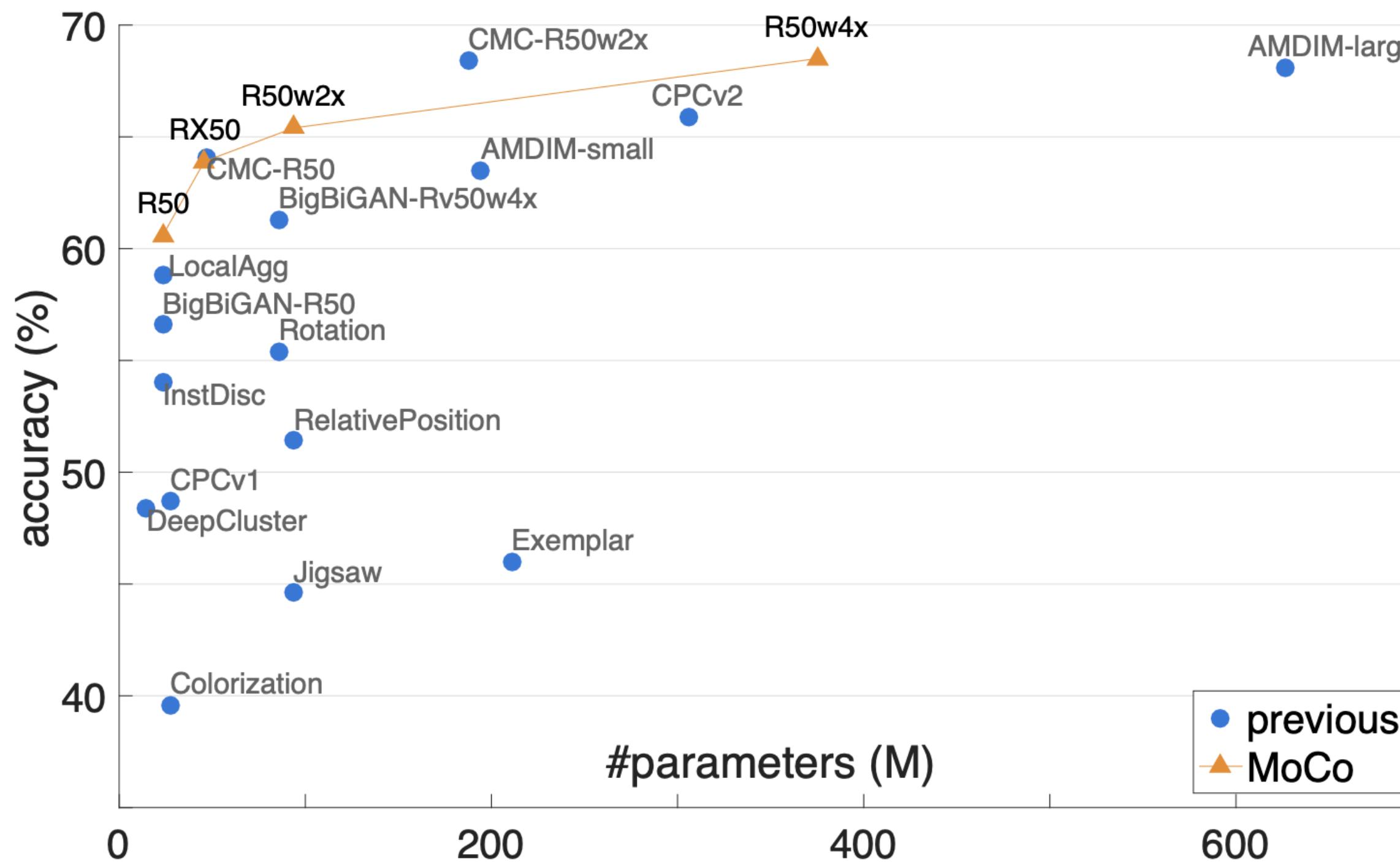
(i) Gaussian blur



(j) Sobel filtering

From [Chen et al.,⁸⁰ SimCLR, 2020]

Performance snapshot



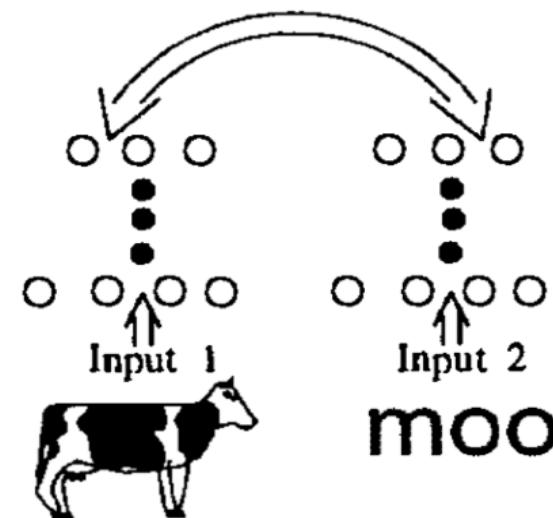
ImageNet linear classification

pre-train	AP ₅₀
random init.	52.5
super. IN-1M	80.8
MoCo IN-1M	81.4 (+0.6)
MoCo IG-1B	82.1 (+1.3)

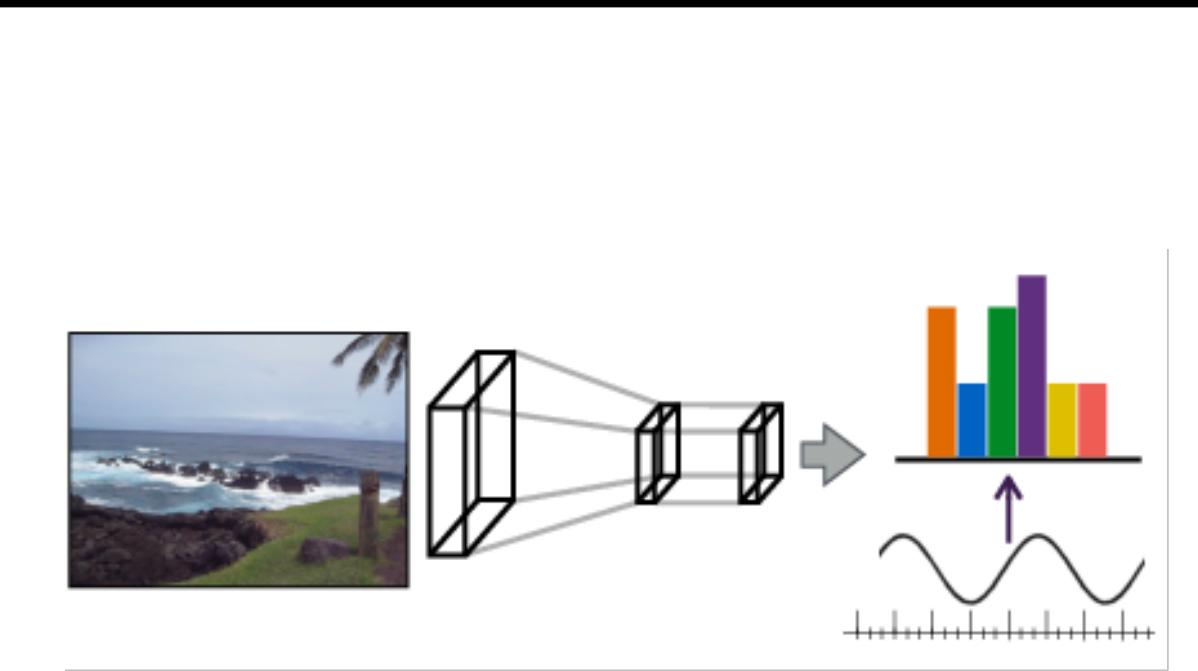
Object detection finetuning

Comparable in many cases to supervised pretraining!

Audio

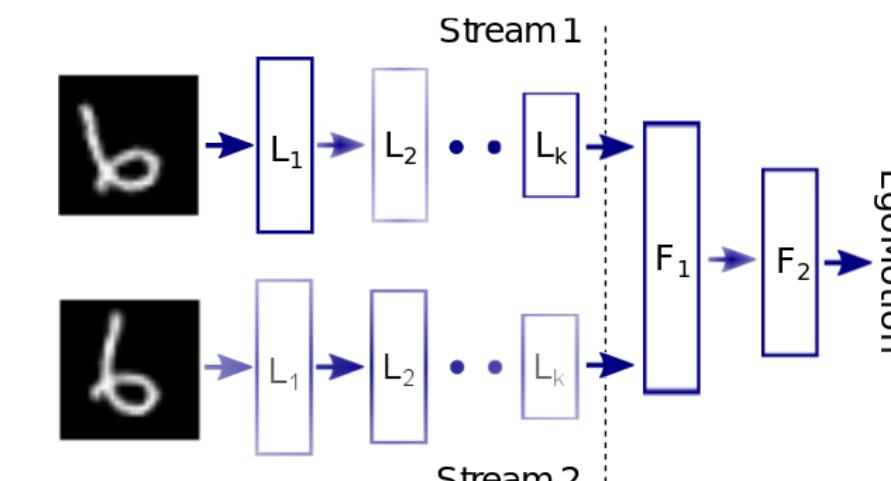


de Sa. NIPS 1994.

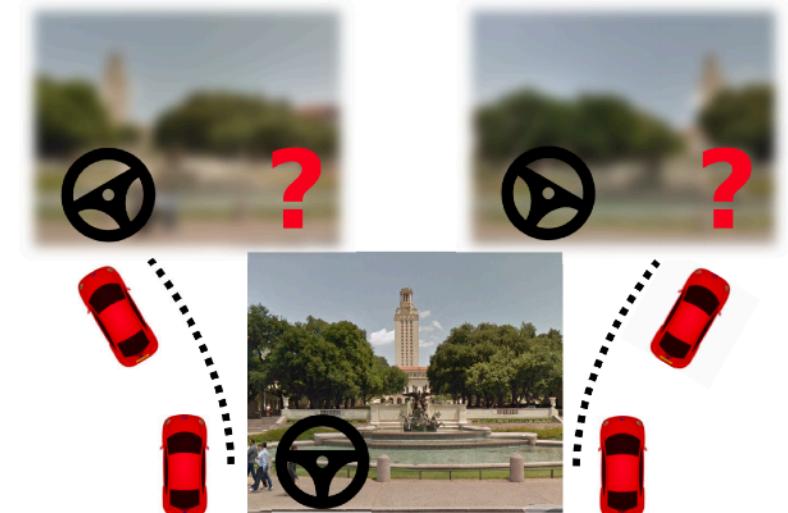


Owens et al. ECCV 2016.

Egomotion



Agrawal et al. ICCV 2015.

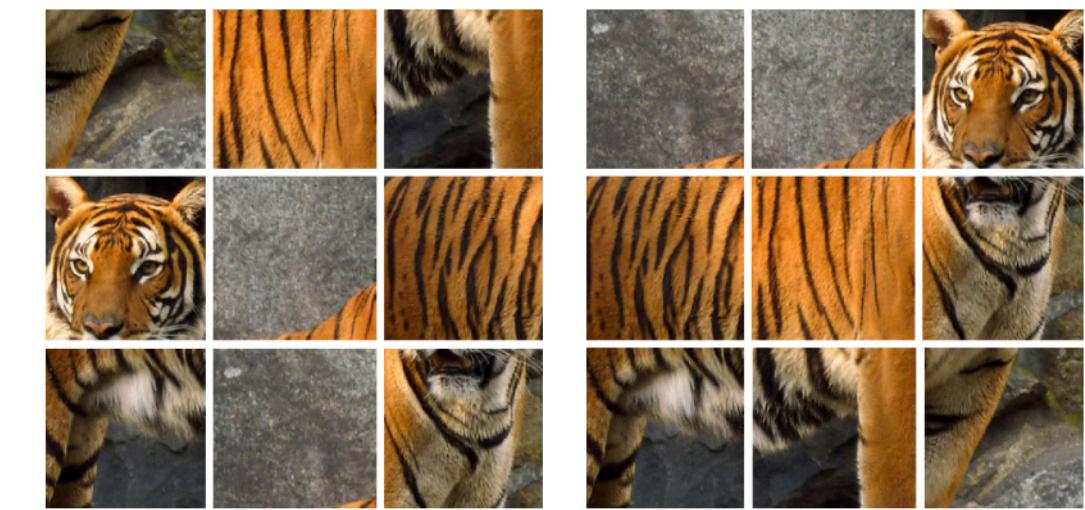


Jayaraman et al. ICCV 2015.

Context

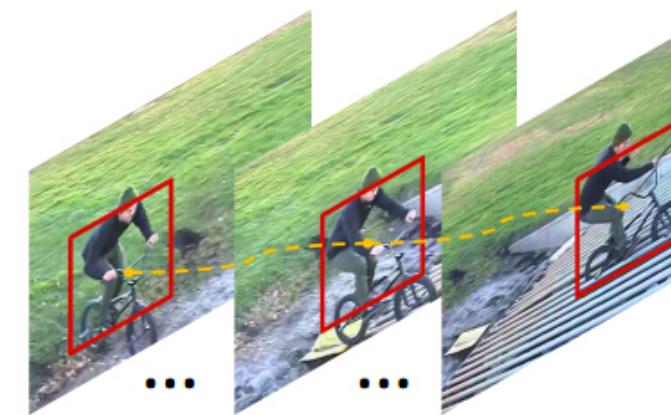


Pathak et al. CVPR 2016.

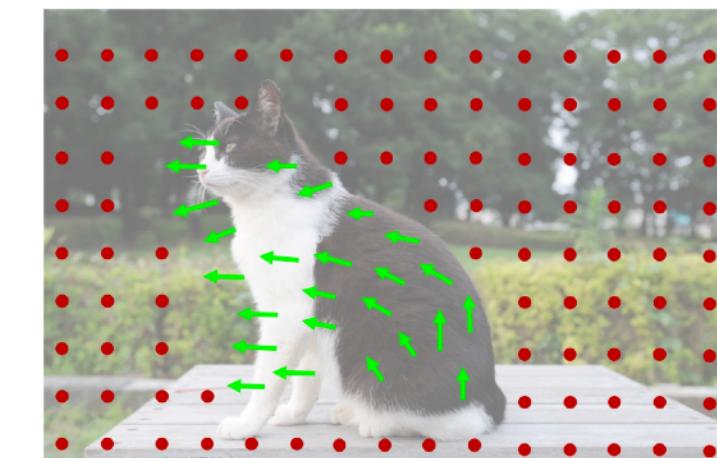


Noroozi and Favaro. ECCV 2016.
Doersch et al. ICCV 2015.

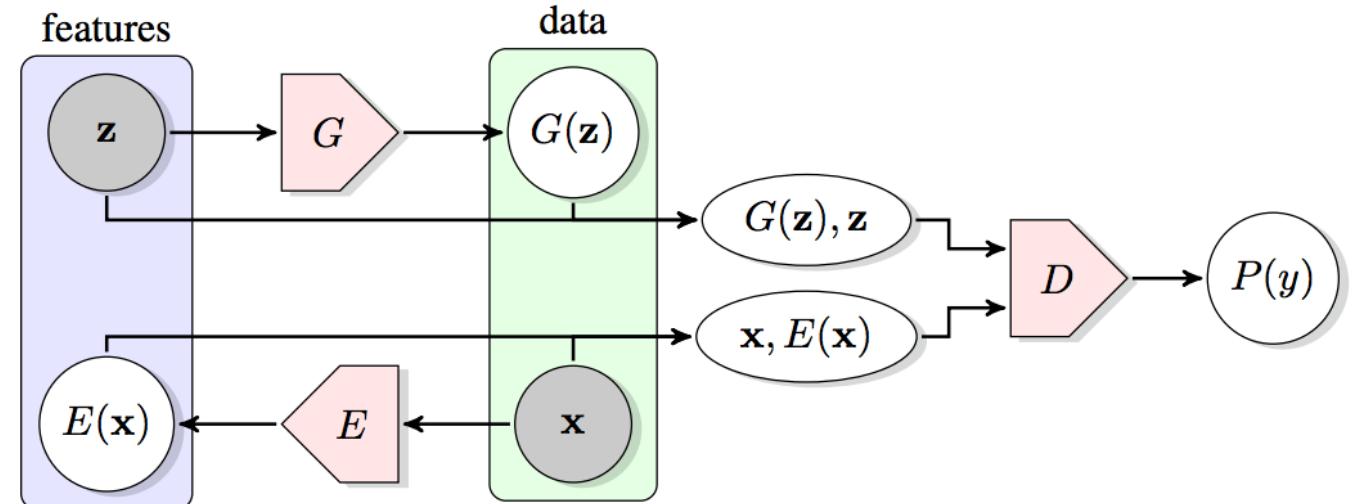
Video



Wang et al. ICCV 2015. Pathak et al. CVPR 2017.
Misra et al. ECCV 2016.

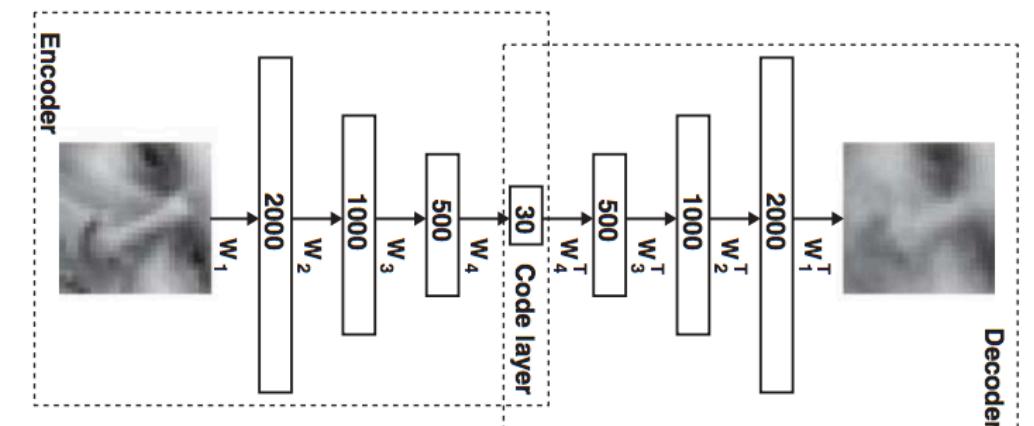


Generative Modeling



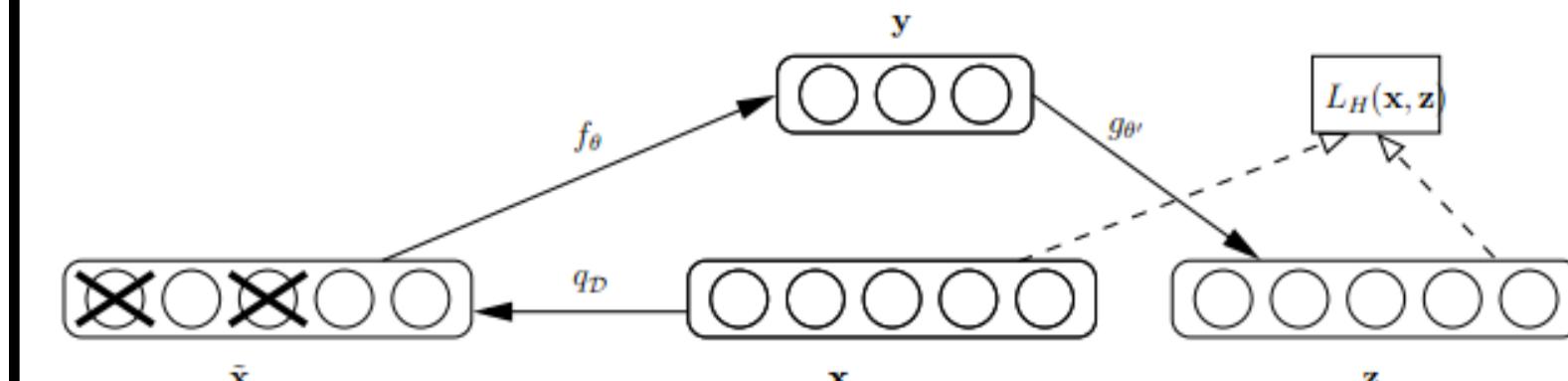
Donahue et al. Dumoulin et al. ICLR 2017.

Autoencoders



Hinton & Salakhutdinov.
Science 2006.

Denoising Autoencoders



Vincent et al. ICML 2008.

Goal: Set up a pre-training scheme to induce a “useful” representation

[Slide credit: Richard Zhang]

Summary

1. Deep nets learn *representations*
2. This is useful because representations transfer – they act as prior knowledge that enables quick learning on new tasks
3. Representations can also be learned without labels
4. Without labels there are many ways to learn representations. We saw:
 1. representations as compressed codes
 2. representations that are predictive of missing data