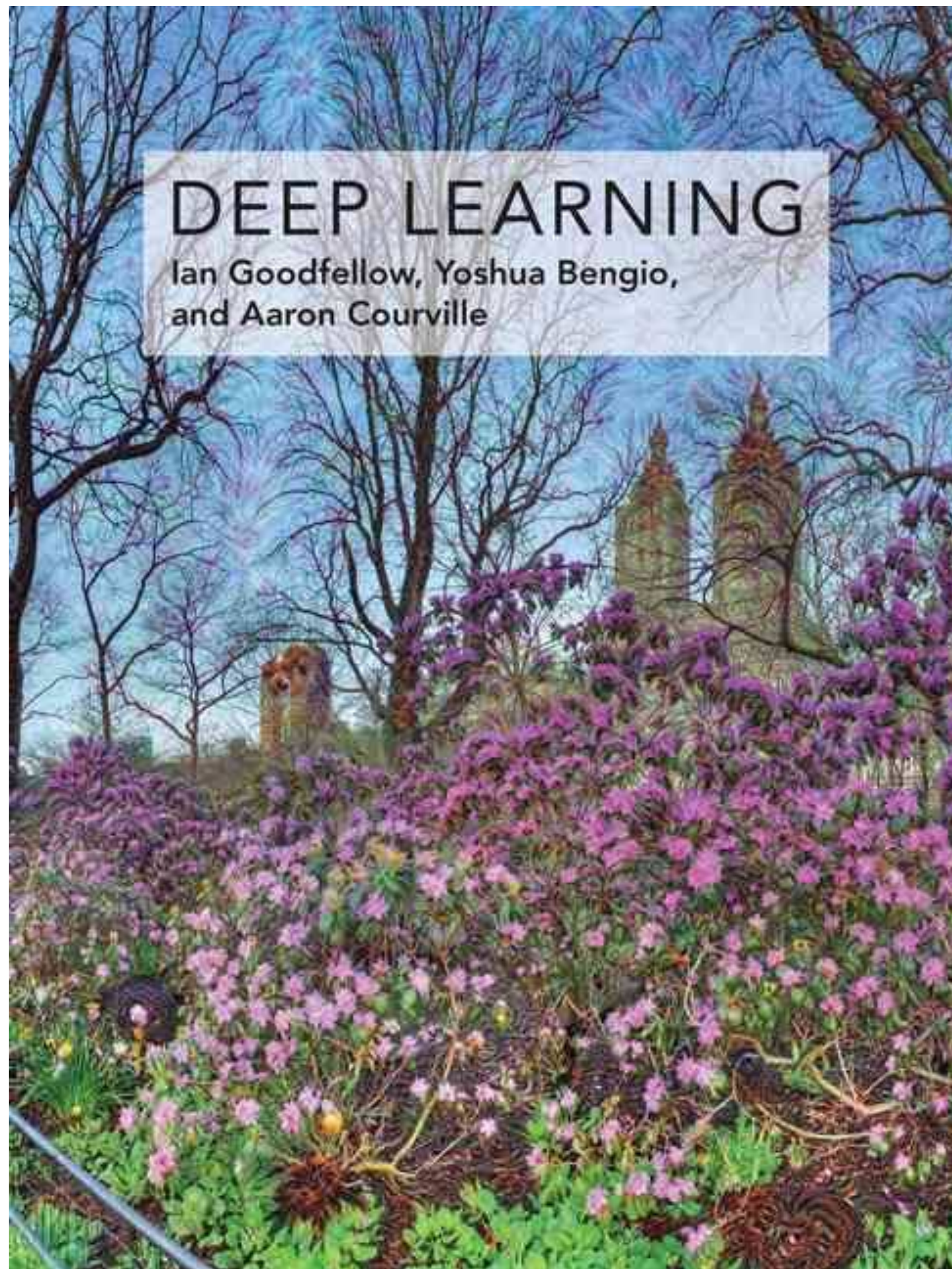




Visual representation learning

Bill Freeman, Antonio Torralba, Phillip Isola
6.819 / 6.869



<http://www.deeplearningbook.org/>

By Ian Goodfellow, Yoshua Bengio and Aaron Courville

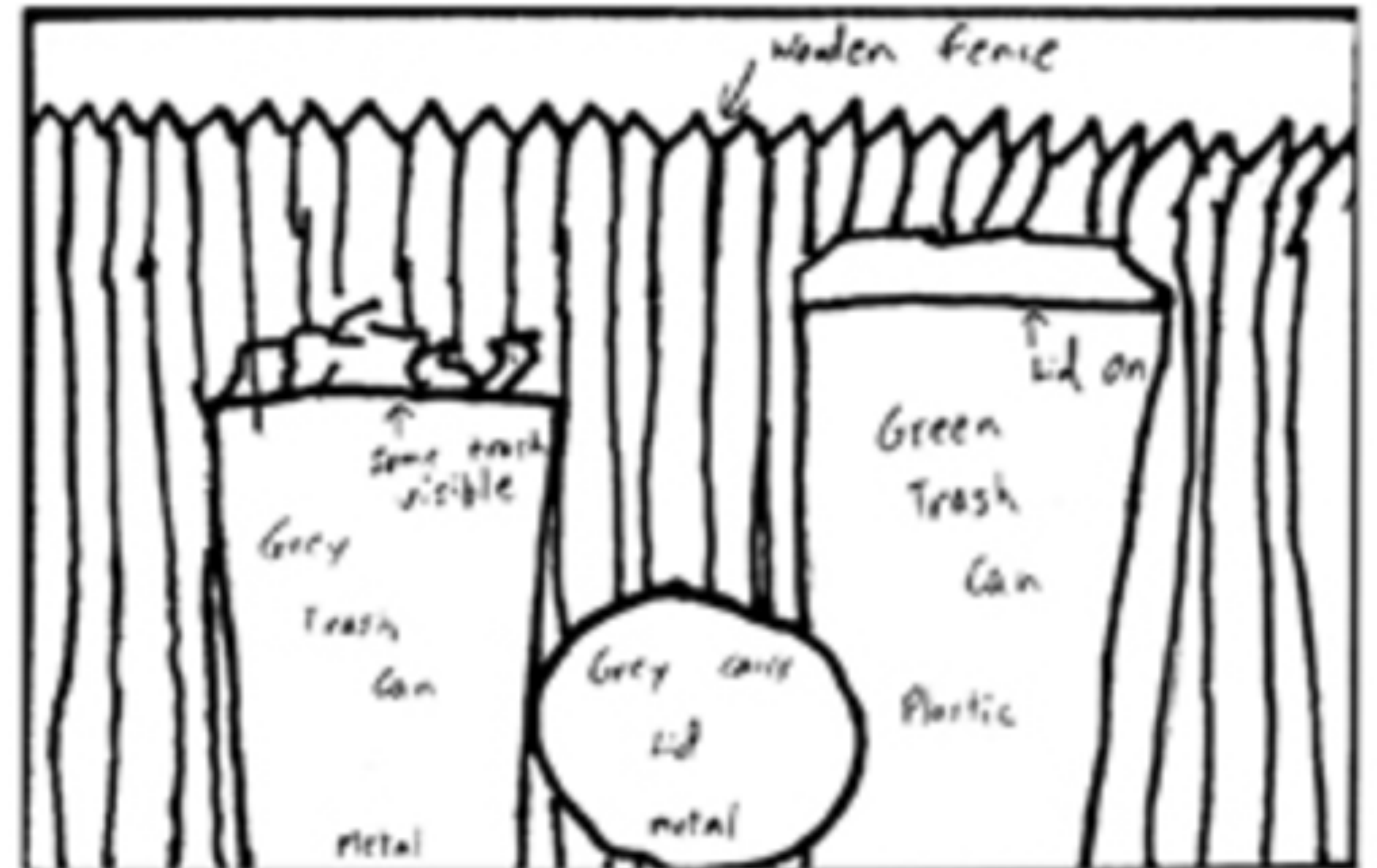
November 2016

Today: parts of chapters 14 and 15 (but this lecture is mostly a departure from the book)

Observed image



Drawn from memory



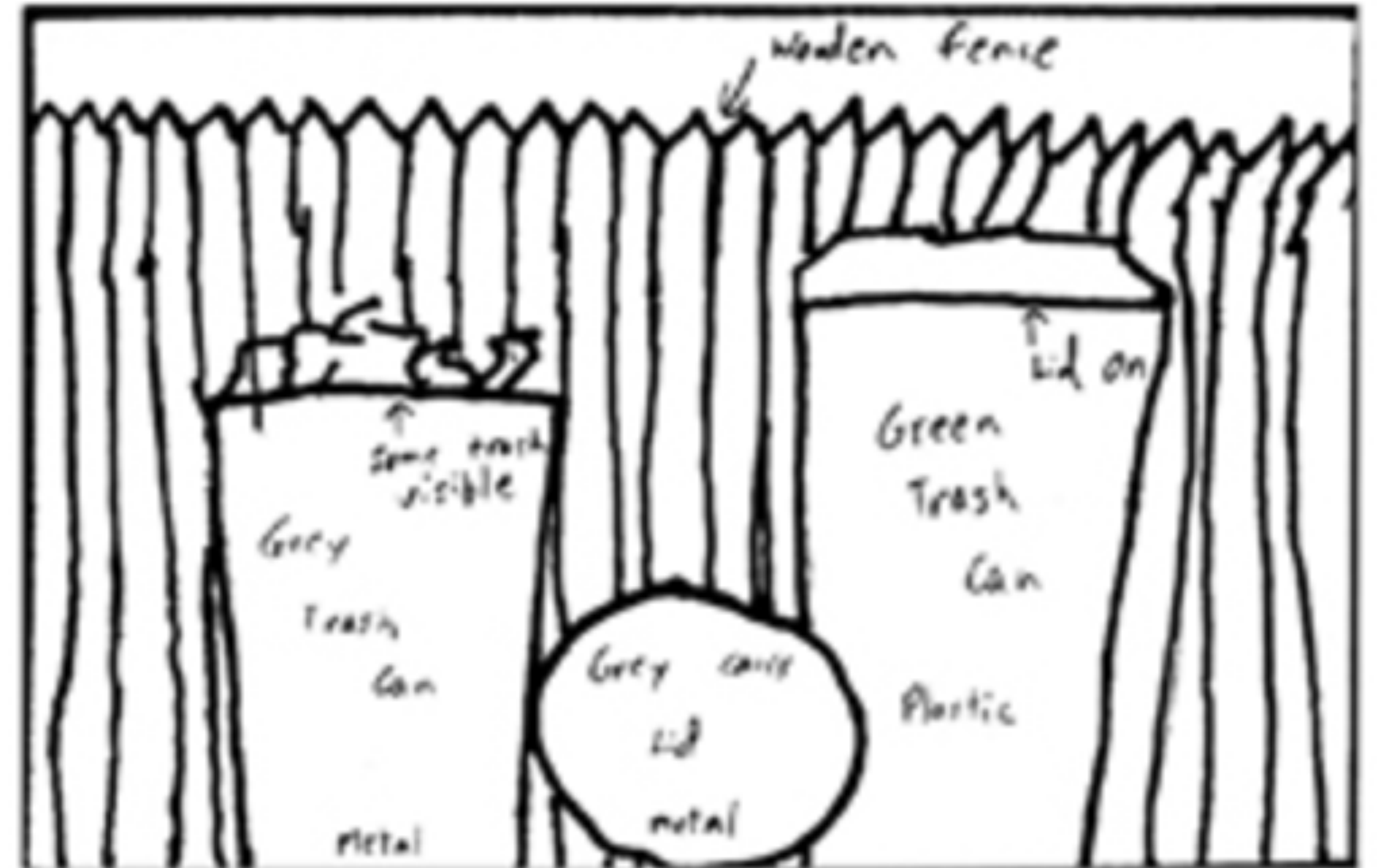
[Bartlett, 1932]

[Intraub & Richardson, 1989]

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."

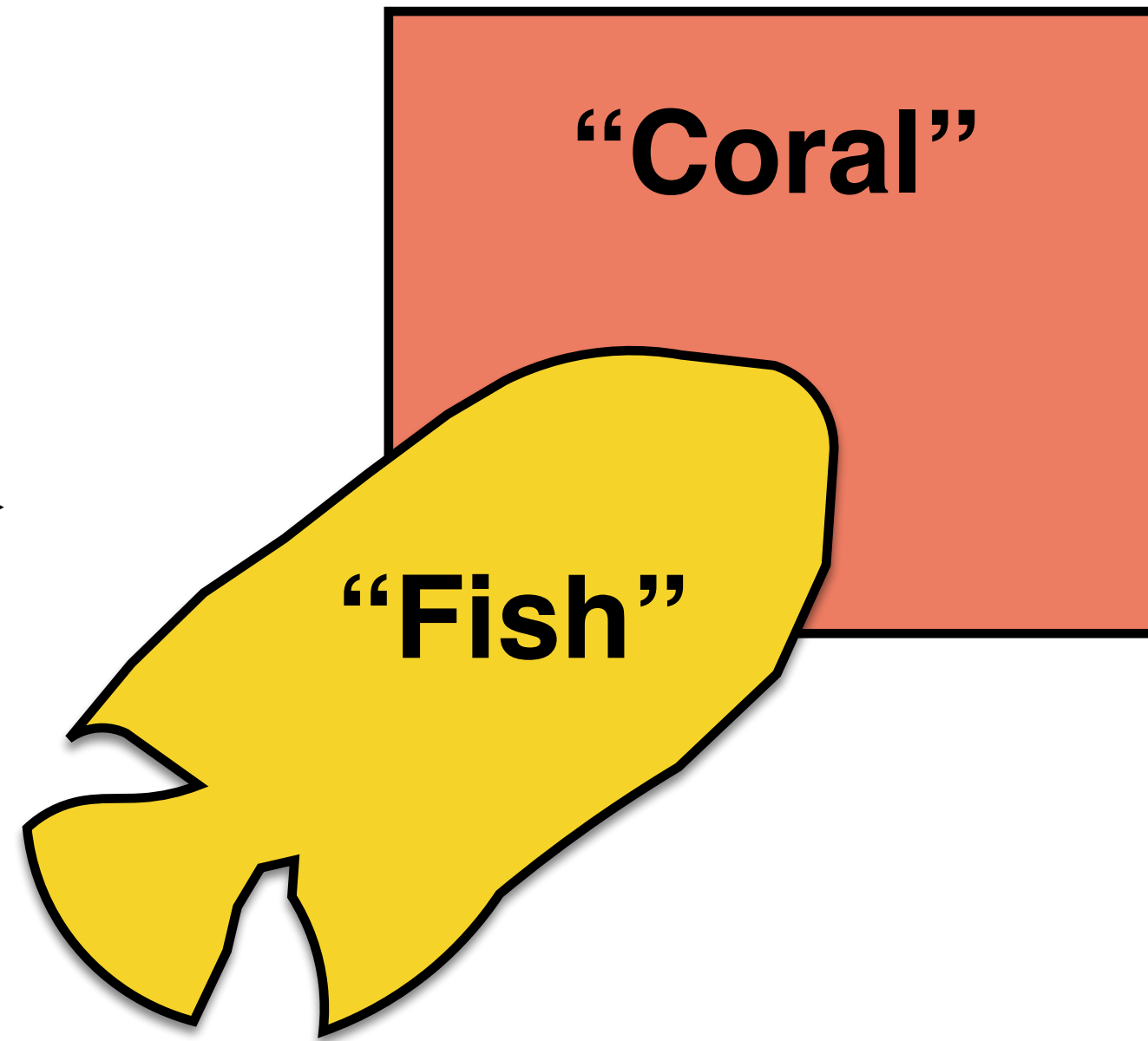
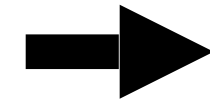
— Max Wertheimer, 1923

Representation learning

X

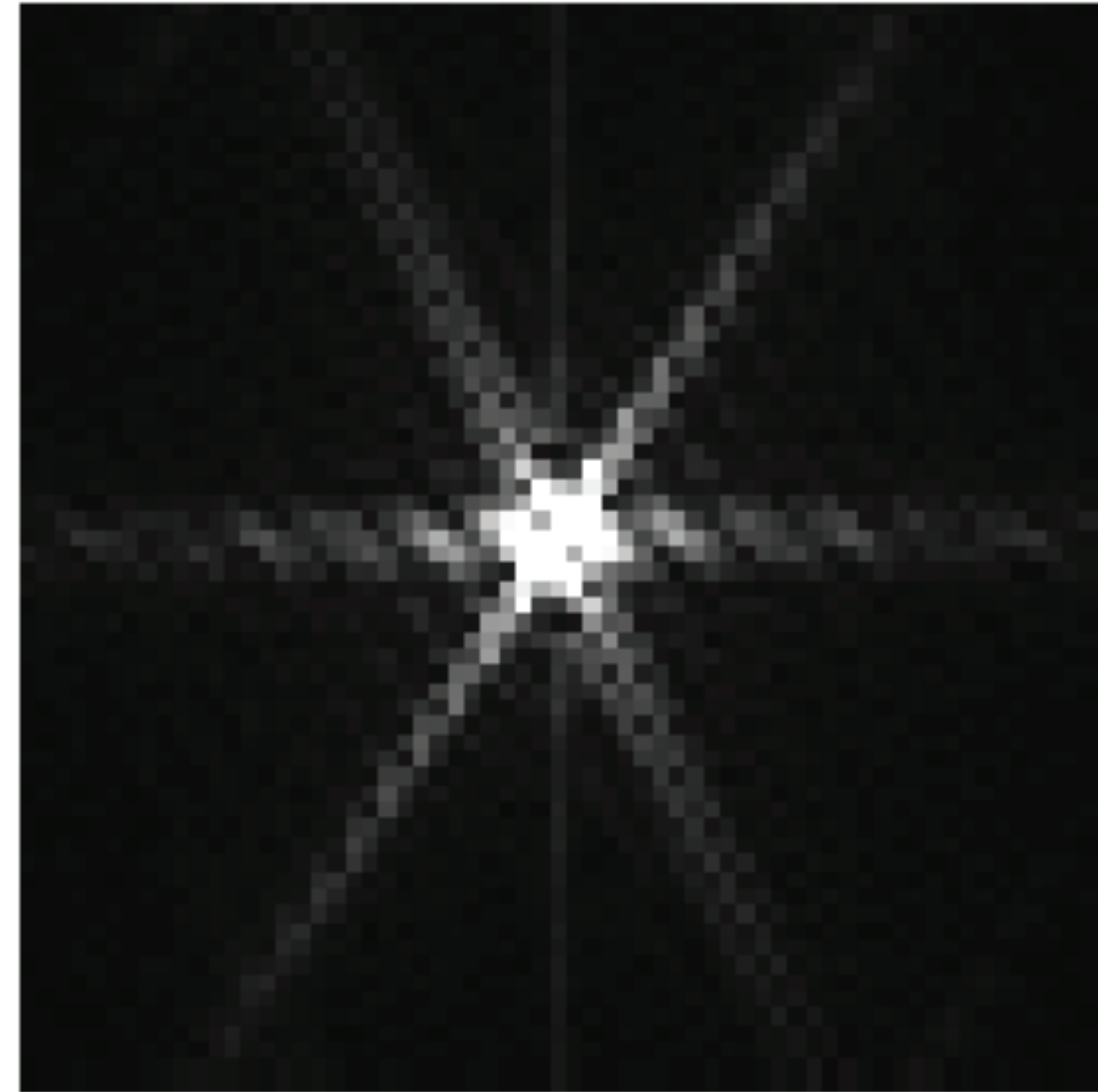


Image



Compact mental
representation

Representation learning

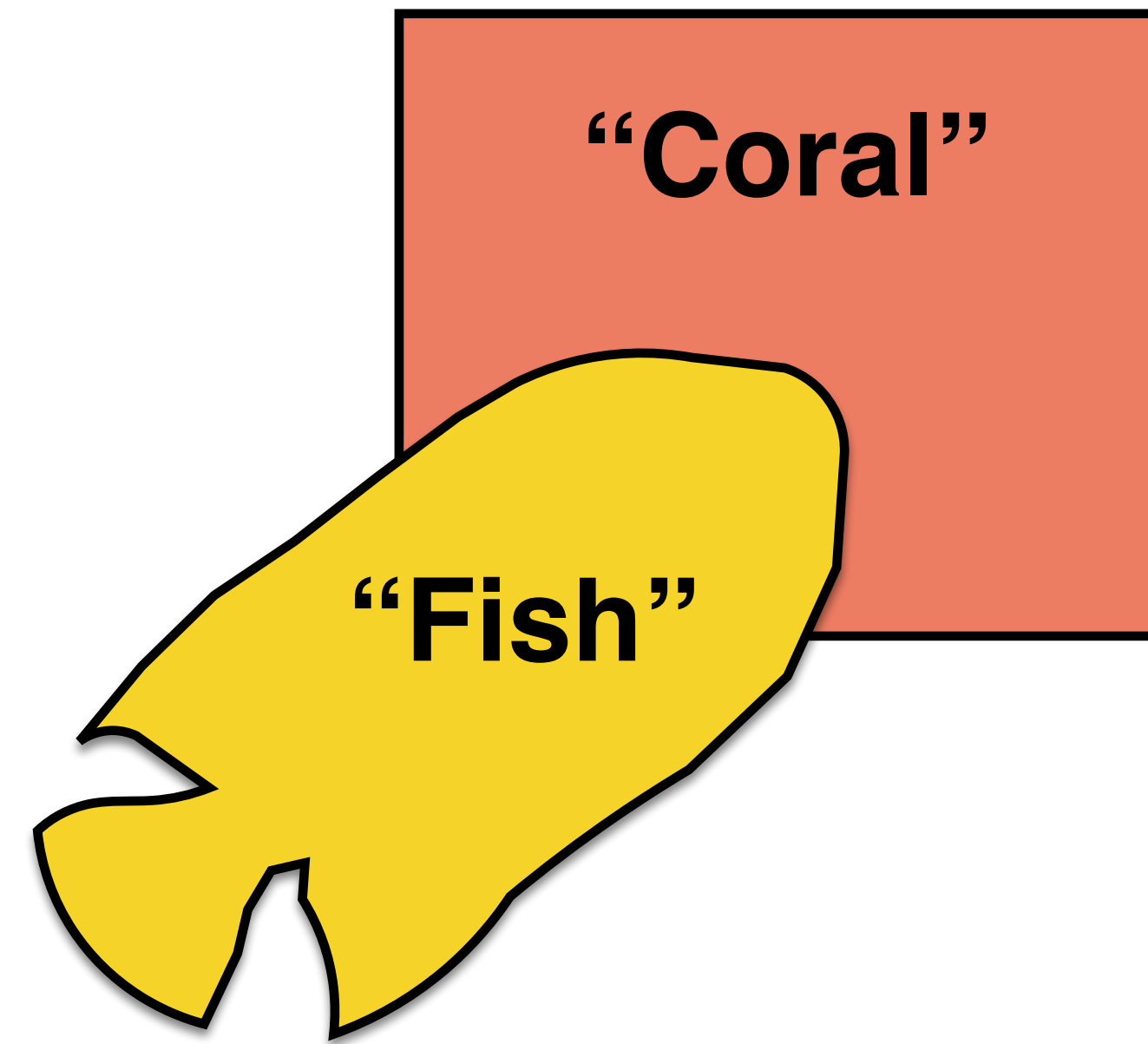


Convolution is pointwise multiplication in the frequency domain.

Representation learning

Good representations are:

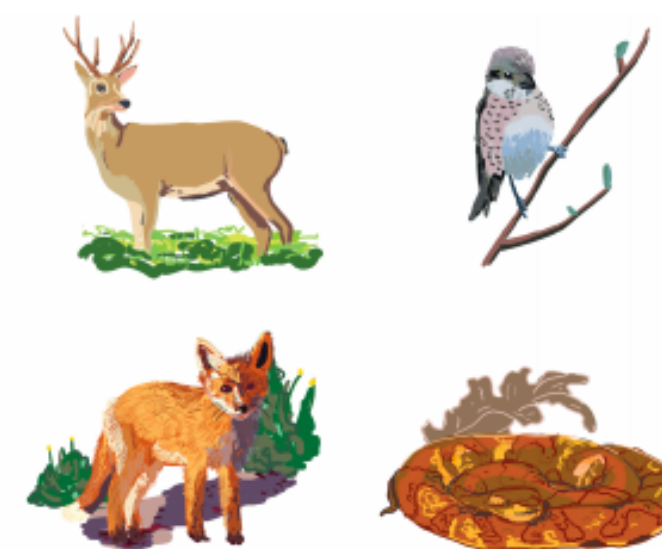
1. Compact
2. Explanatory
3. Disentangled
4. Interpretable



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



Classification units



PIT/AIT



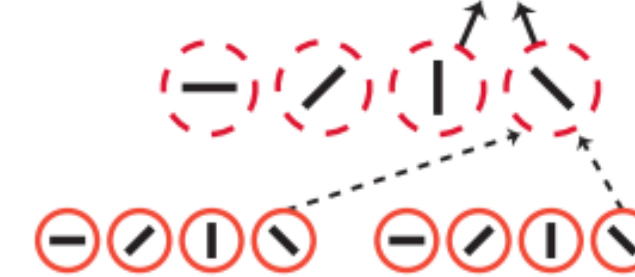
V4/PIT



V2/V4

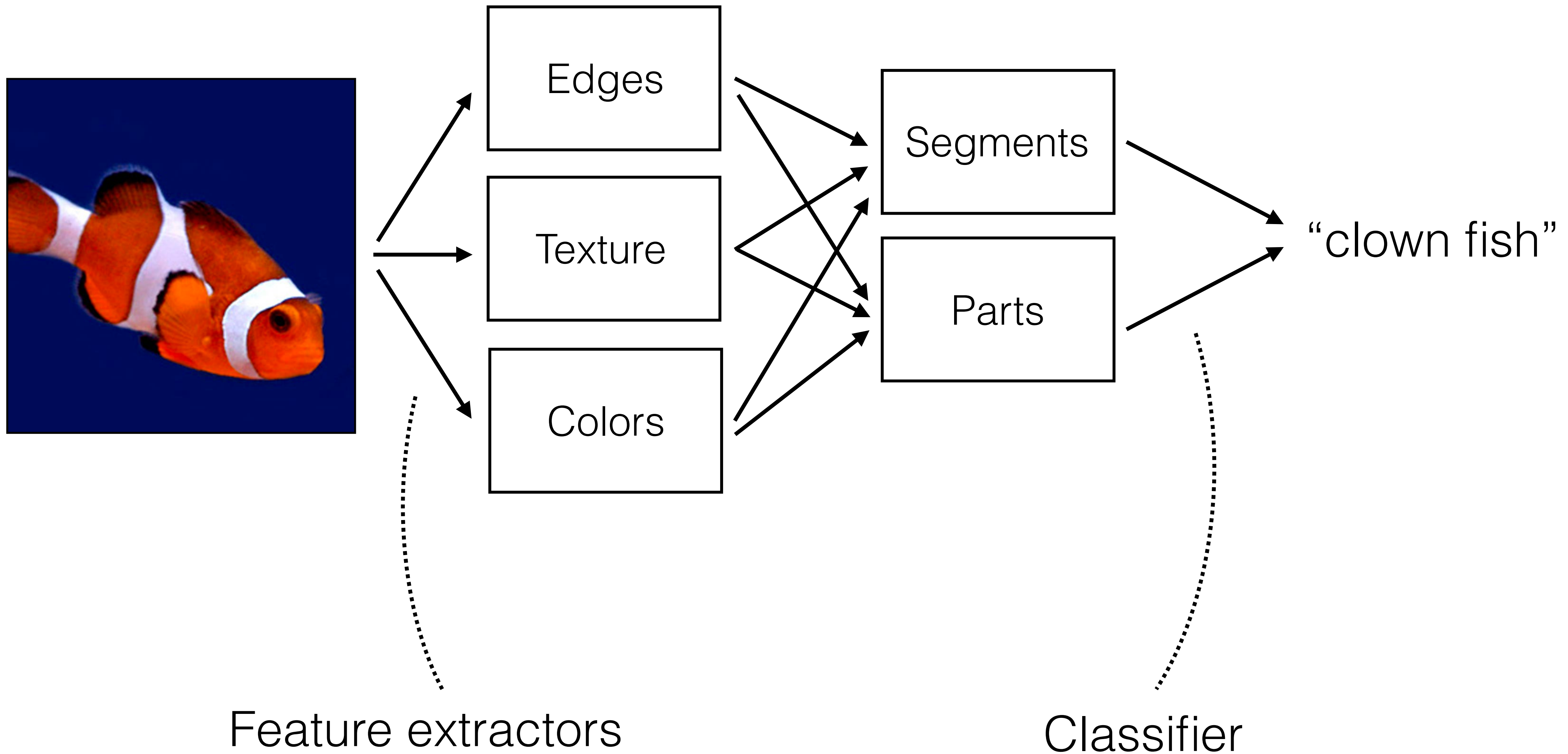


V1/V2

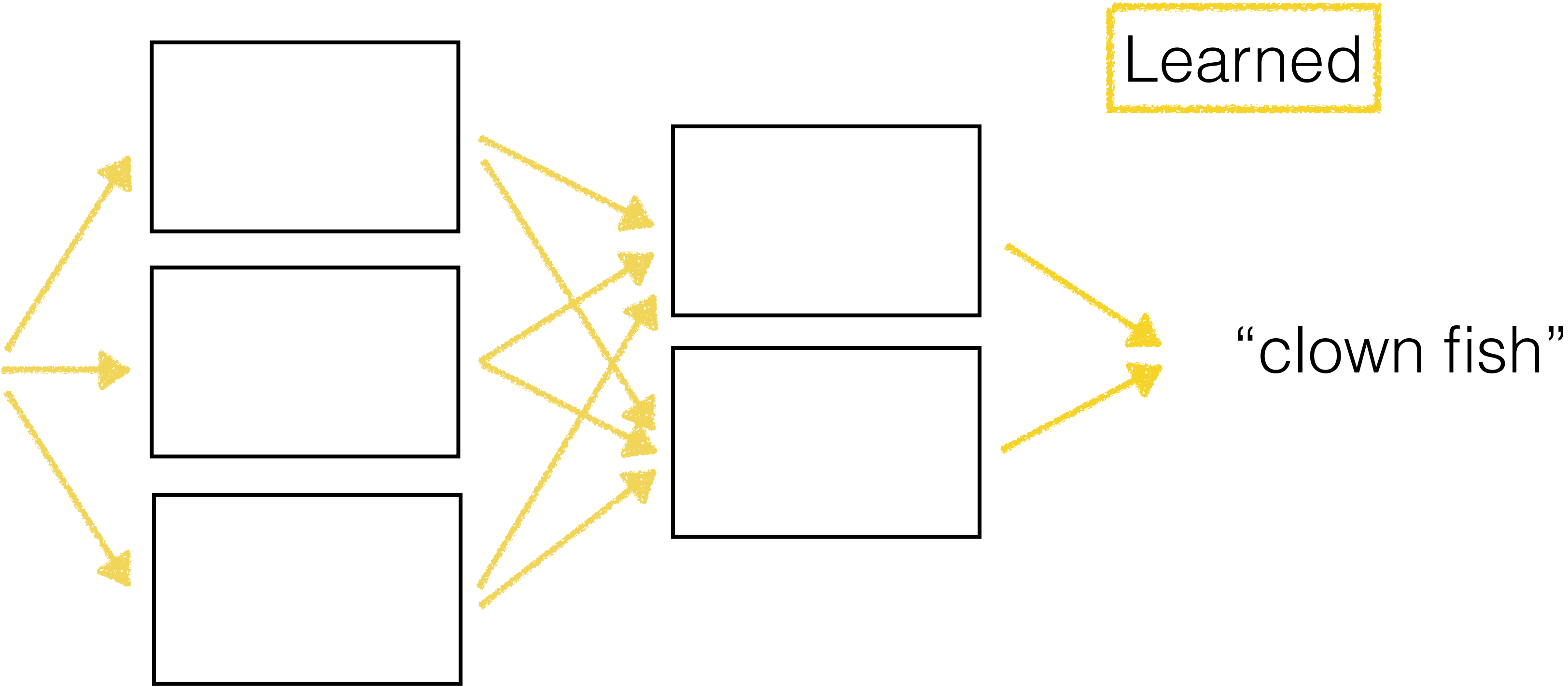


[Serre, 2014]

Classical object recognition



Deep learning

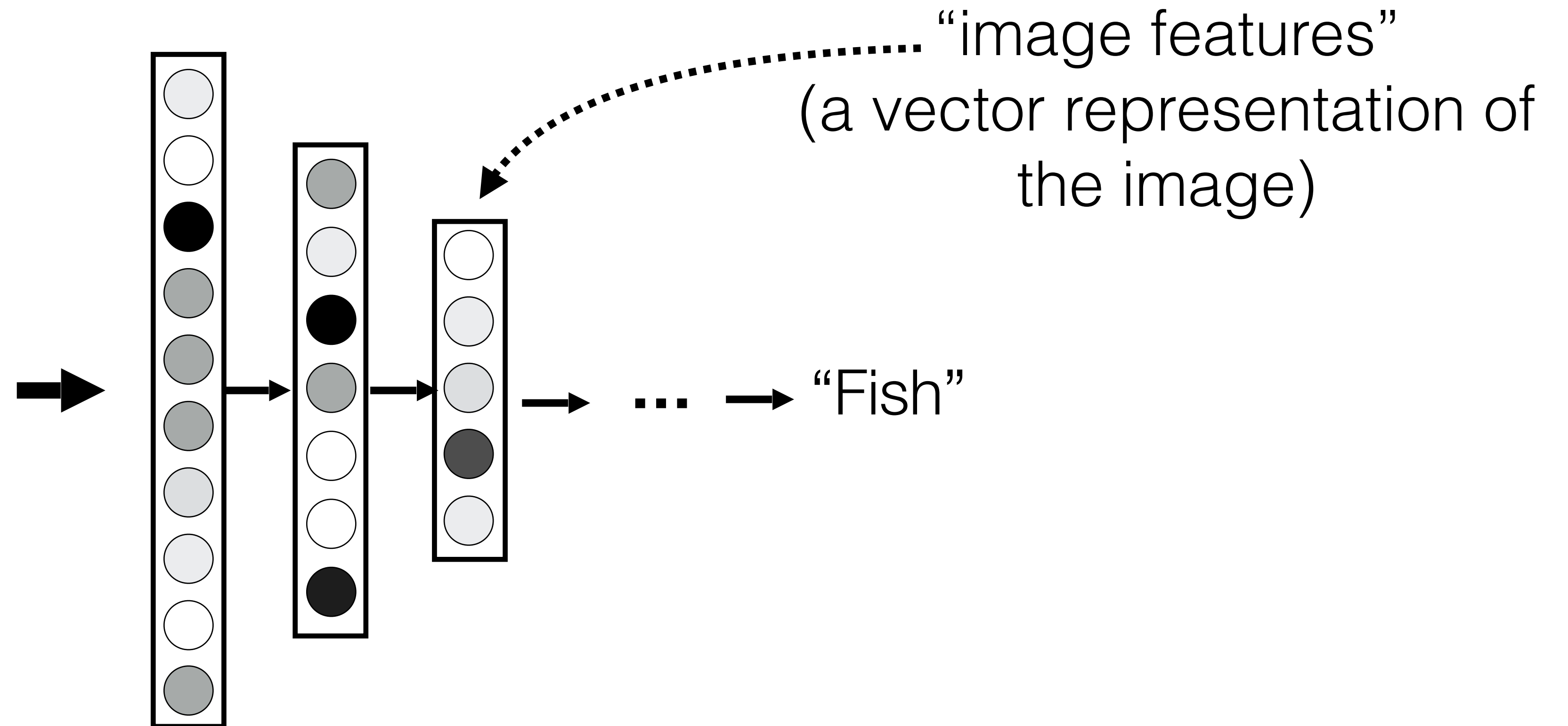


What do deep nets internally learn?

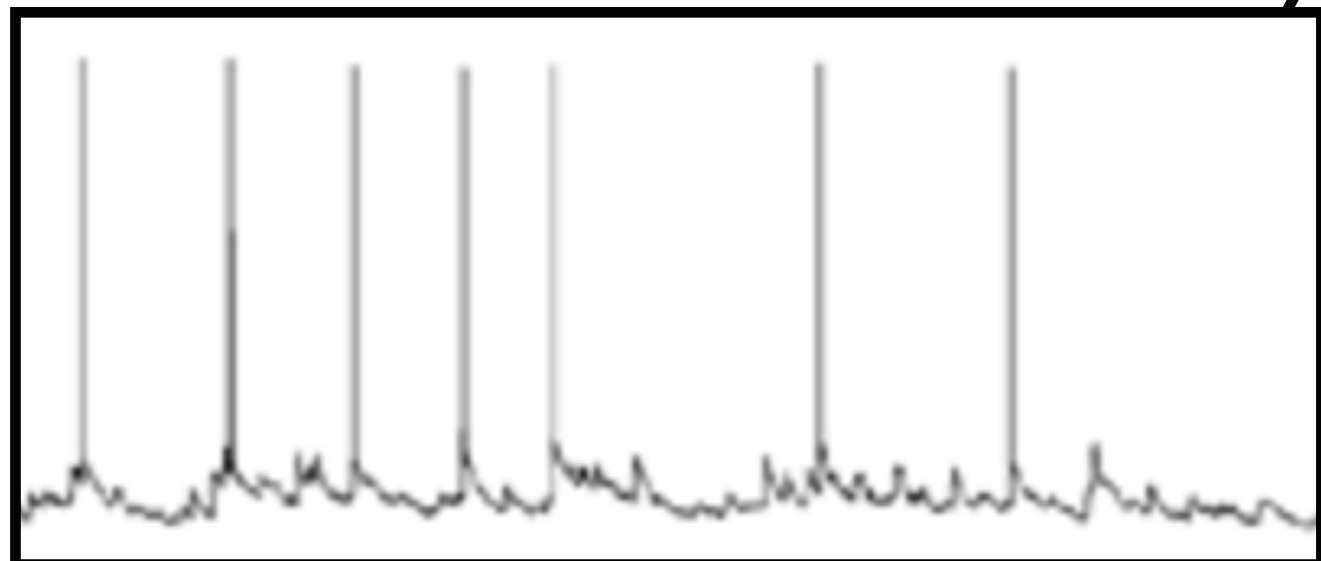
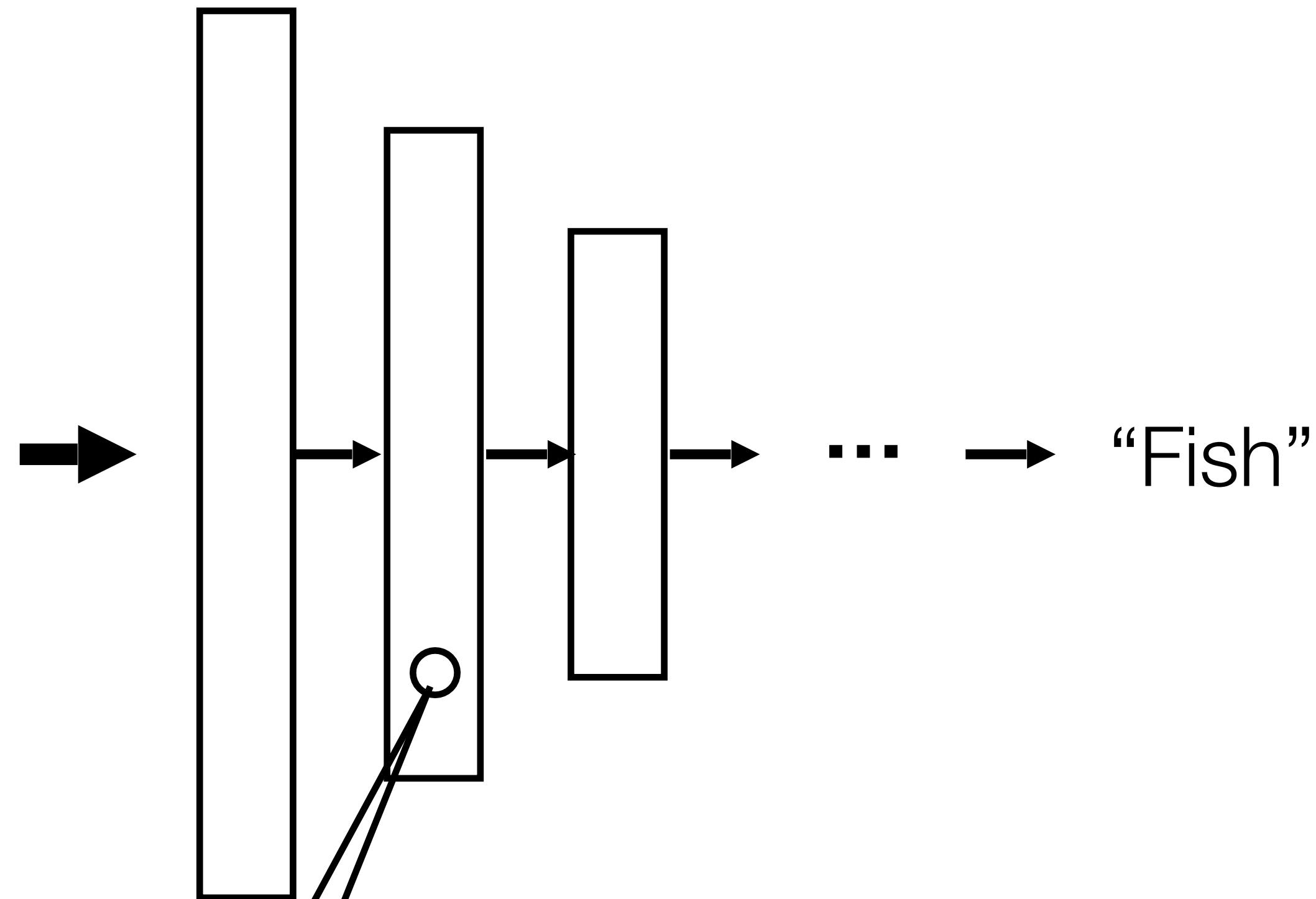
X



Image



Deep Net “Electrophysiology”



[Zeiler & Fergus, ECCV 2014]

[Zhou et al., ICLR 2015]

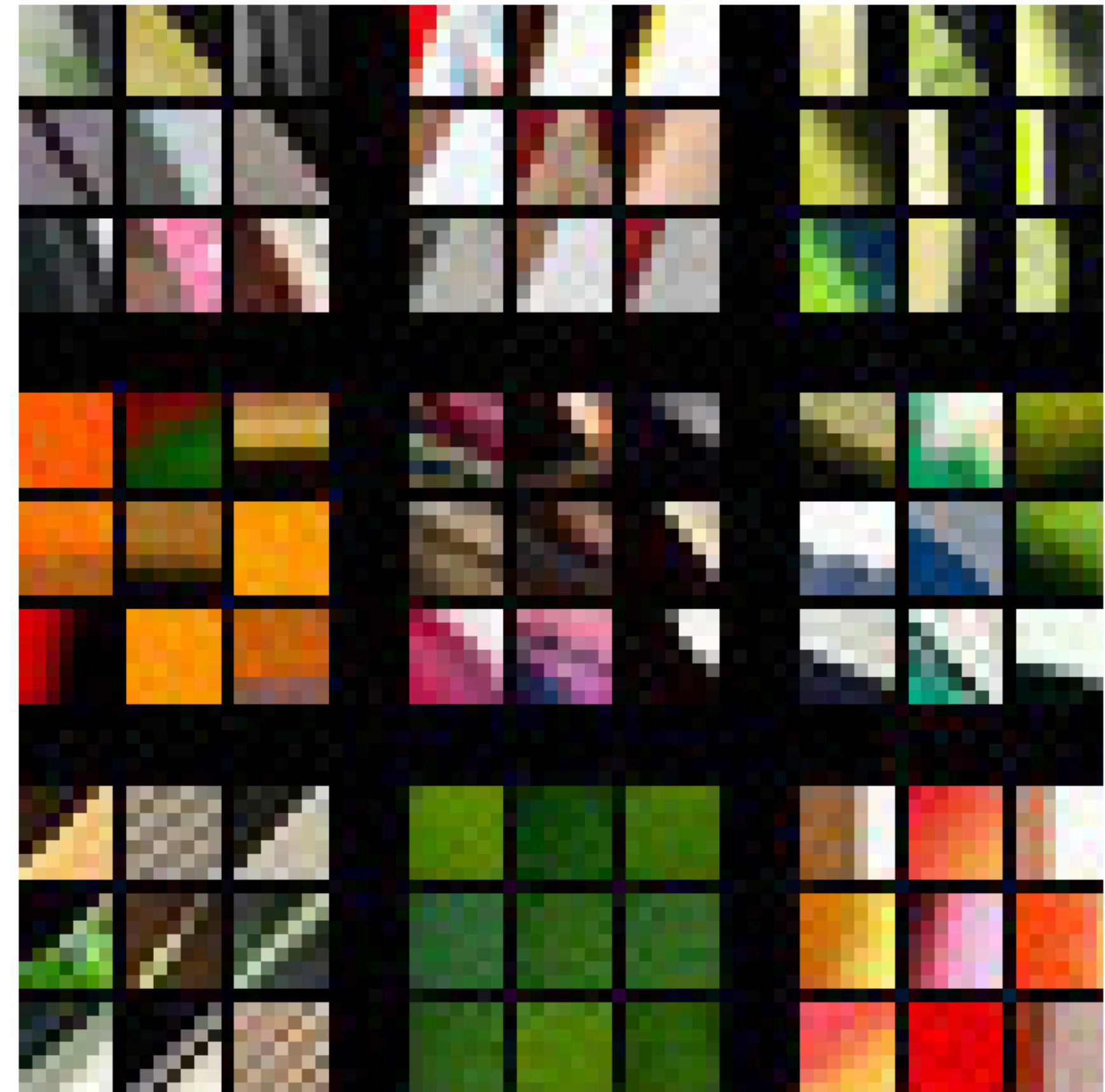
Visualizing and Understanding CNNs

[Zeiler and Fergus, 2014]

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

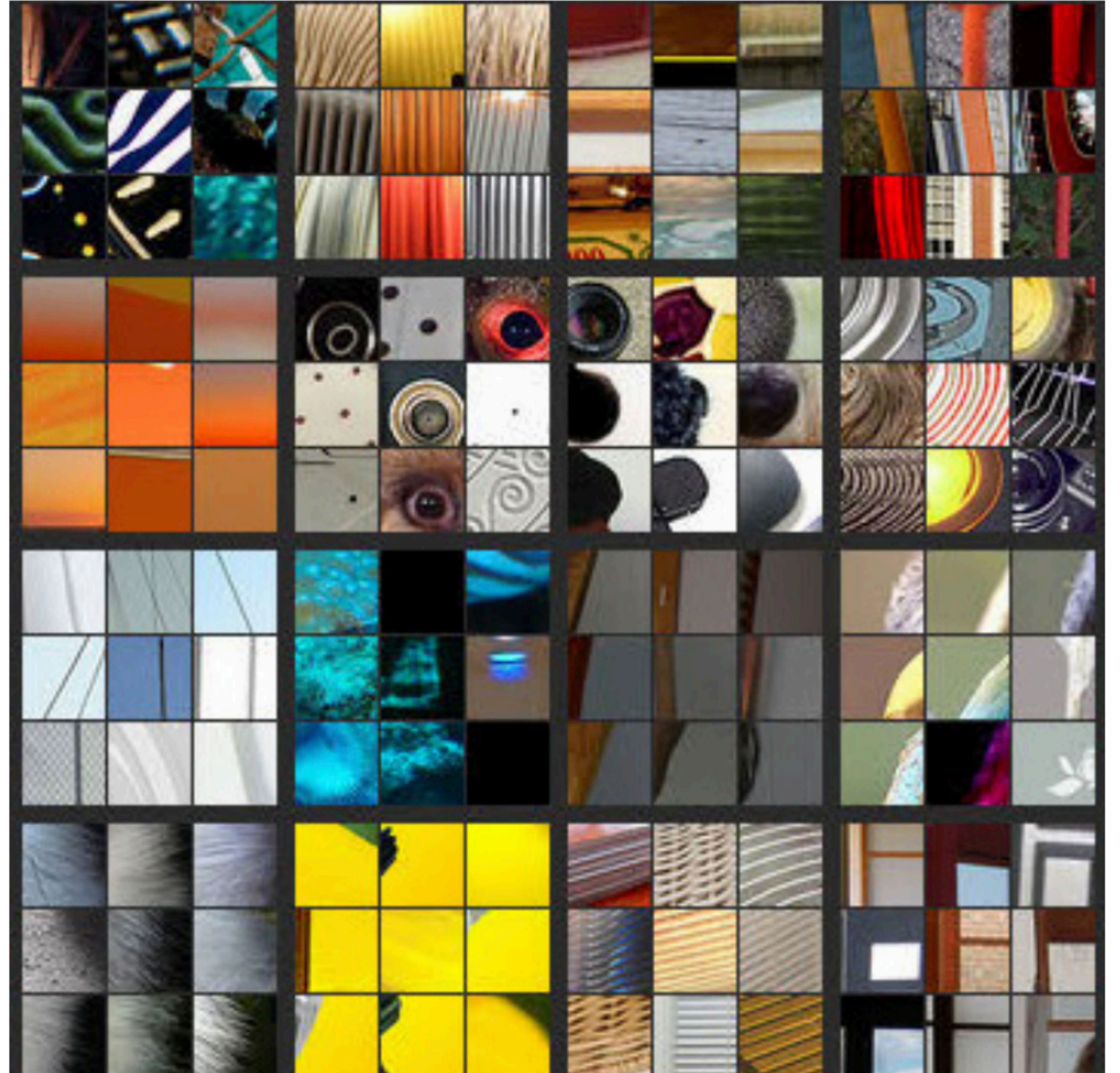


Image patches that activate each of the **layer 1** filters most strongly



[Zeiler and Fergus, 2014]

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



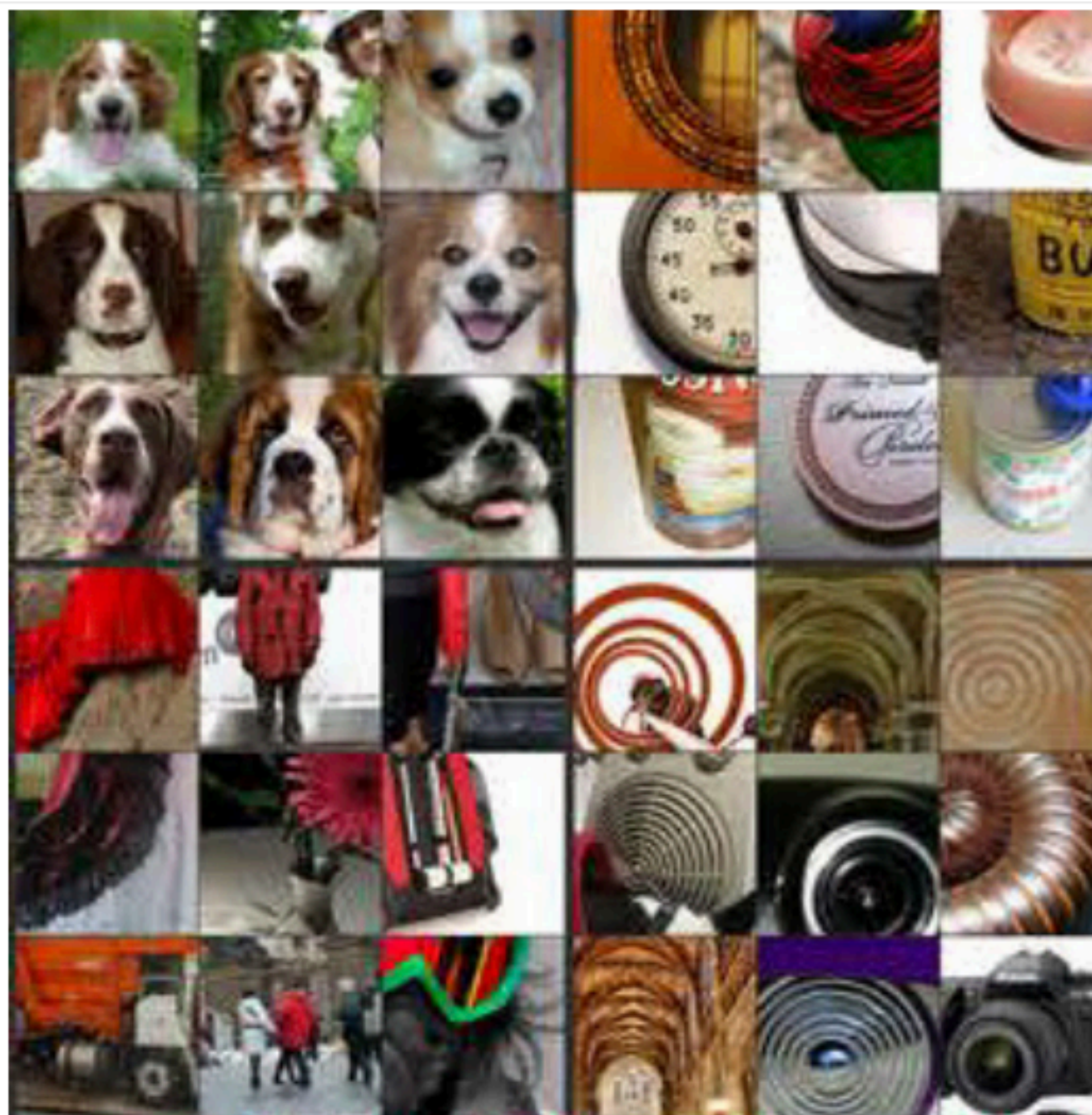
[Zeiler and Fergus, 2014]



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

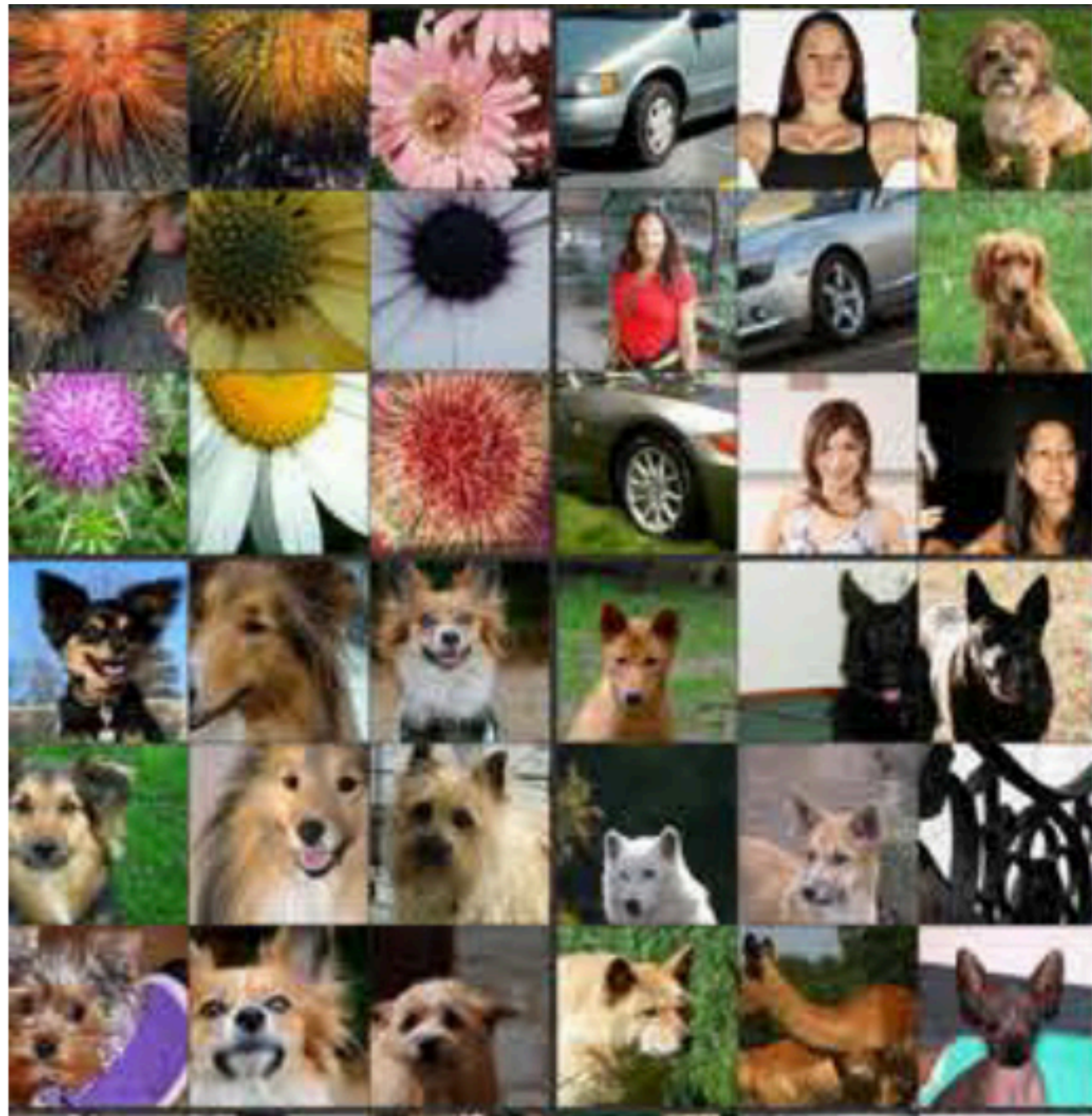
[Zeiler and Fergus, 2014]

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

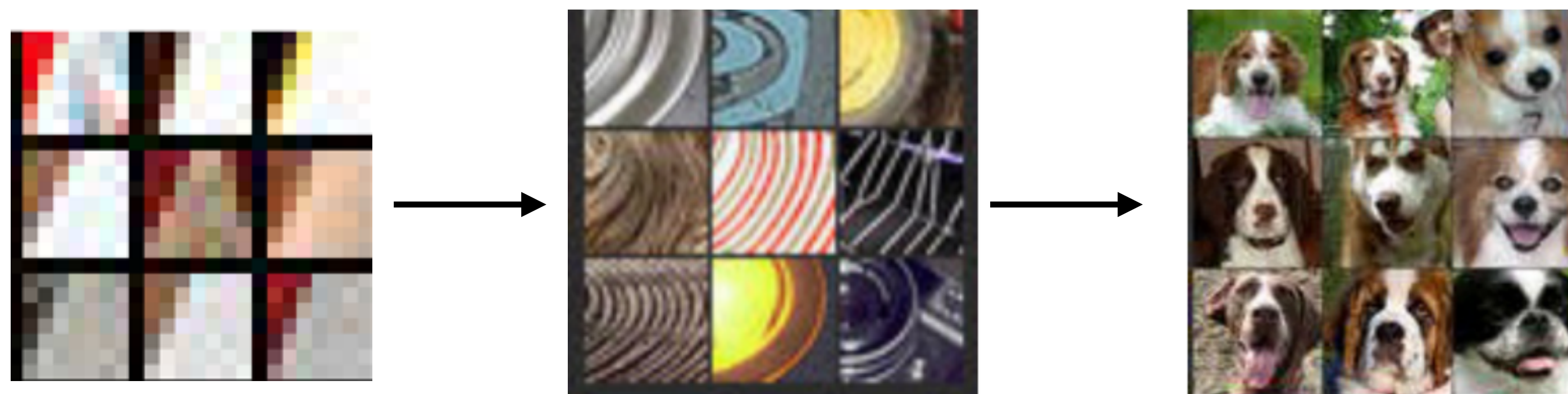
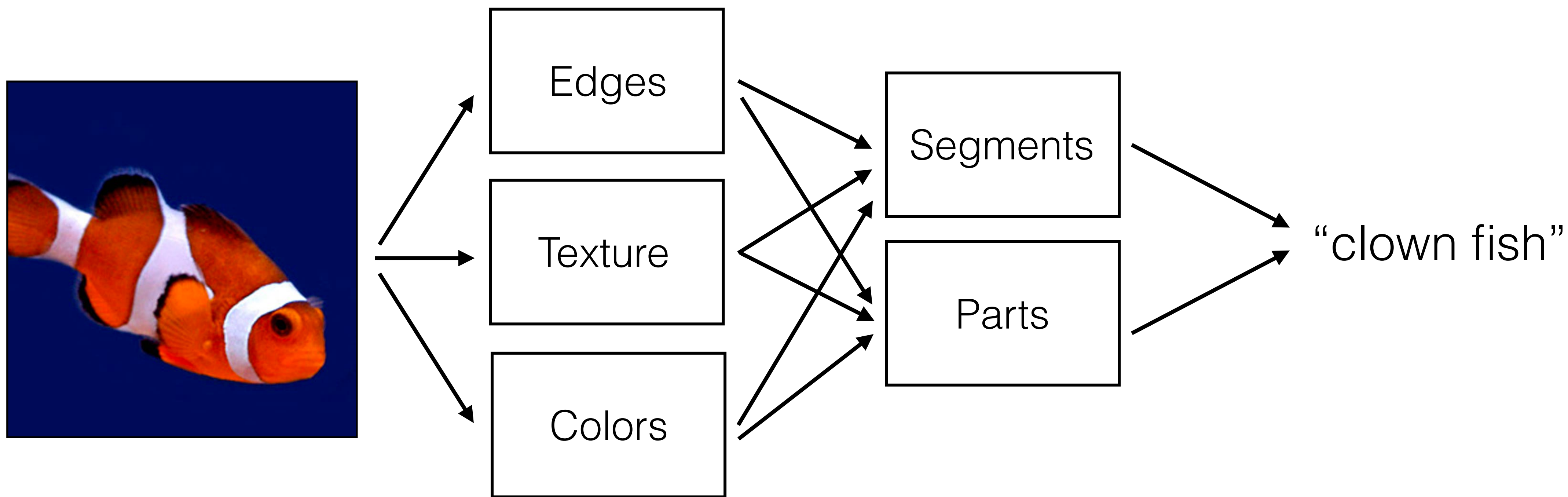


[Zeiler and Fergus, 2014]

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

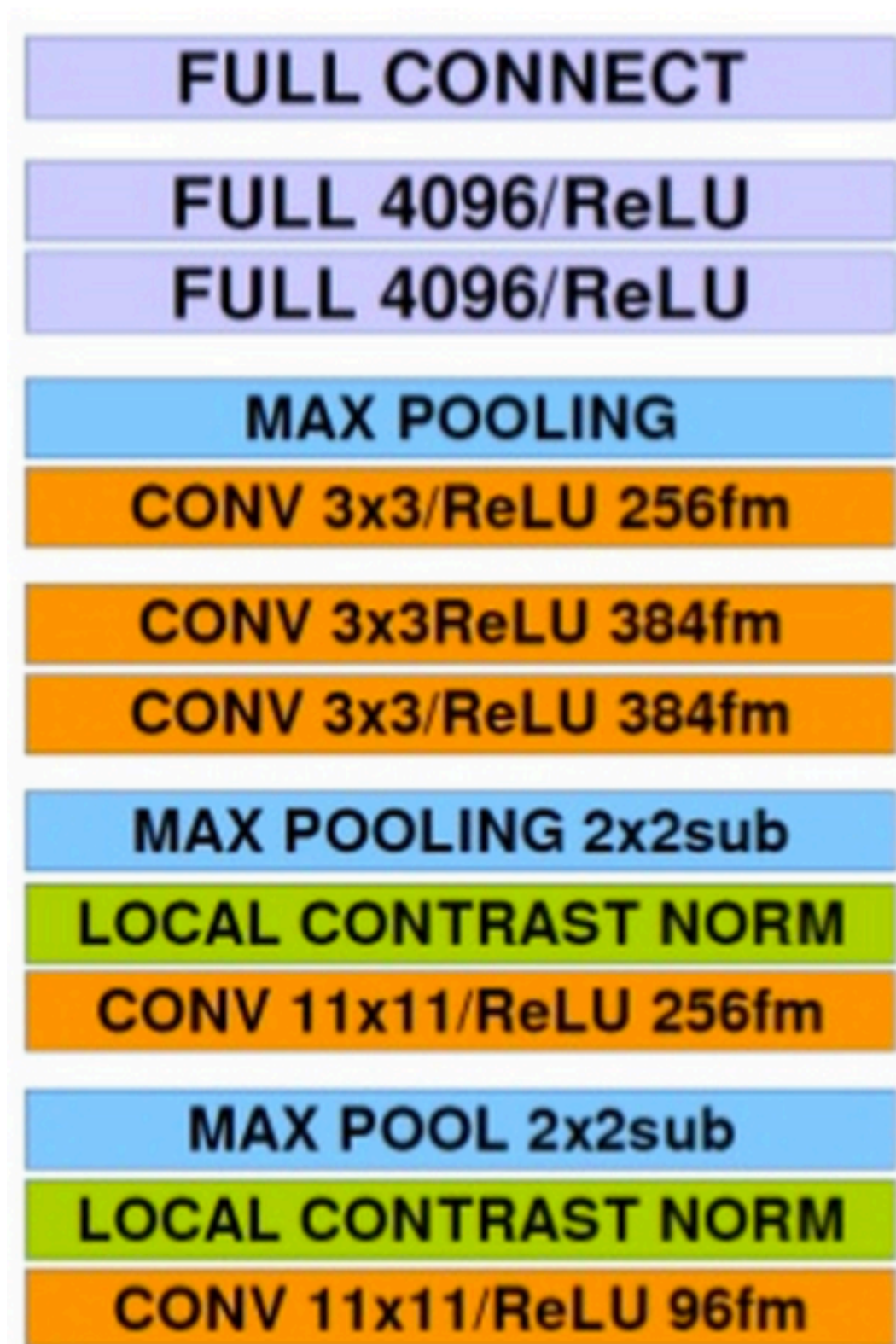


CNNs *learned* the classical visual recognition pipeline!



Object Detectors Emergence in Deep Scene CNNs

[Zhou et al., ICLR 2015]



AlexNet

- For each unit (neuron) in network, find which images it is most selective for (cause it to have highest activation)
- Find which pixels in these images are responsible by occluding regions and seeing which pixels, when occluded, cause activation to change the most

Object Detectors Emergence in Deep Scene CNNs

[Zhou et al., ICLR 2015]

pool 1



Object Detectors Emergence in Deep Scene CNNs

[Zhou et al., ICLR 2015]

pool 2



Object Detectors Emergence in Deep Scene CNNs

[Zhou et al., ICLR 2015]

conv 4



Object Detectors Emergence in Deep Scene CNNs

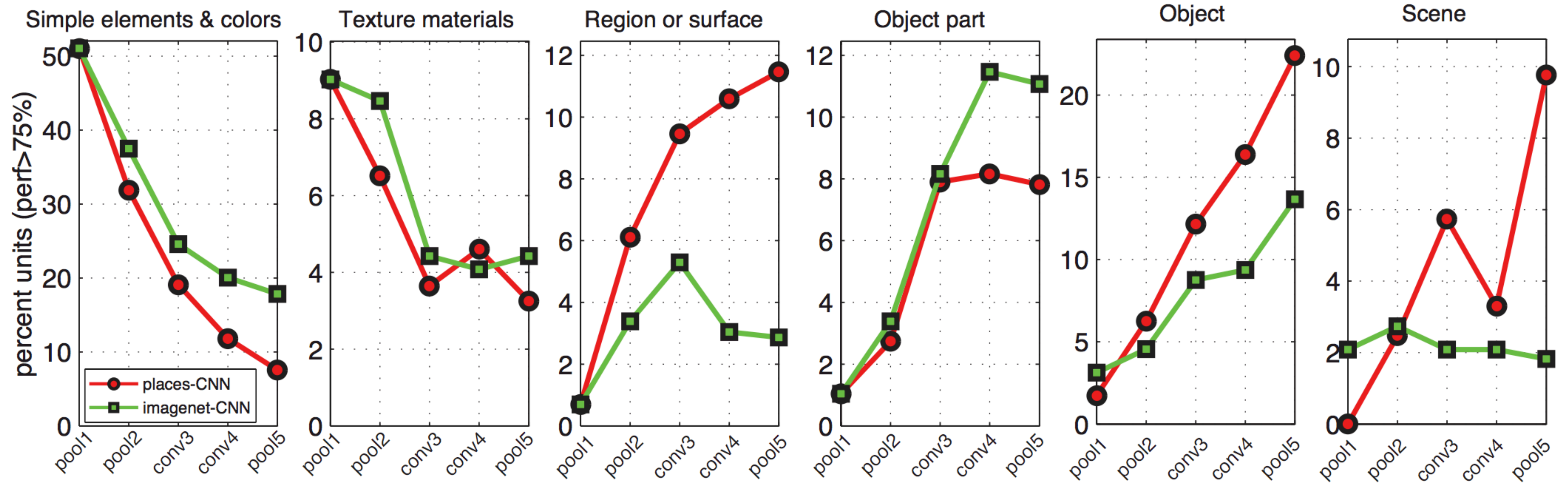
[Zhou et al., ICLR 2015]

pool 5



Object Detectors Emergence in Deep Scene CNNs

[Zhou et al., ICLR 2015]

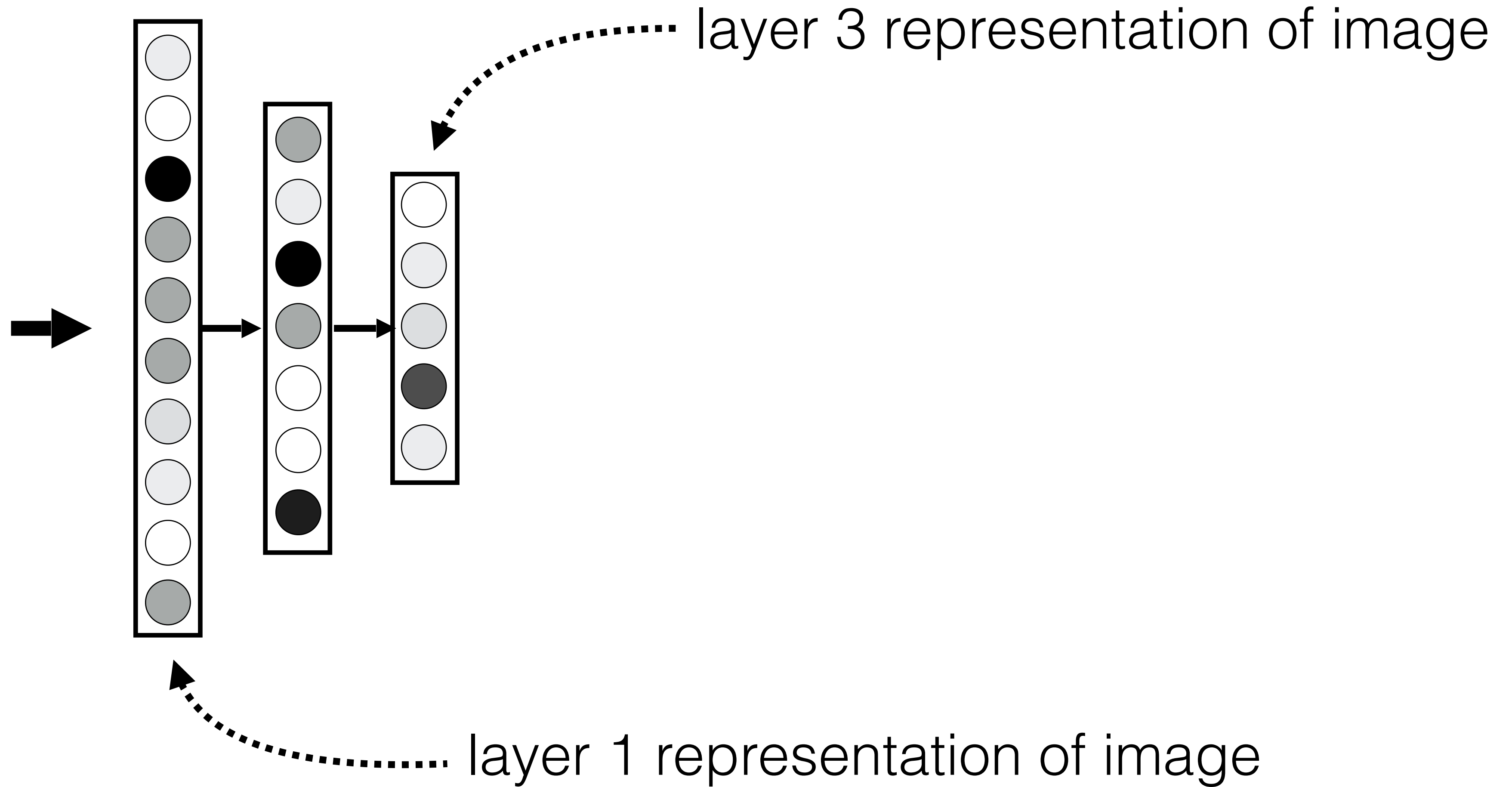


im2vec

X



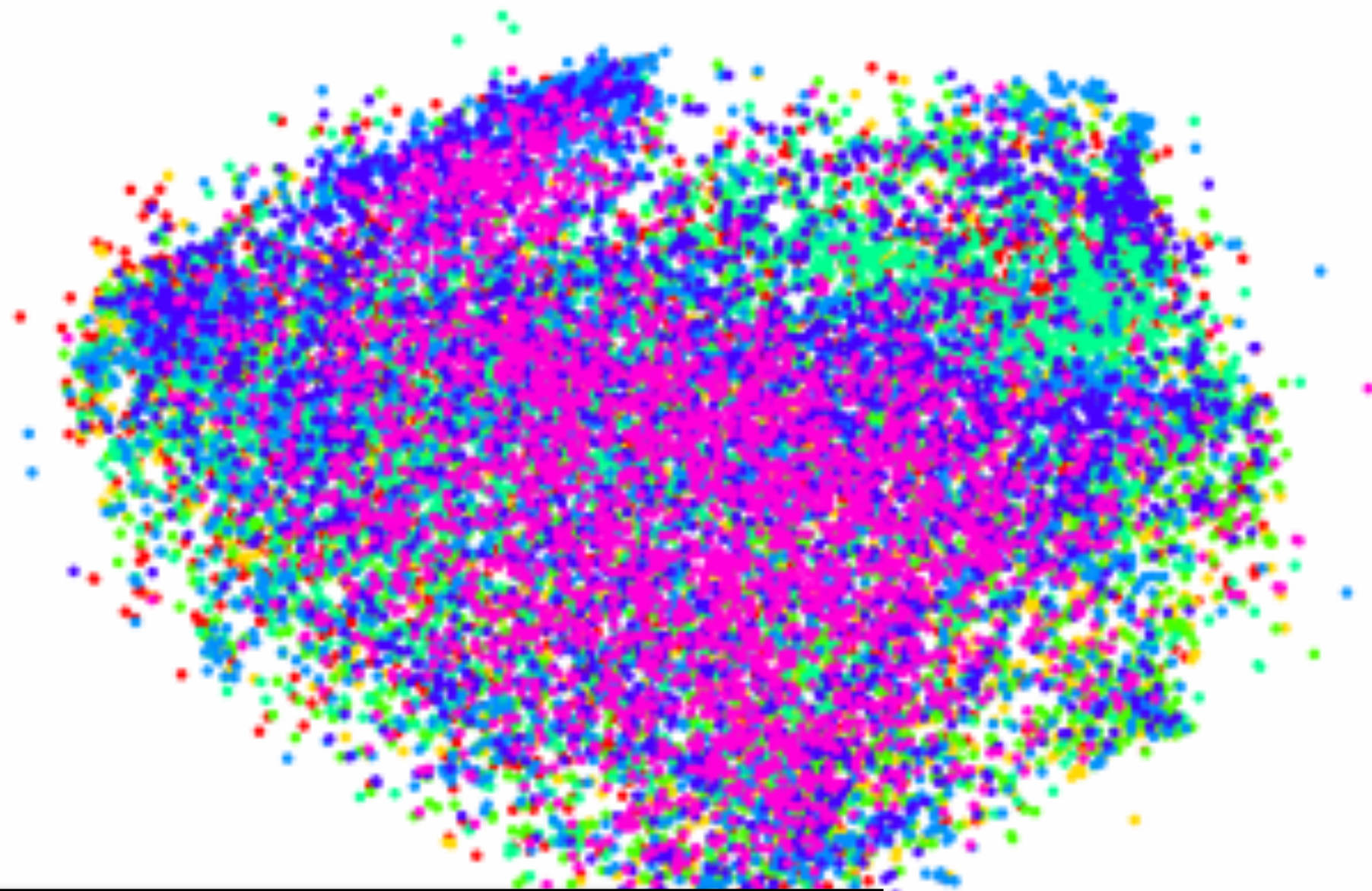
Image



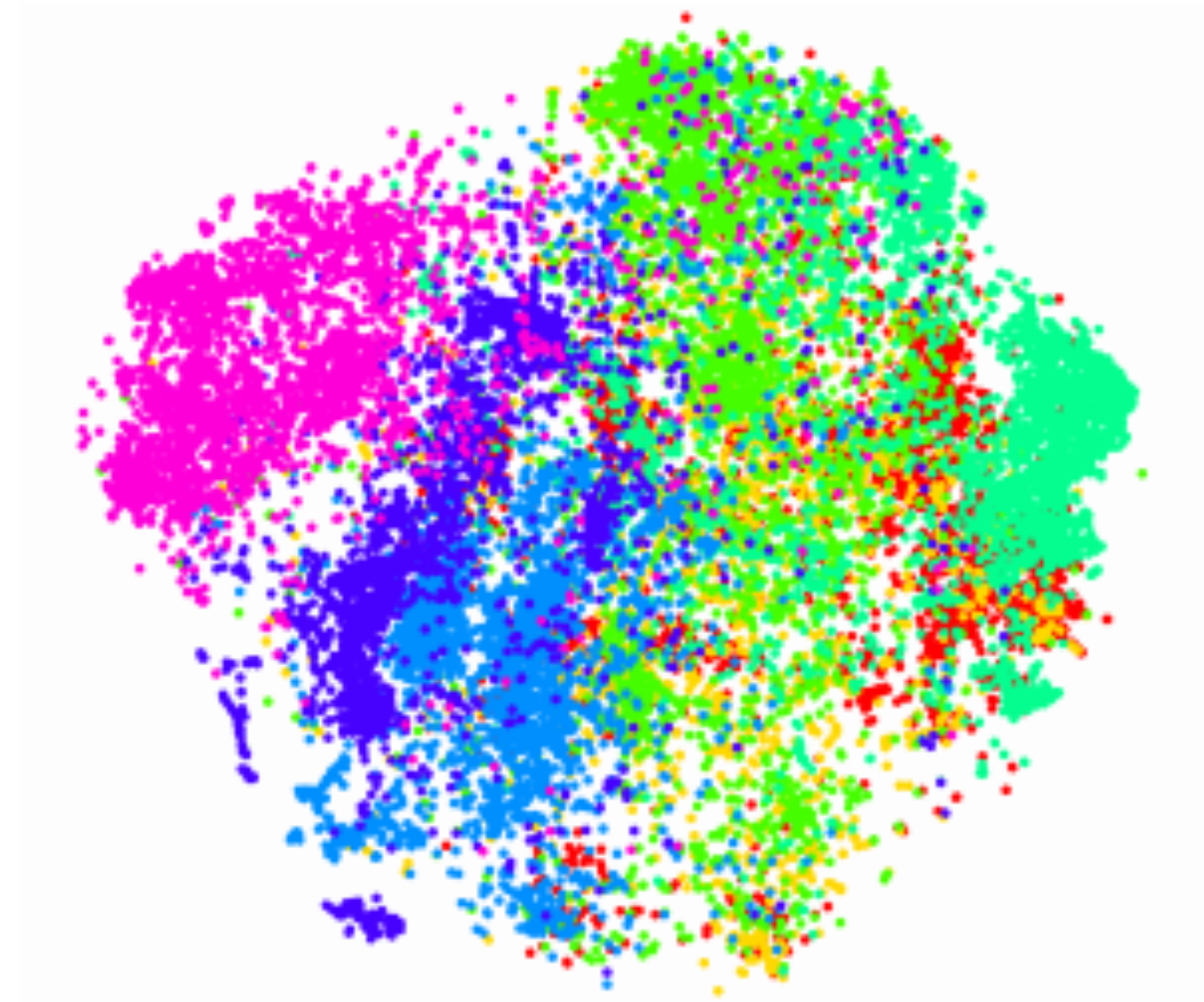
Represent image as a vector of neural activations
(perhaps representing a vector of detected texture patterns or object parts)

Disentangling

Layer 1 representation



Layer 6 representation



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

[DeCAF, Donahue, Jia, et al. 2013]

[Visualization technique : t-sne, van der Maaten & Hinton, 2008]

Investigating a representation via similarity analysis

How similar are these two images?

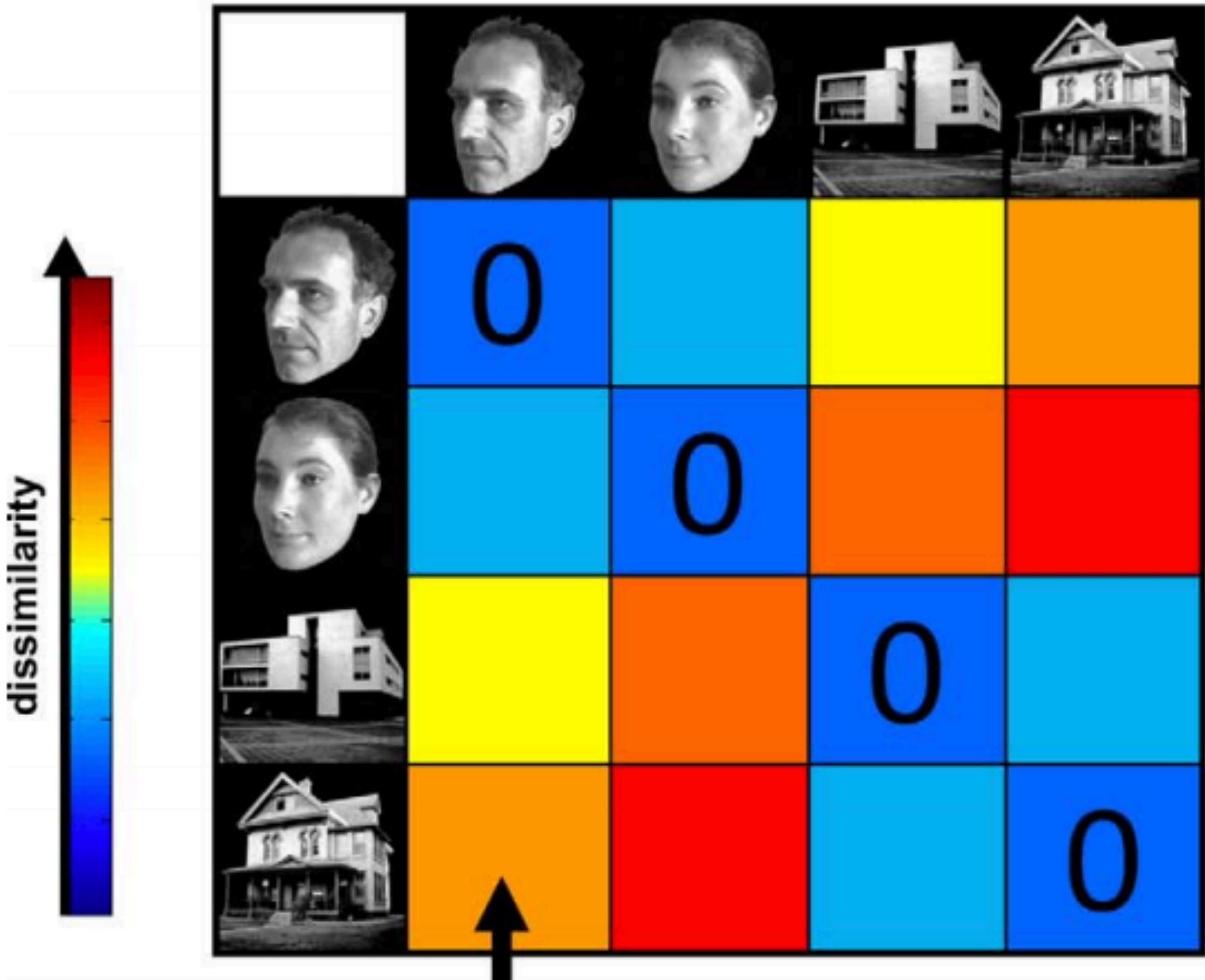


How about these two?



Investigating a representation via similarity analysis

Representational Dissimilarity Matrix



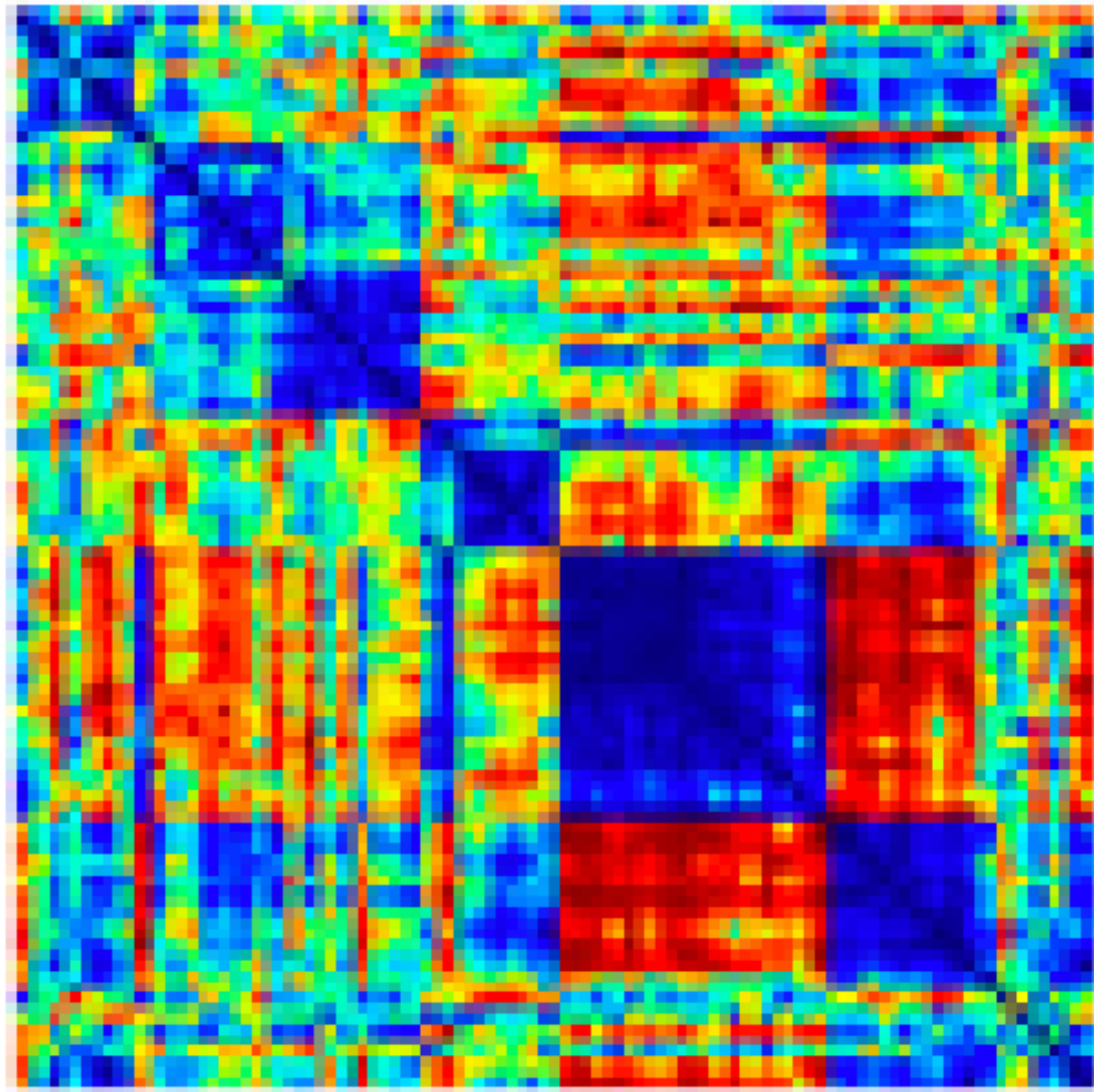
$$\|\mathbf{h}_i - \mathbf{h}_j\|$$

Neural activation vector

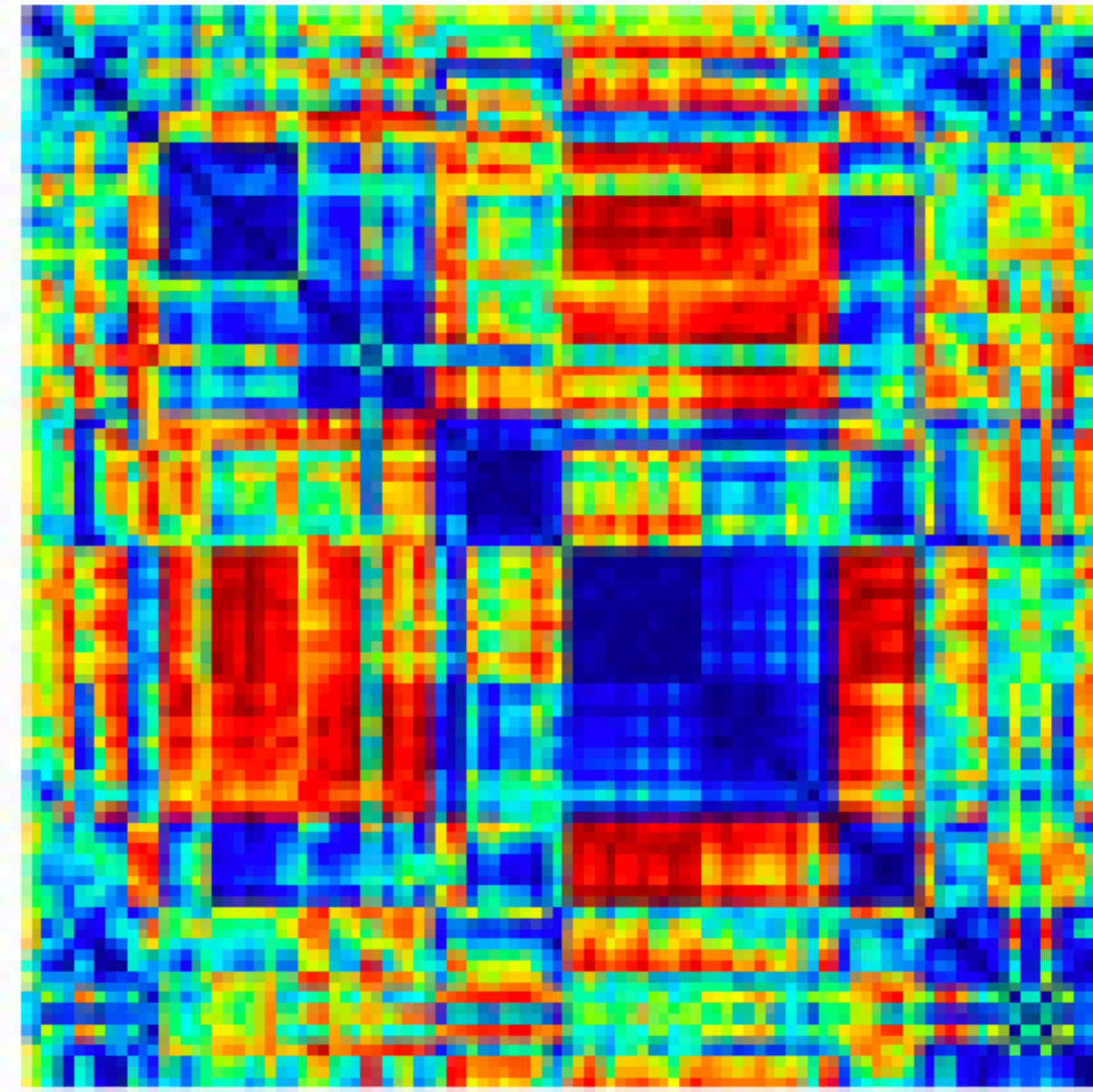
[Kriegeskorte et al. 2008]

Investigating a representation via similarity analysis

IT Neuronal Units



Deep net (in particular, HMO)

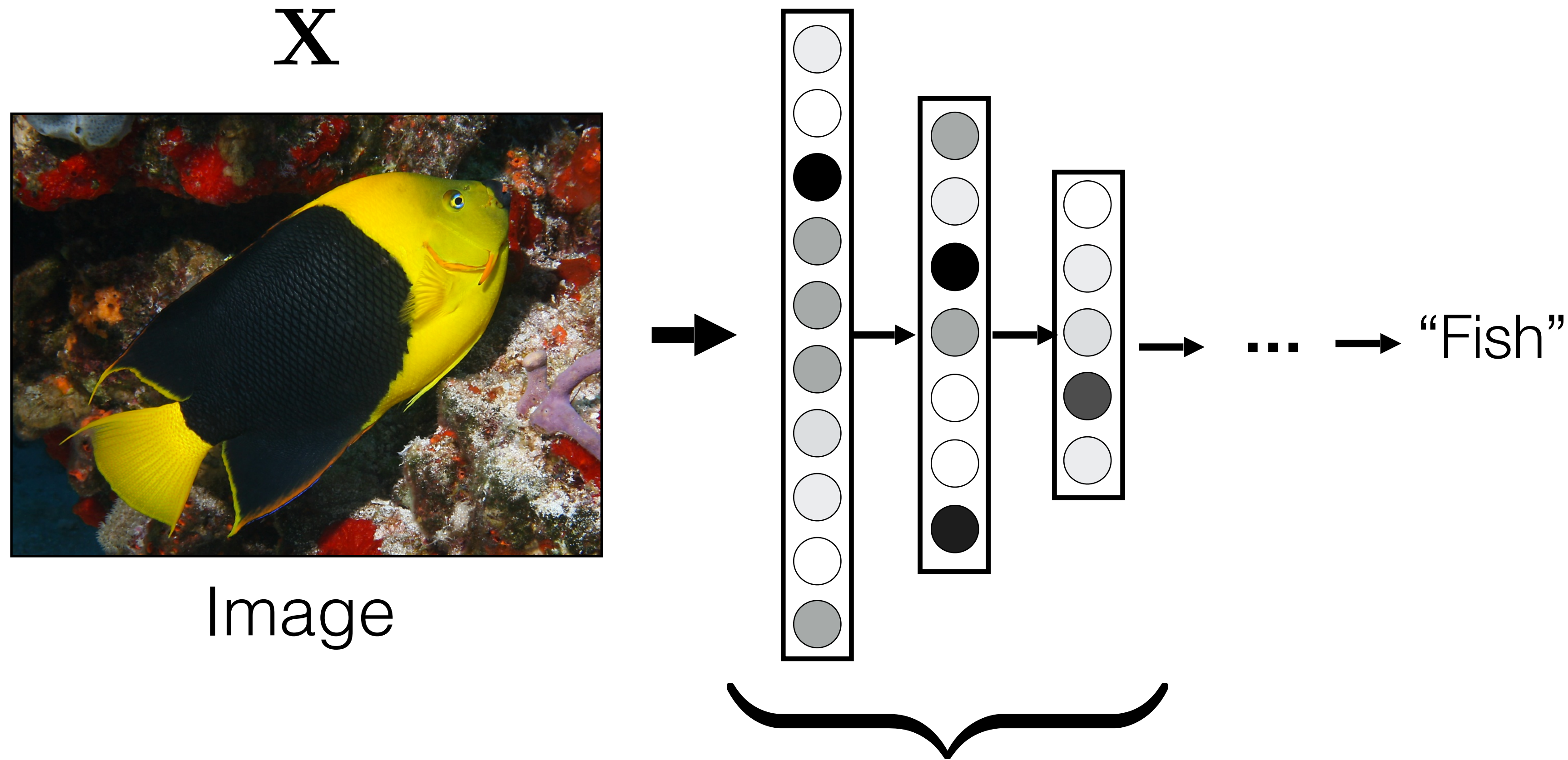


Investigating a representation via similarity analysis

Deep nets and the primate brain both learn similar metric spaces.

Deep nets organize visual information similarly to how our brains do!

What do deep nets internally learn?



A CNN is a multiscale, hierarchical representation of data

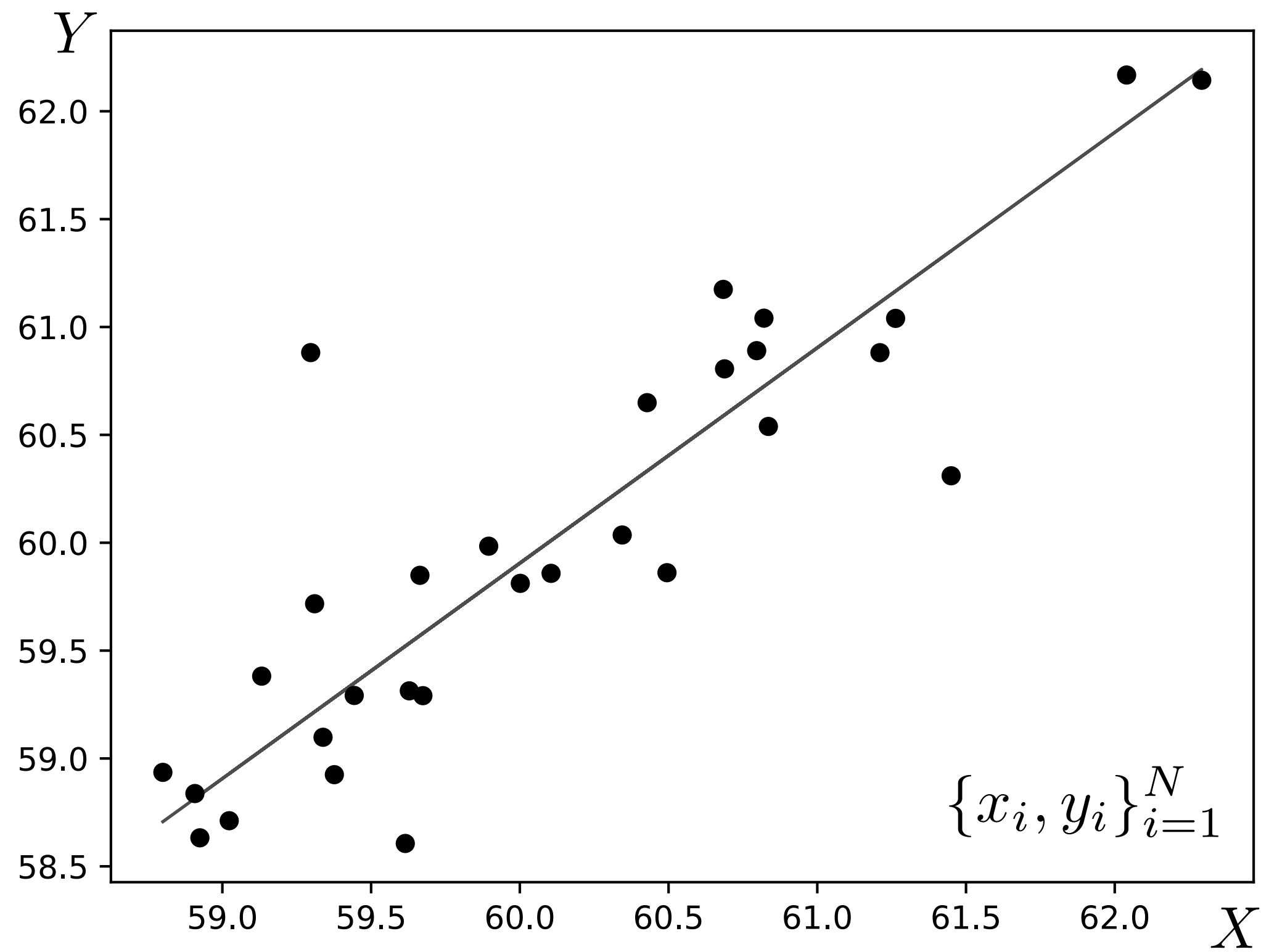
Transfer learning

“Generally speaking, a good representation is one that makes a subsequent learning task easier.” — Deep Learning textbook

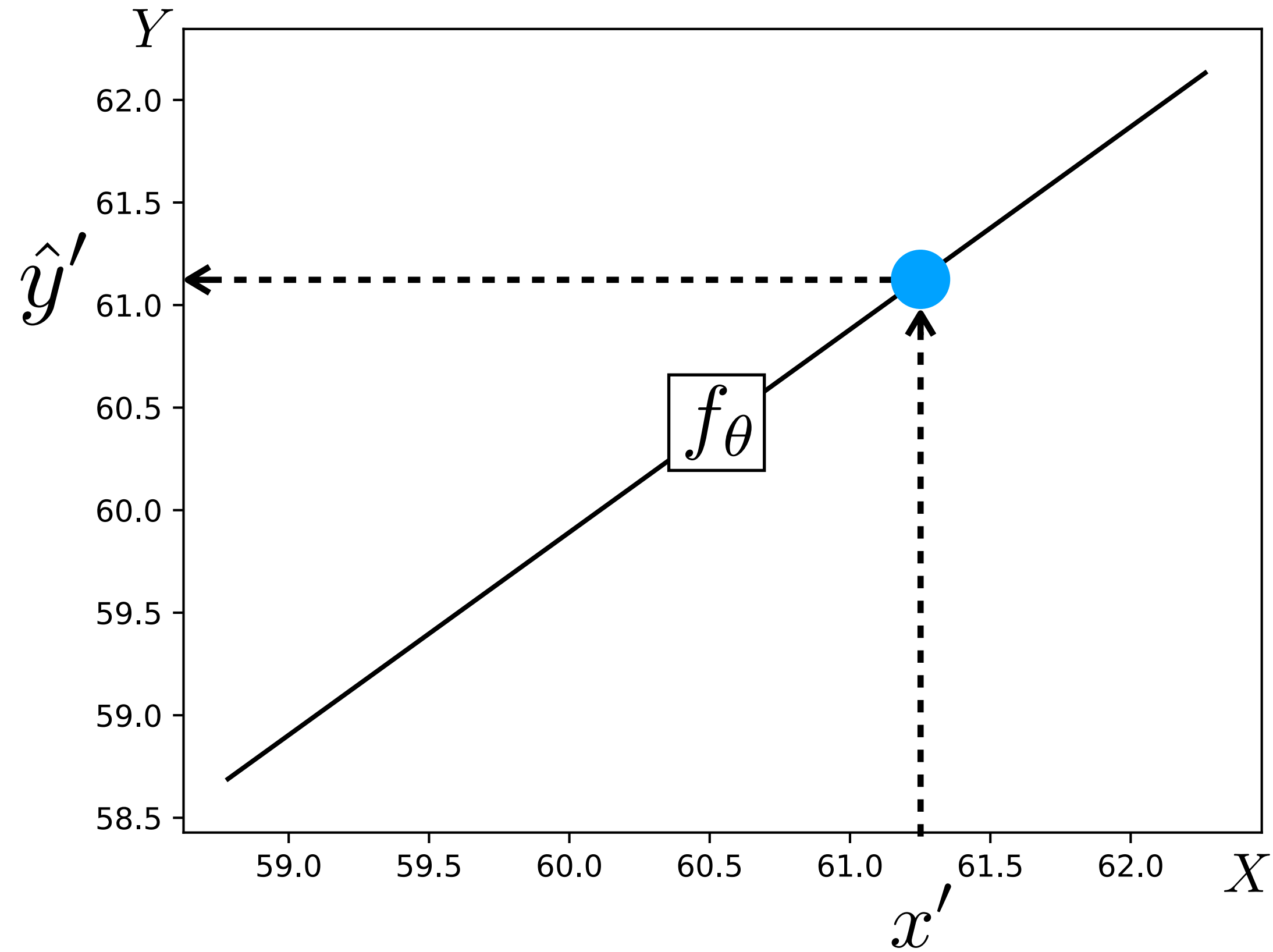


?

Training

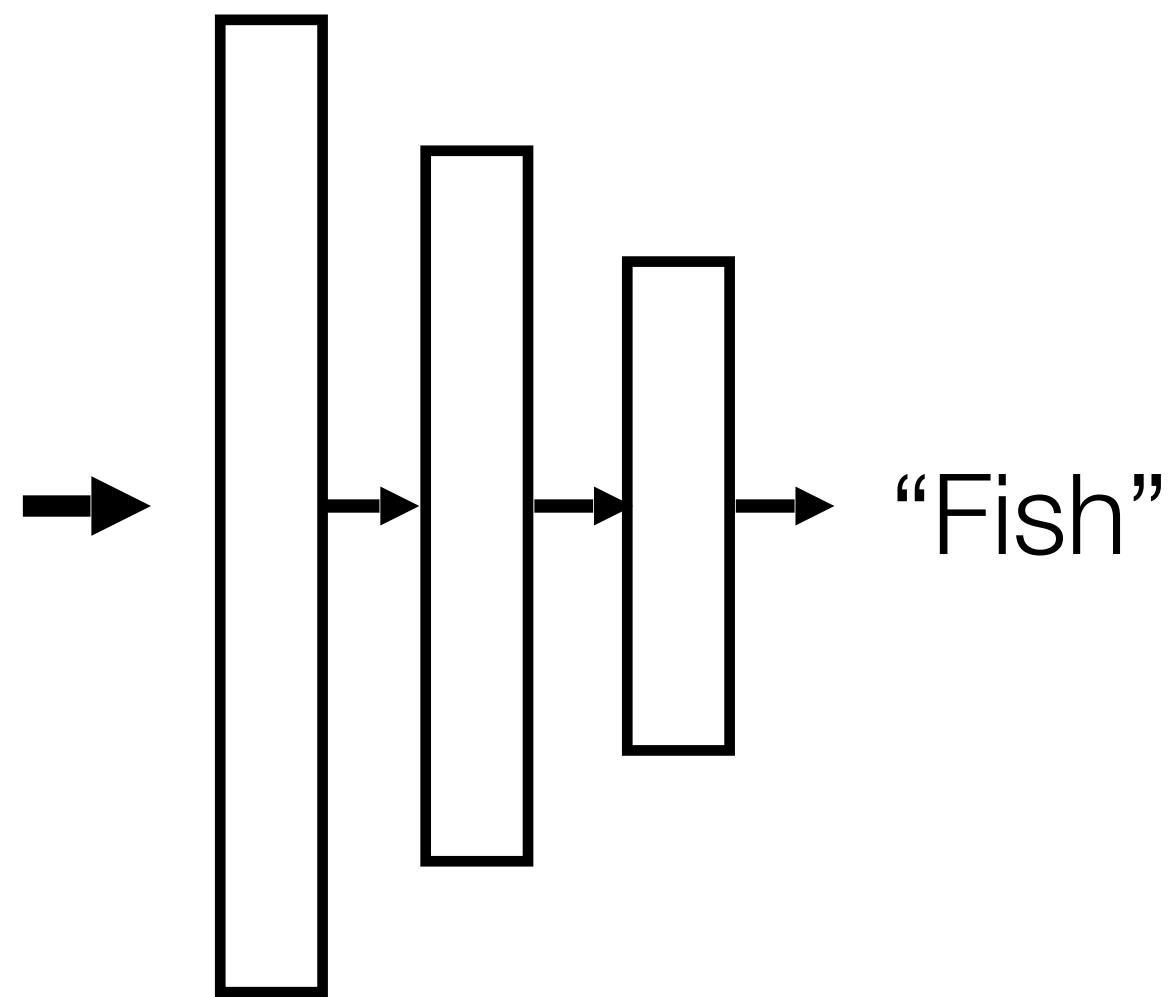


Testing



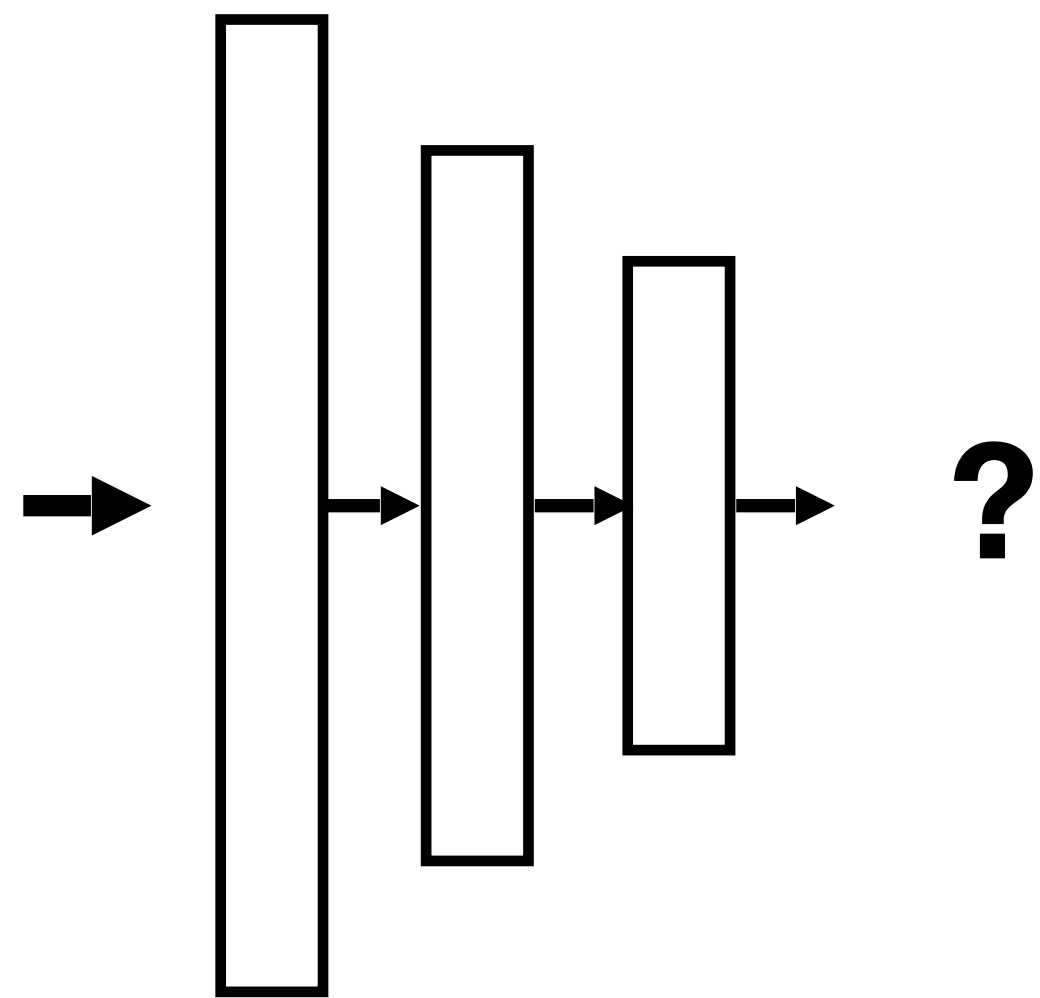
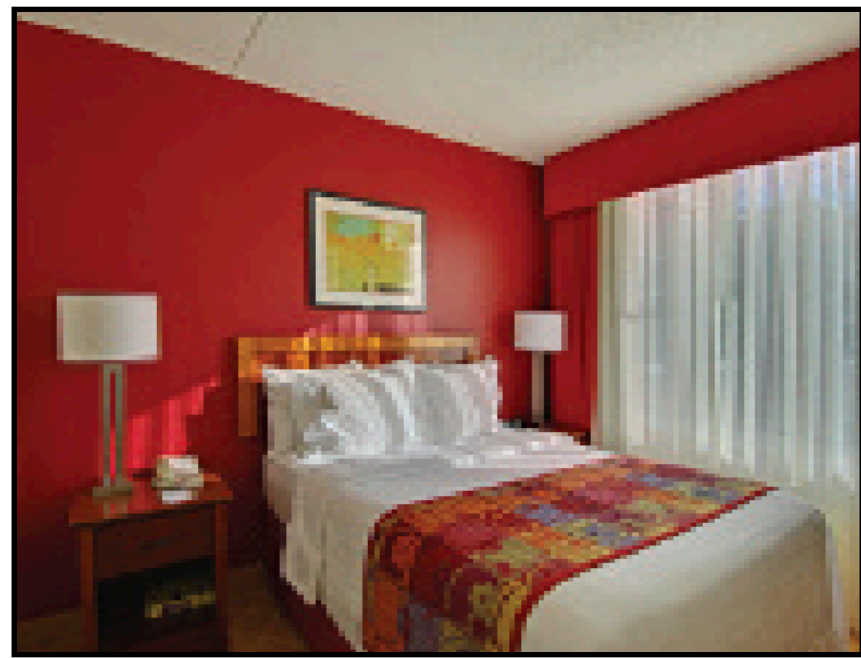
Training

Object recognition



Testing

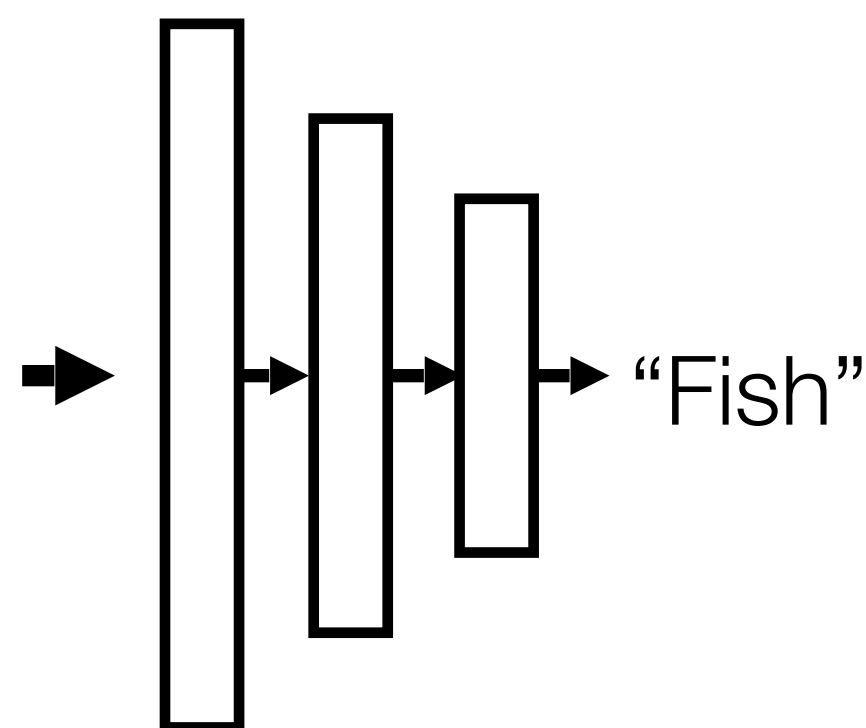
Place recognition



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

Pretraining

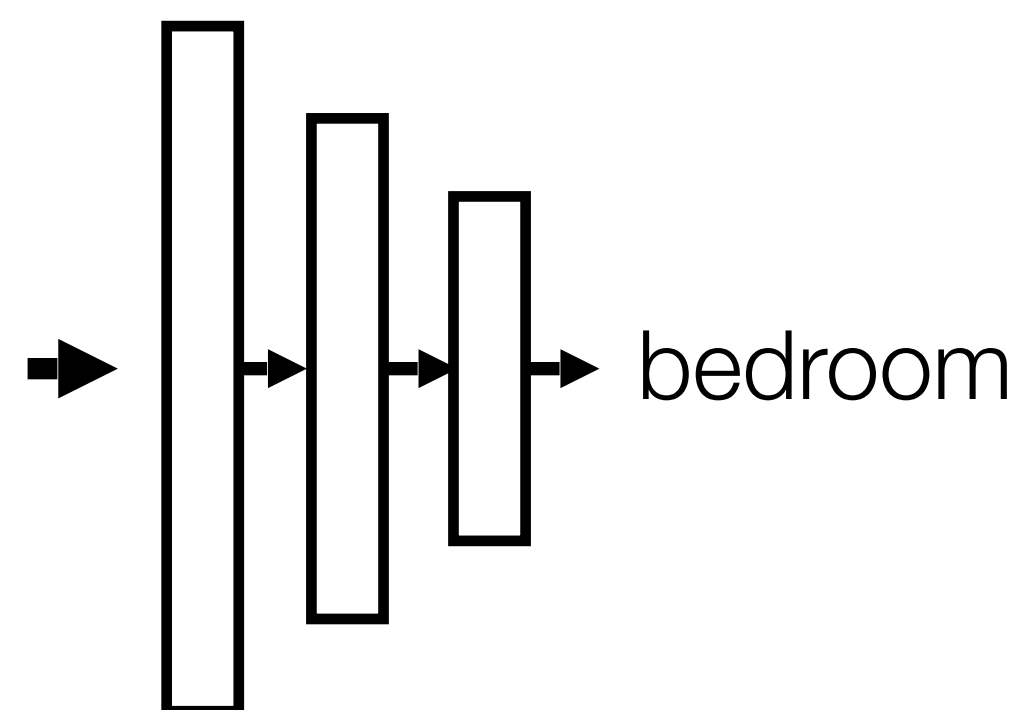
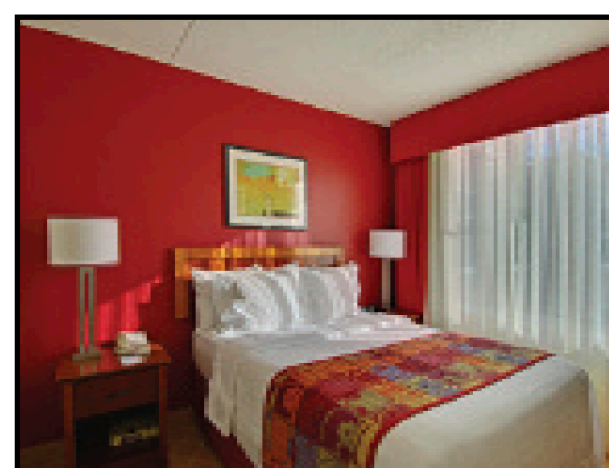
Object recognition



A lot of data

Finetuning

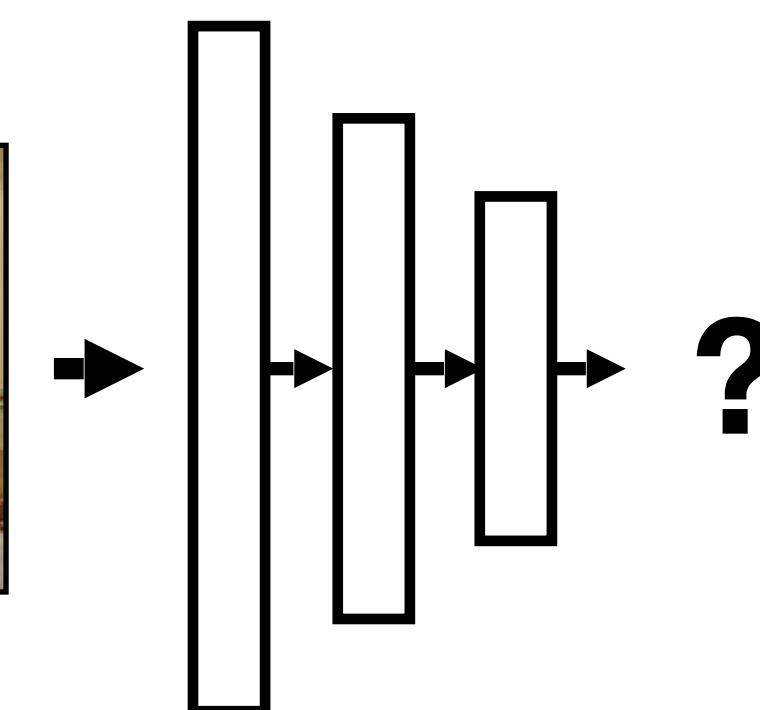
Place recognition



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.

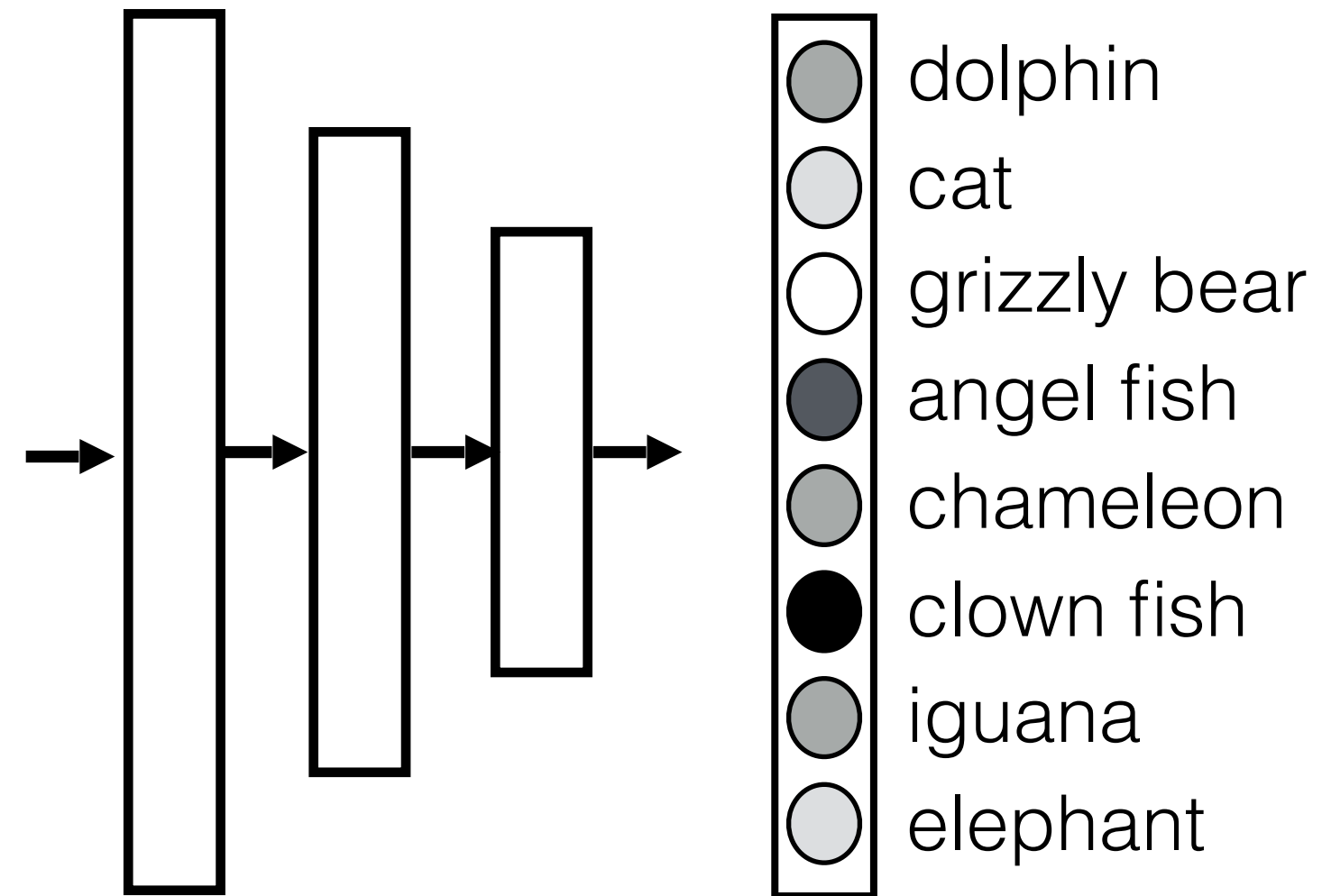
Finetuning in practice

- Pretrain a network on task A (often object recognition), resulting in parameters **W** and **b**
- Initialize a second network with some all of **W** and **b**
- Train the second network on task B, resulting in parameters **W'** and **b'**

Finetuning in practice

Pretraining

Object recognition



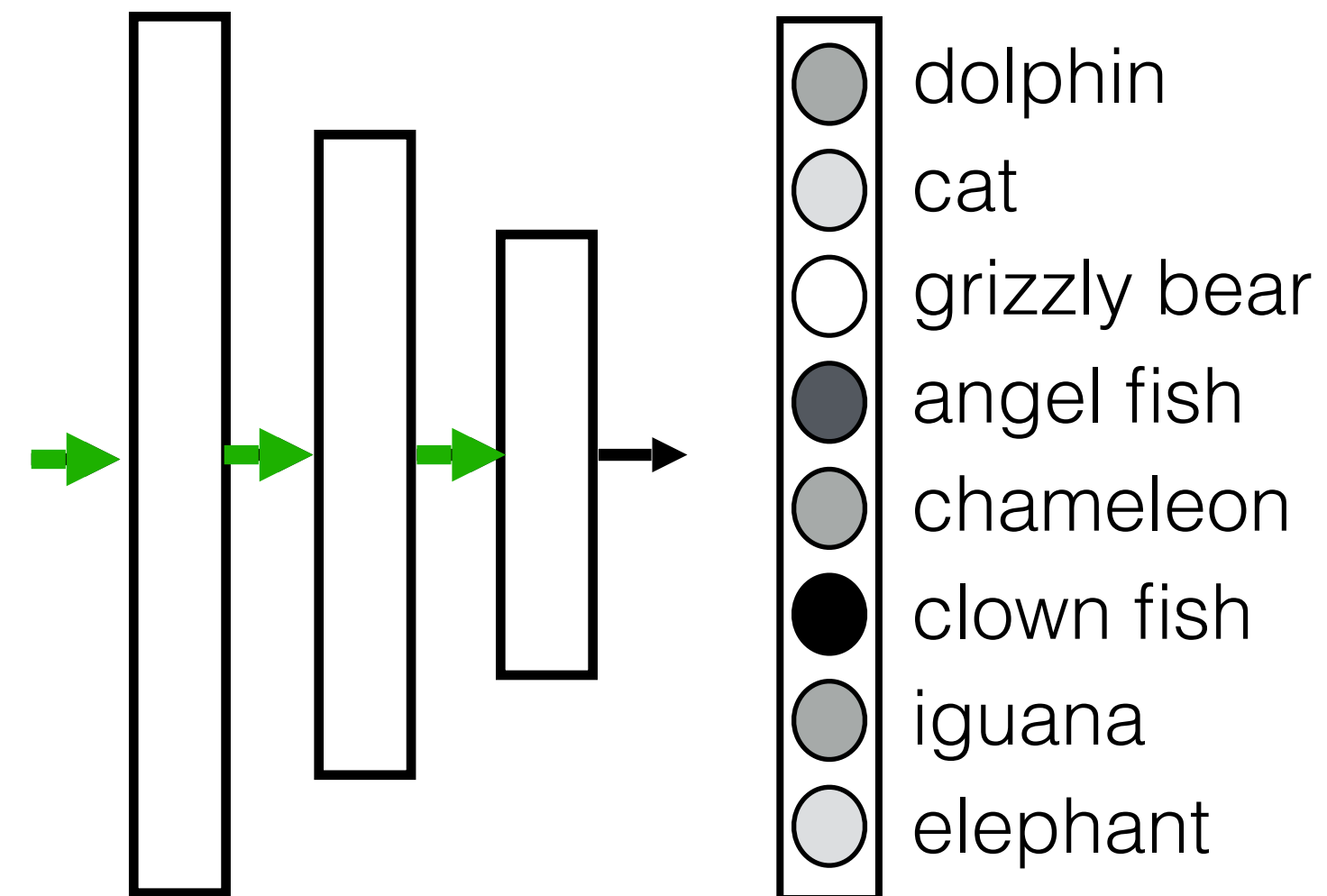
Finetuning

Place recognition

Finetuning in practice

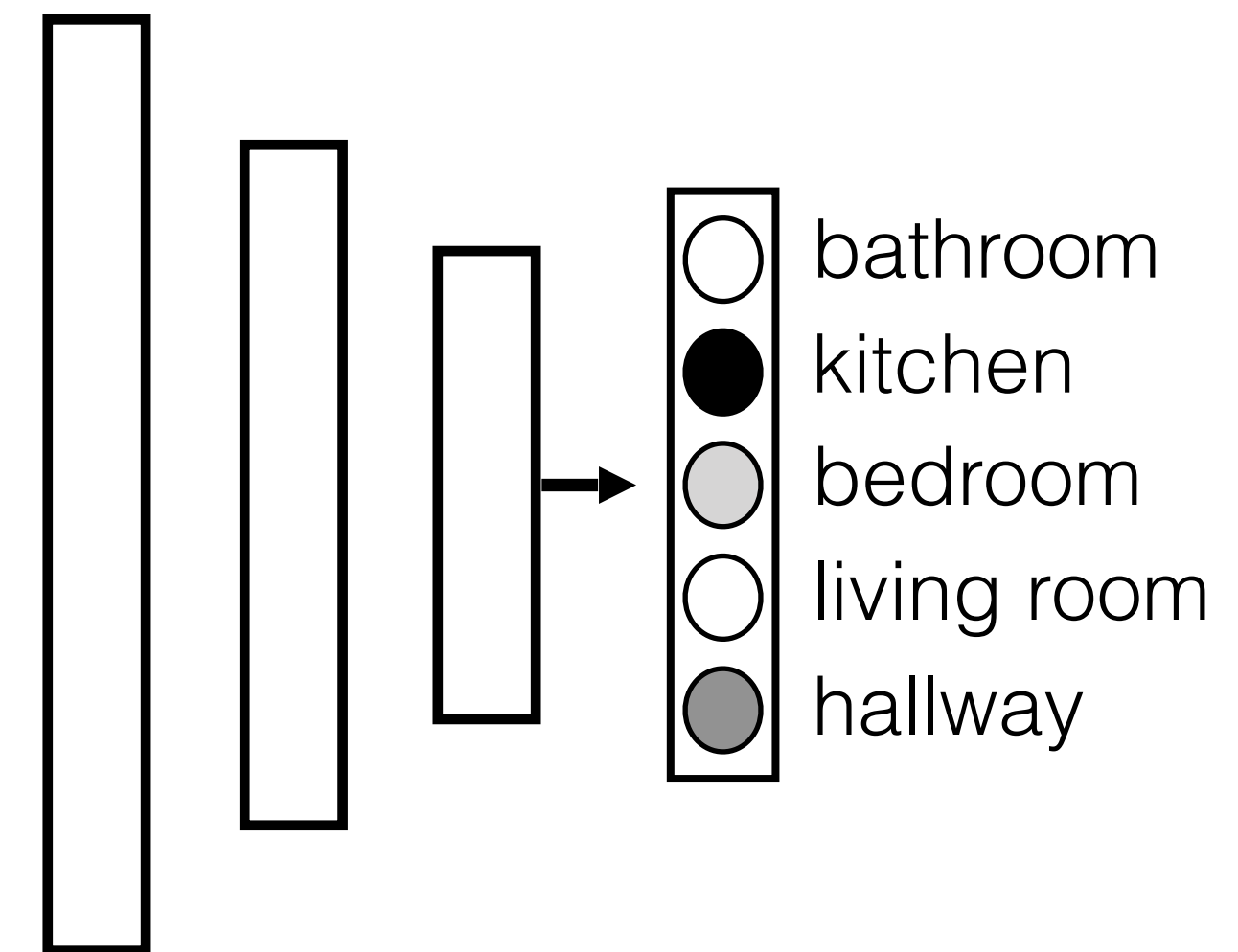
Pretraining

Object recognition



Finetuning

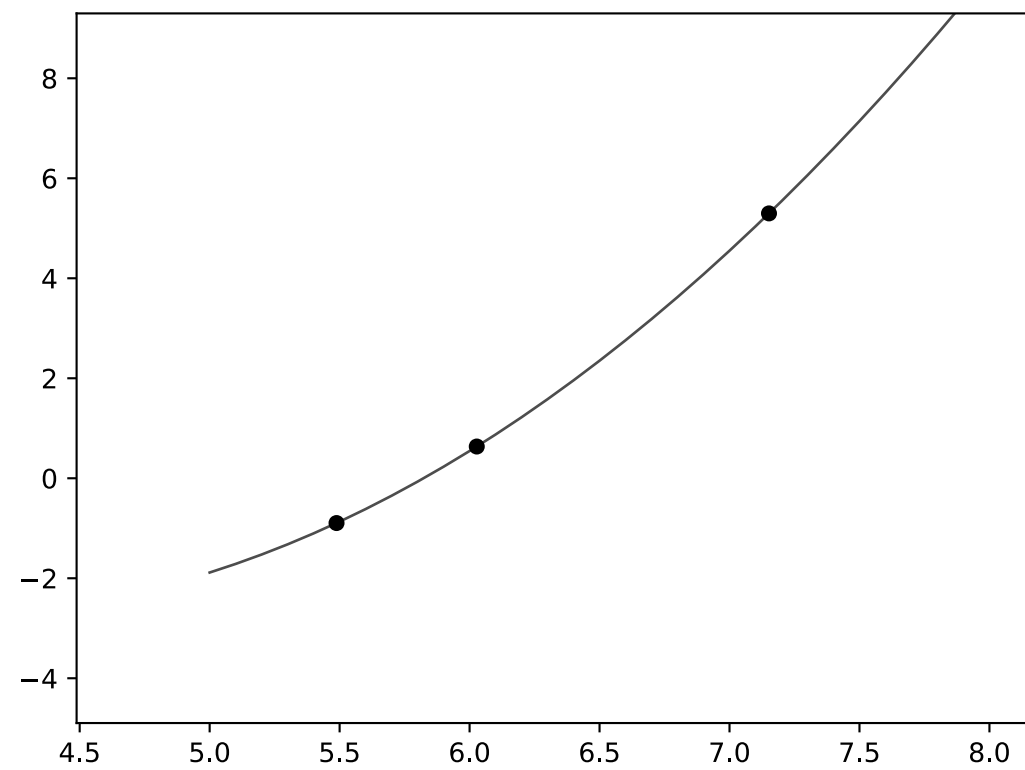
Place recognition



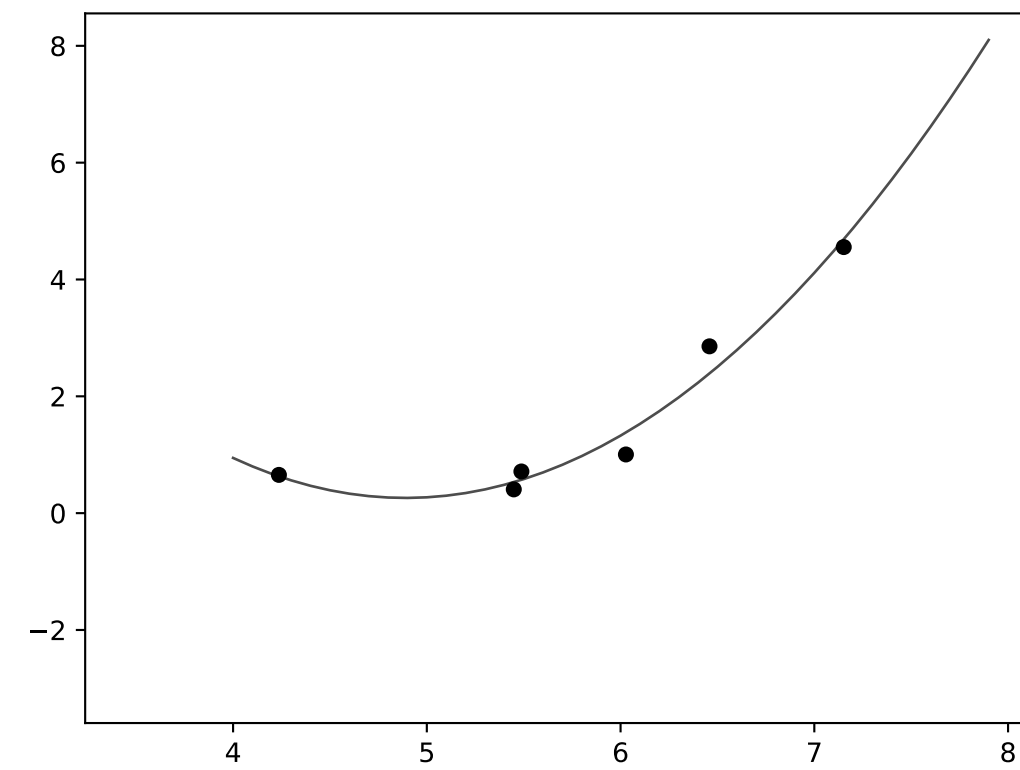
The “learned representation” is just the weights and biases, so that’s what we transfer

If we keep on finetuning for every new datapoint or task that comes our way, we get **online learning**. Humans seem to do this, we never stop learning.

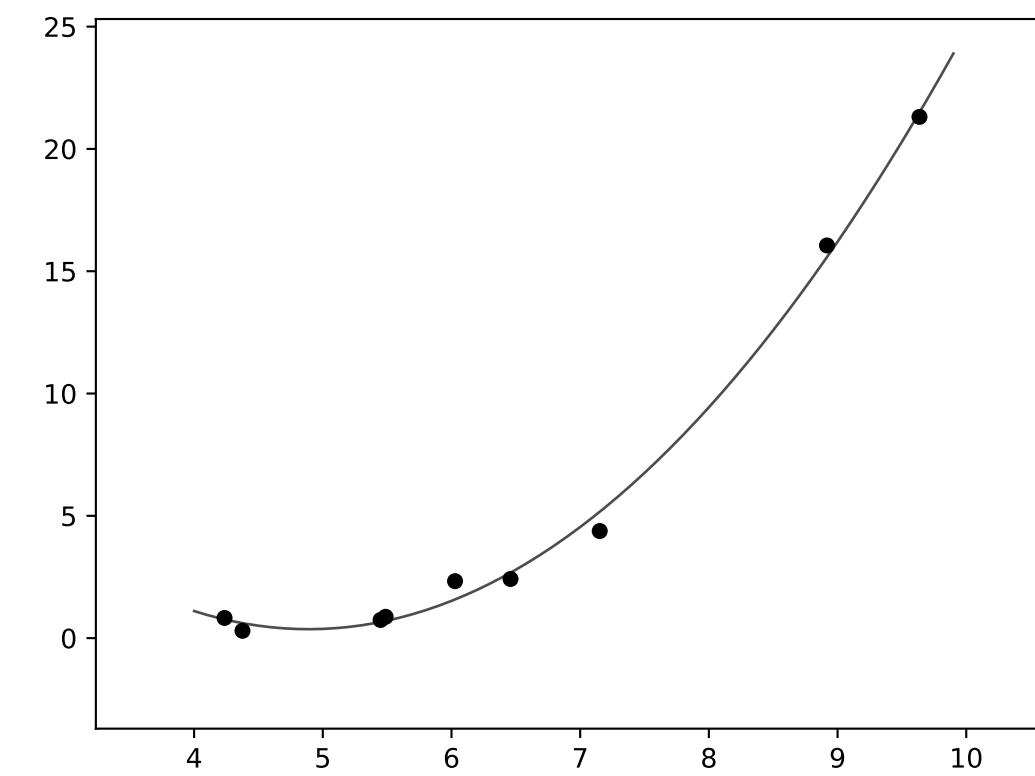
Training



More training



More training



...

Supervised object recognition

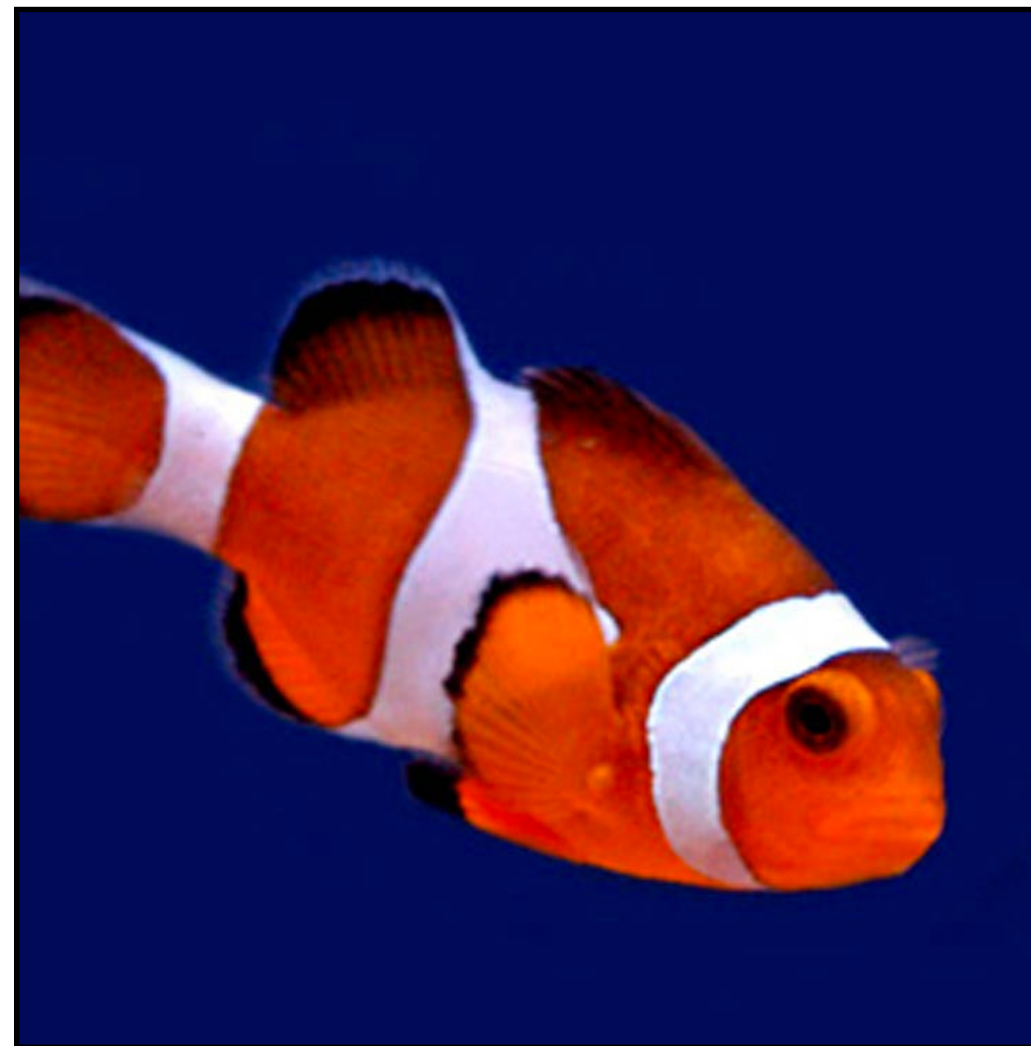
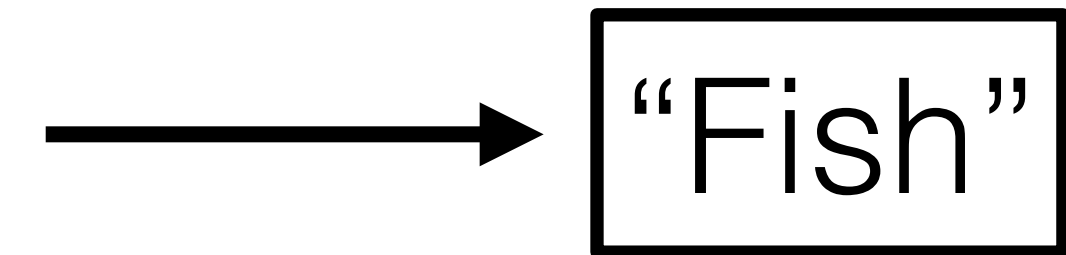


image X

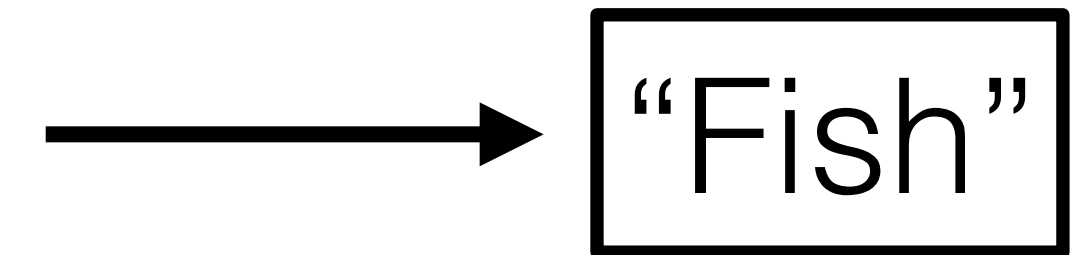


label Y

Supervised object recognition



image X

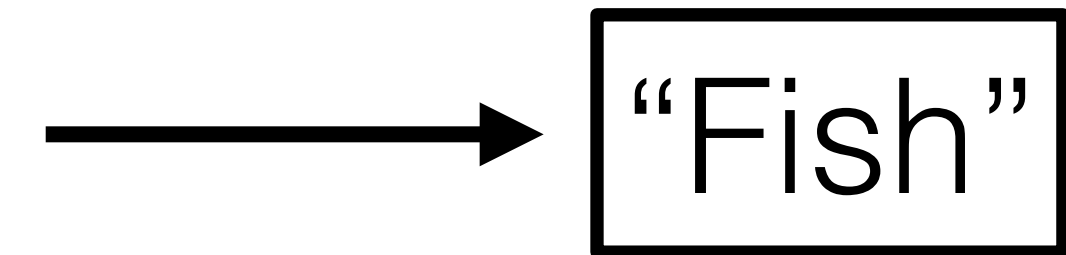


label Y

Supervised object recognition



image X



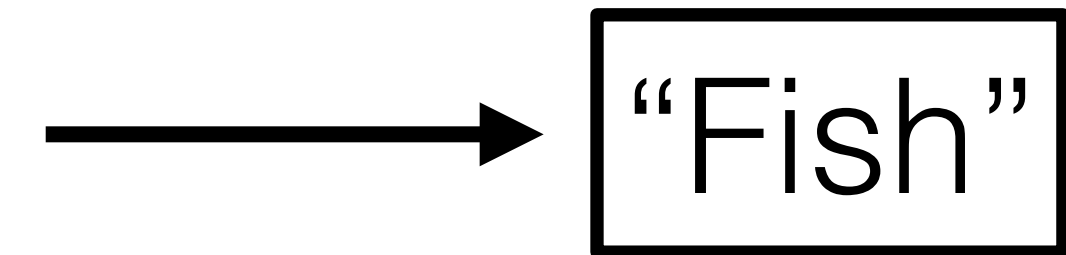
label Y

Supervised object recognition



⋮

image X



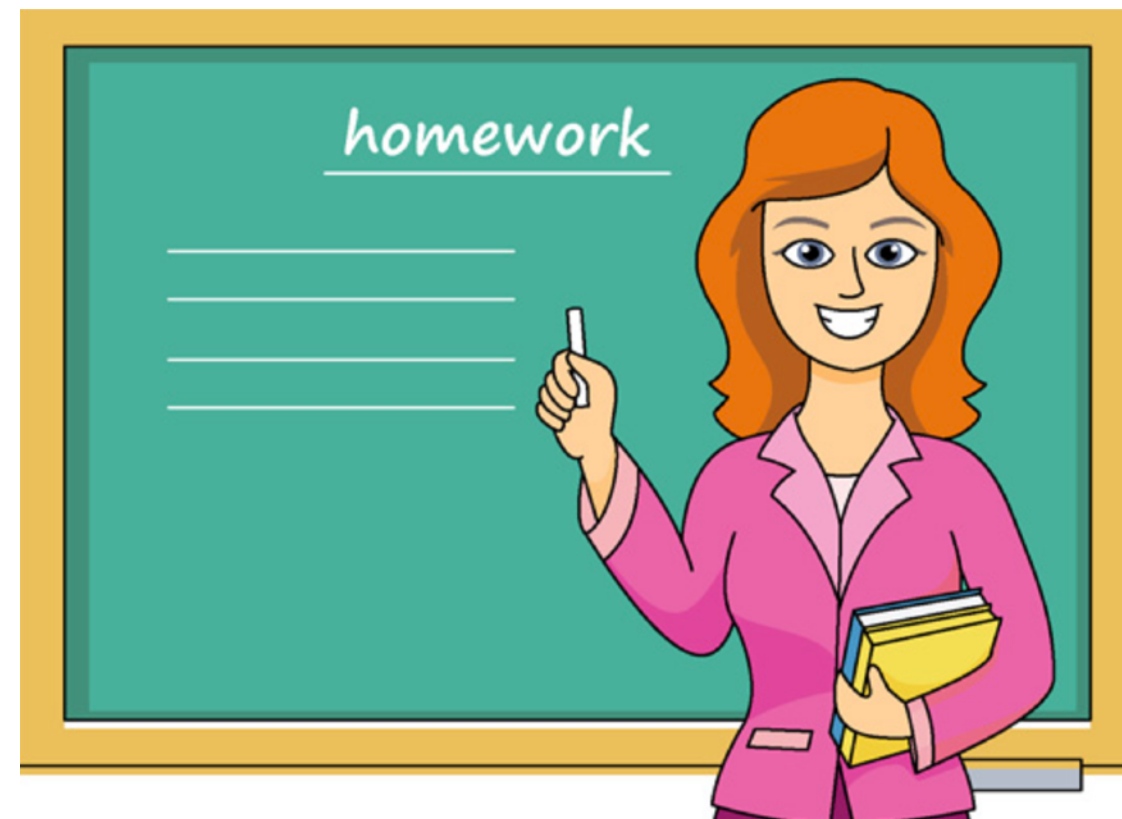
label Y



Supervised computer vision

Hand-curated training data

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



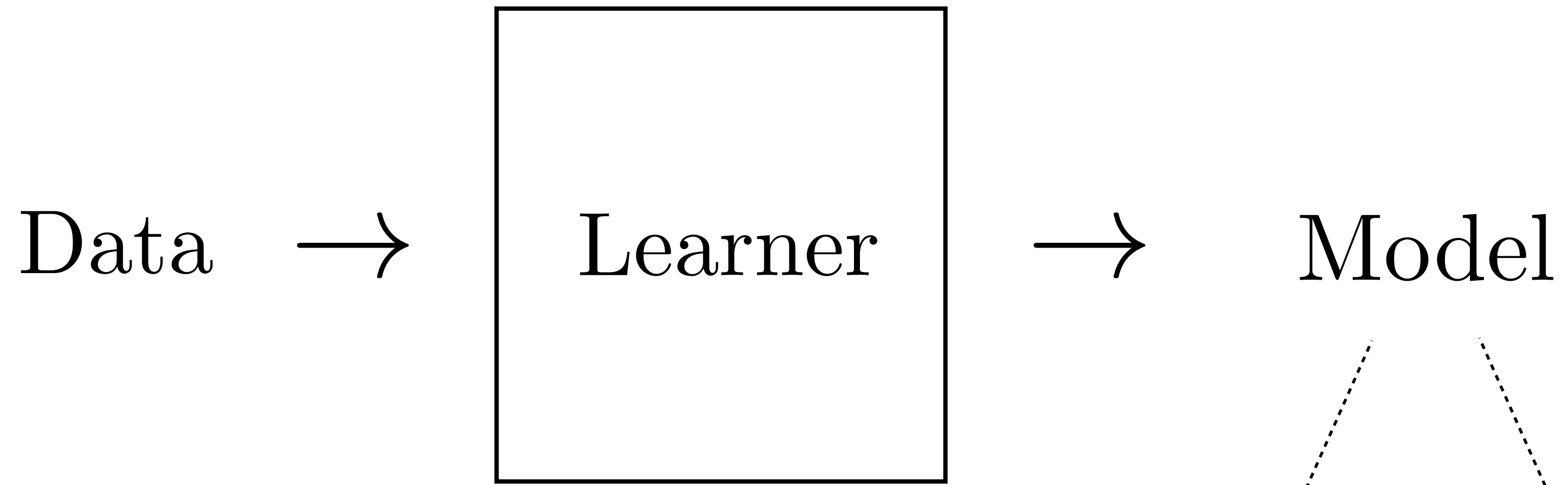
Vision in nature

Raw unlabeled training data

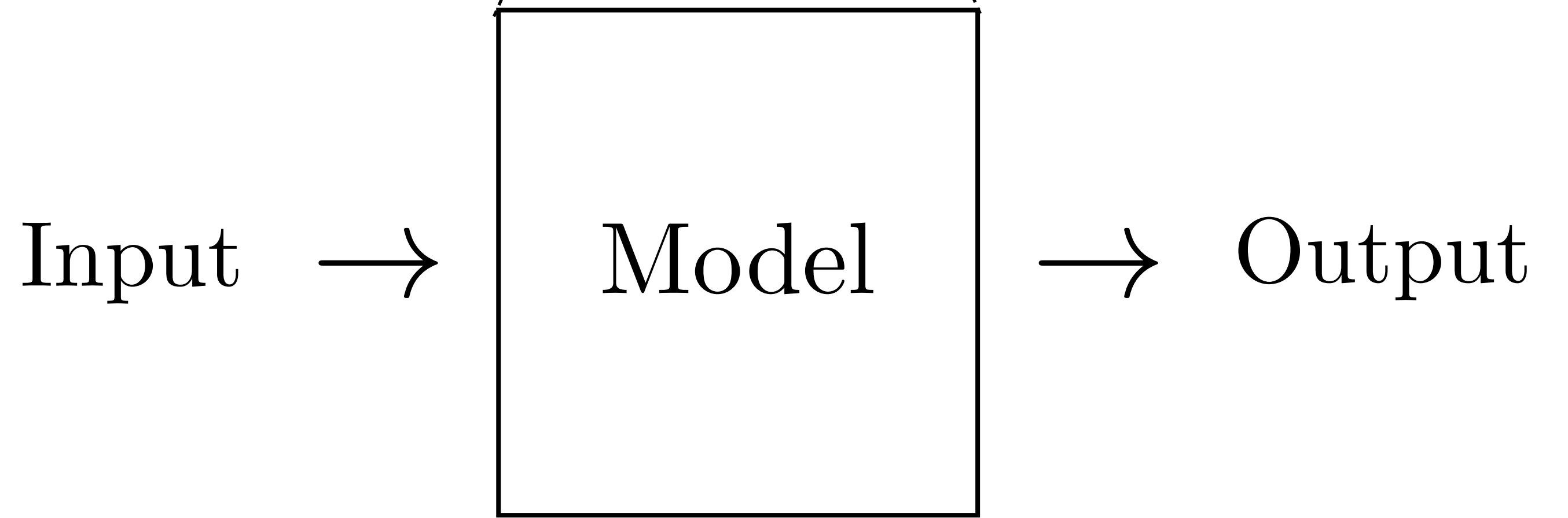
- + Cheap
- Noisy
- Harder to interpret



Learning



Inference



Learning from examples

(aka **supervised learning**)

Training data

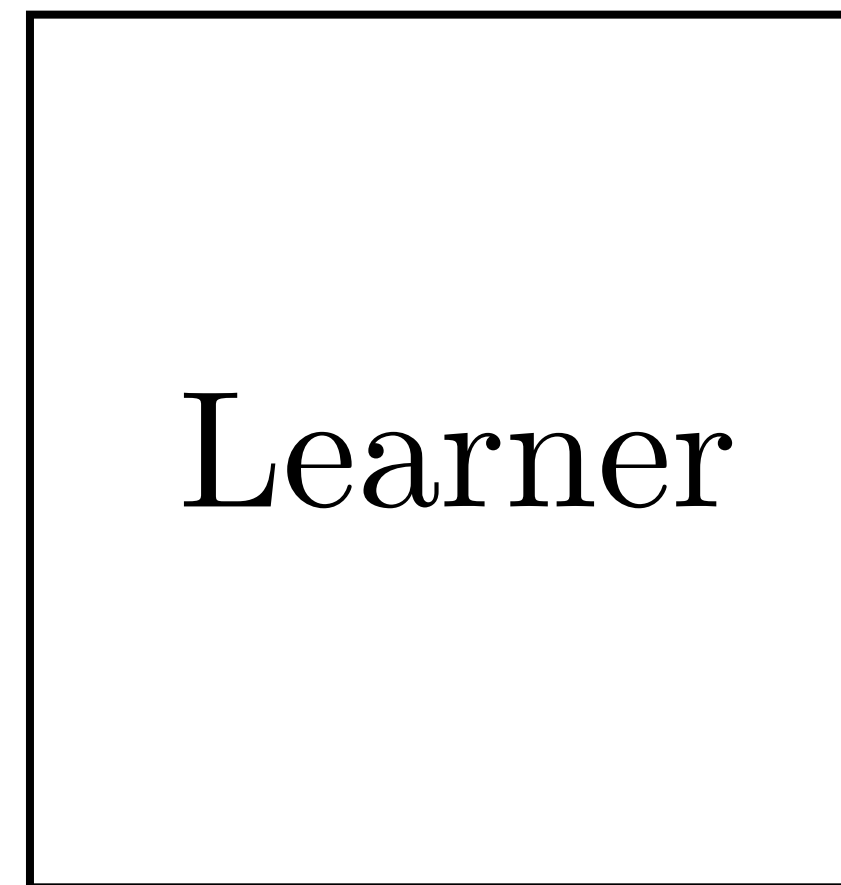
$\{x_1, y_1\}$

$\{x_2, y_2\}$

$\{x_3, y_3\}$

...

→



→

$f : X \rightarrow Y$

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

Learning without examples

(includes **unsupervised learning** and **reinforcement learning**)

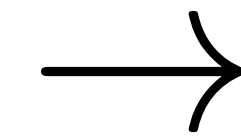
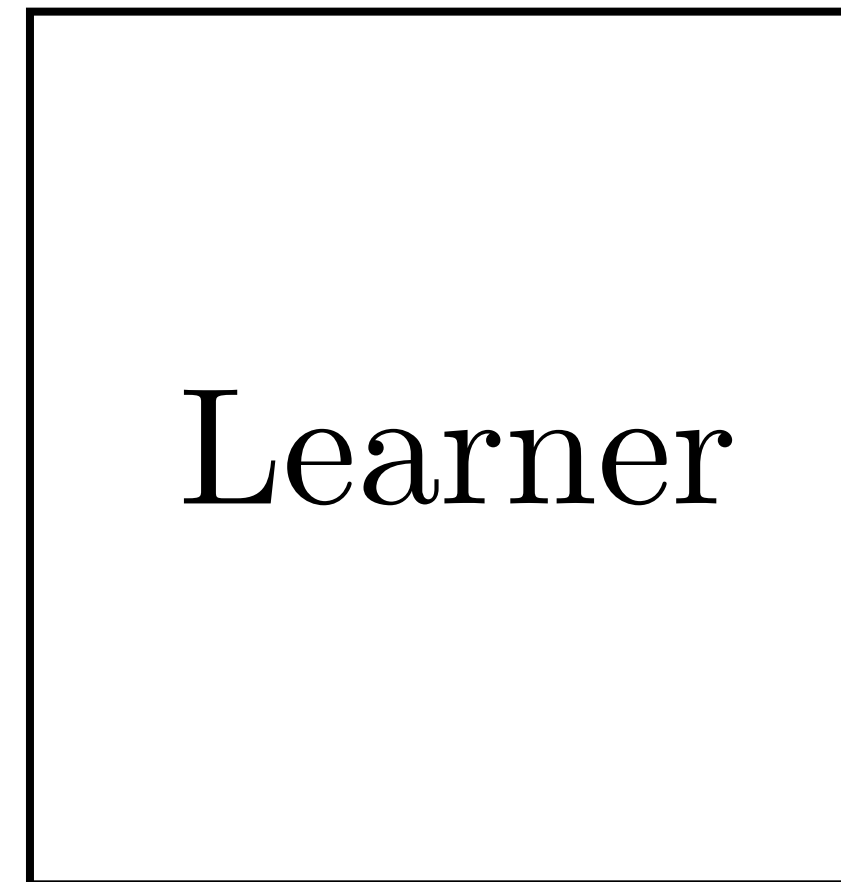
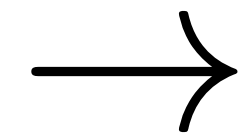
Data

$\{x_1\}$

$\{x_2\}$

$\{x_3\}$

...



?

Unsupervised Representation Learning

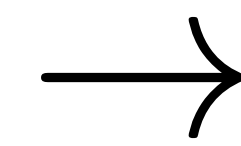
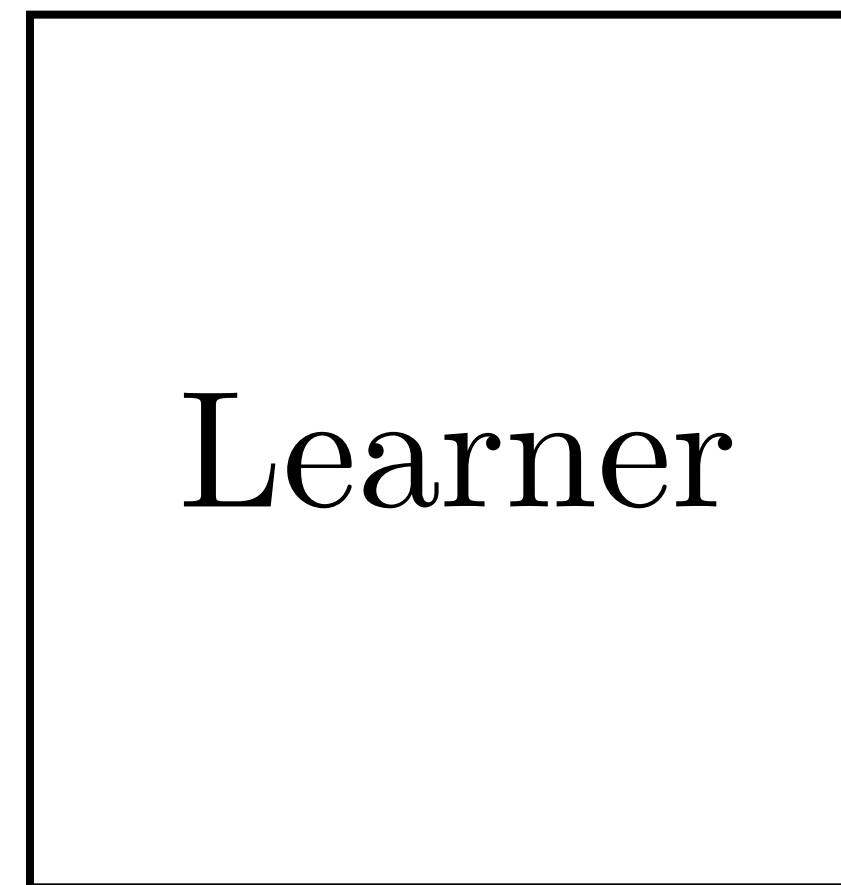
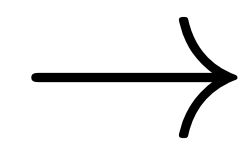
Data

$\{x_1\}$

$\{x_2\}$

$\{x_3\}$

...



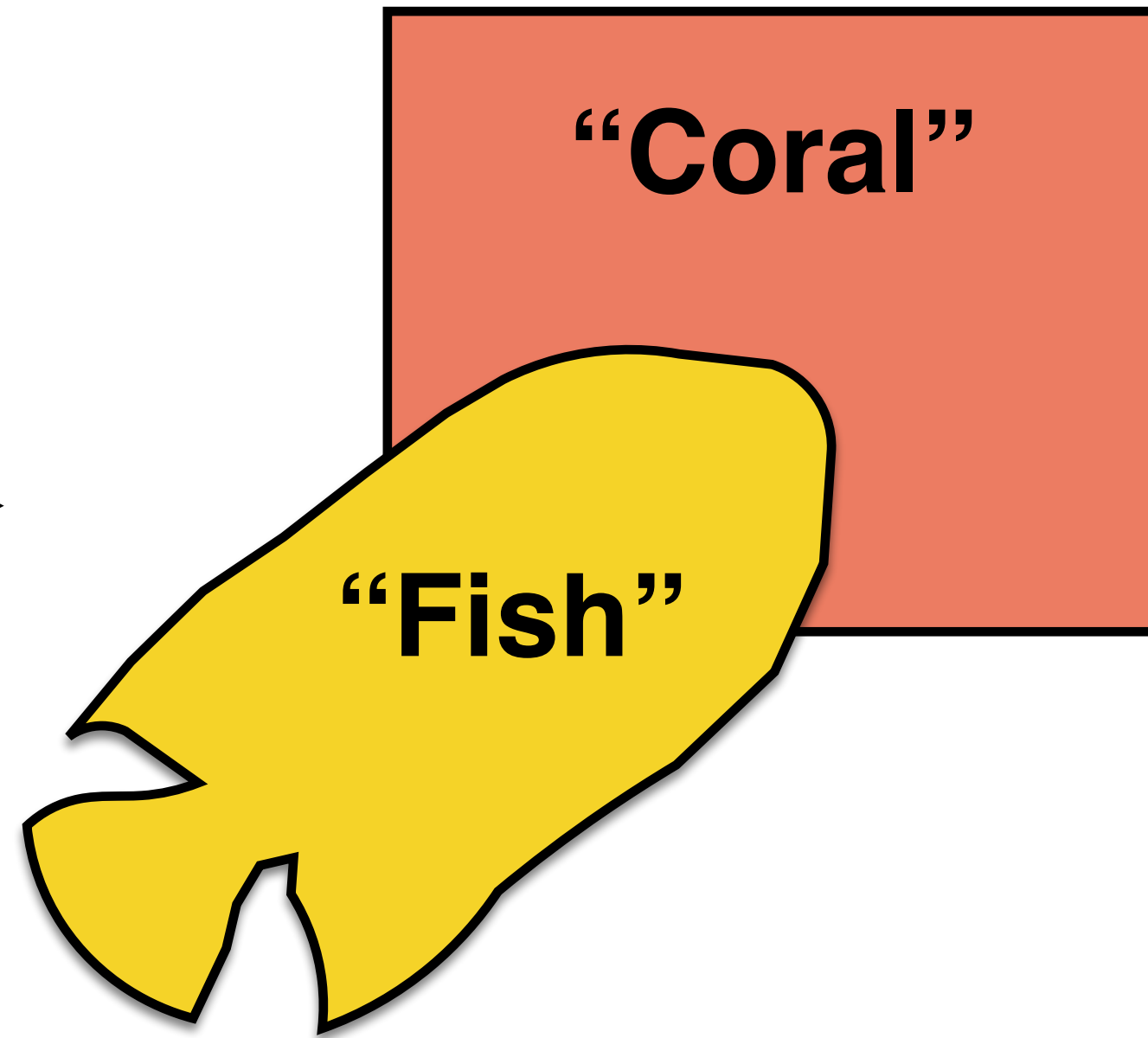
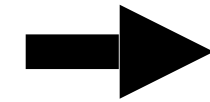
Representations

Unsupervised Representation Learning

X

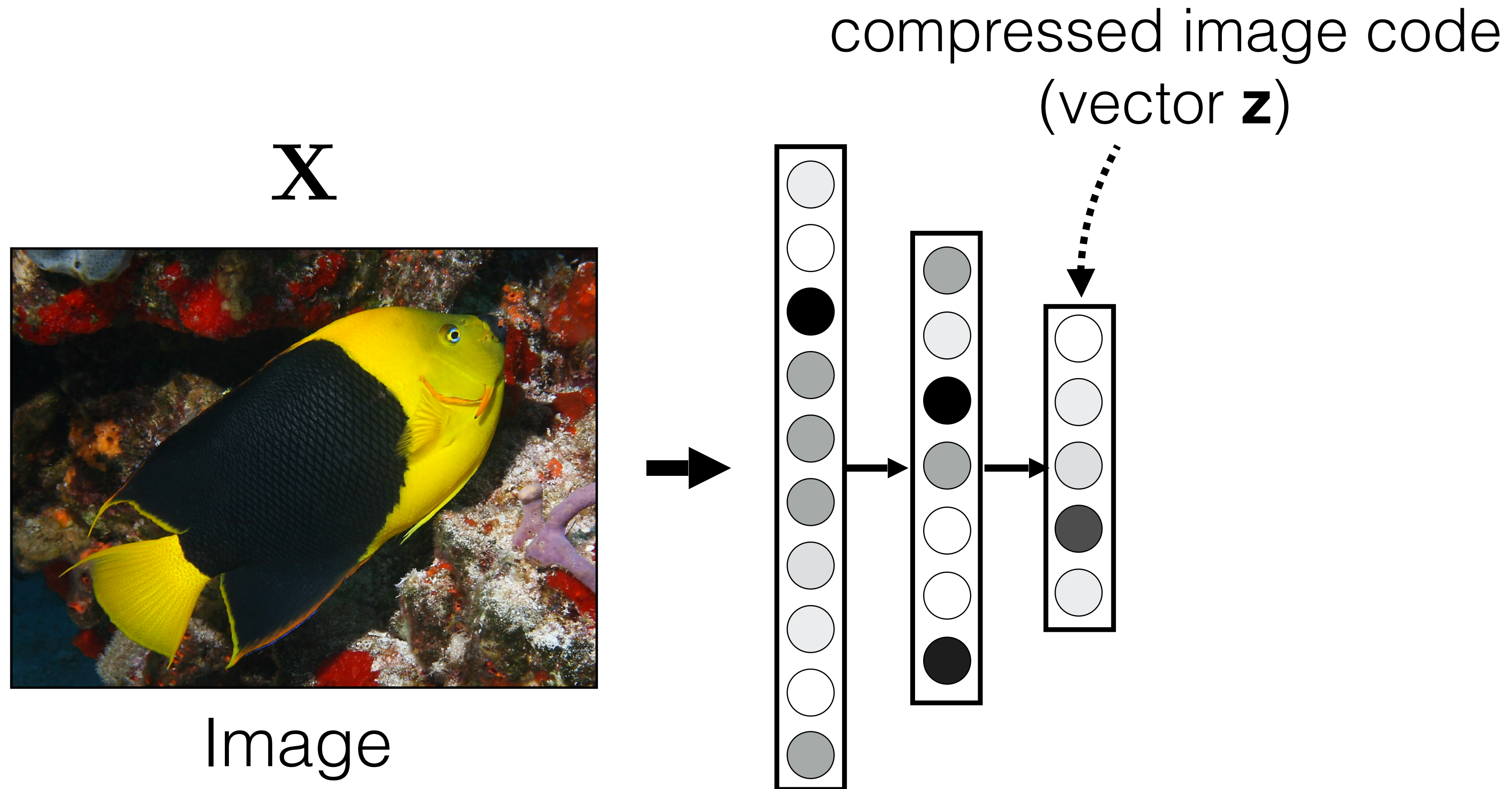


Image

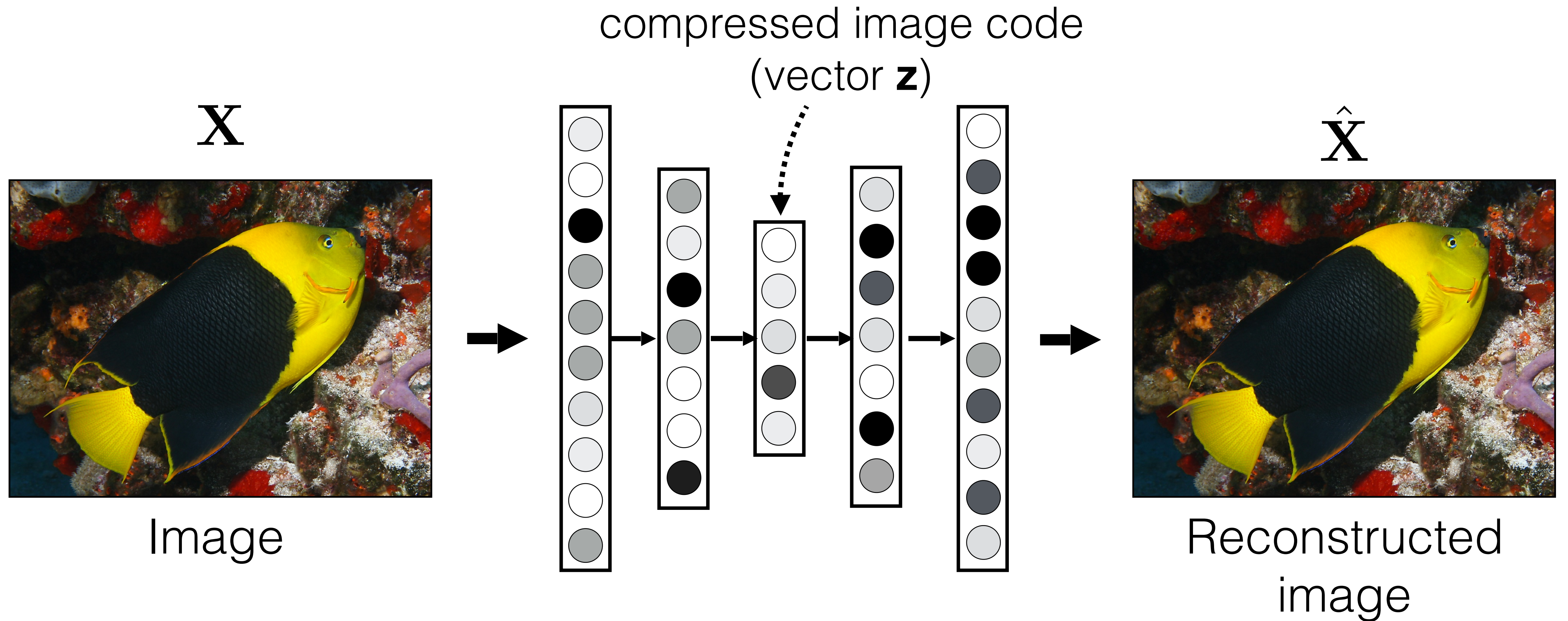


Compact mental
representation

Unsupervised Representation Learning



Unsupervised Representation Learning



“Autoencoder”

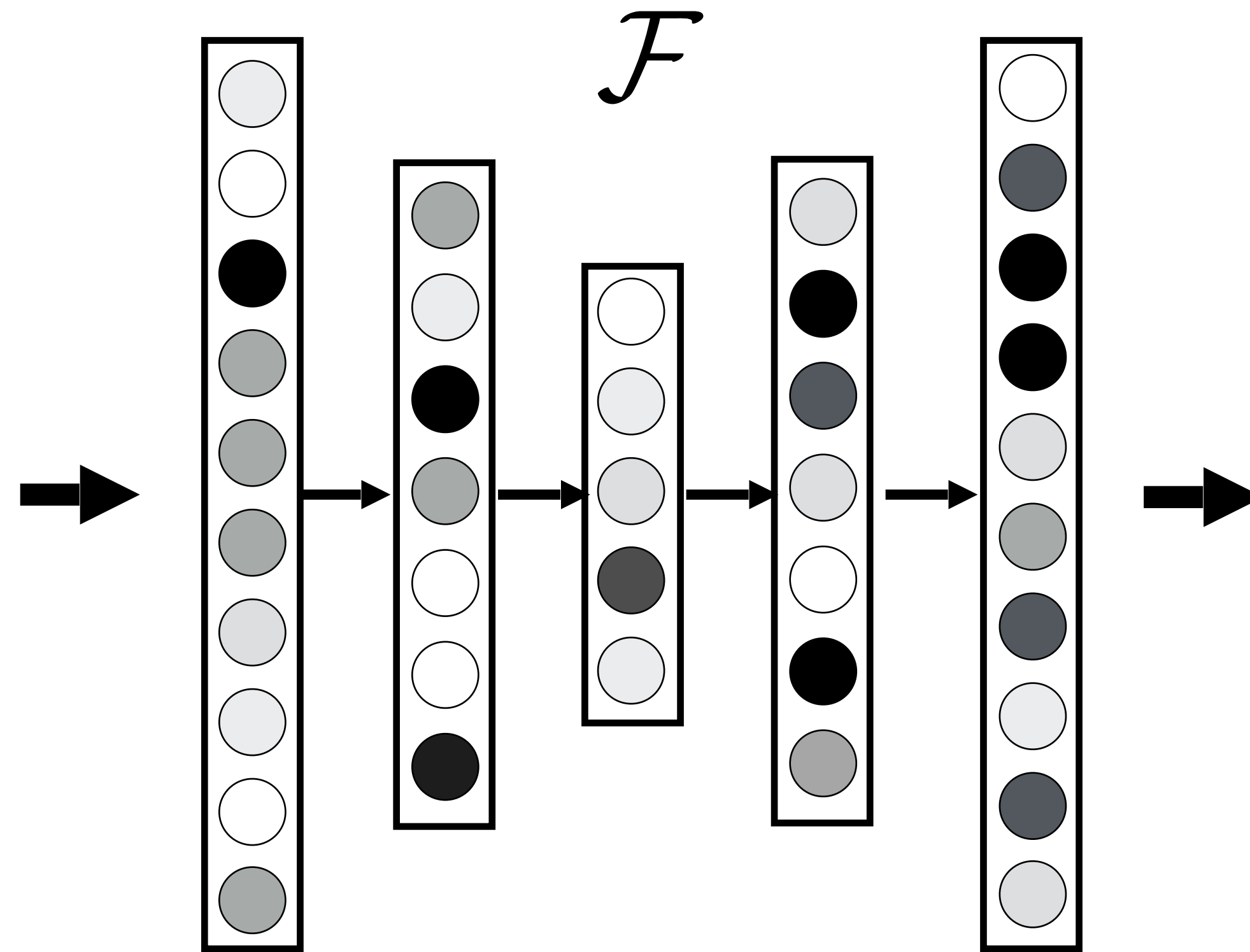
[e.g., Hinton & Salakhutdinov, Science 2006]

Autoencoder

\mathbf{X}



Image



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



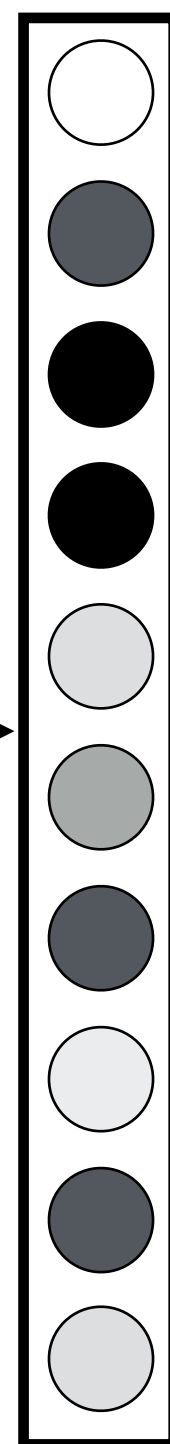
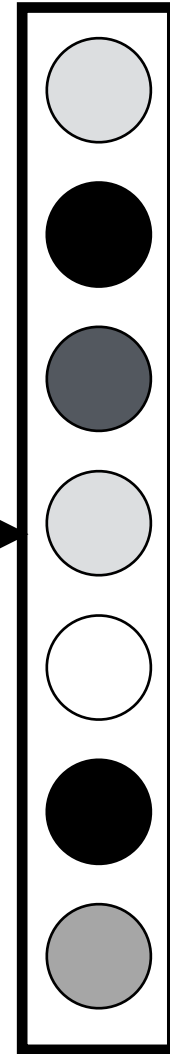
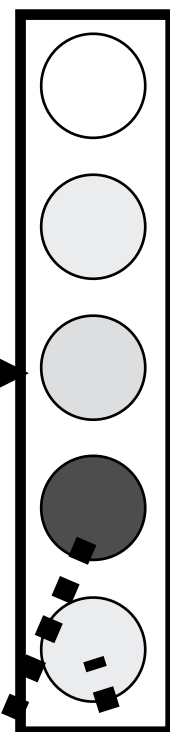
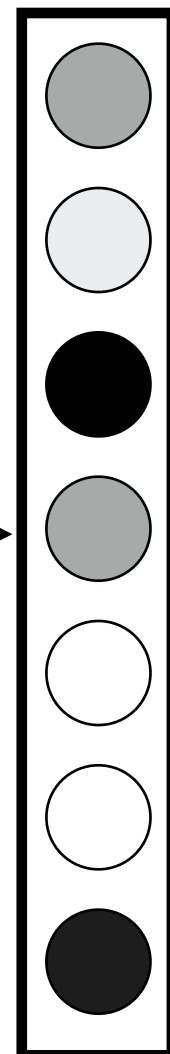
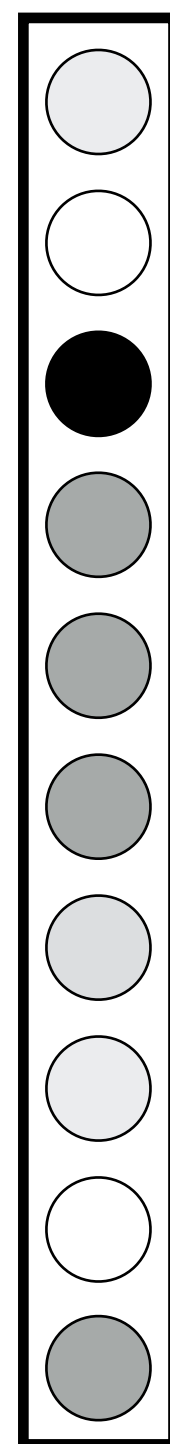
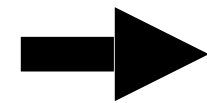
Reconstructed
image

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

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

\mathbf{X} 

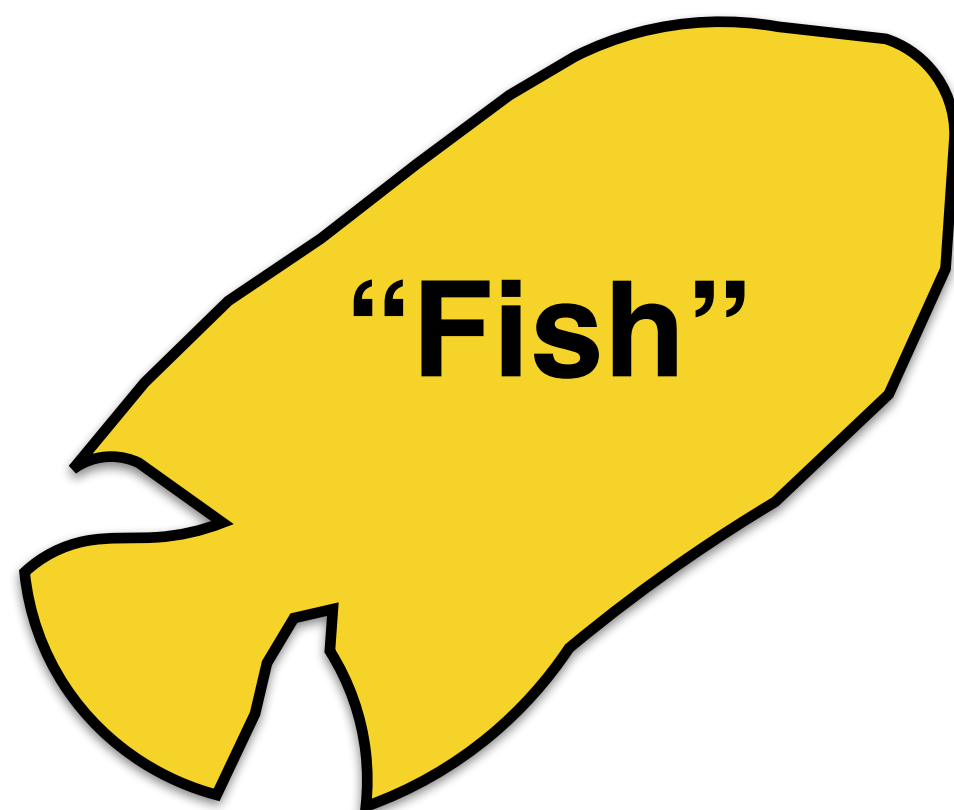
Image



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

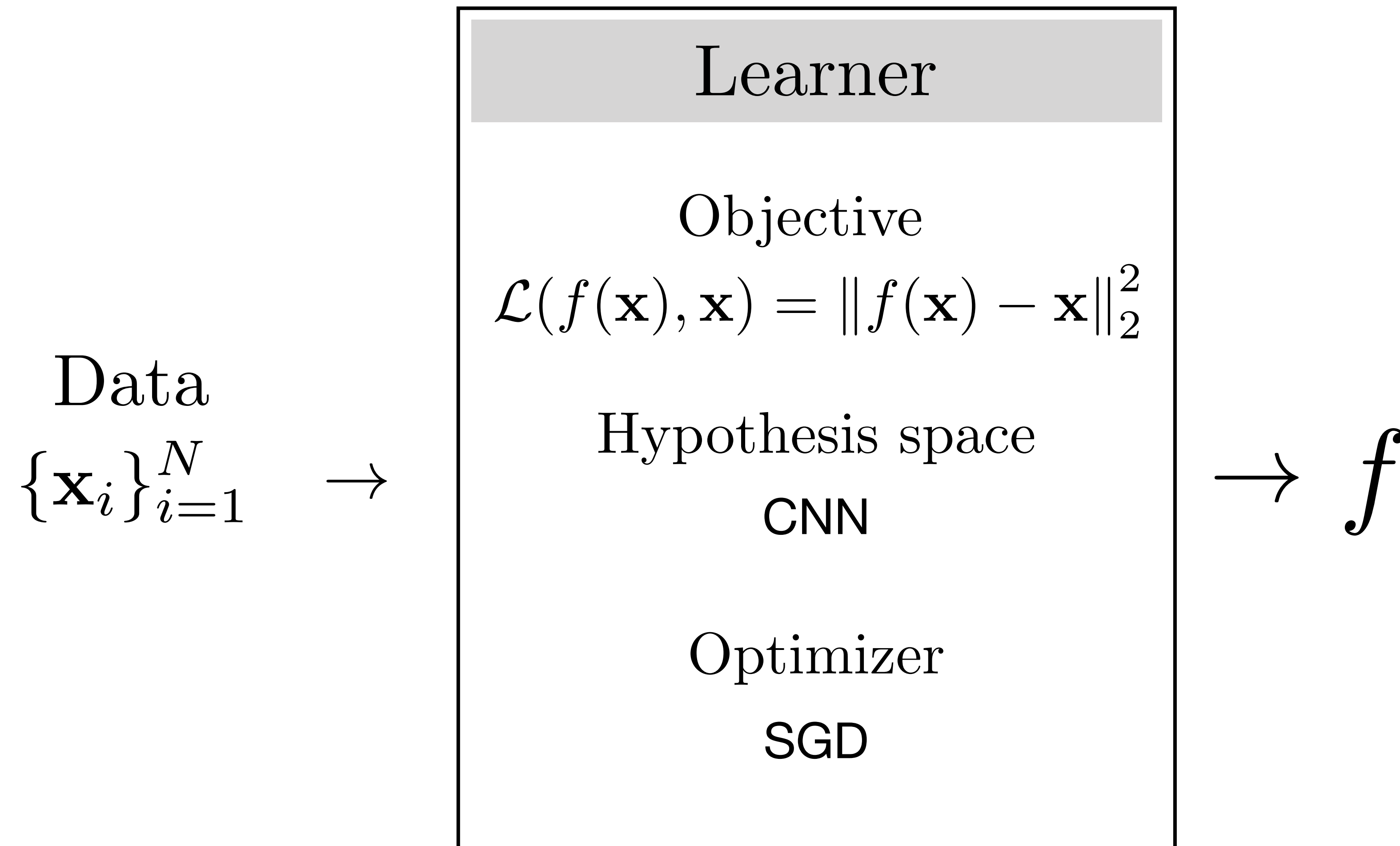


Reconstructed image

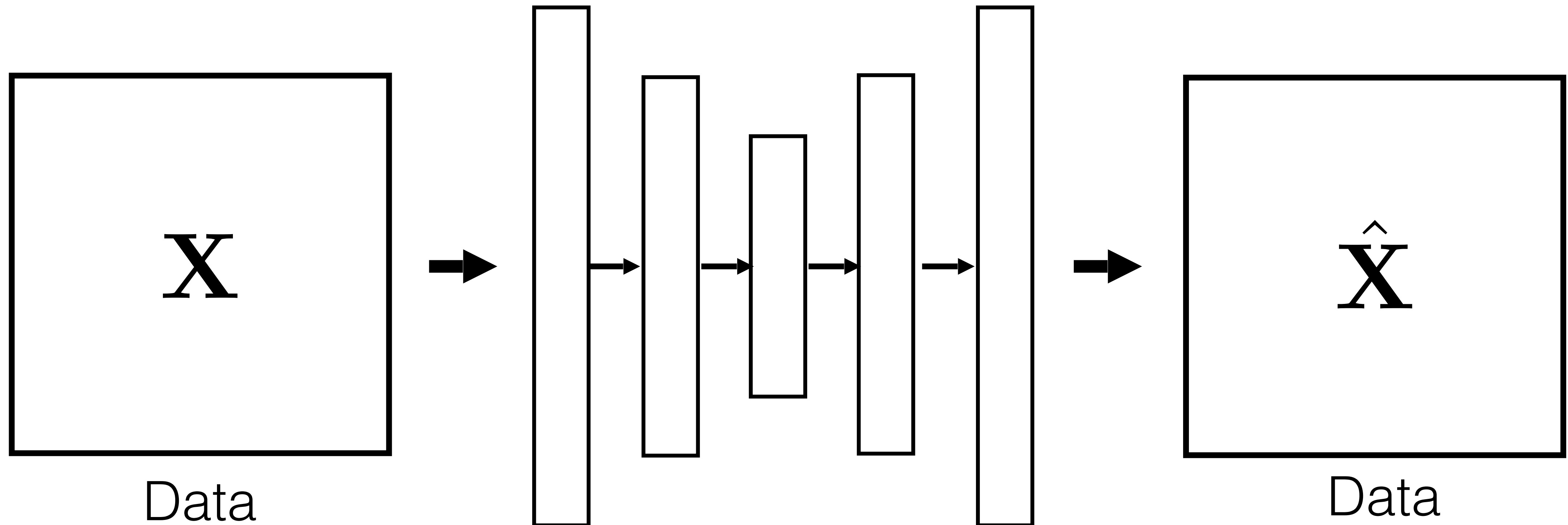


[e.g., Hinton & Salakhutdinov, Science 2006]

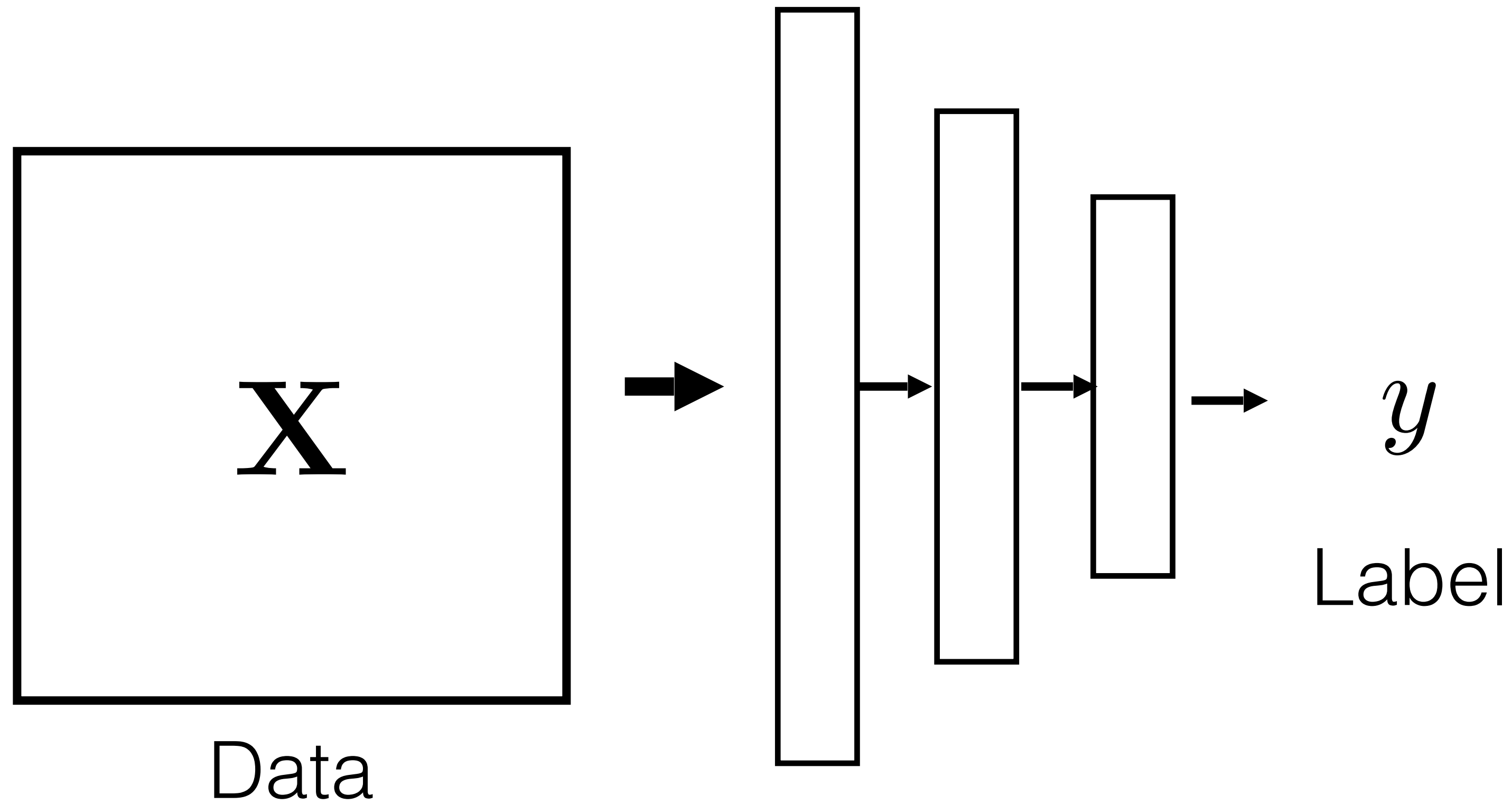
Autoencoder



Data compression

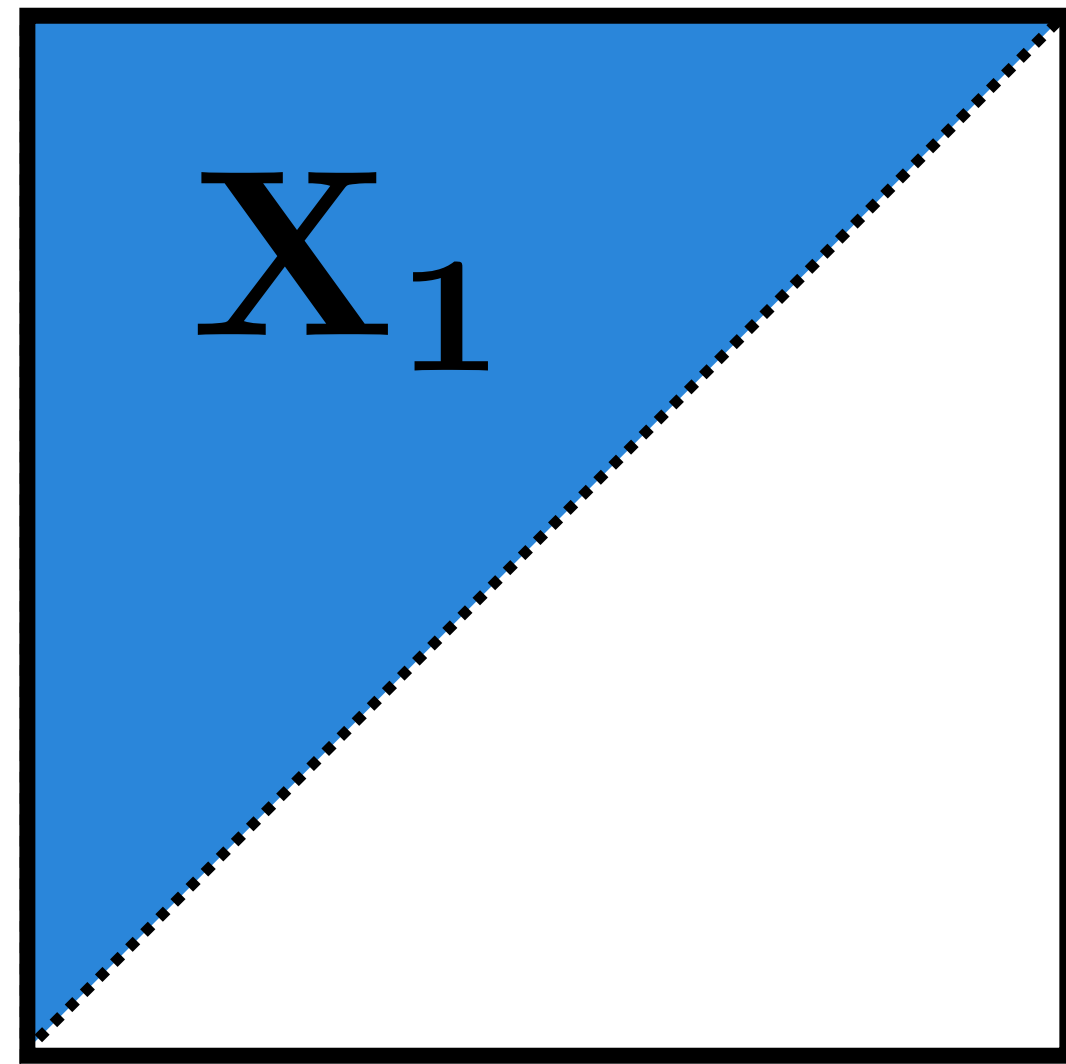


Label prediction

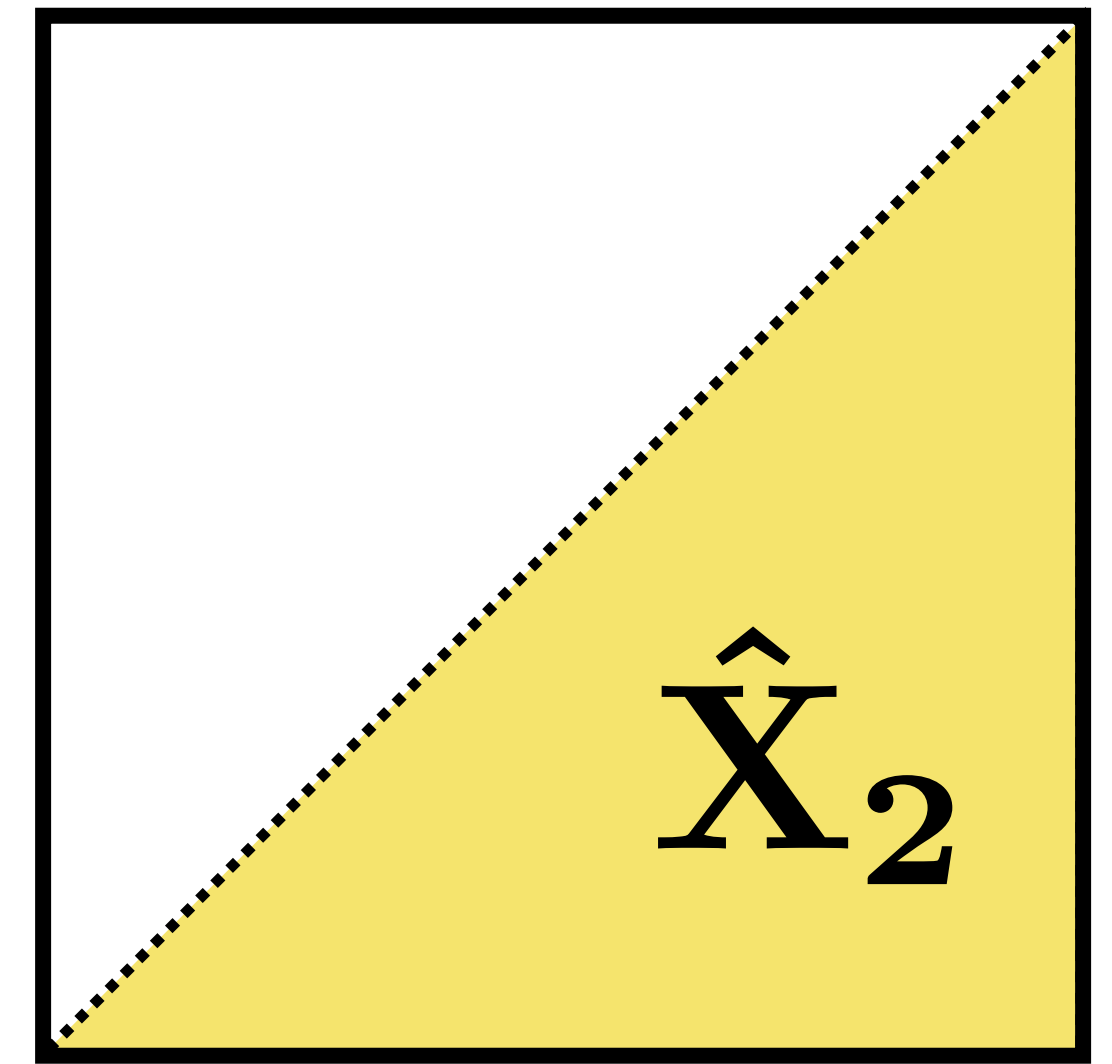
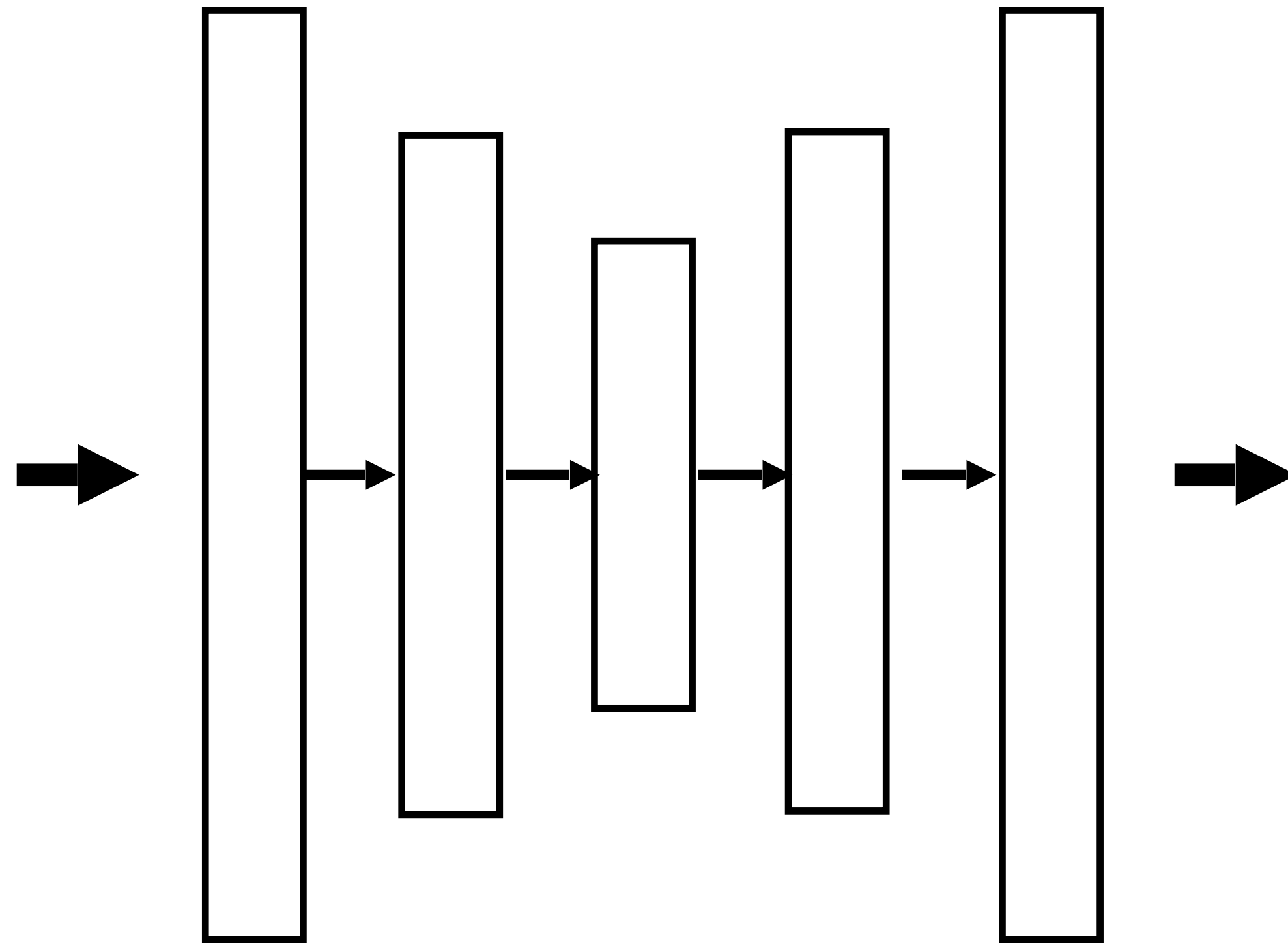


e.g., image classification

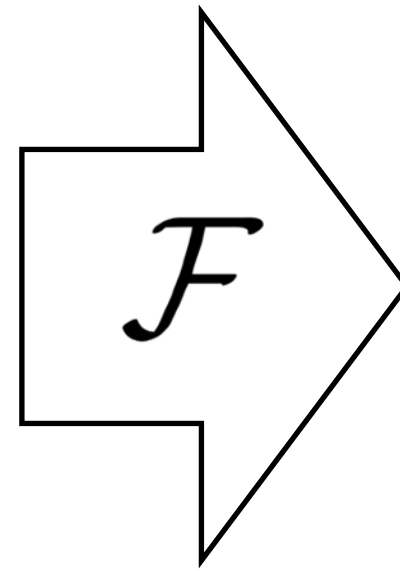
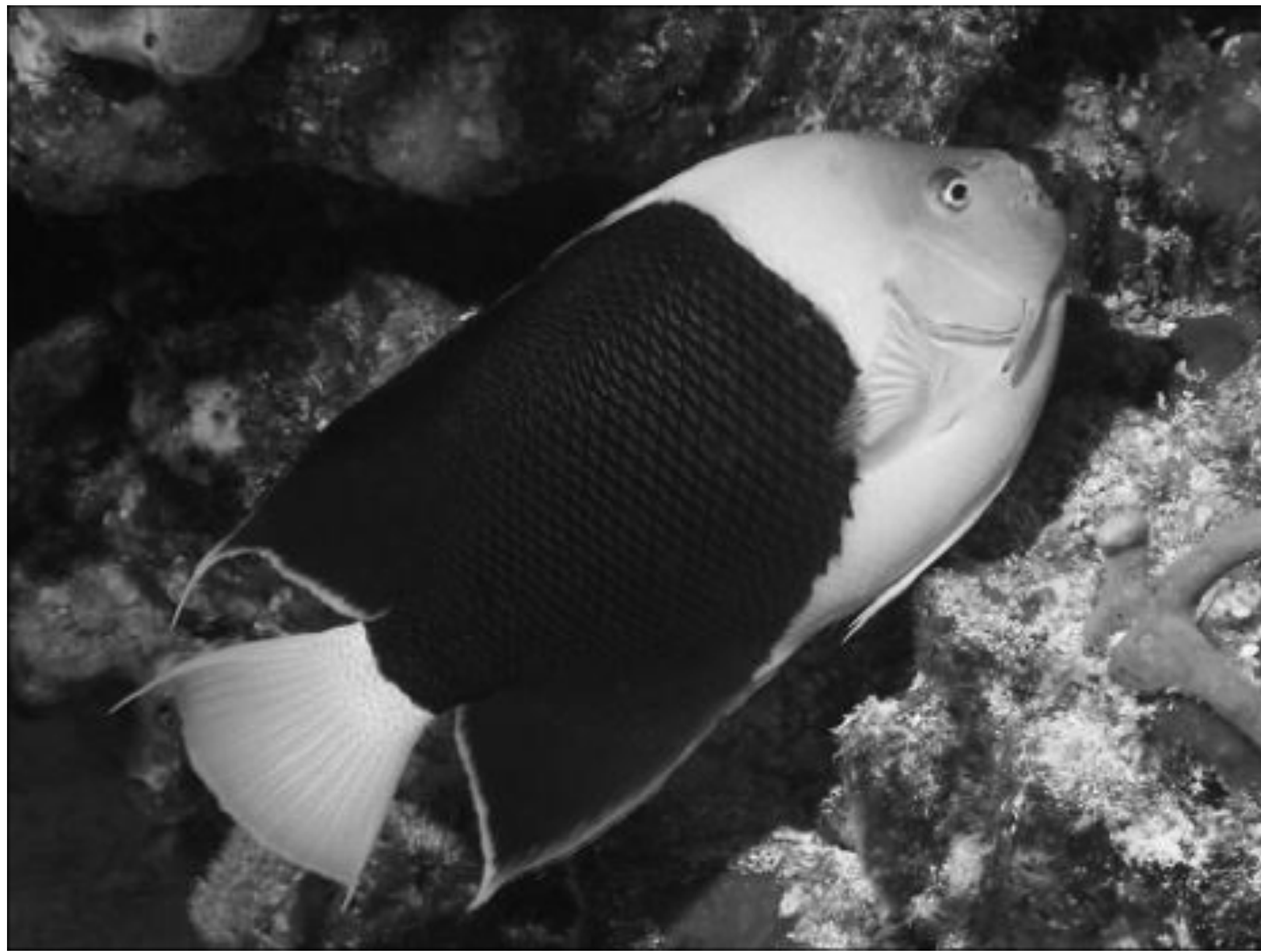
Data prediction



Some data



Other data

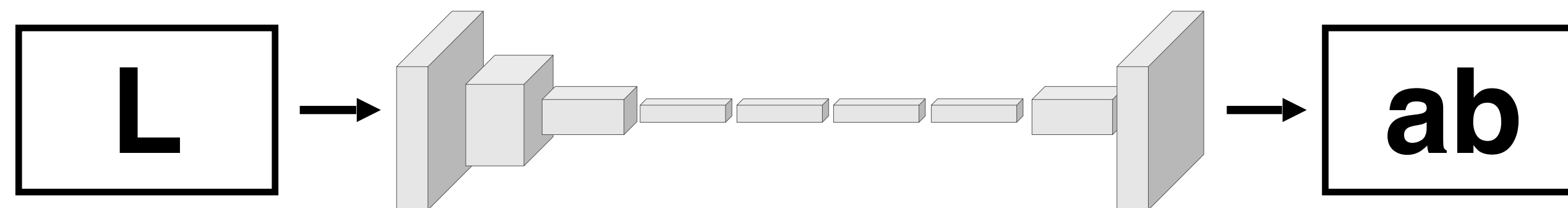


Grayscale image: L channel

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

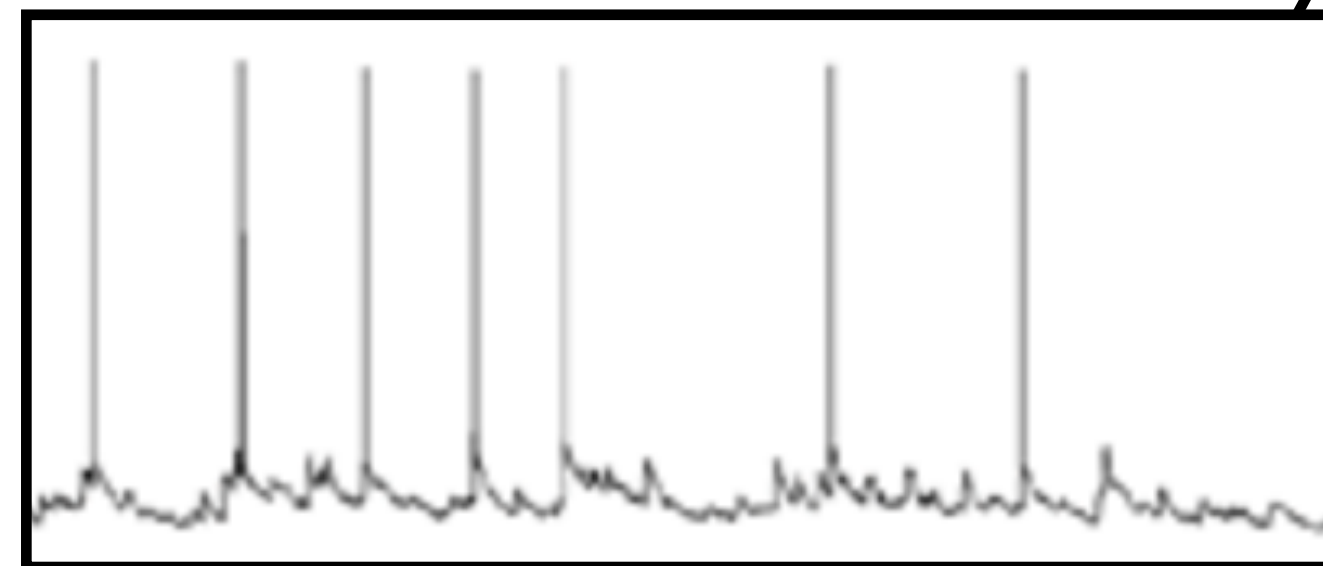
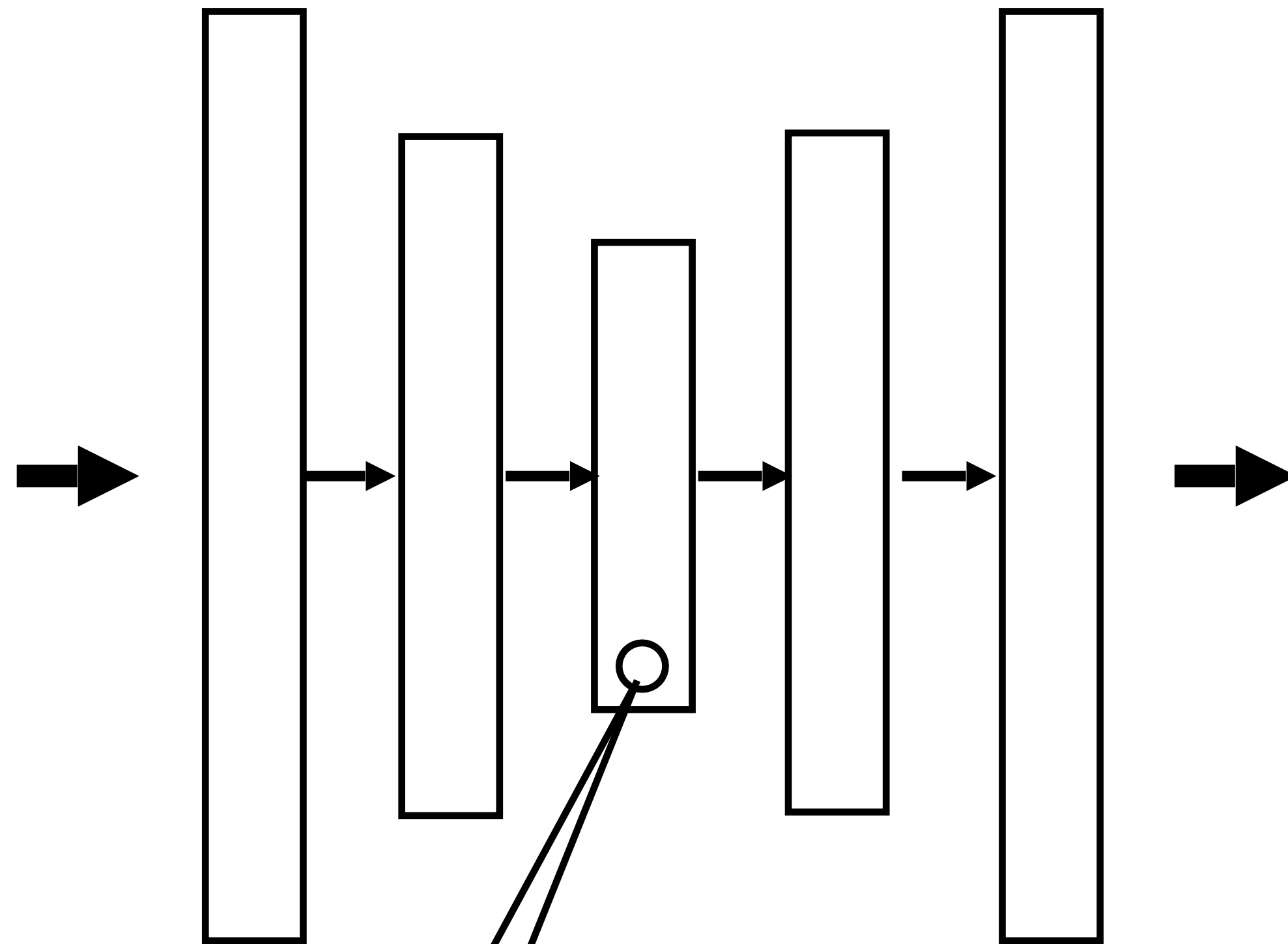
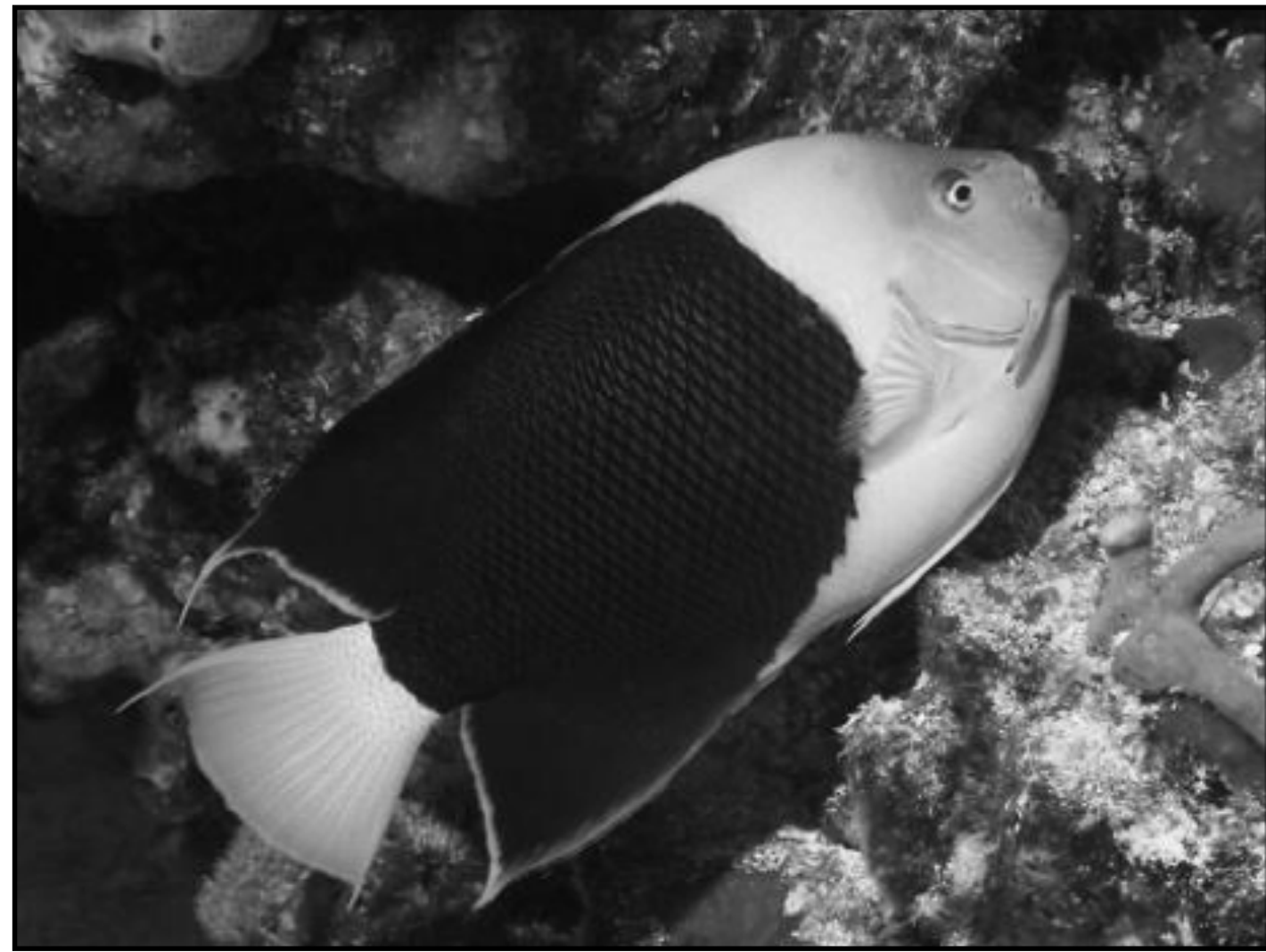
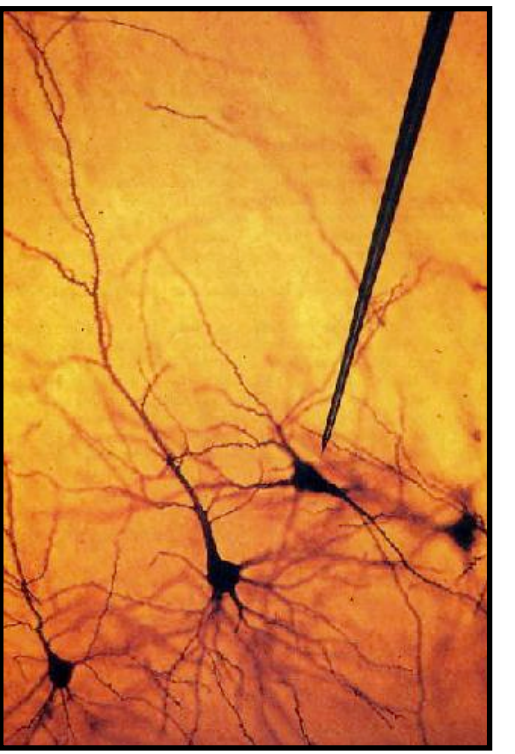
Color information: ab channels

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



[Zhang, Isola, Efros, ECCV 2016]

Deep Net “Electrophysiology”

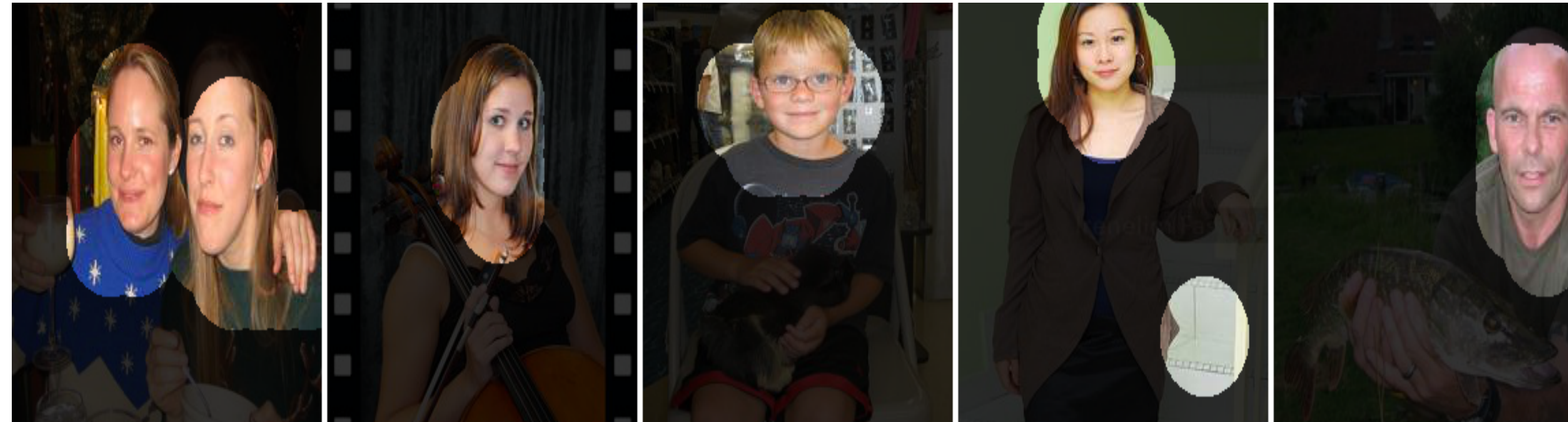


[Zeiler & Fergus, ECCV 2014]

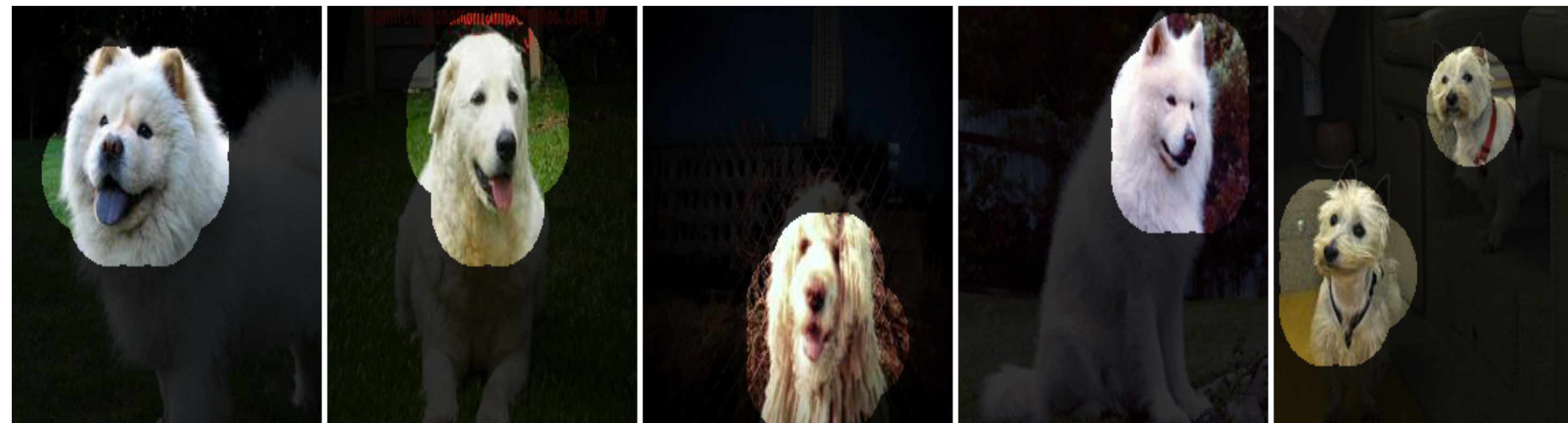
[Zhou et al., ICLR 2015]

Stimuli that drive selected neurons (conv5 layer)

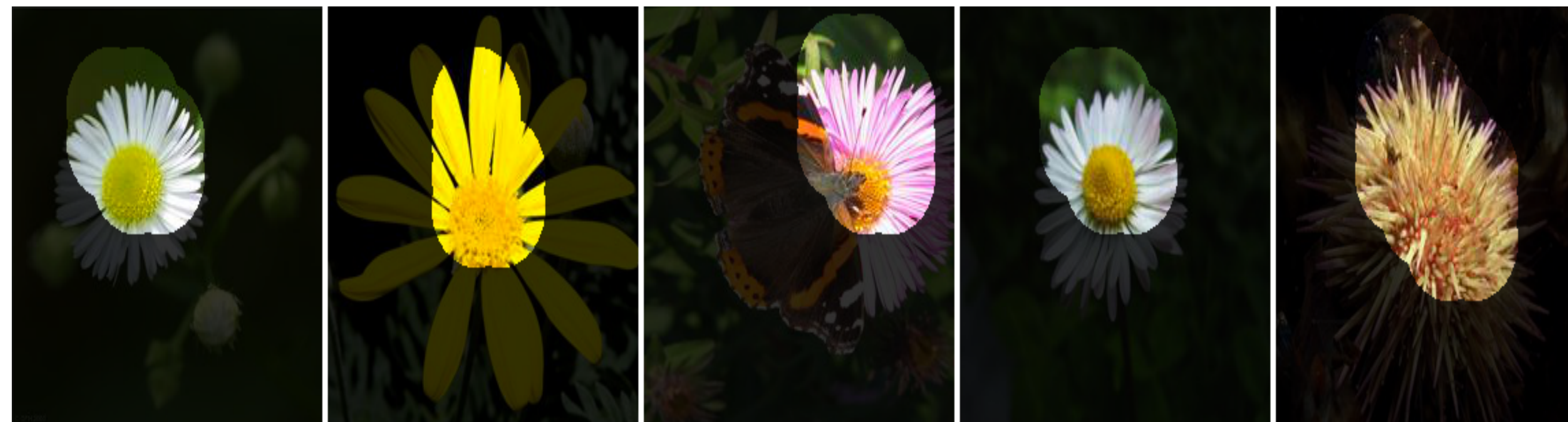
faces



dog
faces

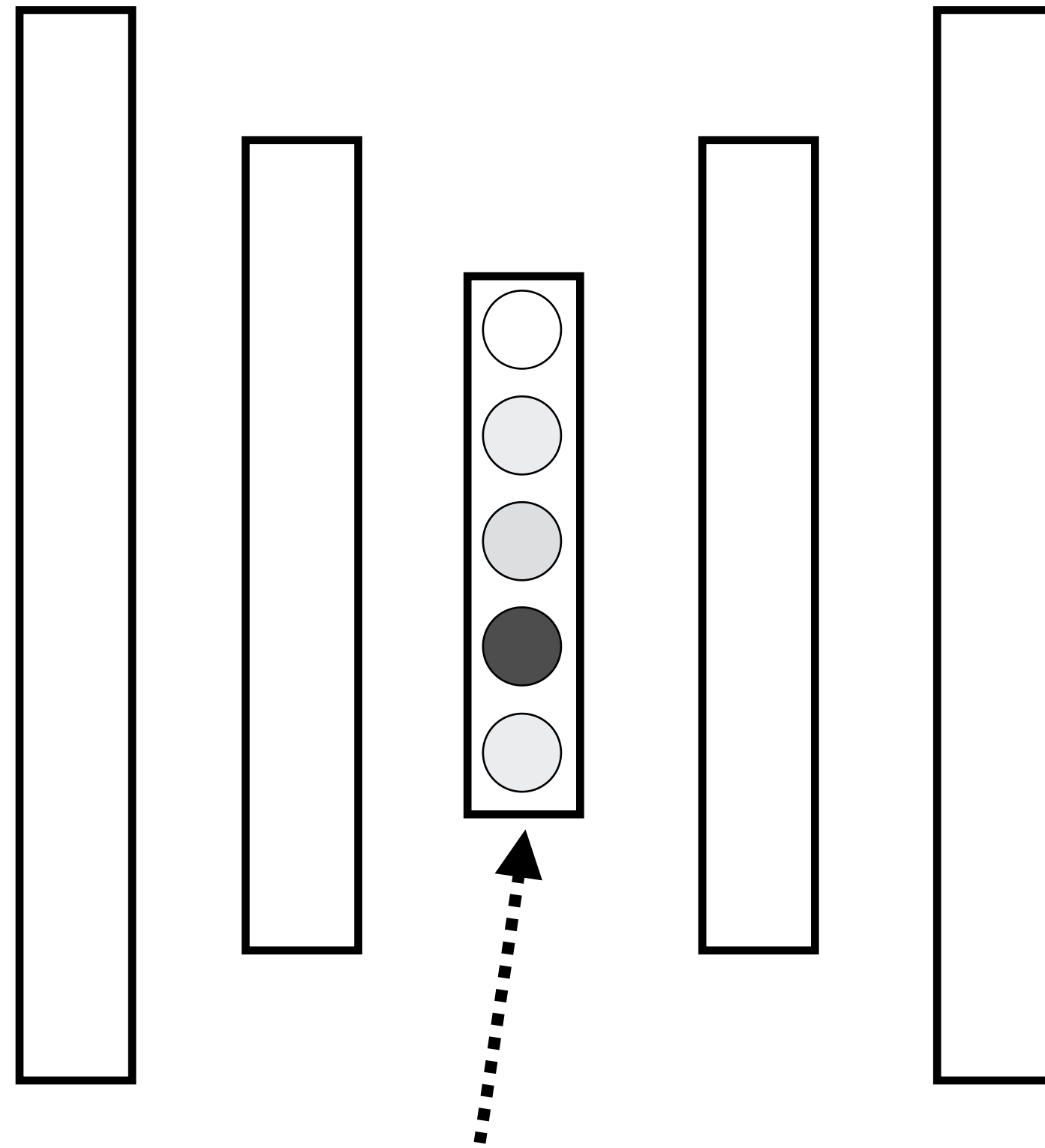


flowers



\mathbf{X} 

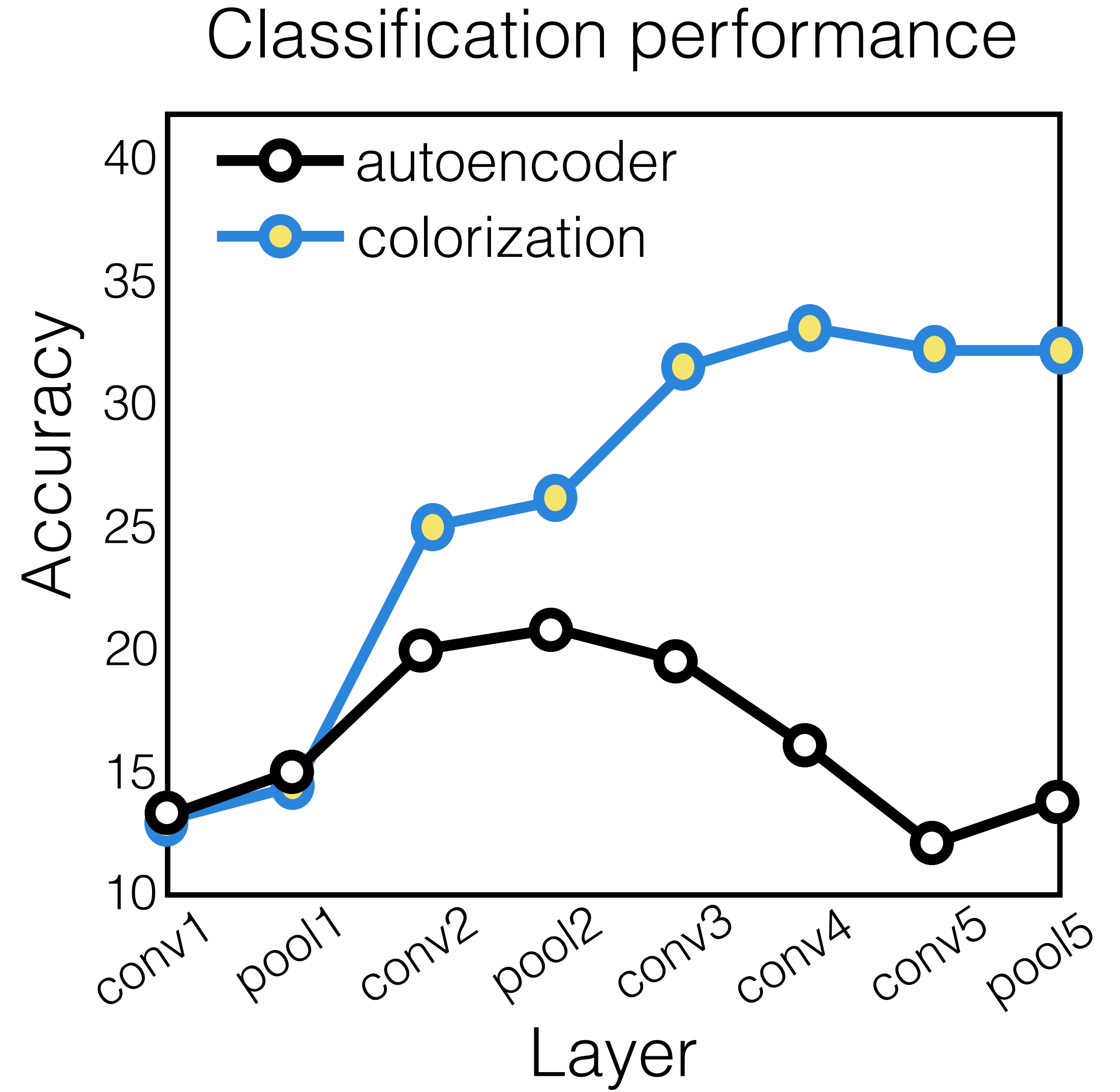
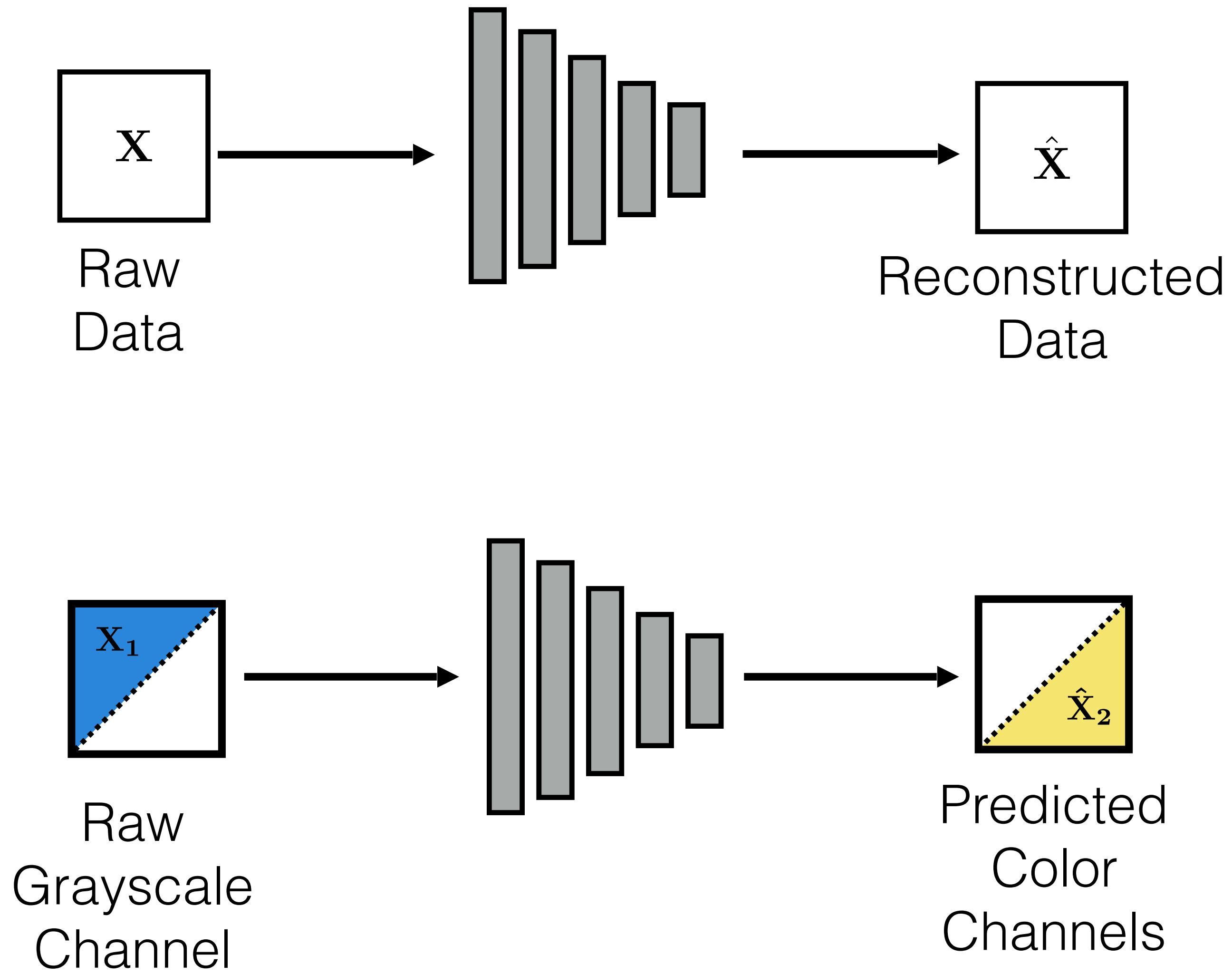
Image

compressed image code
(vector \mathbf{z}) $\hat{\mathbf{X}}$ Reconstructed
image

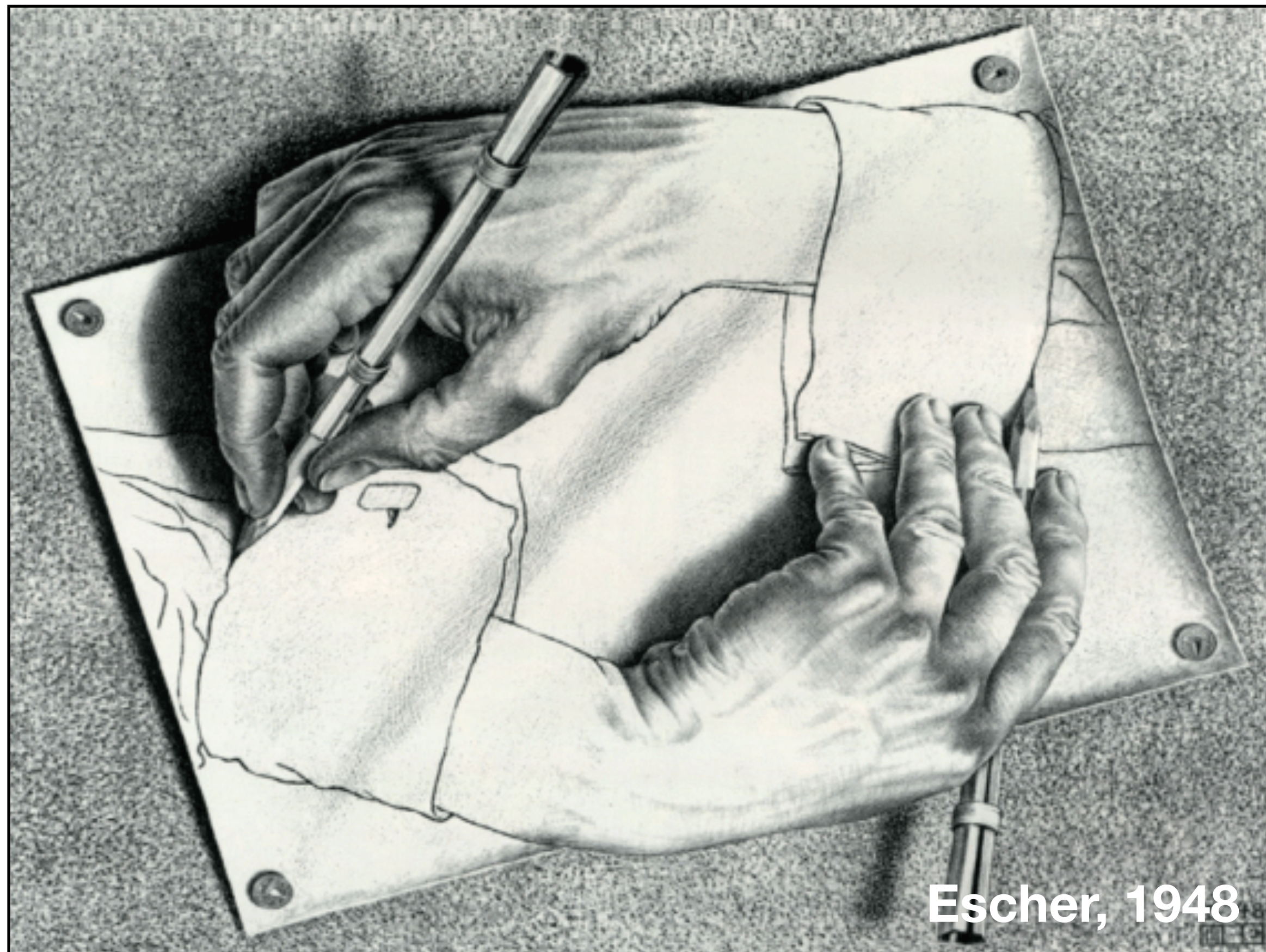
Is the code informative about
object class y ?

Logistic regression:

$$y = \sigma(\mathbf{W}\mathbf{z} + \mathbf{b})$$



Self-supervised learning



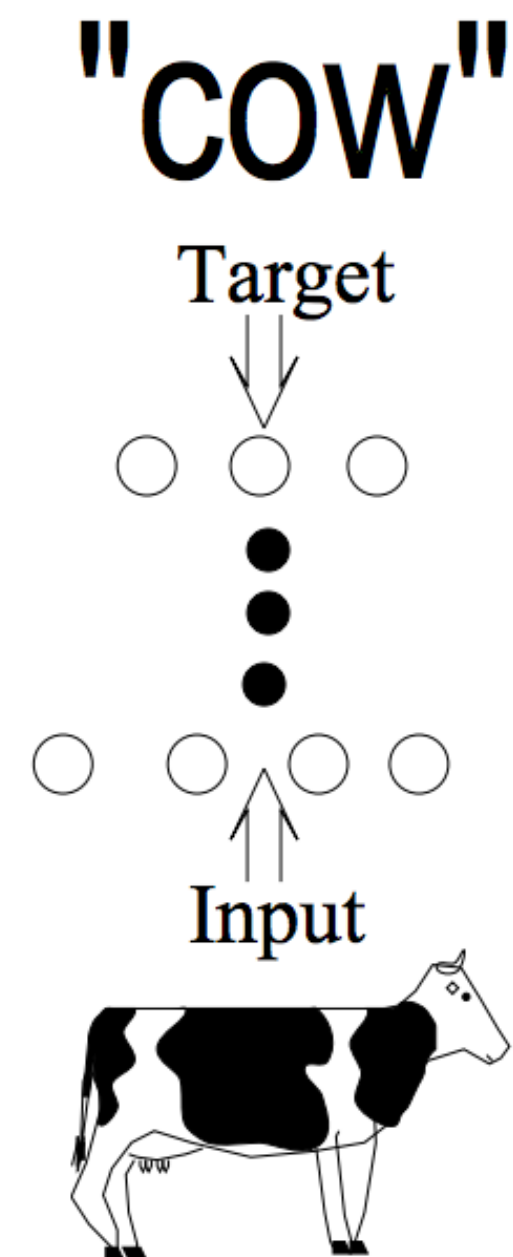
Common trick:

- Convert “unsupervised” problem into “supervised” empirical risk minimization
- Do so by cooking up “labels” (prediction targets) from the raw data itself

Multisensory self-supervision

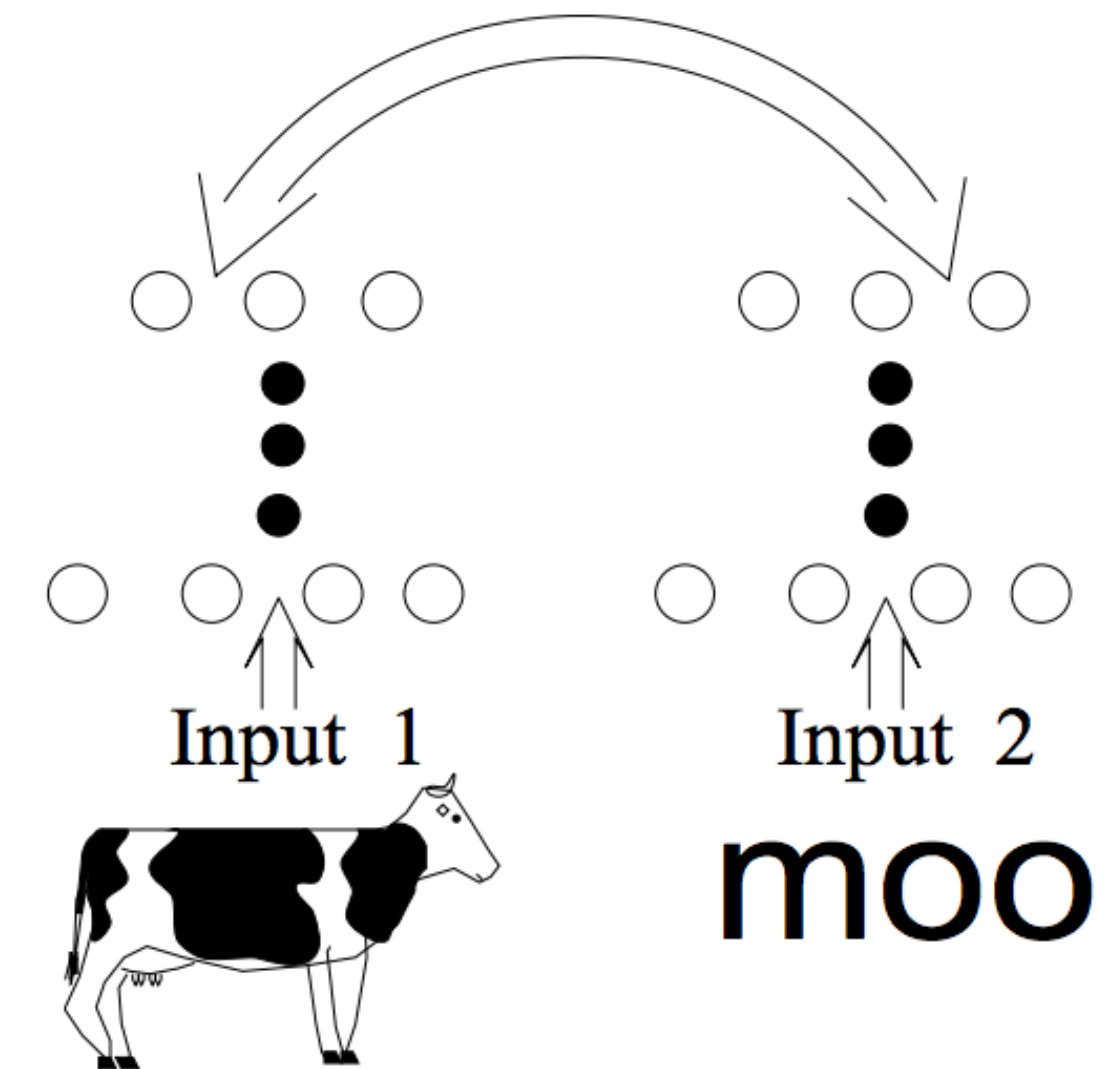
Supervised

- implausible label



Self-Supervised

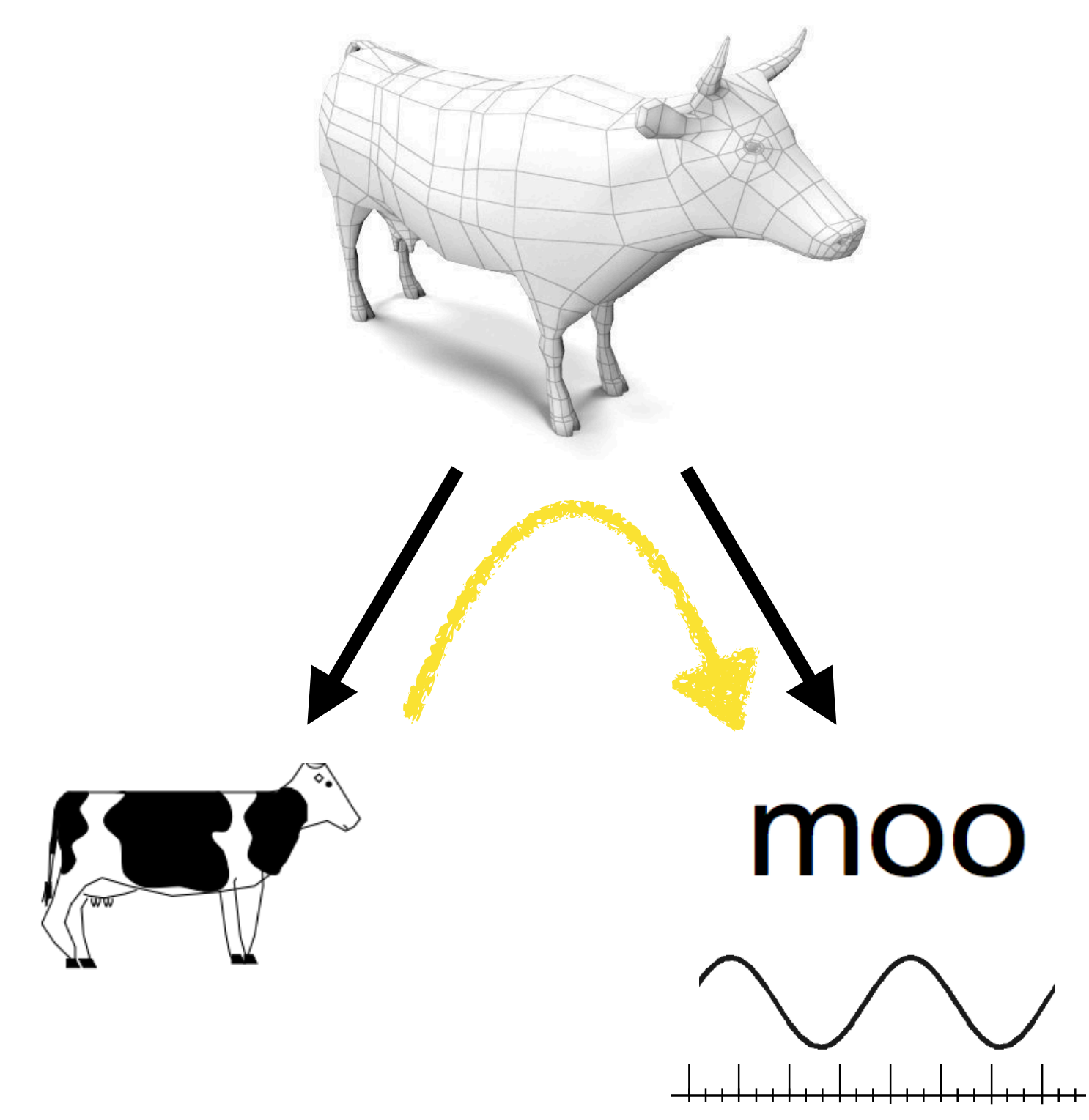
- derives label from a co-occurring input to another modality



Virginia de Sa. *Learning Classification with Unlabeled Data*. NIPS 1994.

[see also "Six lessons from babies", Smith and Gasser 2005]

The allegory of the cave



*Maxima pars hominum cecis immersa tenebris
 Volvitur assidue, et studio letatur inani:
 Adspice ut obiectis obtutis in hereat umbris,
 Ut VERI simulacra omnes mirentur amentq,*

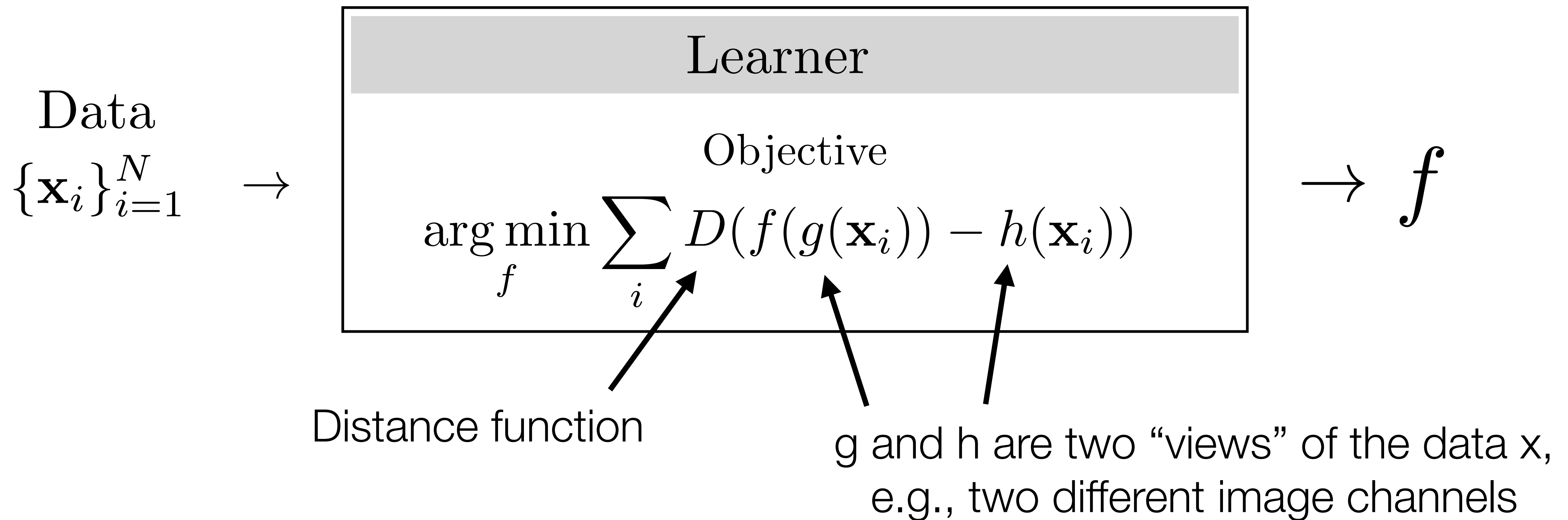
*Et solidi Dana ludantur imagine rerum.
 Quam pauci meliore luto, qui in lumine puro
 Secreti a solidi turba, ludibria cernunt
 Rerum umbras rectas, expendunt omnia lauce:*

*Hi posita erroris nebula dignoscere possunt
 Vera bona, atque alios ceca sub nocte latentes
 Extrahere in claram lucem conantur, ac illis
 Nullus amor lucis, tanta est rationis eges fas.*

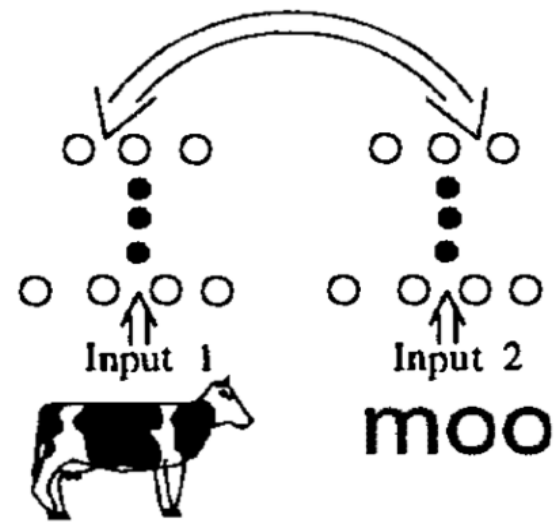
C.C. Harlemensis Inv.
 Janredam Sculpsit.
 Henr. Hondius excudit.
 1604.

H.L. SPIEGEL FIGVRARI ET SCVLPI CVRAVIT. AC DOCTISS. ORNATISS. QZDPET, PAAW IN LVGDVN. ACAD. PROFESSORI MEDICO D.D.

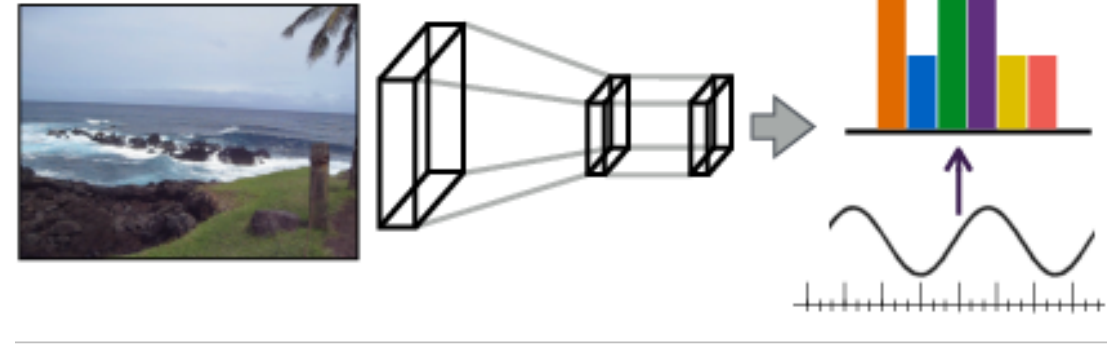
“Multiview” self-supervised learning



Audio

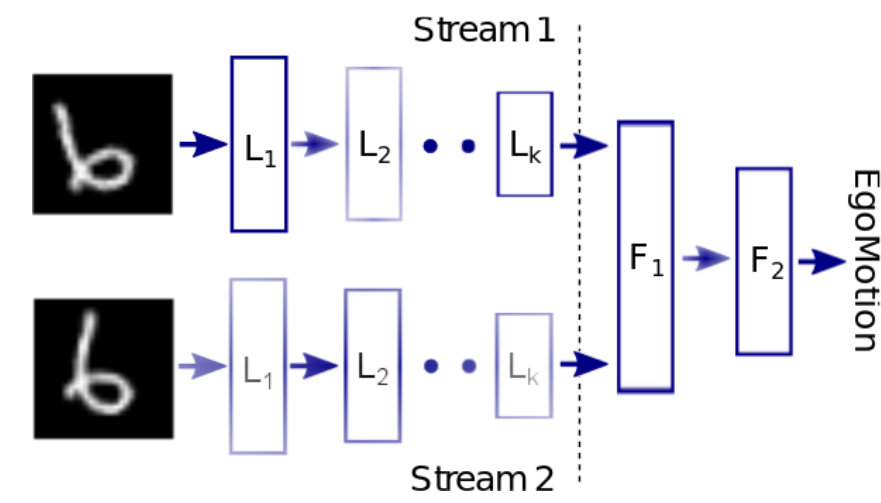


de Sa. NIPS 1994.

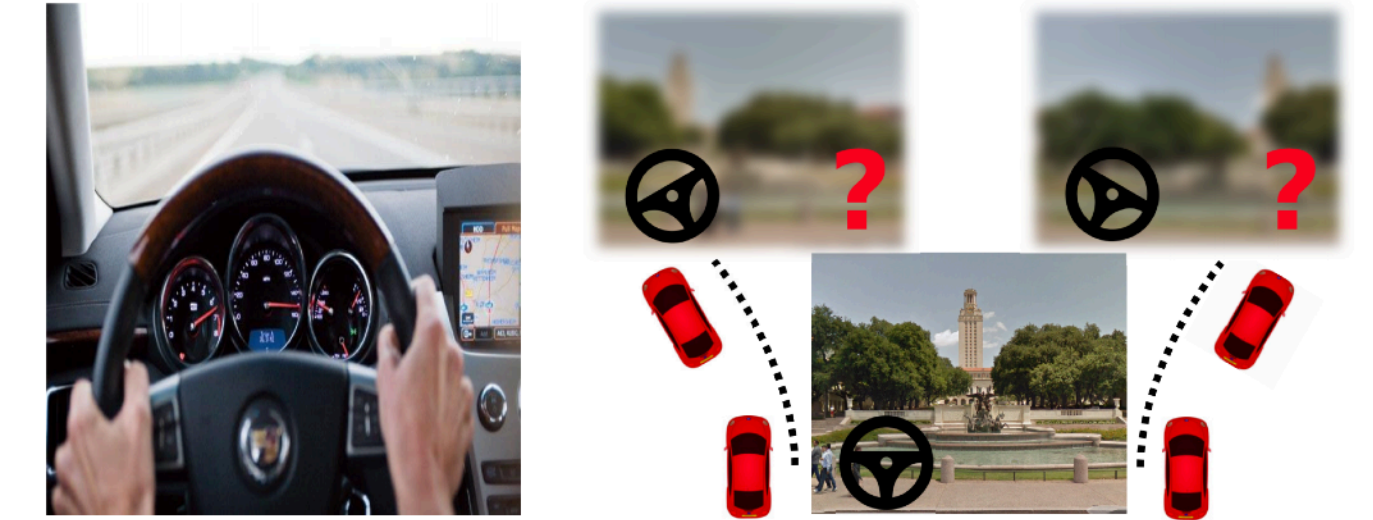


Owens et al. ECCV 2016.

Egomotion



Agrawal et al. ICCV 2015.

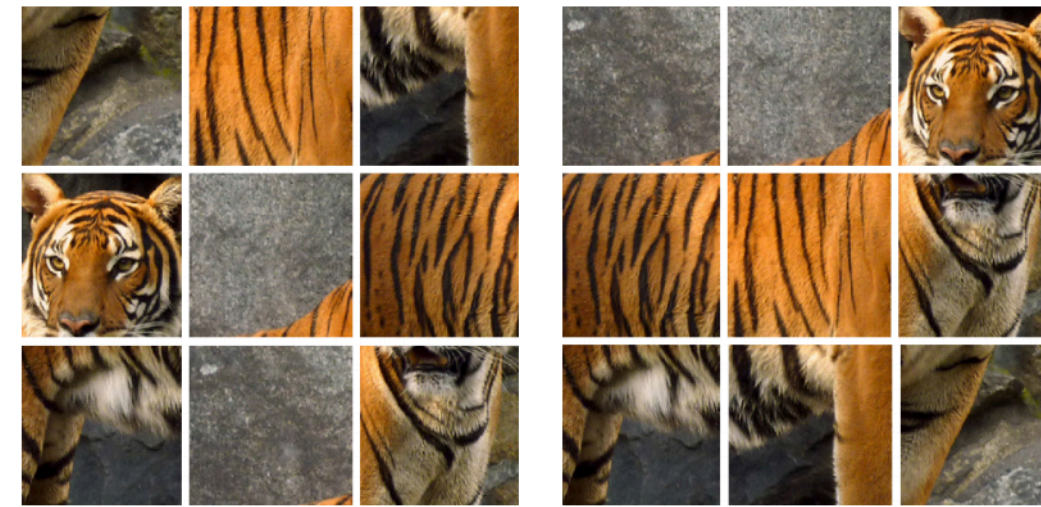


Jayaraman et al. ICCV 2015.

Context

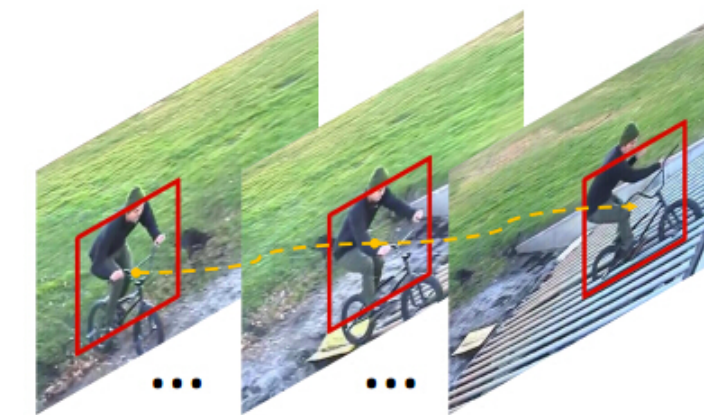


Pathak et al. CVPR 2016.

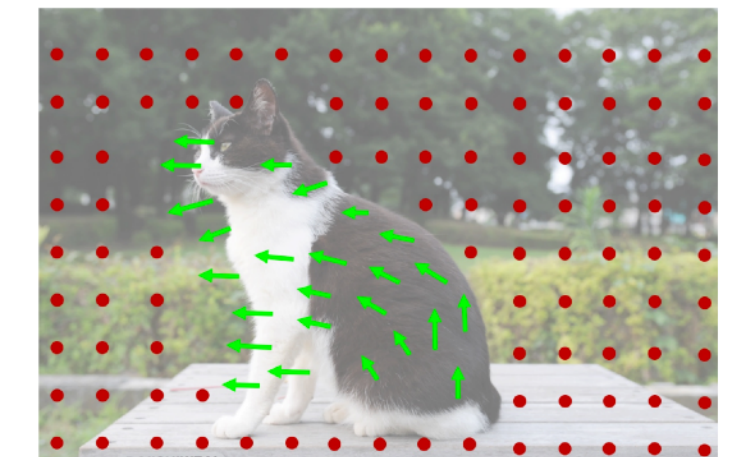


Noroozi and Favaro. ECCV 2016.
Doersch et al. ICCV 2015.

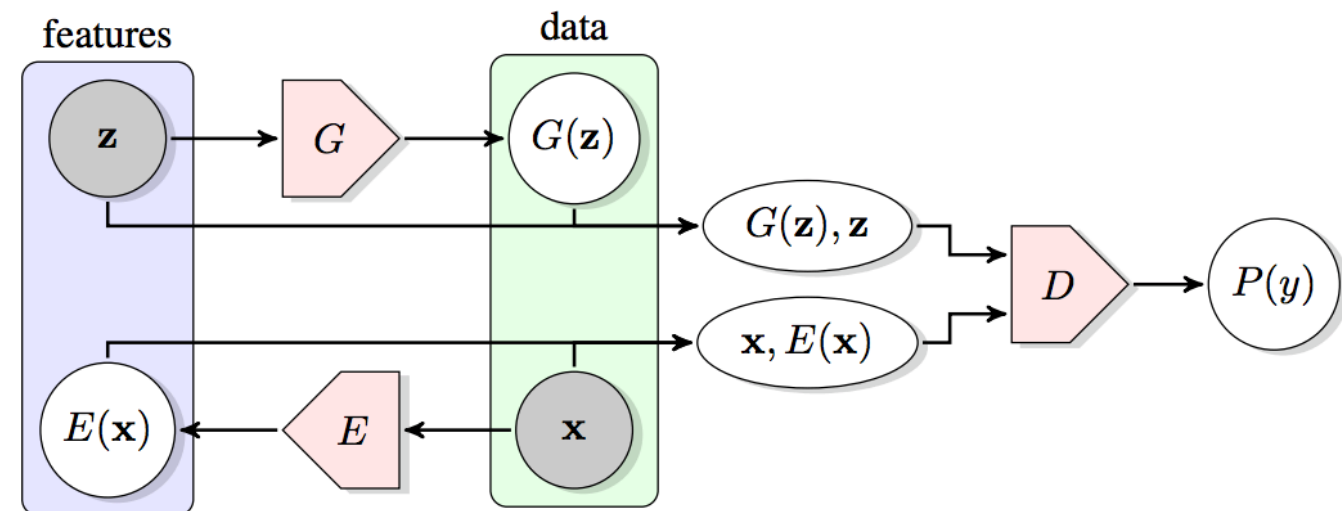
Video



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

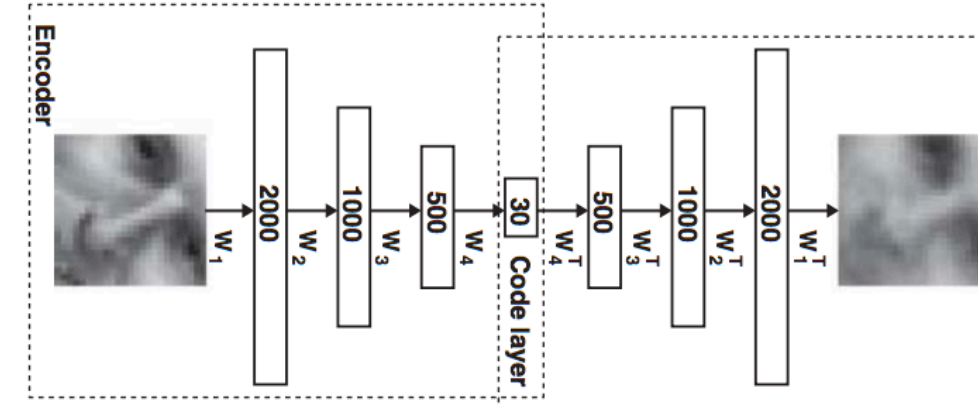


Generative Modeling



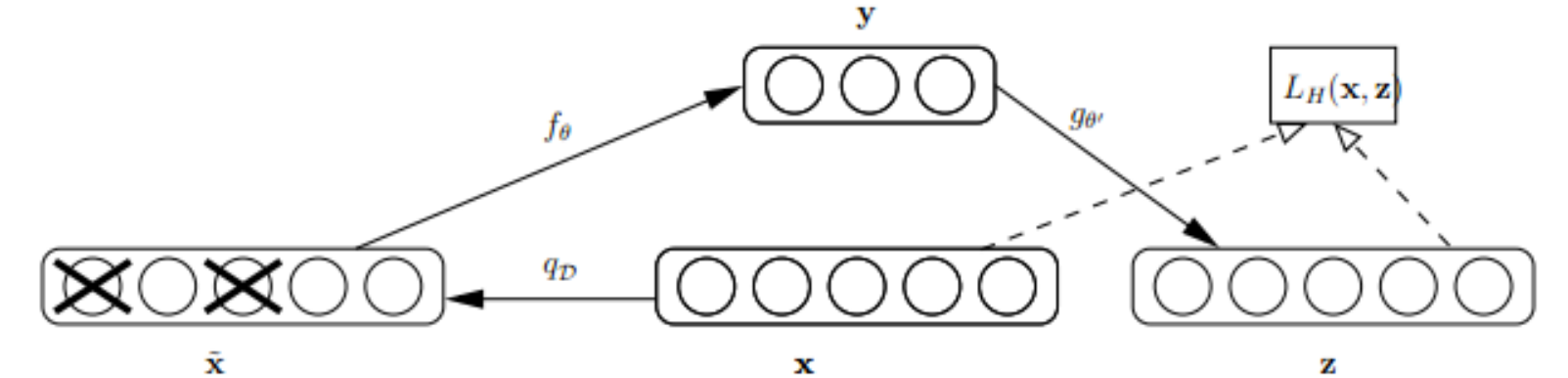
Donahue et al. Dumoulin et al. ICLR 2017.

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]

Unsupervised visual representation learning by context prediction

[Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015]

Context as Supervision

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

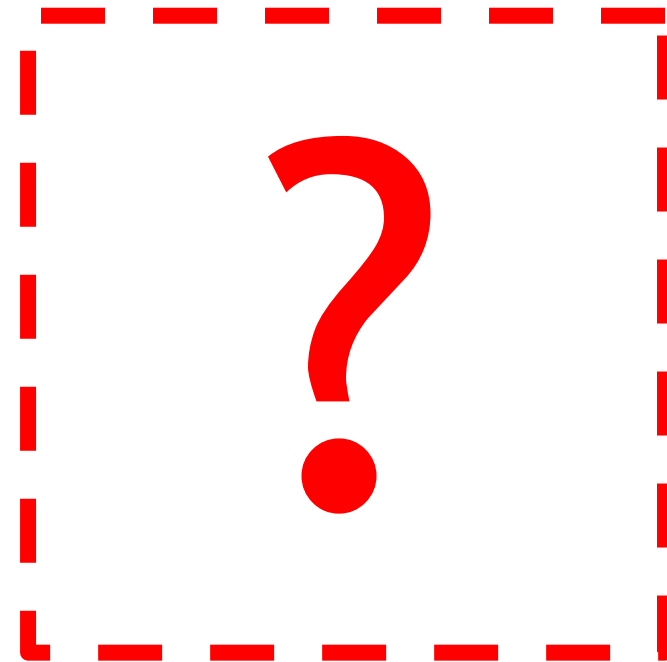
house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high 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

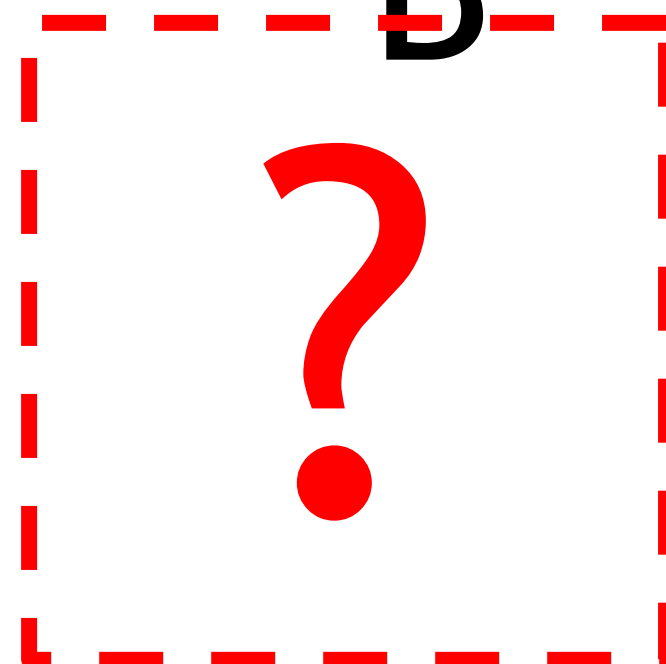
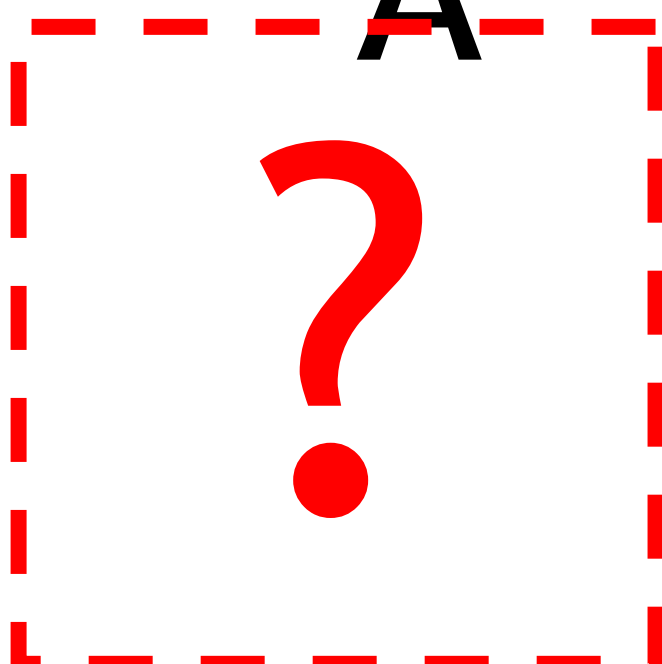
[Slide credit: Carl Doersch]

Context Prediction for Images

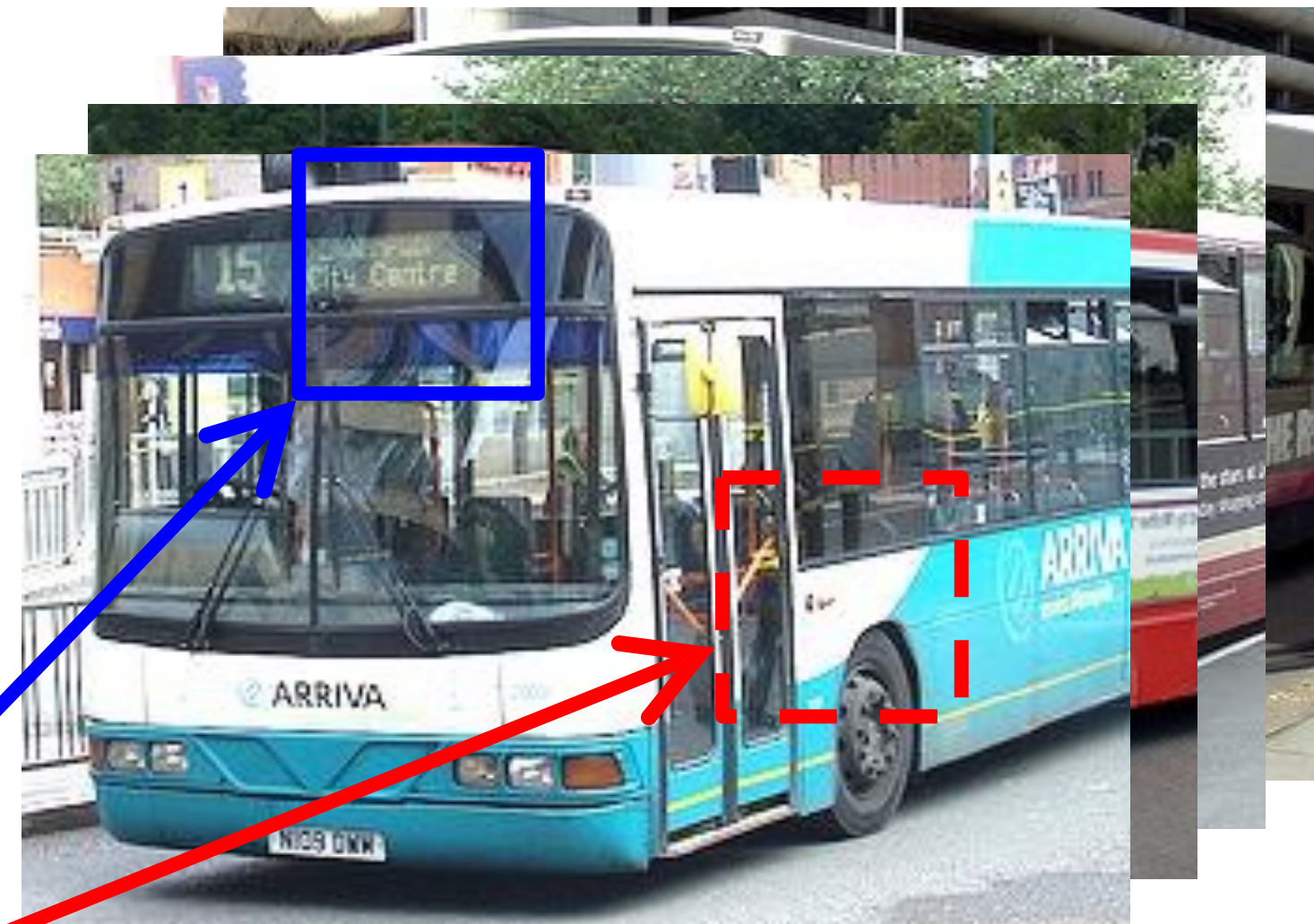


A

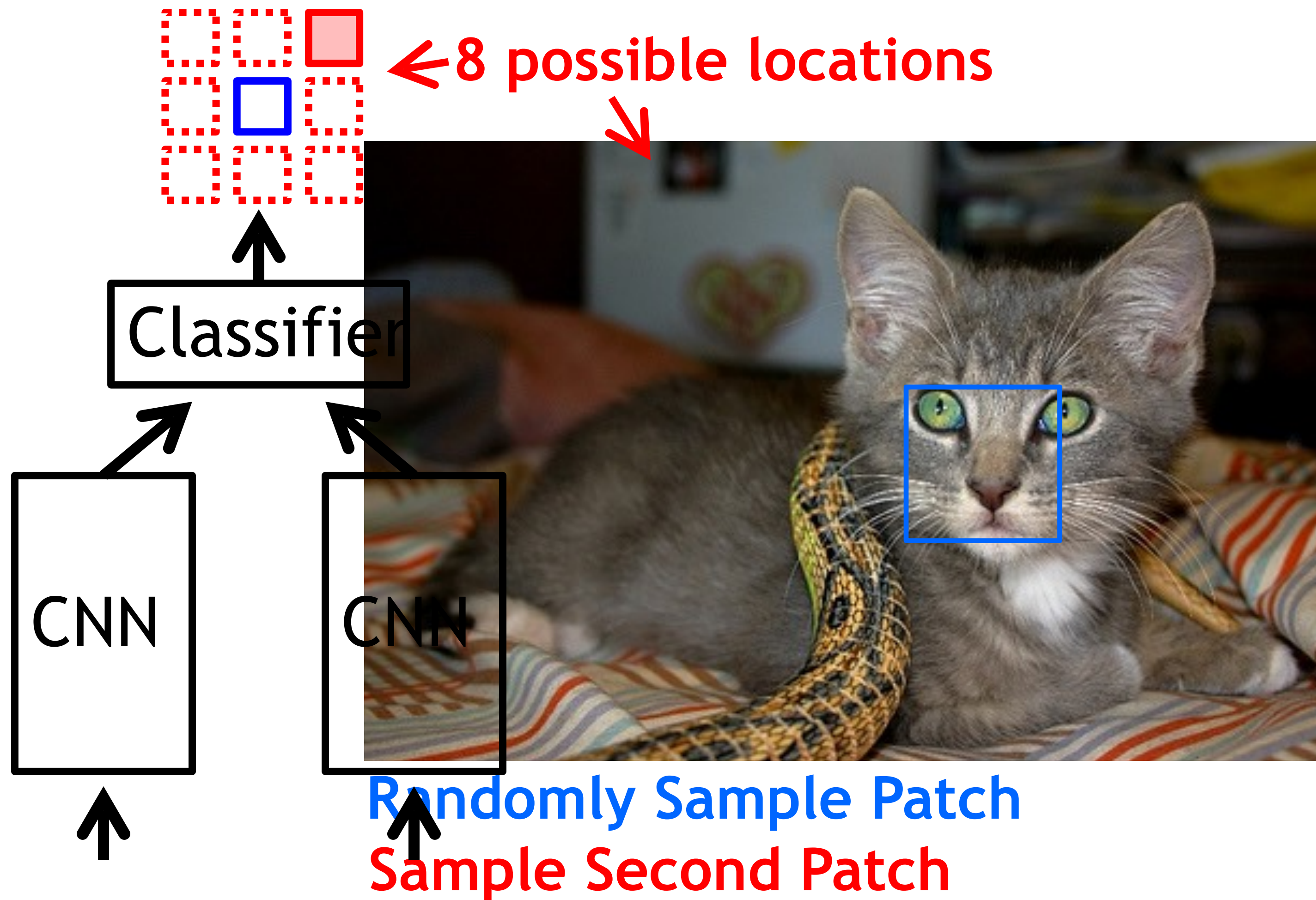
B

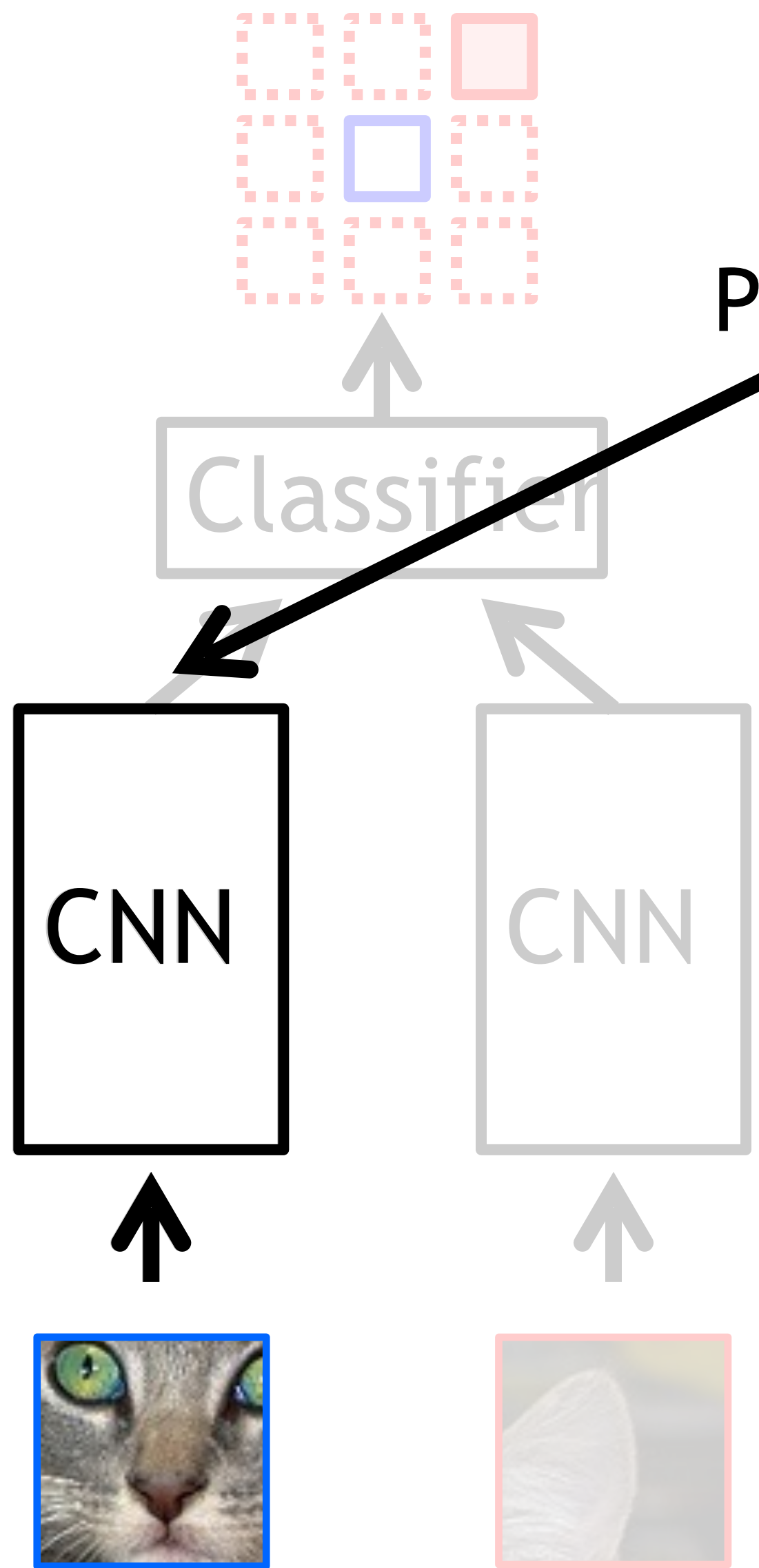


Semantics from a non-semantic task



Relative Position Task





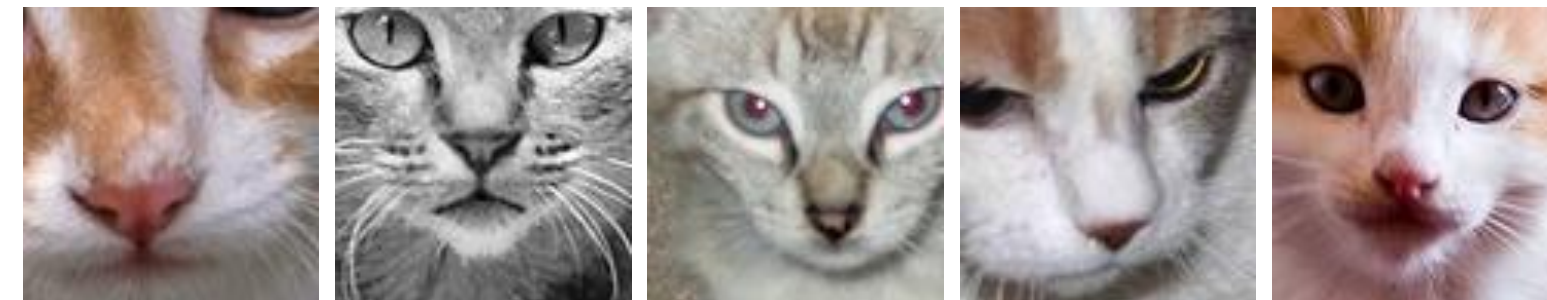
Patch Embedding (representation)

Input



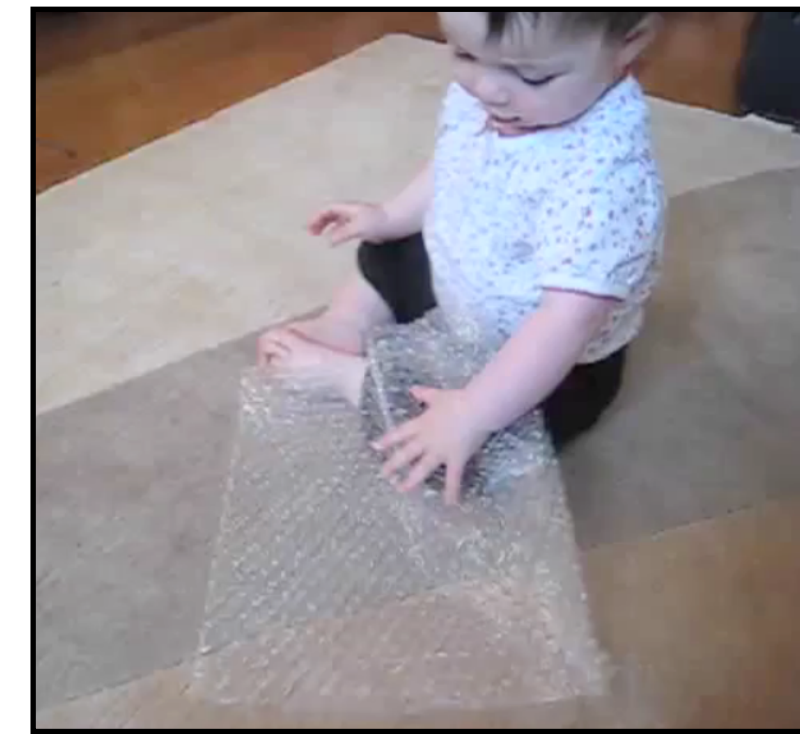
!

Nearest Neighbors



Note: connects *across* instances!

Prediction hypothesis



Henri Cartier-Bresson

1. To survive, biological agents are constantly trying to anticipate, to predict sensations
2. This trains up representations useful for prediction — surfaces, objects, events!

Summary

1. Deep nets learn *representations*, just like our brains do
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, which is great since labels are expensive and limiting
4. Without labels there are many ways to learn representations. We saw:
 1. representations as compressed codes
 2. representations that are shared across sensory modalities
 3. representations that are predictive of the future



Yann LeCun's cake:

1. Cake is unsupervised representation learning
2. Frosting is supervised transfer learning
3. Cherry on top is reinforcement learning (model-based RL)