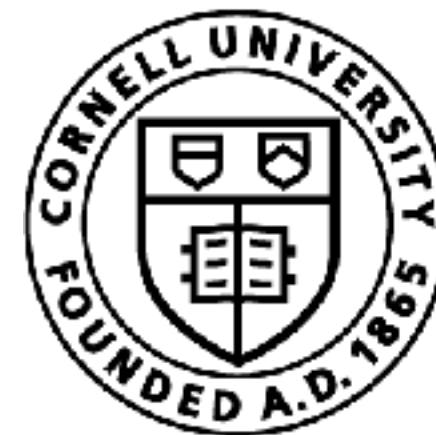
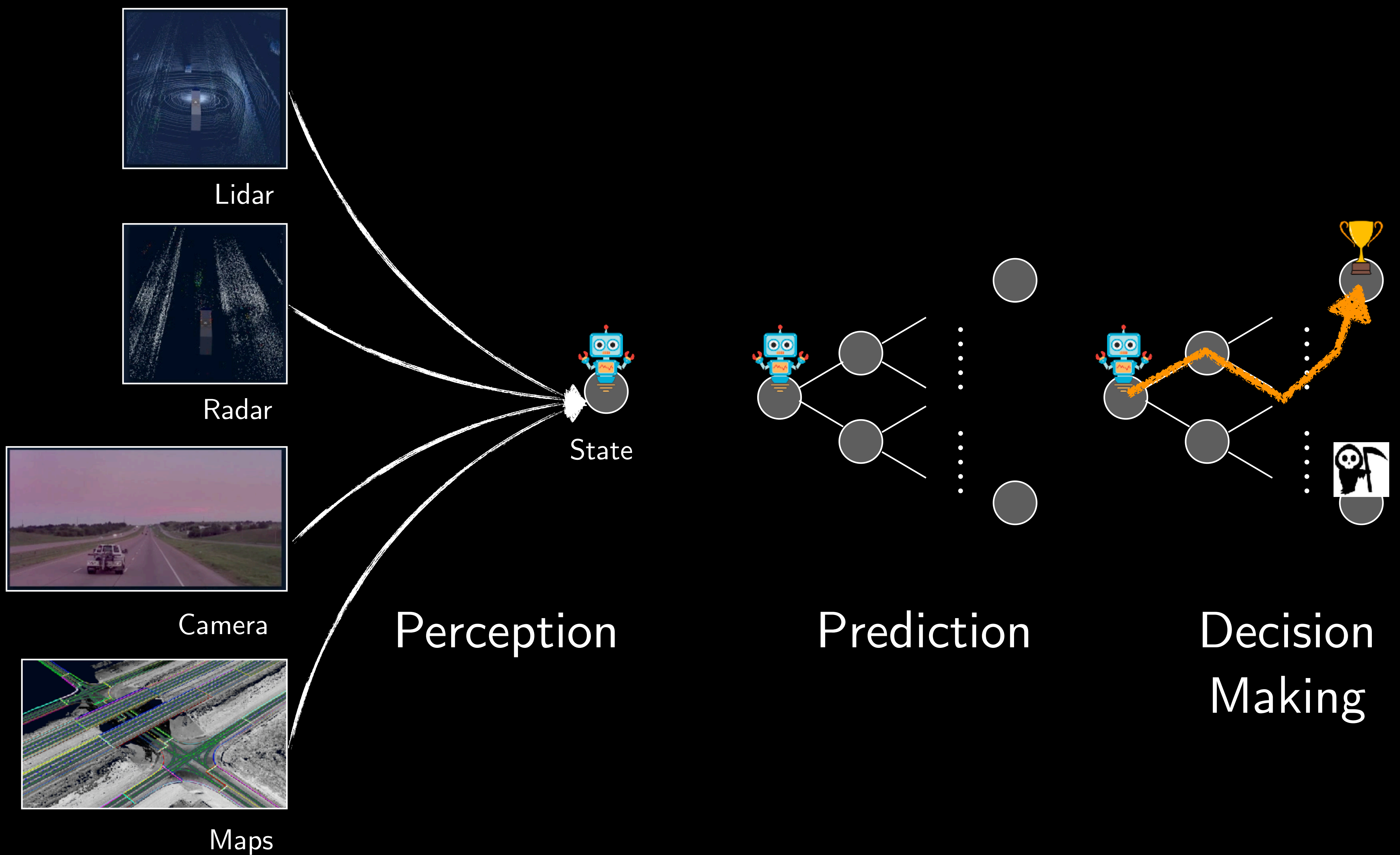


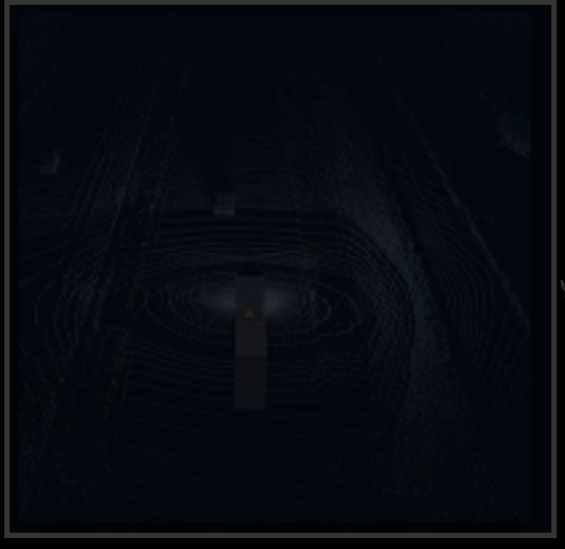
# Visual Representation Learning

Sanjiban Choudhury



Cornell Bowers CIS  
**Computer Science**





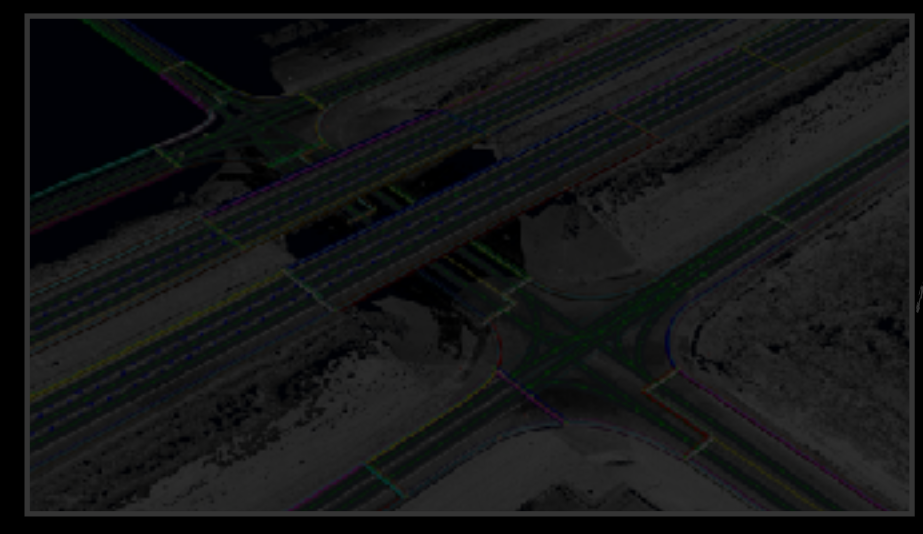
Lidar



Radar

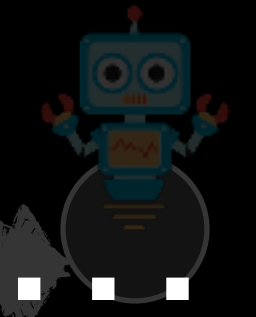


Camera



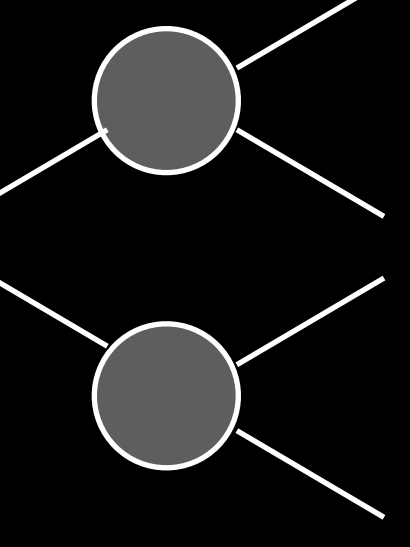
Maps

# Past classes

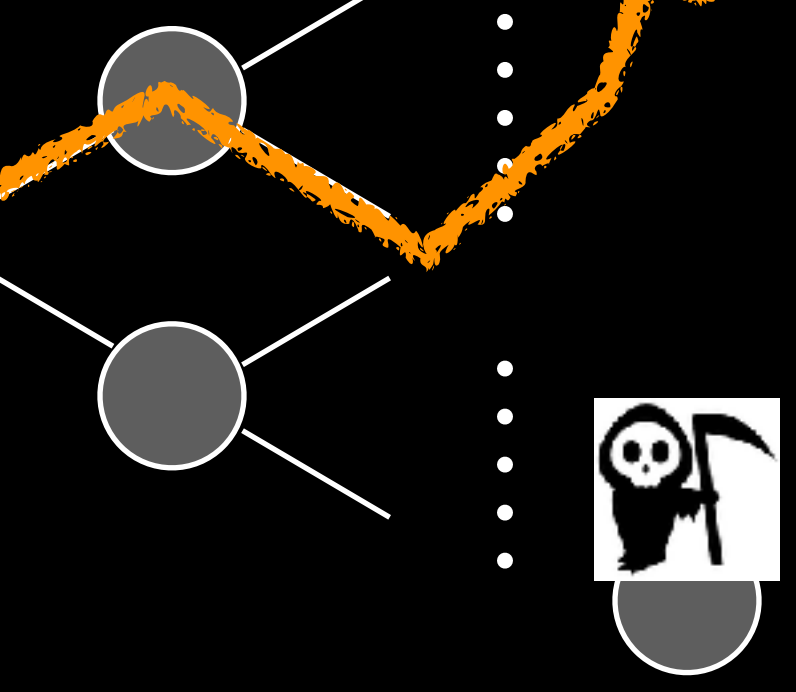
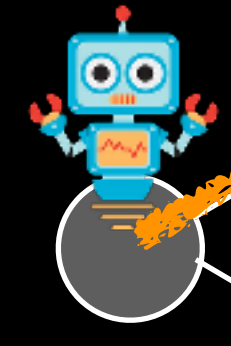


State

## Perception



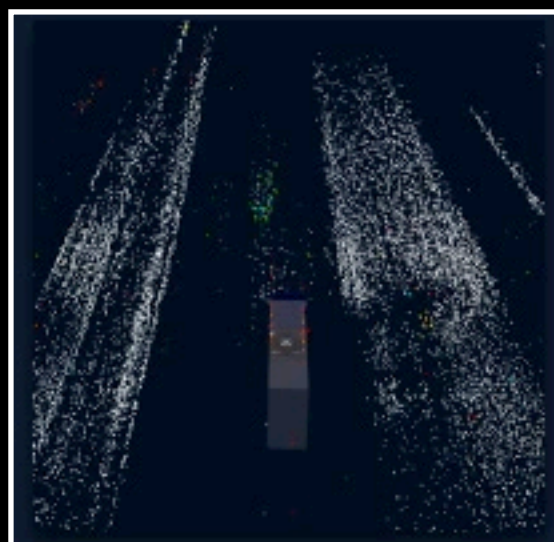
## Prediction



## Decision Making



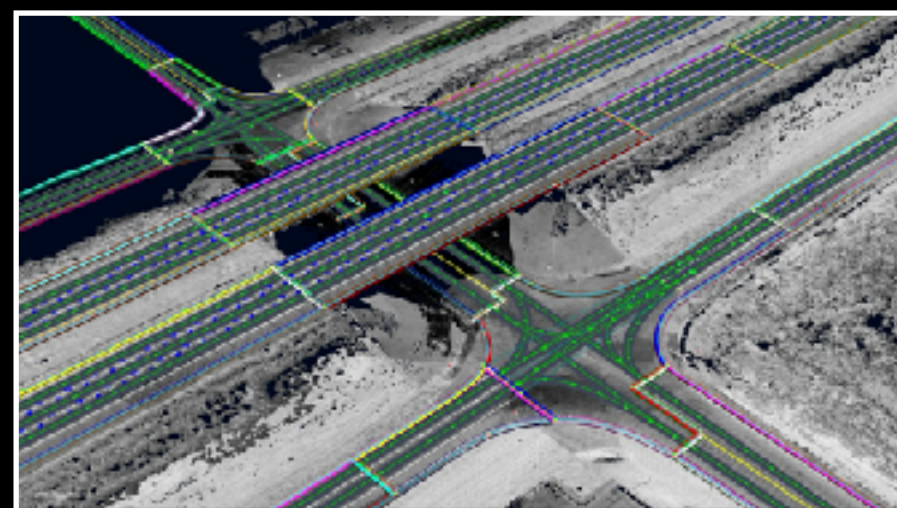
Lidar



Radar



Camera



Maps



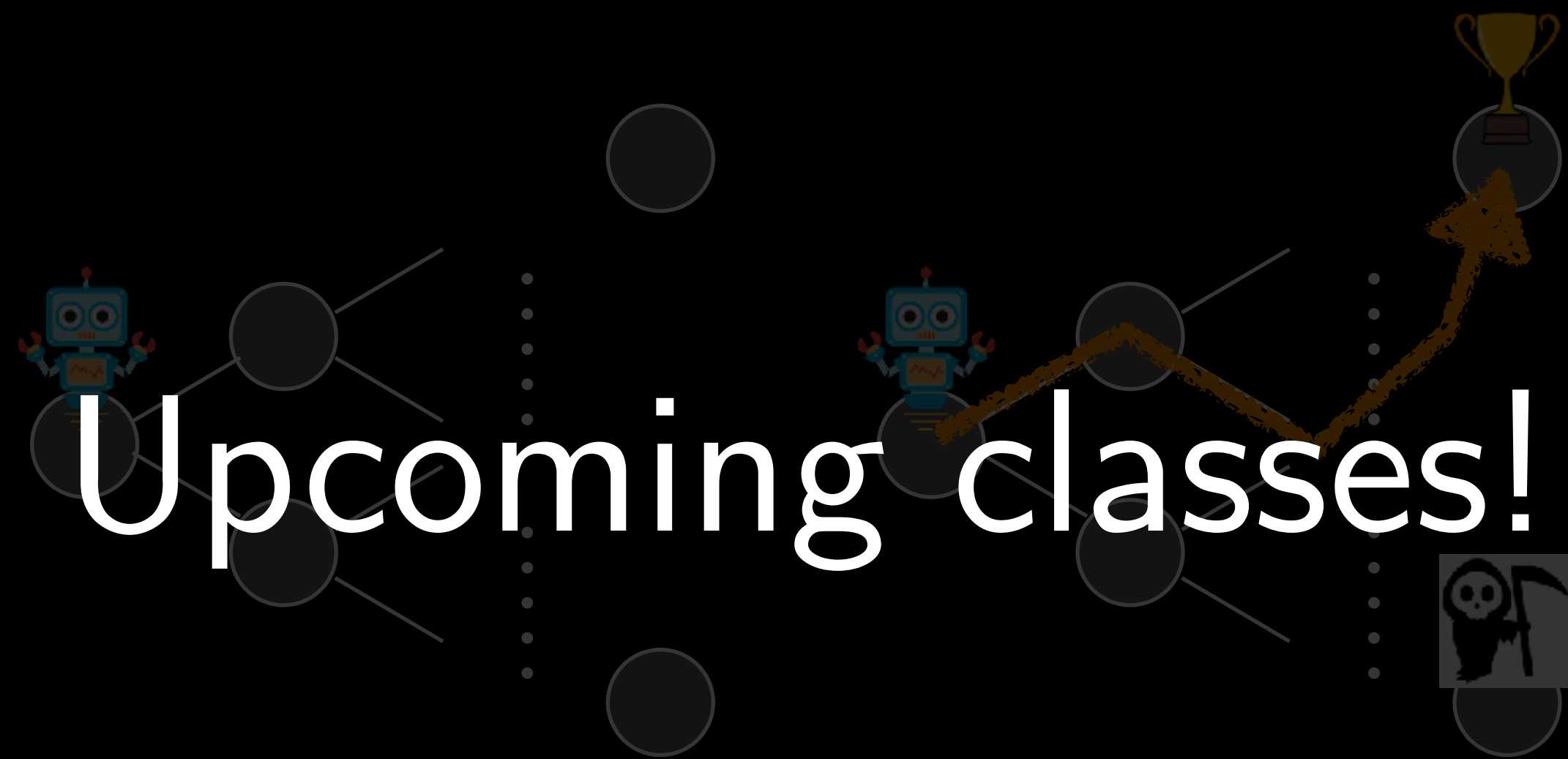
State

Perception

# Upcoming classes!

Prediction

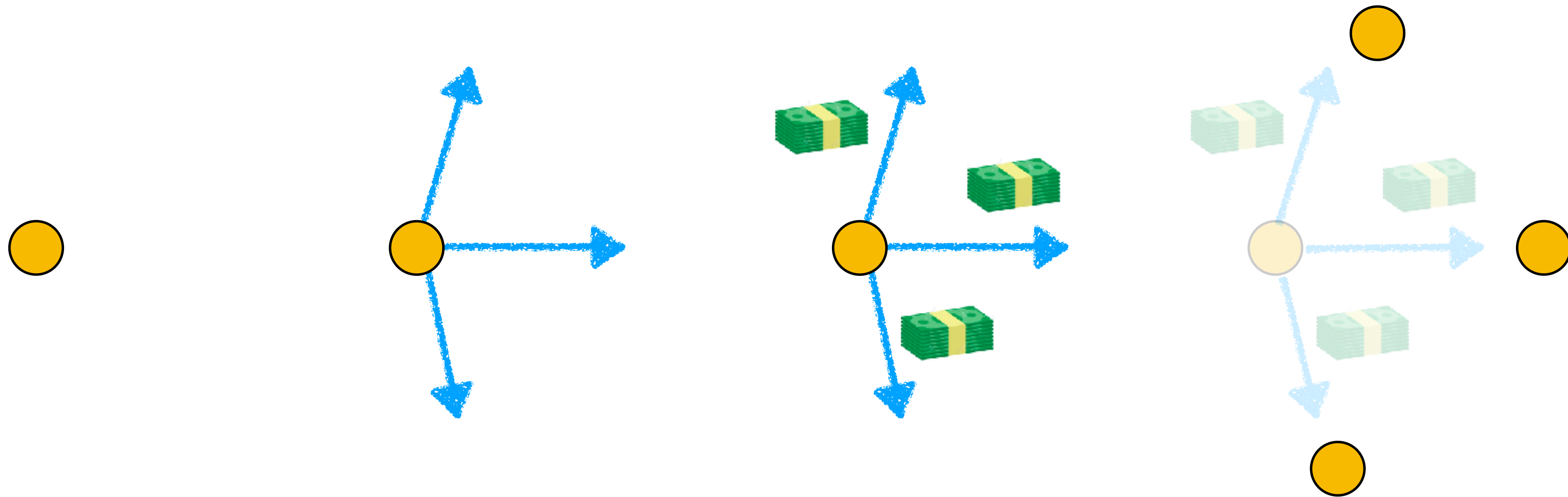
Decision Making



# Markov Decision Process

*A mathematical framework for modeling sequential decision making*

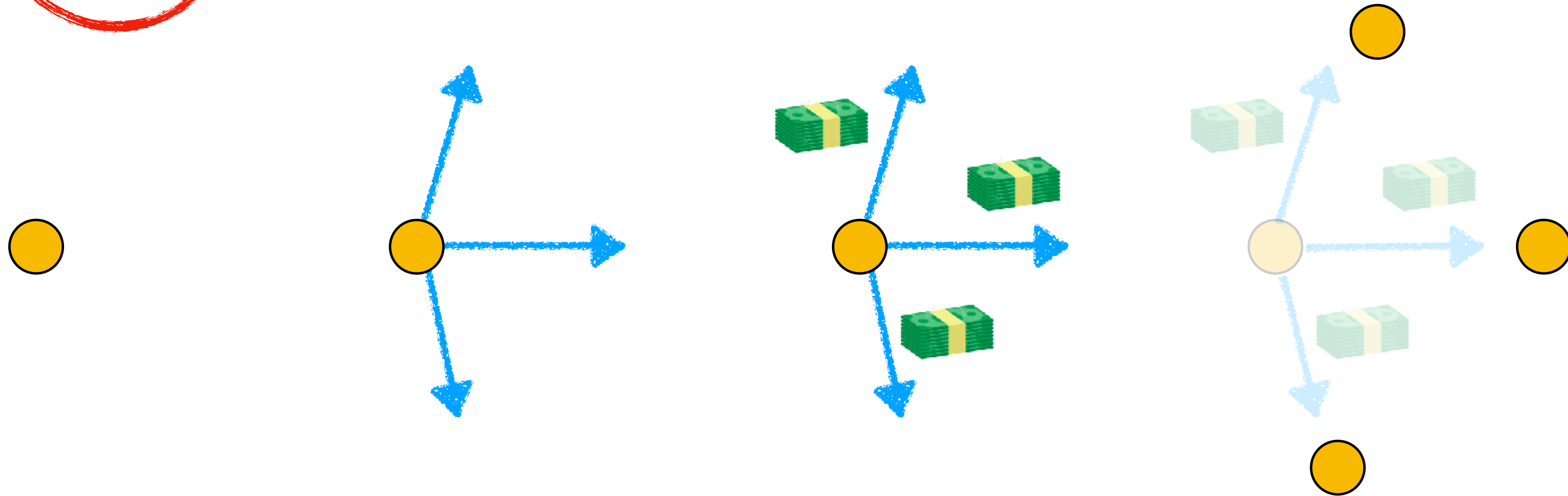
$\langle S, A, C, \mathcal{P} \rangle$



# How do we estimate state?

*A mathematical framework for modeling sequential decision making*

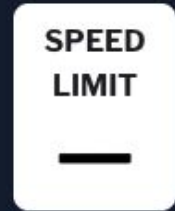
$\langle S, A, C, \mathcal{F} \rangle$



# Self-driving



—  
m/s



# How does a robot build up state?







0.0  
m/s

SPEED  
LIMIT  
**11**

ACTIVE





0.0  
m/s

SPEED  
LIMIT  
**11**

ACTIVE





0.0  
m/s

SPEED  
LIMIT  
**11**

ACTIVE





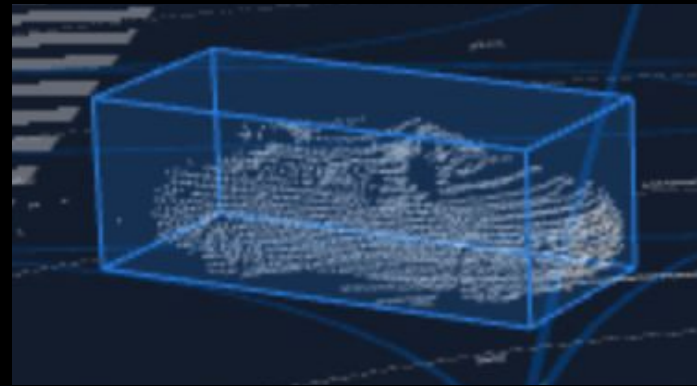
1.2  
m/s

SPEED  
LIMIT  
11

ACTIVE



pose  
 $(x, y, \psi)$



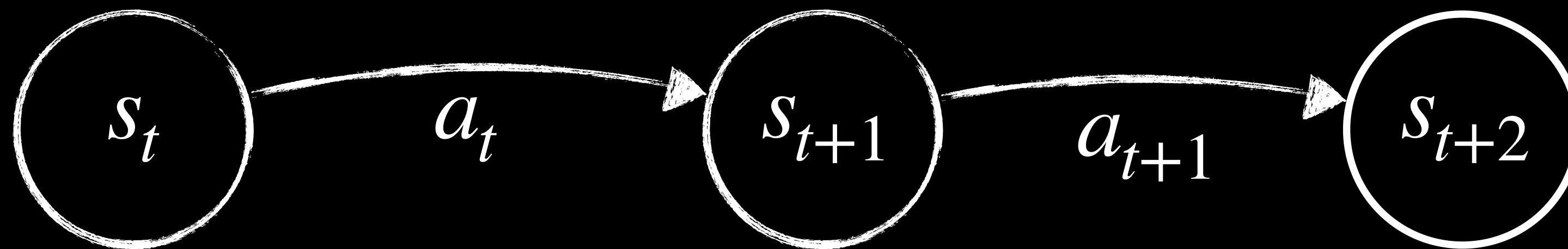
vel  
 $(\dot{x}, \dot{y}, \dot{\psi})$



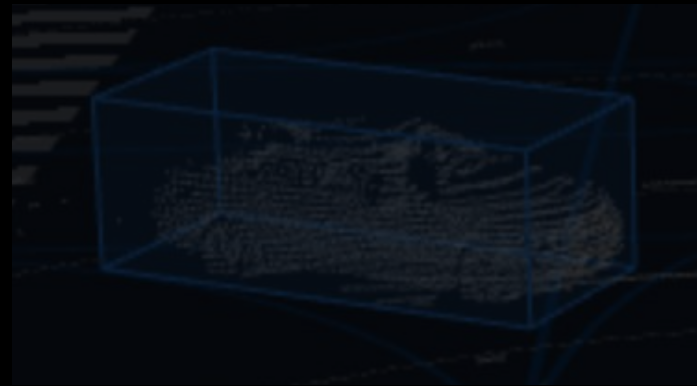
type

(pedestrian, car, cyclist)

But we do not observe  
these directly!



pose  
 $(x, y, \psi)$



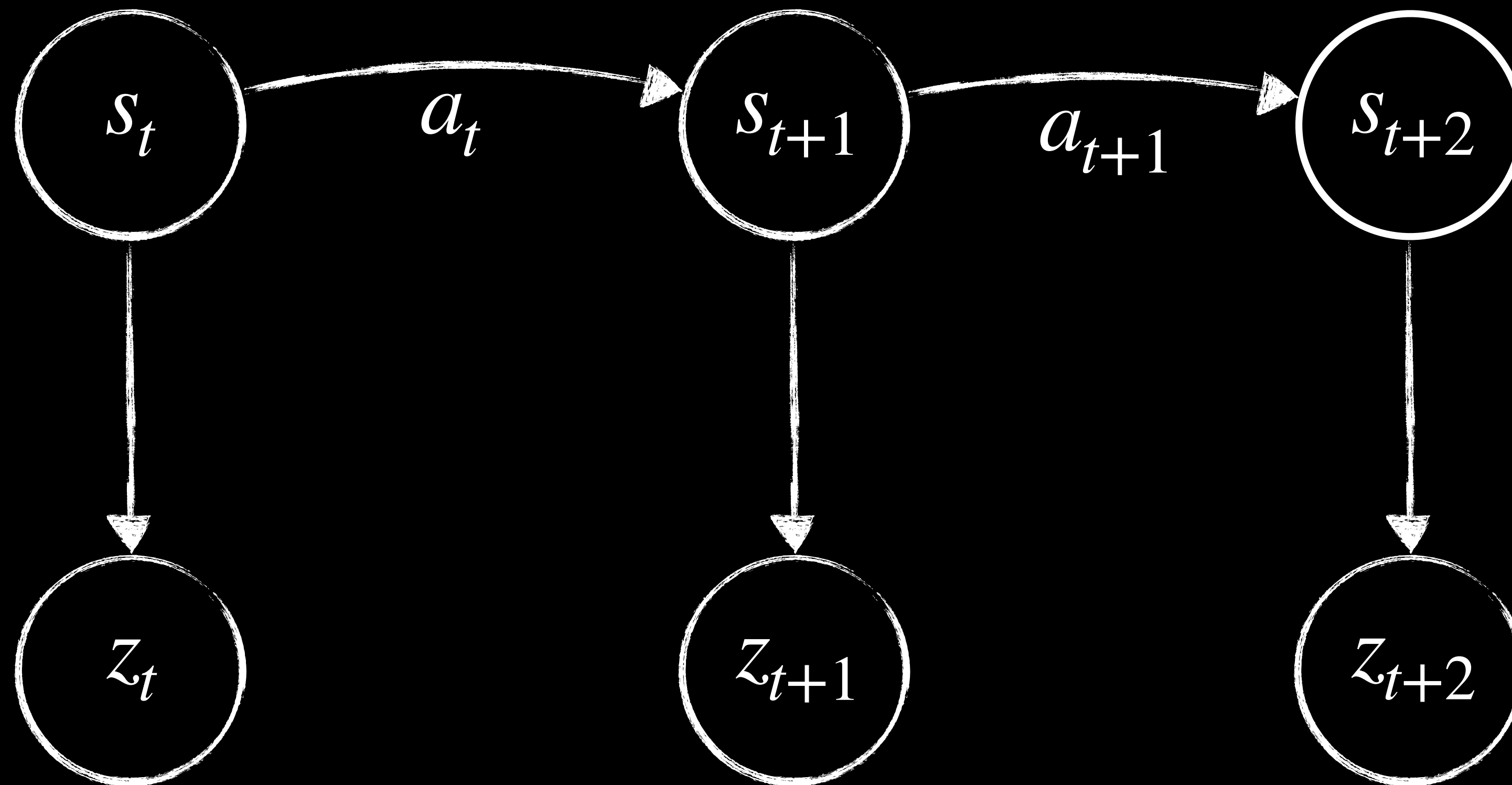
vel  
 $(\dot{x}, \dot{y}, \dot{\psi})$



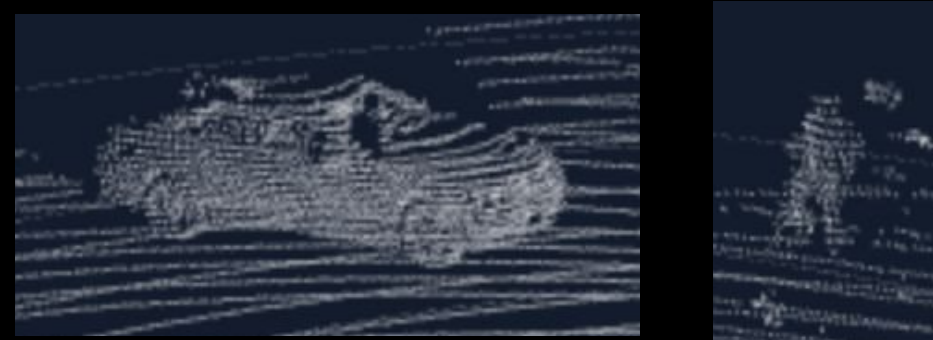
type

(pedestrian, car, cyclist)

# Estimate state from observations

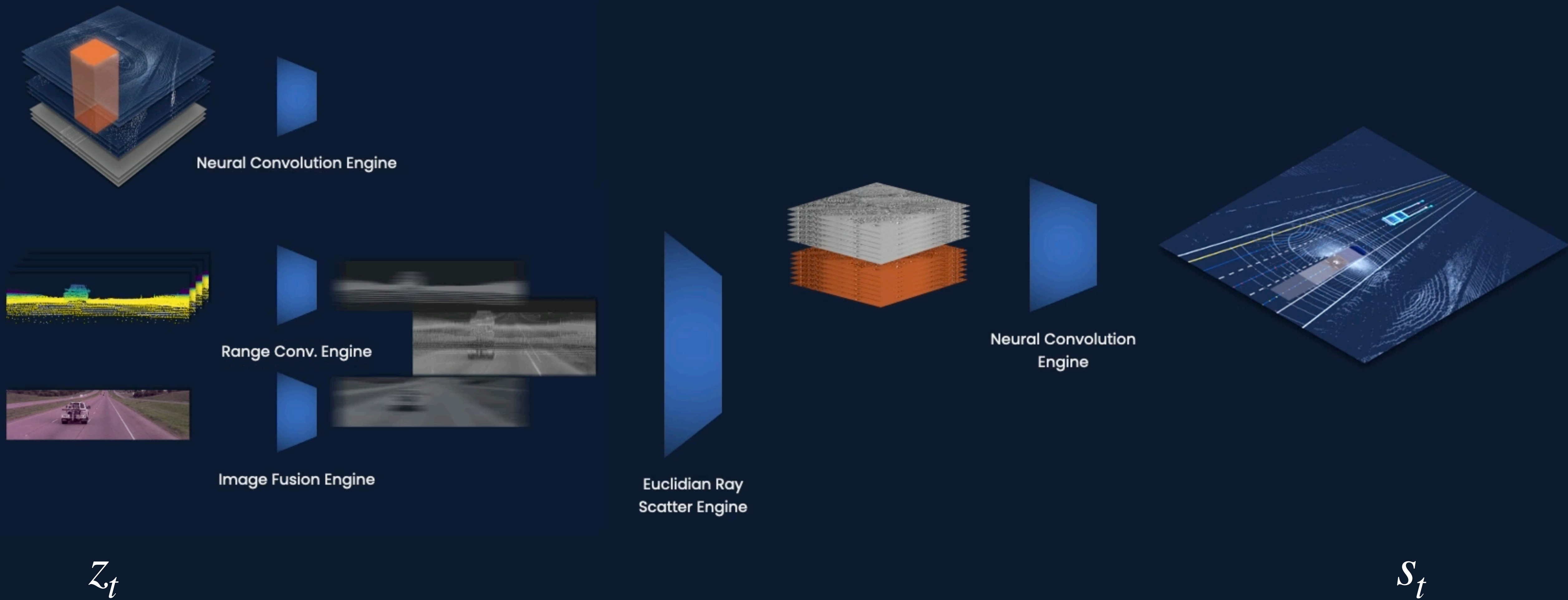


camera



lidar

# Learning a state representation



Goal: Learn a state representation

$$P(s_t | z_t, z_{t-1}, z_{t-2}, \dots)$$



# Another example: Manipulation

$$P(s_t | z_t, z_{t-1}, z_{t-2}, \dots)$$

What is observation  $z_t$  ?

What is state  $s_t$  ?



From Jeannette Bohg: CS231A Computer Vision: From 3D Reconstruction to Recognition

# Another example: Manipulation

$$P(s_t | z_t, z_{t-1}, z_{t-2}, \dots)$$

What is observation  $z_t$  ?

RGB Image + Depth (Kinect)

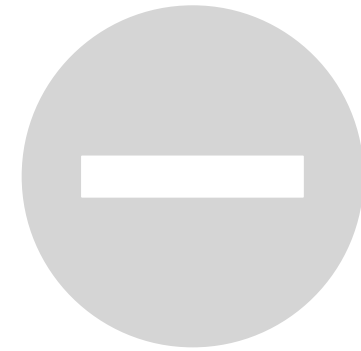
What is state  $s_t$  ?

6D Pose of the object



From Jeannette Bohg: CS231A Computer Vision: From 3D Reconstruction to Recognition

# Today's Class



What is a **visual** representation?  
What makes for a **good** representation?



How do we learn a visual representation?



How do we train a representation by  
**unsupervised** learning and **self-supervised** learning?

# Today's Class



What is a **visual** representation?  
What makes for a **good** representation?



How do we learn a visual representation?

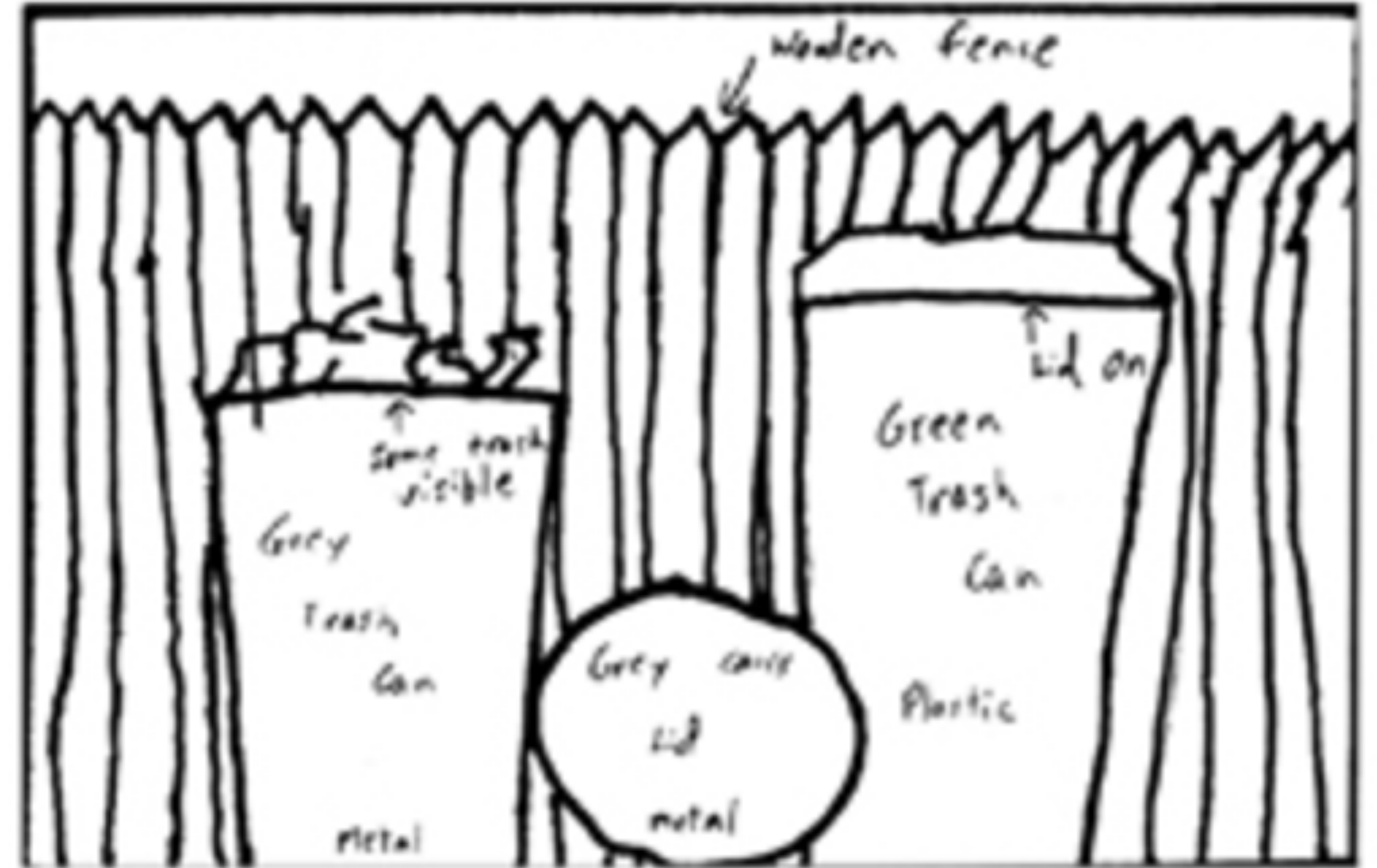


How do we train a representation by  
**unsupervised** learning and **self-supervised** learning?

Observed image



Drawn from memory



[Bartlett, 1932]

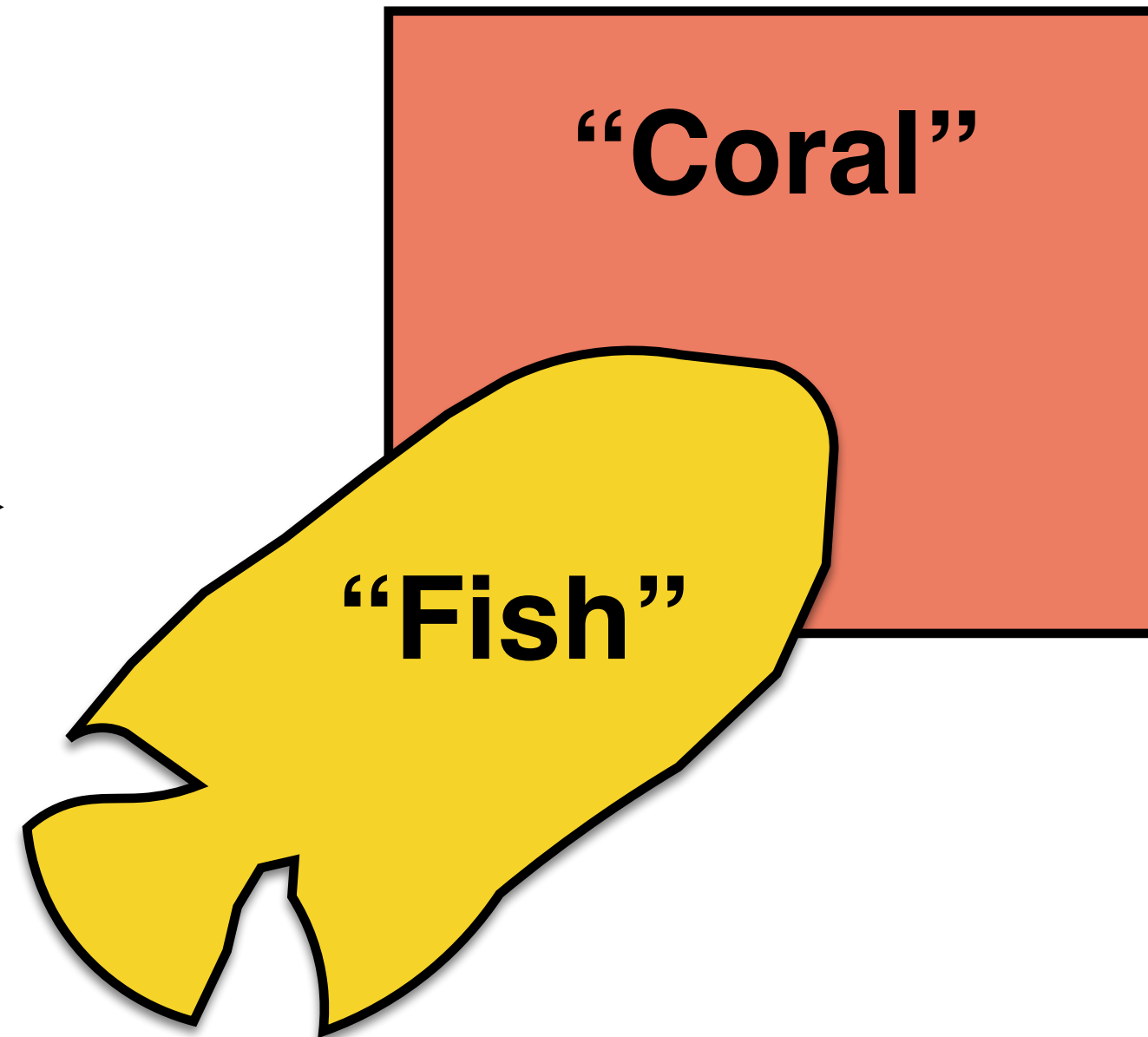
[Intraub & Richardson, 1989]

# Representation learning

$X$



Image



Compact mental  
representation

Activity!

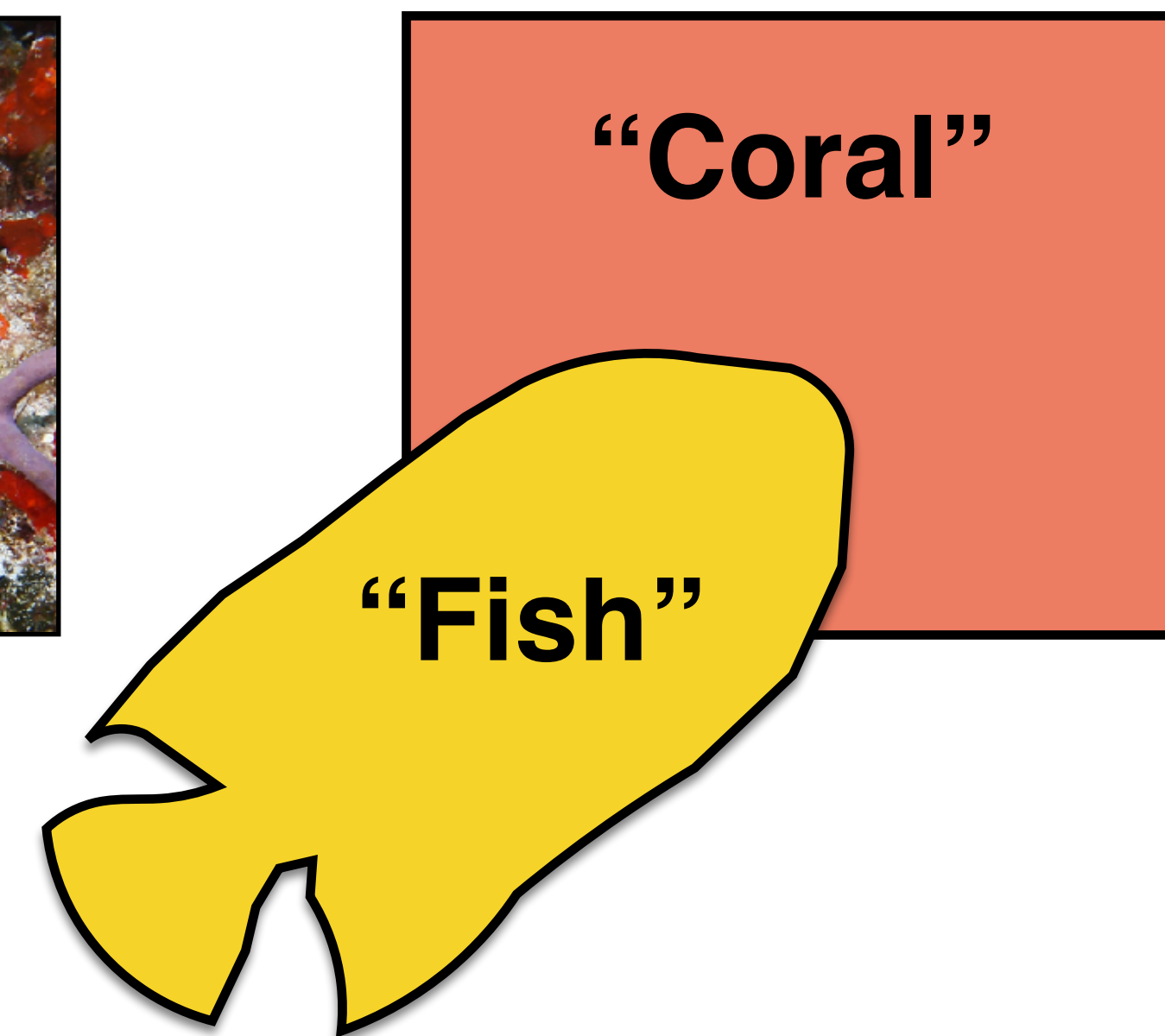


# Think-Pair-Share!

Think (30 sec): What makes a representation good?

Pair: Find a partner

Share (45 sec): Partners exchange ideas



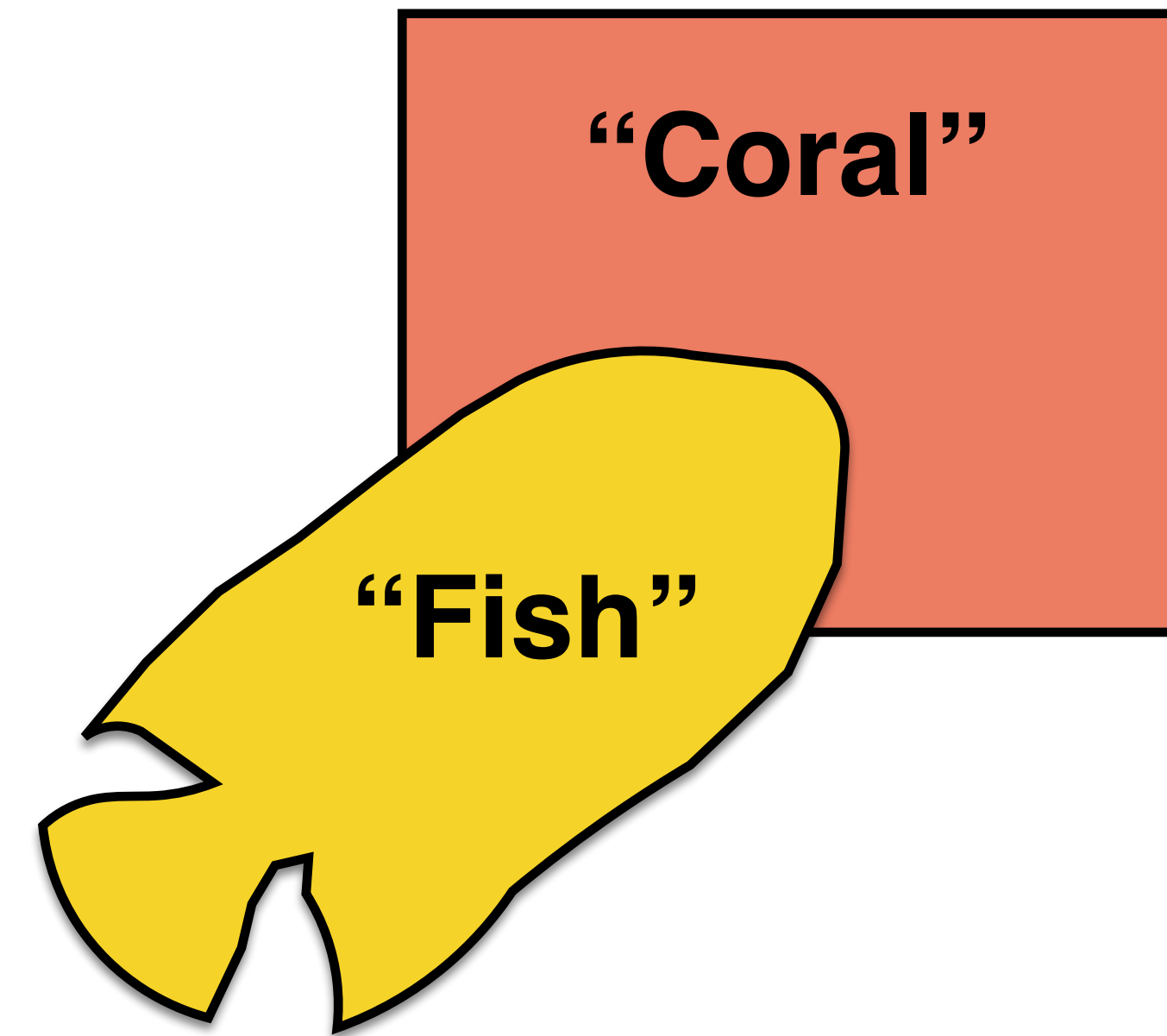


# Representation learning

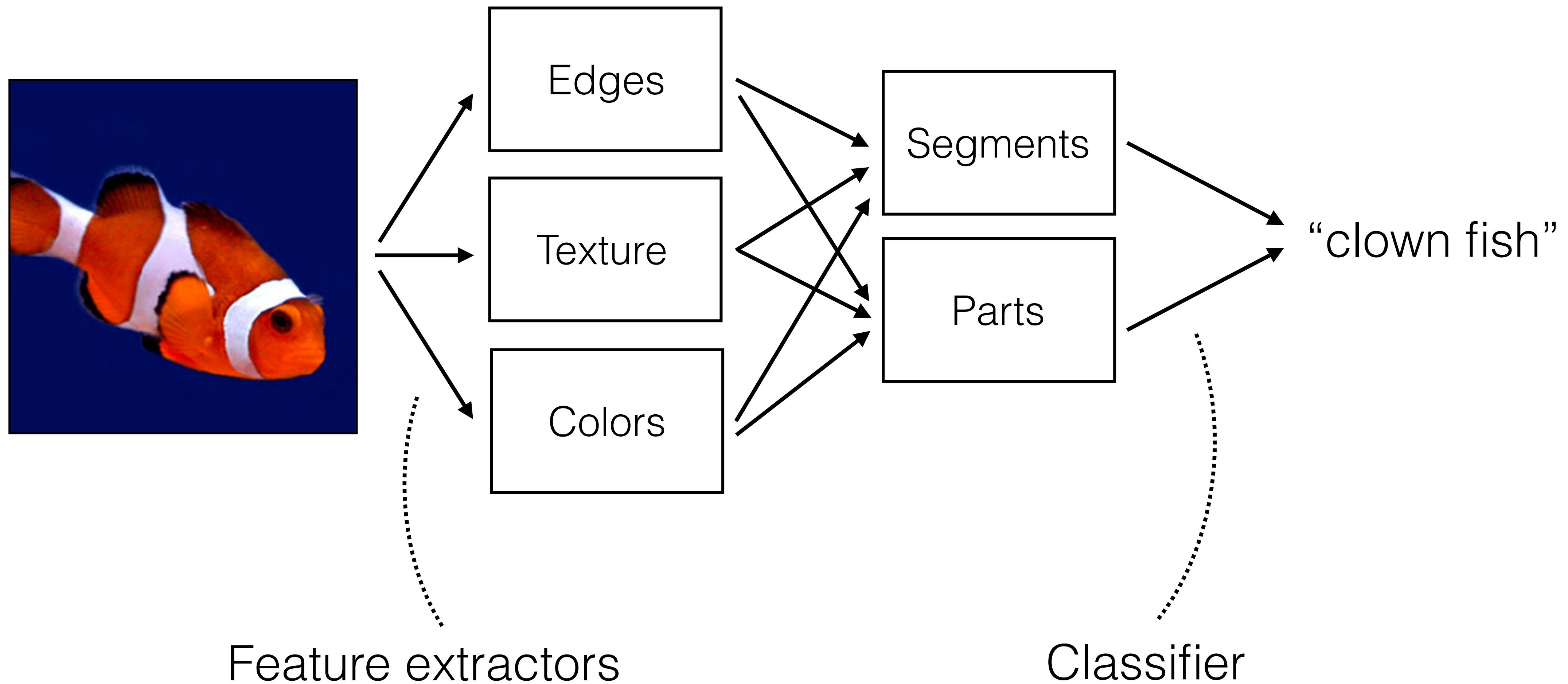
Good representations are:

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

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



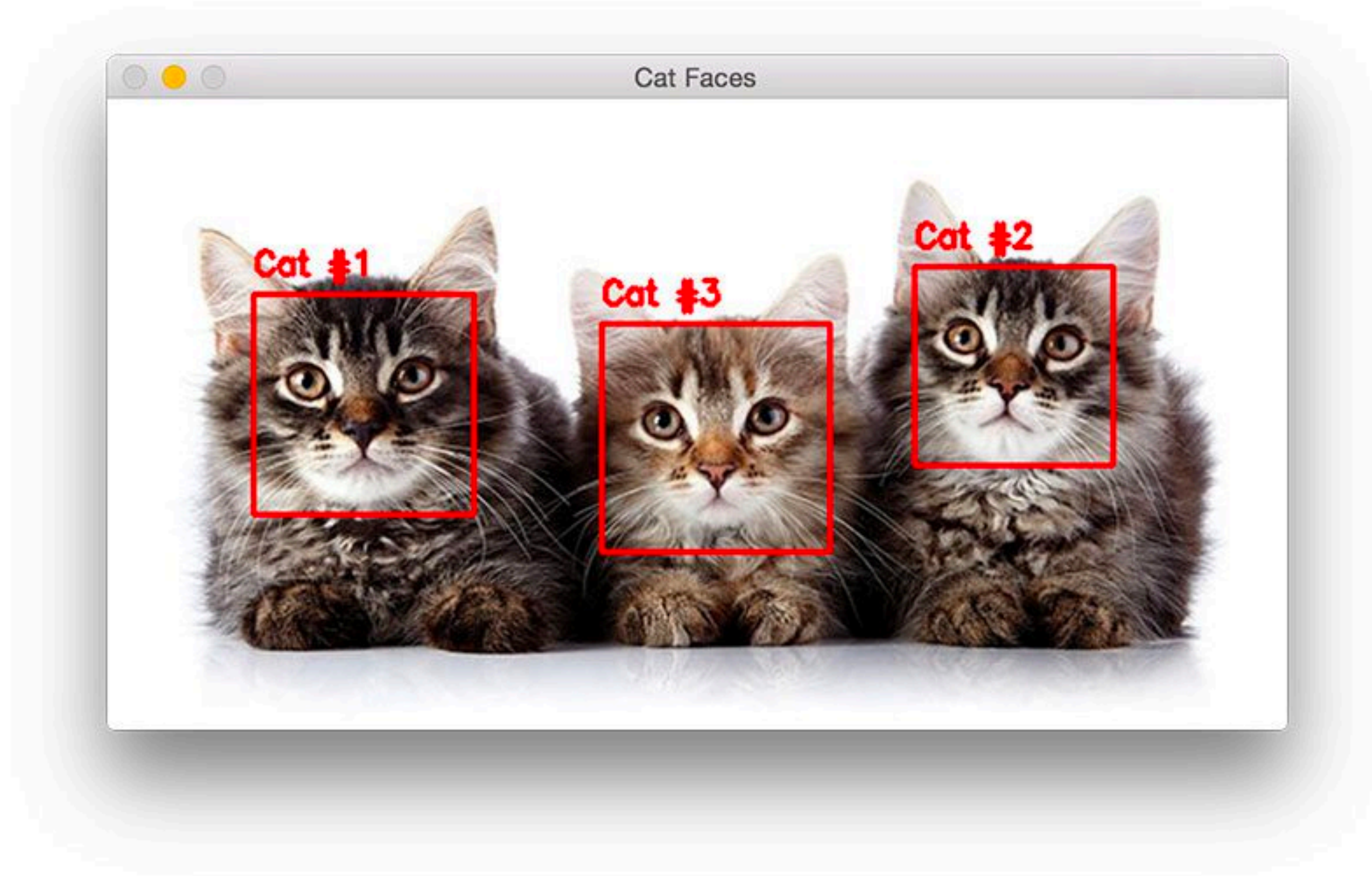
# Traditional CV Pipeline



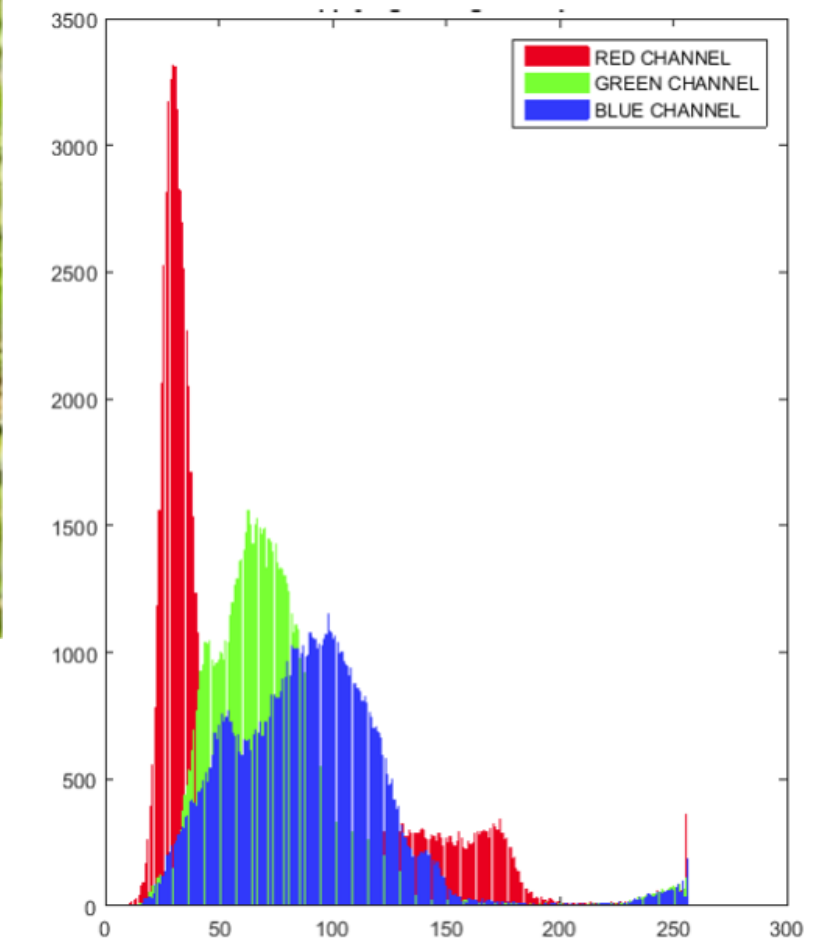
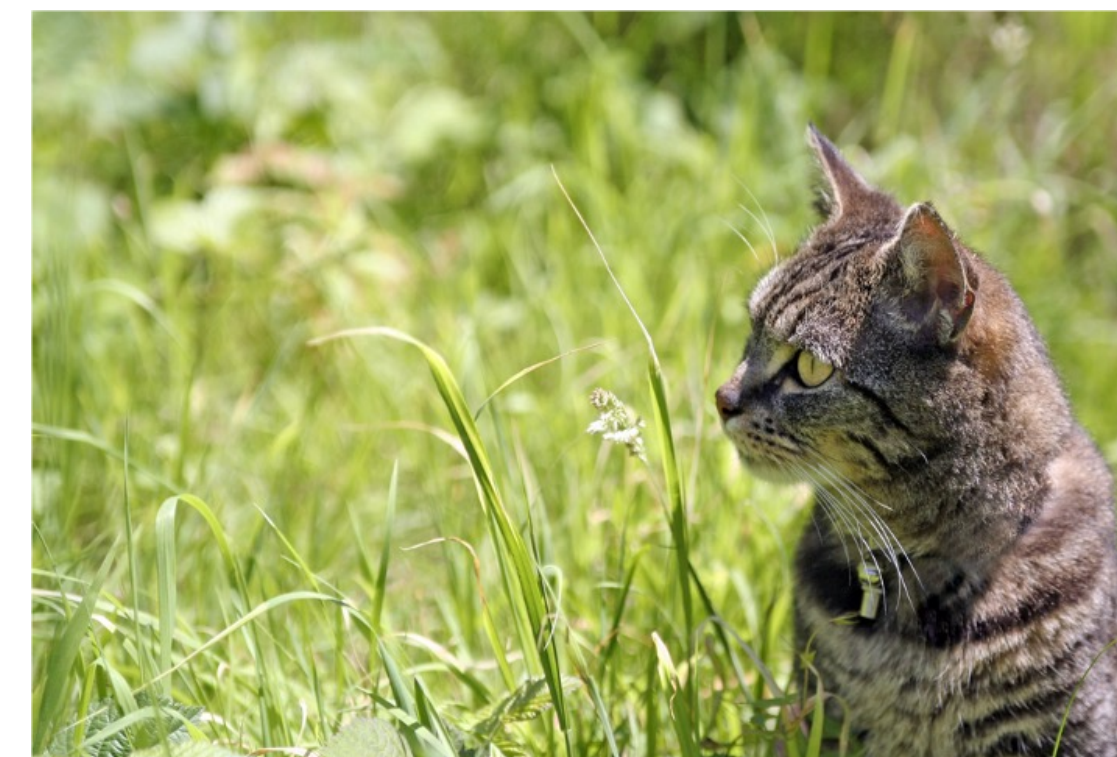
Can there be more than  
one representation that is  
useful?



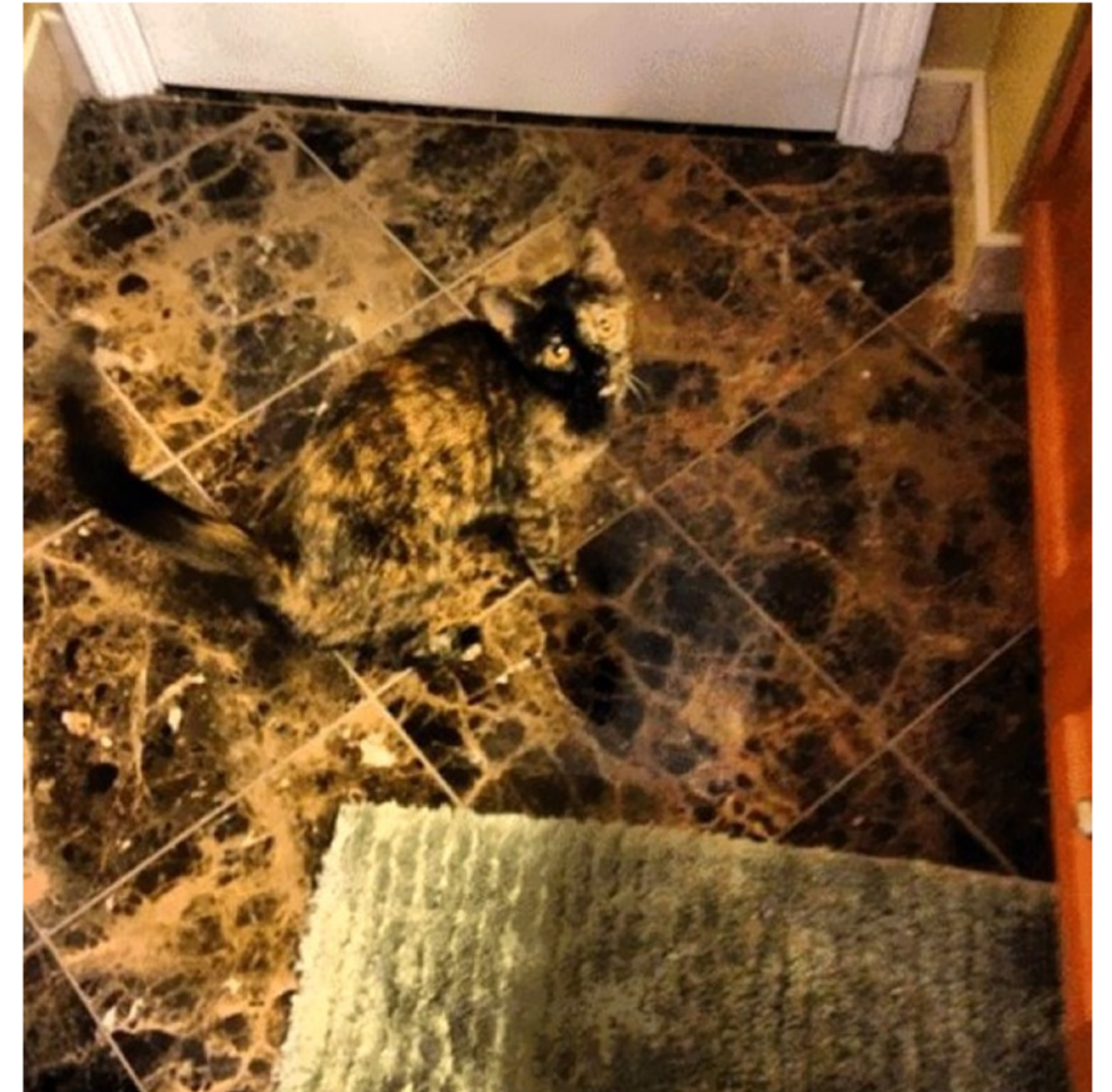
# Exercise: Let's think about cats!



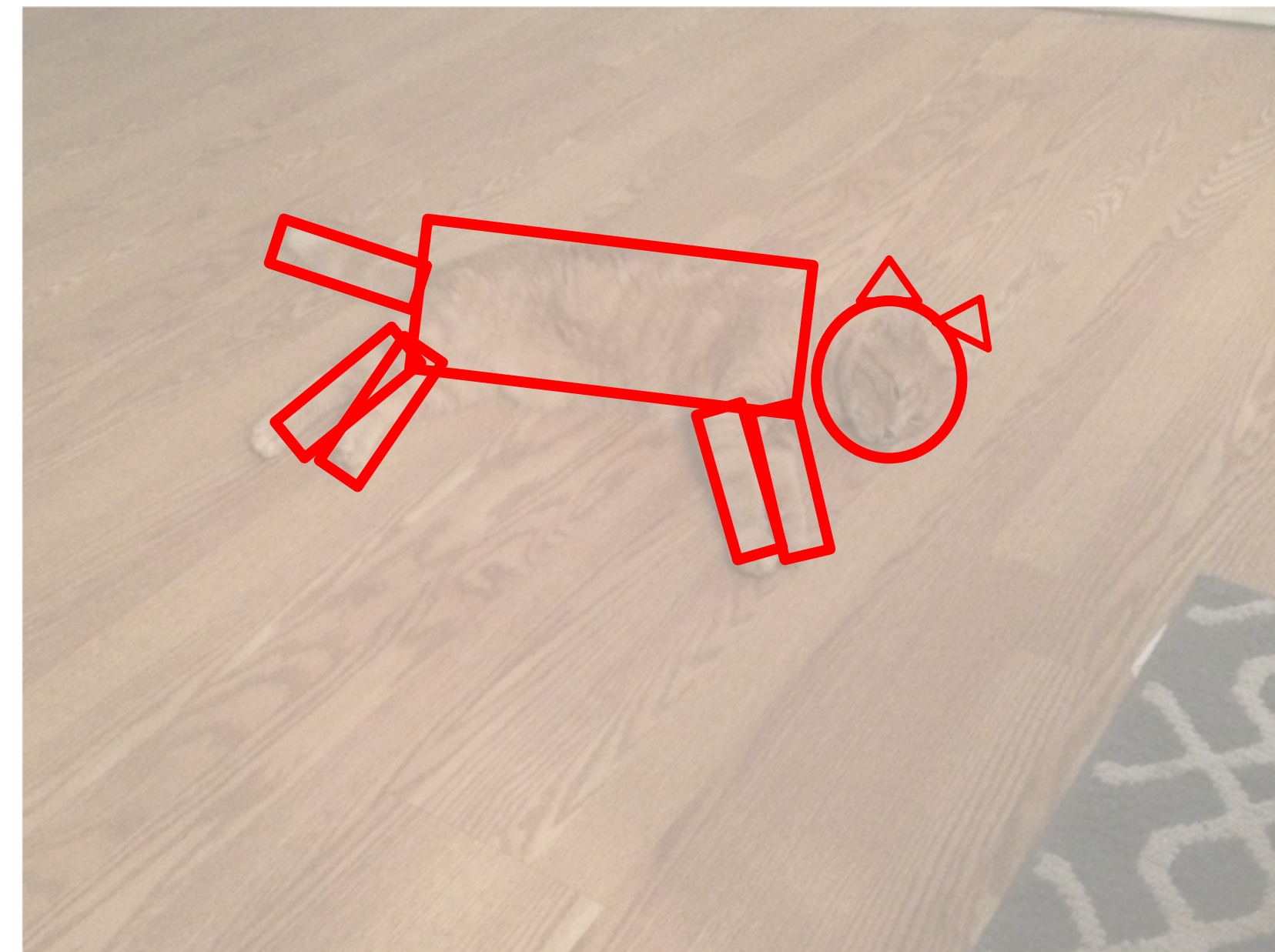
# Represent these cats with a cat detector



# Represent these cats with a cat detector (II)

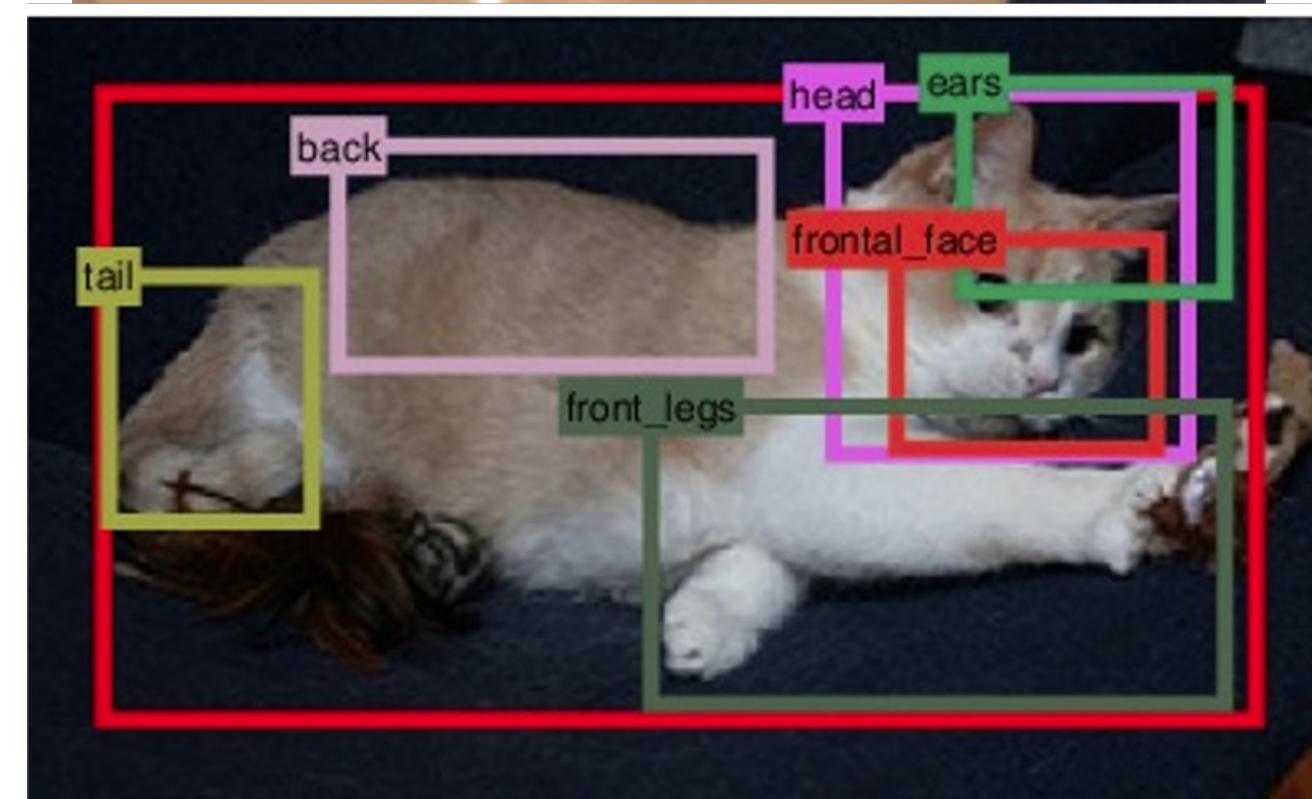


# Represent these cats with a cat detector (II)



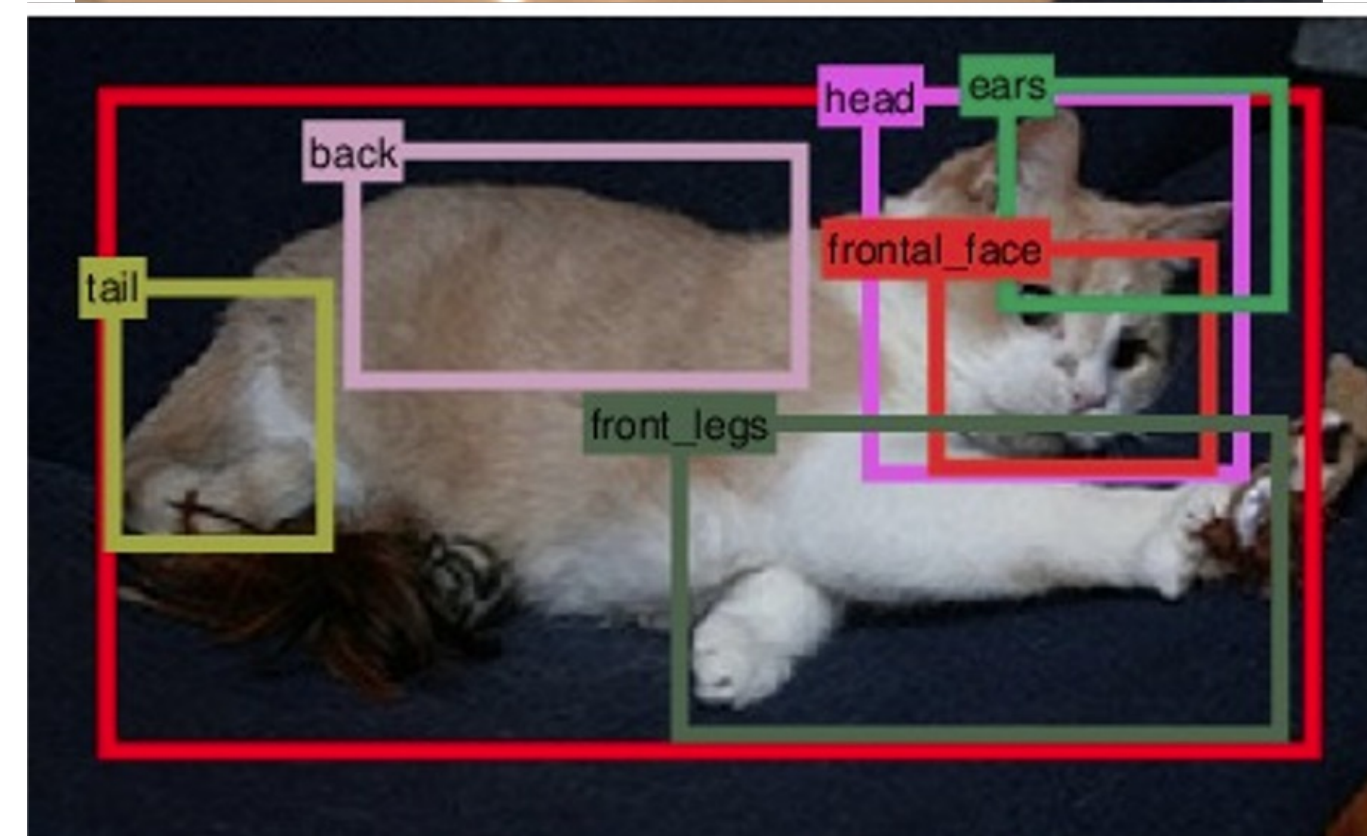
Example from CS331B: Representation Learning in Computer Vision

# Represent these cats with a cat detector (III)

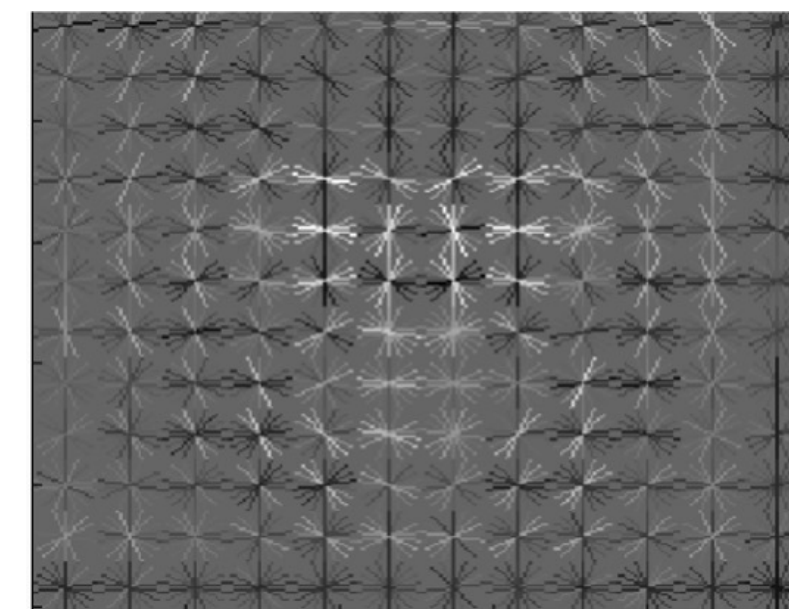
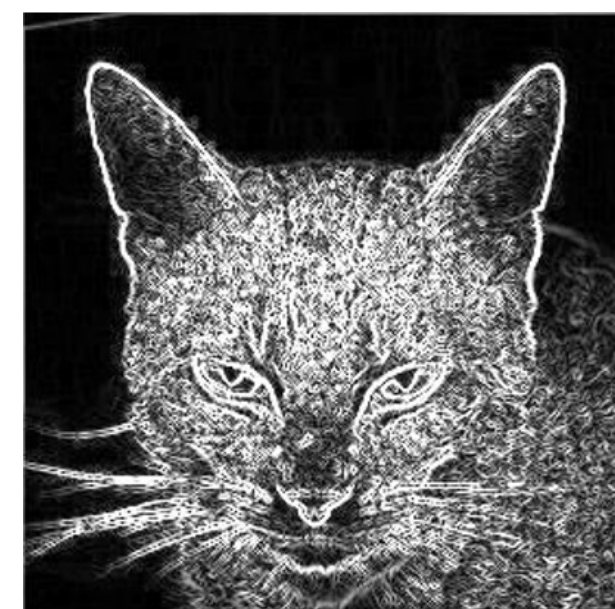
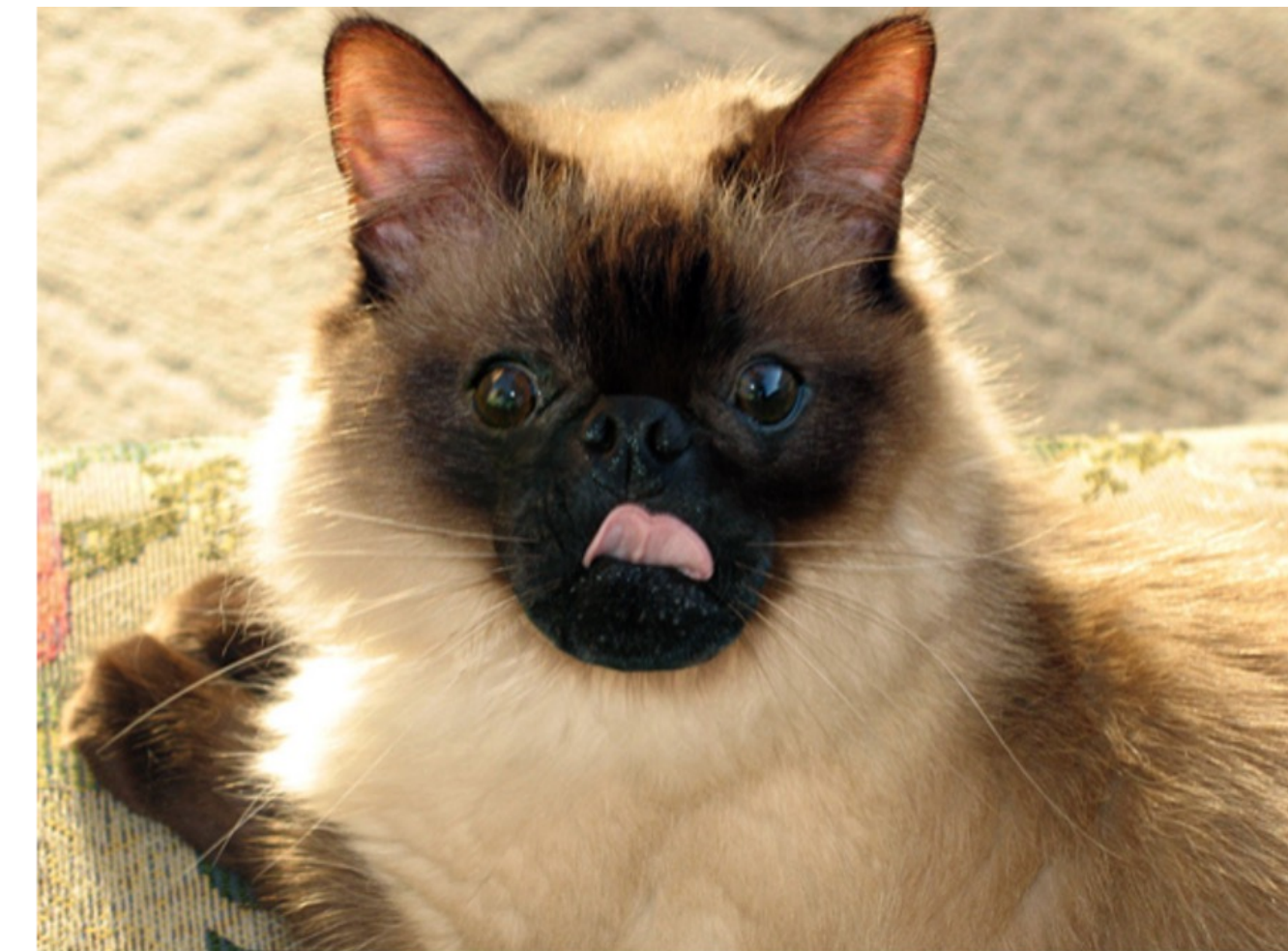




# Represent these cats with a cat detector (III)



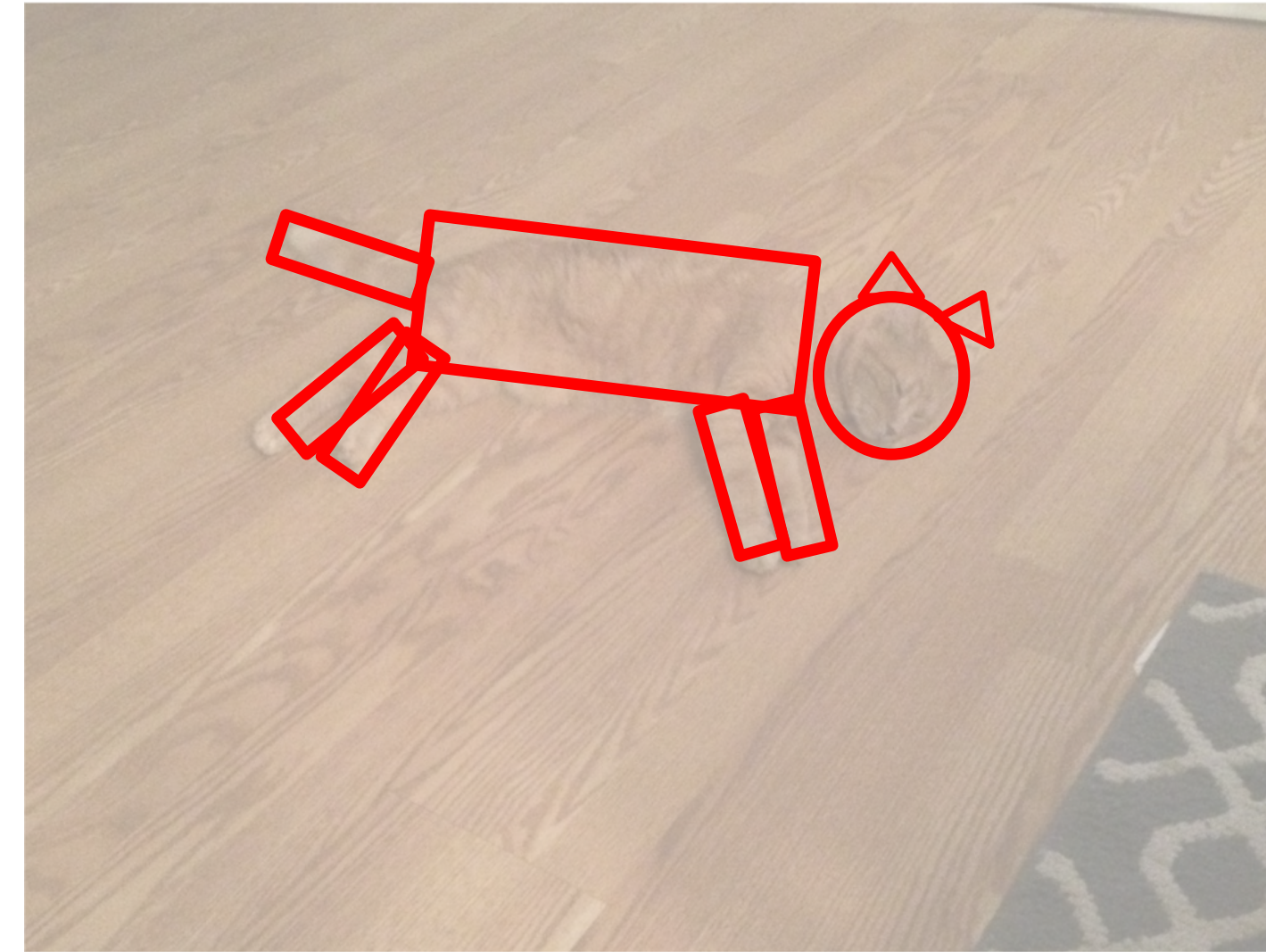
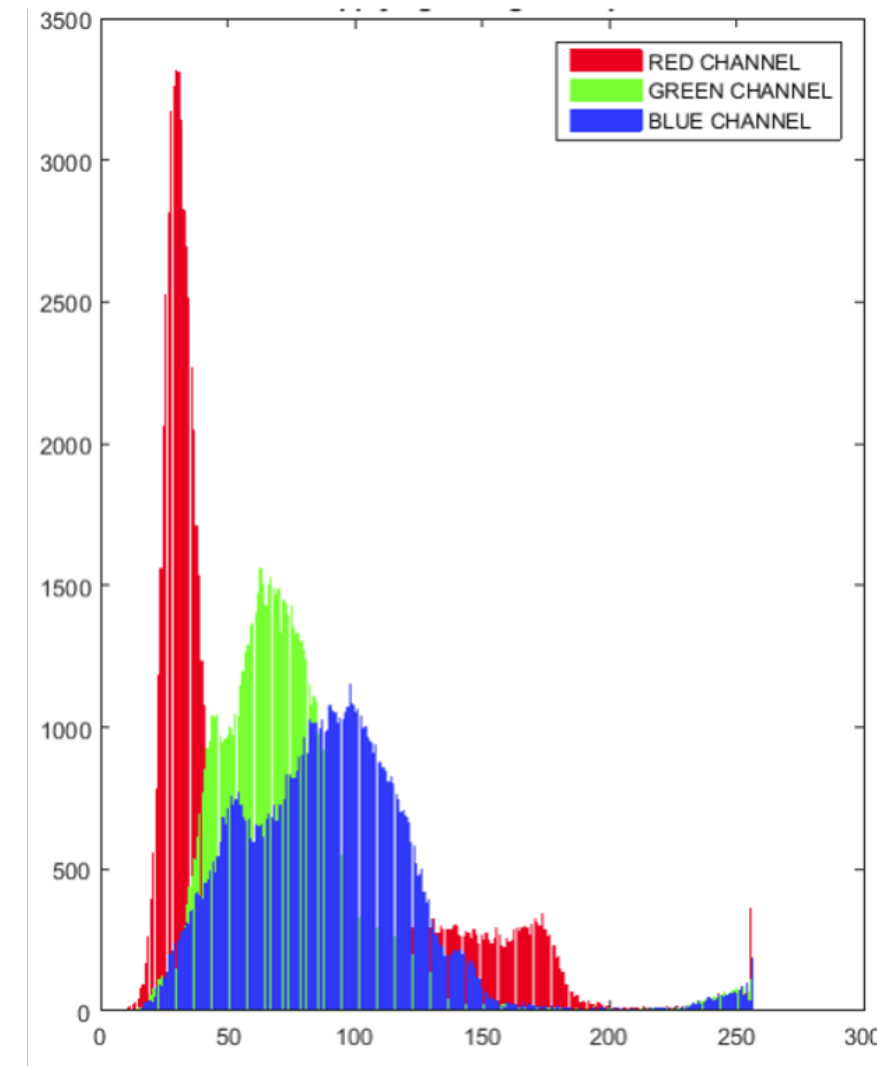
# Represent these cats with a cat detector (IV)



Example from CS331B: Representation Learning in Computer Vision

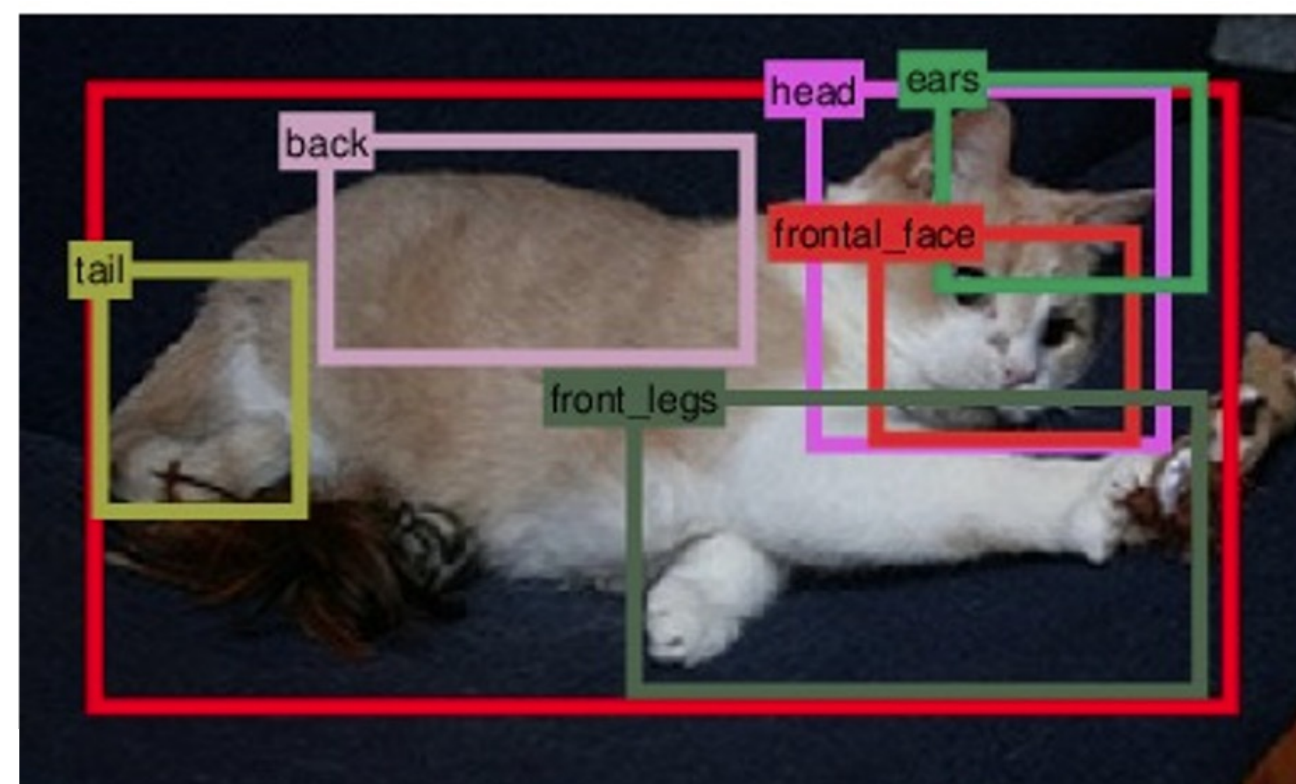
# Summary of representations

Color Histograms

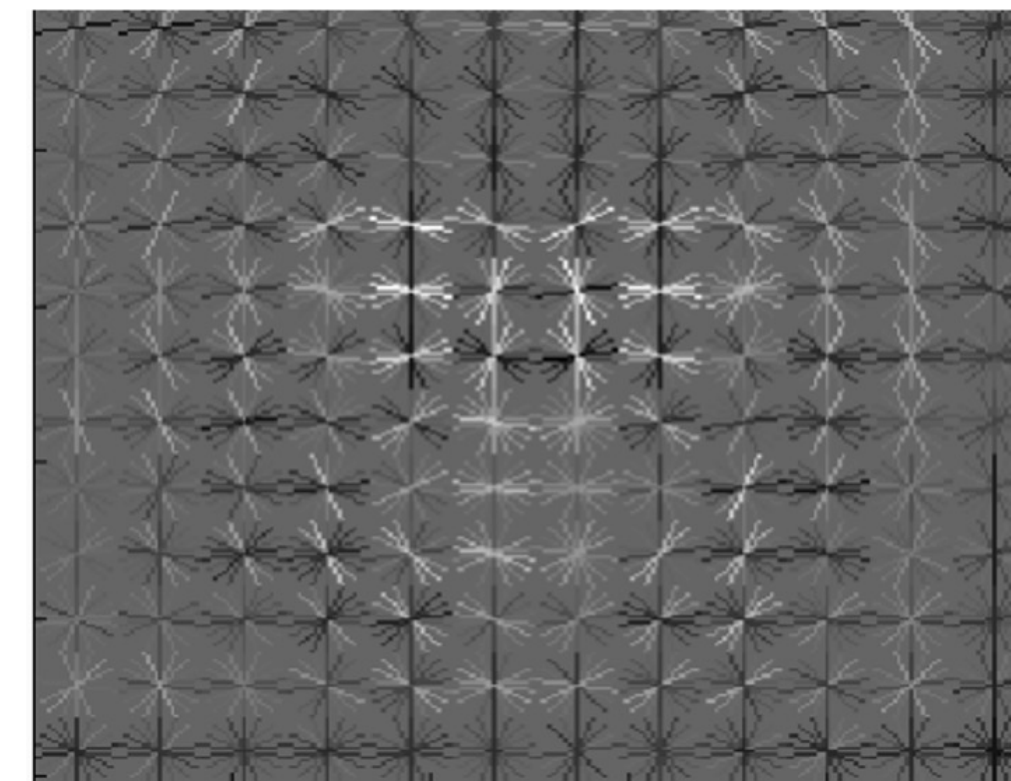


Model based Shapes

Deformable Part based Models (DPM)



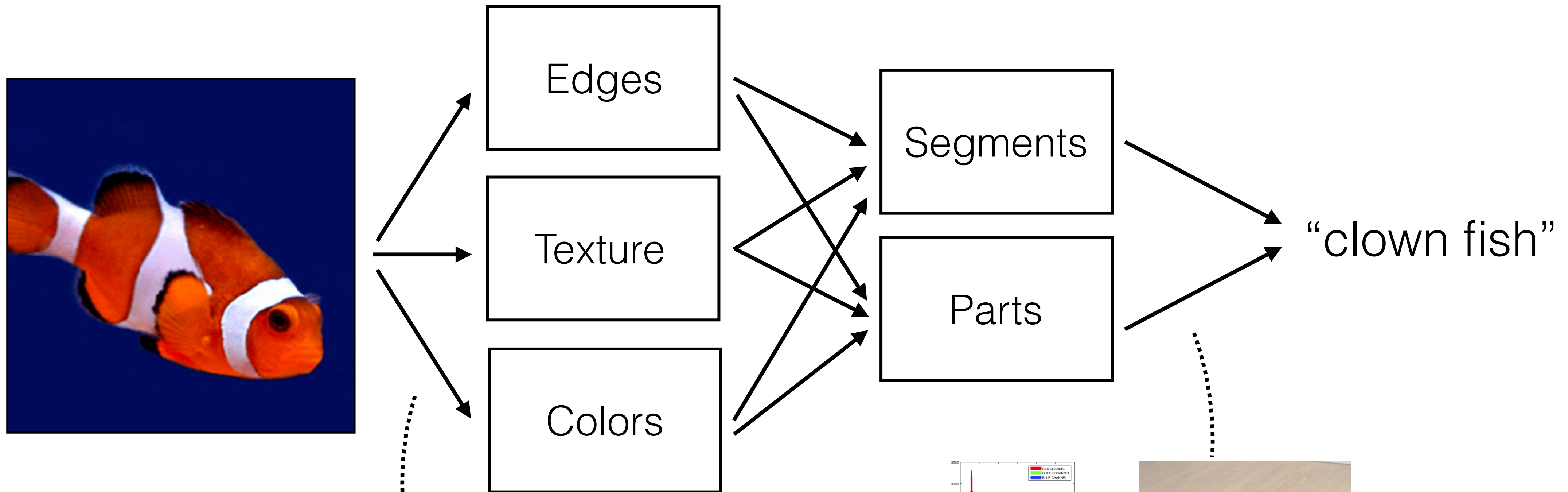
Felzenszwalb et al. 2010.  
Dalal and Triggs, 2005.  
Beis and Lowe, 1997.



Histogram of Gradients (HOG)

Example from CS331B: Representation Learning in Computer Vision

# Traditional CV Pipeline



Feature extractors

Color Histograms

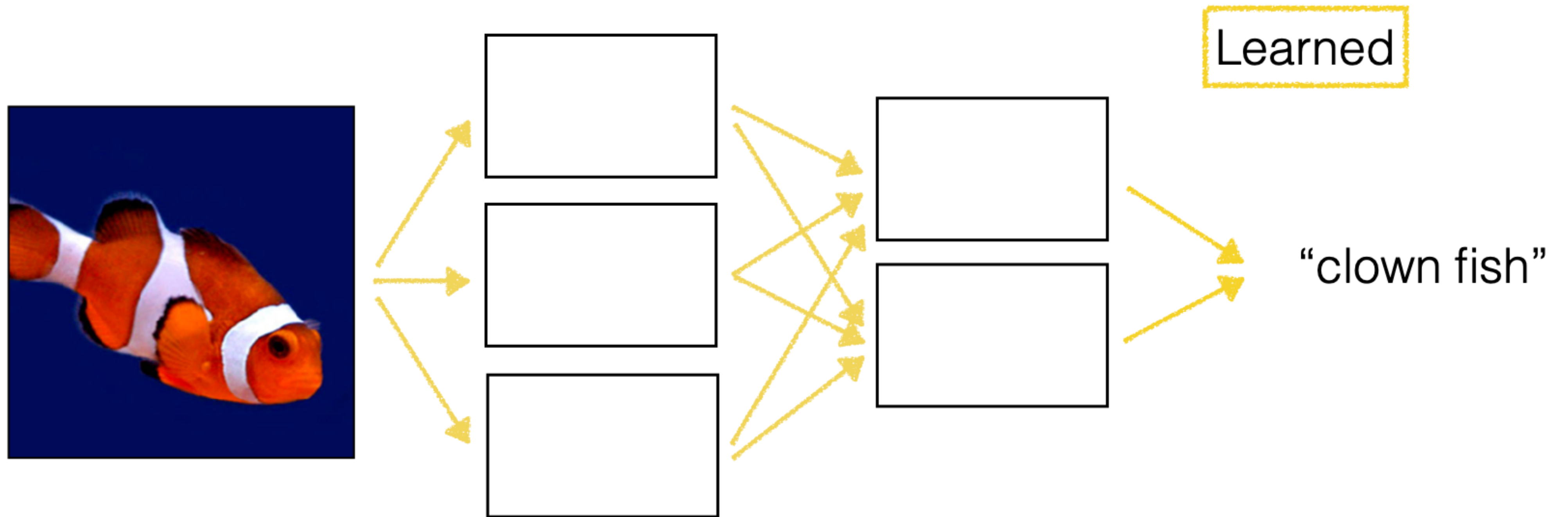
Model based Shapes

Deformable Part based Models (DPM)

Histogram of Gradients (HOG)

Felzenszwalb et al. 2010.  
Dalal and Triggs, 2005.  
Beis and Lowe, 1997.

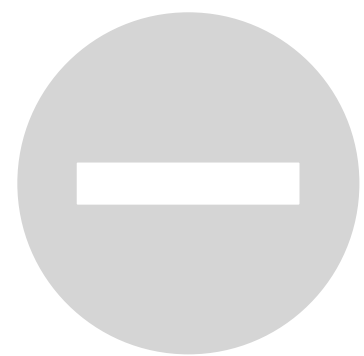
# Learned CV Pipeline



# Today's Class



What is a **visual** representation?  
What makes for a **good** representation?



How do we learn a visual representation?

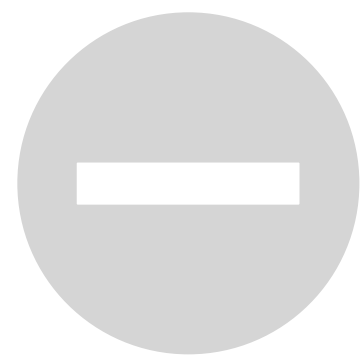


How do we train a representation by  
**unsupervised** learning and **self-supervised** learning?

# Today's Class



What is a **visual** representation?  
What makes for a **good** representation?



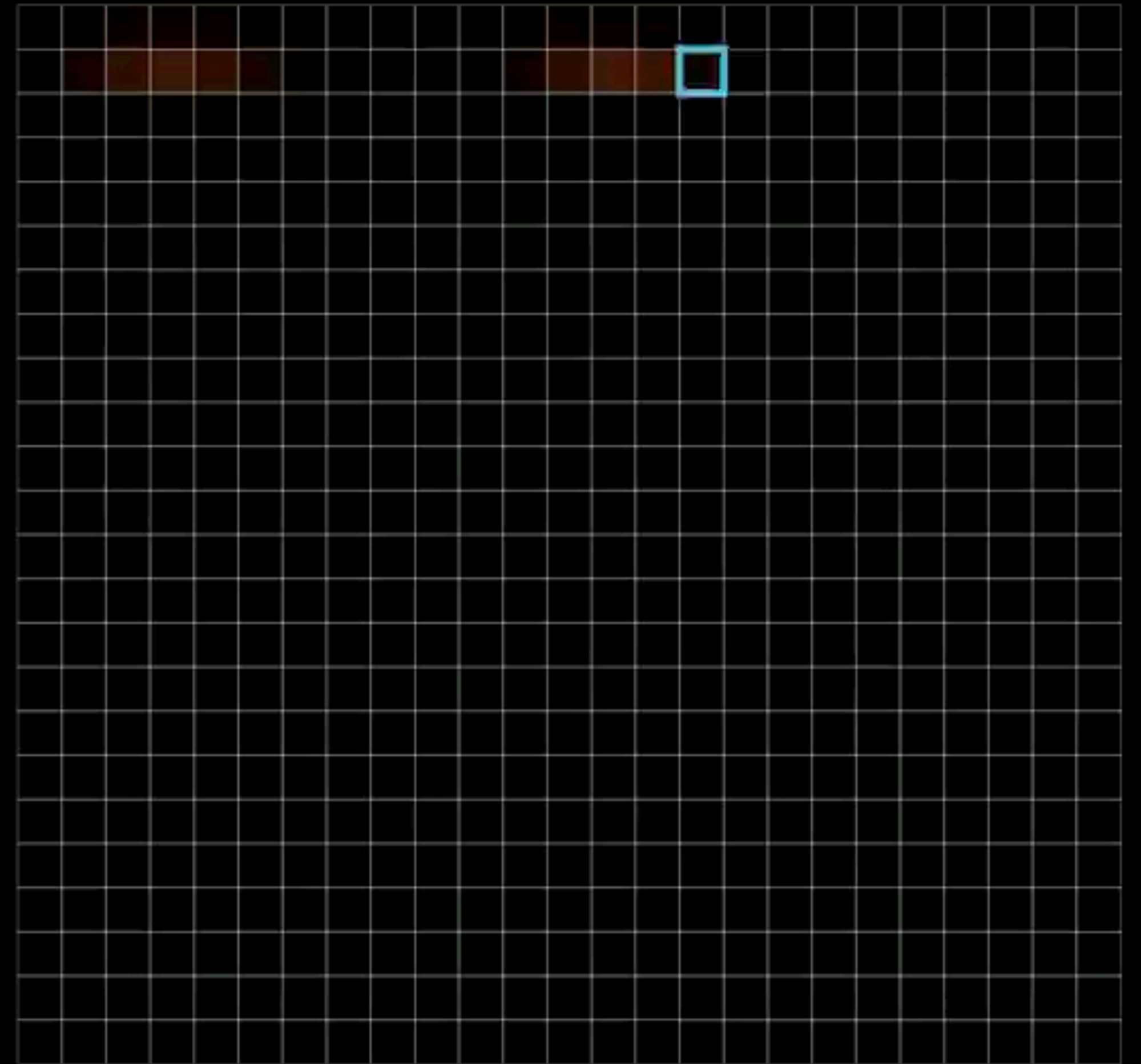
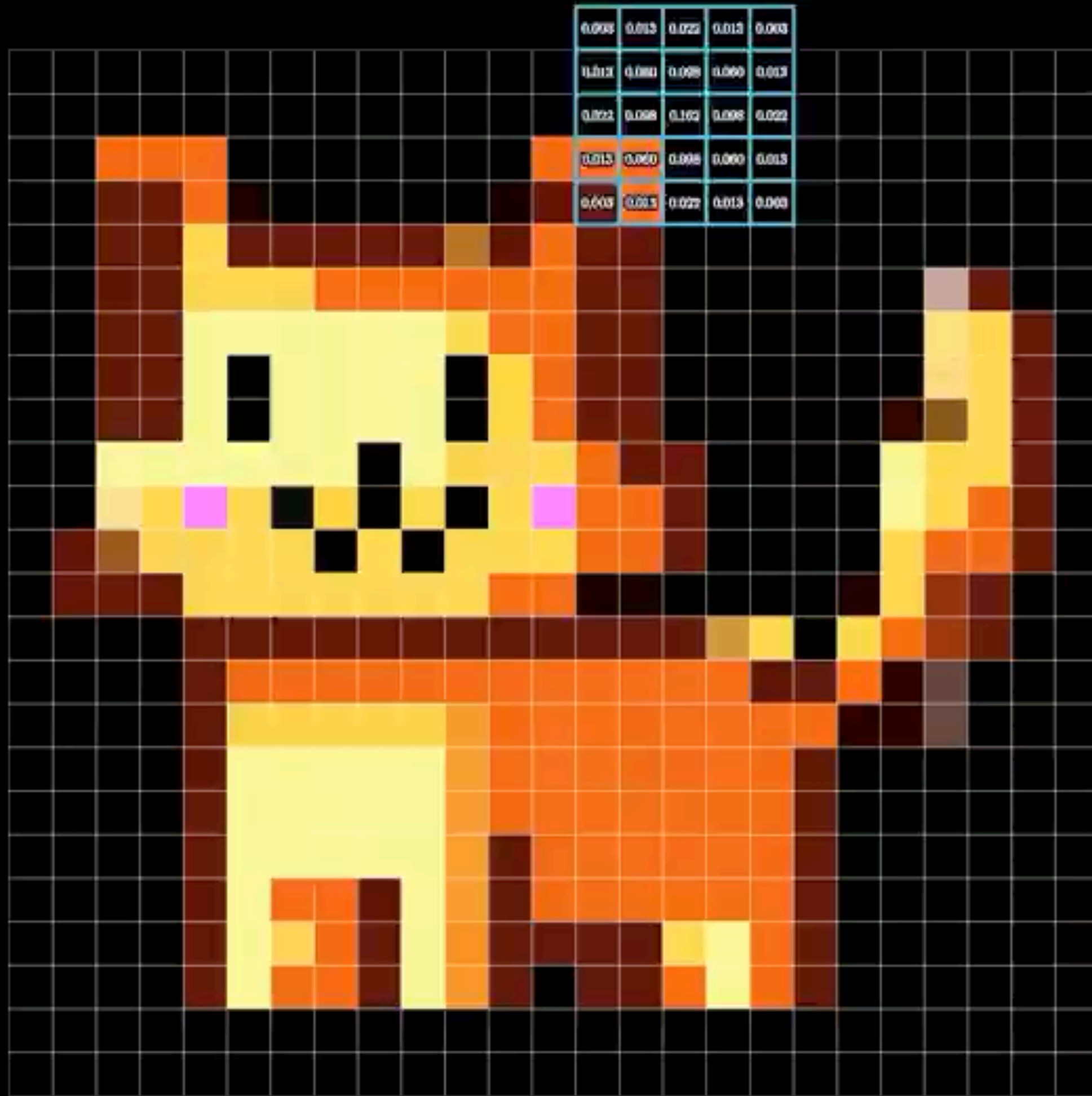
How do we learn a visual representation?



How do we train a representation by  
**unsupervised** learning and **self-supervised** learning?

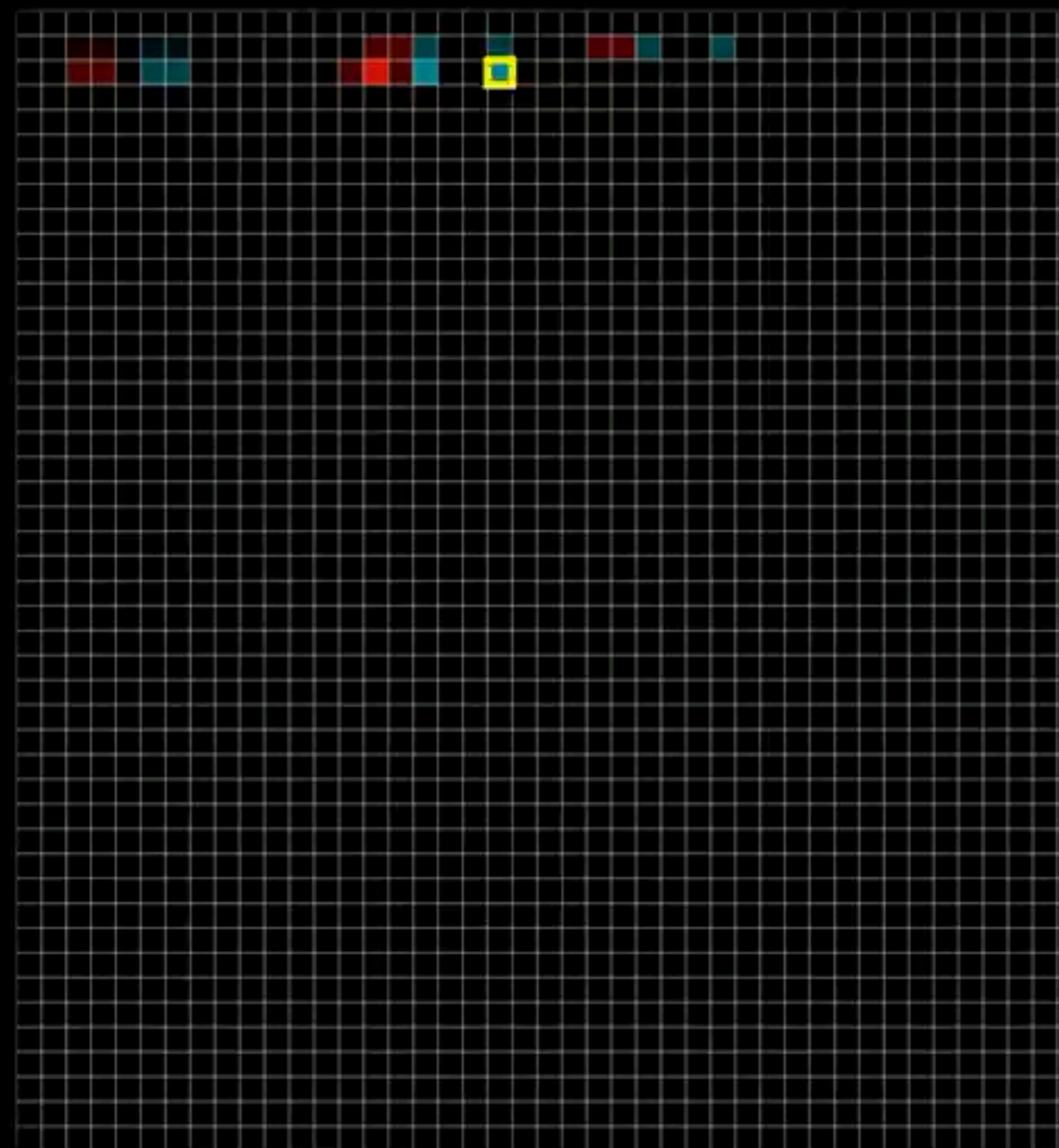
# Convolutional Neural Network





Video from ThreeBlue1Brown: <https://www.youtube.com/watch?v=KuXjwB4LzSA&t=738s>

0.25	0.00	-0.25
0.50	0.00	-0.50
0.25	0.00	-0.25



Video from ThreeBlue1Brown: <https://www.youtube.com/watch?v=KuXjwB4LzSA&t=738s>



# Training a CNN Classifier

[https://pytorch.org/tutorials/  
beginner/blitz/cifar10\\_tutorial.html](https://pytorch.org/tutorials/beginner/blitz/cifar10_tutorial.html)

# CIFAR-10 Dataset

**airplane**



**automobile**



**bird**



**cat**



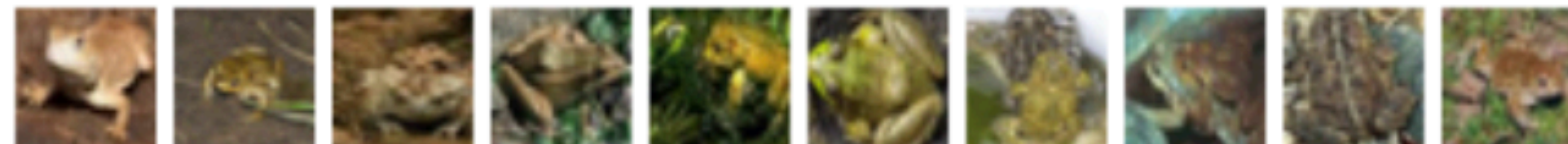
**deer**



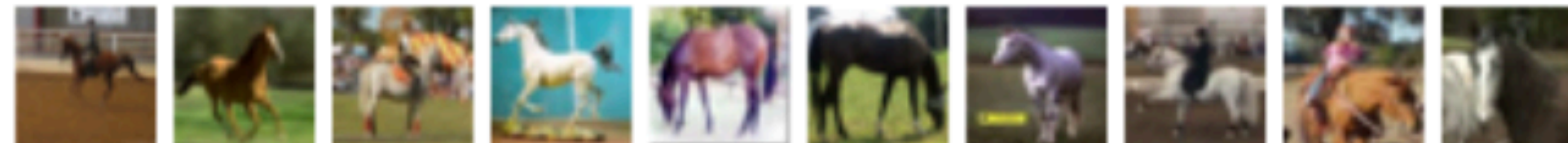
**dog**



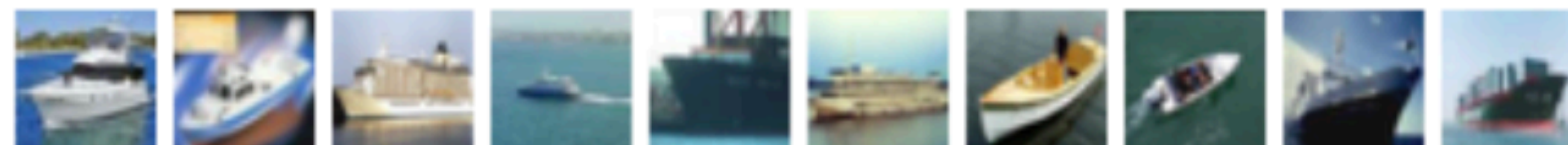
**frog**



**horse**



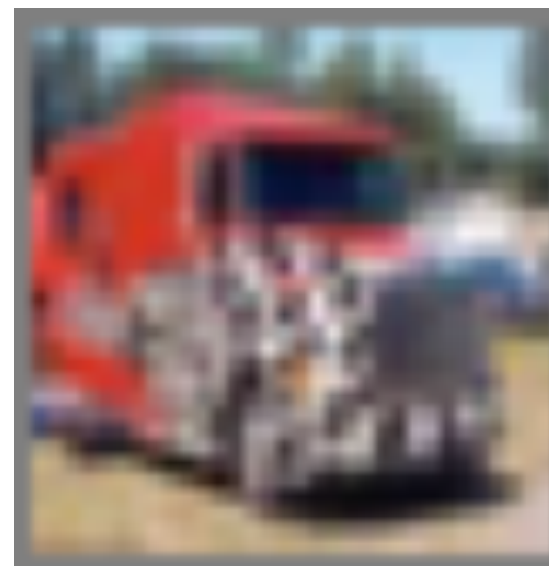
**ship**



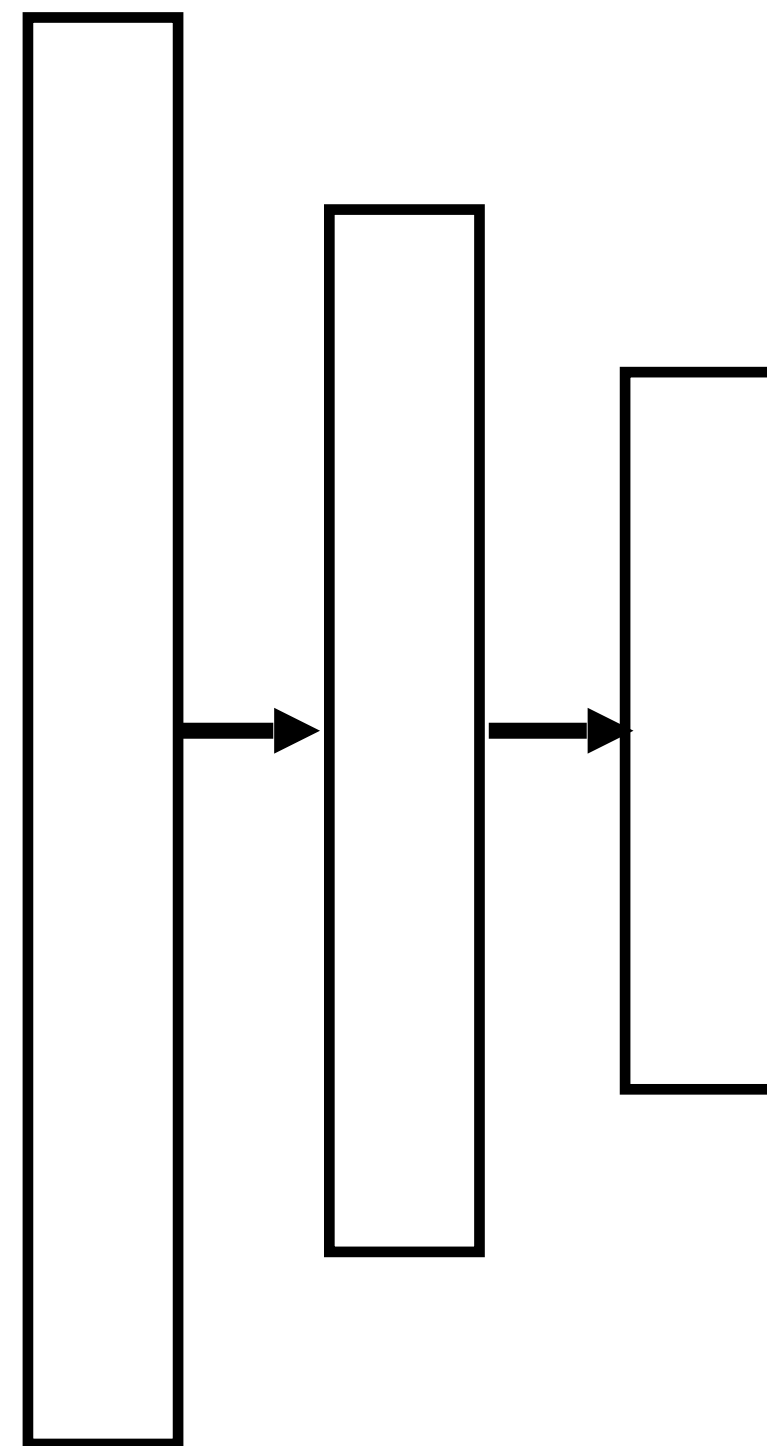
**truck**



# CIFAR-10 Dataset



Input: Image



0.0

0.1

0.0

0

0.07

0.01

0.0

0.01

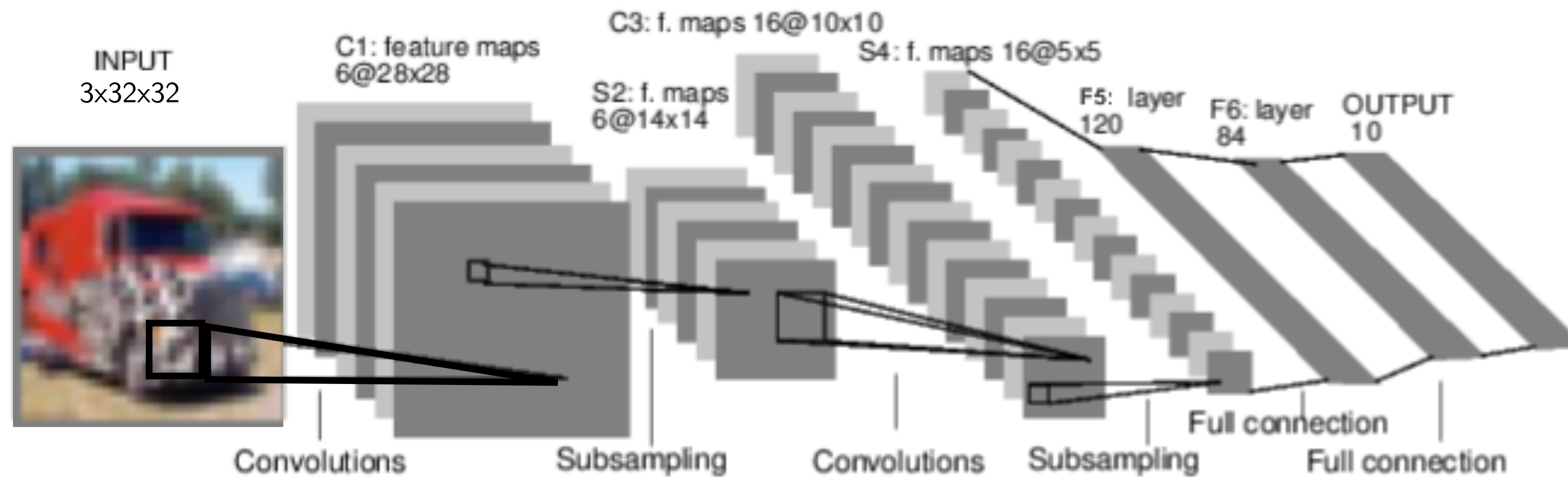
**0.8**

0.01

Truck!

Output: Class Probabilities

# Model Architecture

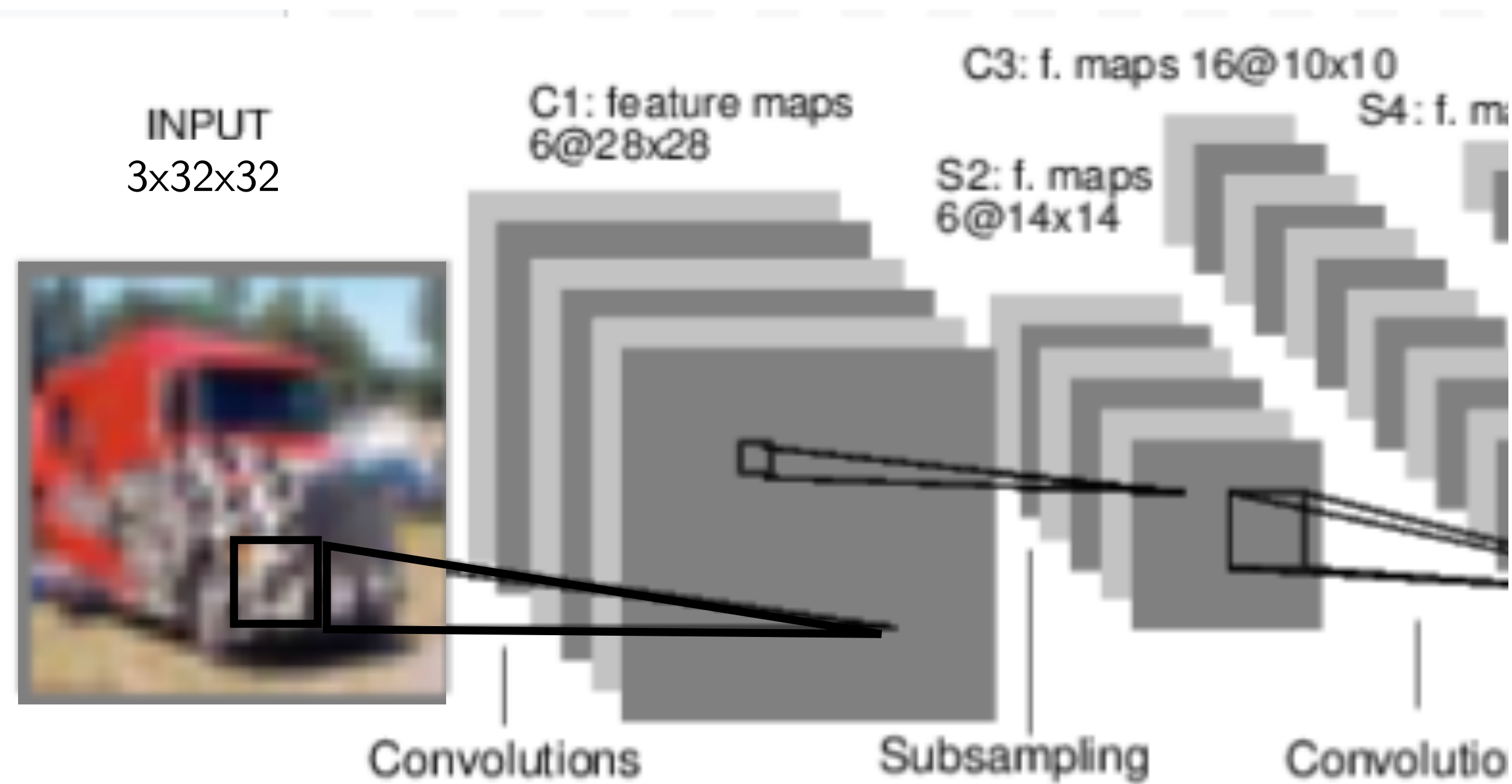


Input:  
Image

Output:  
Logits for class  
(1-10)

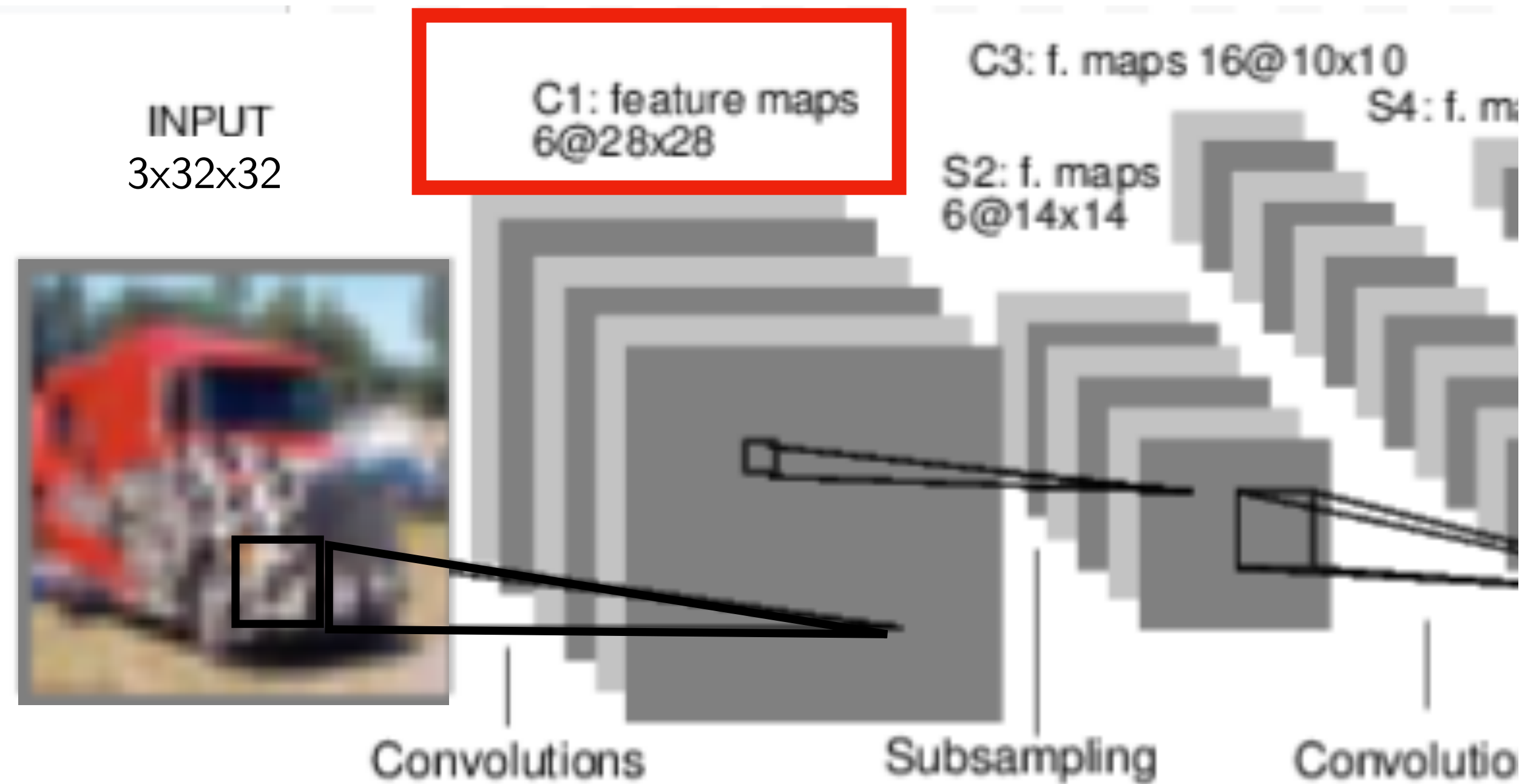
classes = ('plane', 'car', 'bird', 'cat',  
'deer', 'dog', 'frog', 'horse', 'ship', 'truck')

# Model Architecture



```
class Net(nn.Module):  
    def __init__(self):  
        super().__init__()
```

# Model Architecture

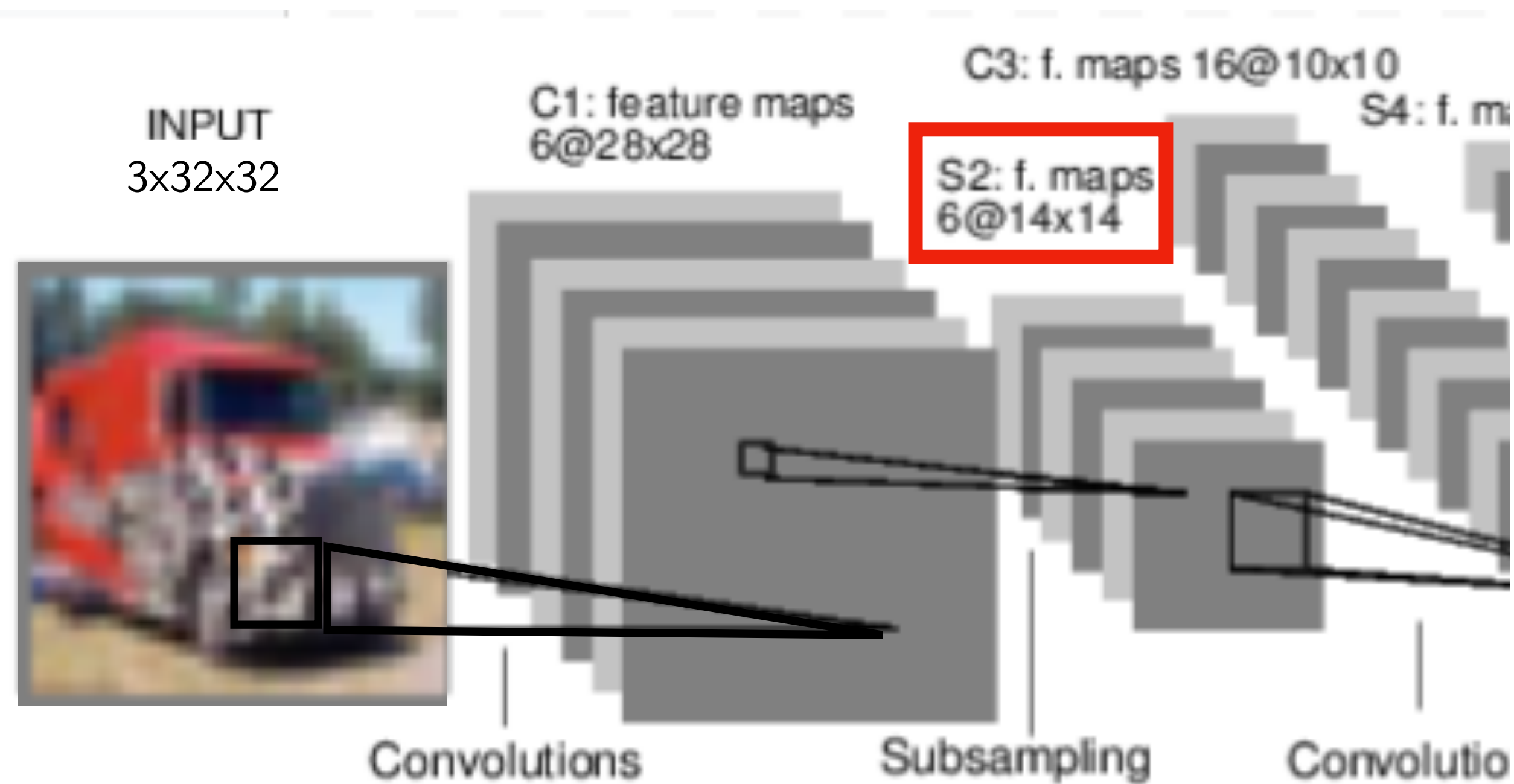


```
class Net(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.conv1 = nn.Conv2d(3, 6, 5)
```

Input Channel    Output Channel    Kernel Size



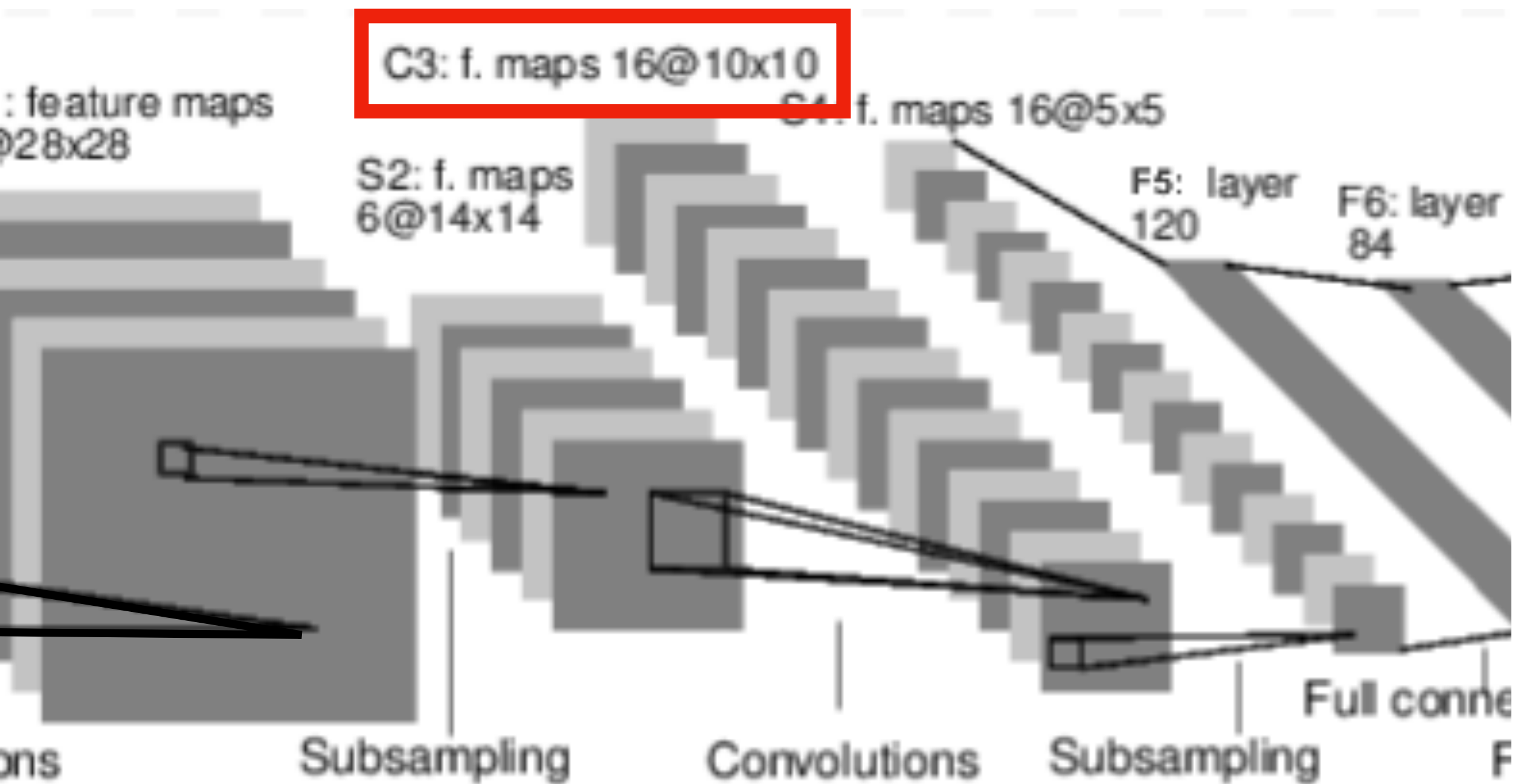
# Model Architecture



```
class Net(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.conv1 = nn.Conv2d(3, 6, 5)  
        self.pool = nn.MaxPool2d(2, 2)
```

Kernel  
Size

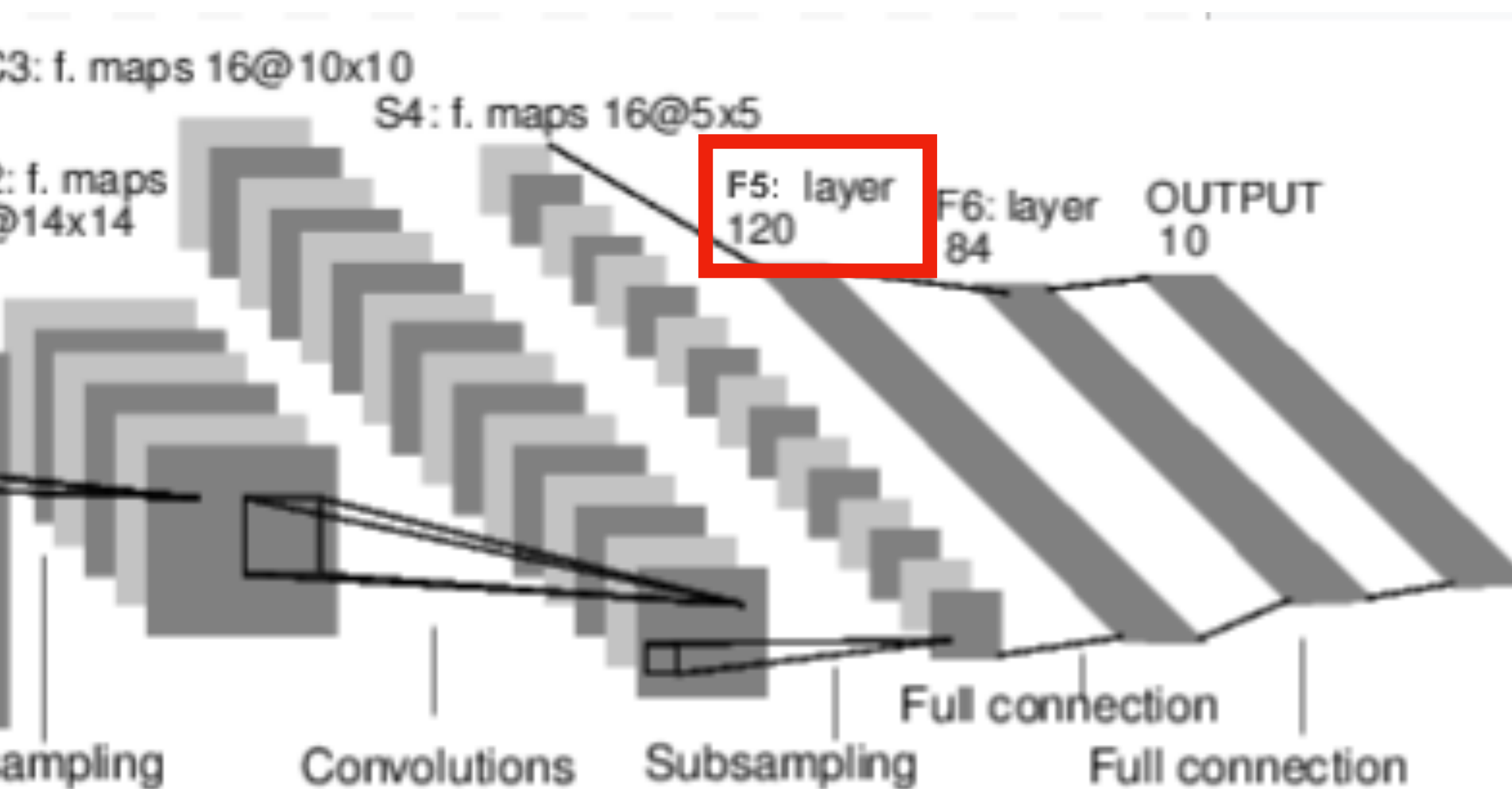
# Model Architecture



```
class Net(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.conv1 = nn.Conv2d(3, 6, 5)  
        self.pool = nn.MaxPool2d(2, 2)  
        self.conv2 = nn.Conv2d(6, 16, 5)
```

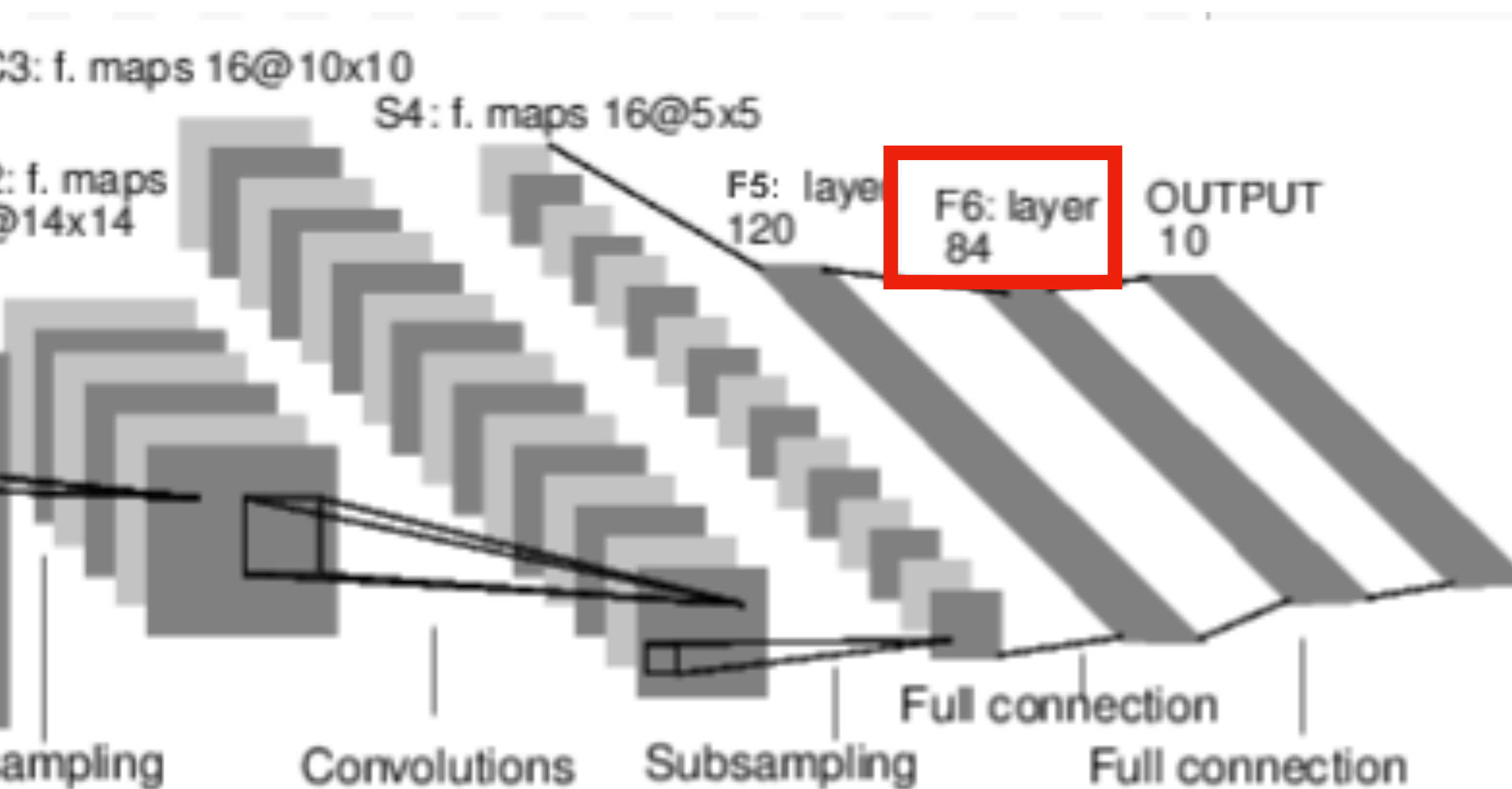
Input Channel      Kernel Size  
Output Channel

# Model Architecture



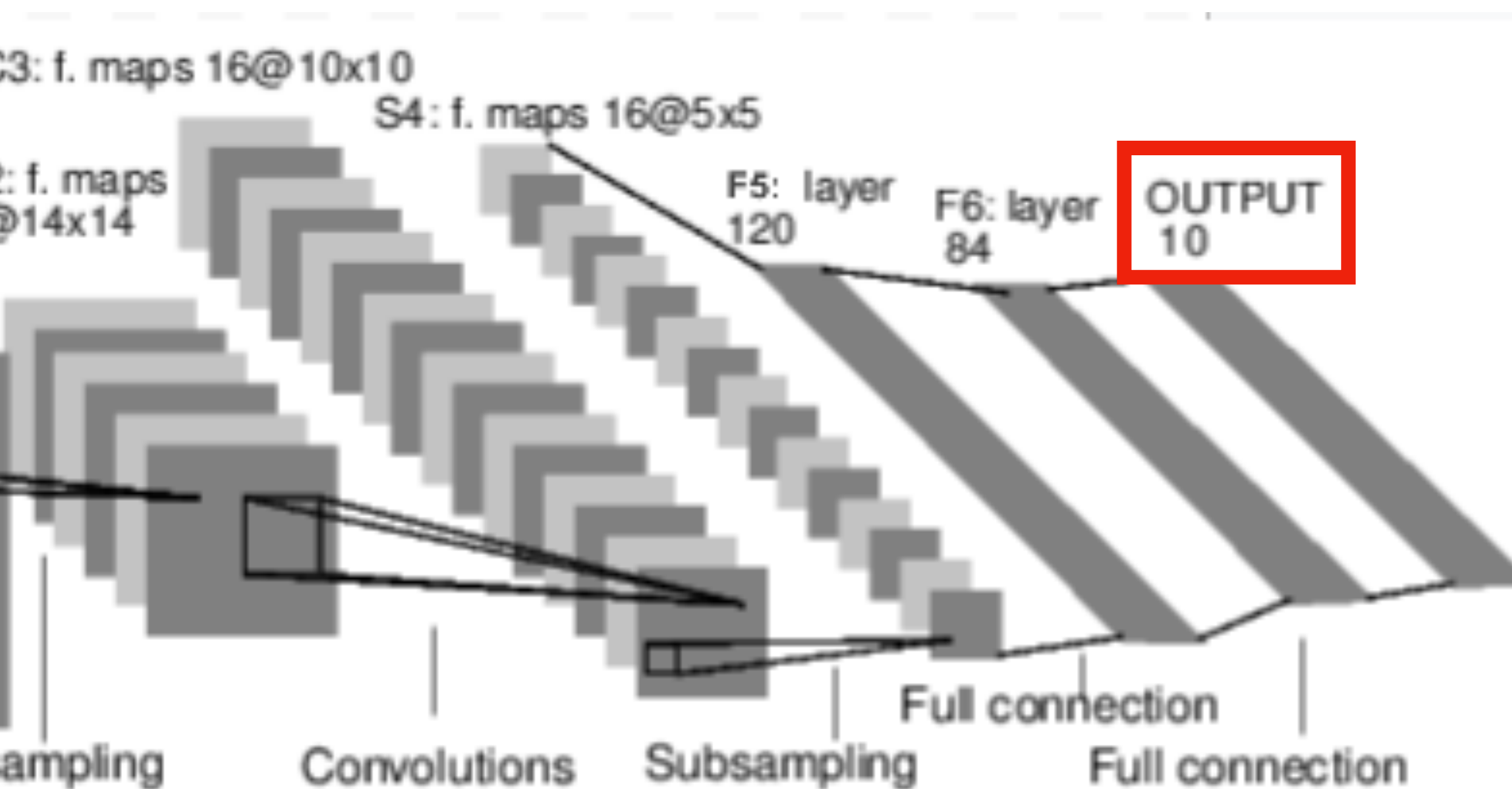
```
class Net(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.conv1 = nn.Conv2d(3, 6, 5)  
        self.pool = nn.MaxPool2d(2, 2)  
        self.conv2 = nn.Conv2d(6, 16, 5)  
        self.fc1 = nn.Linear(16 * 5 * 5,  
                              120)
```

# Model Architecture



```
class Net(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.conv1 = nn.Conv2d(3, 6, 5)  
        self.pool = nn.MaxPool2d(2, 2)  
        self.conv2 = nn.Conv2d(6, 16, 5)  
        self.fc1 = nn.Linear(16 * 5 * 5,  
                              120)  
        self.fc2 = nn.Linear(120, 84)
```

# Model Architecture



```
class Net(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.conv1 = nn.Conv2d(3, 6, 5)  
        self.pool = nn.MaxPool2d(2, 2)  
        self.conv2 = nn.Conv2d(6, 16, 5)  
        self.fc1 = nn.Linear(16 * 5 * 5,  
                              120)  
        self.fc2 = nn.Linear(120, 84)  
        self.fc3 = nn.Linear(84, 10)
```

Go through rest of the tutorial,  
train a classifier, and try to improve  
accuracy  
(Hint: Increase the width of the  
network)



[https://pytorch.org/tutorials/  
beginner/blitz/cifar10\\_tutorial.html](https://pytorch.org/tutorials/beginner/blitz/cifar10_tutorial.html)

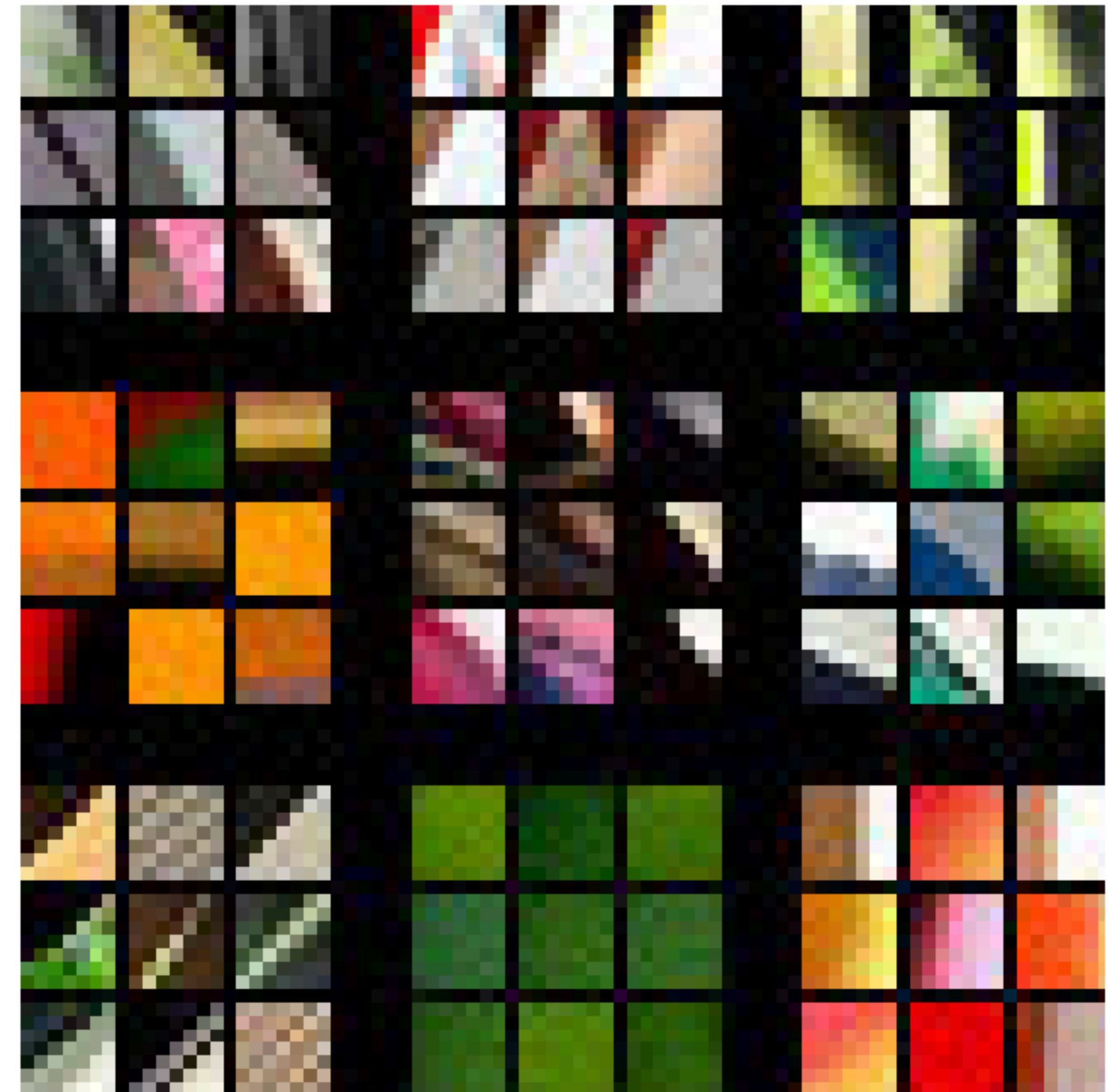
# Visualizing and Understanding CNNs

[Zeiler and Fergus, 2014]

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

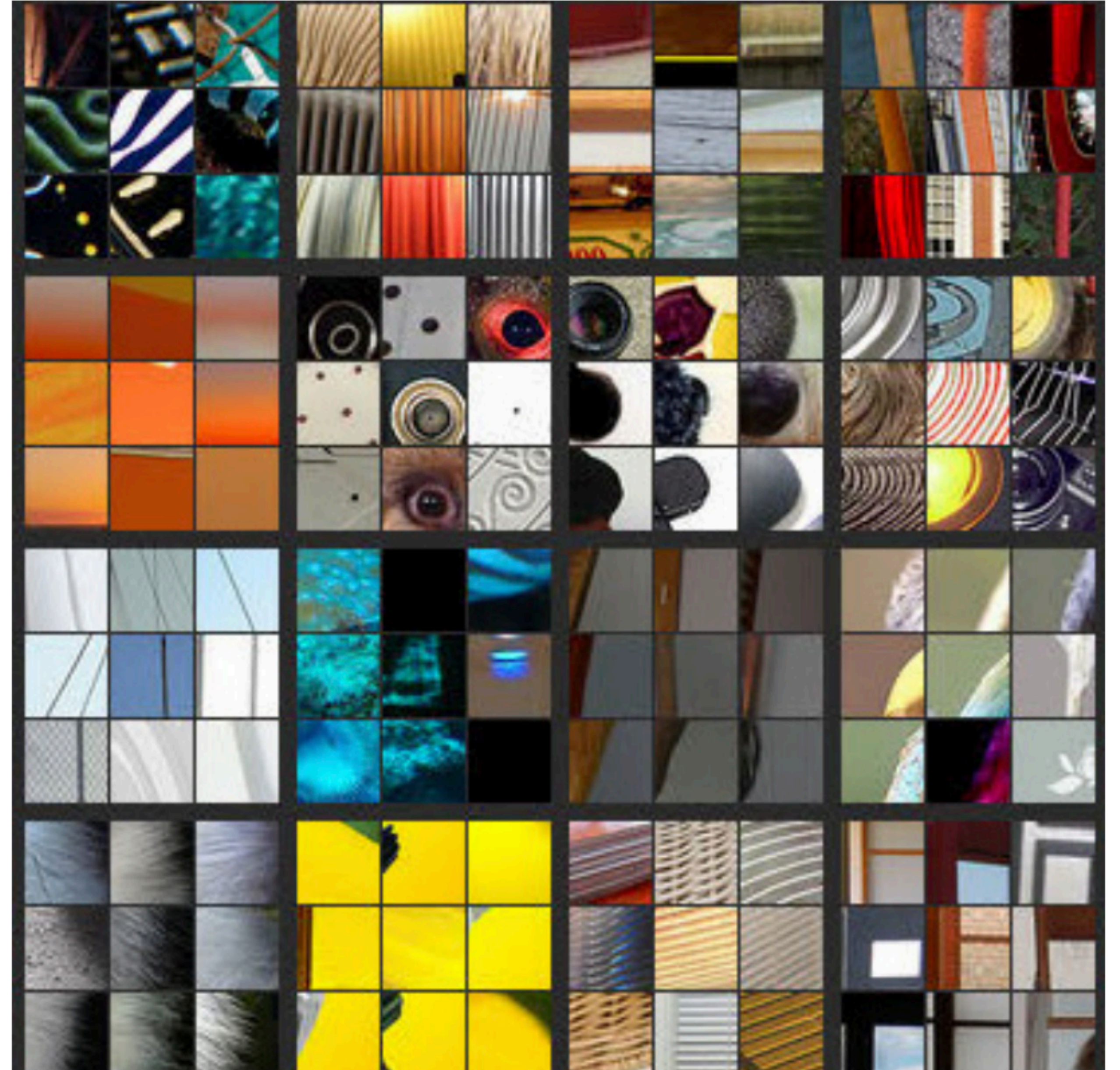


Image patches that activate each of the **layer 1** filters most strongly



[Zeiler and Fergus, 2014]

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





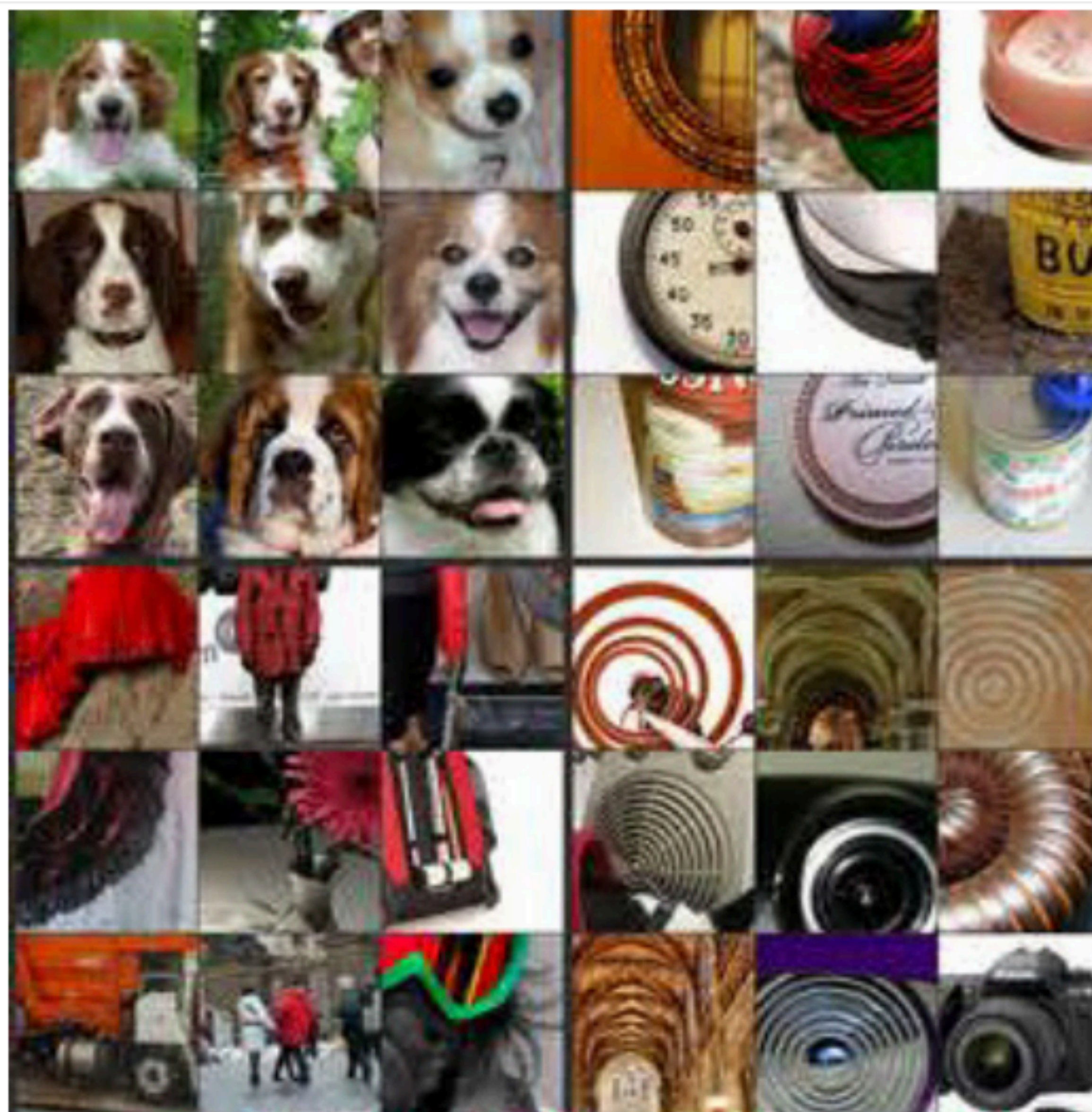
[Zeiler and Fergus, 2014]

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



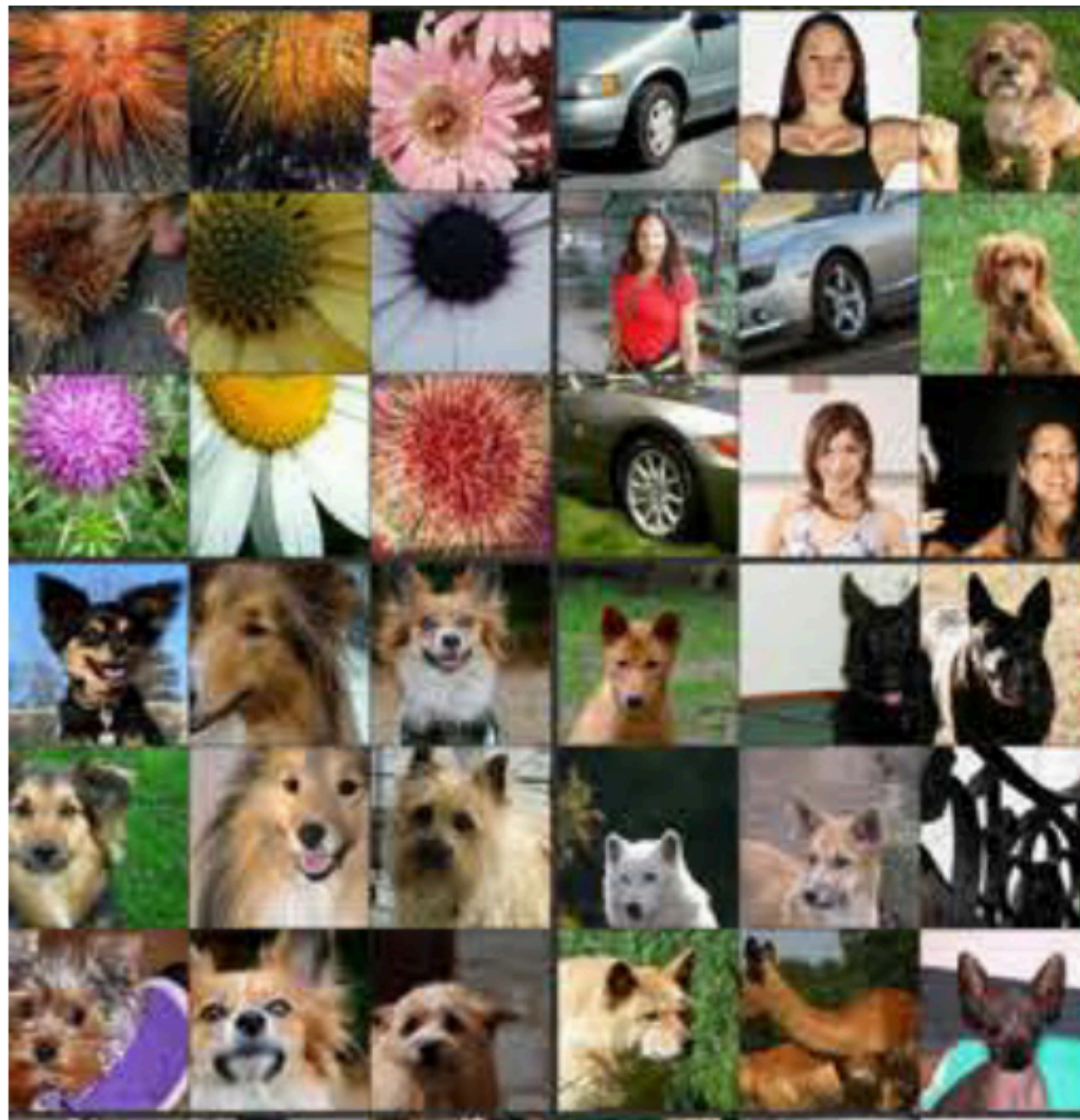
[Zeiler and Fergus, 2014]

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

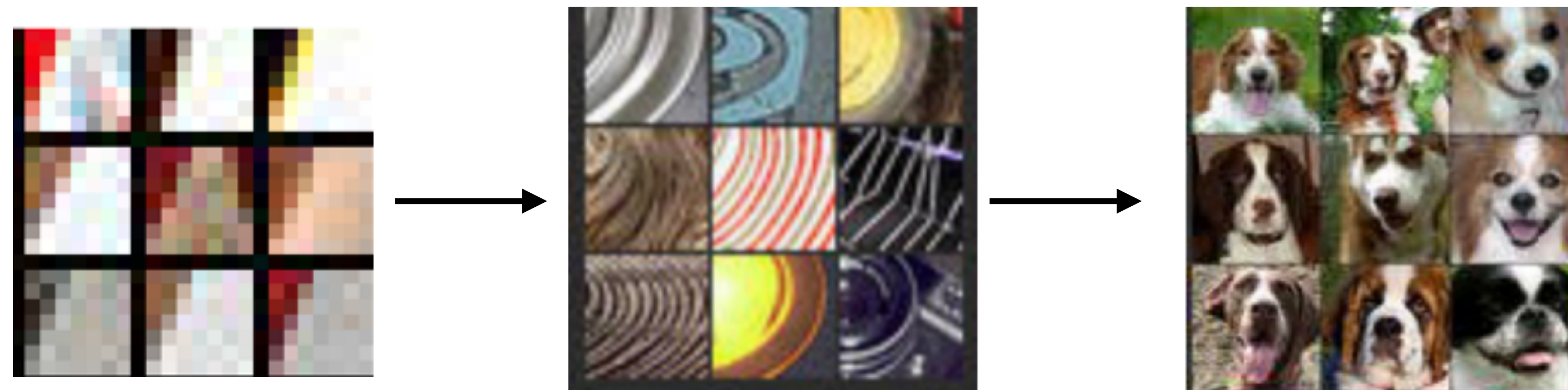
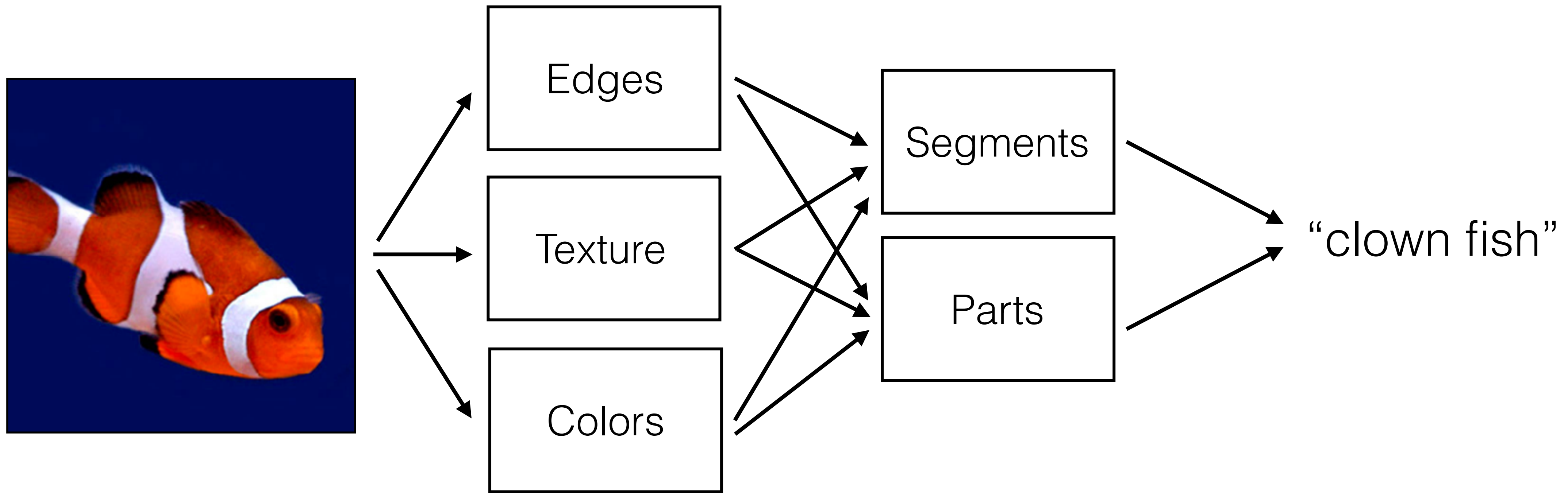


[Zeiler and Fergus, 2014]

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



# CNNs *learned* the classical visual recognition pipeline!



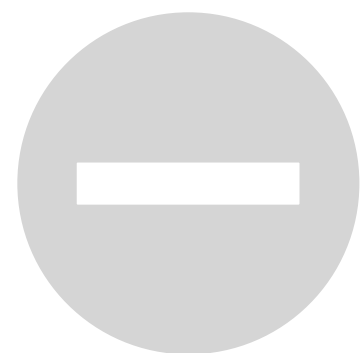
# Today's Class



What is a **visual** representation?  
What makes for a **good** representation?



How do we learn a visual representation?



How do we train a representation by  
**unsupervised** learning and **self-supervised** learning?

# Today's Class



What is a **visual** representation?  
What makes for a **good** representation?



How do we learn a visual representation?



How do we train a representation by  
**unsupervised** learning and **self-supervised** learning?

# Supervised object recognition

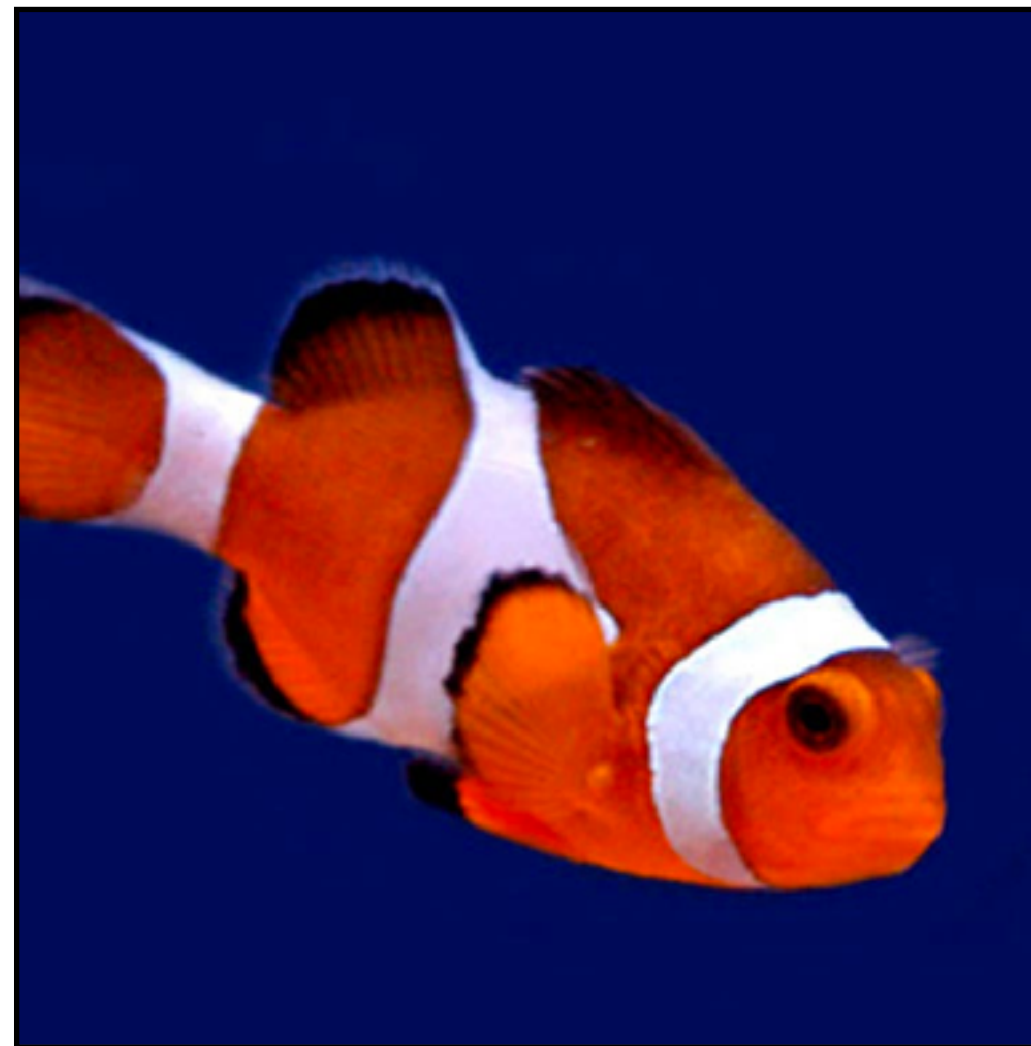
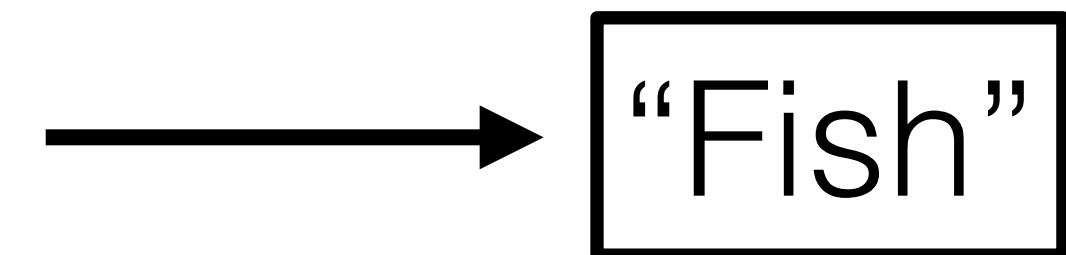


image X

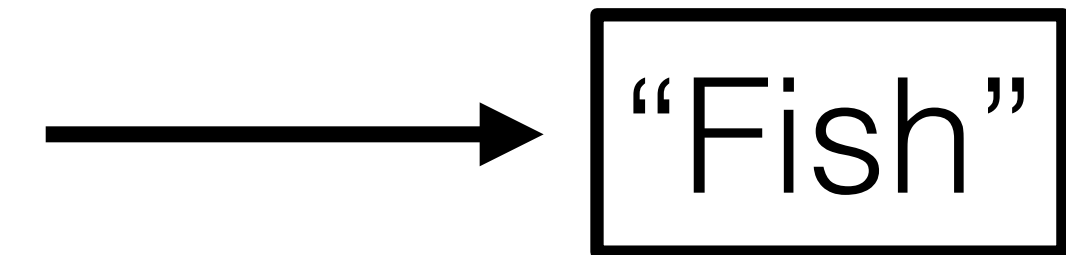


label Y

# Supervised object recognition



image X



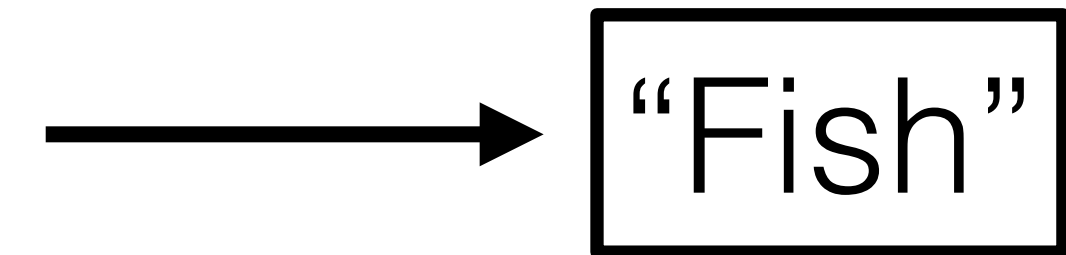
label Y



# Supervised object recognition

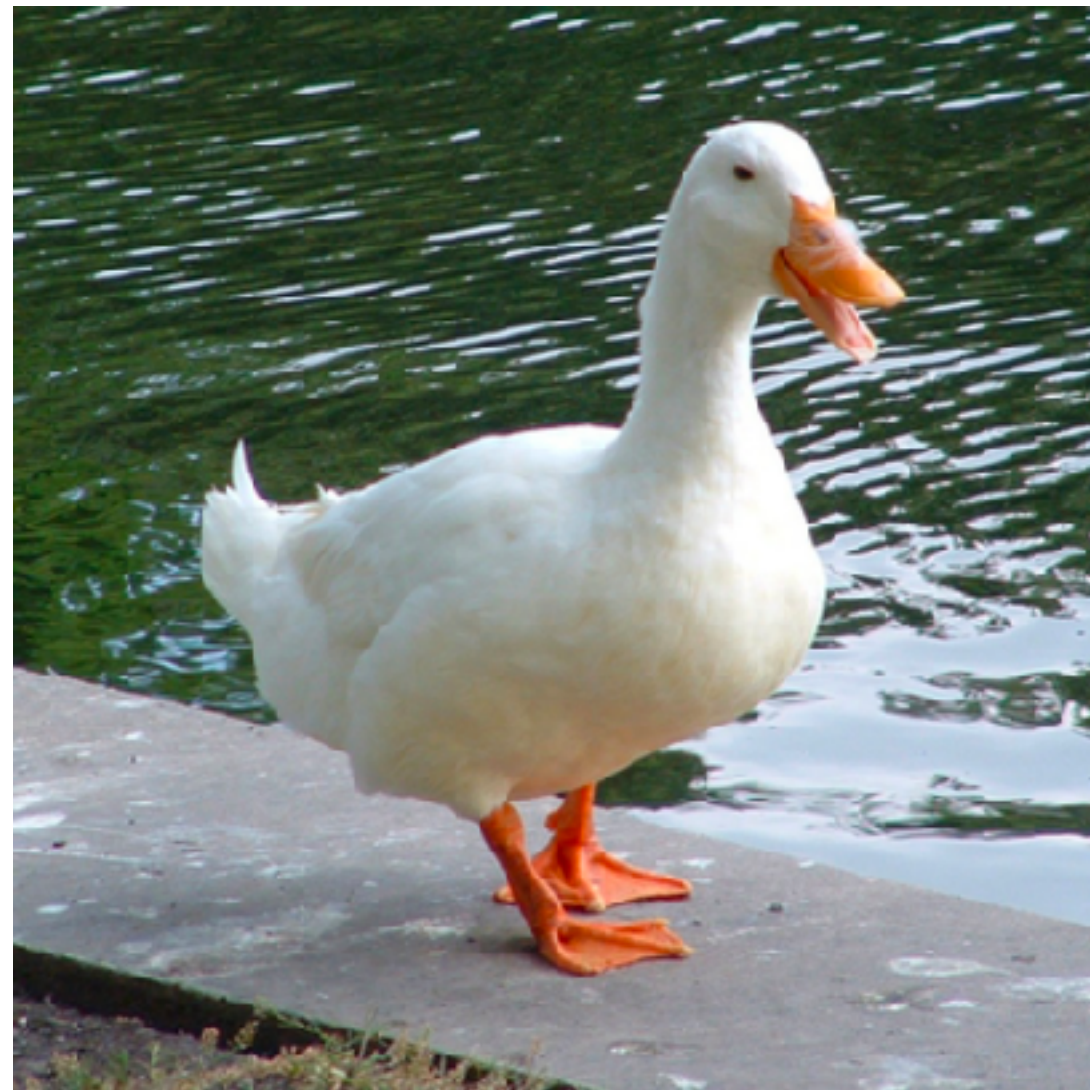


image X



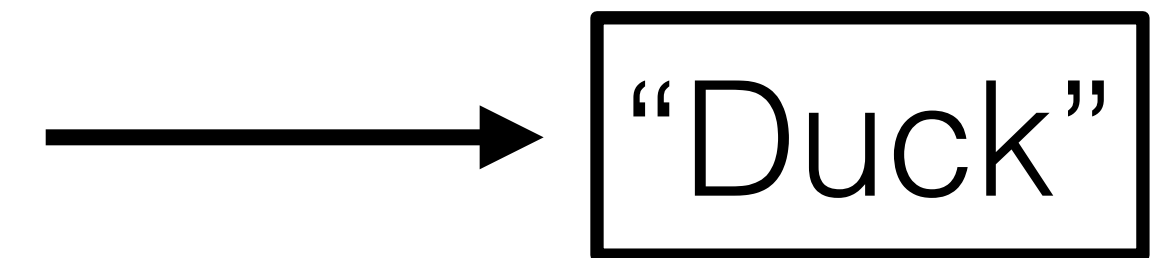
label Y

# Supervised object recognition



⋮

image X



label Y

# The real world is not labeled!



Imagine a robot having to clean this kitchen.

How does it learn about all the objects in the scene?

# 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 labels

(aka **supervised learning**)

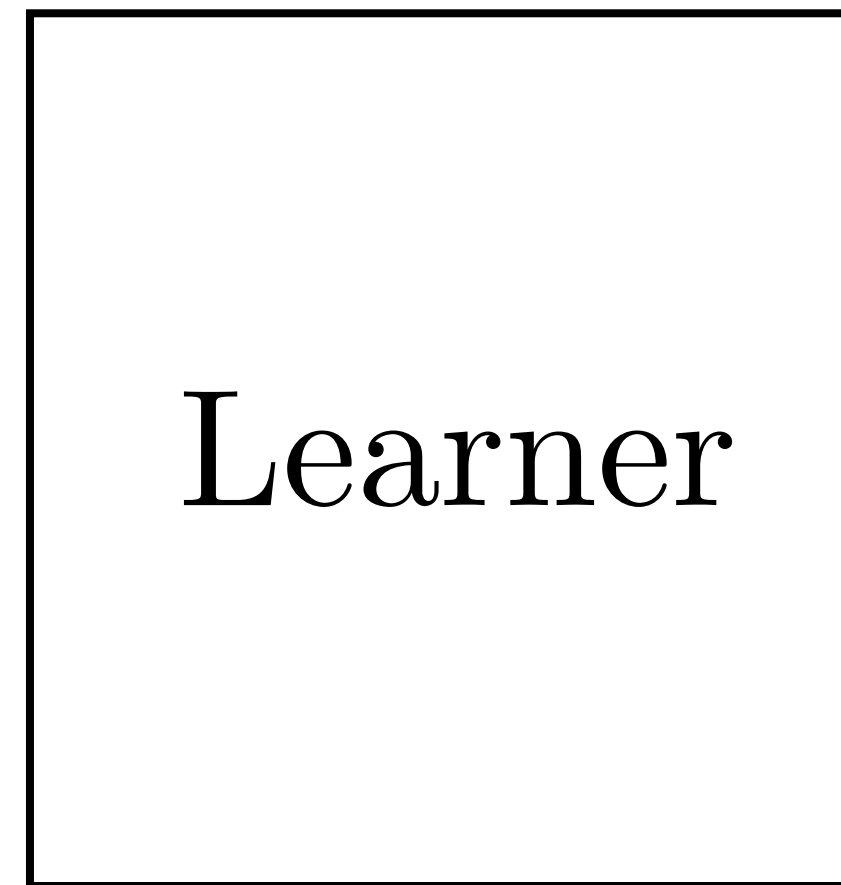
Training data

$\{x^{(1)}, y^{(1)}\}$

$\{x^{(2)}, y^{(2)}\}$

$\{x^{(3)}, y^{(3)}\}$

...



$\rightarrow f : X \rightarrow Y$

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

# Learning without labels

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

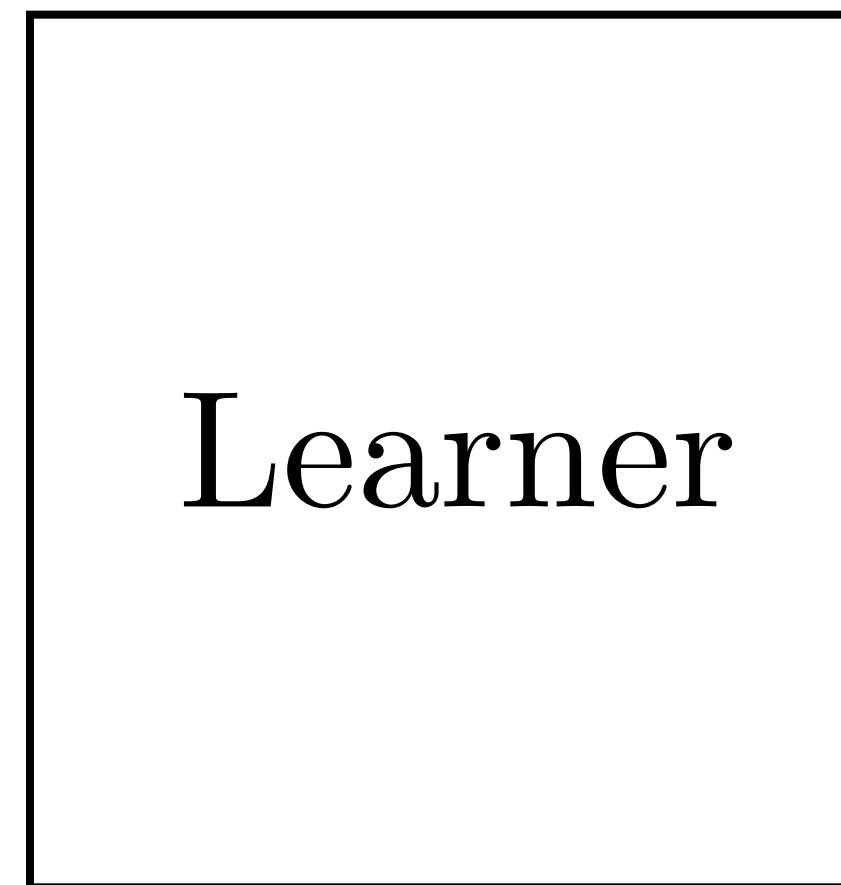
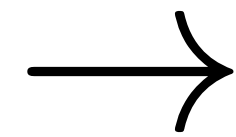
Data

$\{x^{(1)}\}$

$\{x^{(2)}\}$

$\{x^{(3)}\}$

...



?

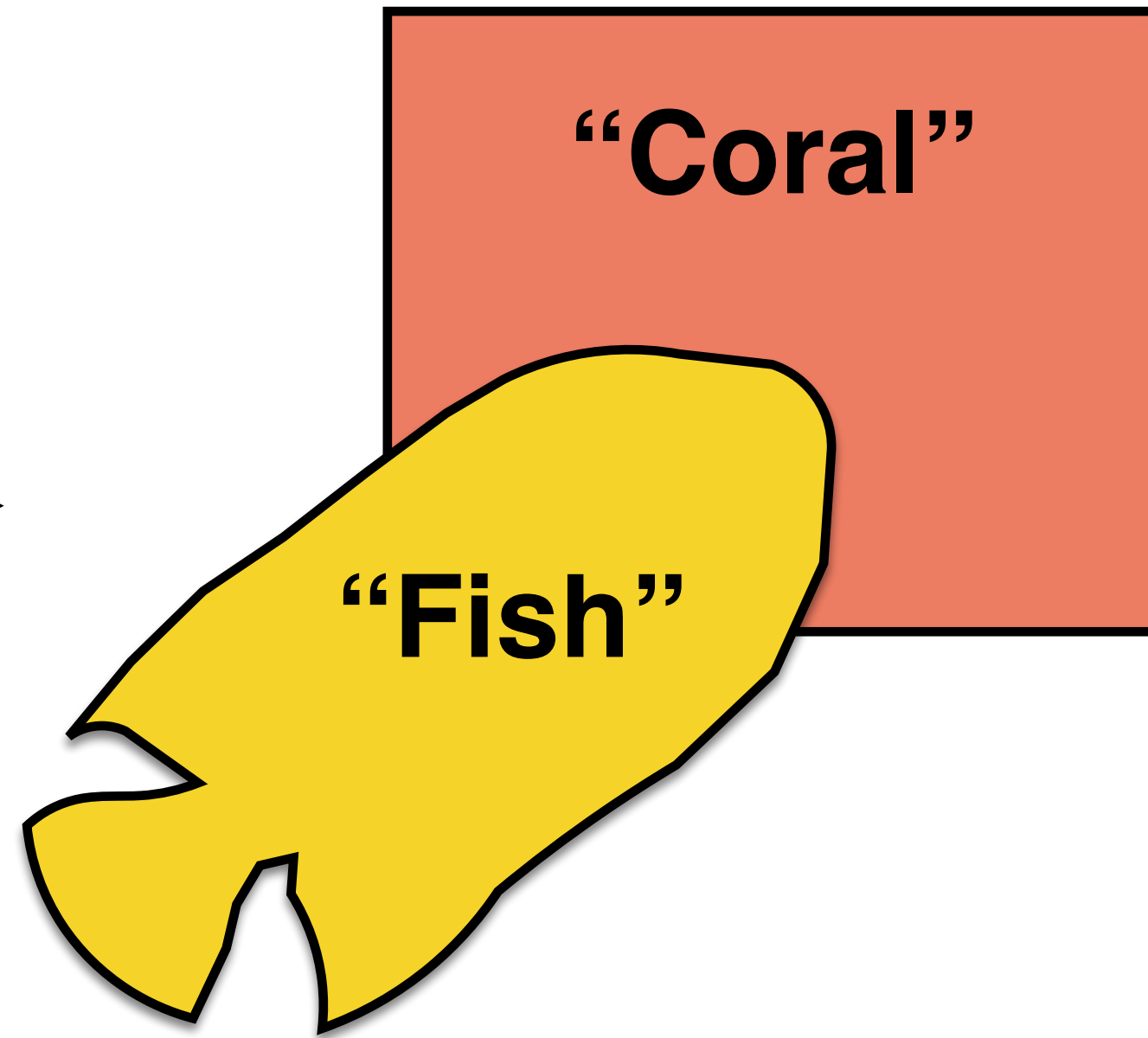
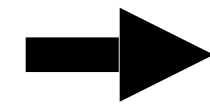
# Unsupervised Learning

# Unsupervised Representation Learning

$X$



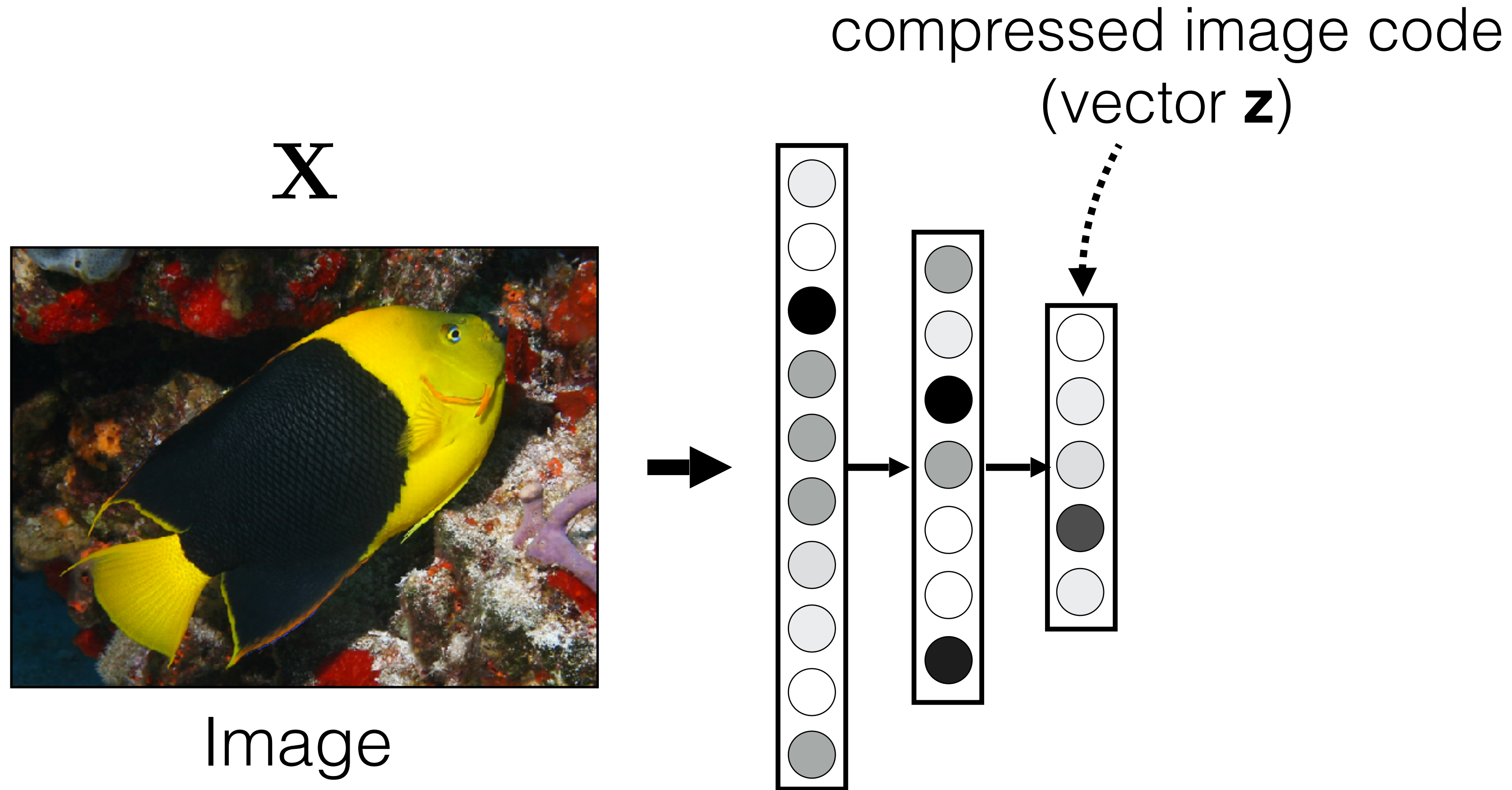
Image



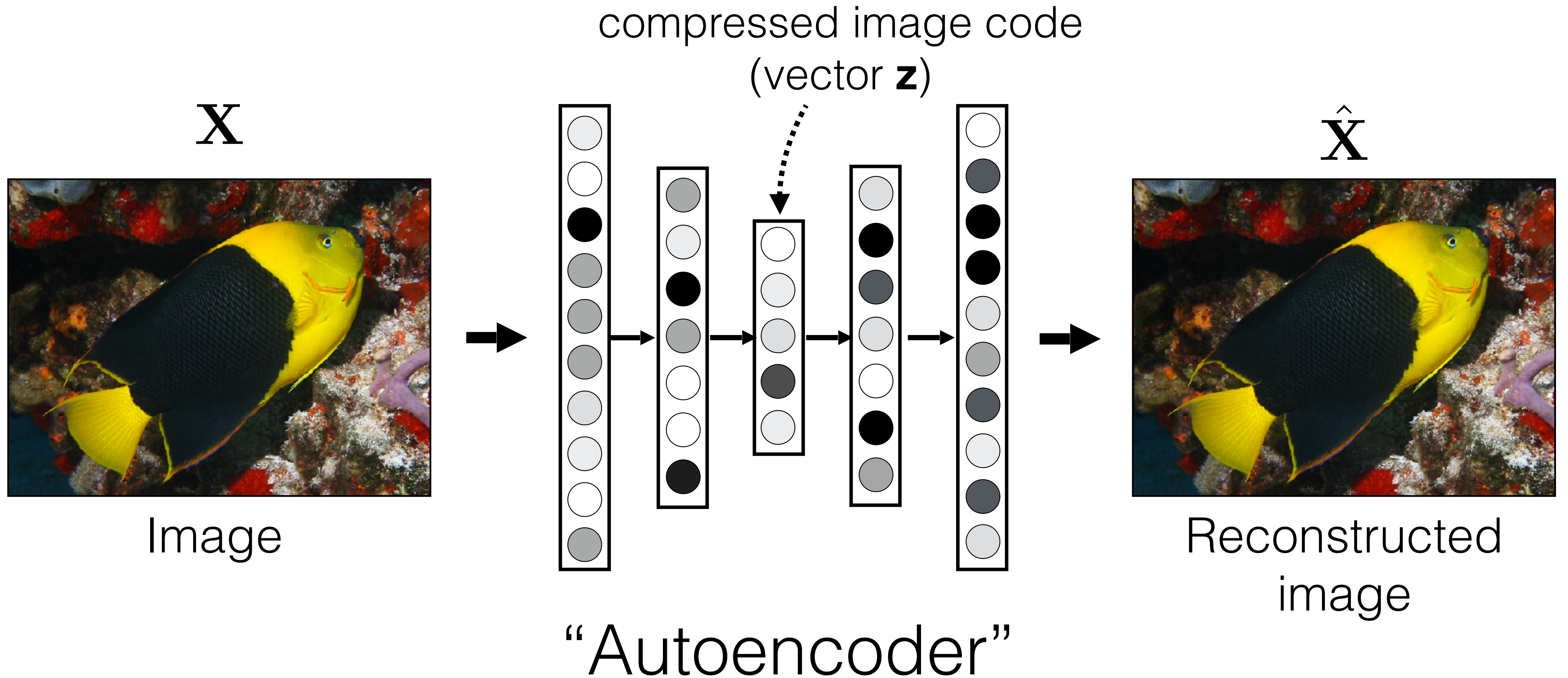
Compact mental  
representation



# Unsupervised Representation Learning



# Unsupervised Representation Learning



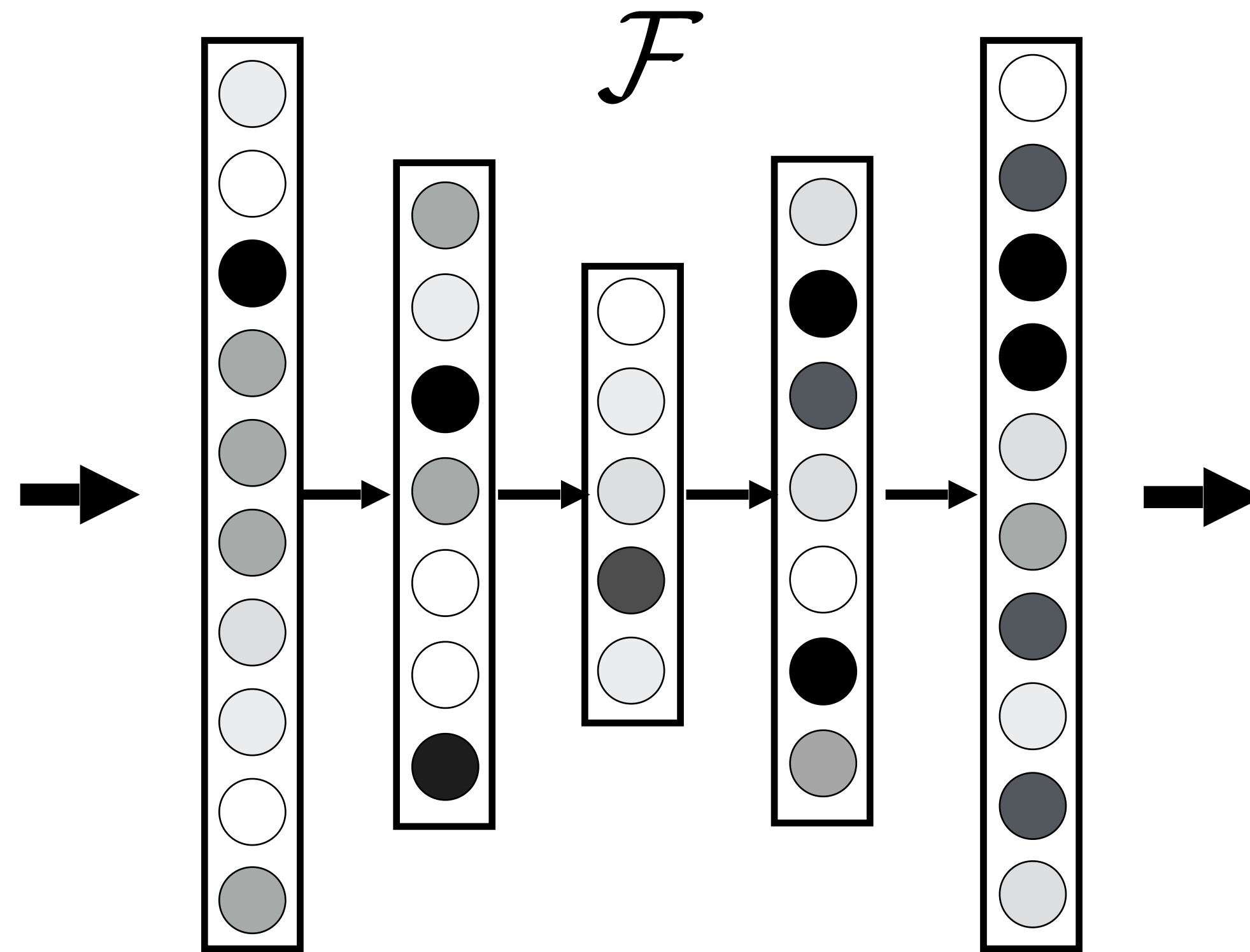
[e.g., Hinton & Salakhutdinov, Science 2006]

# Autoencoder

$\mathbf{X}$



Image



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



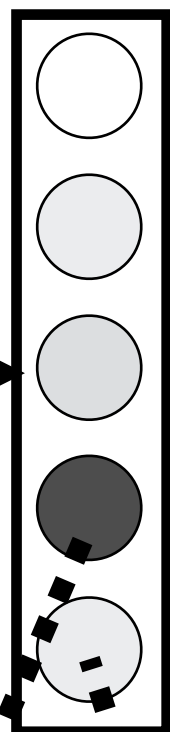
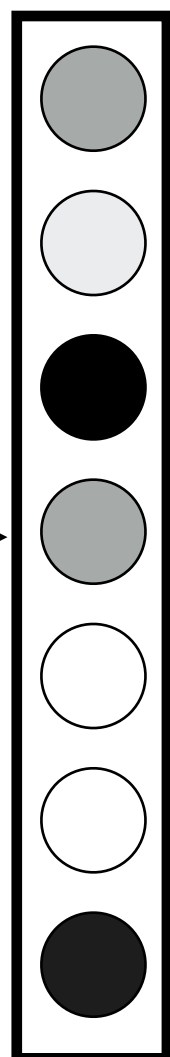
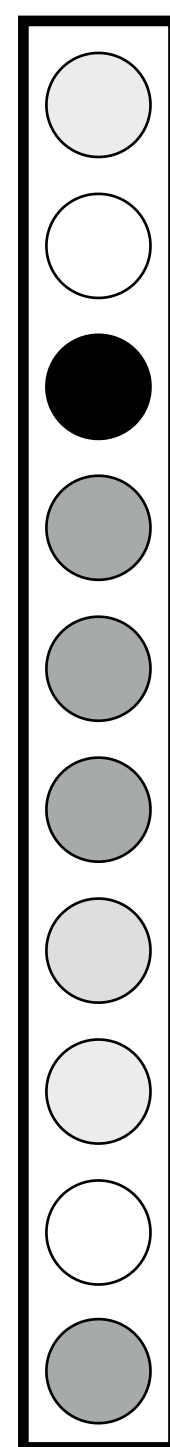
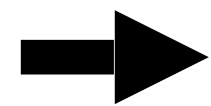
Reconstructed image

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

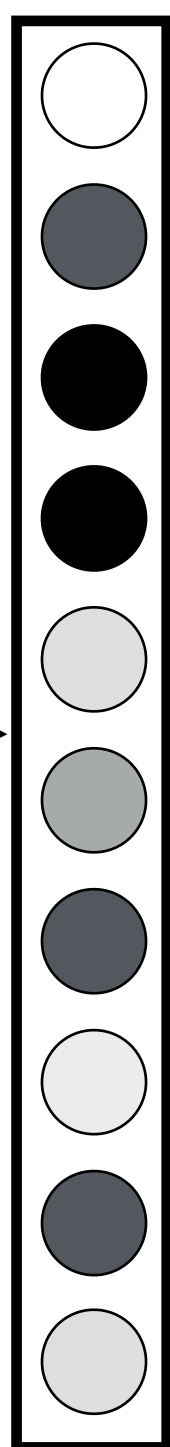
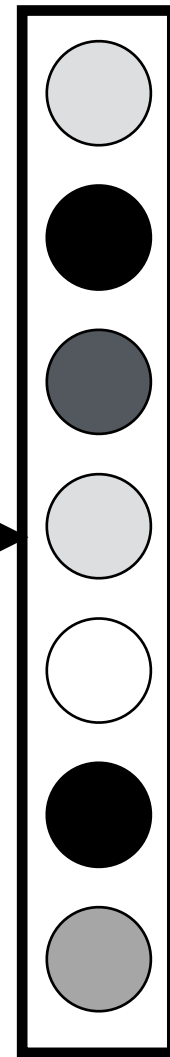
$\mathbf{X}$



Image



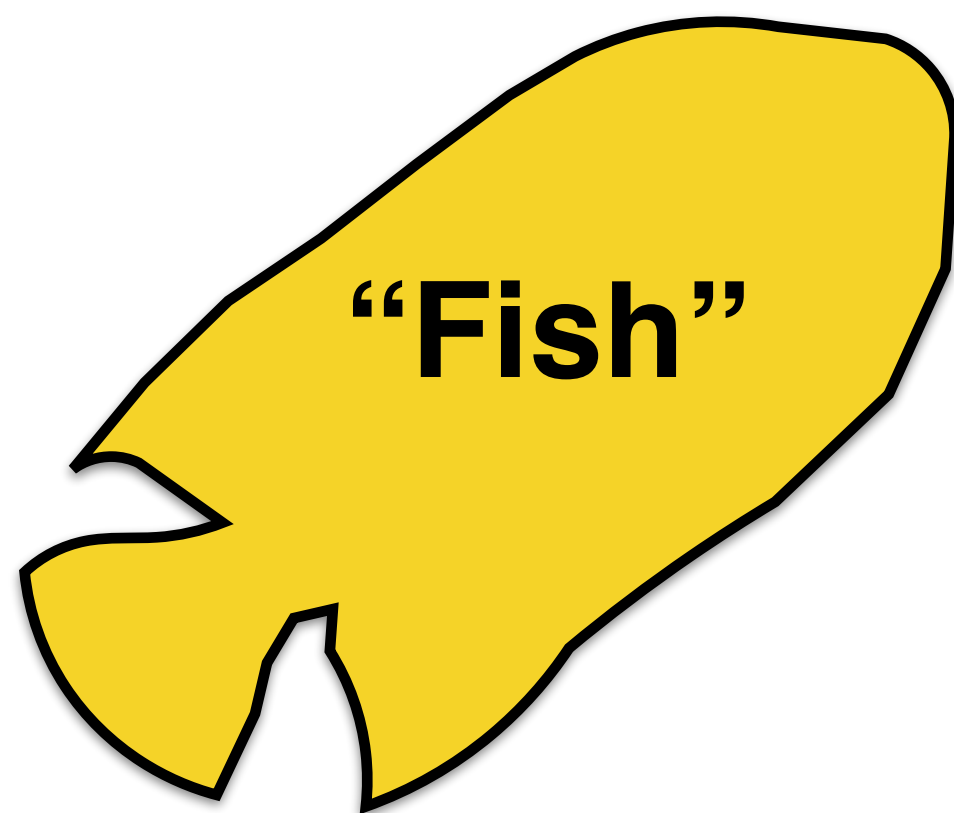
$\mathcal{F}$



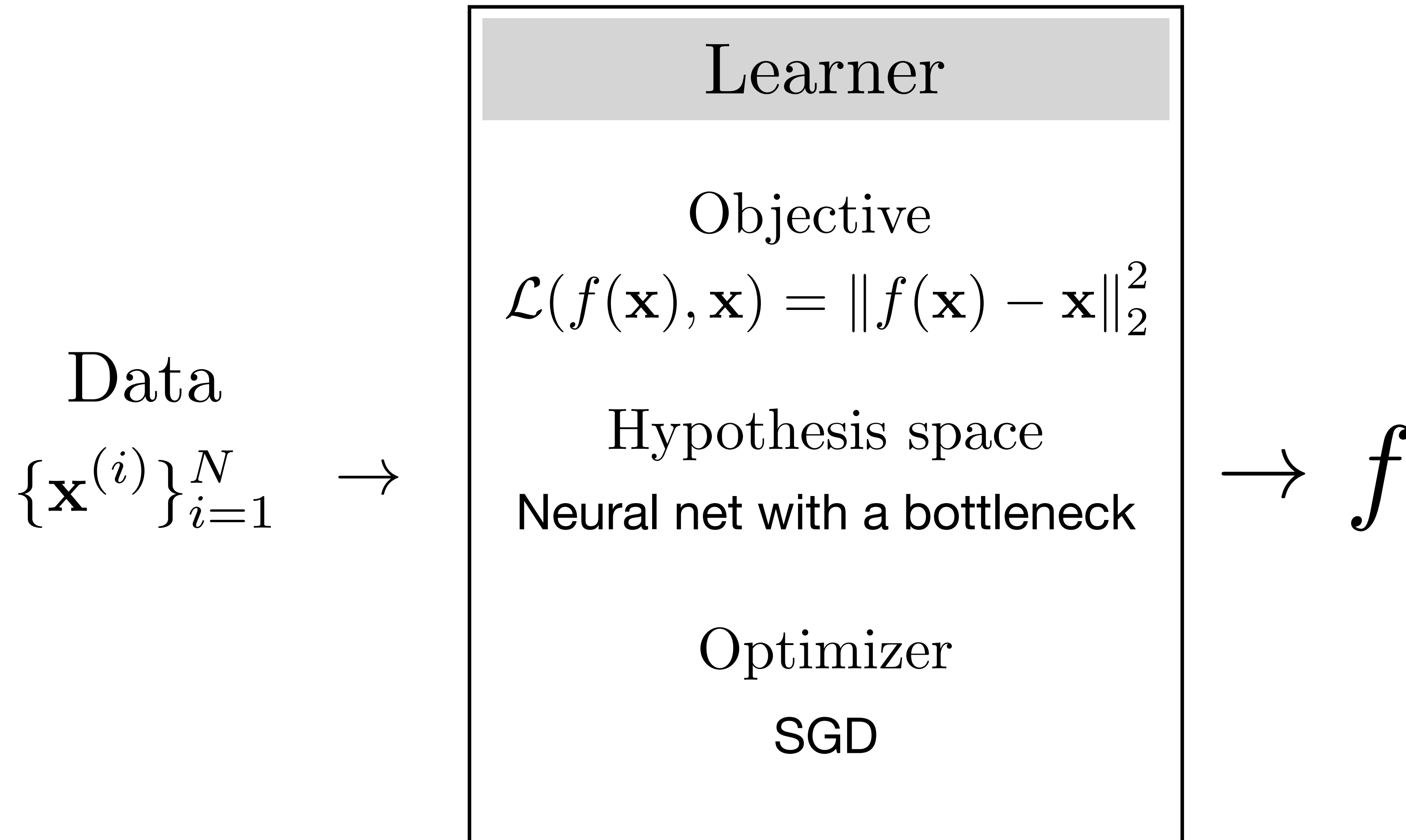
$$\hat{\mathbf{X}} = \mathcal{F}(\mathbf{X})$$



Reconstructed image

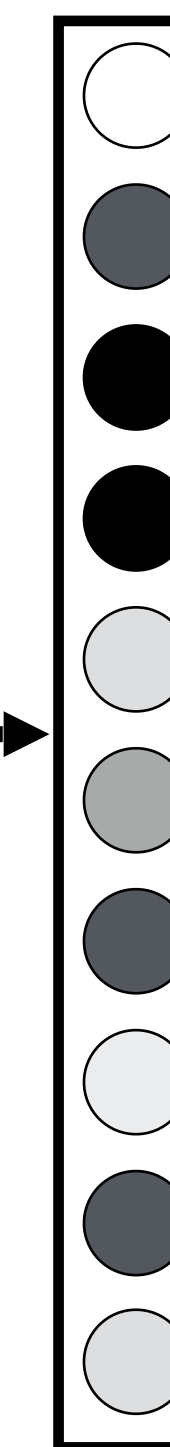
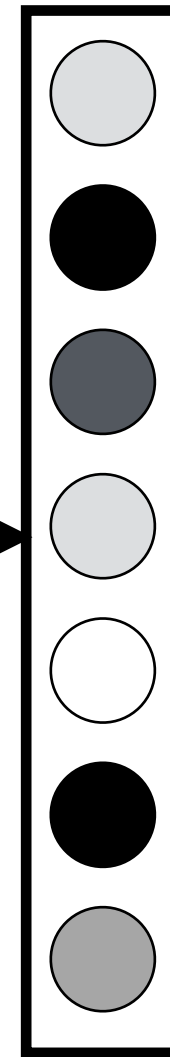
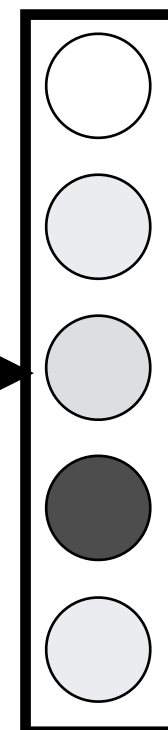
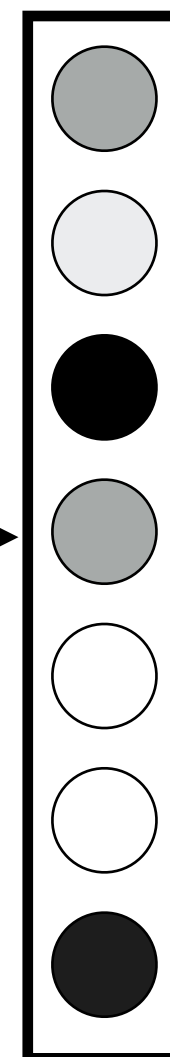
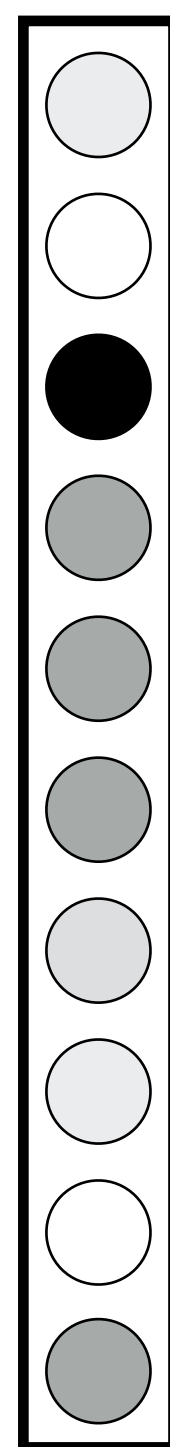
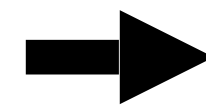


# Autoencoder



$\mathbf{X}$ 

Image



Reconstructed image

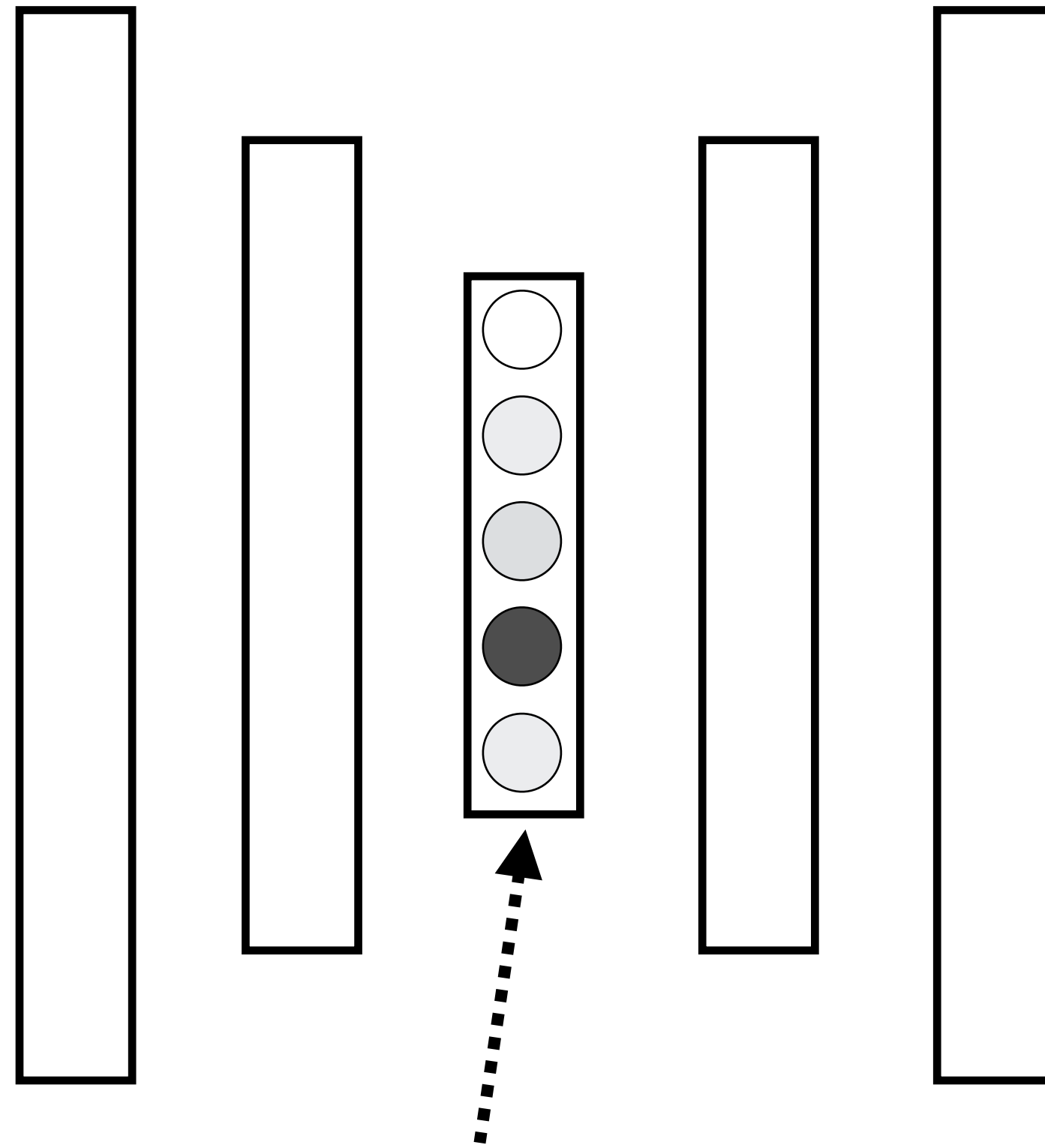
 $\mathcal{F}$ 

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

**Encoder****Decoder**

$\mathbf{X}$ 

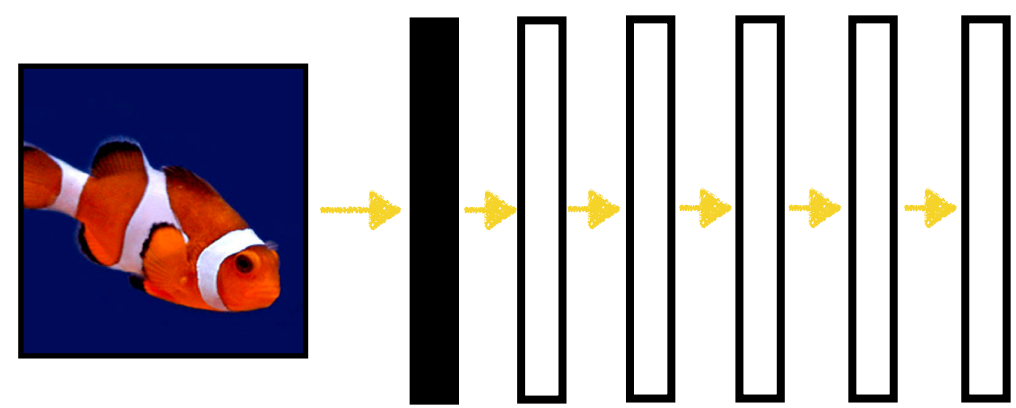
Image

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

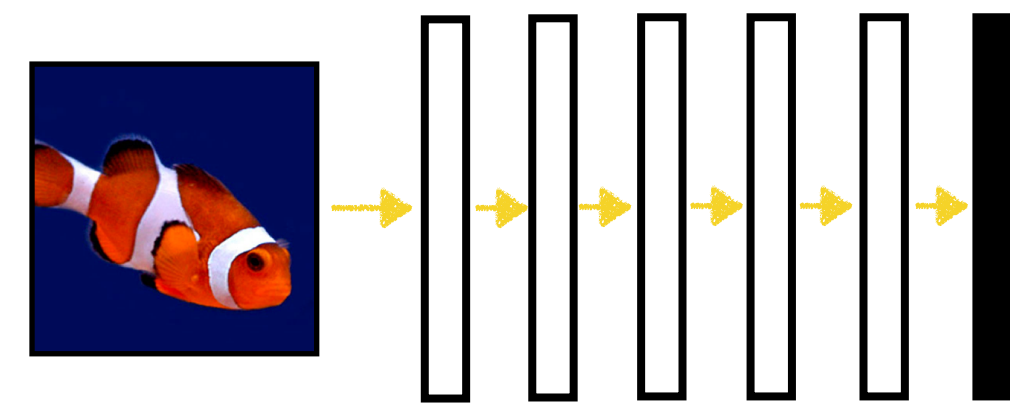
Is the code informative about  
object class  $y$ ?

Logistic regression:

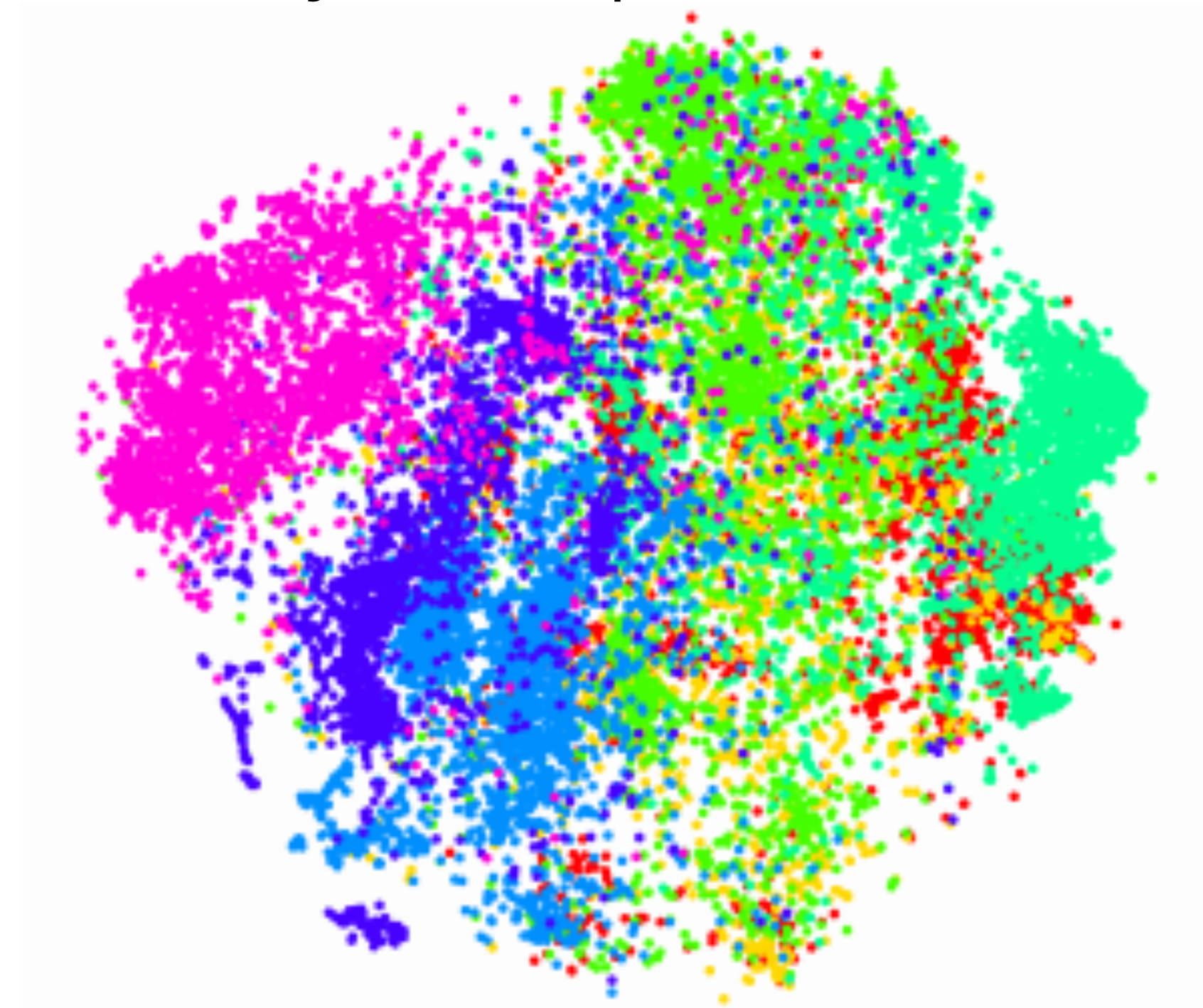
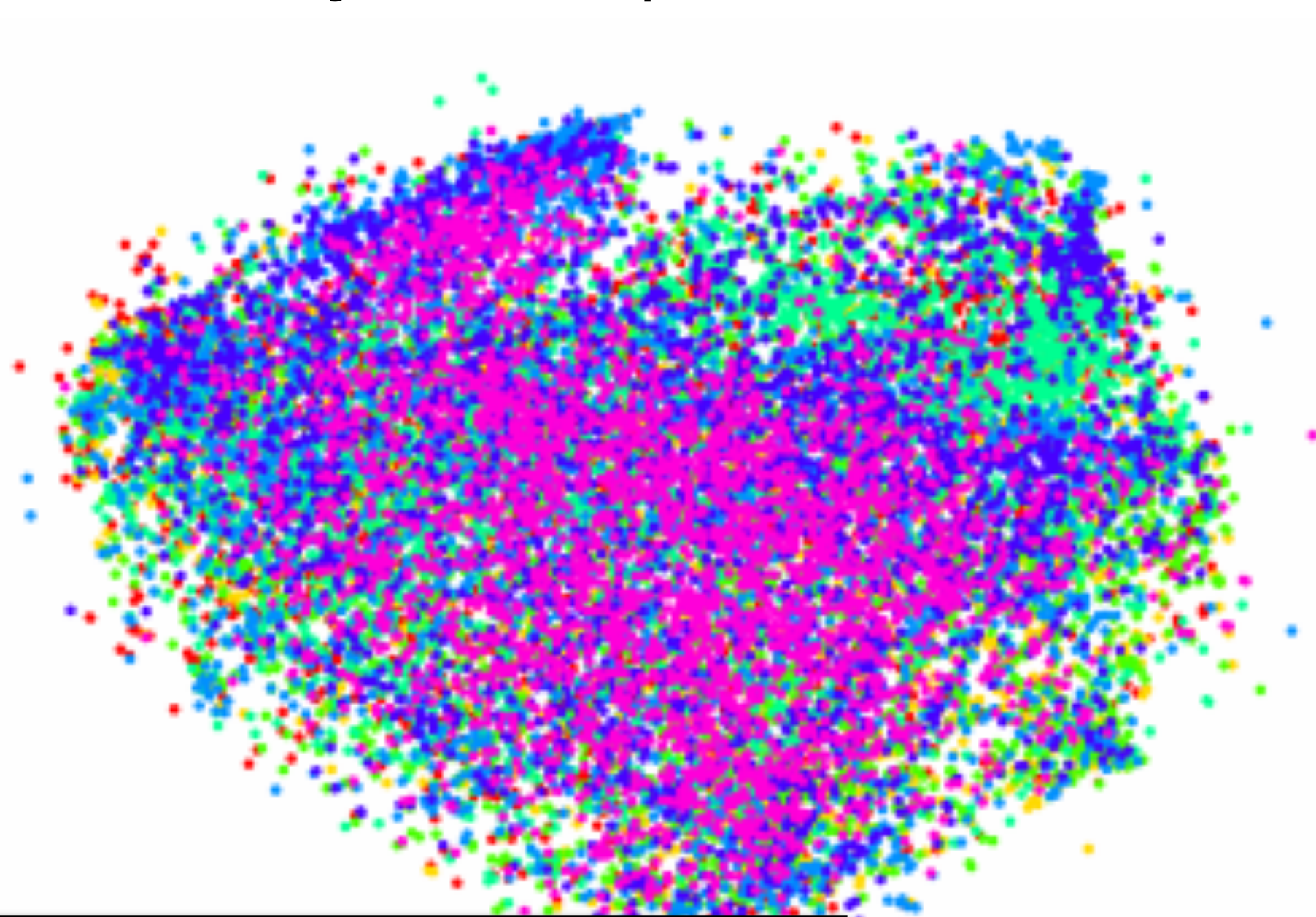
$$y = \sigma(\mathbf{W}\mathbf{z} + \mathbf{b})$$



Layer 1 representation



Layer 6 representation



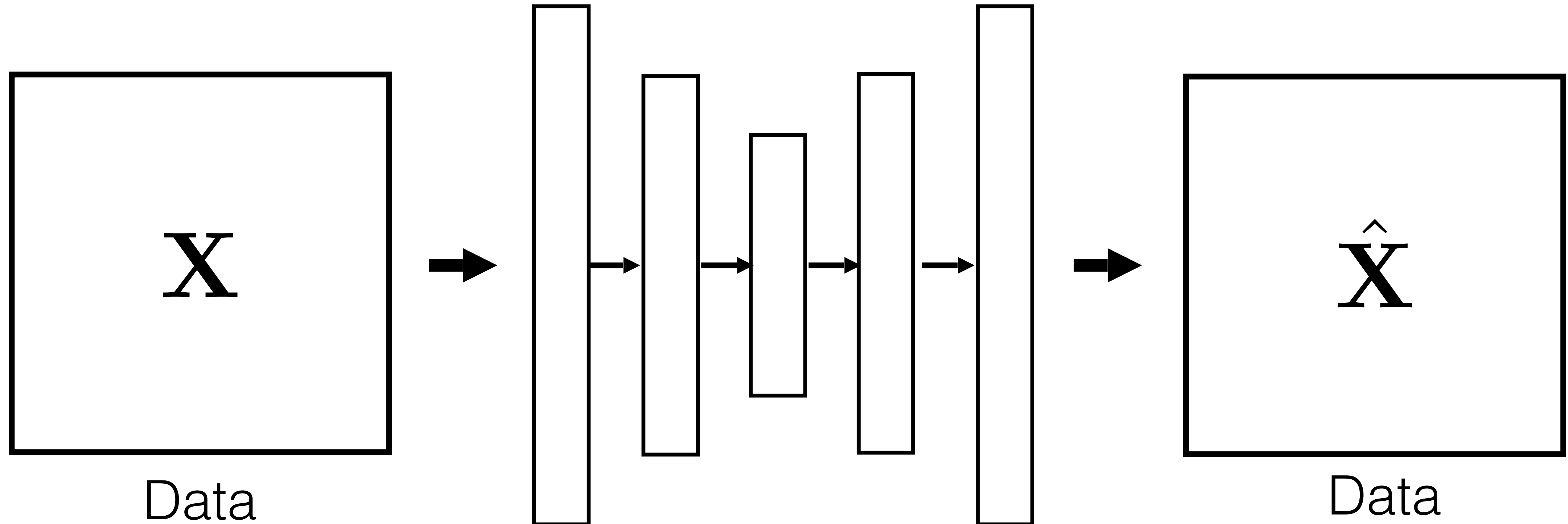
- structure, construction
- covering
- commodity, trade good, good
- conveyance, transport
- invertebrate
- bird
- hunting dog

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

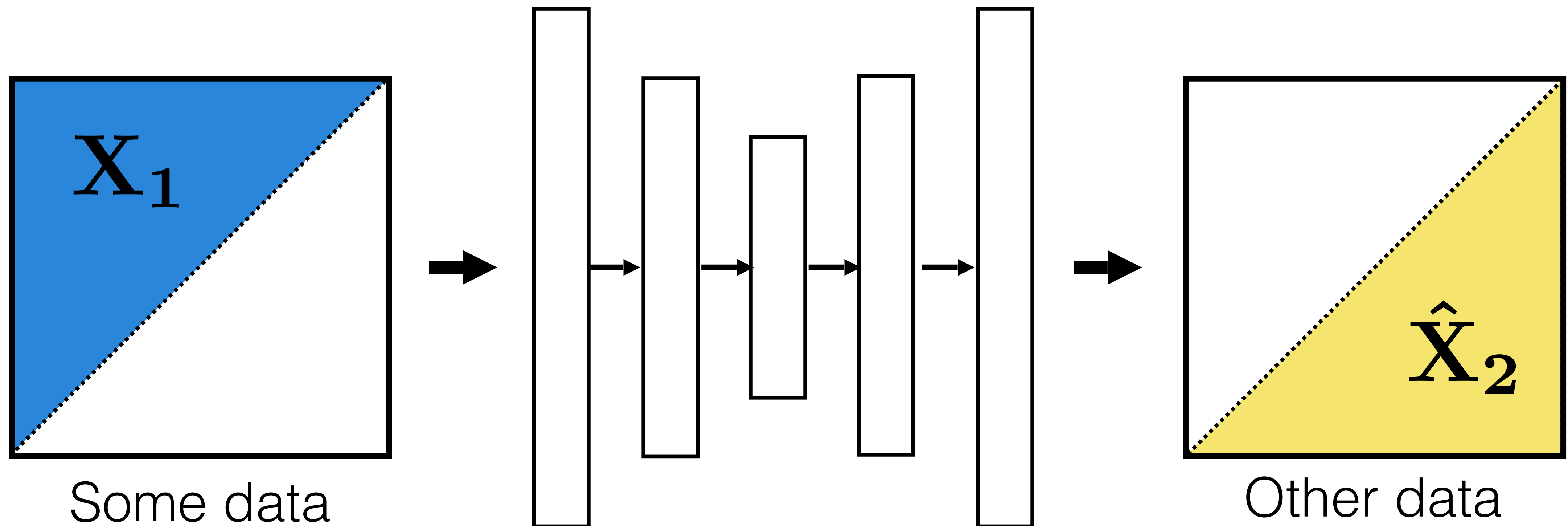
[Visualization technique : t-sne, van der Maaten & Hinton, 2008]



# Data compression



# Data prediction aka “self-supervised learning”



# Today's Class



What is a **visual** representation?  
What makes for a **good** representation?



How do we learn a visual representation?



How do we train a representation by  
**unsupervised** learning and **self-supervised** learning?

# Why learn representations?

## Reinforcement Learning (Cherry)

Predicting a scalar reward given once in a while

A few bits for some samples

## Supervised Learning (Chocolate Coat)

Predicting category or vector of scalars per input as provided by human labels.  
10-10k bits per sample

## Unsupervised / Self-Supervised Learning (Cake)

Predicting parts of observed input or predicting future observations or events  
Millions of bits per sample



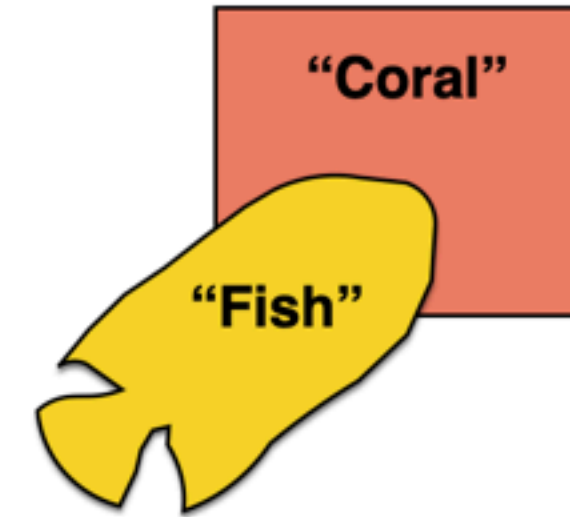
Visualisation Idea by Yann LeCun  
Photo by [Kristina Paukshtite](#) from [Pexels](#)

# tl;dr

## Representation learning

Good representations are:

1. Compact (*minimal*)
2. Explanatory (*sufficient*)
3. Disentangled (*independent factors*)
4. Hierarchical (*feature reuse*)

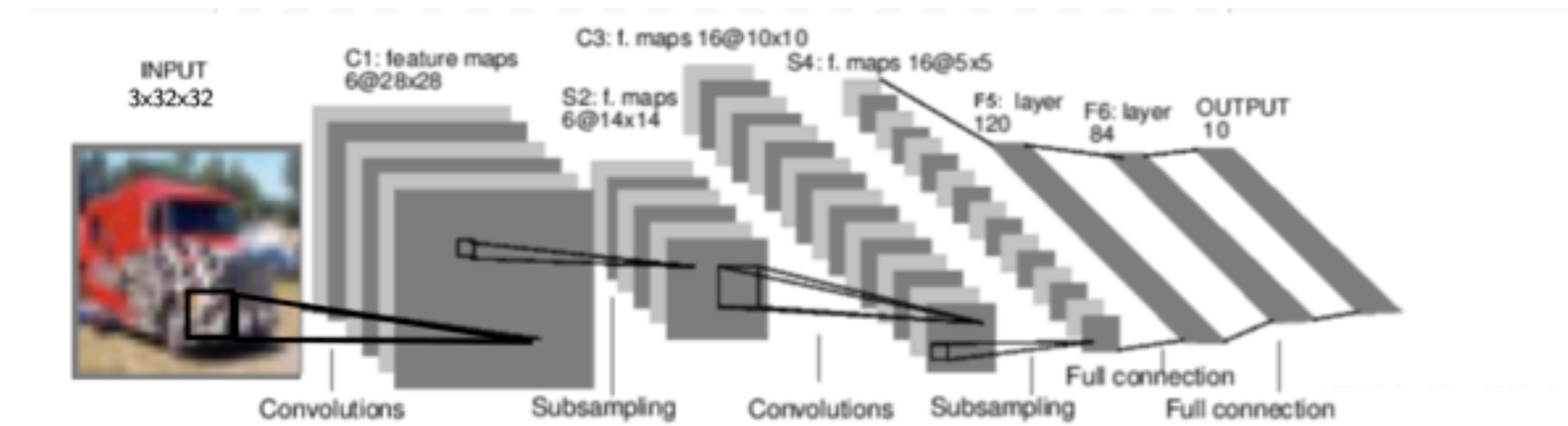


5. *Make subsequent problem solving easy*

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

From MIT 6.8300/6.8301: Advances in Computer Vision

## Model Architecture



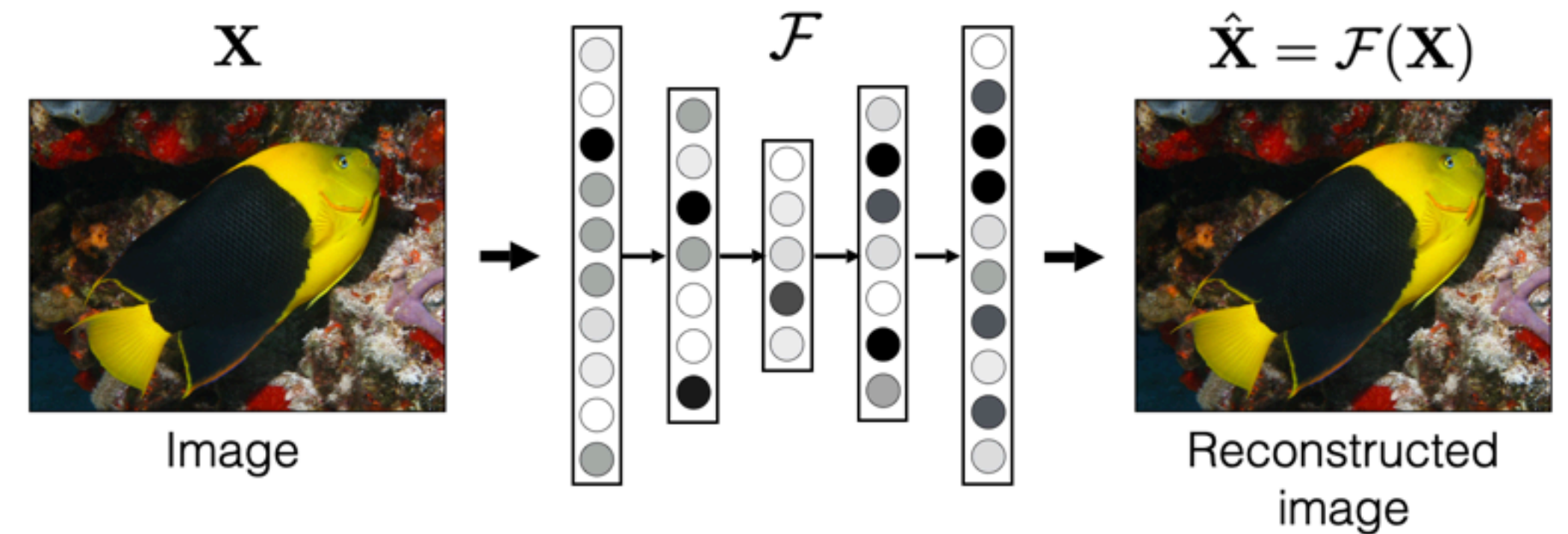
Input:  
Image

Output:  
Logits for class  
(1-10)

classes = ('plane', 'car', 'bird', 'cat',  
'deer', 'dog', 'frog', 'horse', 'ship', 'truck')

44

## Autoencoder



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

From MIT 6.8300/6.8301: Advances in Computer Vision