



THE UNIVERSITY OF BRITISH COLUMBIA

# Topics in AI (CPSC 532S): Multimodal Learning with Vision, Language and Sound

**Lecture 3: Introduction to Computer Vision**

# Computer vs. human vision



objects, scenes, people

**Human** Vision

\*slide from V. Ordonex



# Computer vs. human vision



objects, scenes, people

**Human** Vision

0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0

matrix of numbers

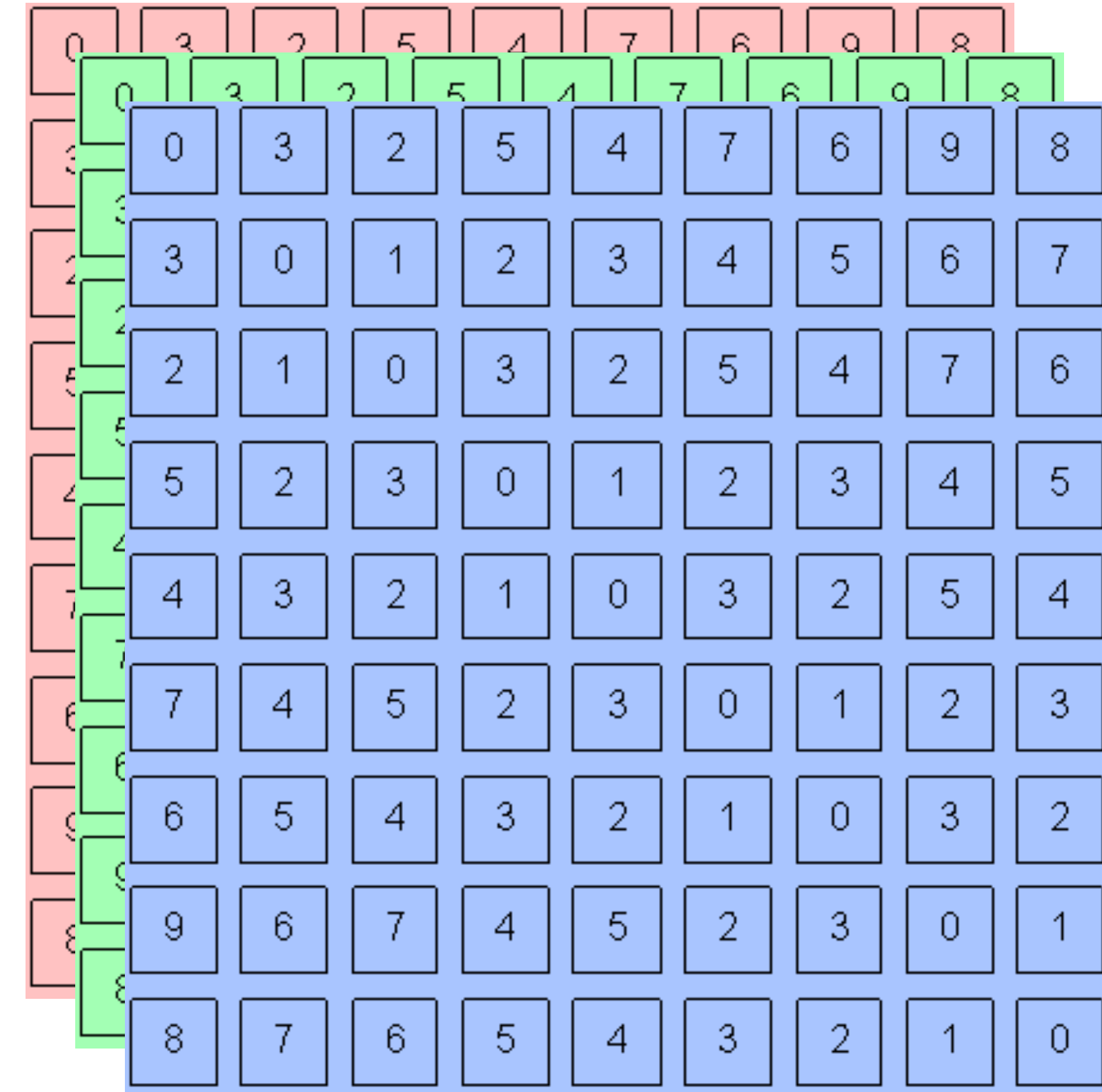
**Computer** Vision

# Computer vs. human vision



objects, scenes, people

**Human** Vision



tensor of numbers

**Computer** Vision



# Computer Vision

Computer vision studies the **tools and theories** that enable the design of machines that can **extract useful information from imagery data** (images and videos) toward the goal of **interpreting the world**

\*curtesy of Peter Meer





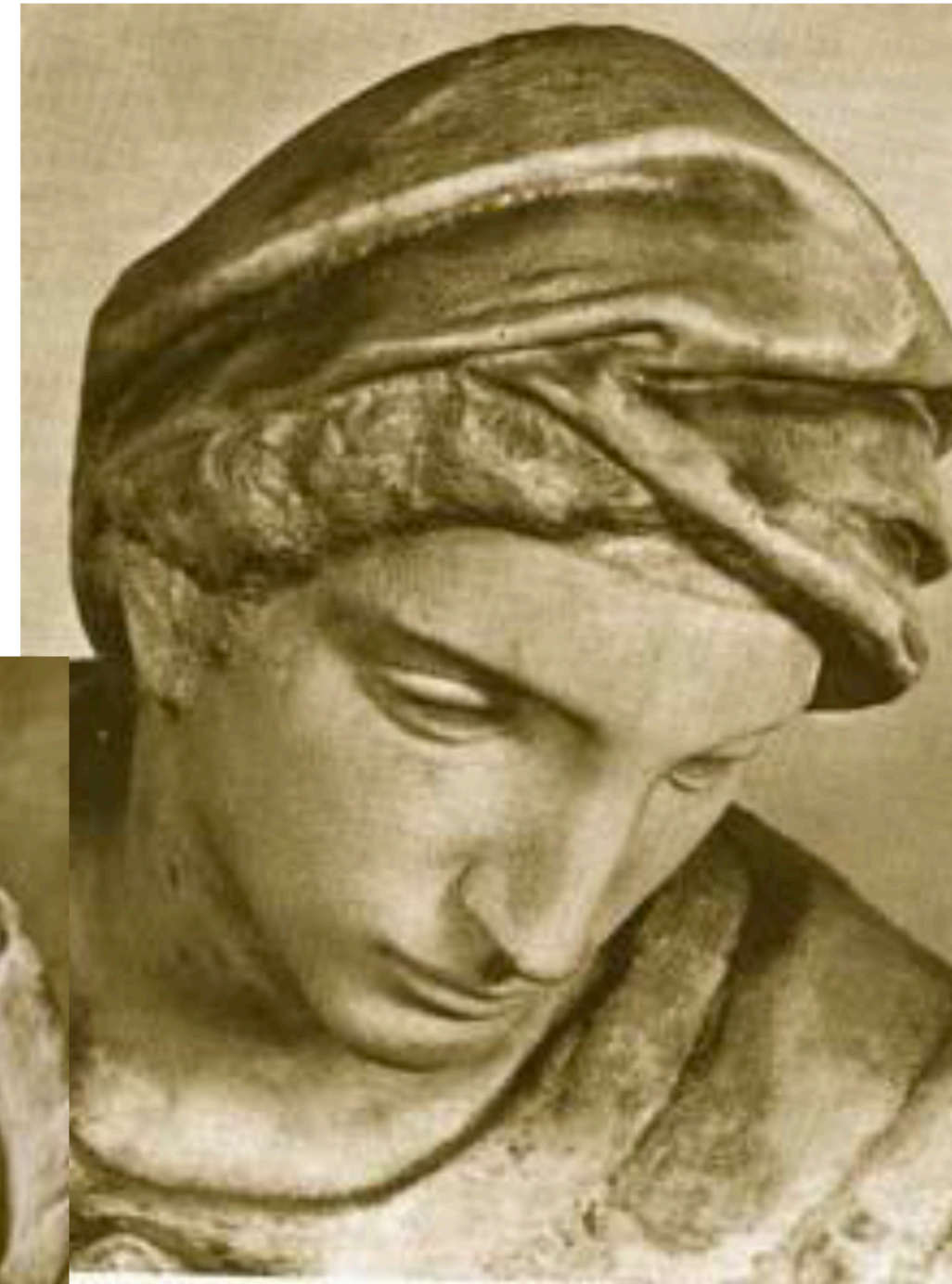
# **Vision** is Amazing Feat of **Natural Intelligence**

~ 55% of **cerebral cortex** in humans (13 billion neurons) are devoted to vision  
more human brain devoted to vision than anything else





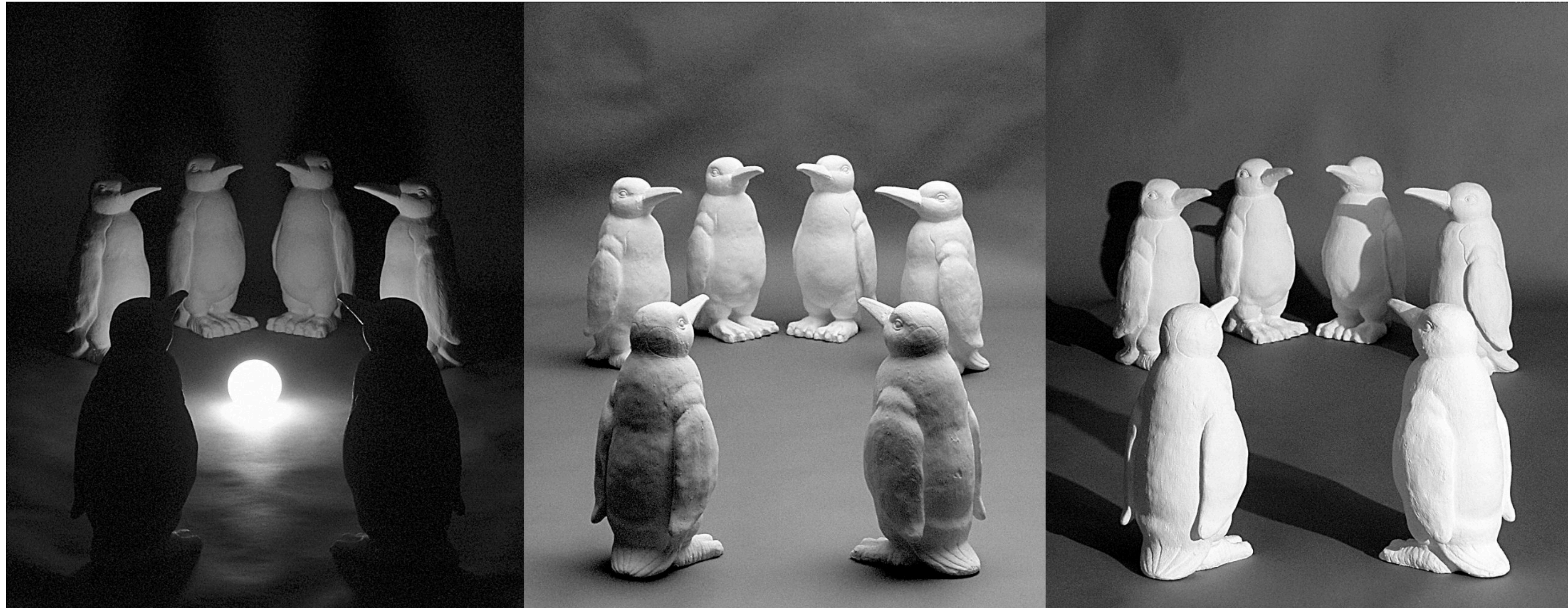
# Challenges: Viewpoint invariance



**Michelangelo** 1475-1564



# Challenges: Lighting



\*image credit J. Koenderink



# Challenges: Scale



\*slide credit Fei-Fei, Fergus & Torralba

# Challenges: Deformation



\*image credit Peter Meer



# Challenges: Occlusions

**Rene Magritte 1965**





# Challenges: Background clutter

**Kilmeny Niland** 1995





# Challenges: Local ambiguity and context



\*image credit Fergus & Torralba

# Challenges: Local ambiguity and context



\*image credit Fergus & Torralba



# Challenges: Motion



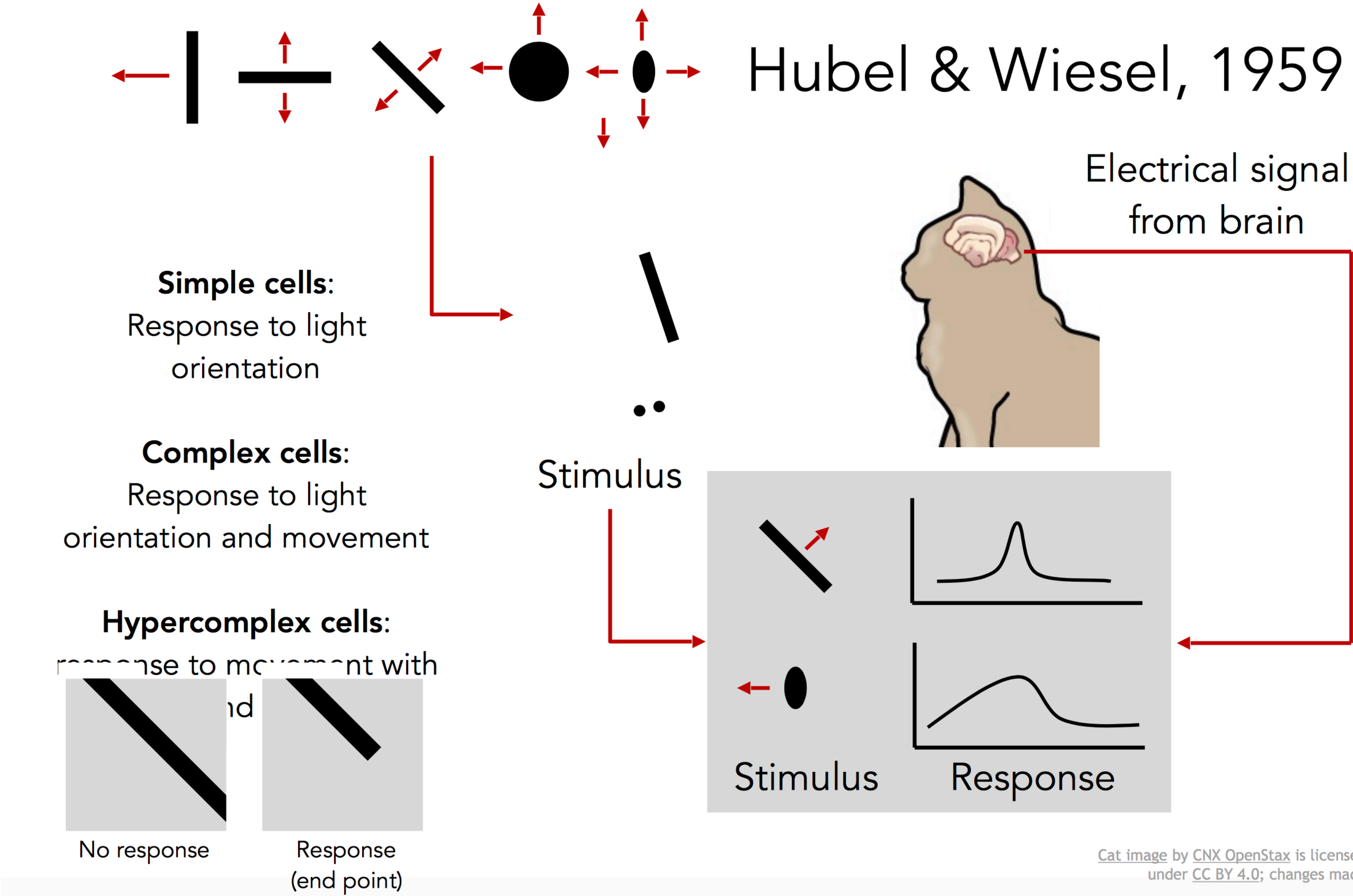
\*image credit Peter Meer



# Challenges: Object inter-class variation



# Human vision ...



\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**



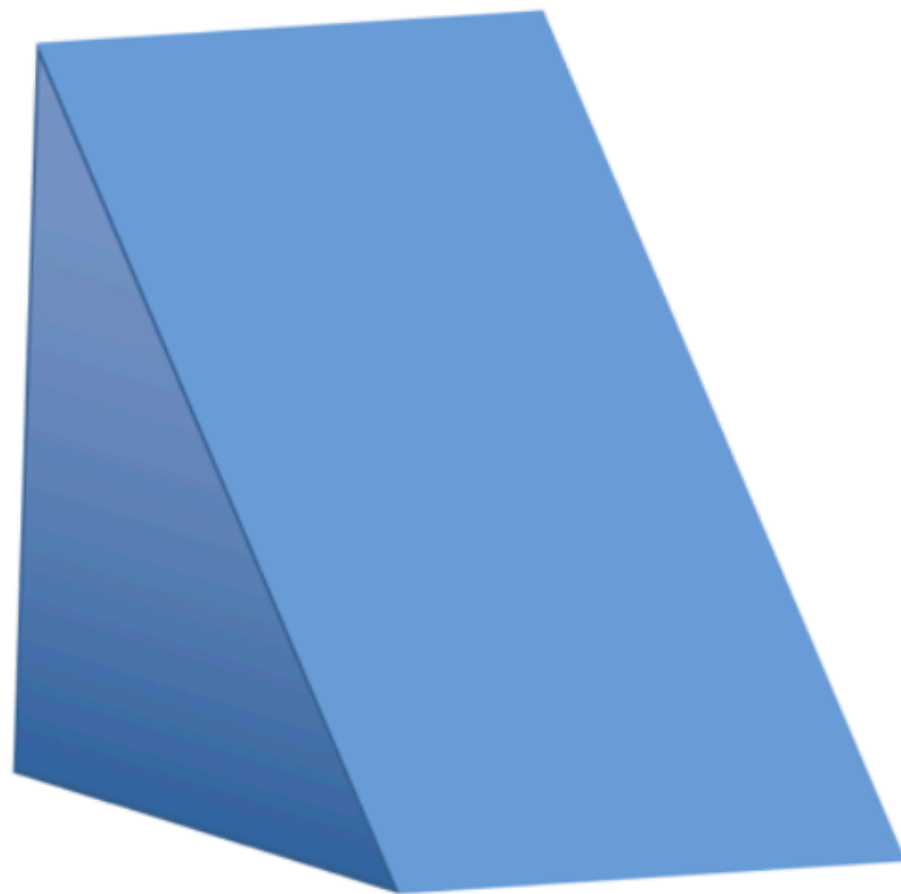
# Computer vision ... the beginning ...



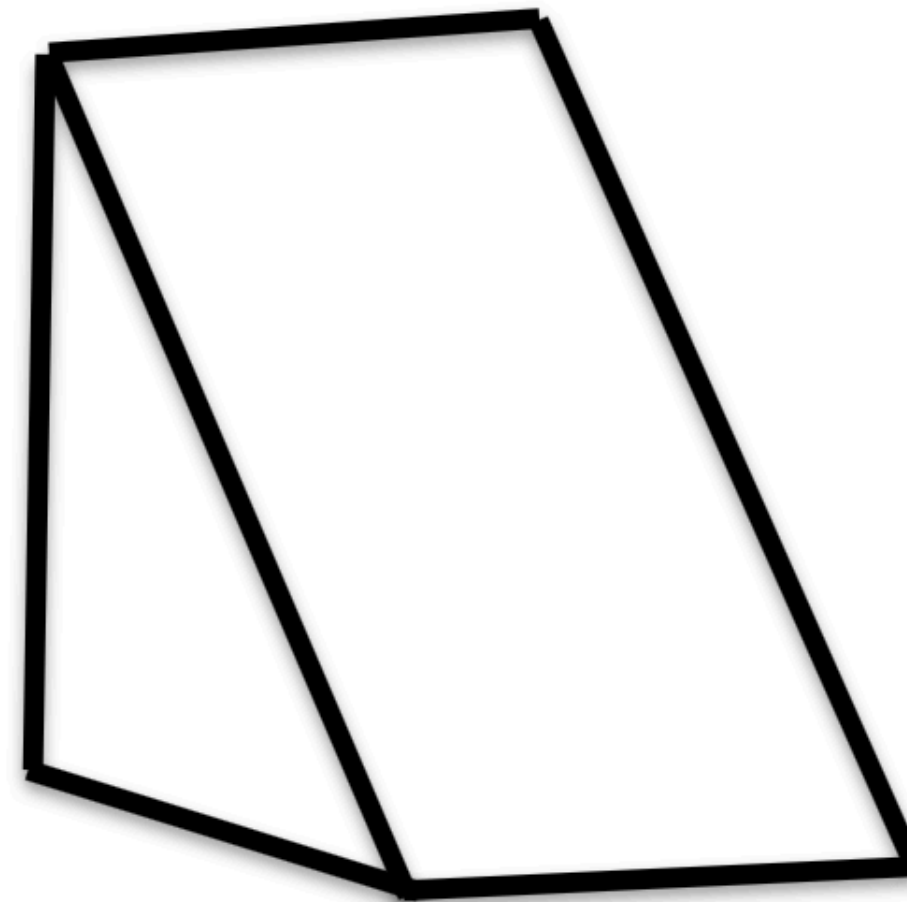
Larry Roberts

**Blocks World.** first thesis in  
computer vision, 1963

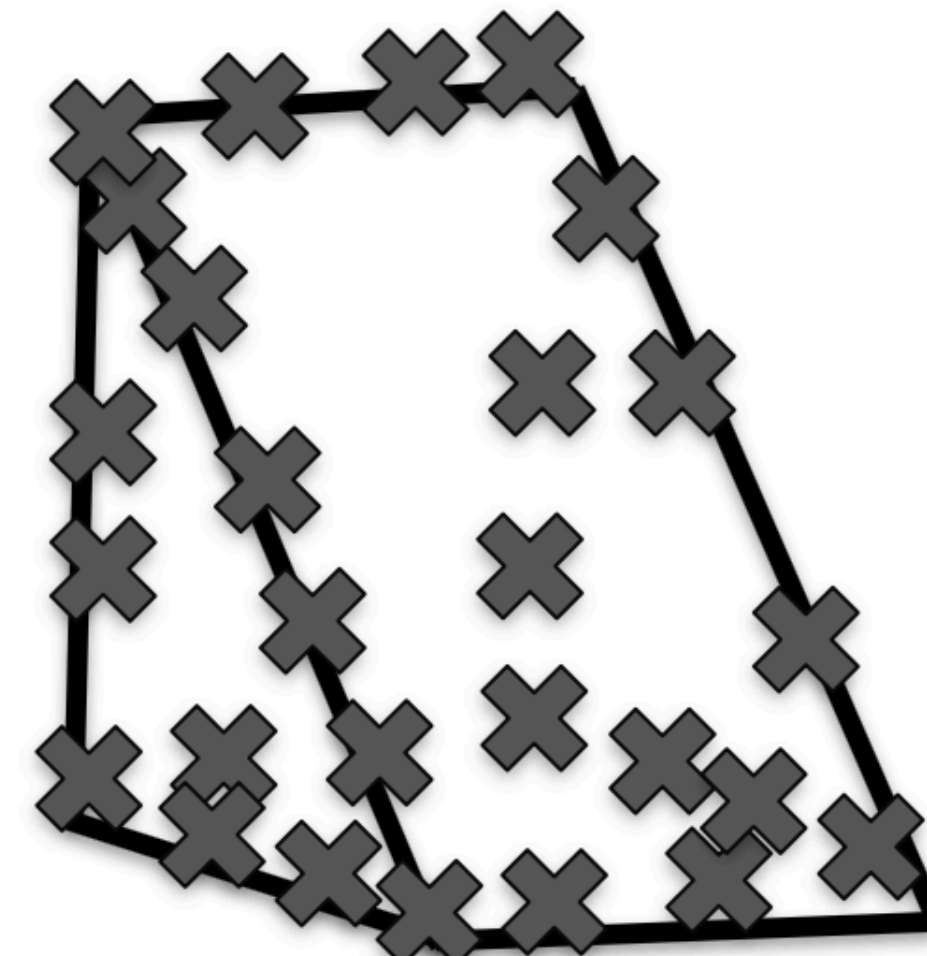
“the perception of **solid objects** is a process which can be based on  
the **properties of three-dimensional** transformations and **the  
laws of nature**”



(a) Original picture



(b) Differentiated picture



(c) Feature points selected



# Computer vision ... the beginning ...



Larry Roberts

**Blocks World.** first thesis in computer vision, 1963

"the perception of **solid objects** is a process which can be based on the **properties of three-dimensional** transformations and **the laws of nature**"

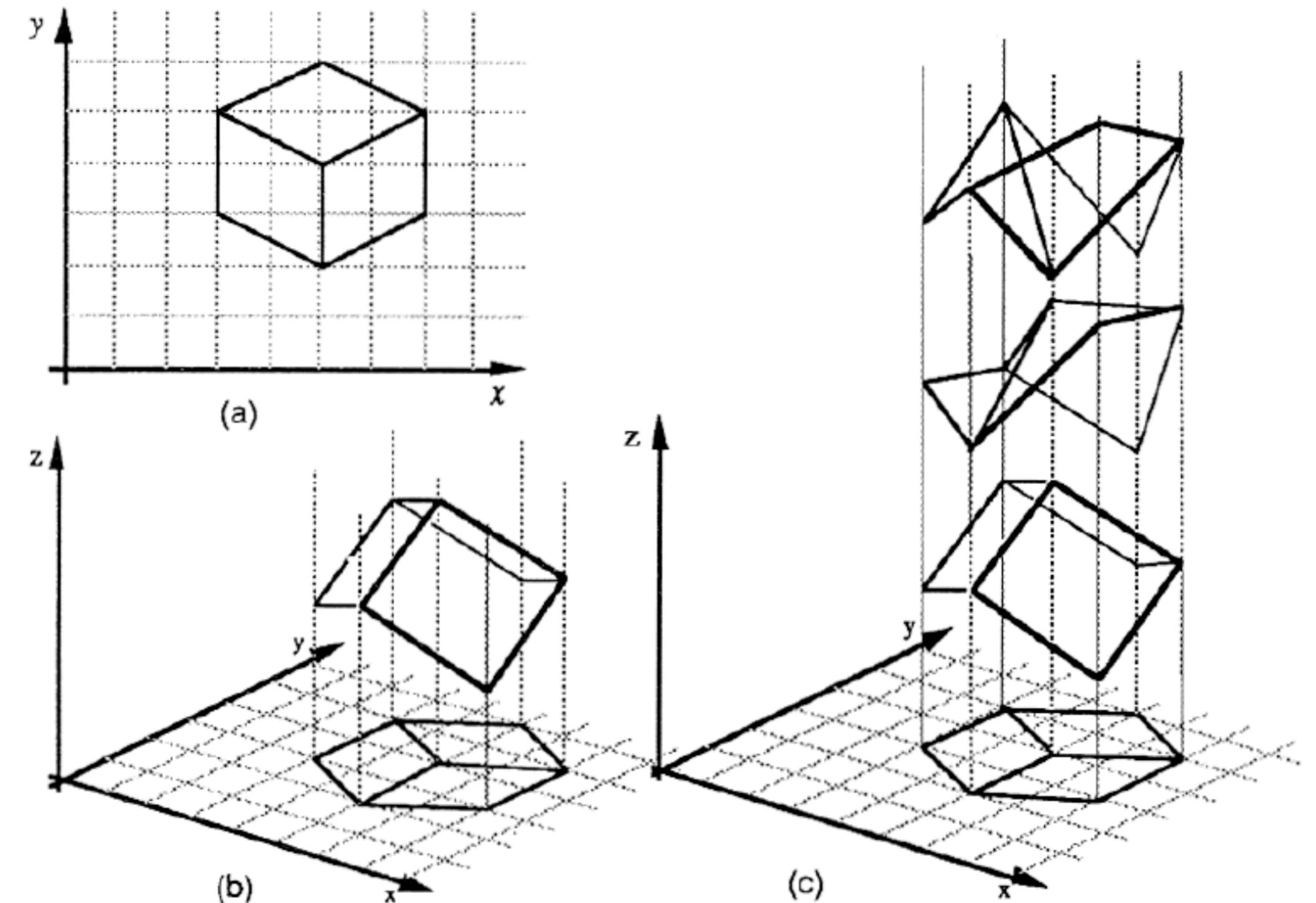


Figure 1. (a) A line drawing provides information only about the  $x, y$  coordinates of points lying along the object contours. (b) The human visual system is usually able to reconstruct an object in three dimensions given only a single 2D projection (c) Any planar line-drawing is geometrically consistent with infinitely many 3D structures.



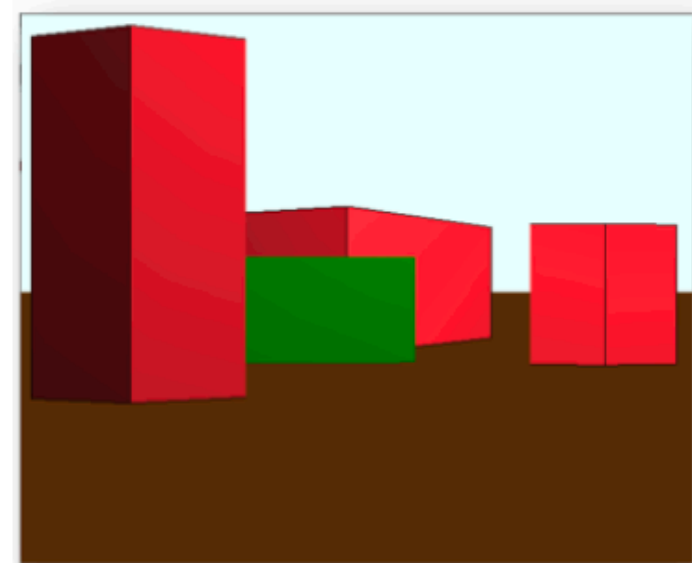
# Computer vision ... the beginning ...



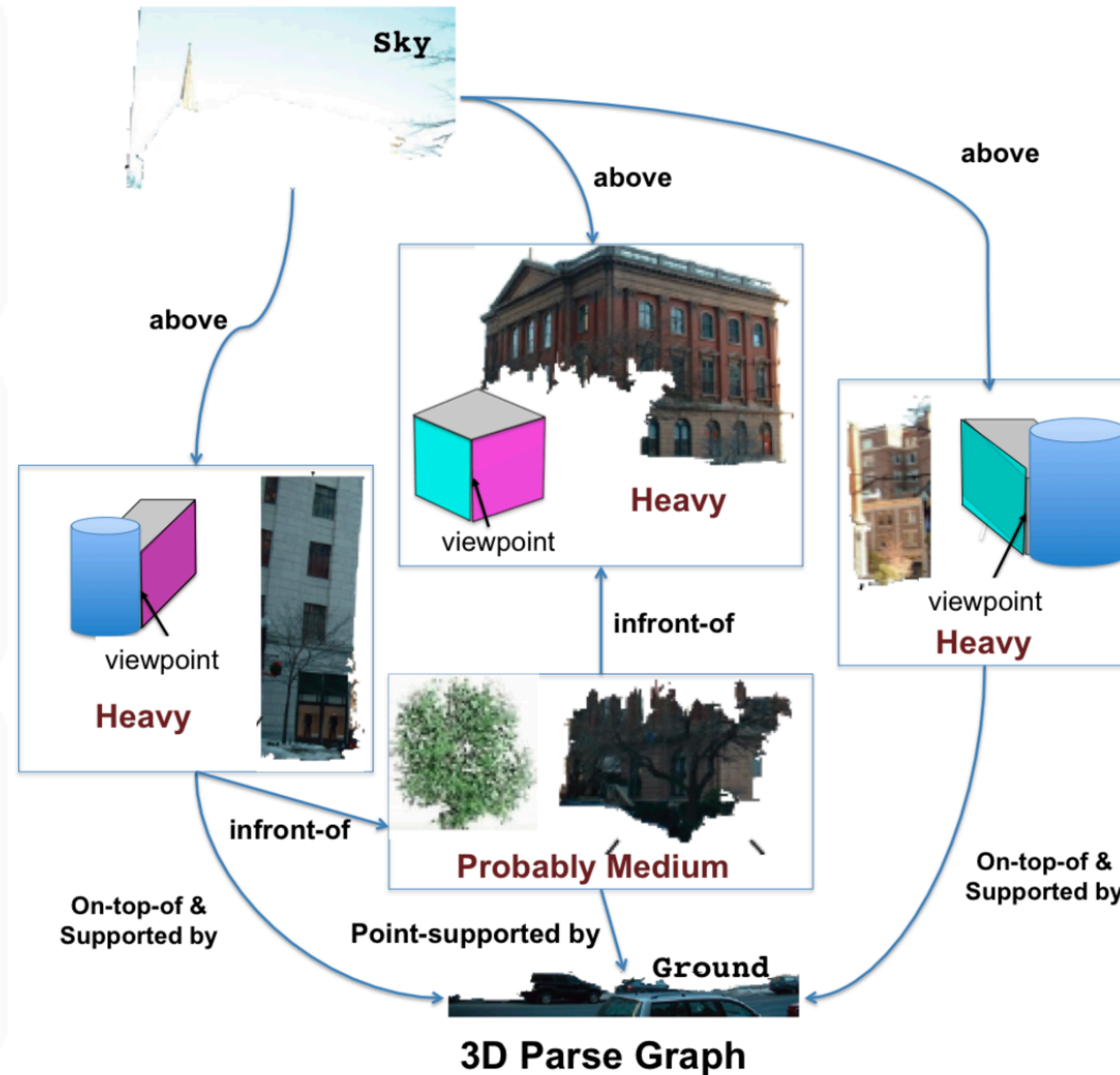
Input Image



Blocks World



3D Rendering



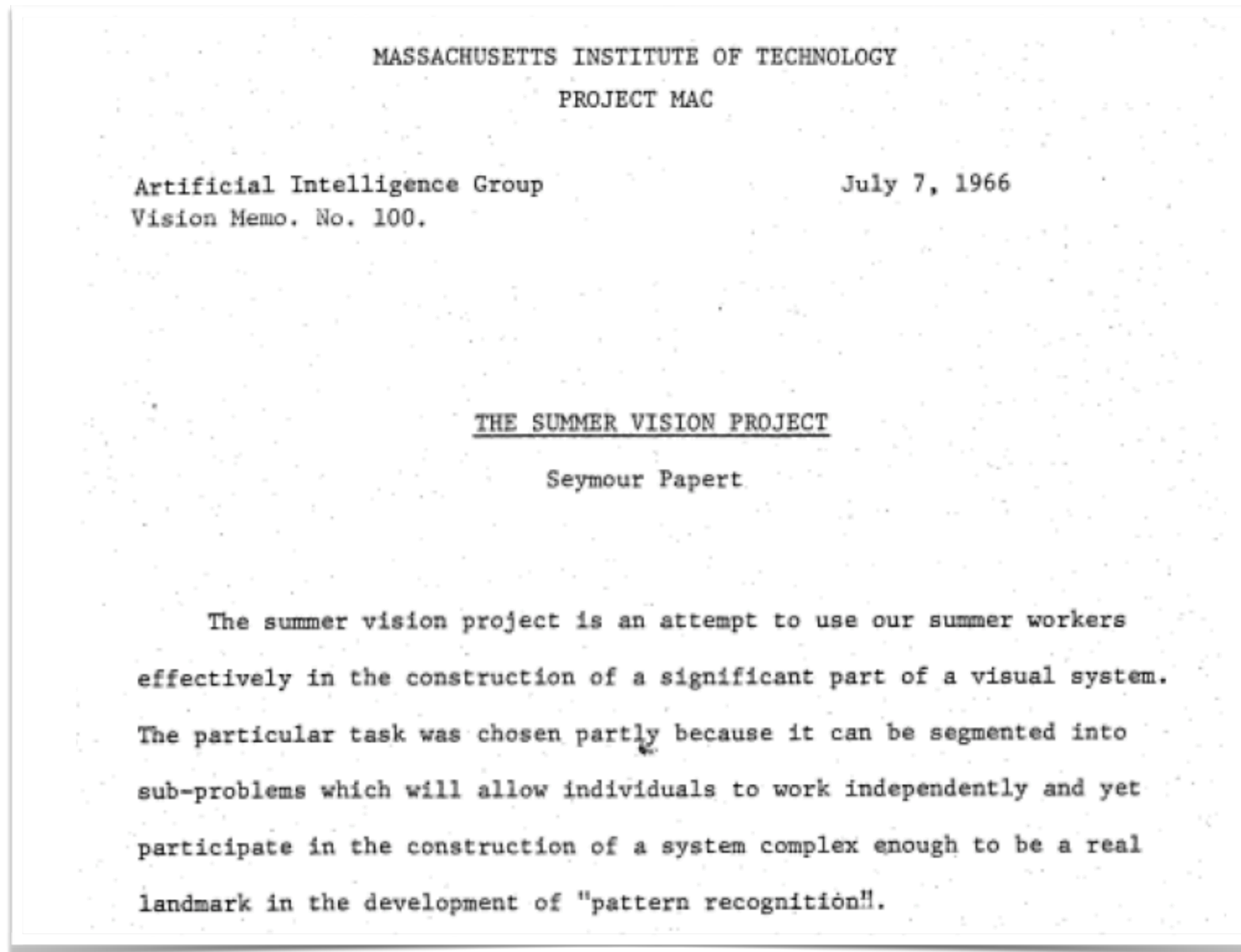
**Static Equilibrium:** Forces and torques acting on a block should cancel each other out.

**Support Force Constraint:** Supporting object should have enough strength to provide contact reactionary forces

**Volumetric Constraints:** All objects in the world must have finite volume & cannot penetrate each other



# Computer vision ... the beginning ...

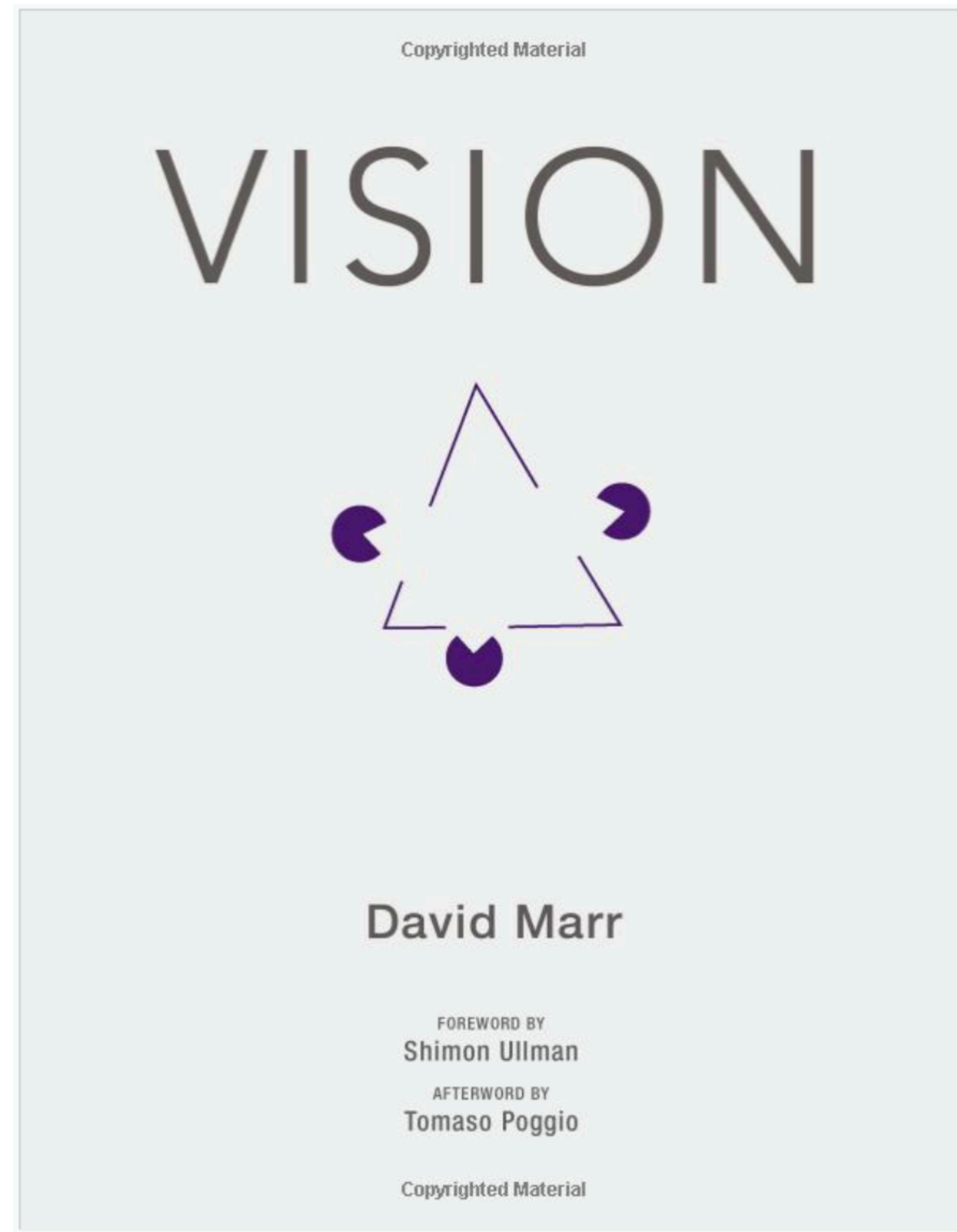


In 1966, Marvin Minsky at MIT asked his undergraduate student Gerald Jay Sussman to “spend the summer linking a camera to a computer and getting the computer to describe what it saw”

[ Szeliski 2009, Computer Vision ]

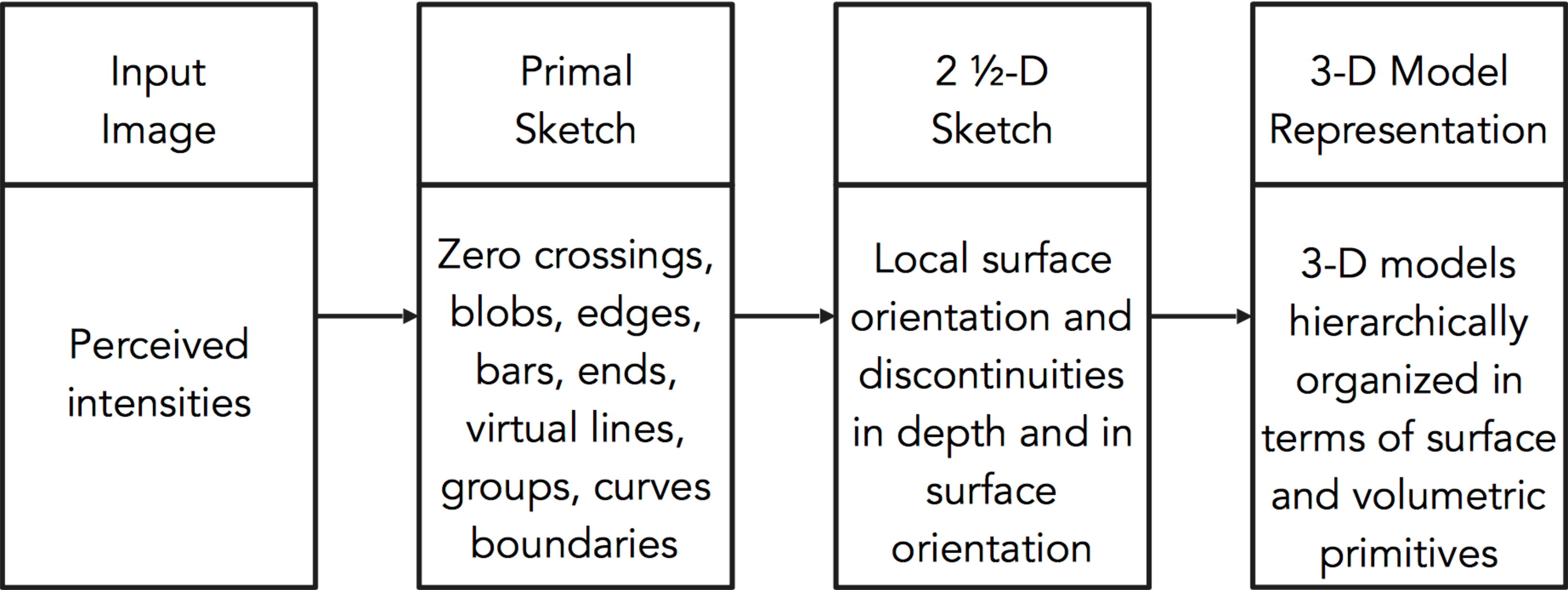
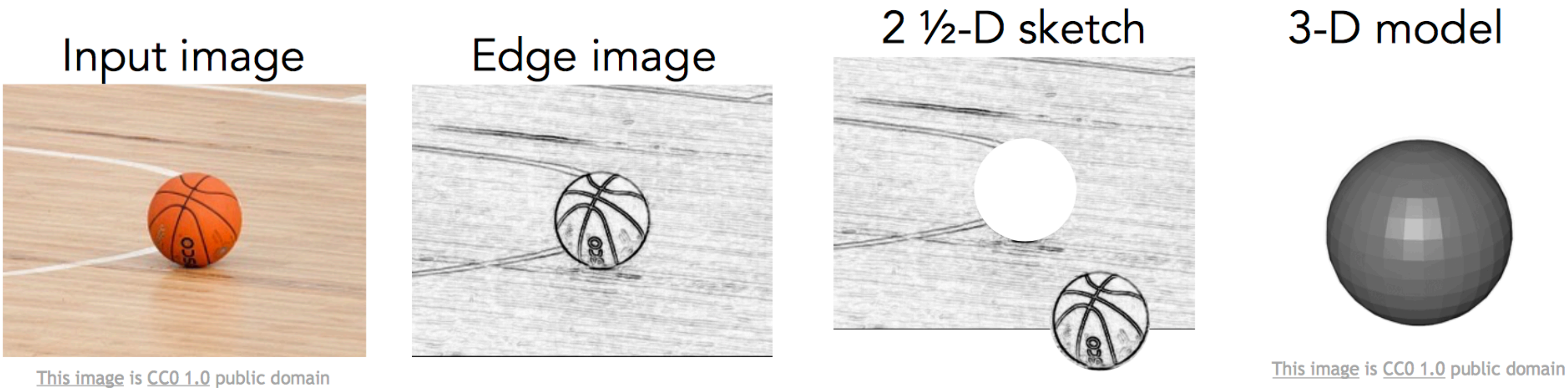


# David **Marr**, 1970s





# David Marr, 1970s



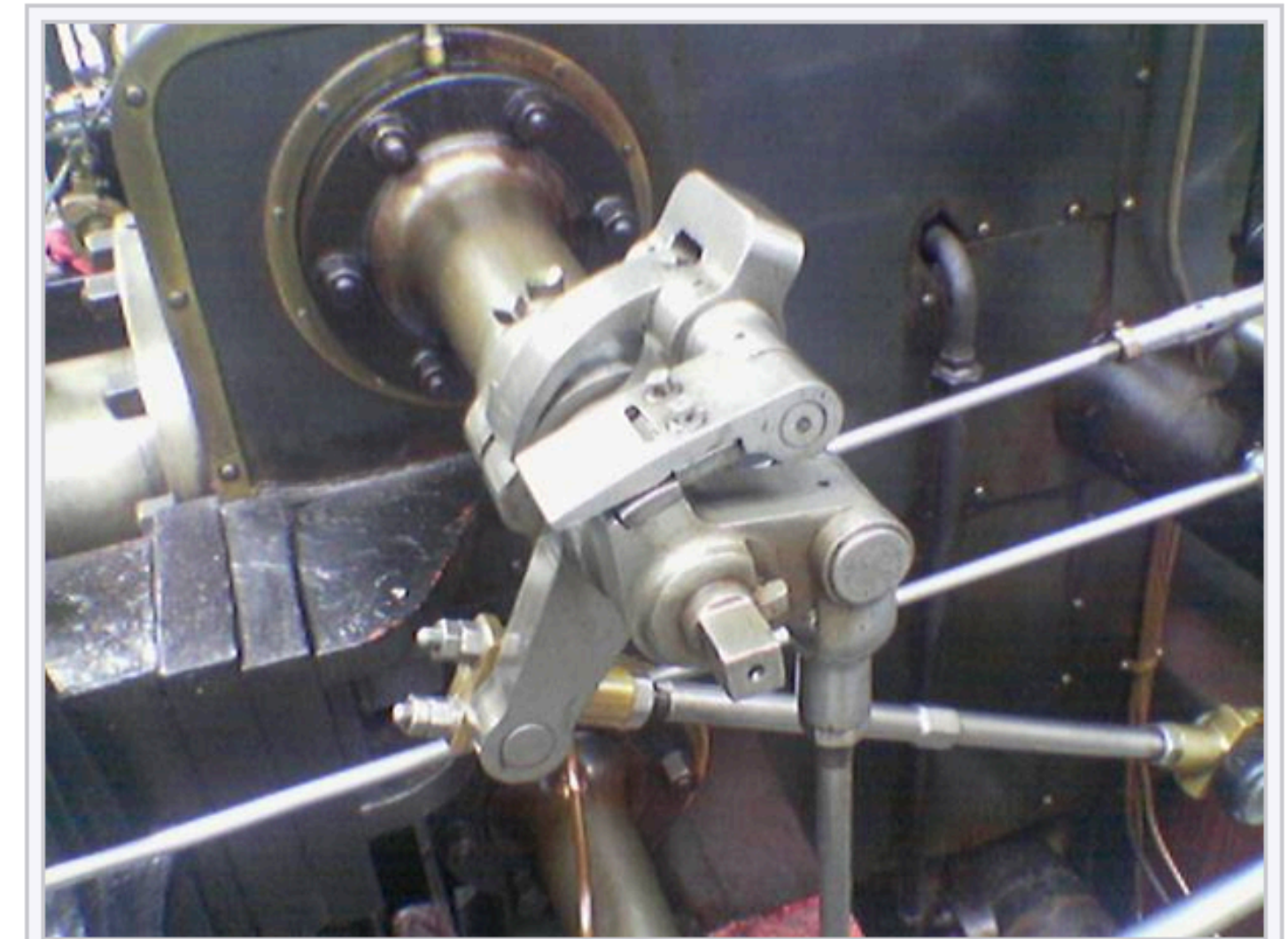
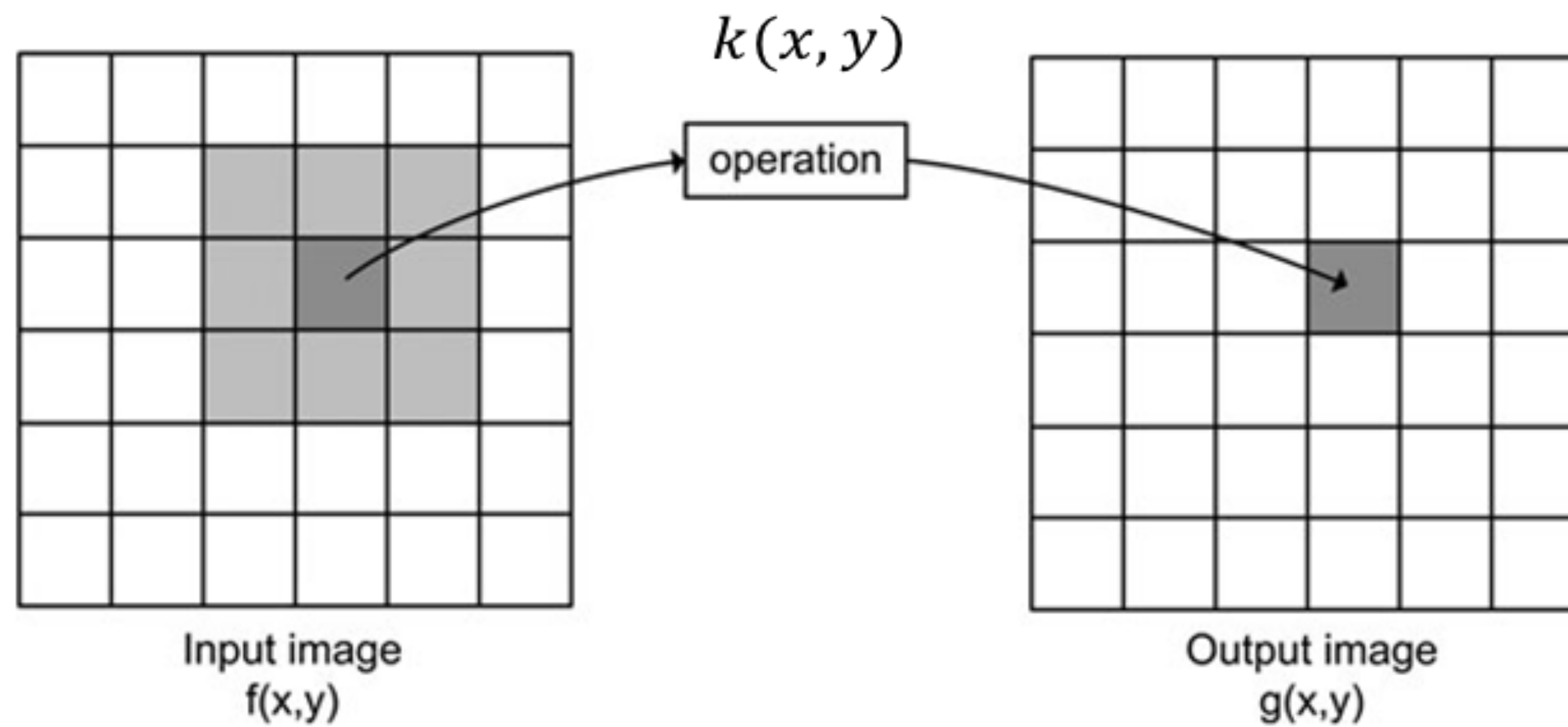
[ Stages of Visual Representation, **David Marr** ]

\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

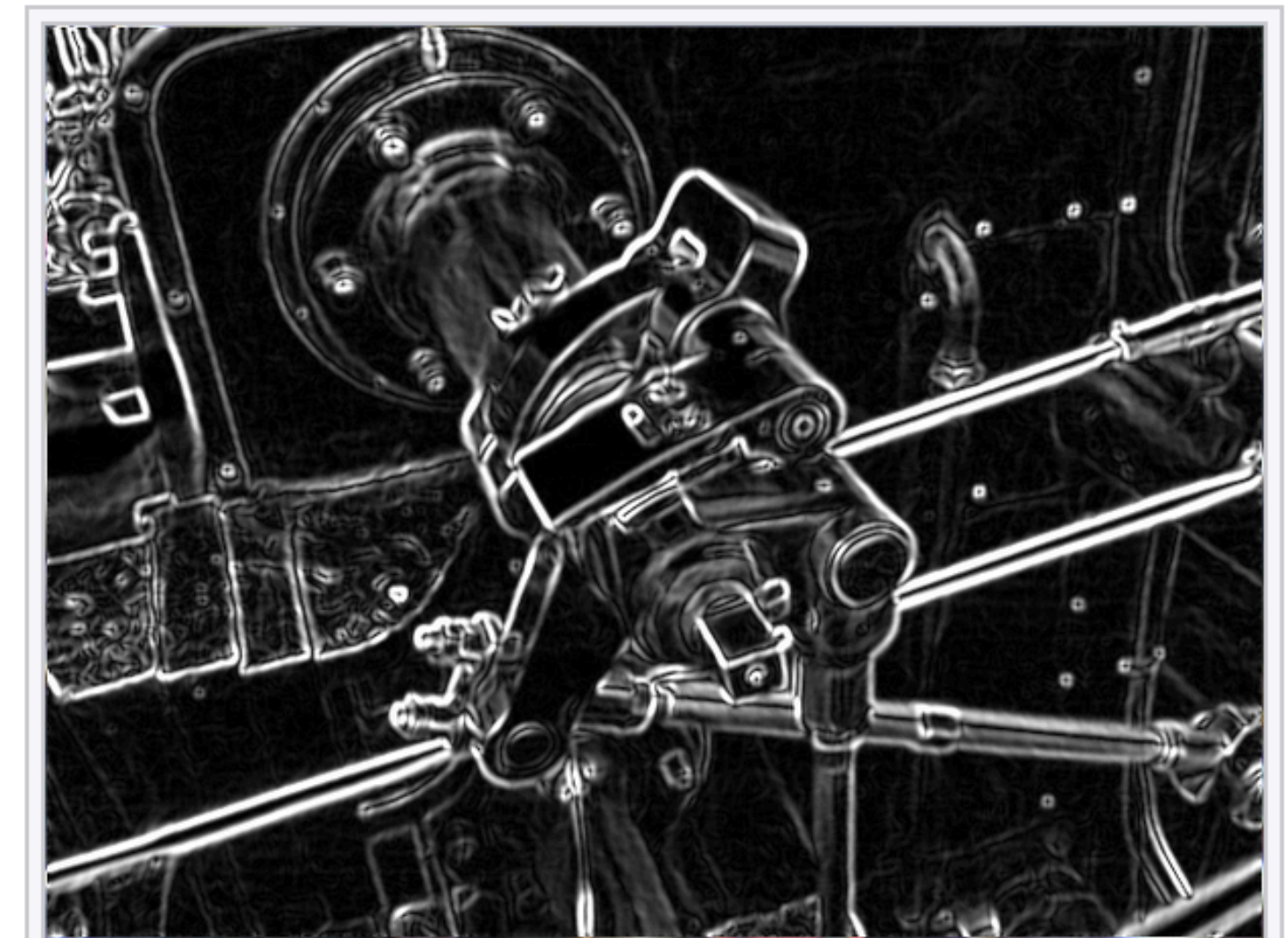


# Edges

1	0	-1
2	0	-2
1	0	-1



A color picture of a steam engine

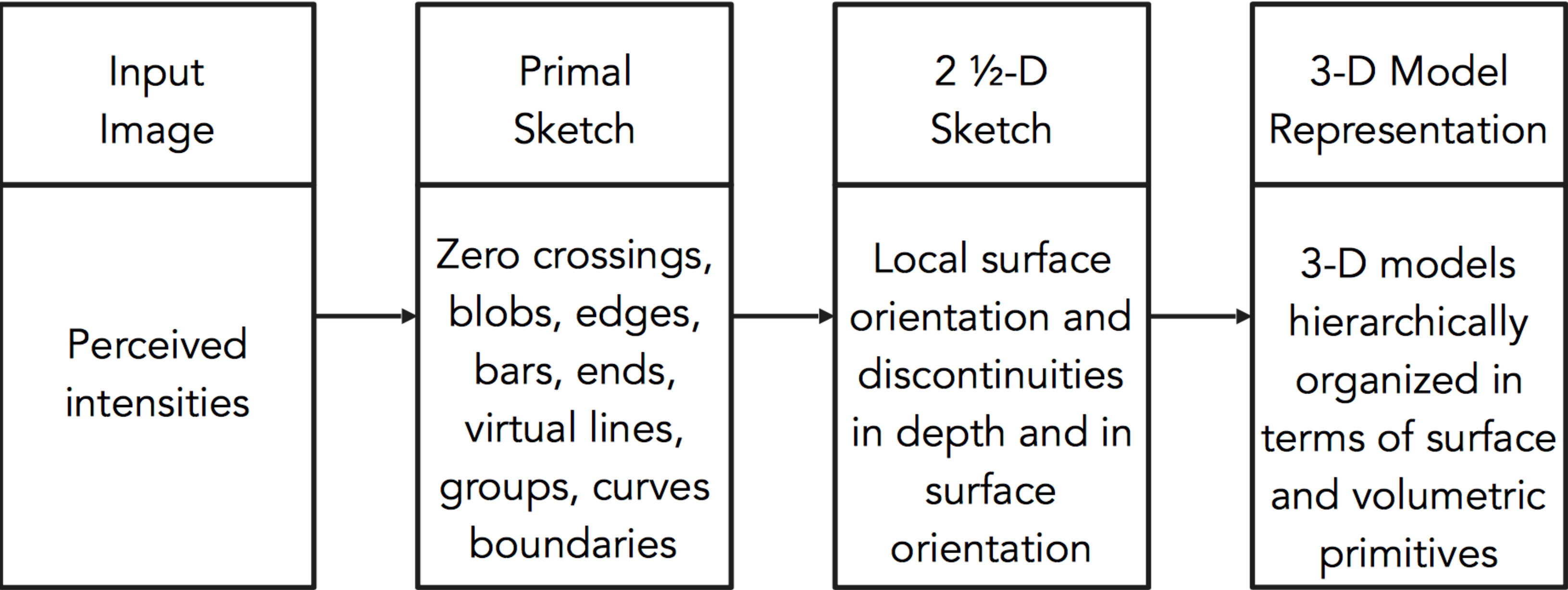
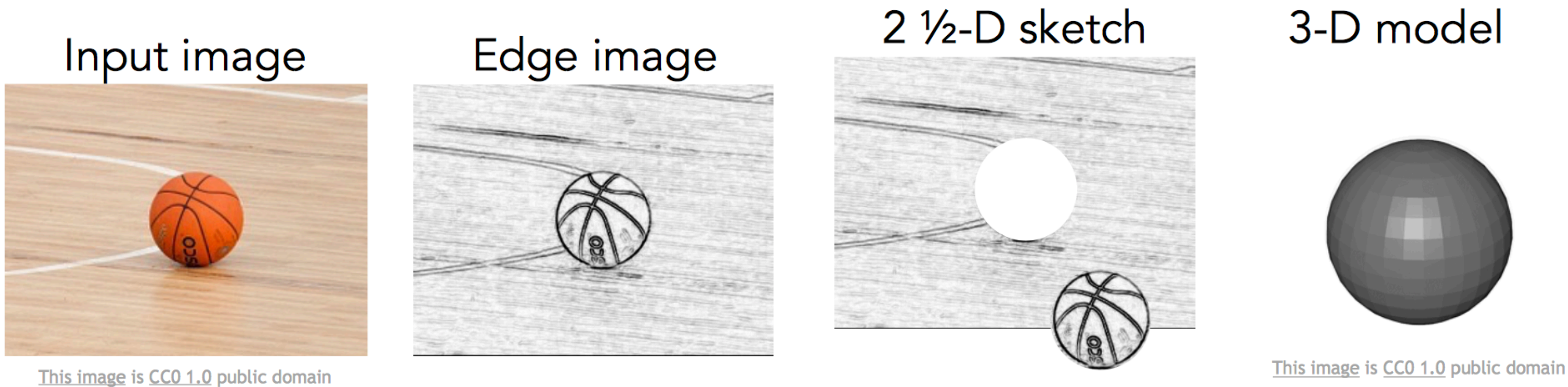


The Sobel operator applied to that image





# David Marr, 1970s



[ Stages of Visual Representation, **David Marr** ]

\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**



# Segmentation - GraphCuts

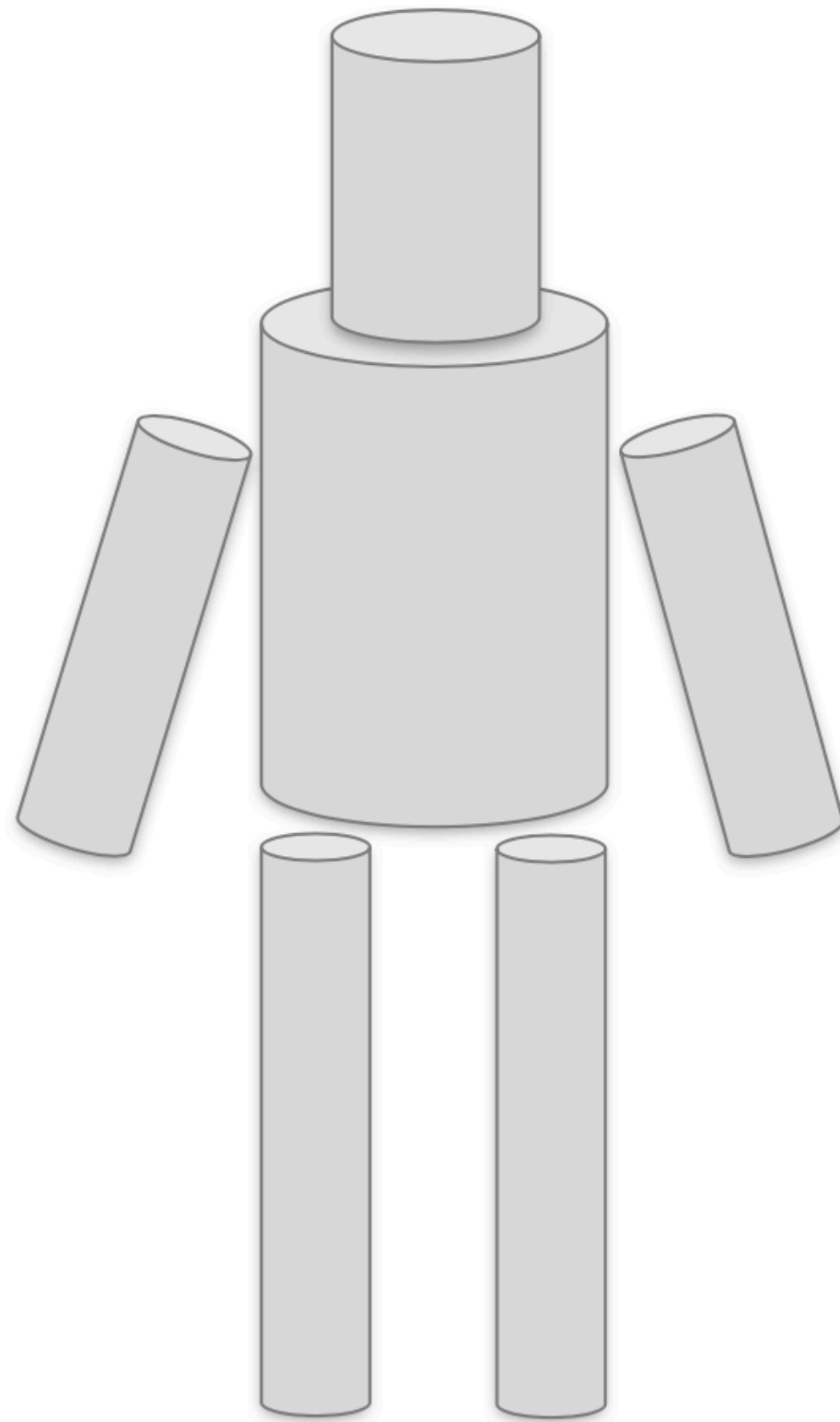


[ Shi & Malik, 2000 ]



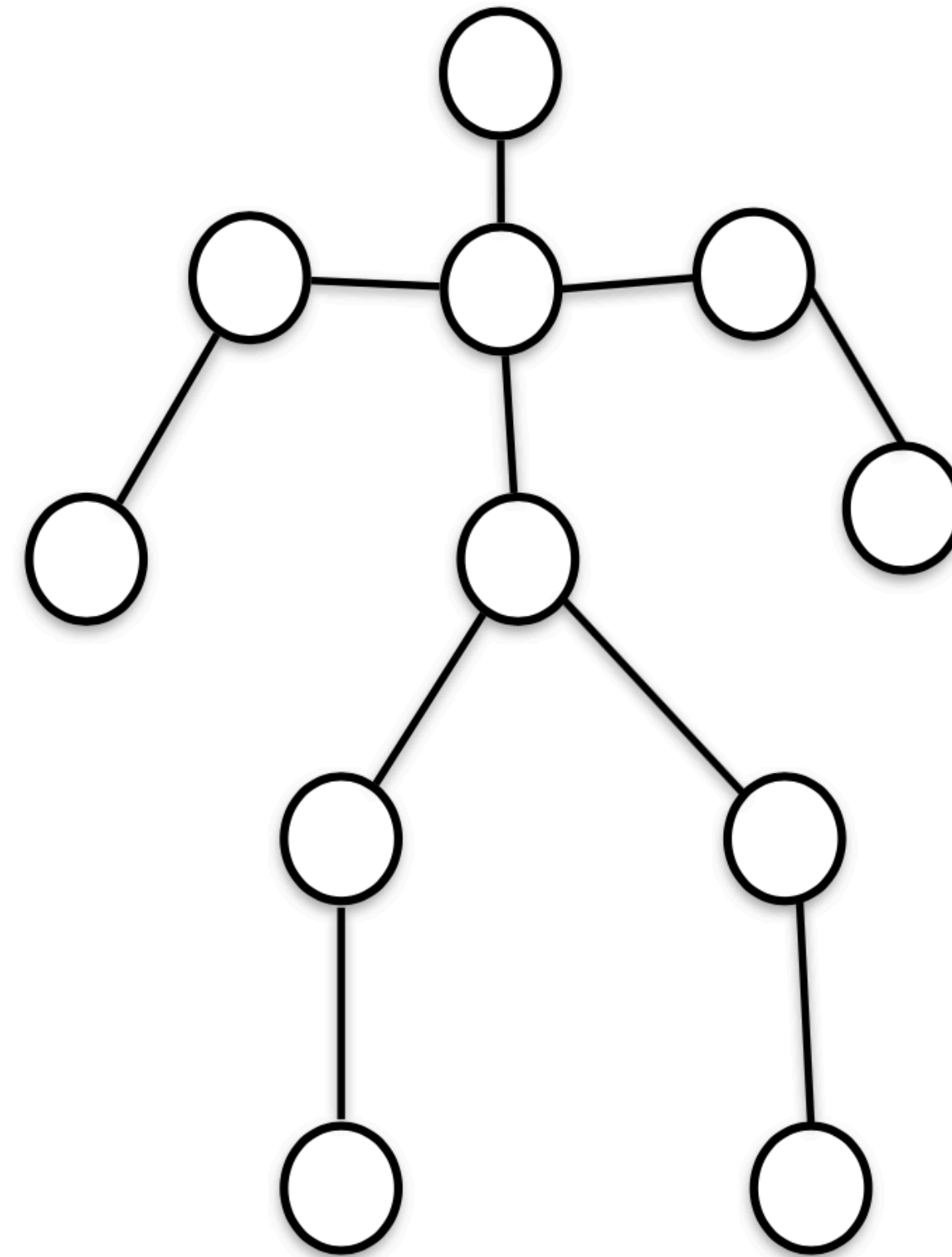
# Part-based Models

## Generalized Cylinders

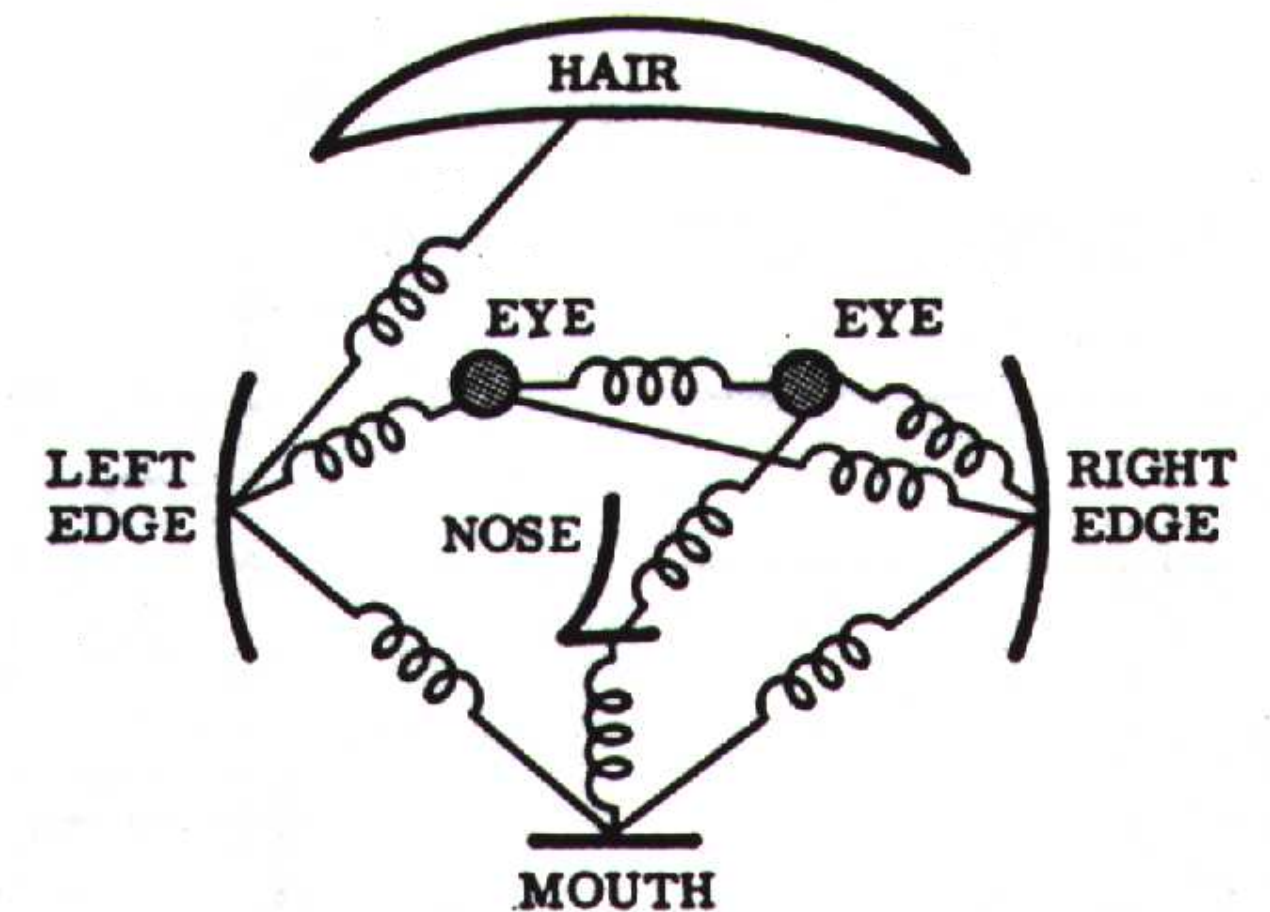


[ Brooks & Binford, 1979 ]

## Pictorial Structures

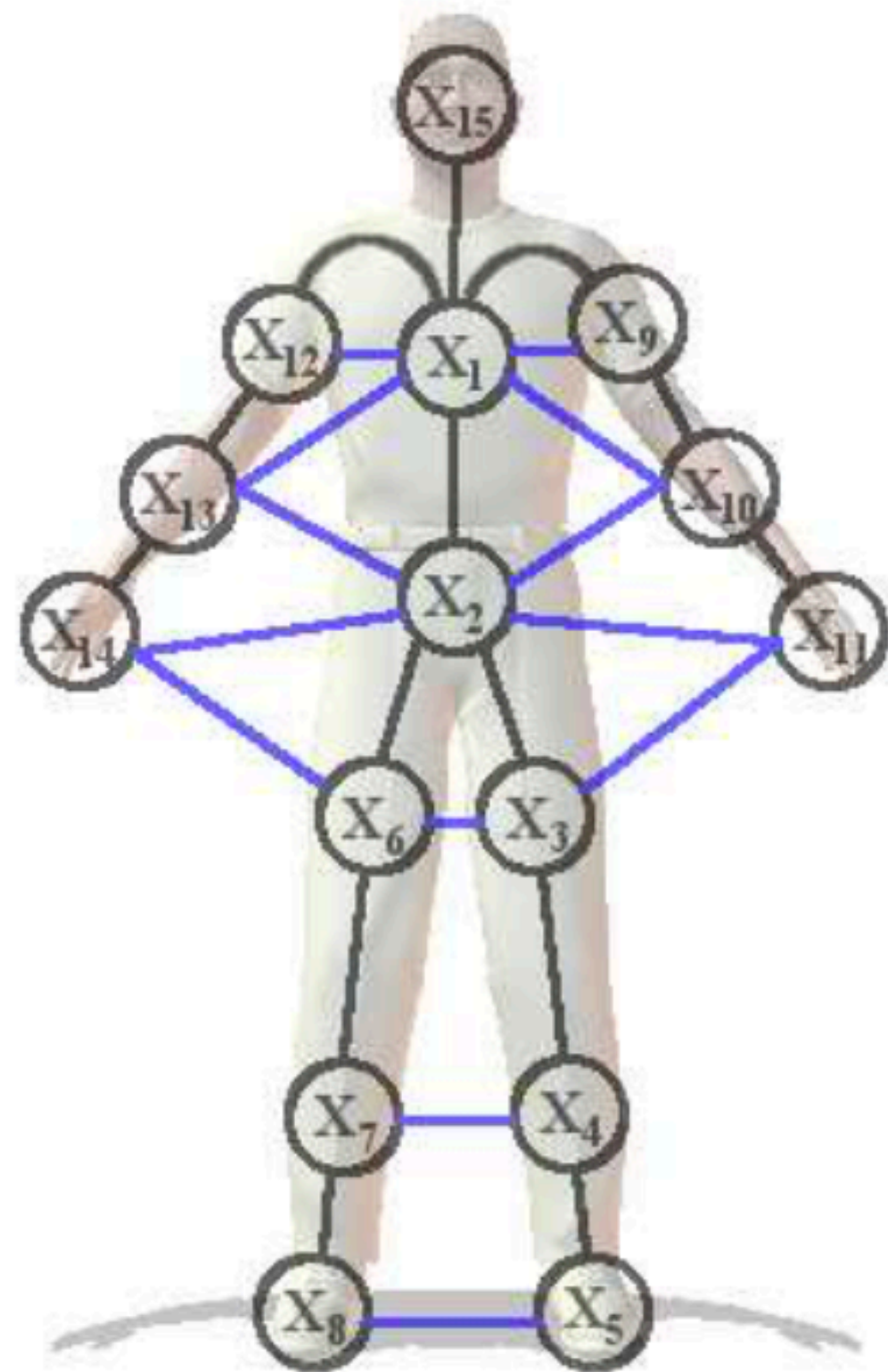


[ Fischler & Elschlager, 1973 ]

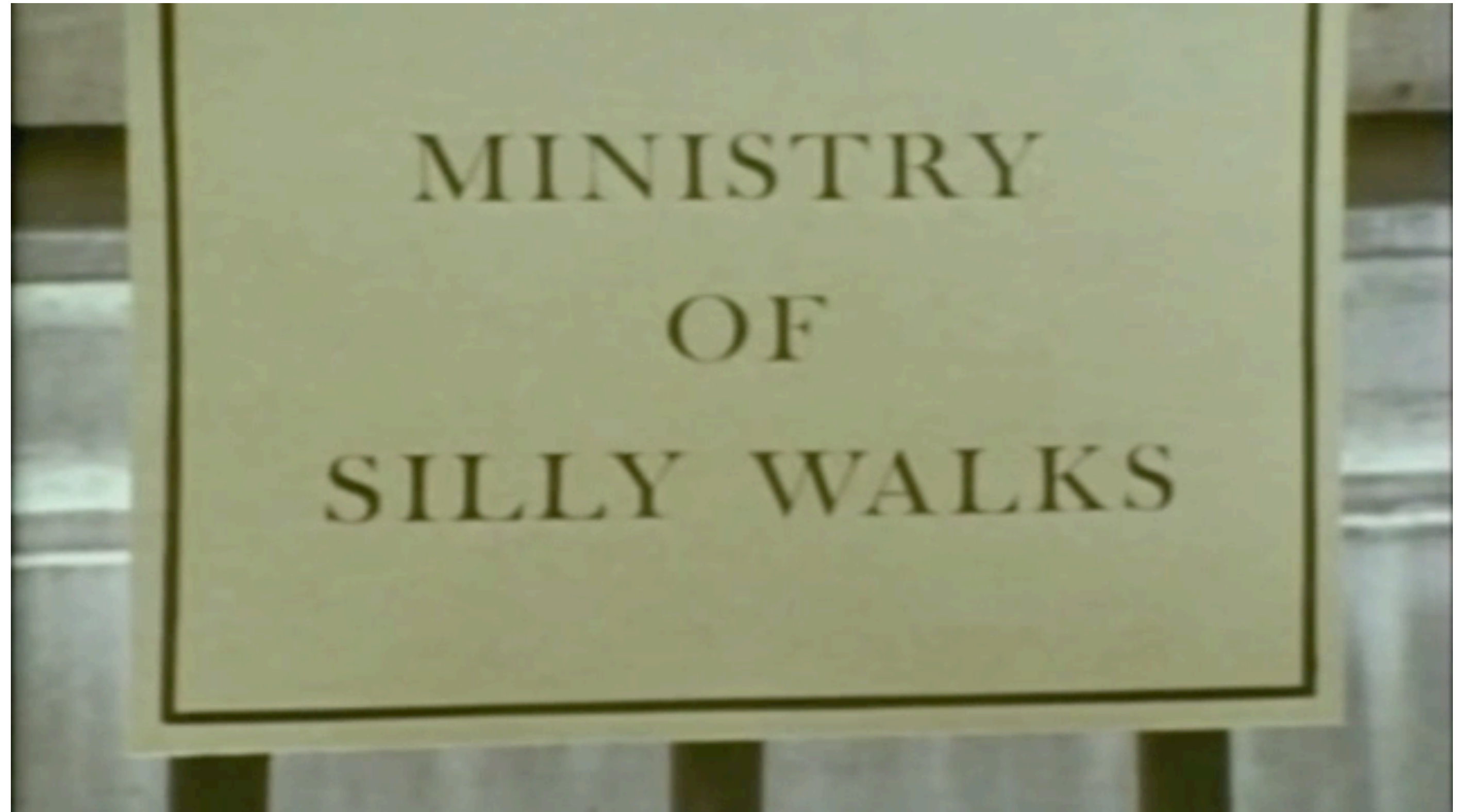




# Part-based Models



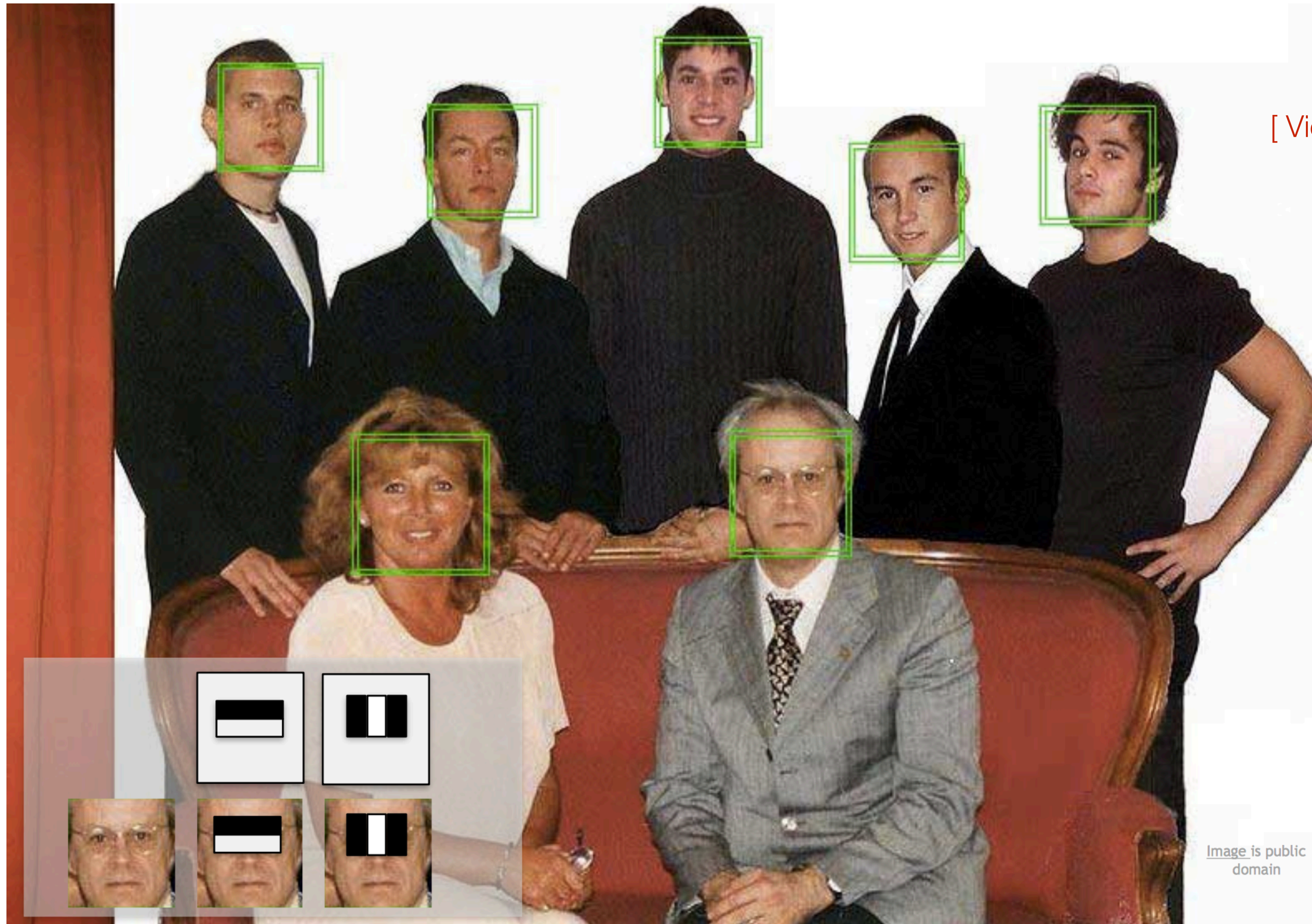
[ Sigal et al. 2004]



Monty Python's **Ministry of Silly Walks**



# Face Detection 1999-2000



[ Viola & Jones, 2001 ]

Image is public domain



# Feature-based Vision



Image is public domain



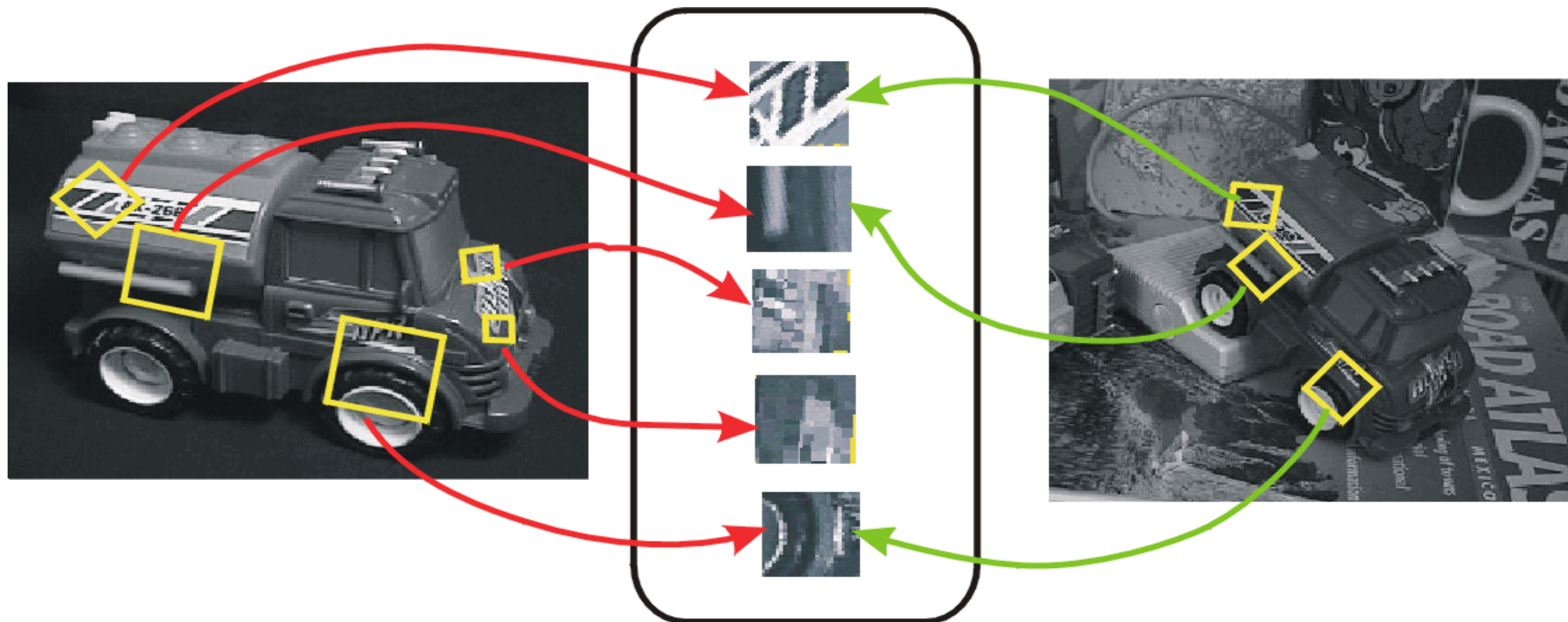
Image is CC BY-SA 2.0

[ **David Lowe**, 1999 ]



# SIFT Idea

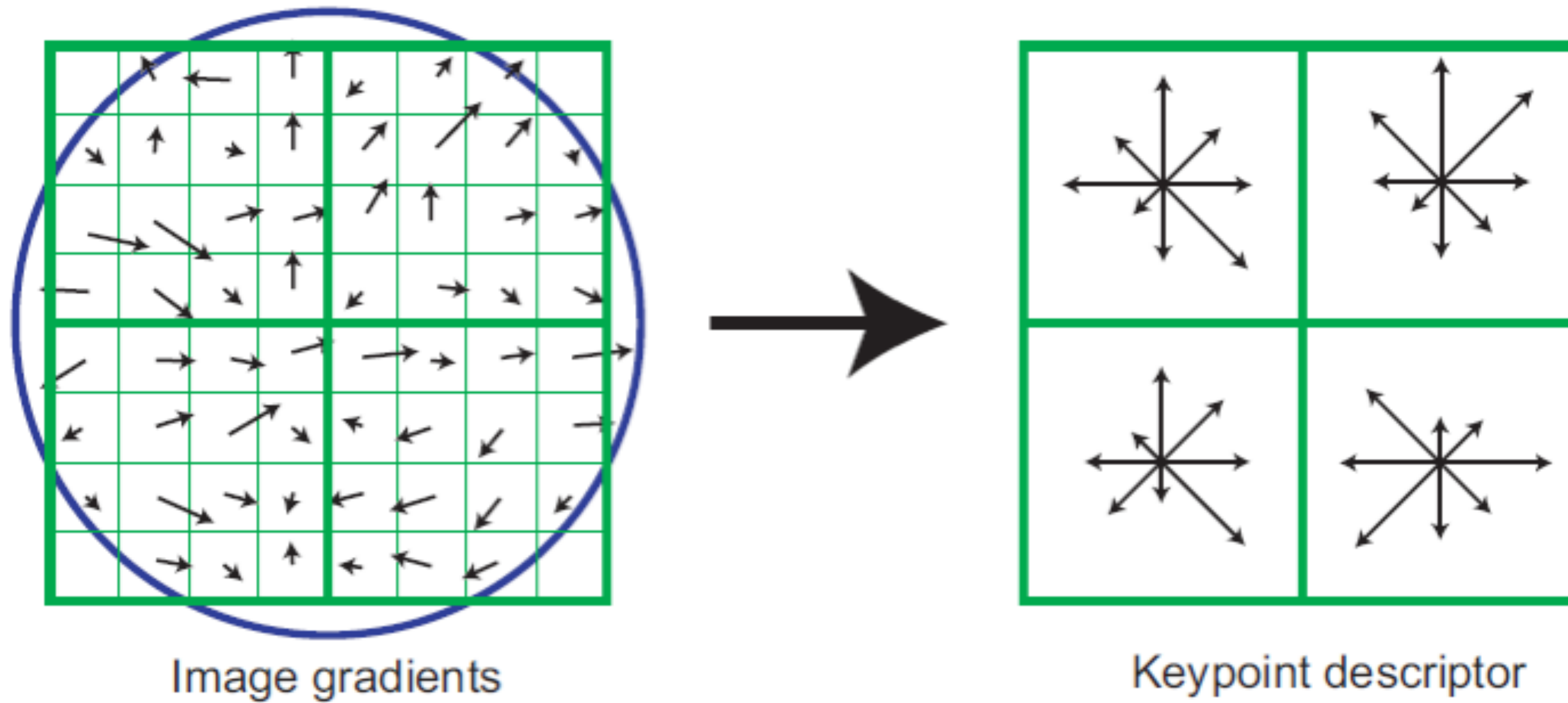
Image content is transformed into local feature coordinates that are **invariant** to translation, rotation, scale and imaging parameters



[ **David Lowe**, 1999 ]



# SIFT Descriptor



[ David Lowe, 1999 ]



# Massive 3D Reconstructions



[ Agarwal, Furukawa, Snavely, Curless, Seitz, Szeliski, 2010 ]



# Bag-of-Words

\*slide credit Li Fei-Fei

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that come from our eyes. For a long time, the visual image was considered as a movie scene. It was discovered that behind the image in the brain is a very complicated system. Hubel and Wiesel have been able to demonstrate that the message about the image falling on the retina undergoes a step-wise analysis in a series of nerve cells stored in columns. In this system, each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.

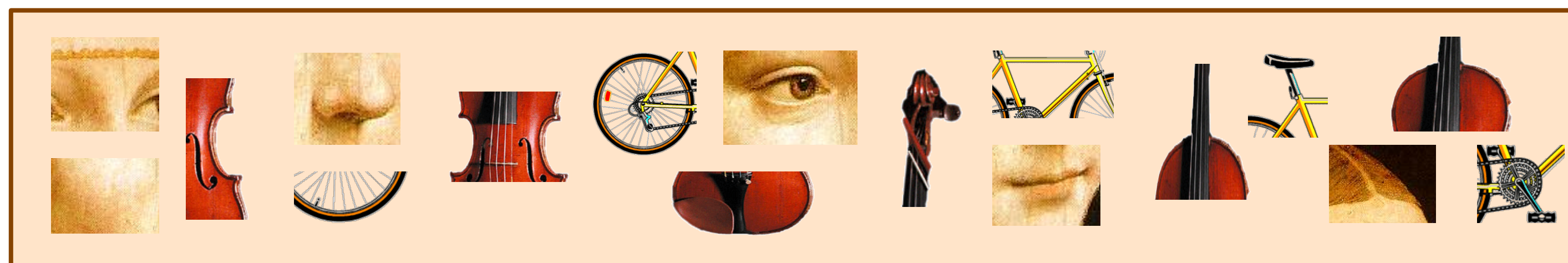
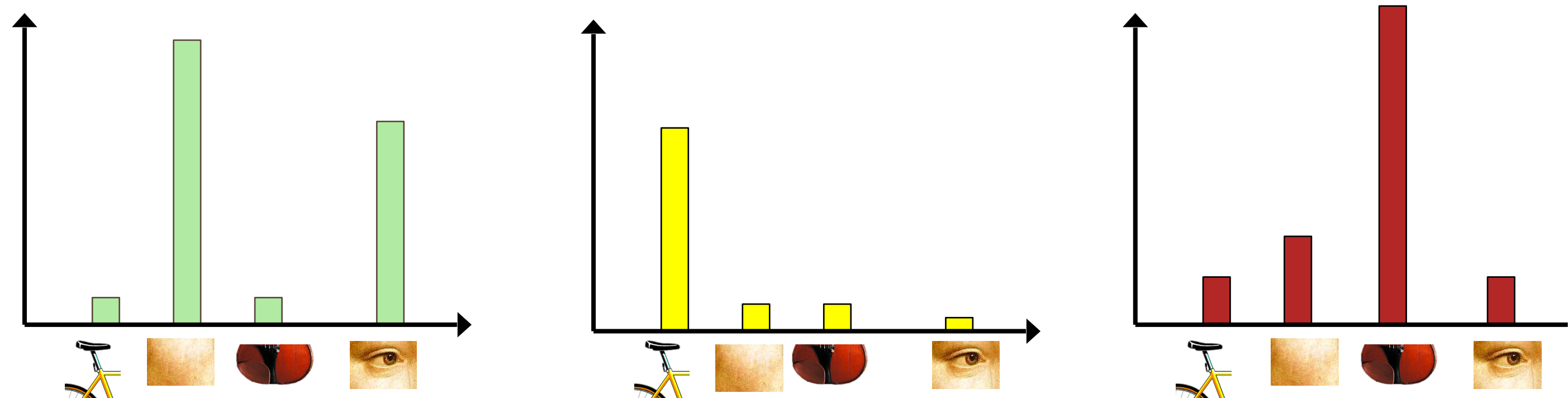
**sensory, brain,  
visual, perception,  
retinal, cerebral cortex,  
eye, cell, optical  
nerve, image  
Hubel, Wiesel**

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus will be created by a predicted 30% increase in exports to \$750bn, compared with \$575bn in 2004. The ministry said the surplus will not annoy the US. China's government has agreed to allow the yuan to rise against the dollar by 2.1% in 2005 and permitted it to trade within a narrow band but the US wants the yuan to be allowed to rise freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

**China, trade,  
surplus, commerce,  
exports, imports, US,  
yuan, bank, domestic,  
foreign, increase,  
trade, value**

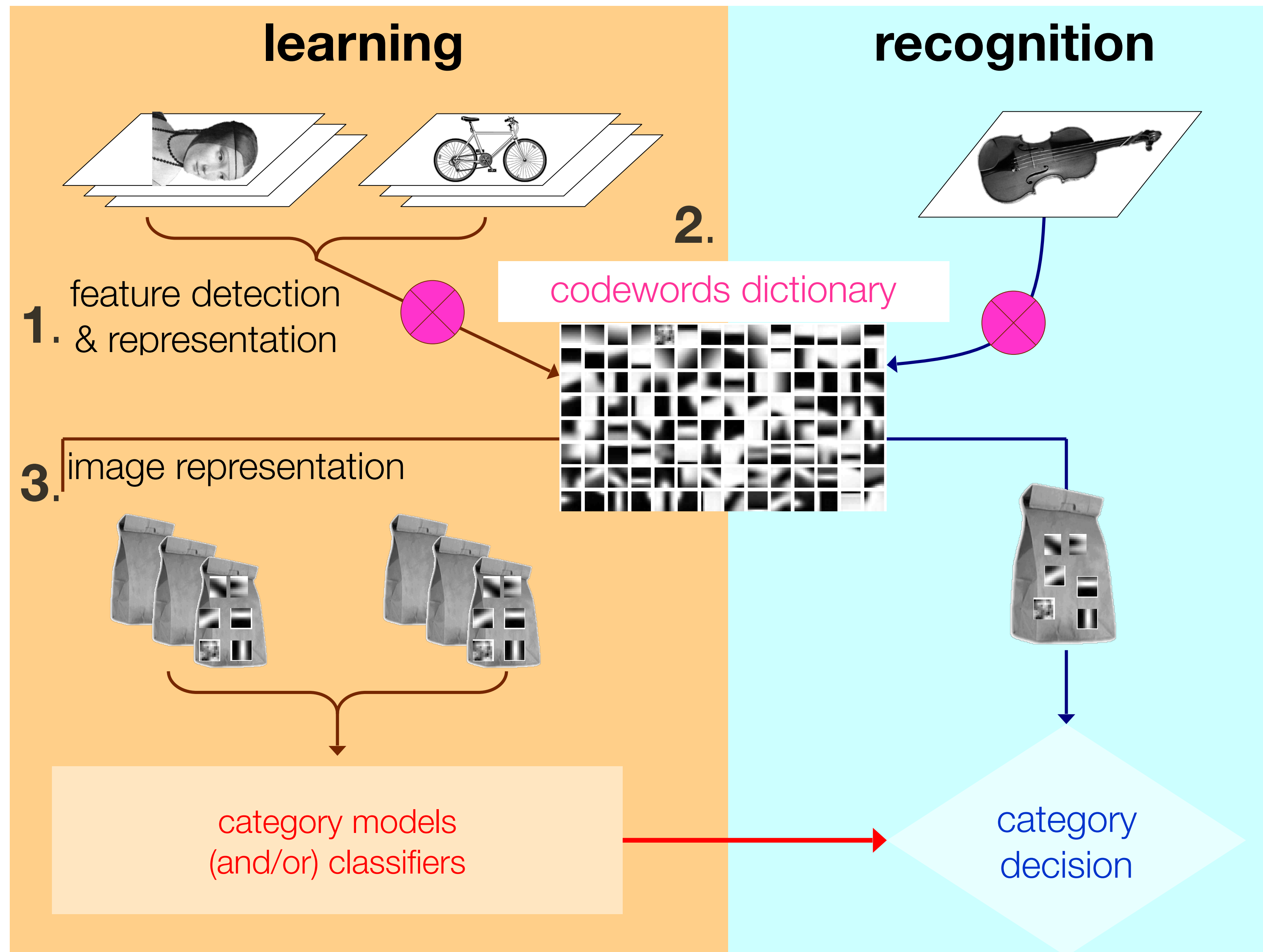


# Bag-of-Visual-Words



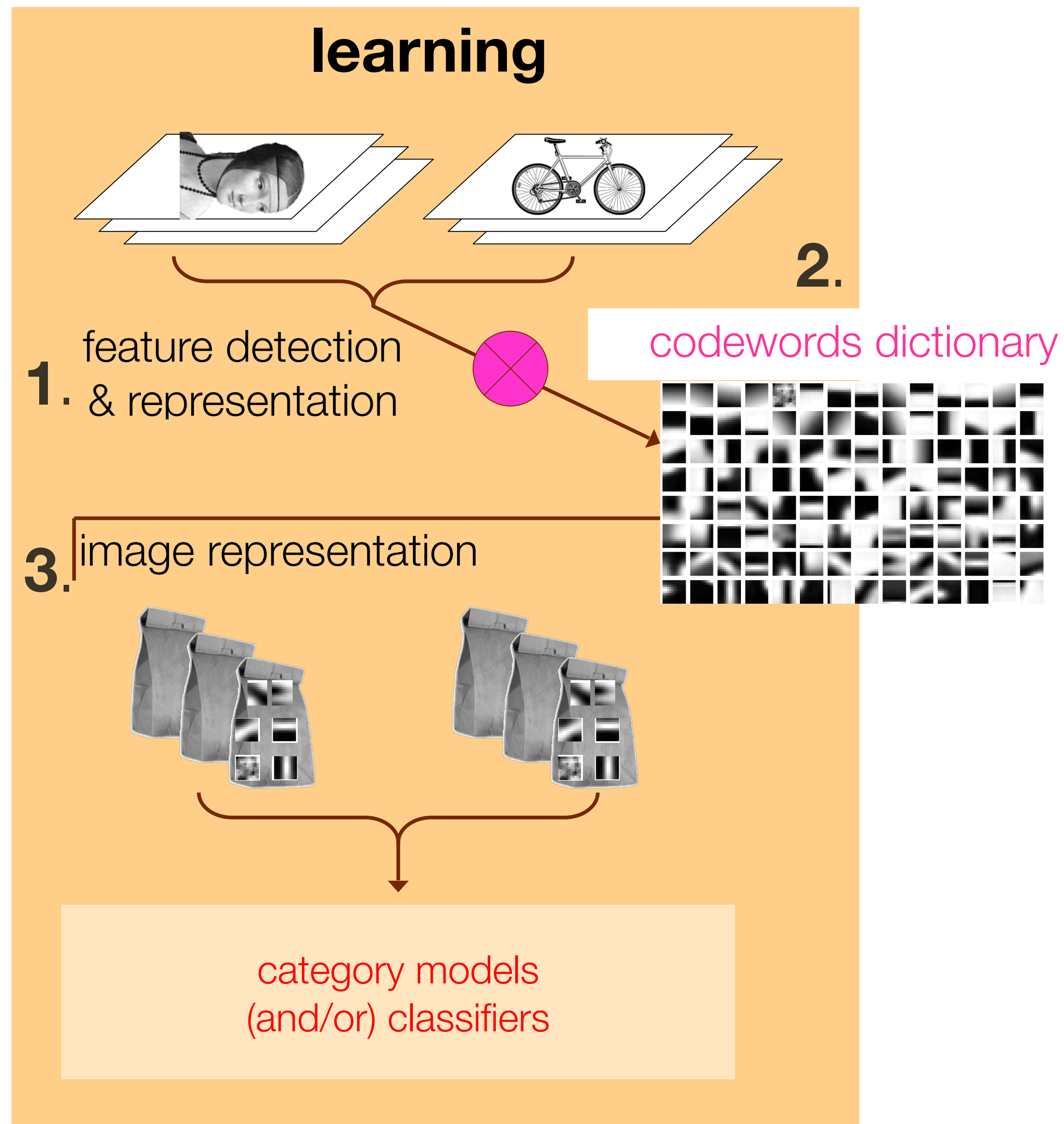


# Bag-of-Visual-Words





# Bag-of-**Visual**-Words: Learning





# Feature Detection & Representation

## Regular Grid

- Vogel et al. 2003
- Fei-Fei et al. 2005





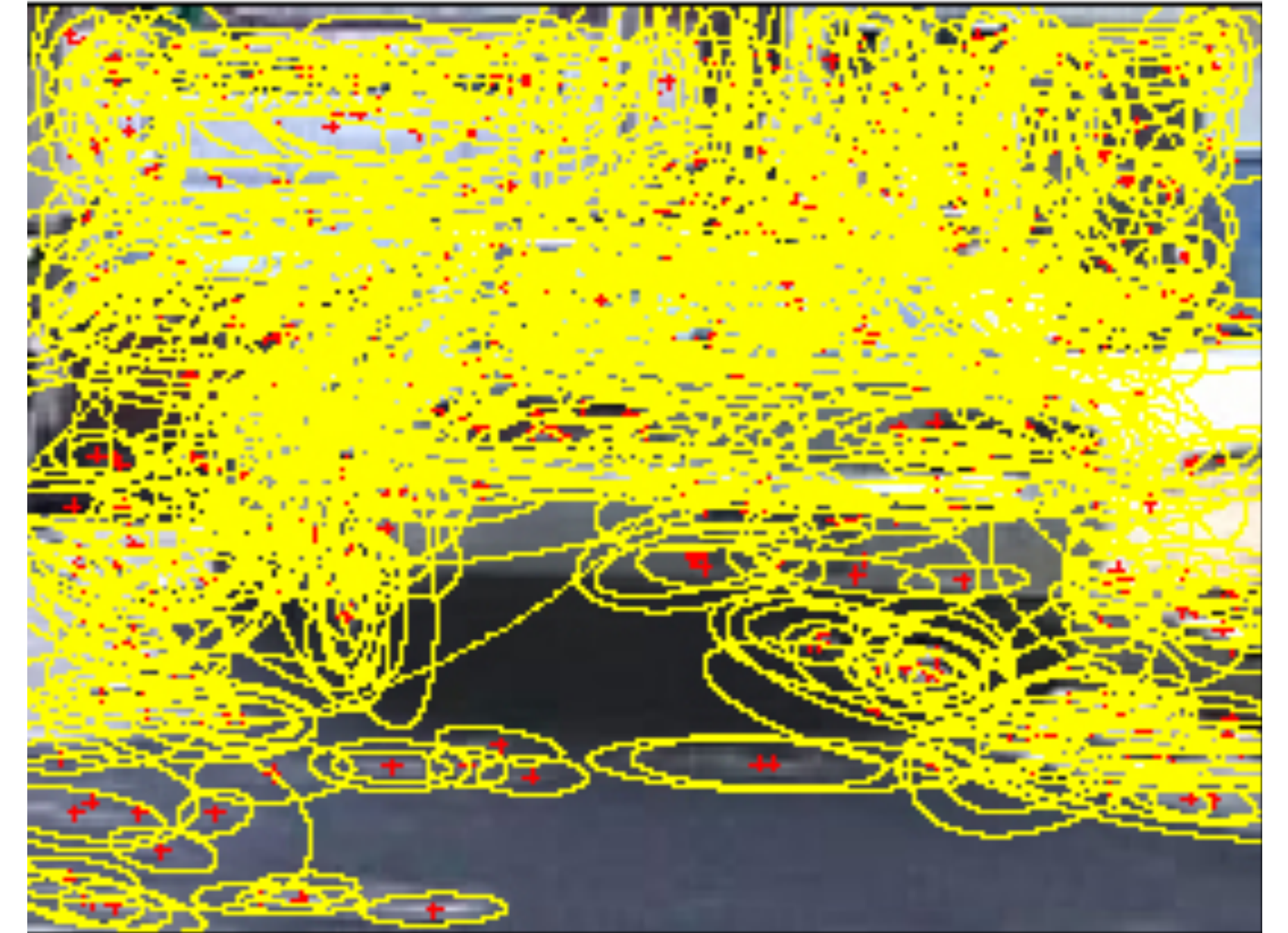
# Feature Detection & Representation

## Regular Grid

- Vogel et al. 2003
- Fei-Fei et al. 2005

## Interest Point Detector

- Csurka et al. 2004
- Fei-Fei et al. 2005
- Sivic et al. 2005





# Feature Detection & Representation

## Regular Grid

- Vogel et al. 2003
- Fei-Fei et al. 2005

## Interest Point Detector

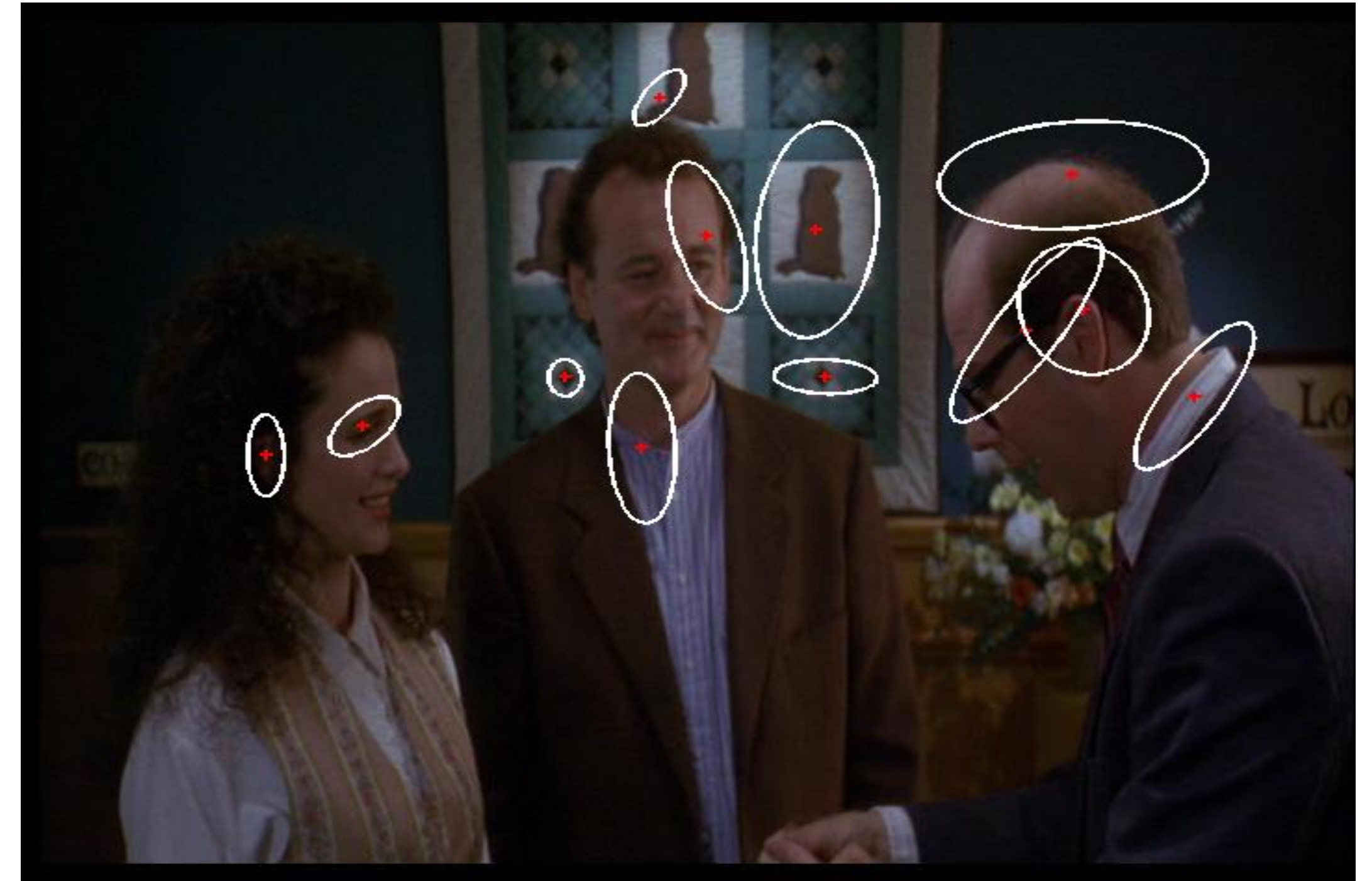
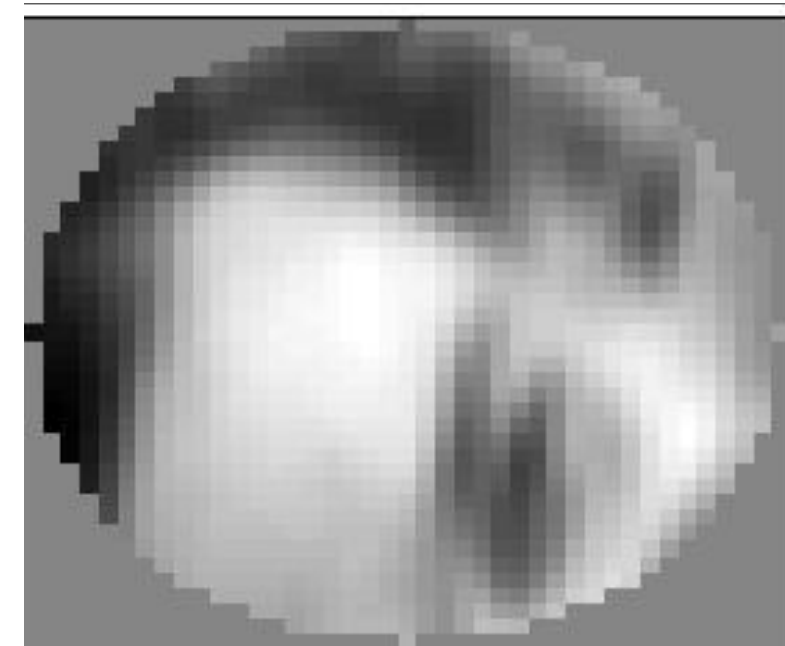
- Csurka et al. 2004
- Fei-Fei et al. 2005
- Sivic et al. 2005

## Other Methods

- Random sampling (Ullman et al. 2002)
- Segmentation based patches (Barnard et al. 2003)

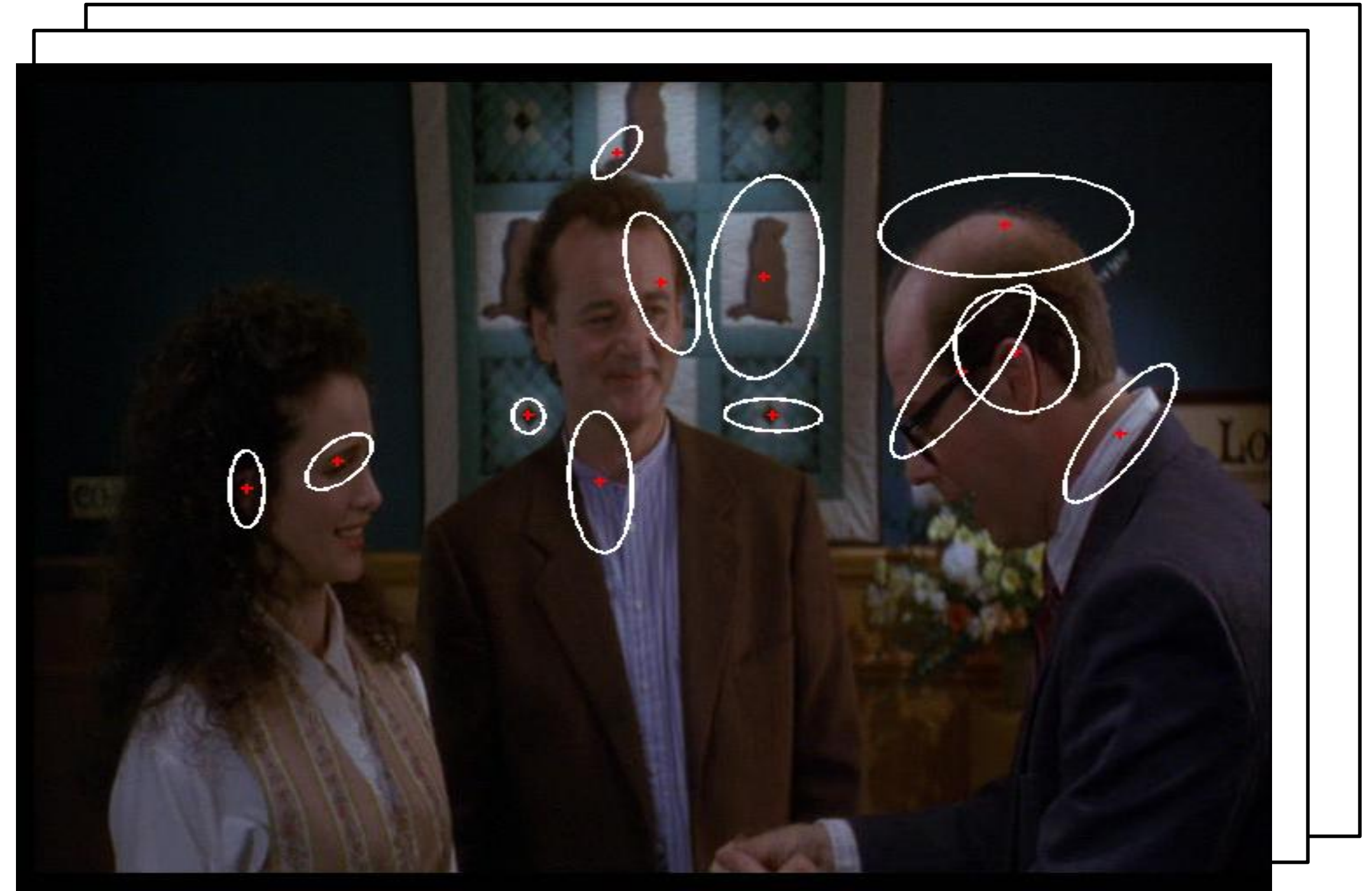
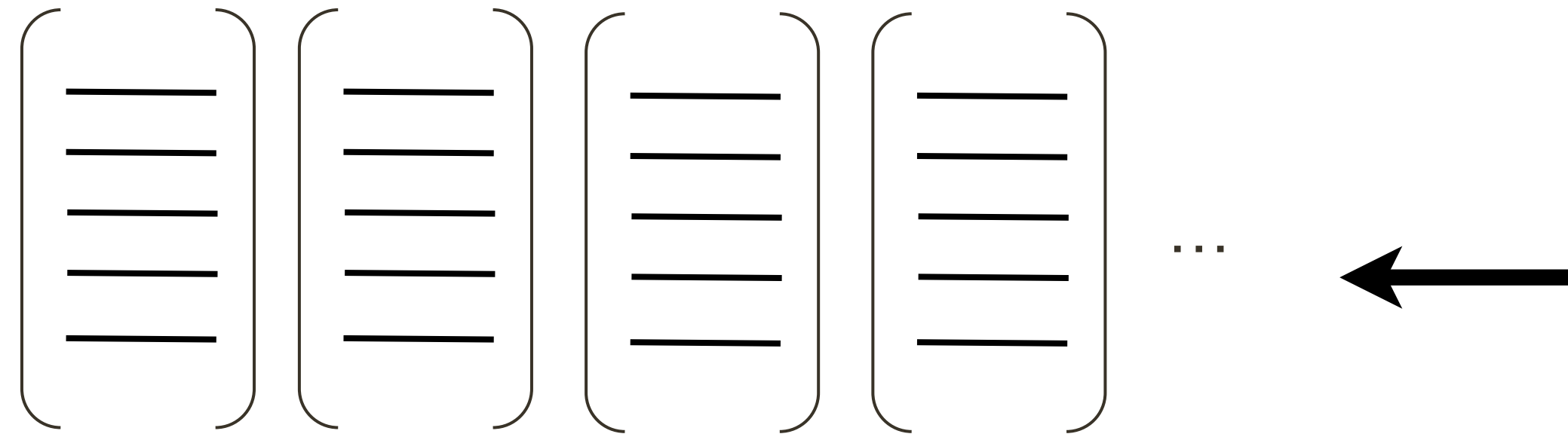


# Feature Detection & Representation



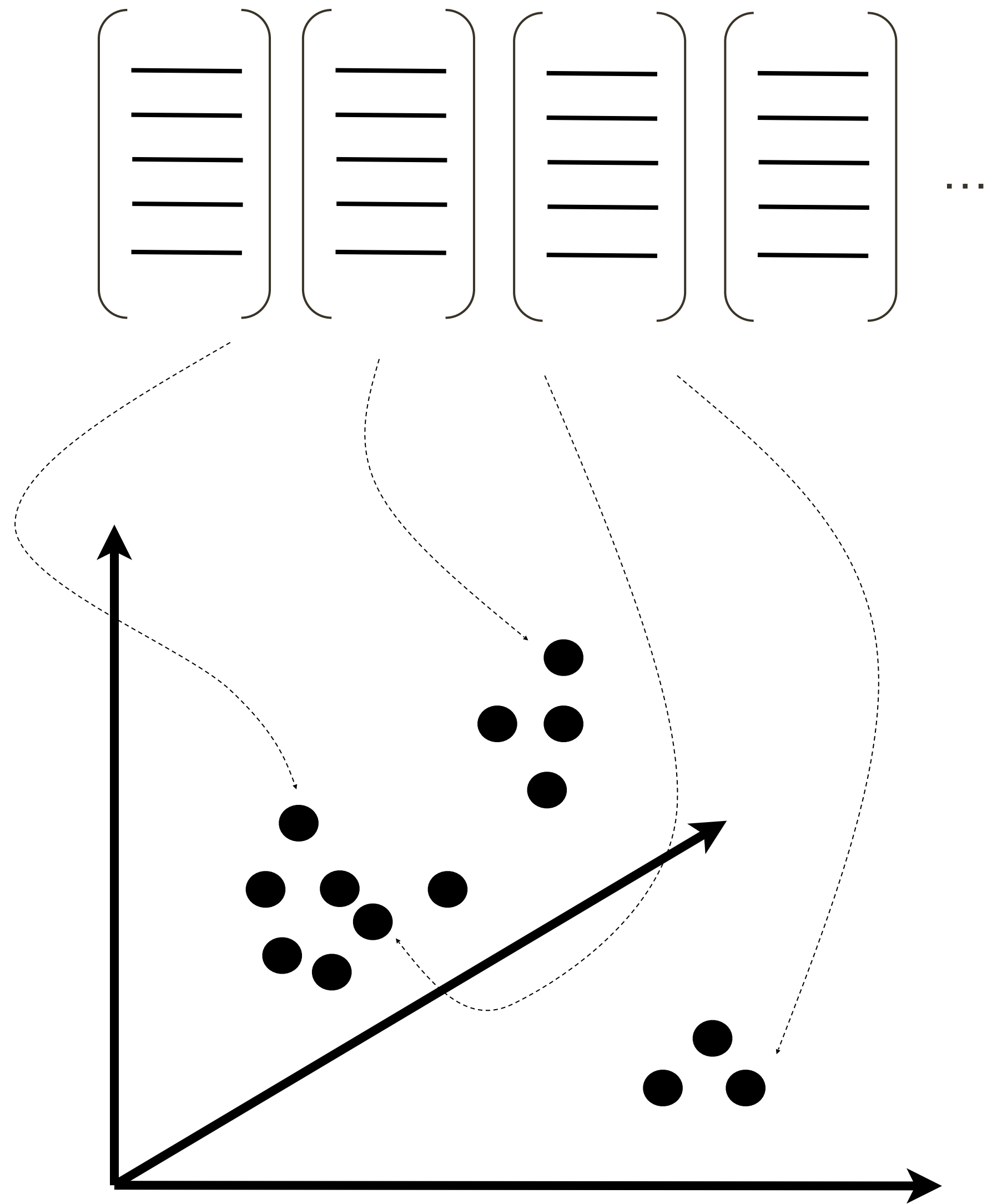


# Feature Detection & Representation



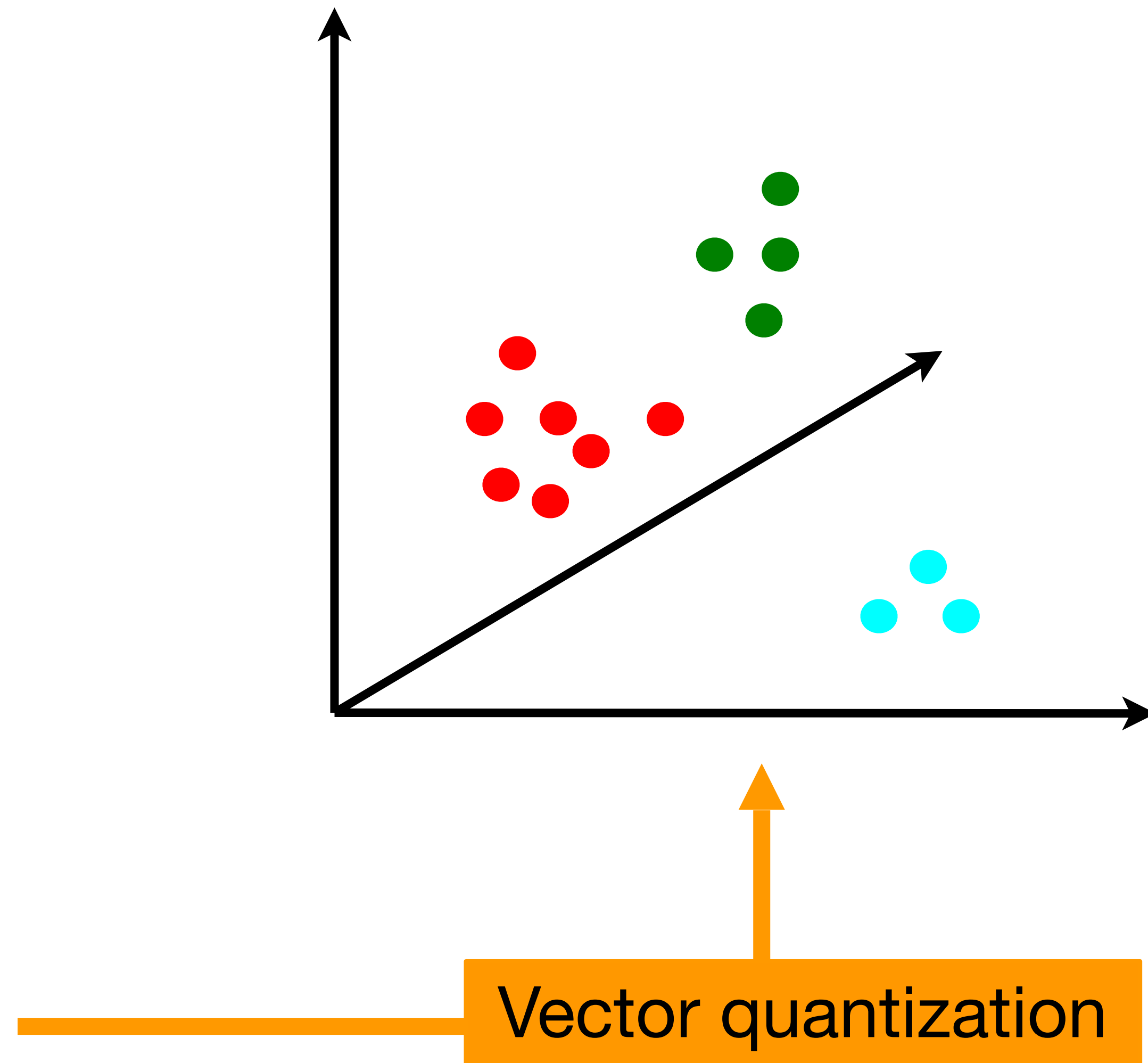
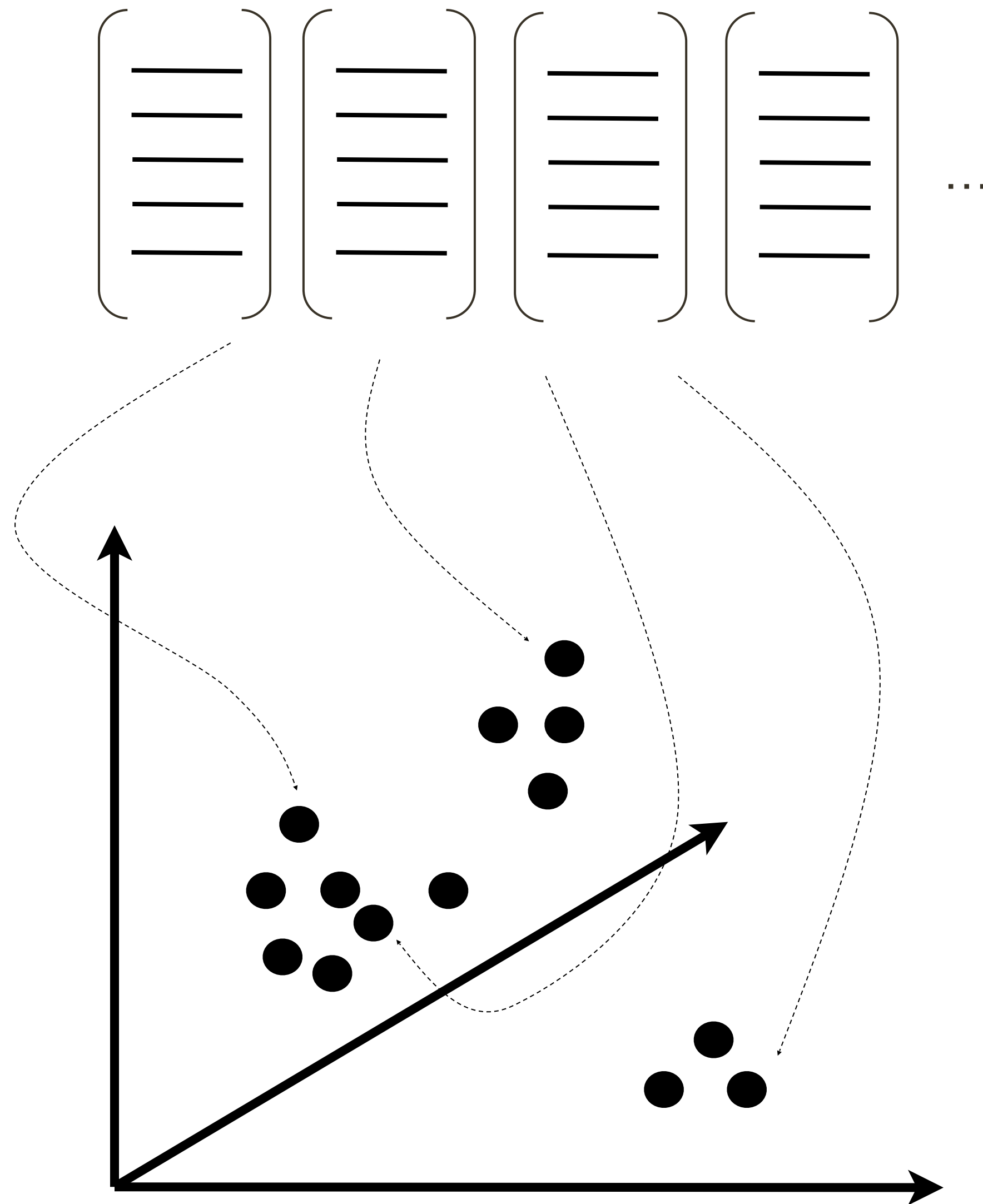


# Feature Detection & Representation

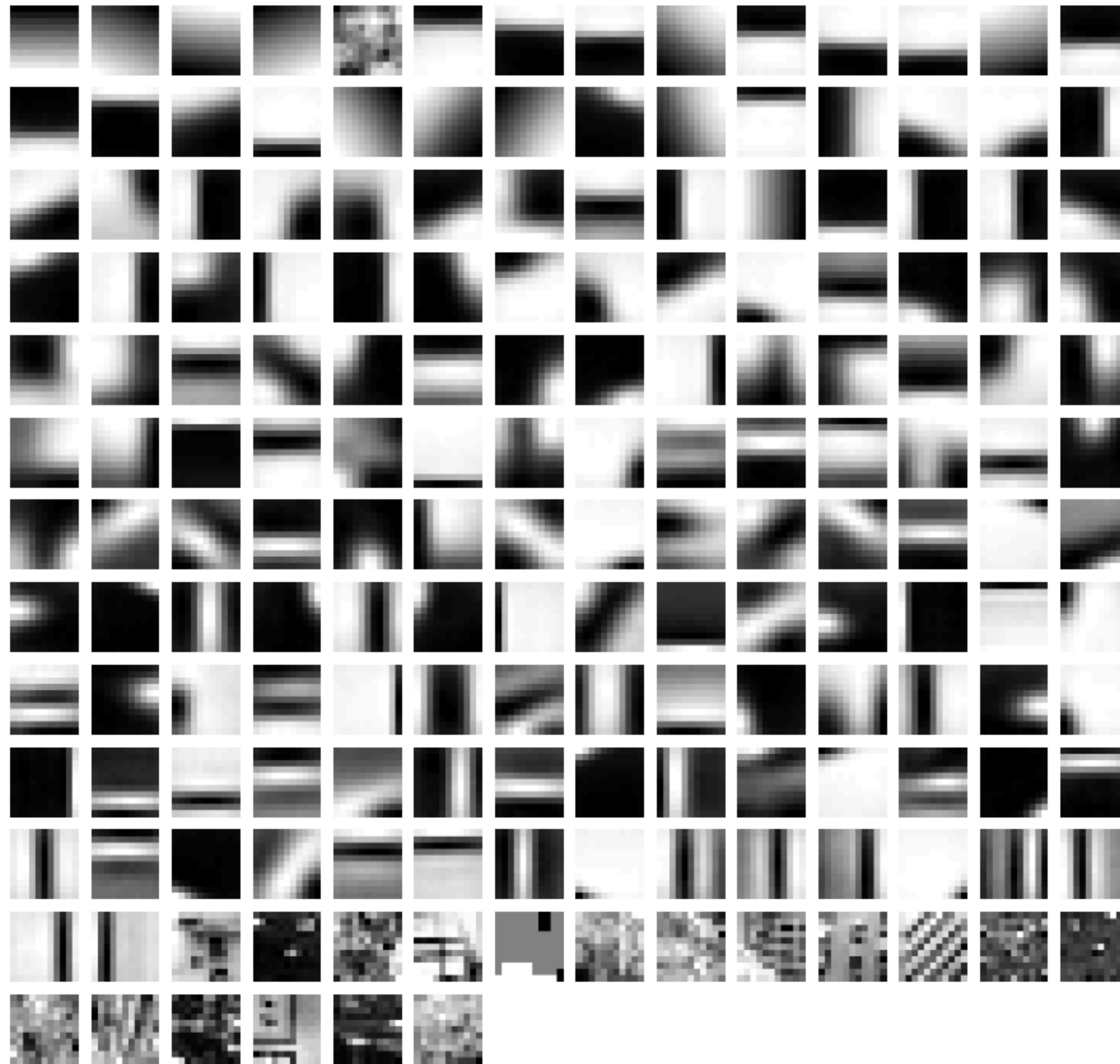




# Feature Detection & Representation



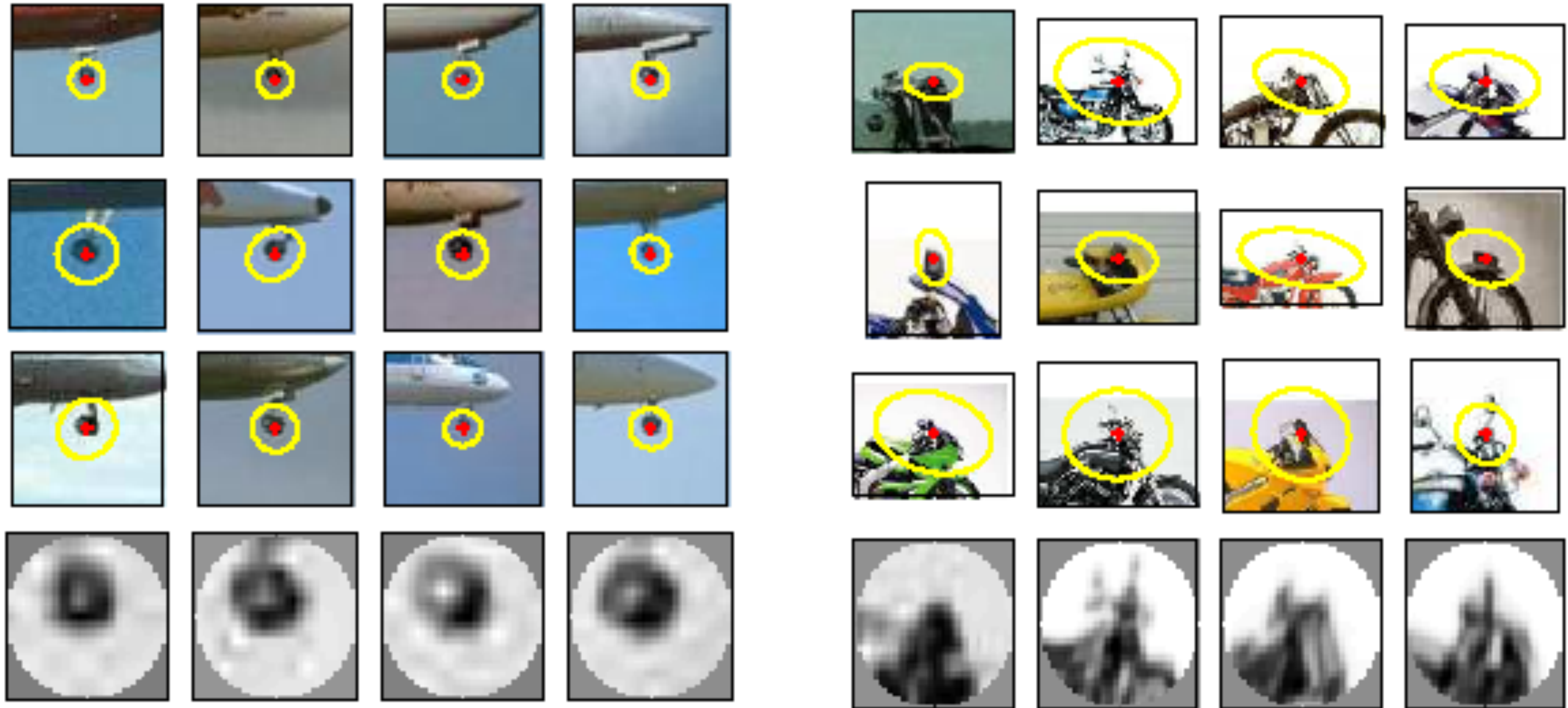
# Codeword Dictionary



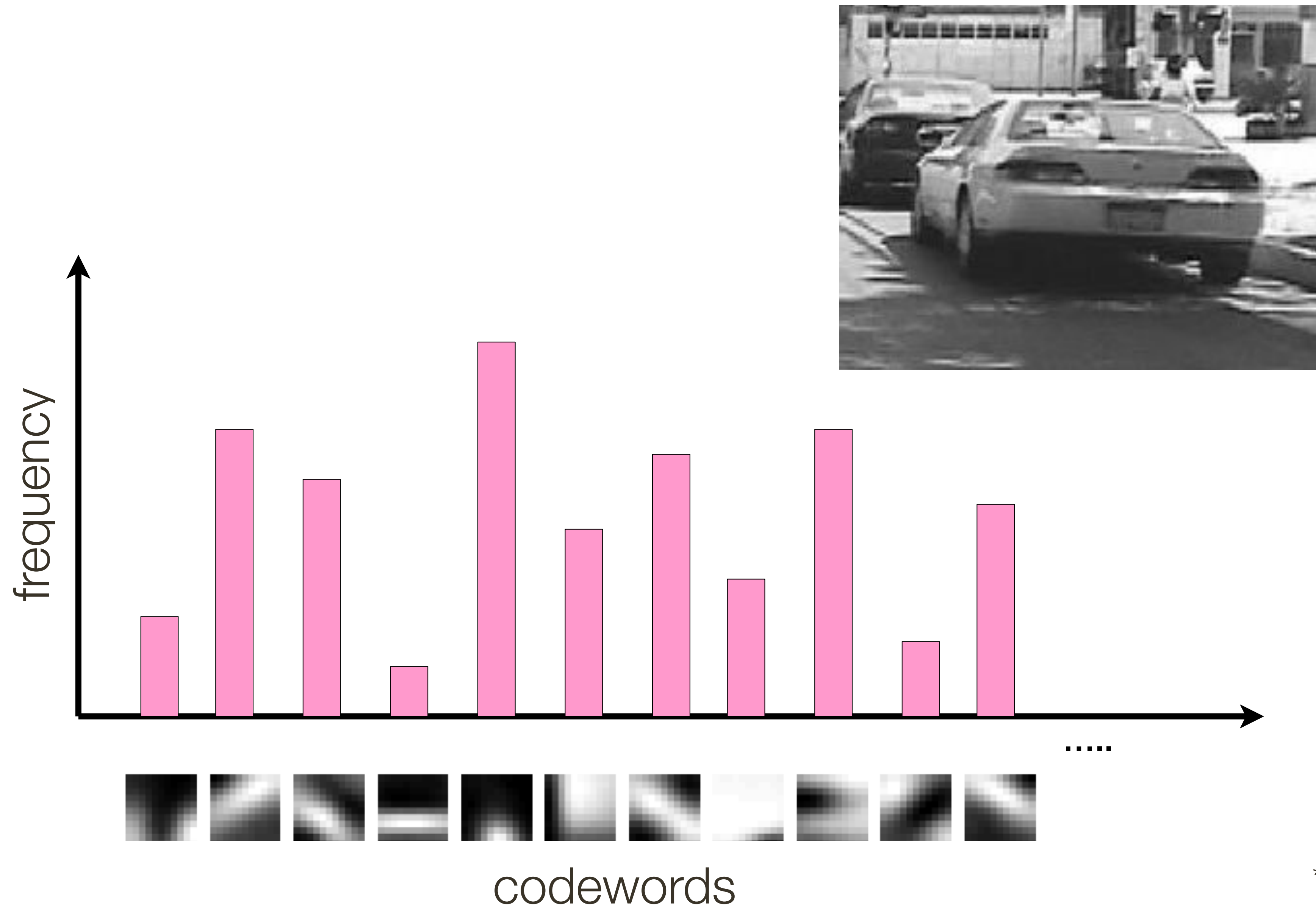
[ Fei-Fei et al., 2005 ]



# Image Patch Examples of Code Words



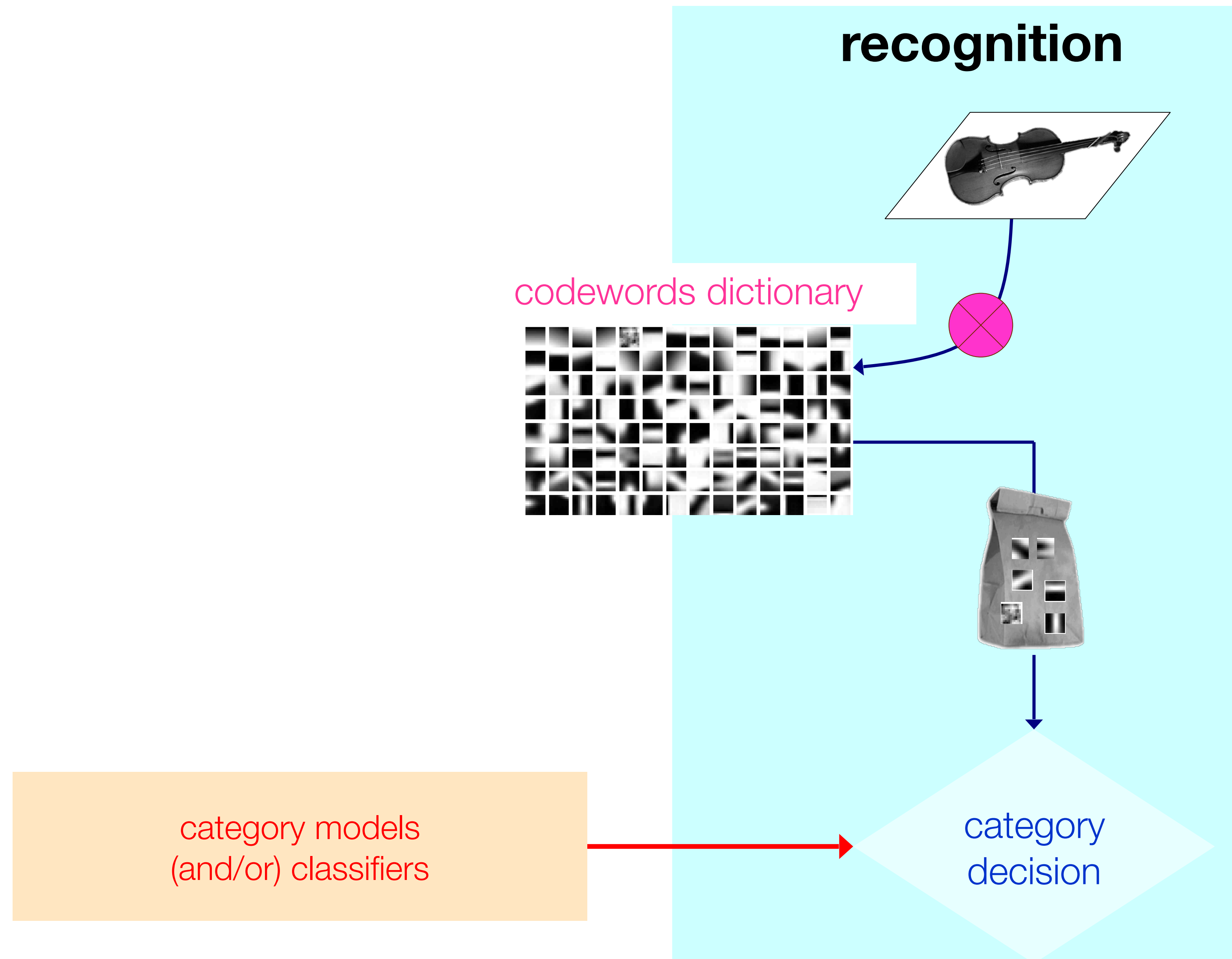
# Image Representation



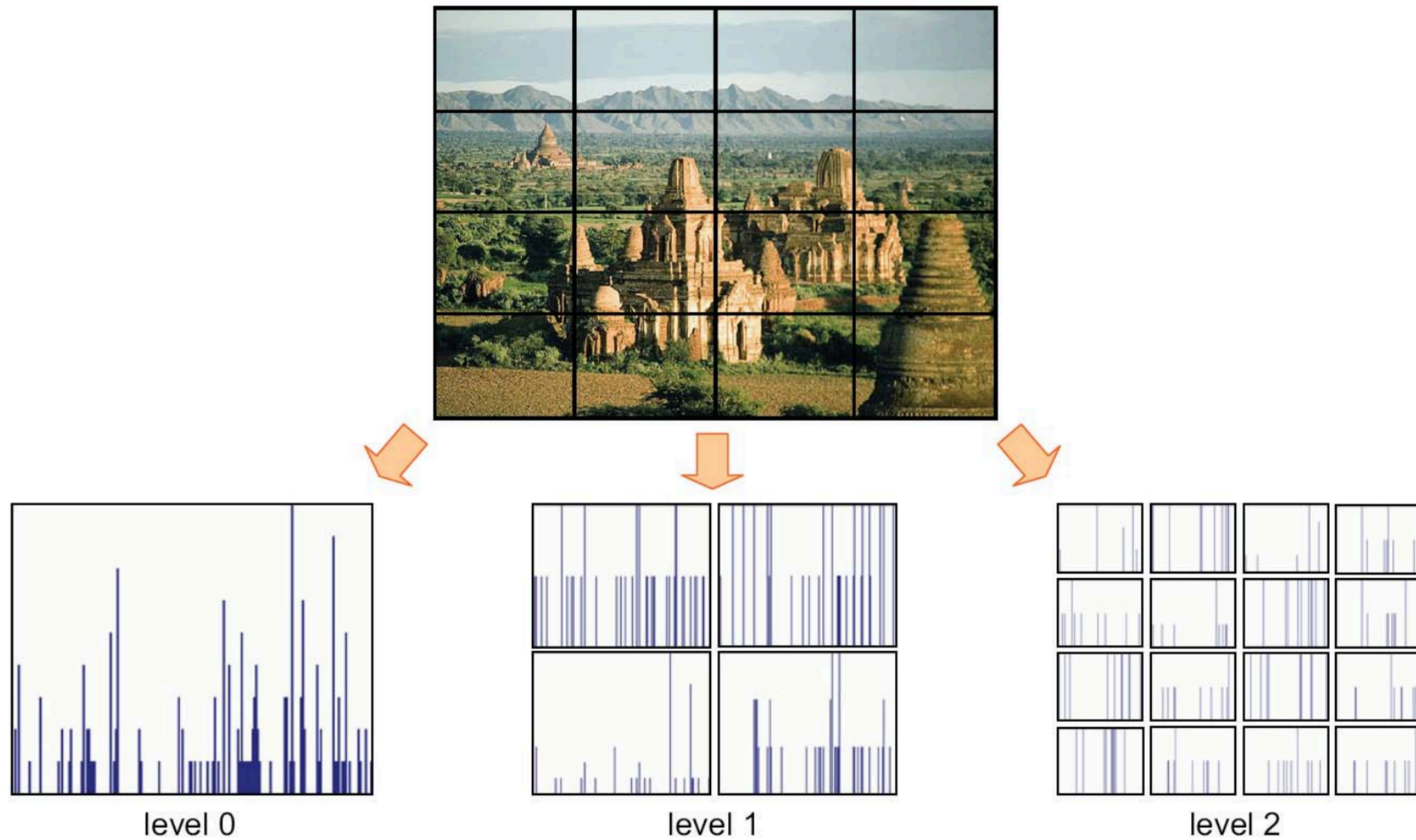
\*slide credit Li Fei-Fei



# Bag-of-**Visual**-Words



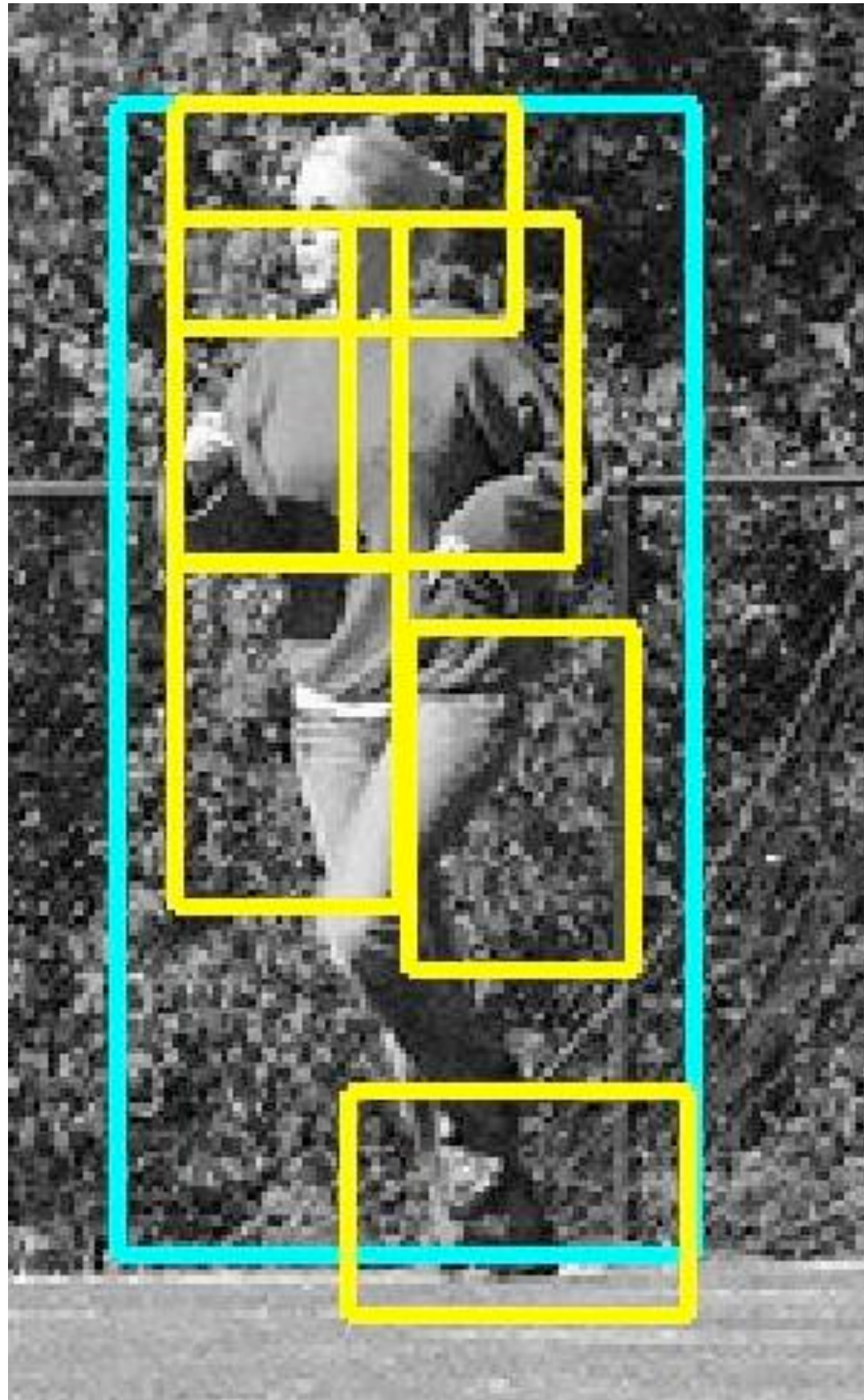
# Beyond Bag of Features



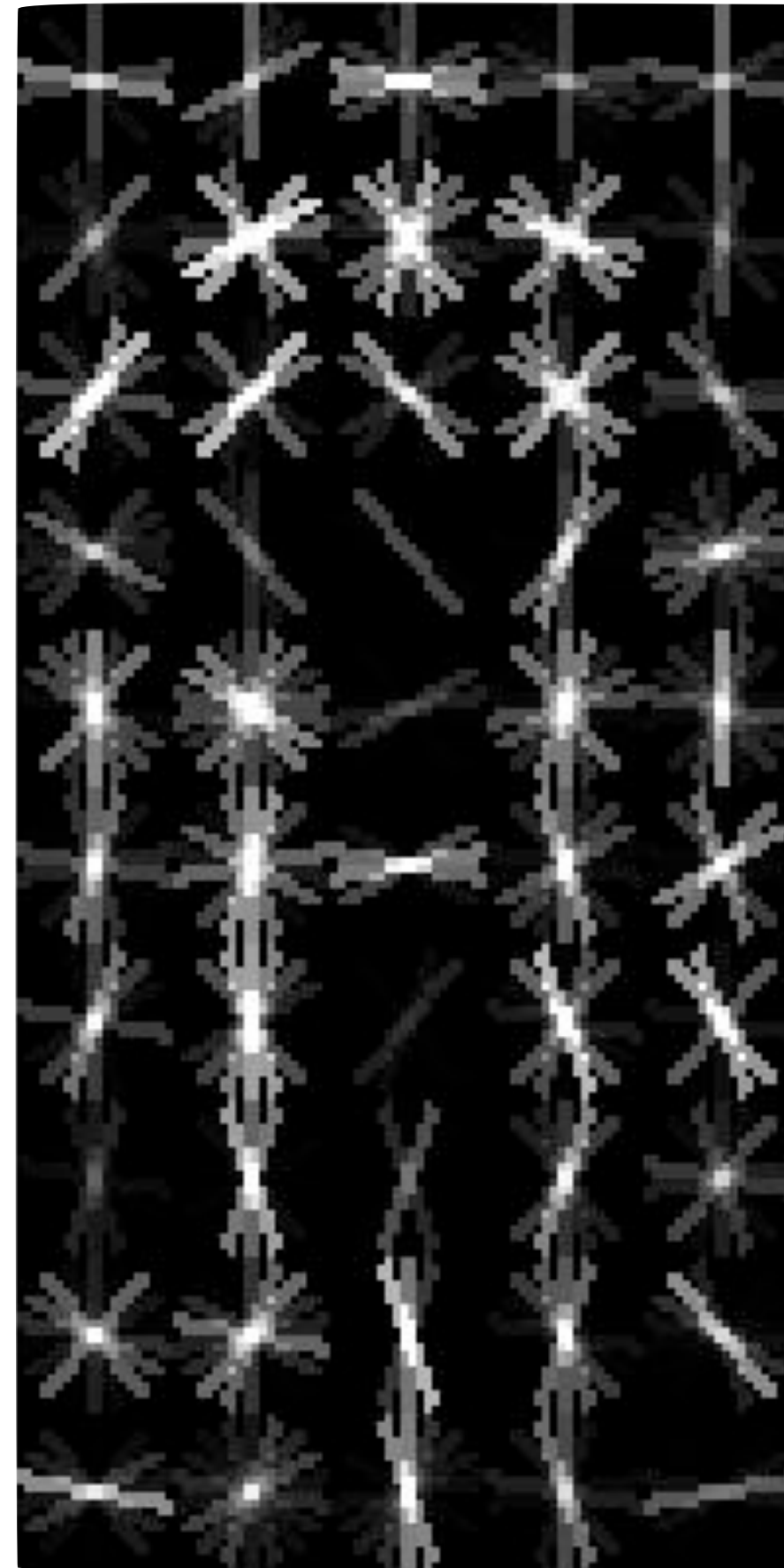
[ Lazebnik, Schmid, Ponce, 2006 ]



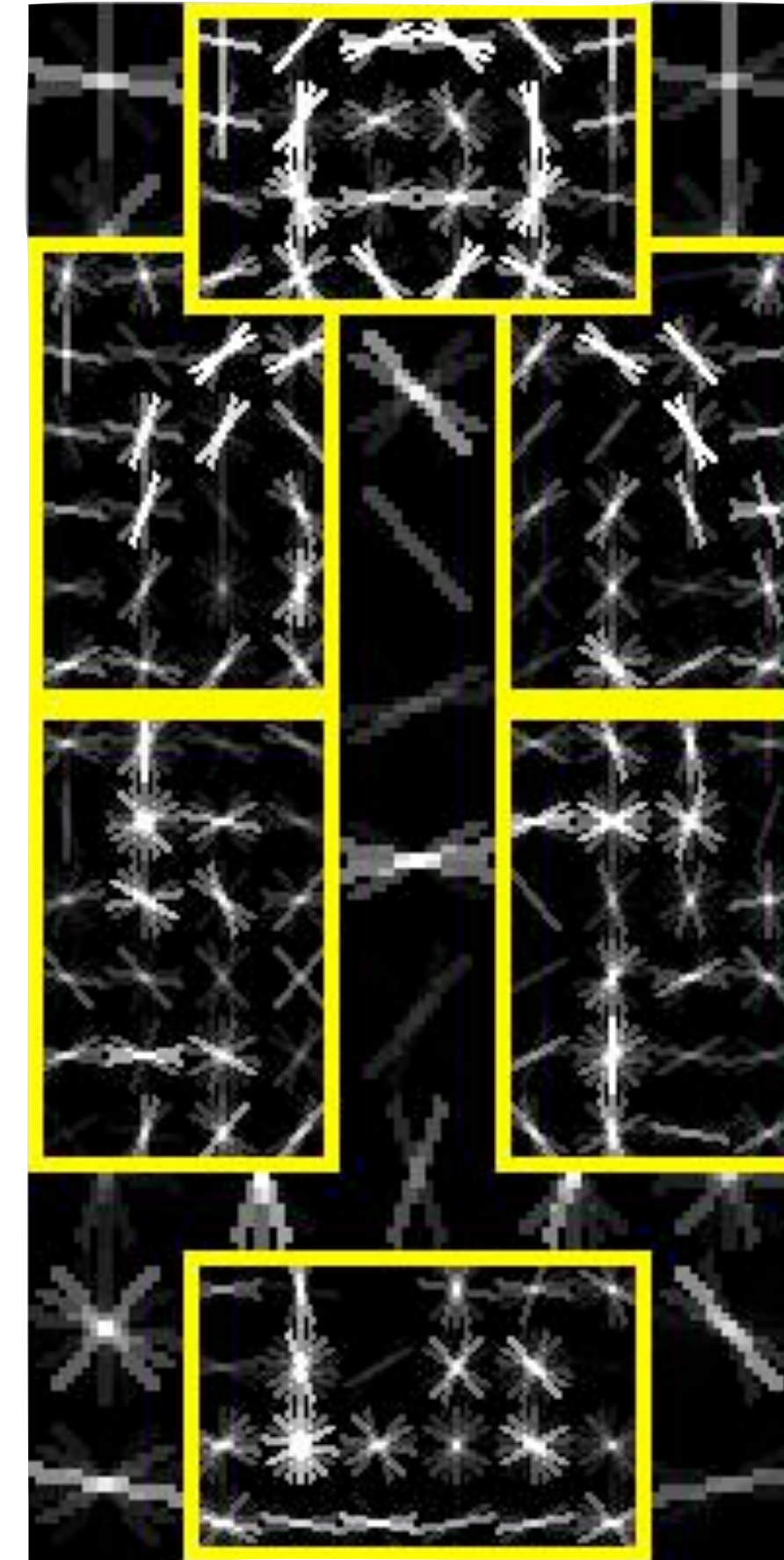
# Deformable Part Models



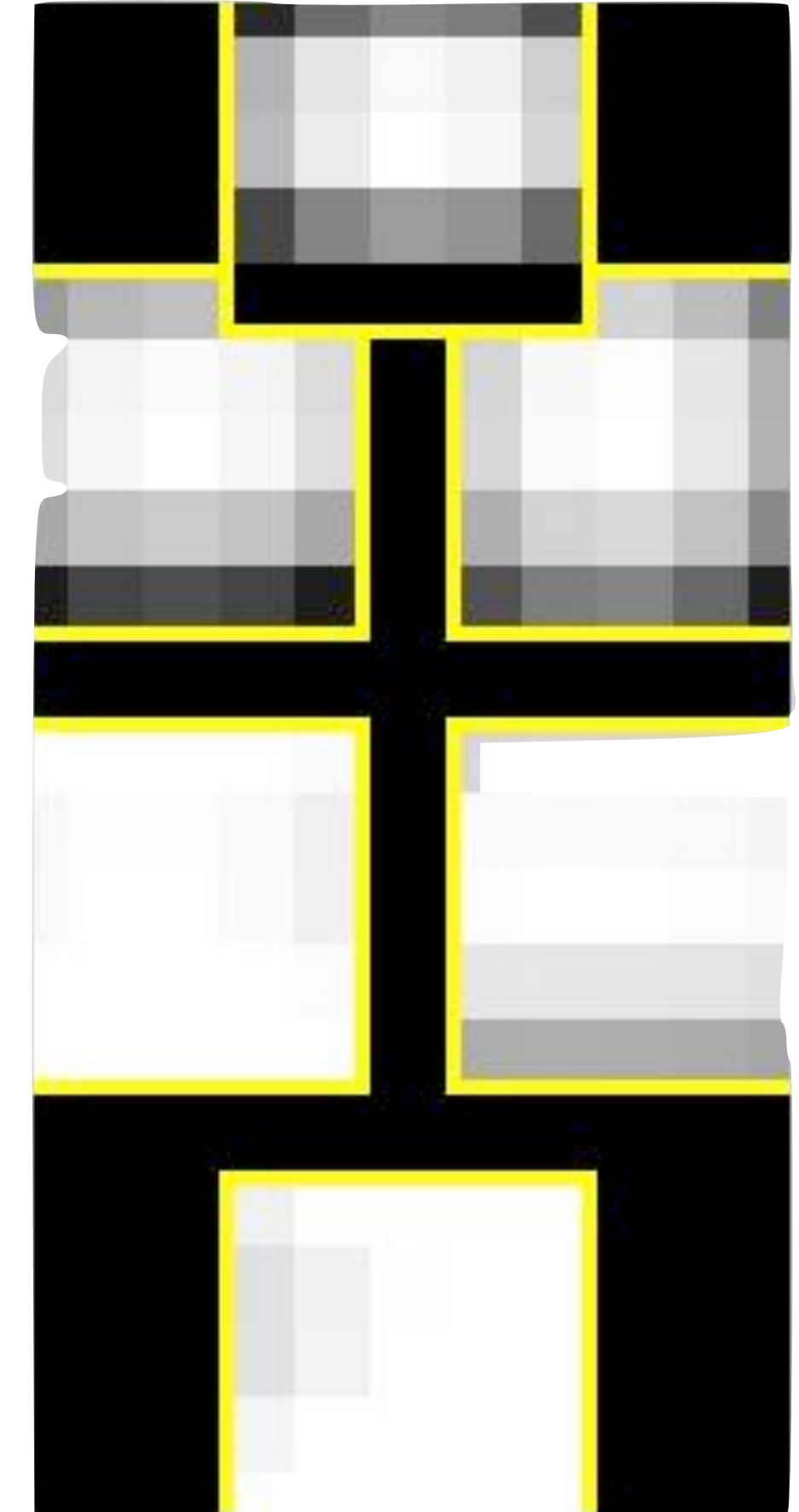
Detection



Root Filter



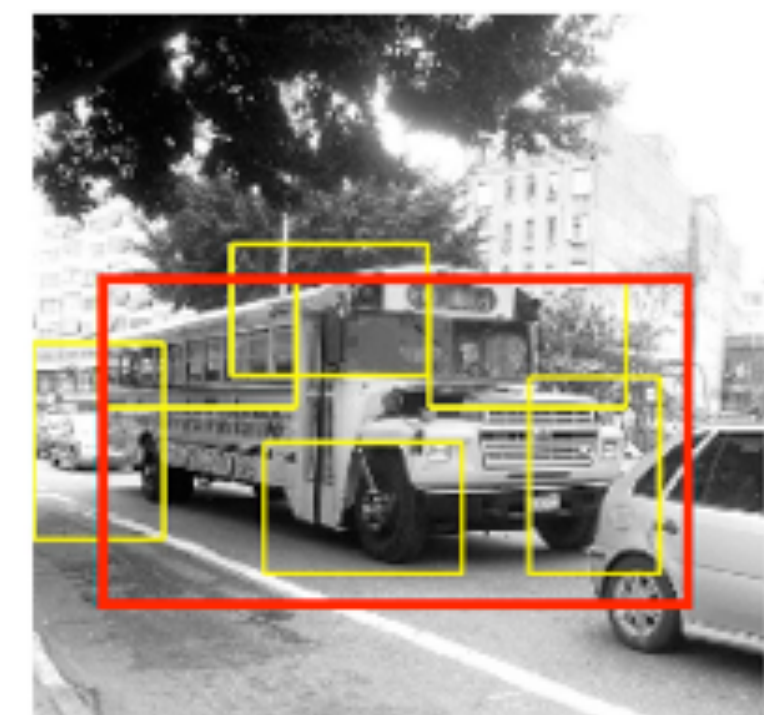
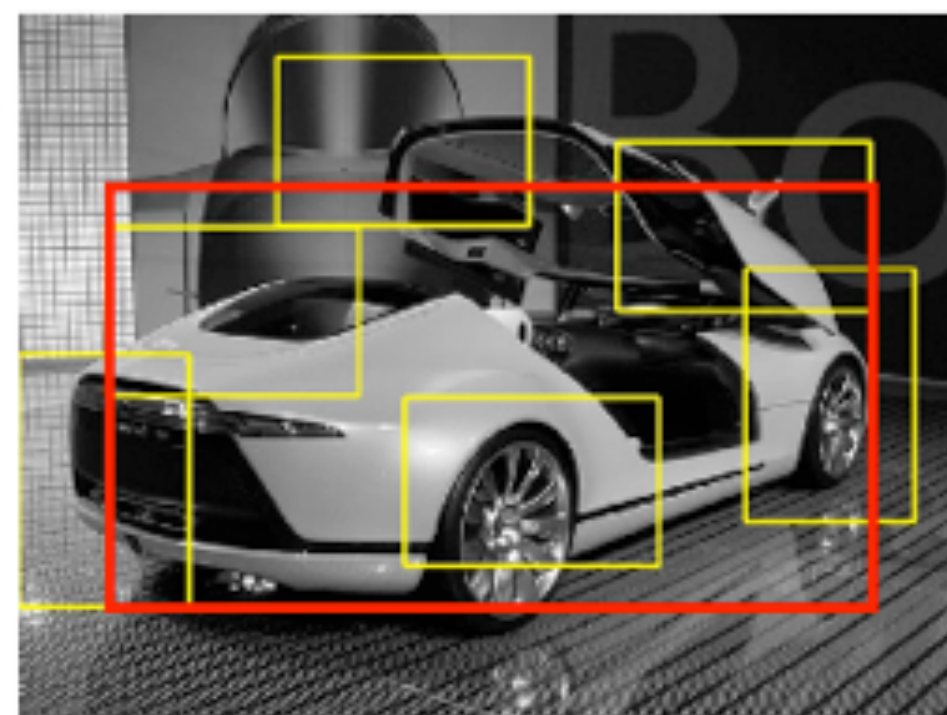
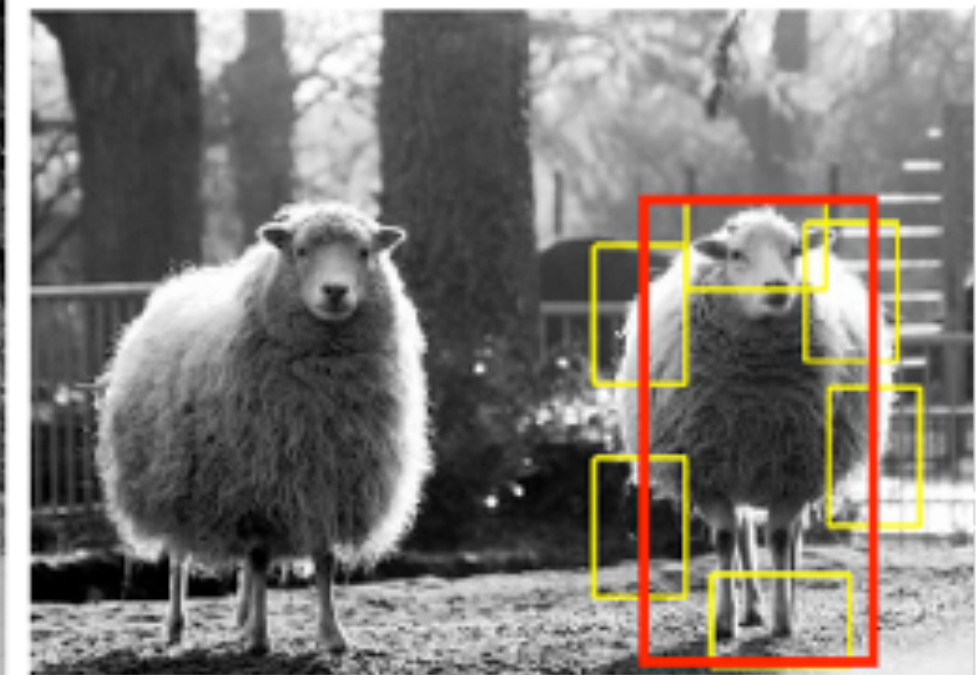
Part Filters



Deformations

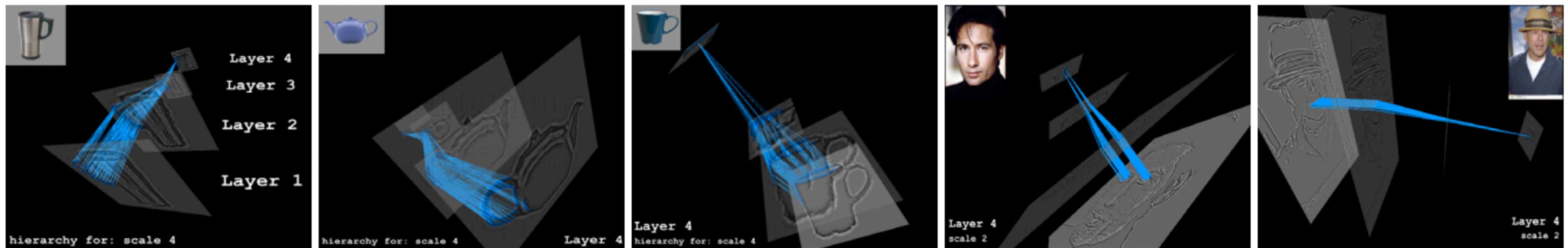
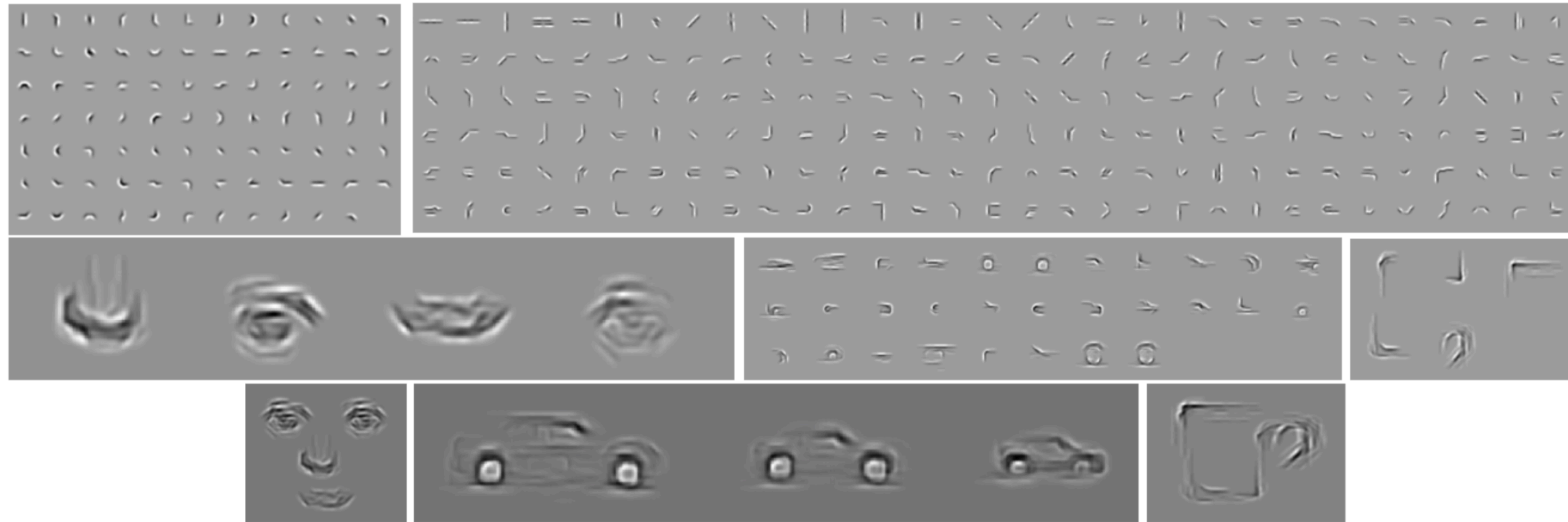


# Deformable Part Models





# Hierarchical Models





# PASCAL Visual Object Challenge (VOC)

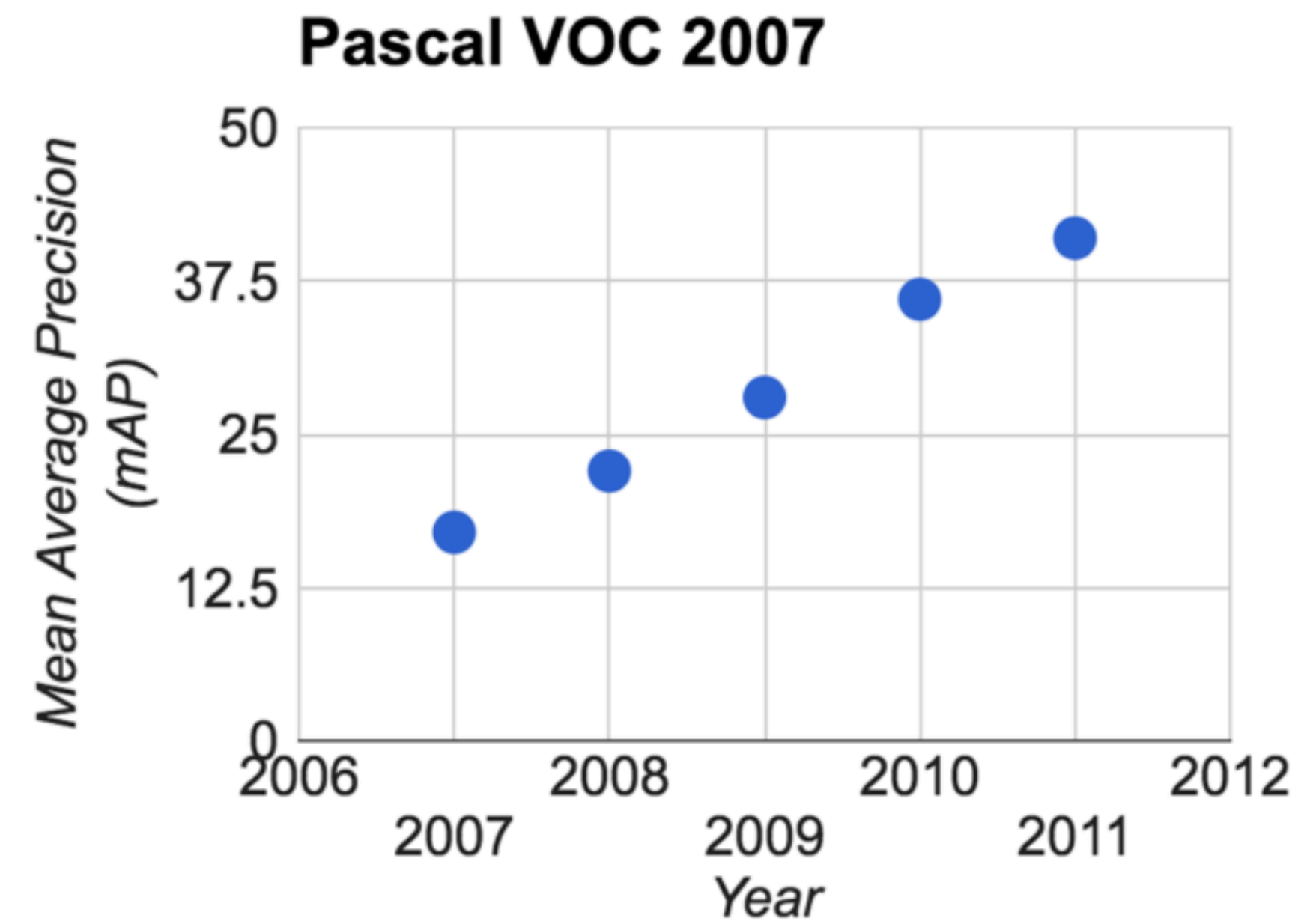
Image is CC BY-SA 3.0



Image is CC0 1.0 public domain



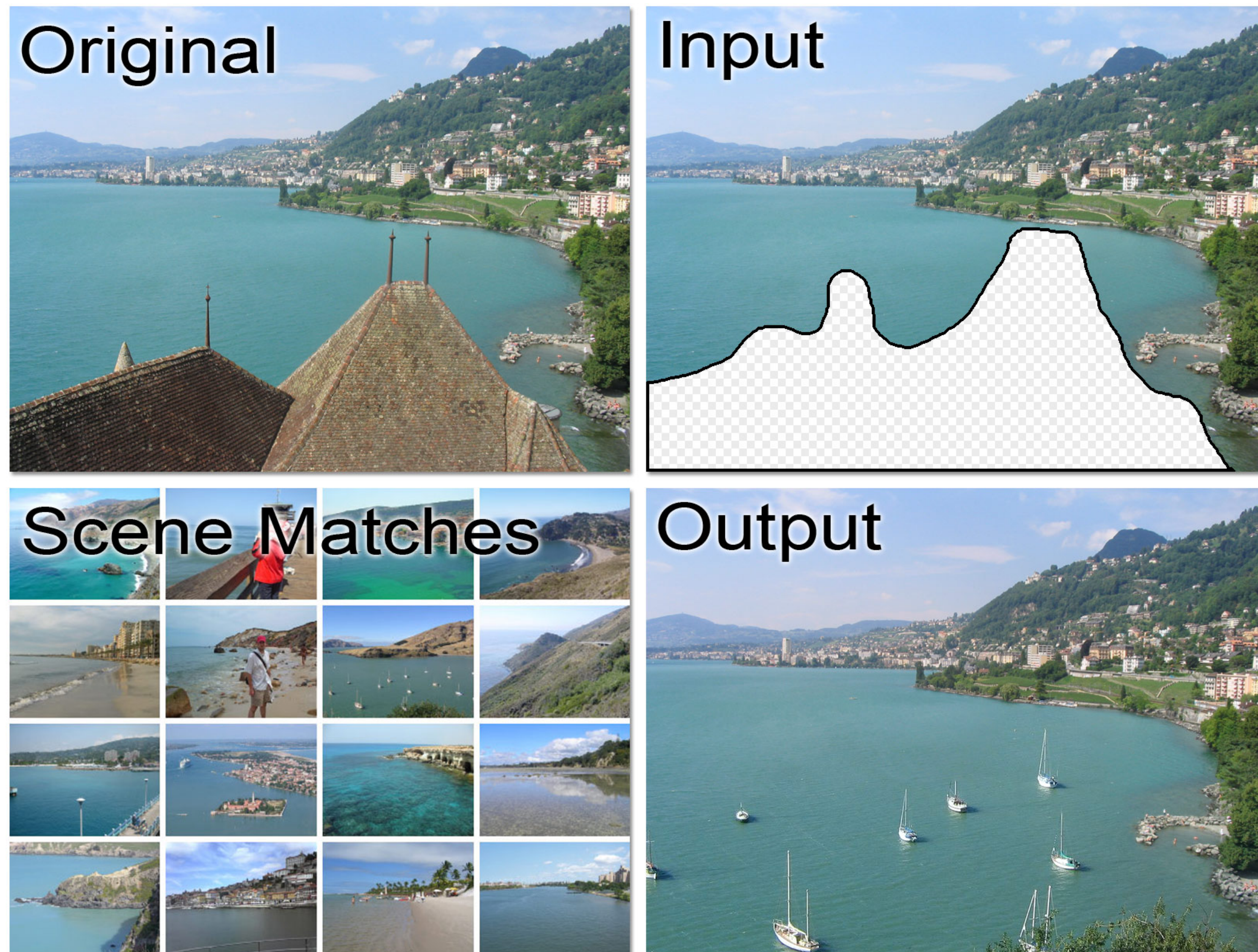
This image is licensed under  
CC BY-SA 2.0; changes made



[ Everingham et al. 2006-2012 ]



# Effectiveness of **Data**




[ Hays, Efros, ACM Siggraph 2007 ]



[ Hays, Efros, CVPR 2008 ]



# ImageNet Bechmark




**IMAGENET**

[www.image-net.org](http://www.image-net.org)

**22K** categories and **14M** images

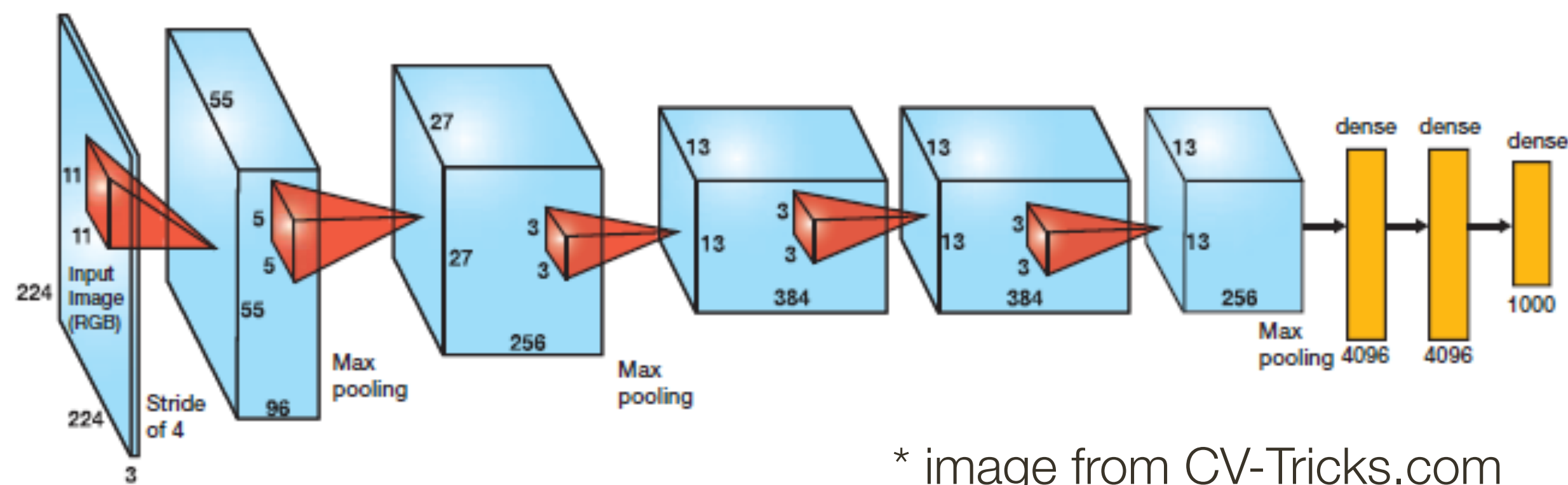
- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plants
  - Tree
  - Flower
  - Food
  - Materials
- Structures
  - Artifact
  - Tools
  - Appliances
  - Structures
- Person
  - Scenes
    - Indoor
    - Geological Formations
  - Sport Activities



Deng, Dong, Socher, Li, Li, & Fei-Fei, 2009



# AlexNet on ImageNet



\* image from [CV-Tricks.com](http://CV-Tricks.com)

Model	Top-1	Top-5
<i>Sparse coding [2]</i>	47.1%	28.2%
<i>SIFT + FVs [24]</i>	45.7%	25.7%
CNN	<b>37.5%</b>	<b>17.0%</b>

## ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky  
University of Toronto  
kriz@cs.utoronto.ca

Ilya Sutskever  
University of Toronto  
ilya@cs.utoronto.ca

Geoffrey E. Hinton  
University of Toronto  
hinton@cs.utoronto.ca

### Abstract

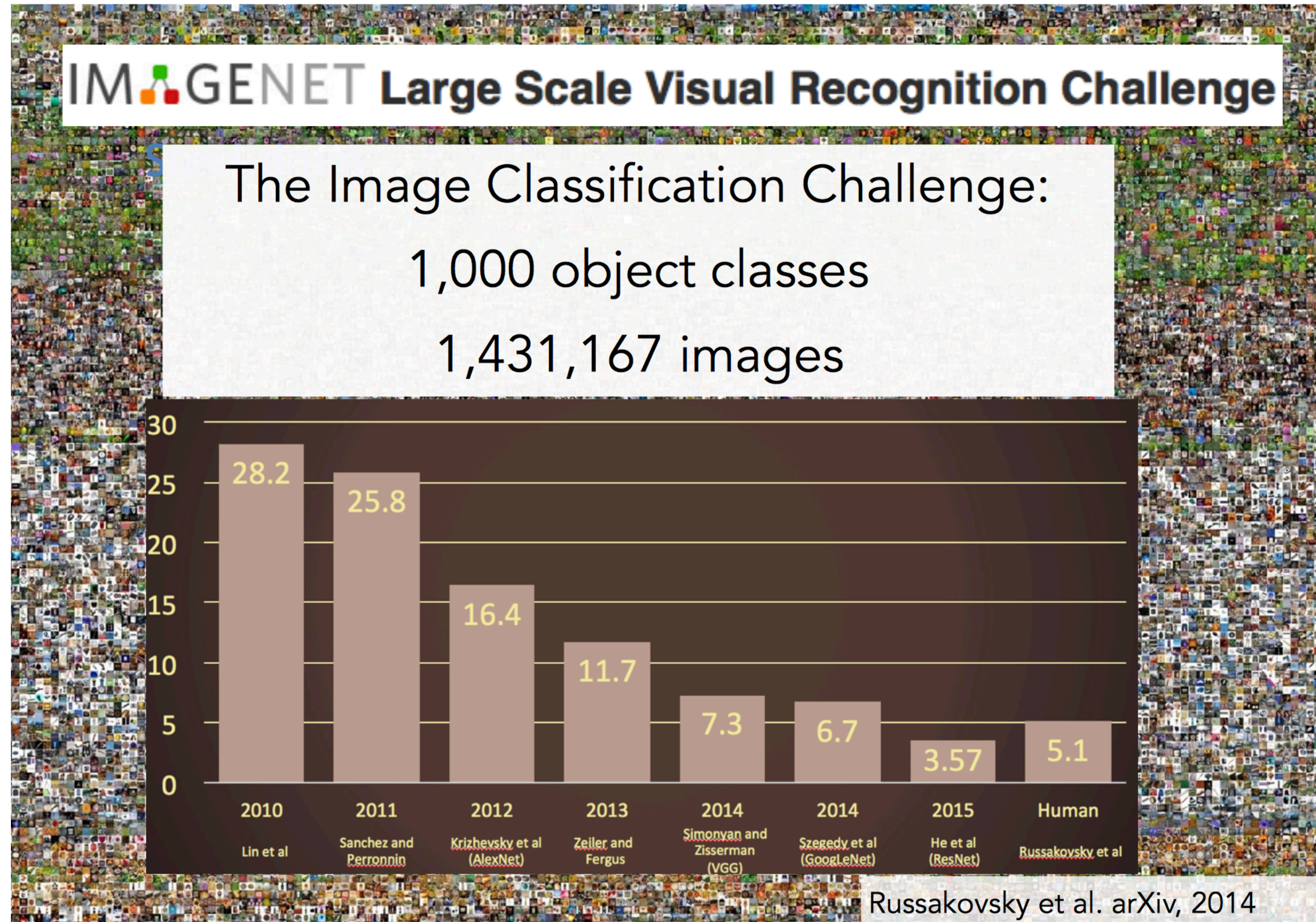
We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet ILSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called “dropout” that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
<i>SIFT + FVs [7]</i>	—	—	26.2%
1 CNN	40.7%	18.2%	—
5 CNNs	38.1%	16.4%	<b>16.4%</b>
1 CNN*	39.0%	16.6%	—
7 CNNs*	36.7%	15.4%	<b>15.3%</b>

[ Krizhevsky, Sutskever, Hinton, NIPS 2012 ]



# Success of **Deep Learning**





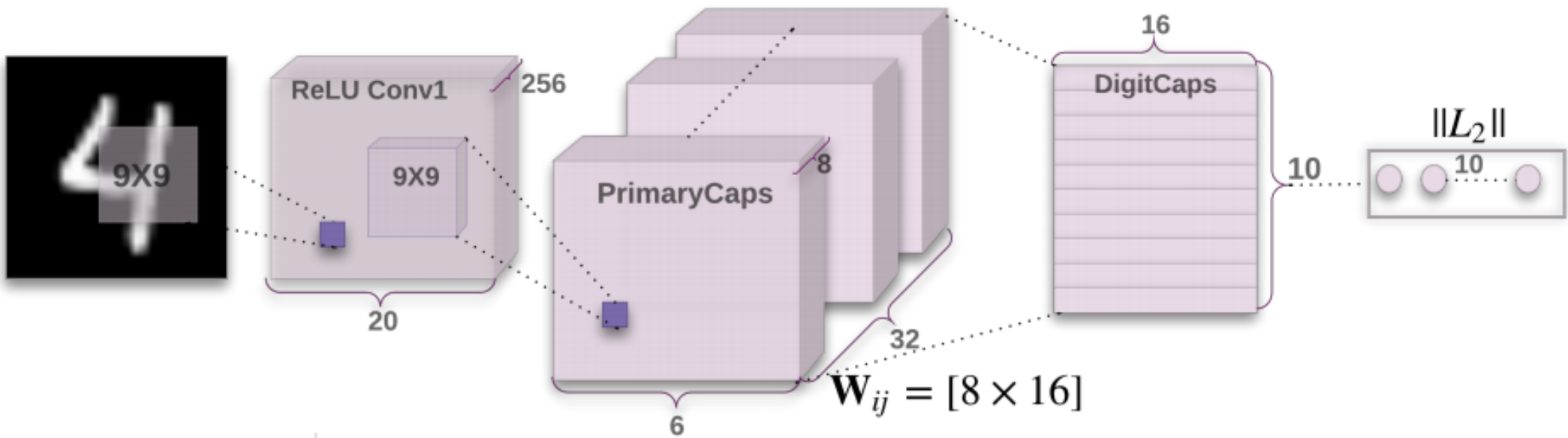
# Final thought ...

- Model based, compositional, primitives, inverse graphics
- Hand-crafted features for given invariances & matching
- Hand-crafted features with learned statistical models on top
- Joint learning of features and statistical models for recognition

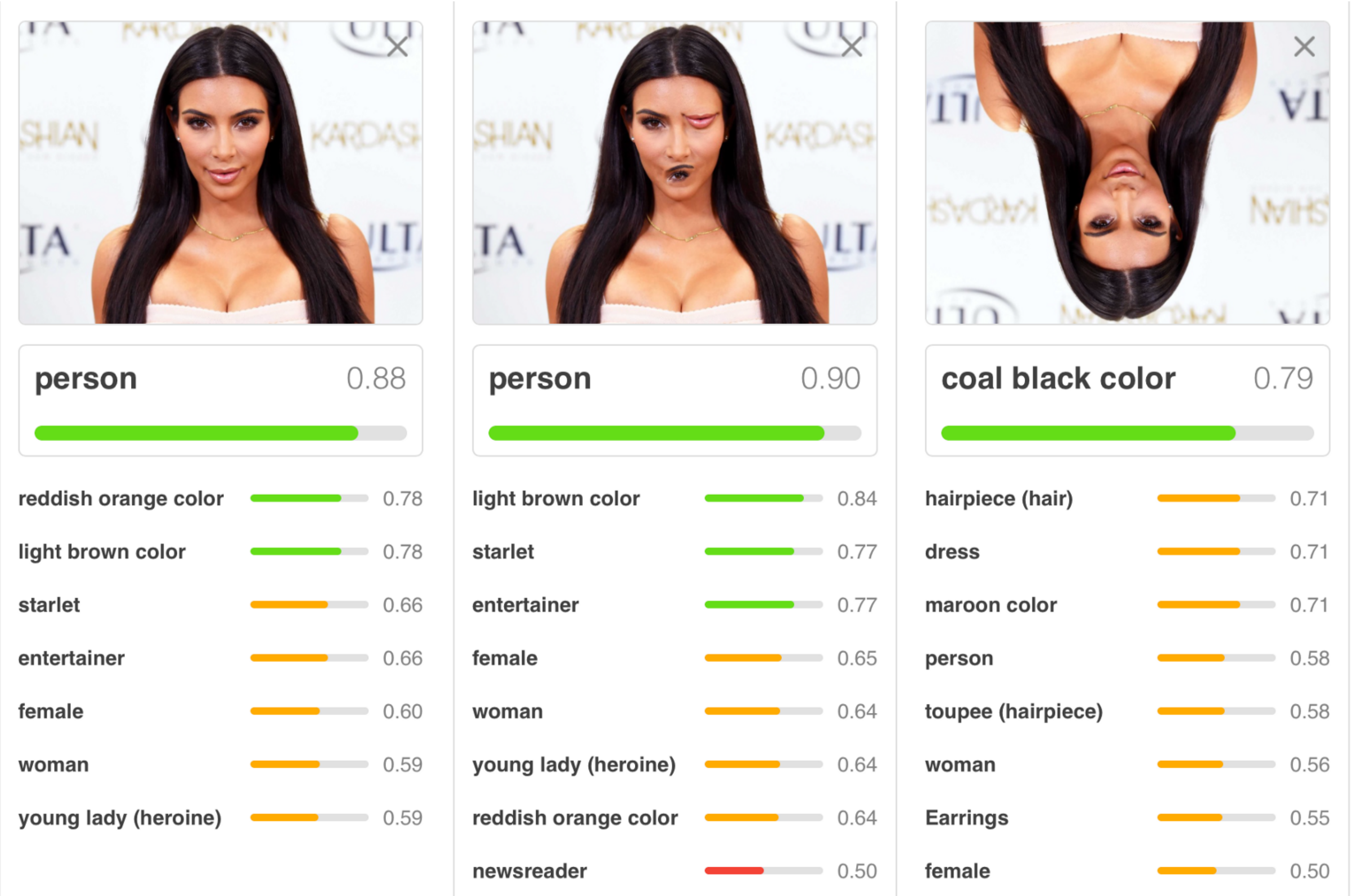


# CapsuleNET

Going **back to inverse** graphics



[ Sabour, Frosst, Hinton, NIPS 2017 ]



\*image credit [medium.com](https://medium.com)



# Neural Modular Networks

