



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

## Lecture 4: Introduction to Computer Vision

# Course **Logistics**

- **Assignment 1** was due 11:59pm today
- **Assignment 2** will be out today (on CNNs) and is due Thursday next week  
(note, it will take computation time)

# Computer vs. human vision



**Human** Vision

# Computer vs. human vision



objects, scenes, people

**Human** Vision

# 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

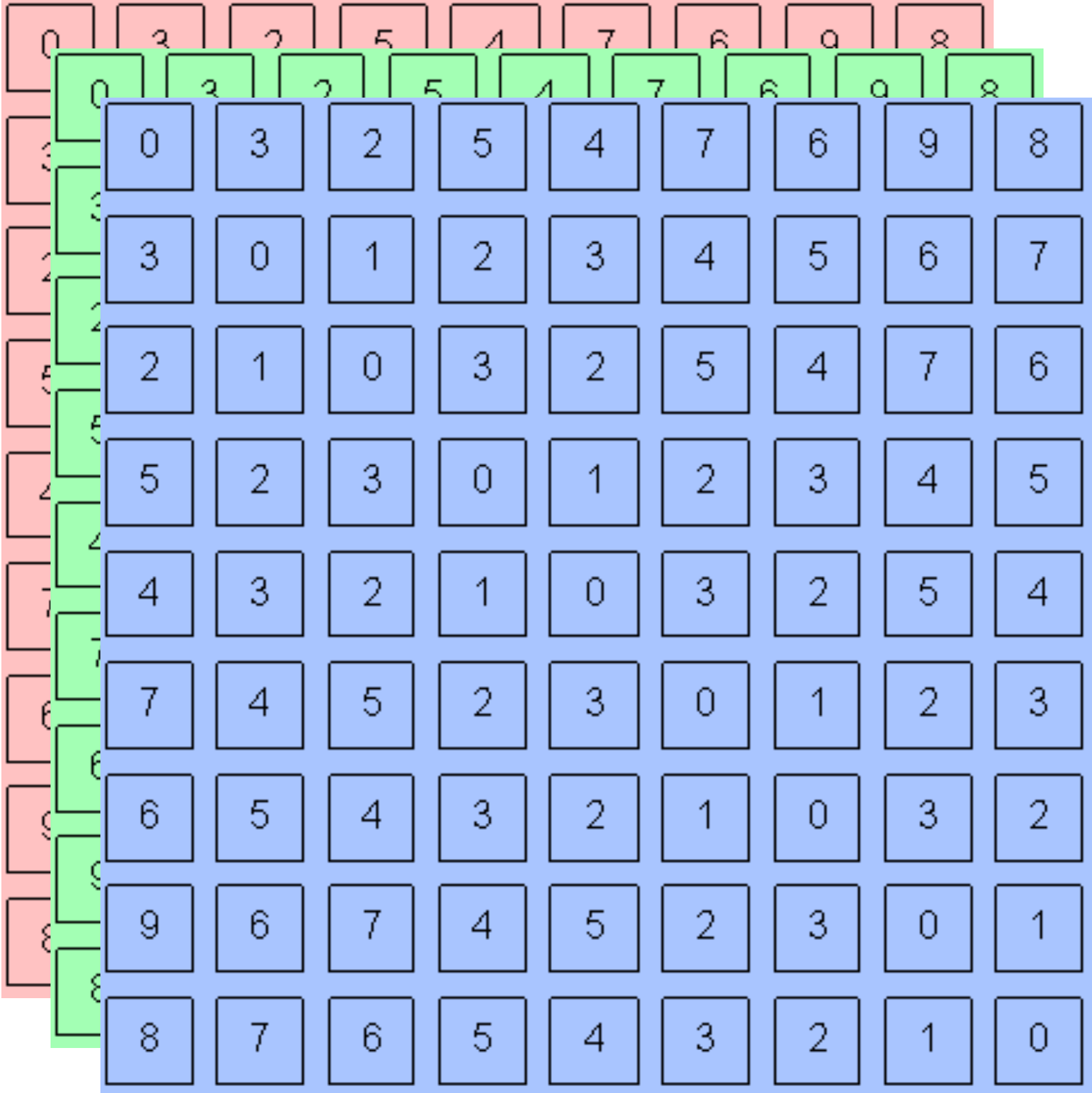
\*slide from V. Ordonex

# Computer vs. human vision



objects, scenes, people

**Human** Vision



tensor of numbers

**Computer** Vision

\*slide from V. Ordonex

# 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**

\*courtesy 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



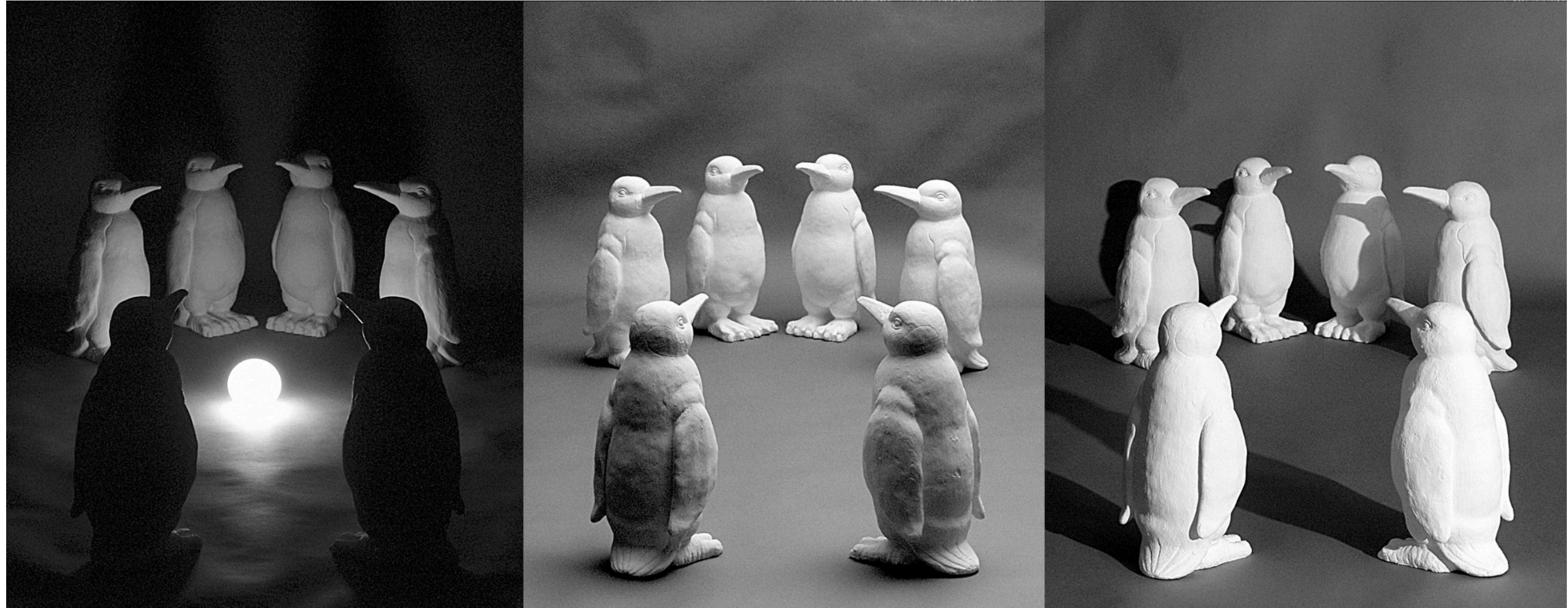


# Challenges: Viewpoint invariance



**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



# Challenges: Local ambiguity and context



\*image credit Fergus & Torralba



# Challenges: Motion

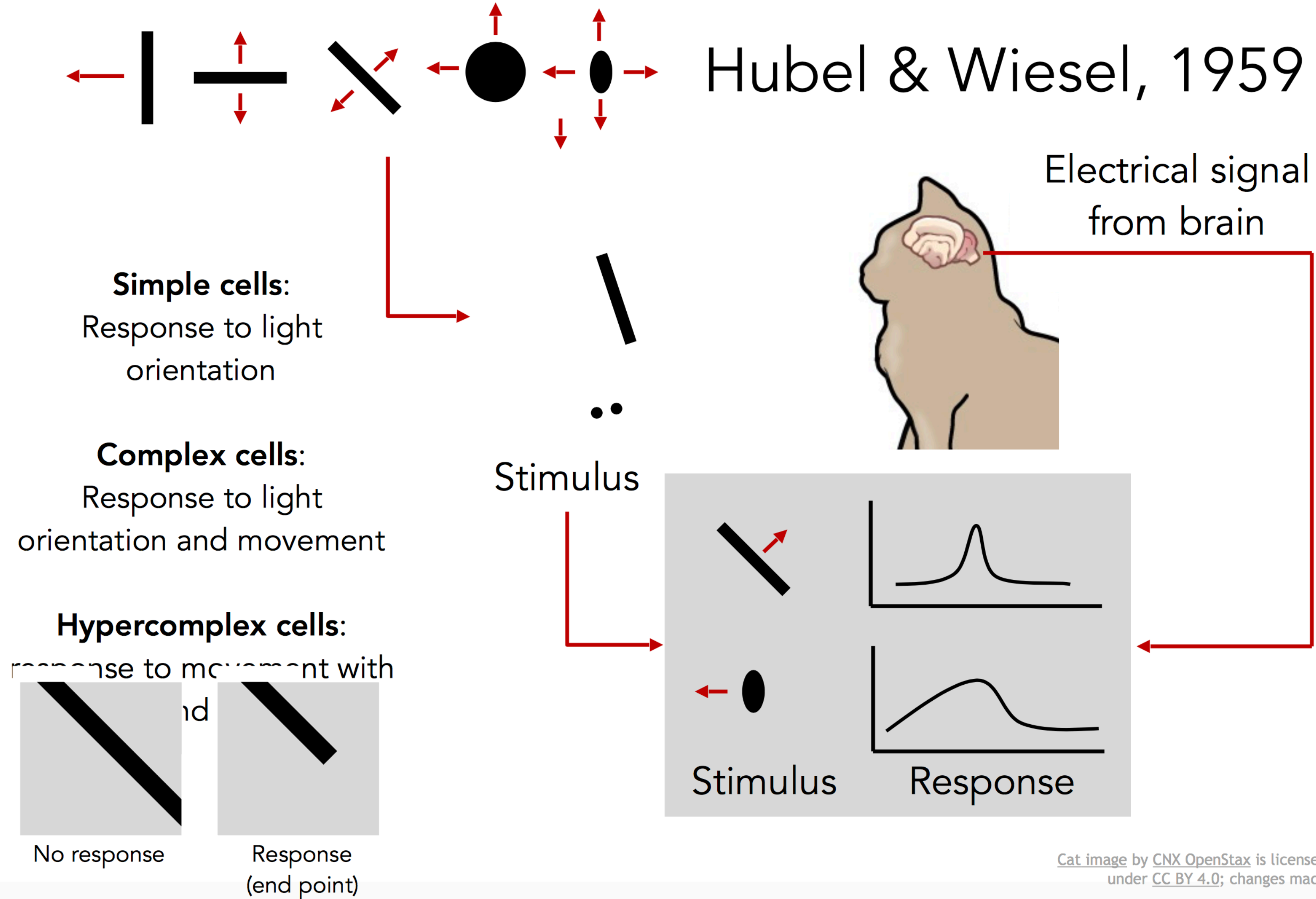


\*image credit Peter Meer

# Challenges: Object inter-class variation



# Human vision ...



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

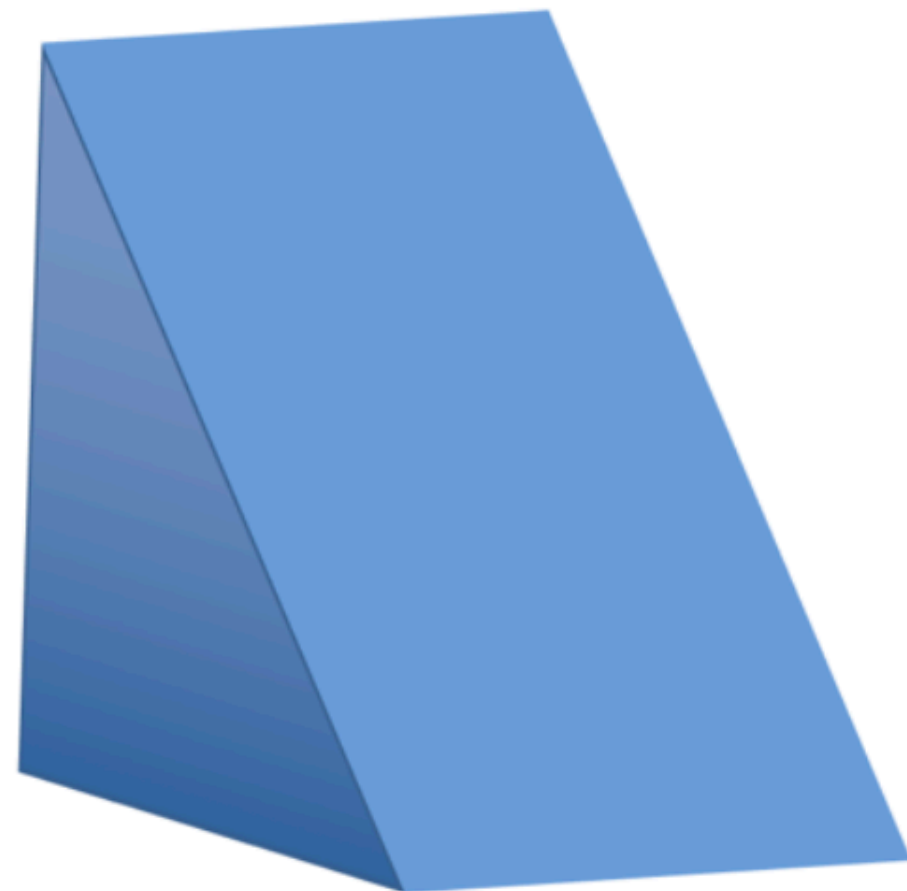
# Computer vision ... the beginning ...



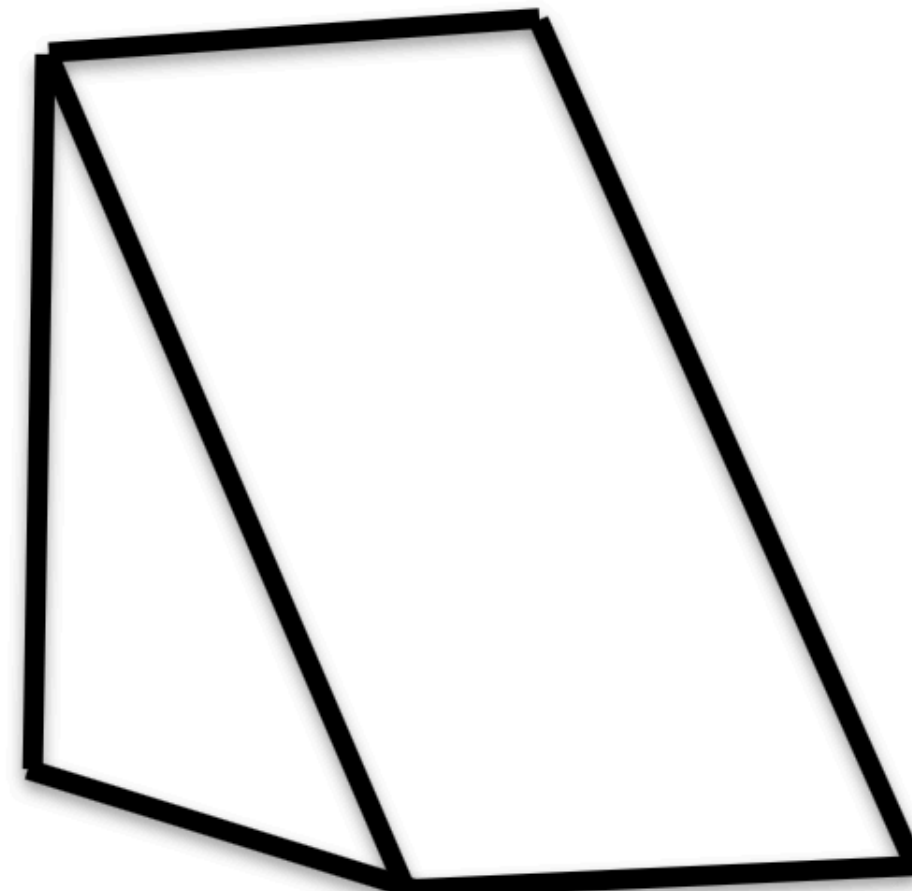
Larry Roberts

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

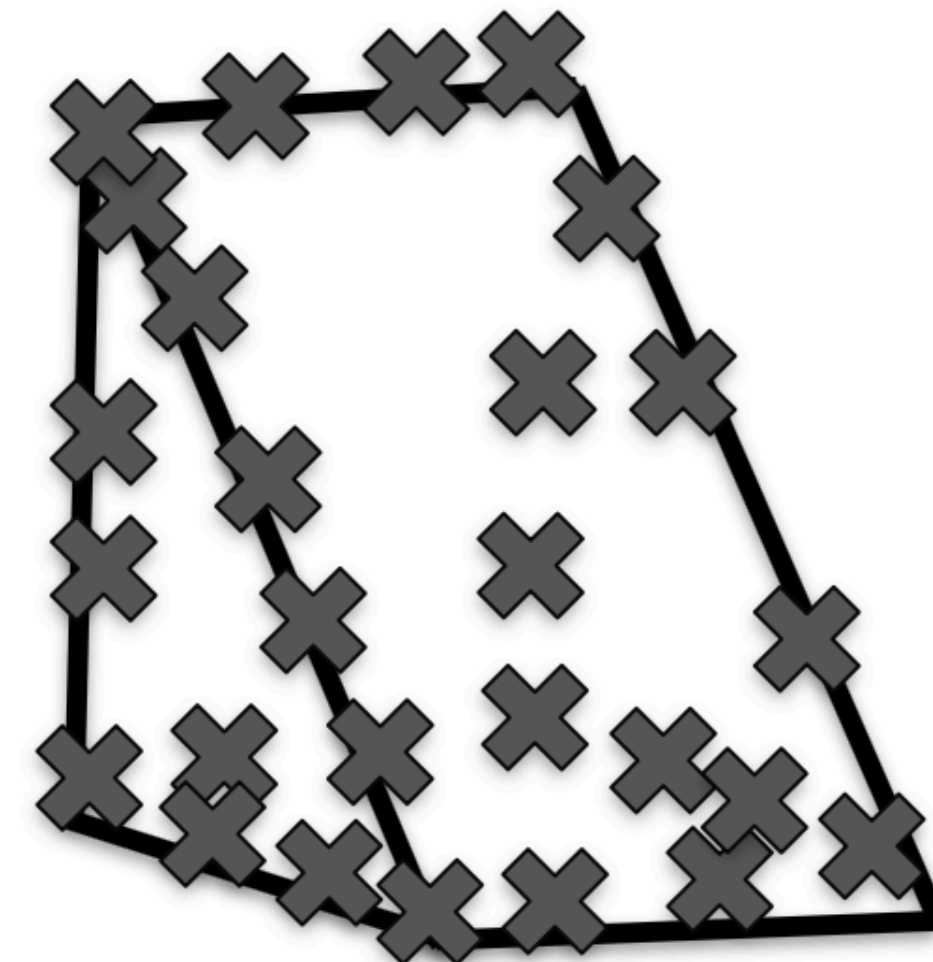
"the perception of **solid objects** is a process which can be based on the **properties of three-dimensional** transformations and **the laws of nature**"



(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**"

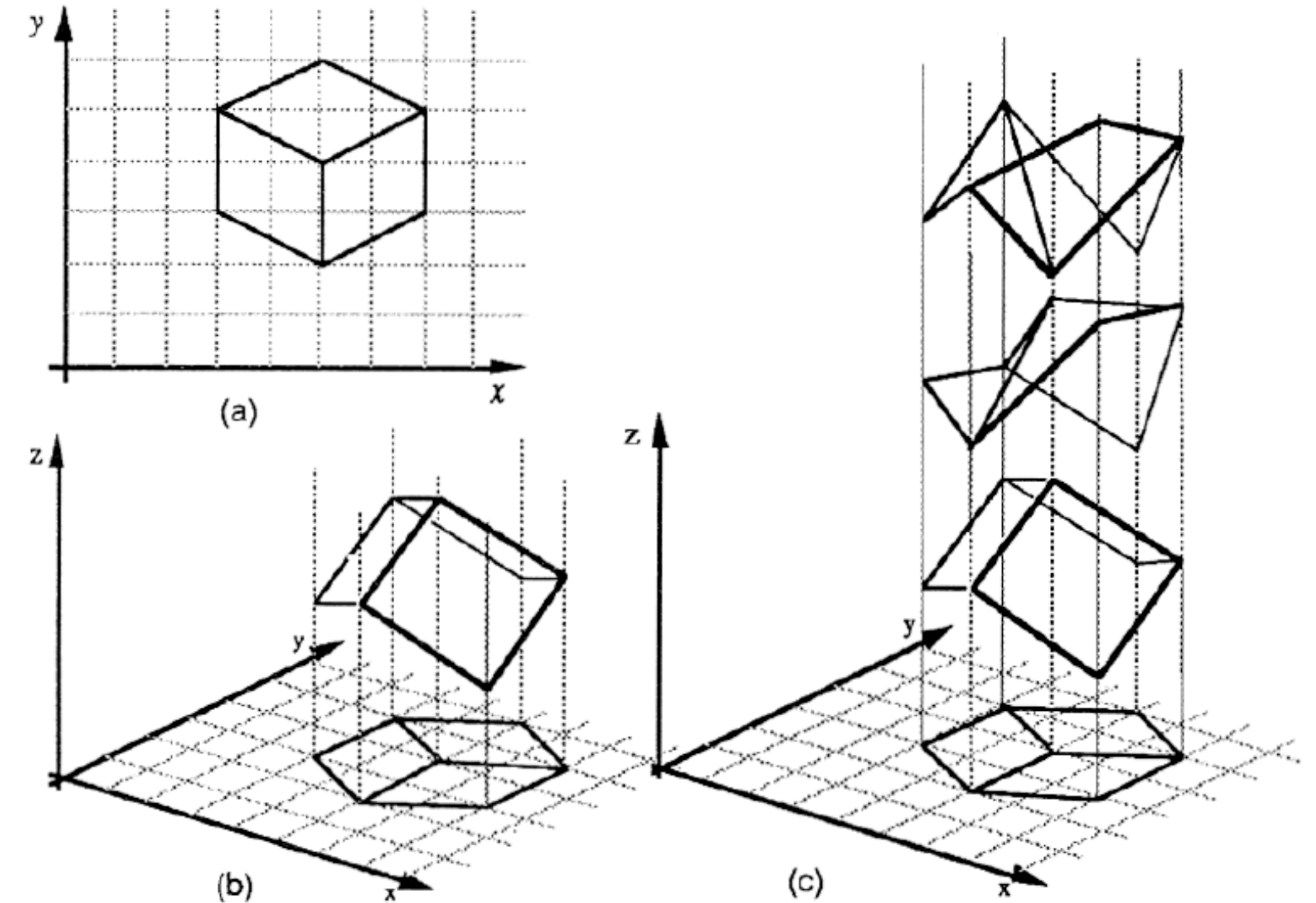


Figure 1. (a) A line drawing provides information only about the  $x, y$  coordinates of points lying along the object contours. (b) The human visual system is usually able to reconstruct an object in three dimensions given only a single 2D projection (c) Any planar line-drawing is geometrically consistent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

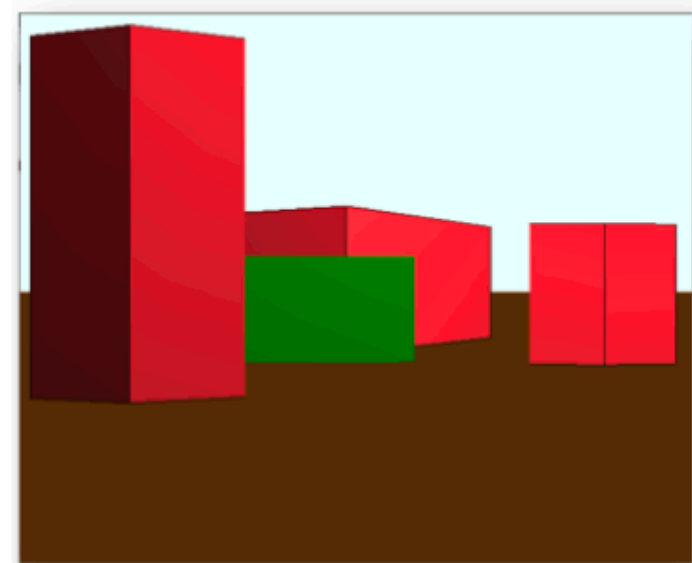
# Computer vision ... the beginning ...



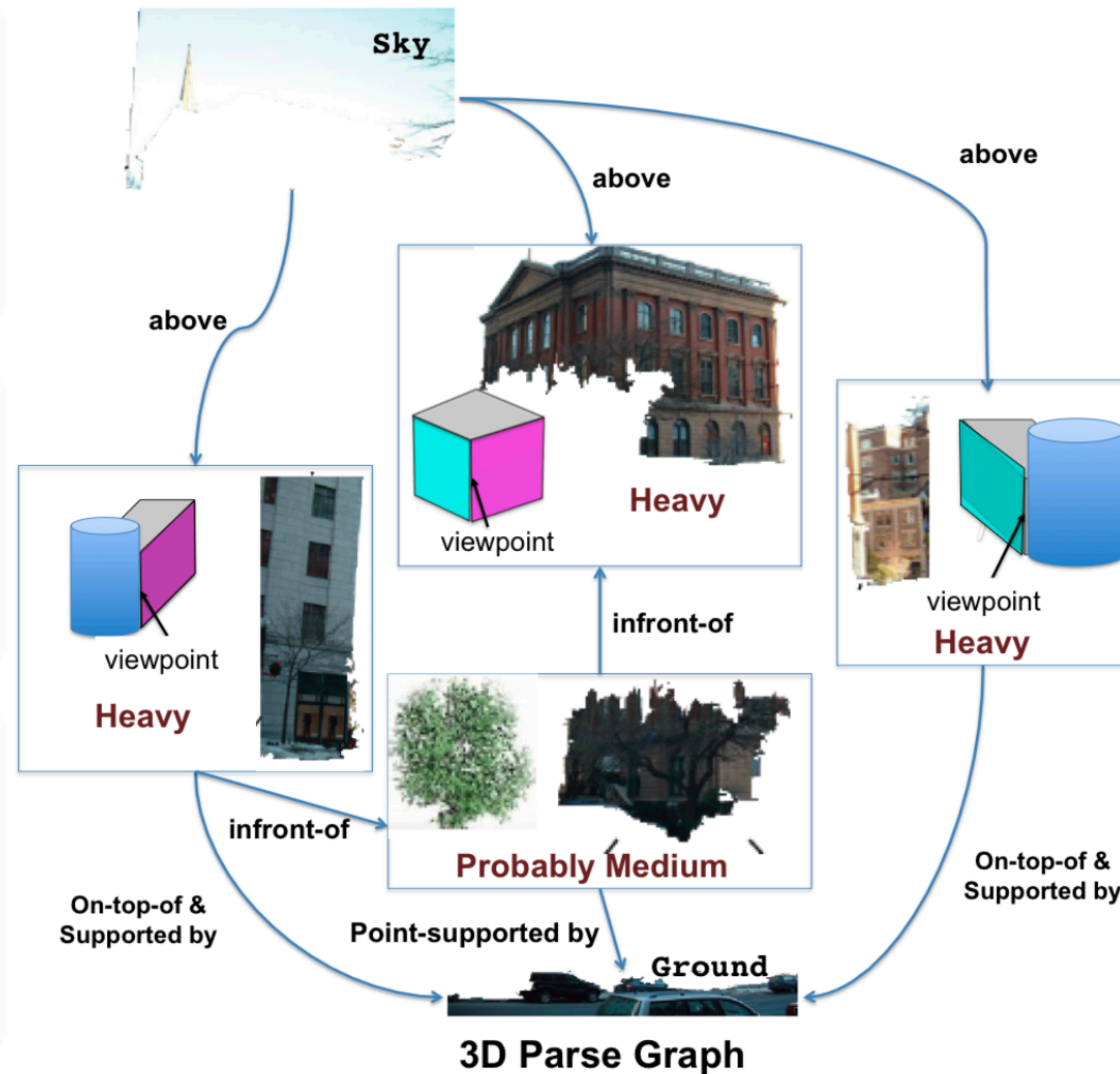
Input Image



Blocks World



3D Rendering

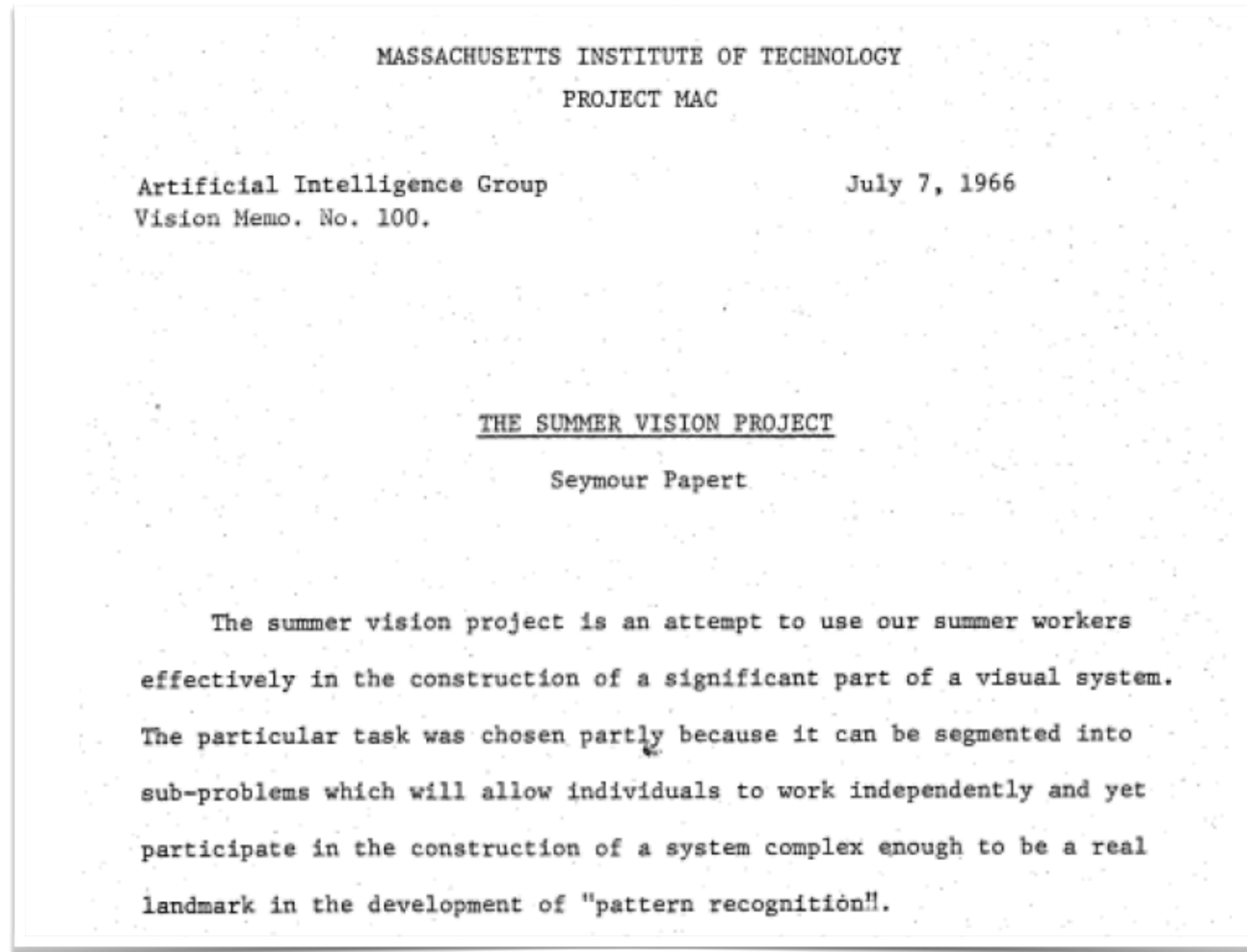


**Static Equilibrium:** Forces and torques acting on a block should cancel each other out.

**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

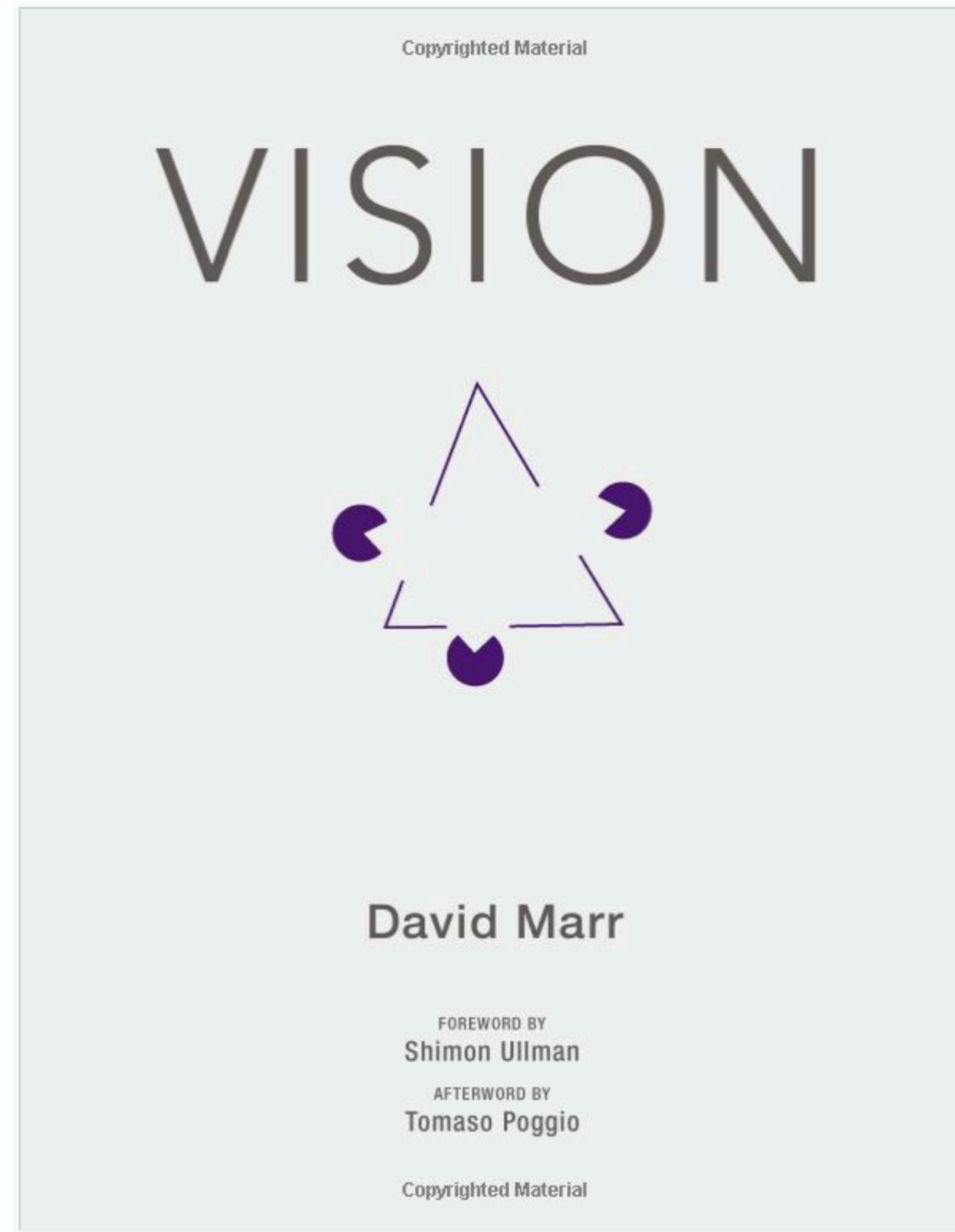
# Computer vision ... the beginning ...



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”

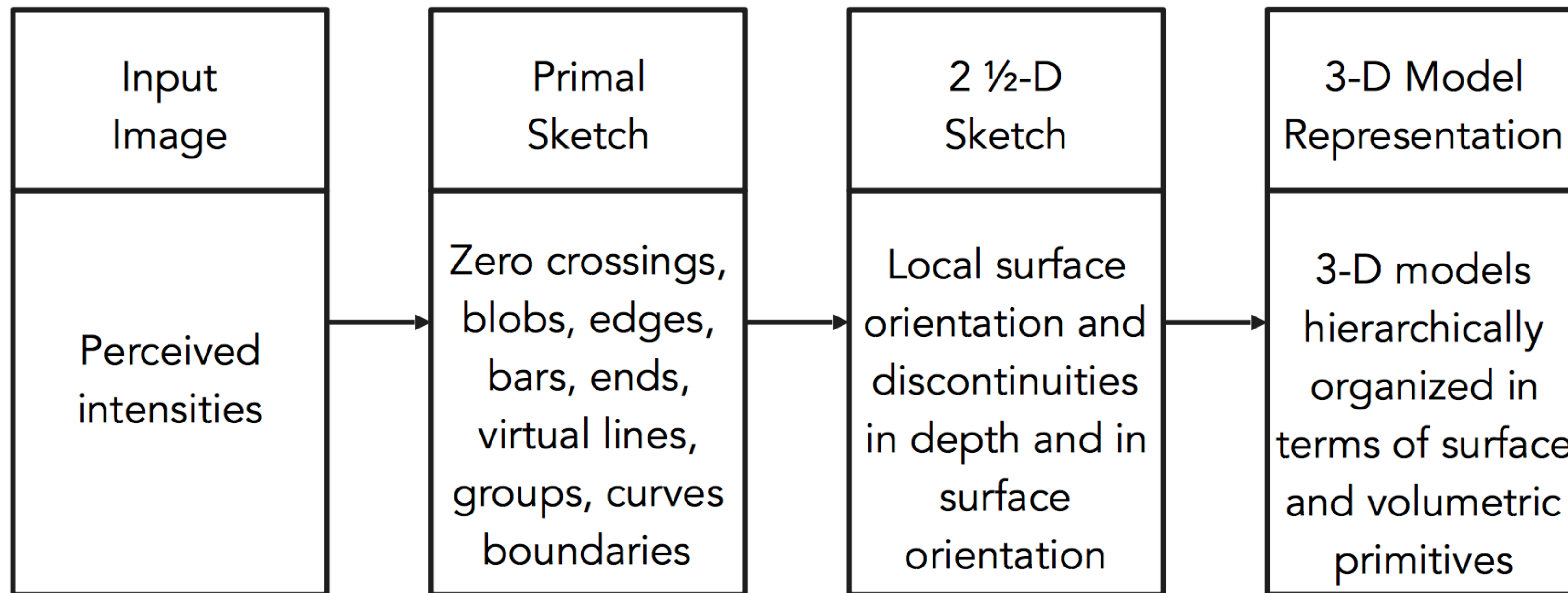
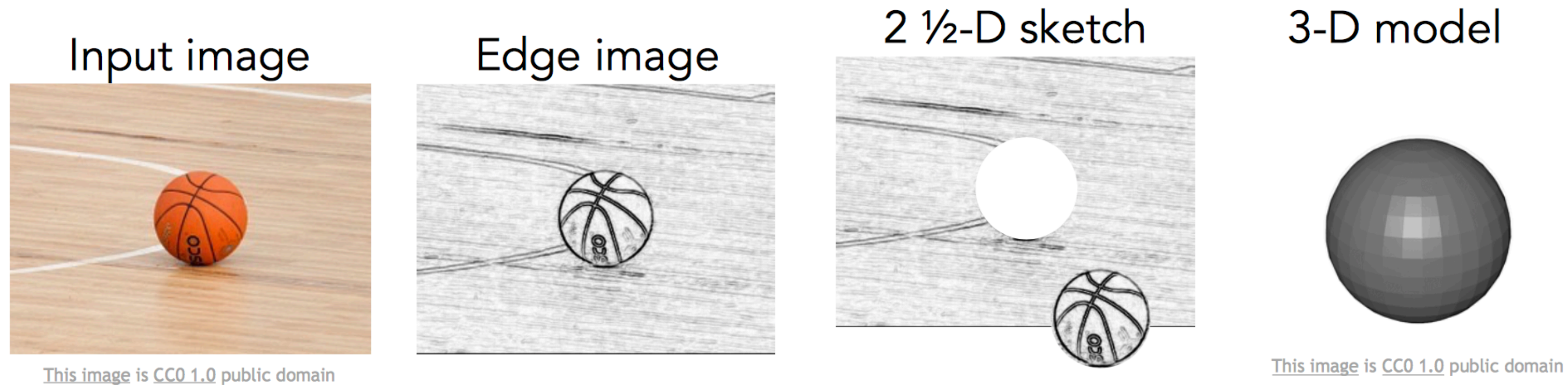
[ Szeliski 2009, Computer Vision ]

# David Marr, 1970s



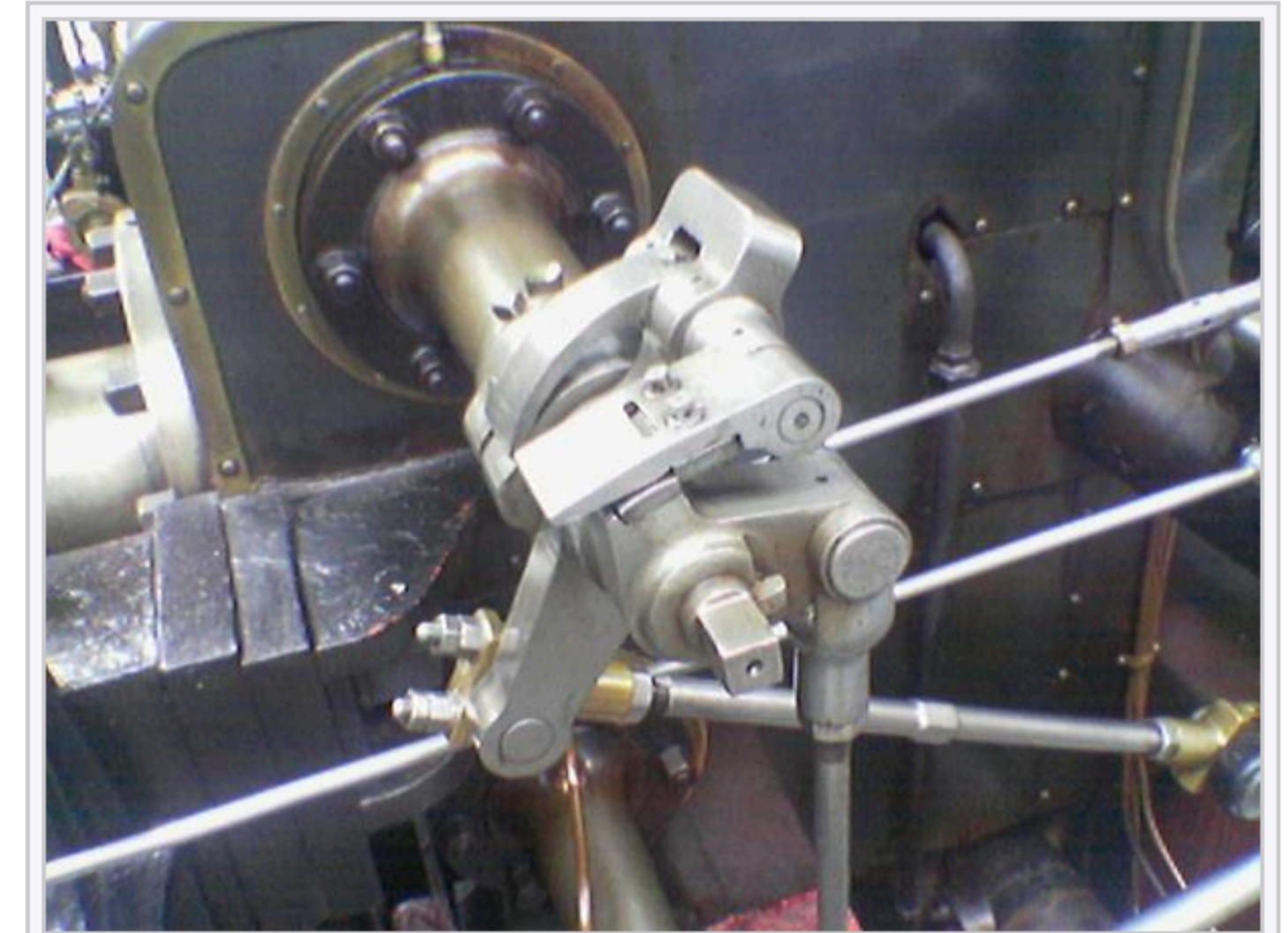
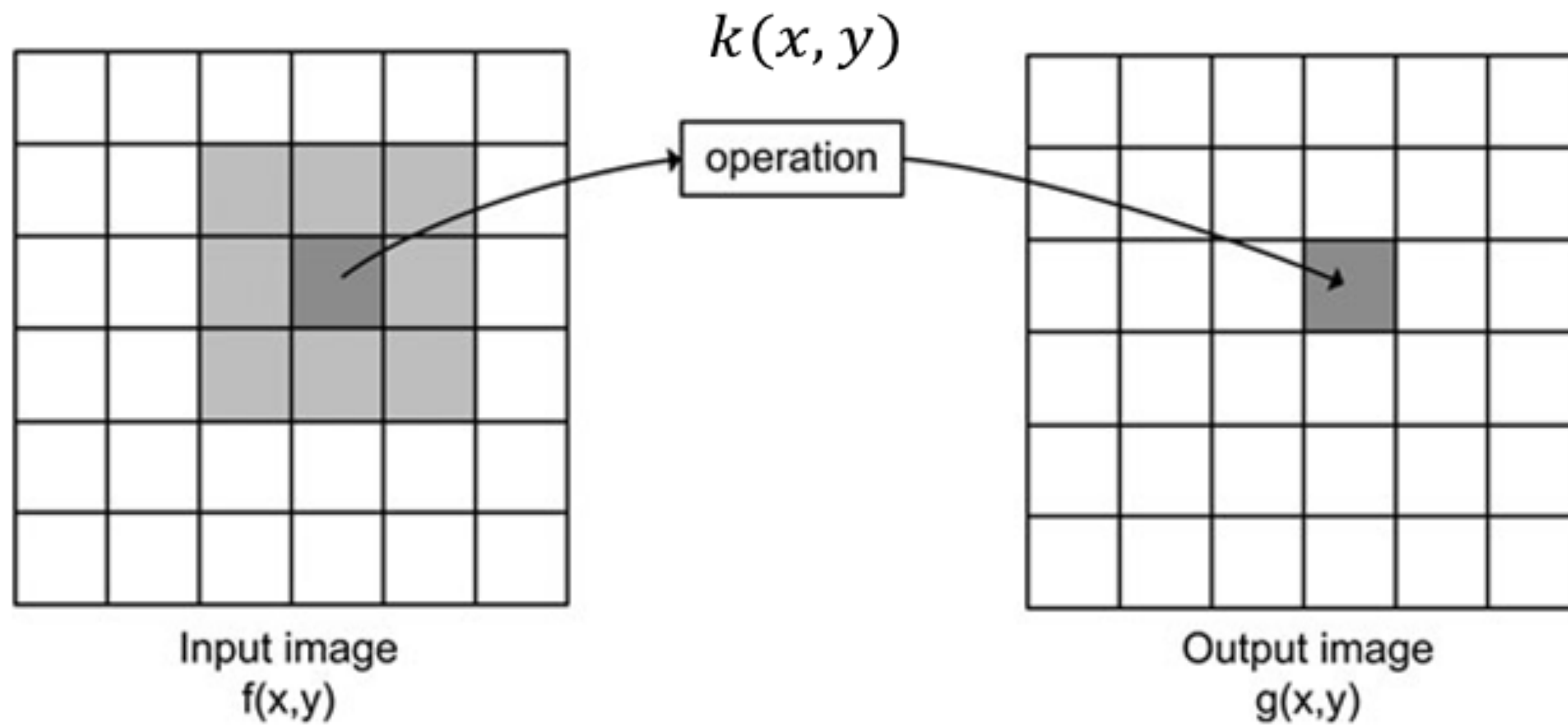


# David Marr, 1970s

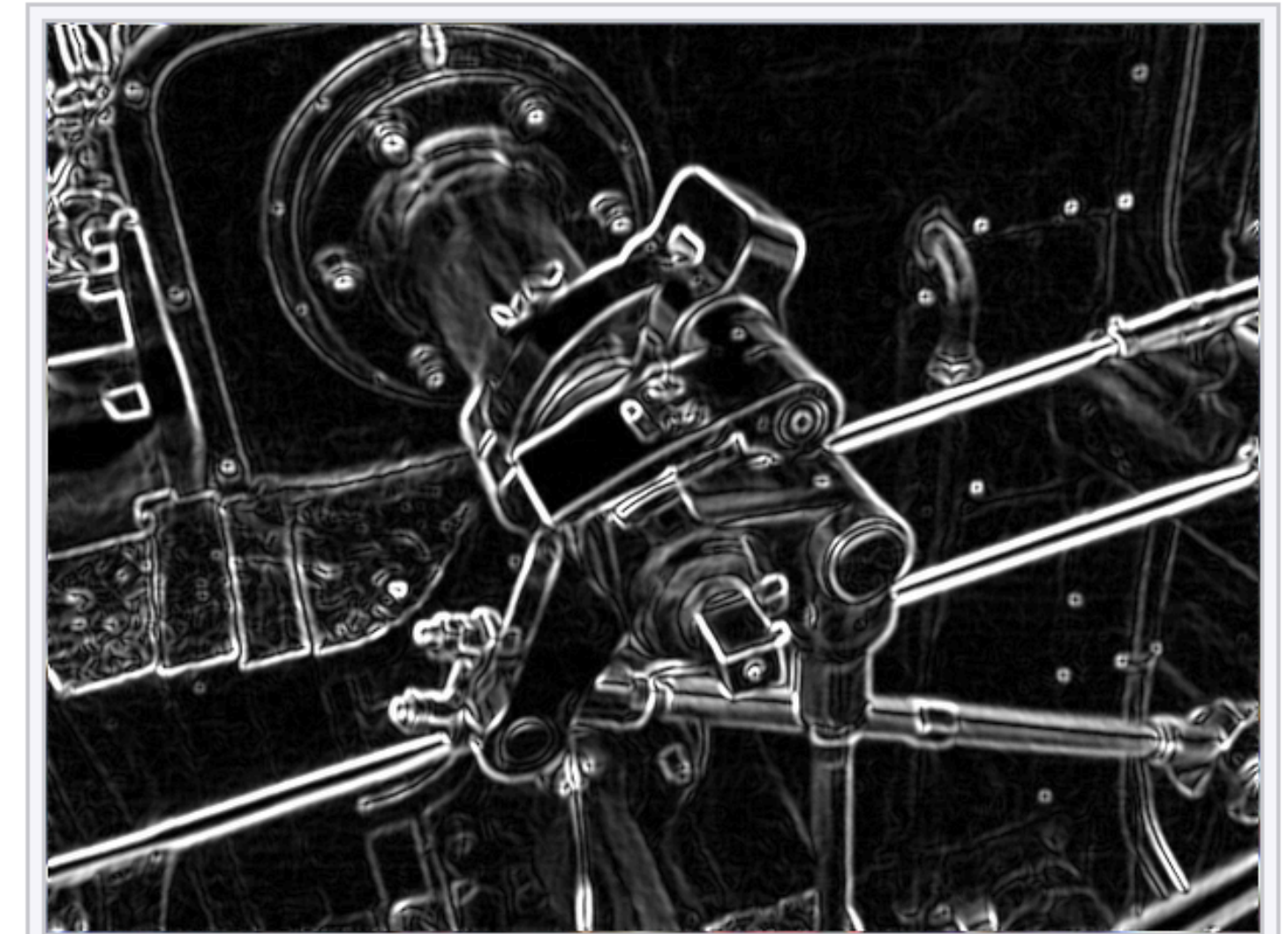


# Edges

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



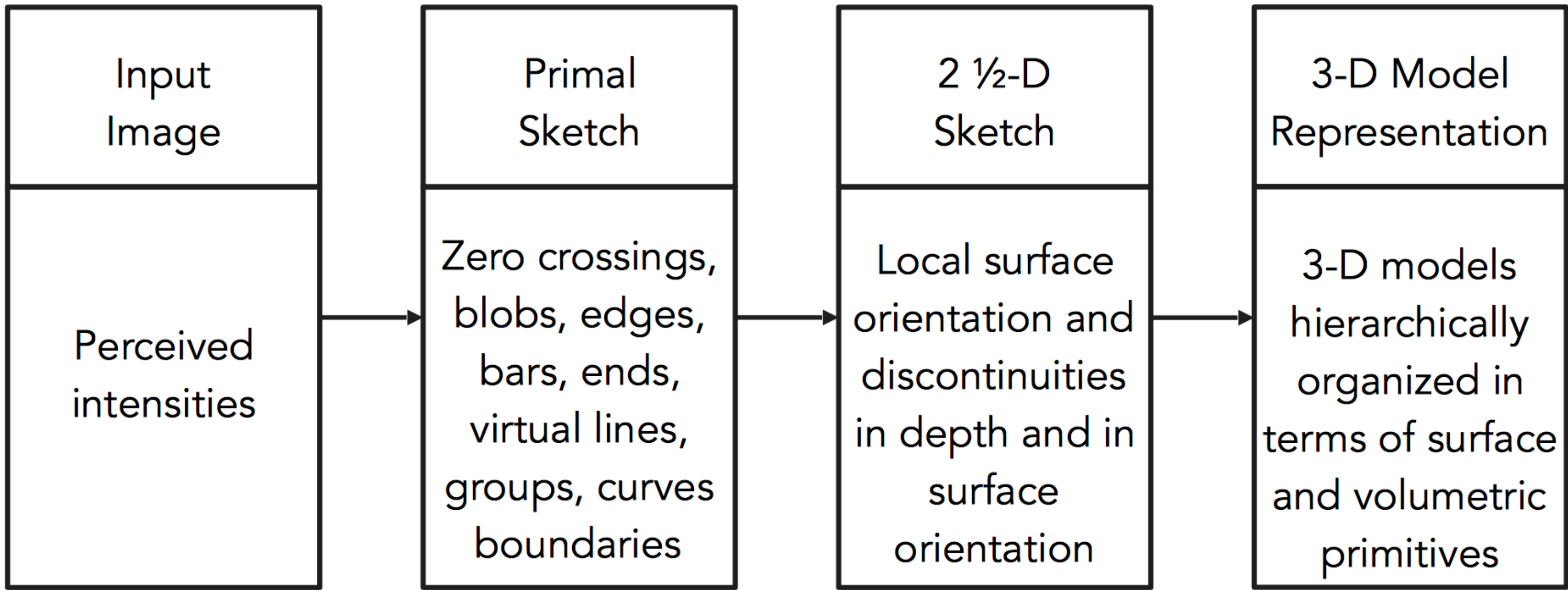
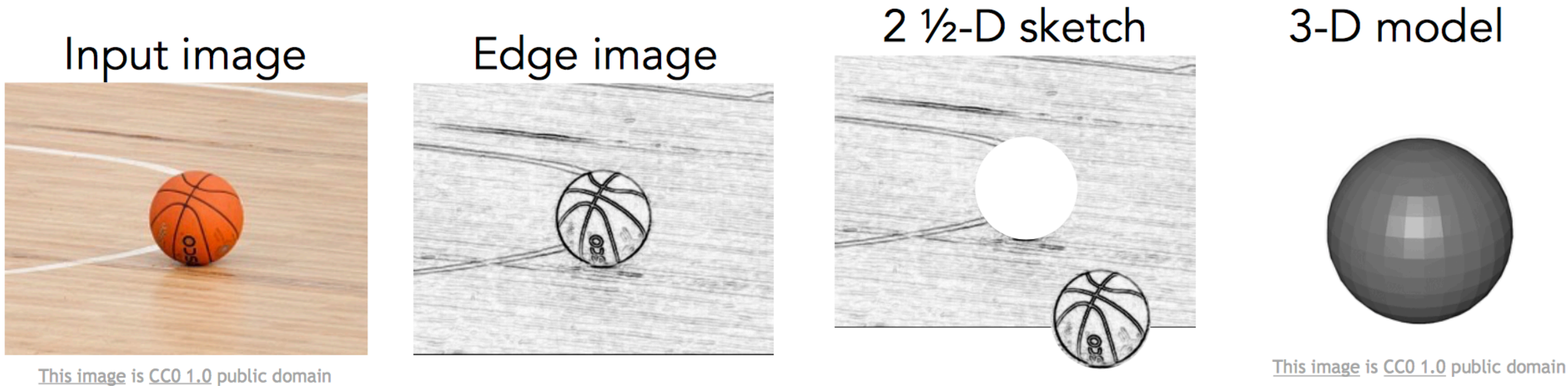
A color picture of a steam engine



The Sobel operator applied to that image



# David Marr, 1970s



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

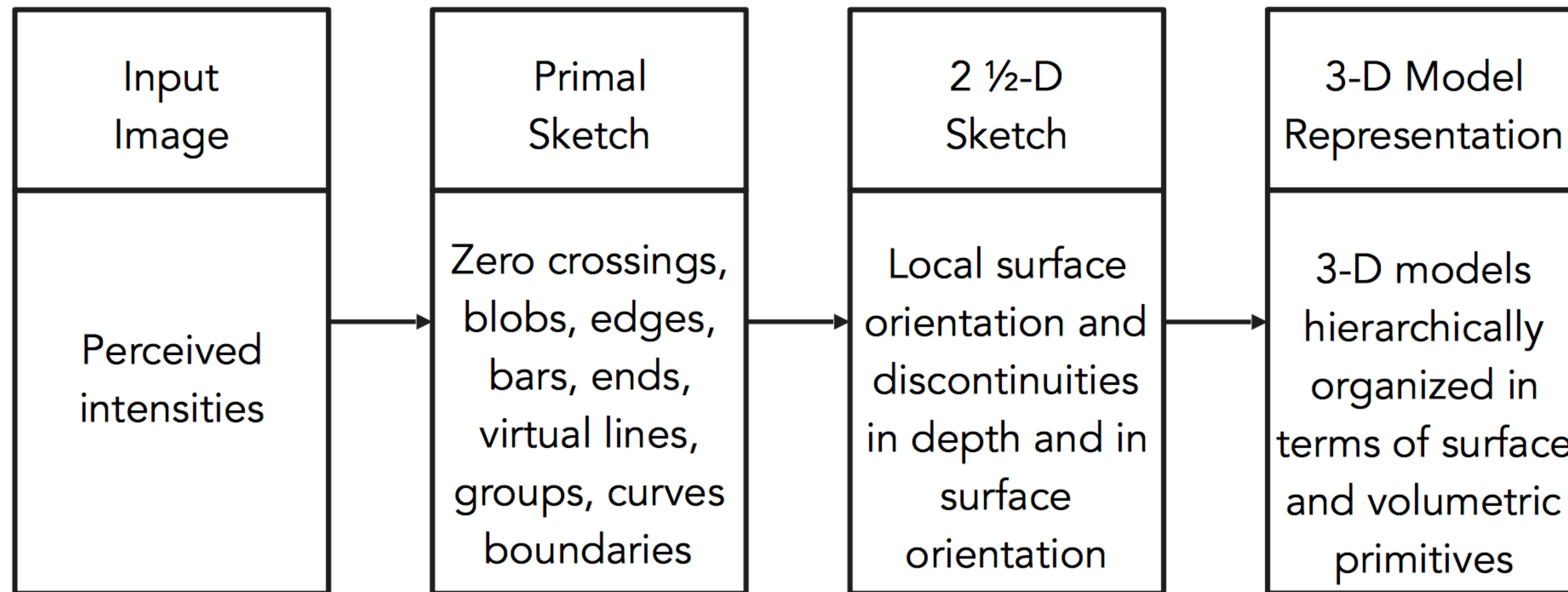
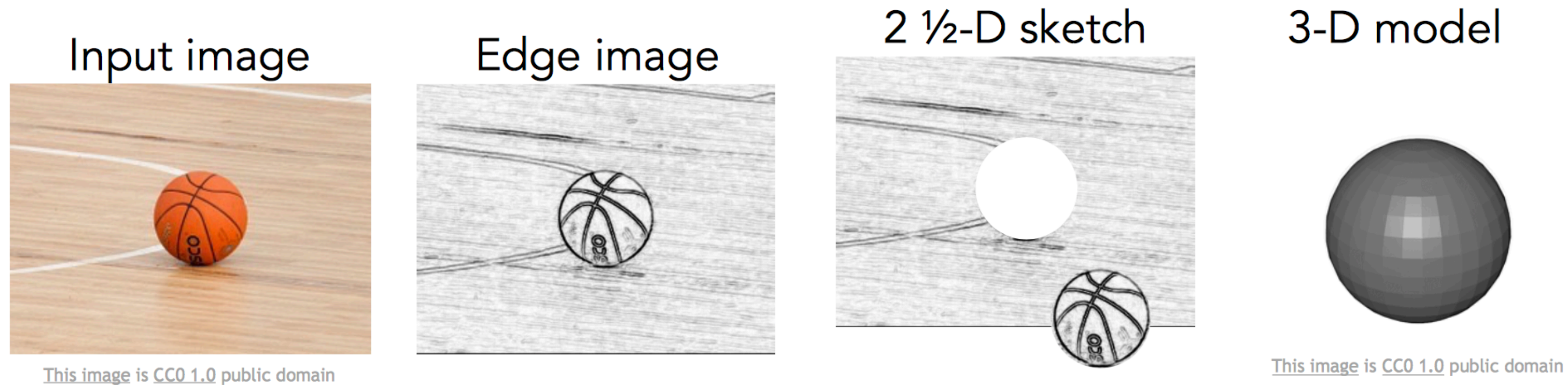
\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# Segmentation - GraphCuts



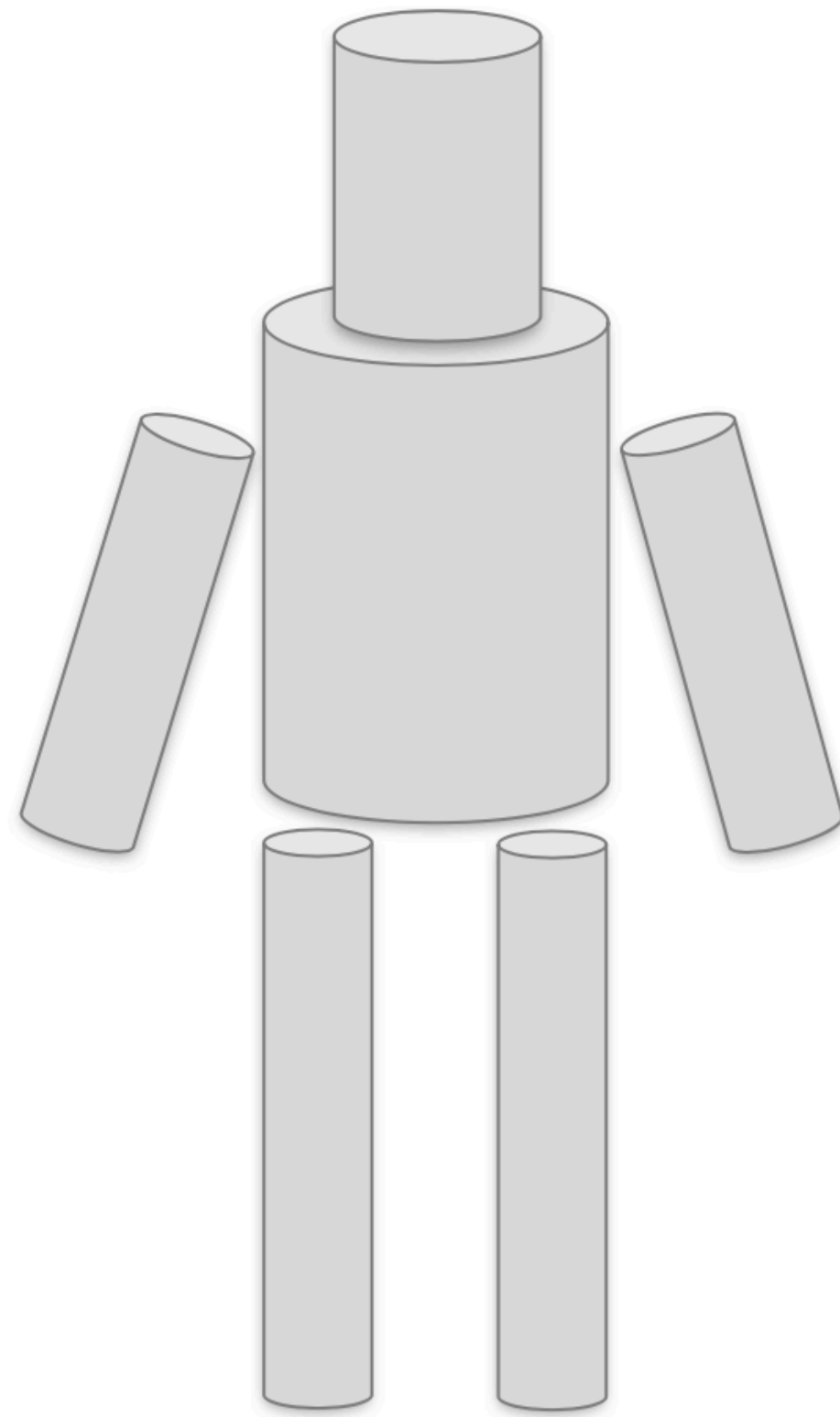
[ Shi & Malik, 2000 ]

# David Marr, 1970s



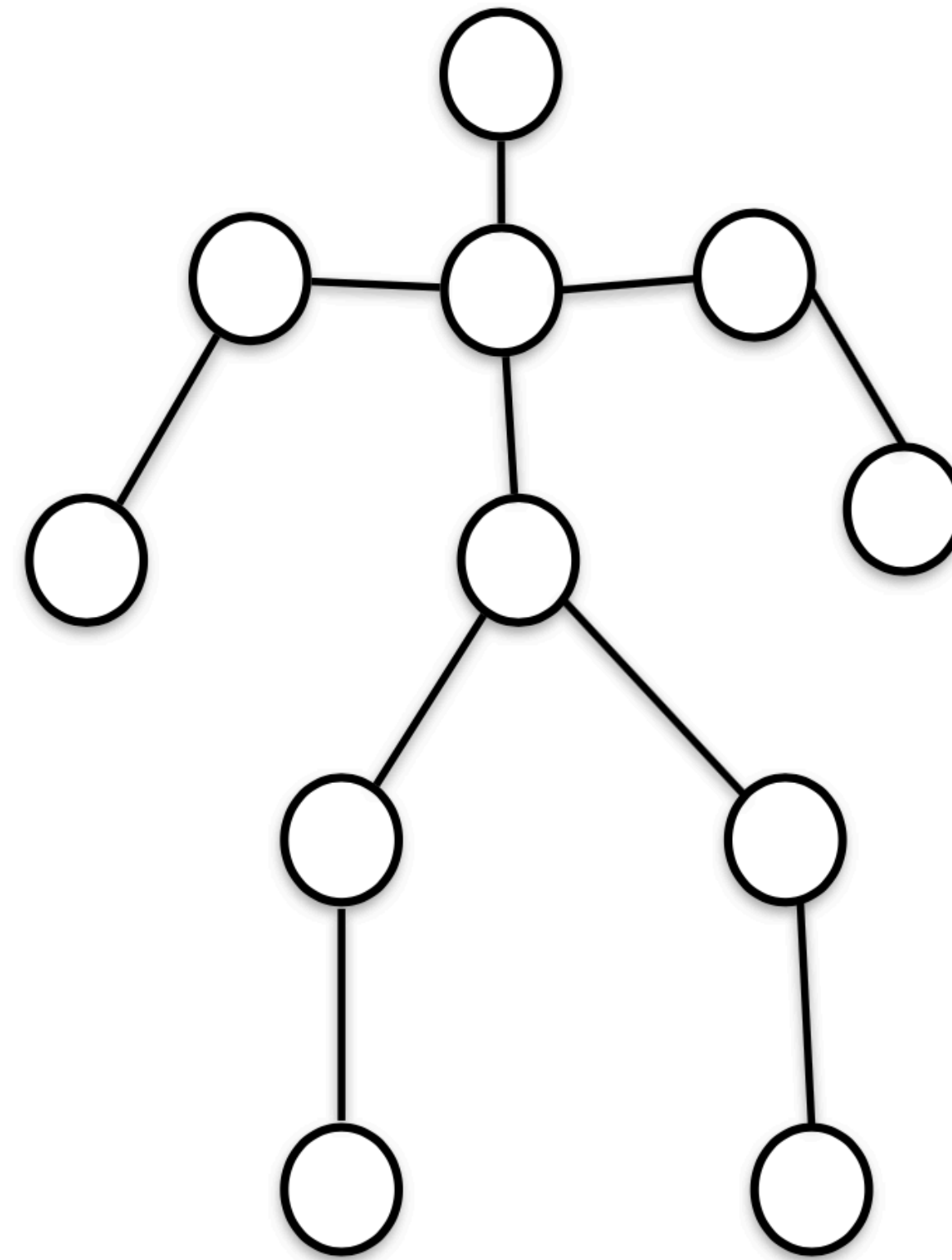
# Part-based Models

## Generalized Cylinders

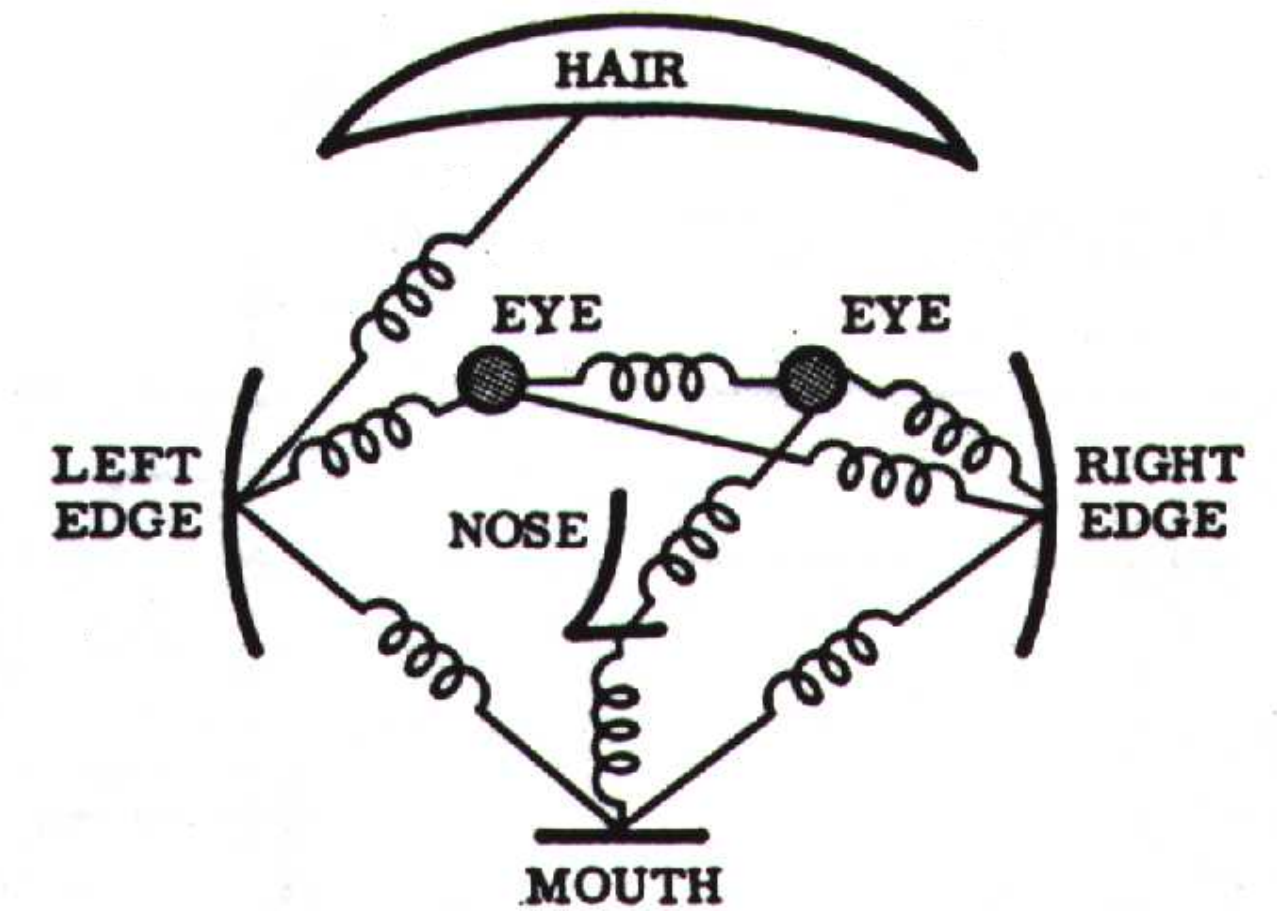


[ Brooks & Binford, 1979 ]

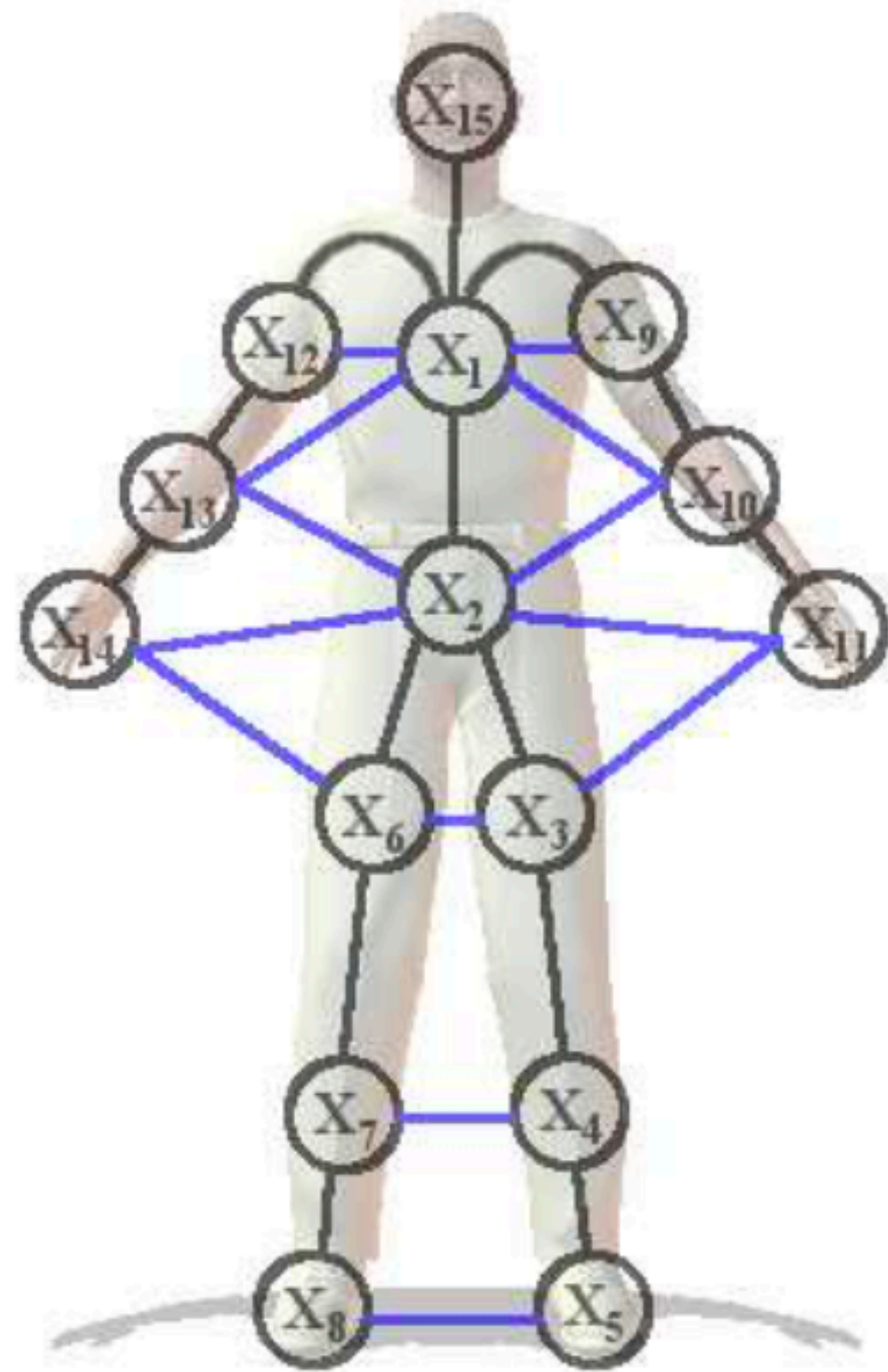
## Pictorial Structures



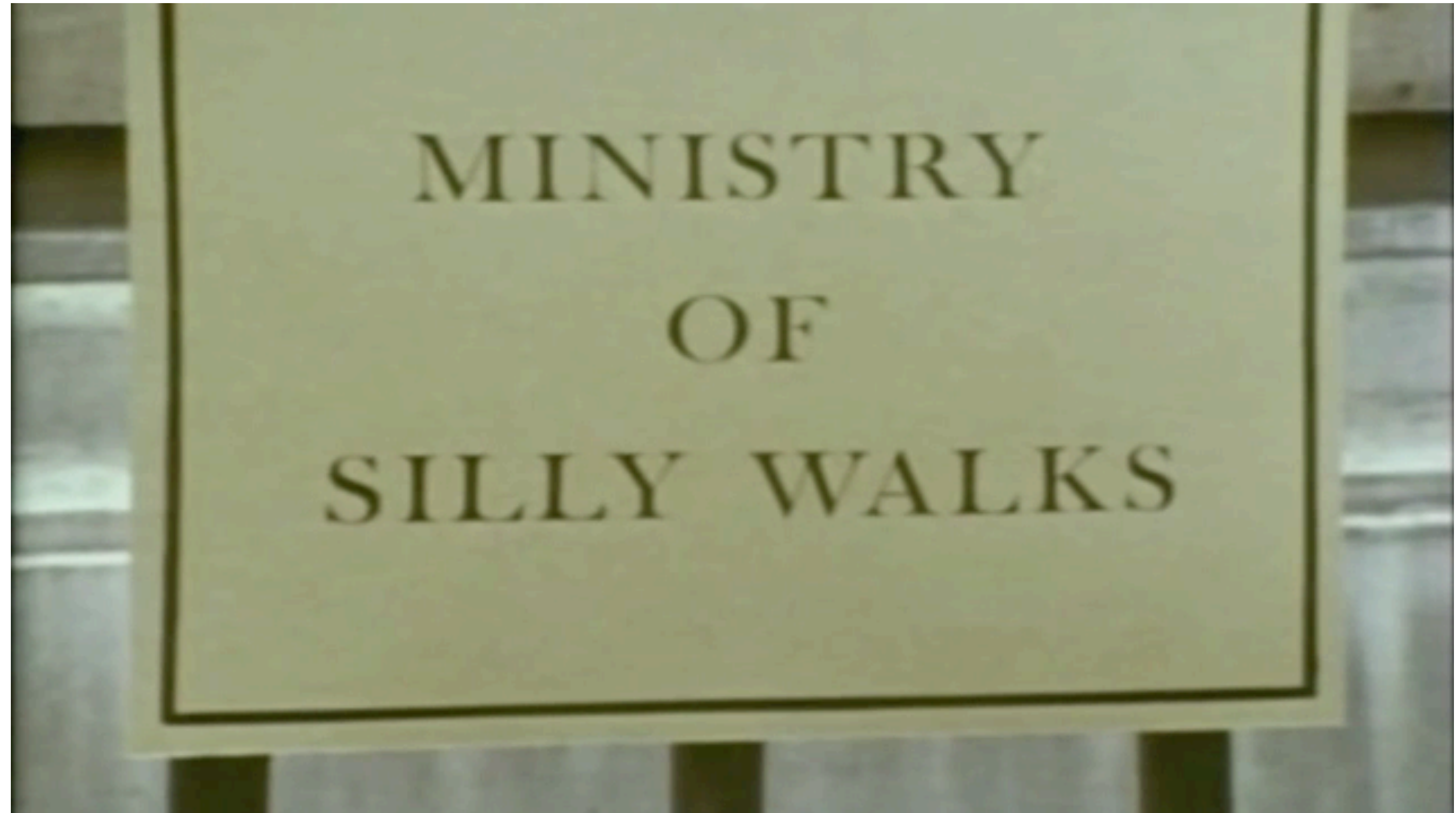
[ Fischler & Elschlager, 1973 ]



# Part-based Models

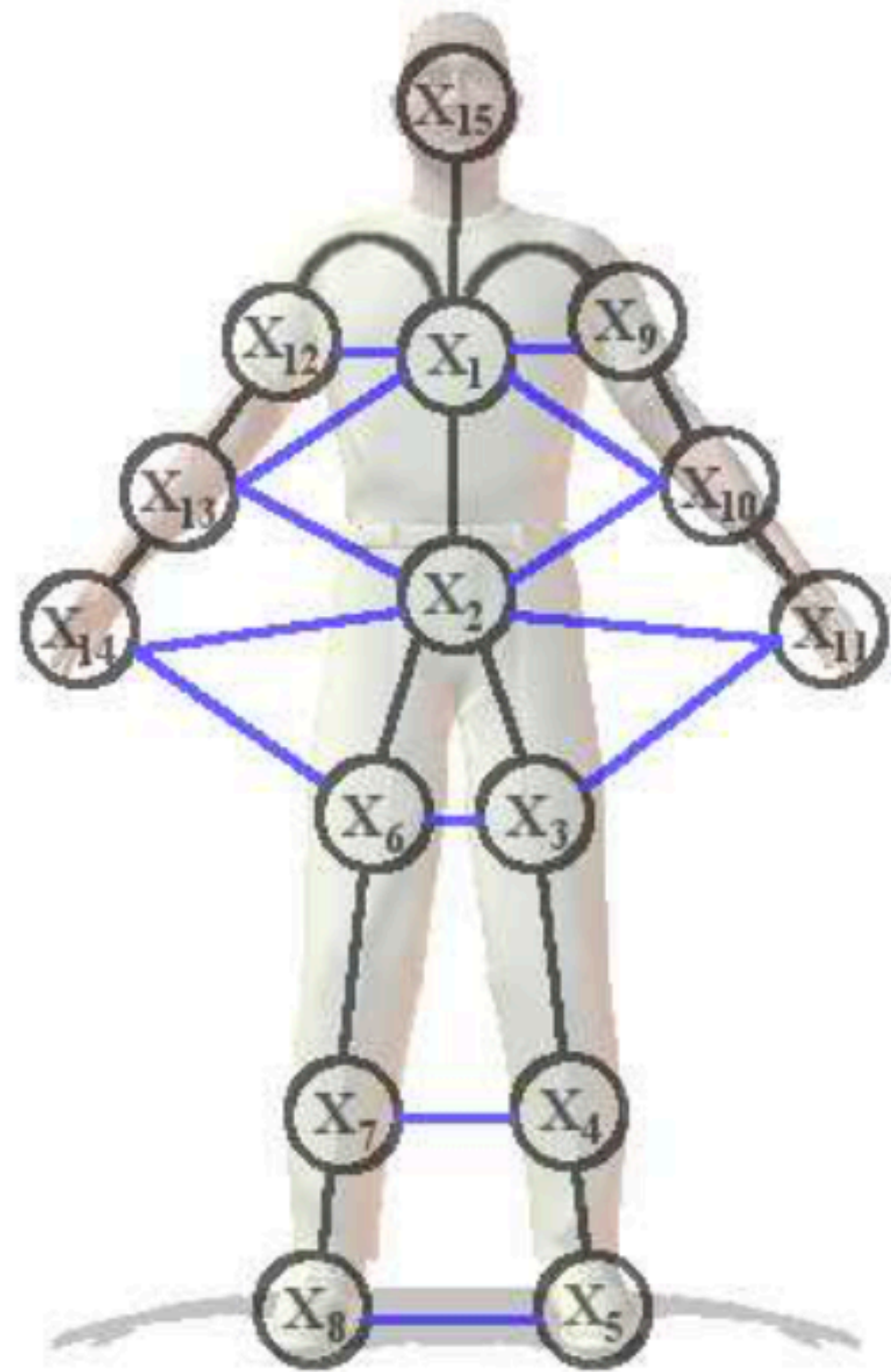


[ Sigal et al. 2004]

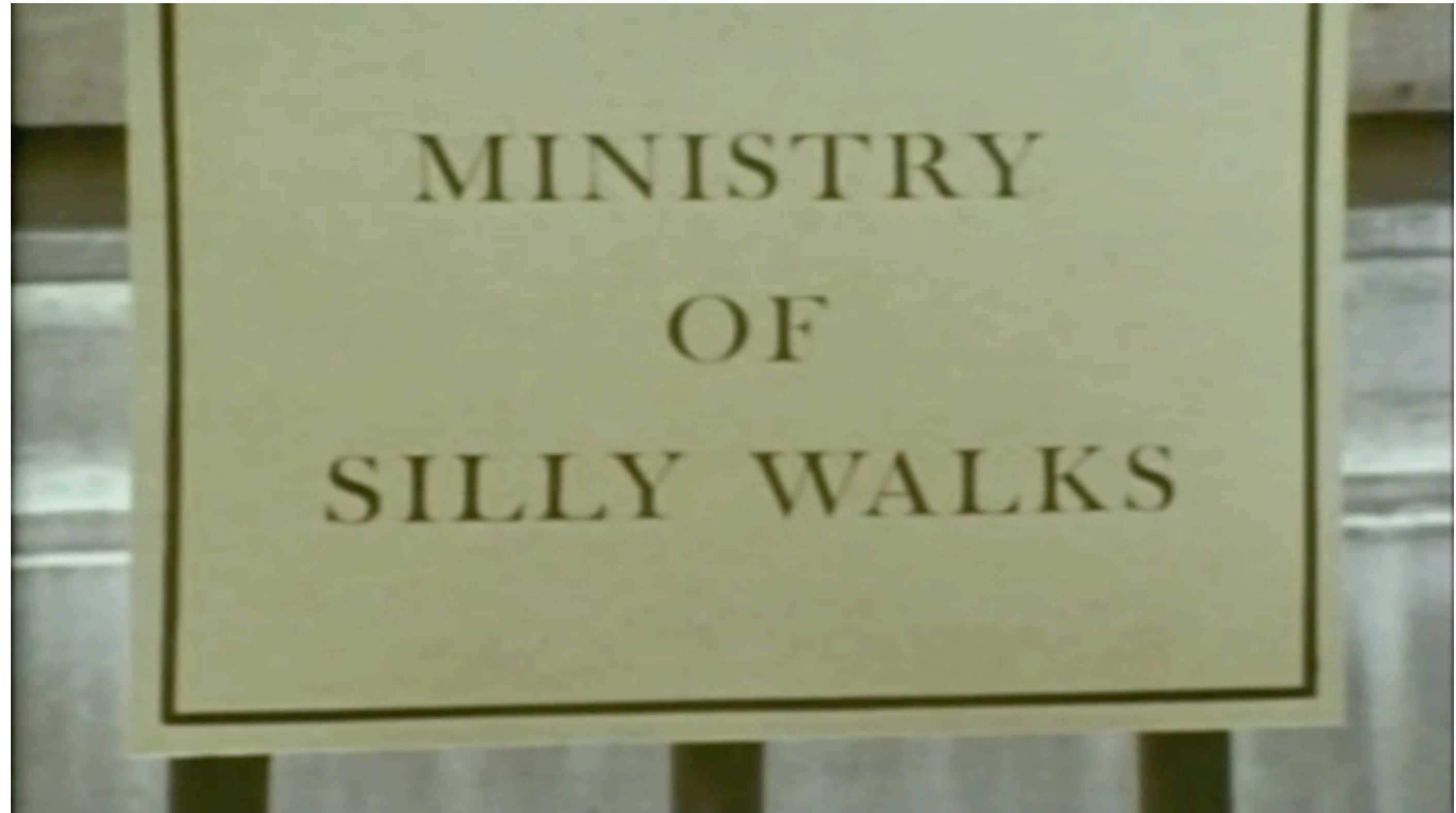


Monty Python's **Ministry of Silly Walks**

# Part-based Models



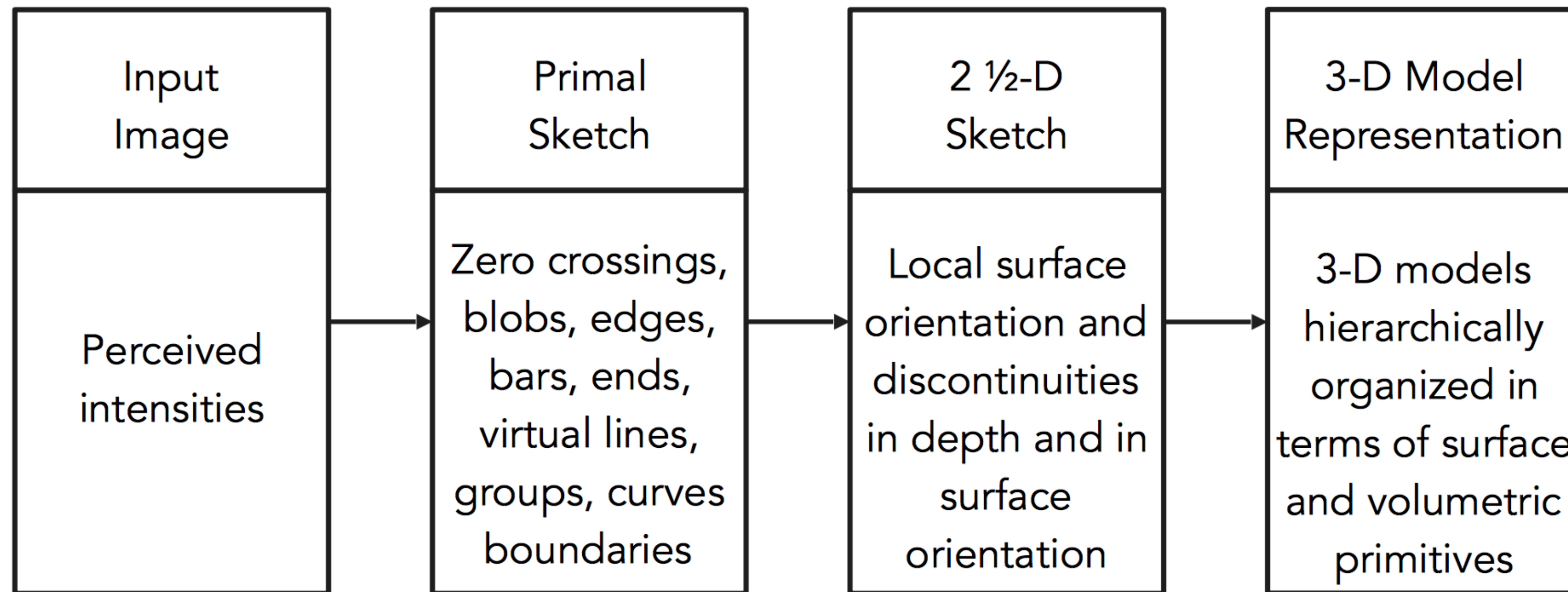
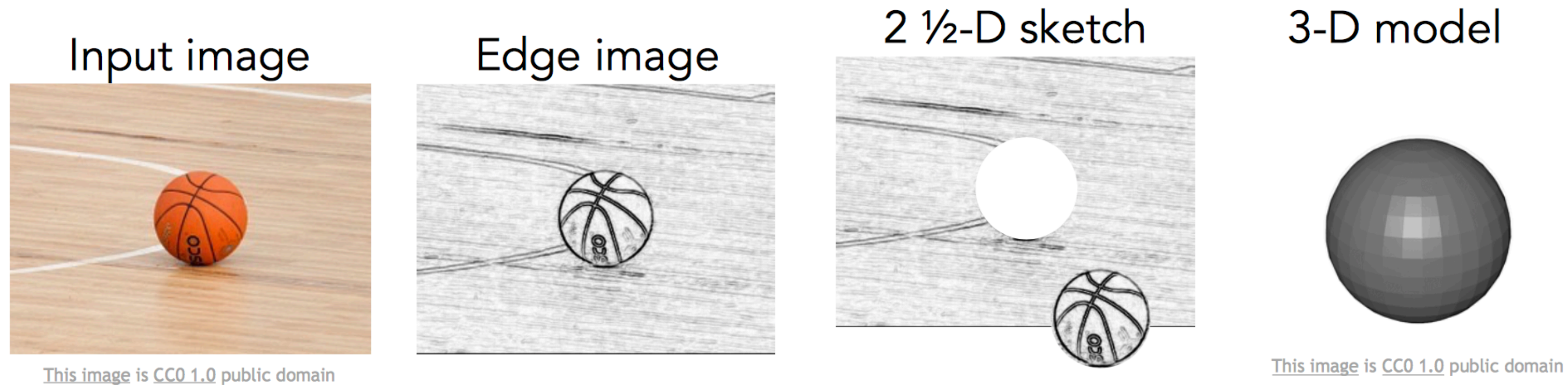
[ Sigal et al. 2004]



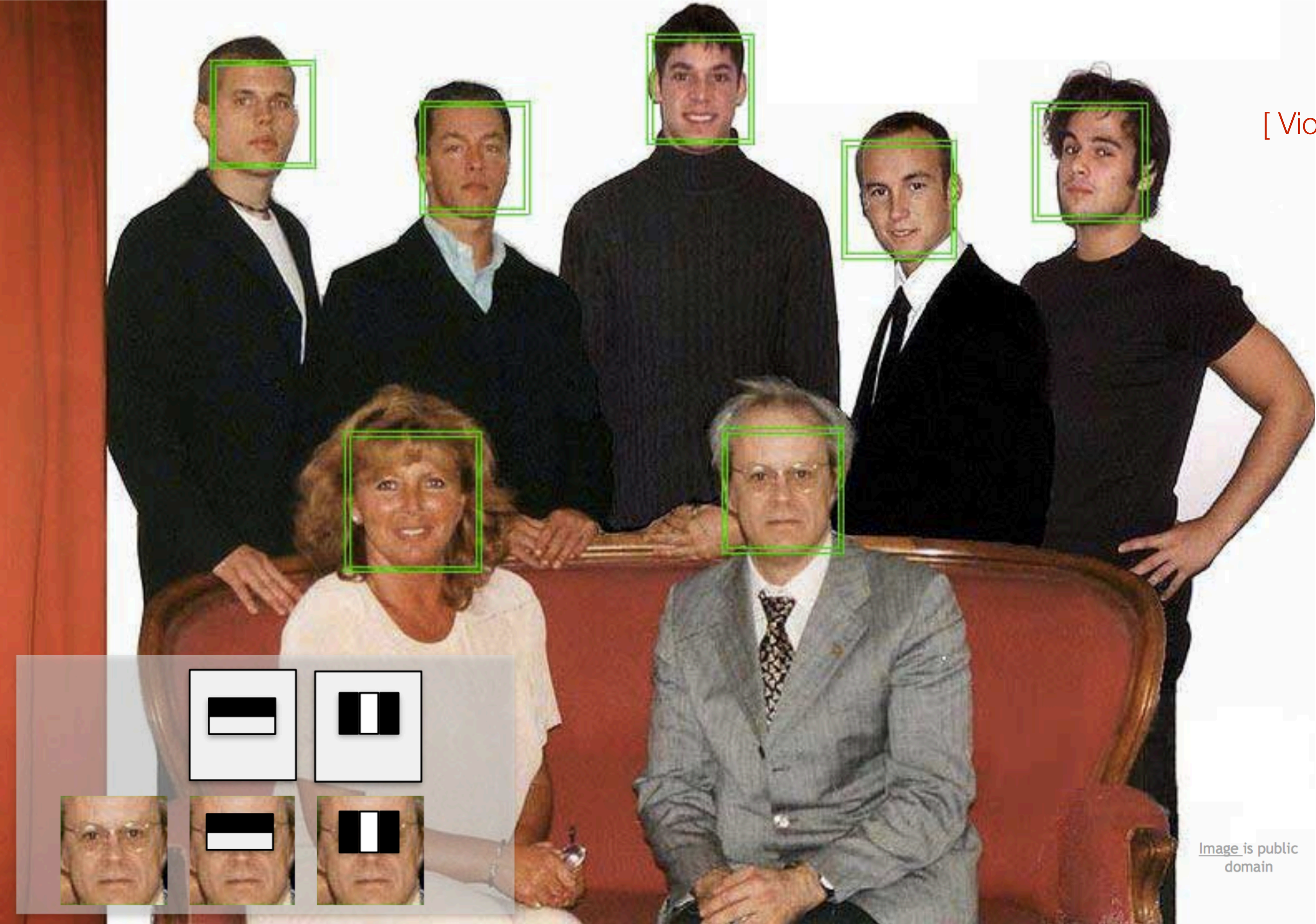
Monty Python's **Ministry of Silly Walks**



# David Marr, 1970s



# Face Detection 1999-2000



[ Viola & Jones, 2001 ]

Image is public domain

# Feature-based Vision



Image is public domain

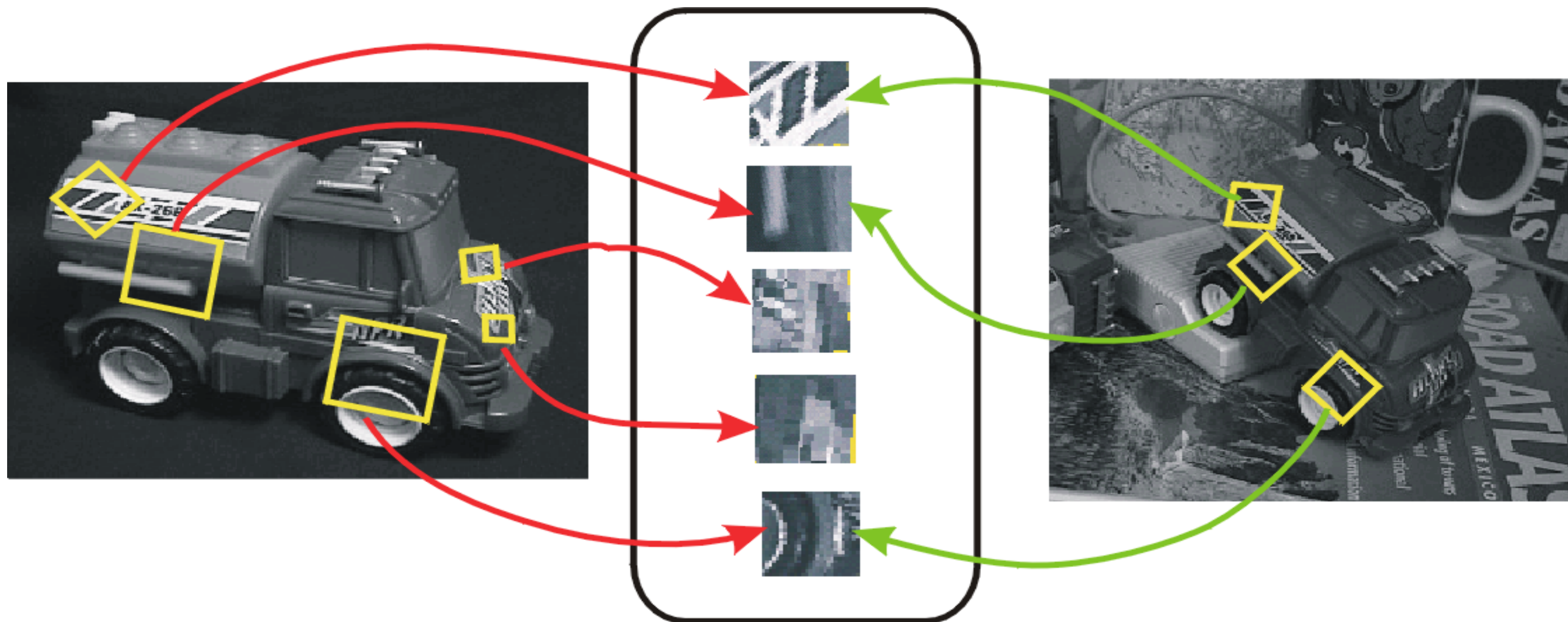


Image is CC BY-SA 2.0

[ **David Lowe**, 1999 ]

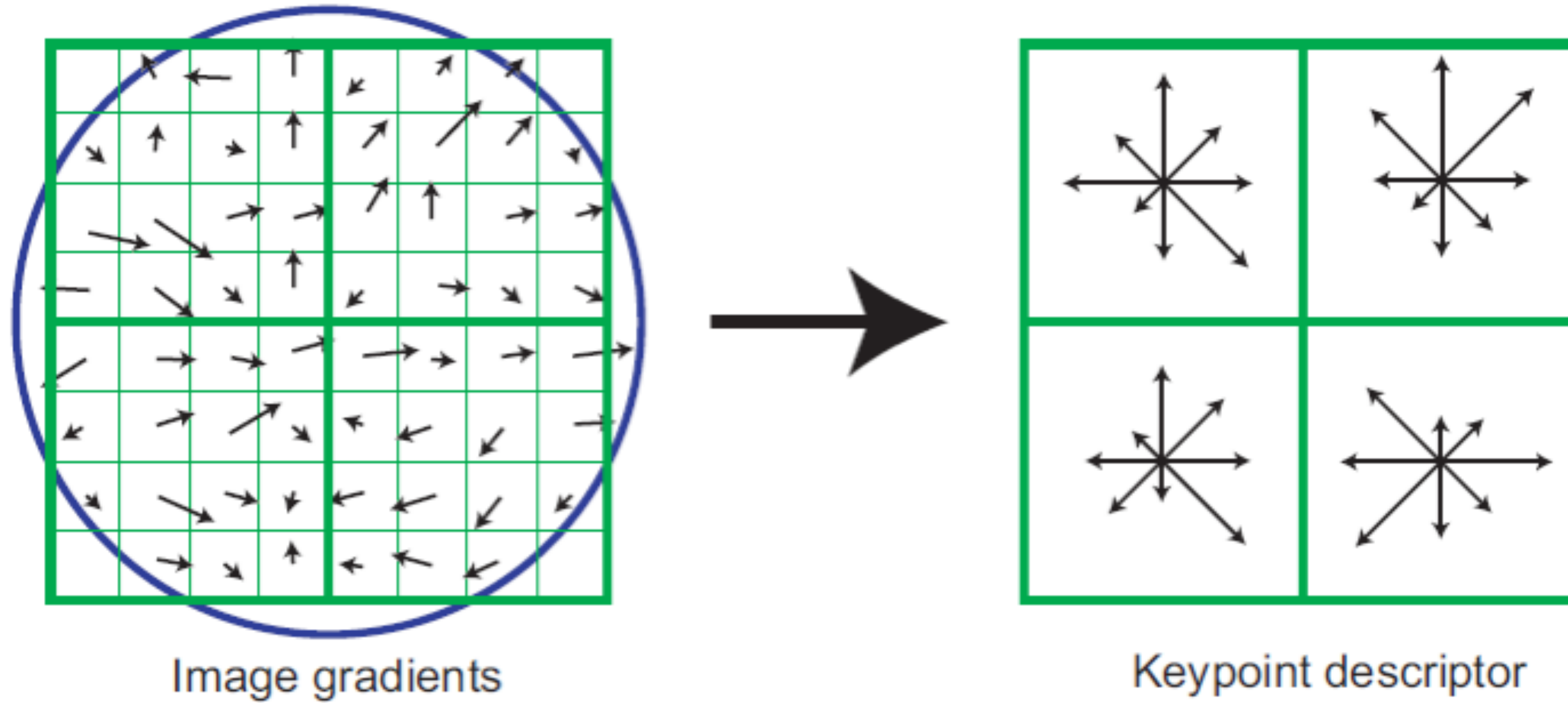
# SIFT Idea

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 ]

# Massive 3D Reconstructions

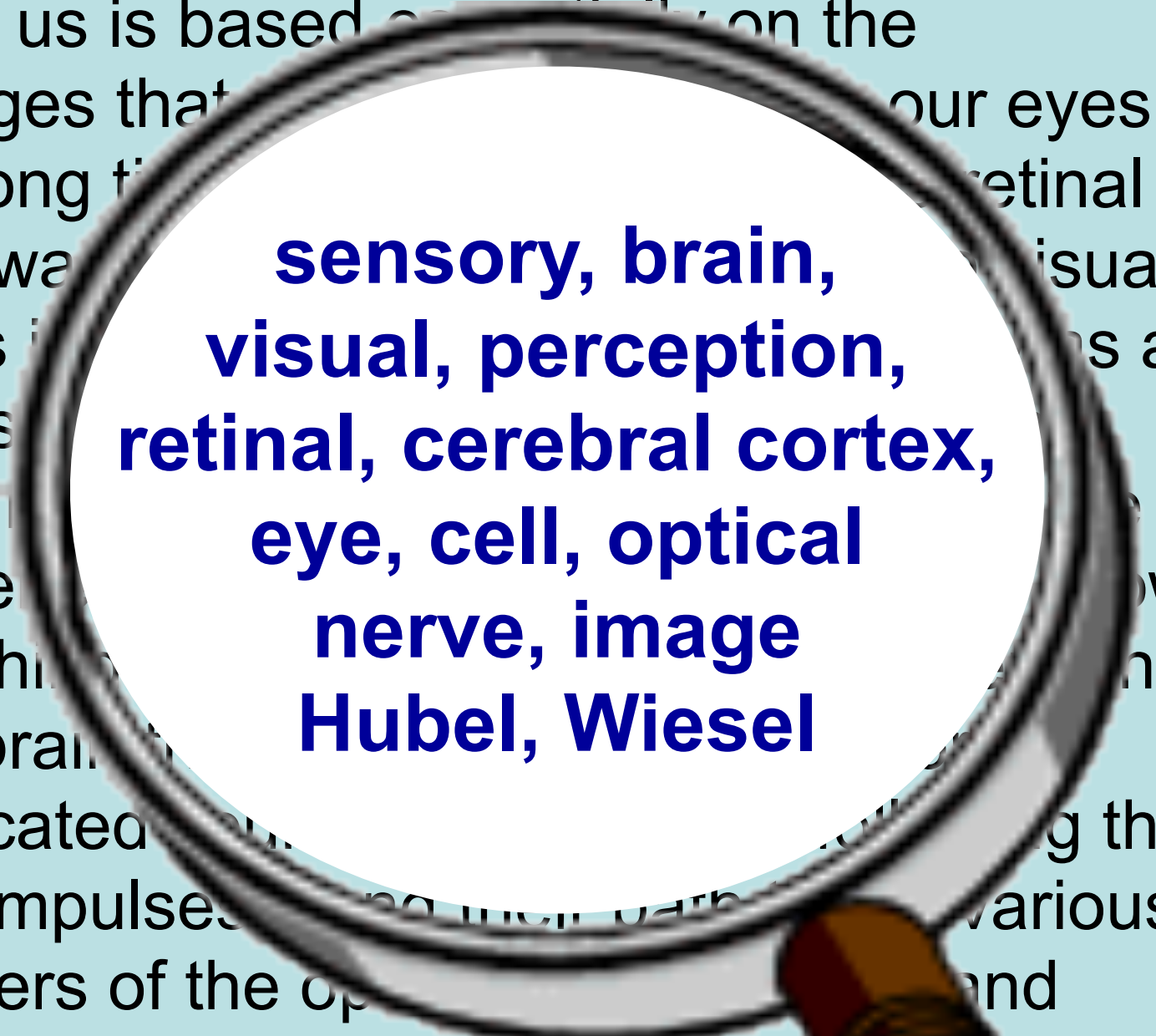


[ Agarwal, Furukawa, Snavely, Curless, Seitz, Szeliski, 2010 ]

# 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 primarily on the messages that come from our eyes. For a long time, the visual image was thought of as a picture that is processed in the visual centers of the brain. However, the discovery of the visual pathway and the work of 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**

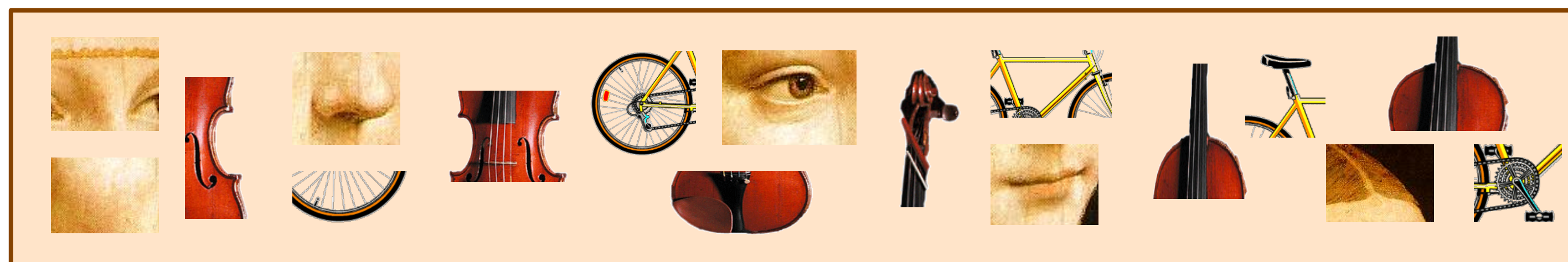
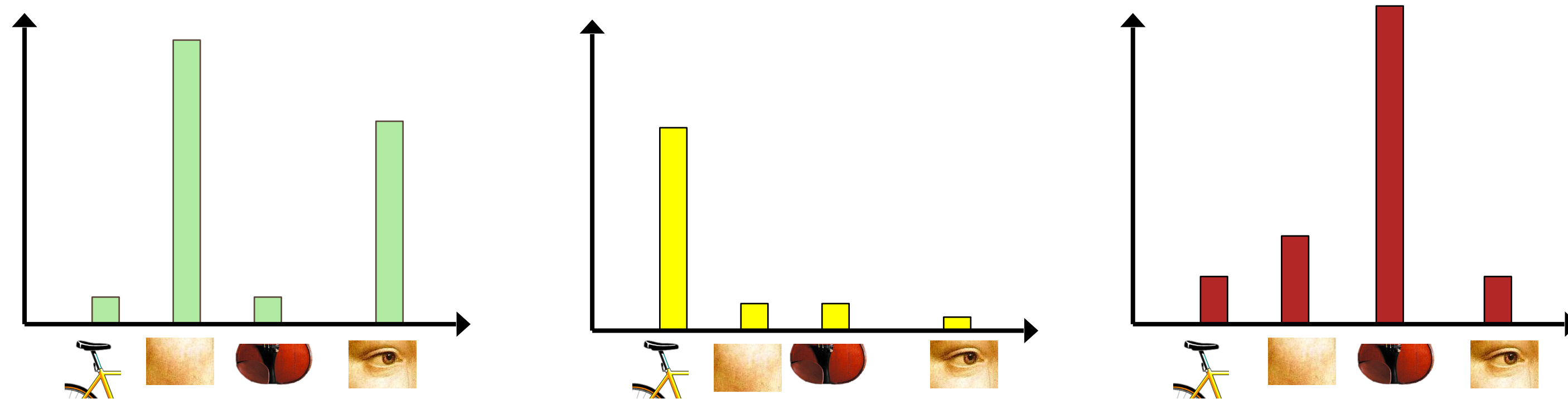
China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus will be created by a predicted 30% increase in exports to \$750bn, compared with \$580bn in 2004. The increase will annoy the US, which has long complained about China's trade surplus. The US government has agreed to a deal with China that the yuan is to be allowed to rise in value. The governor of the People's Bank of China also needed to be convinced that the demand so much of the world's goods from the country. China increased the value of the yuan 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.



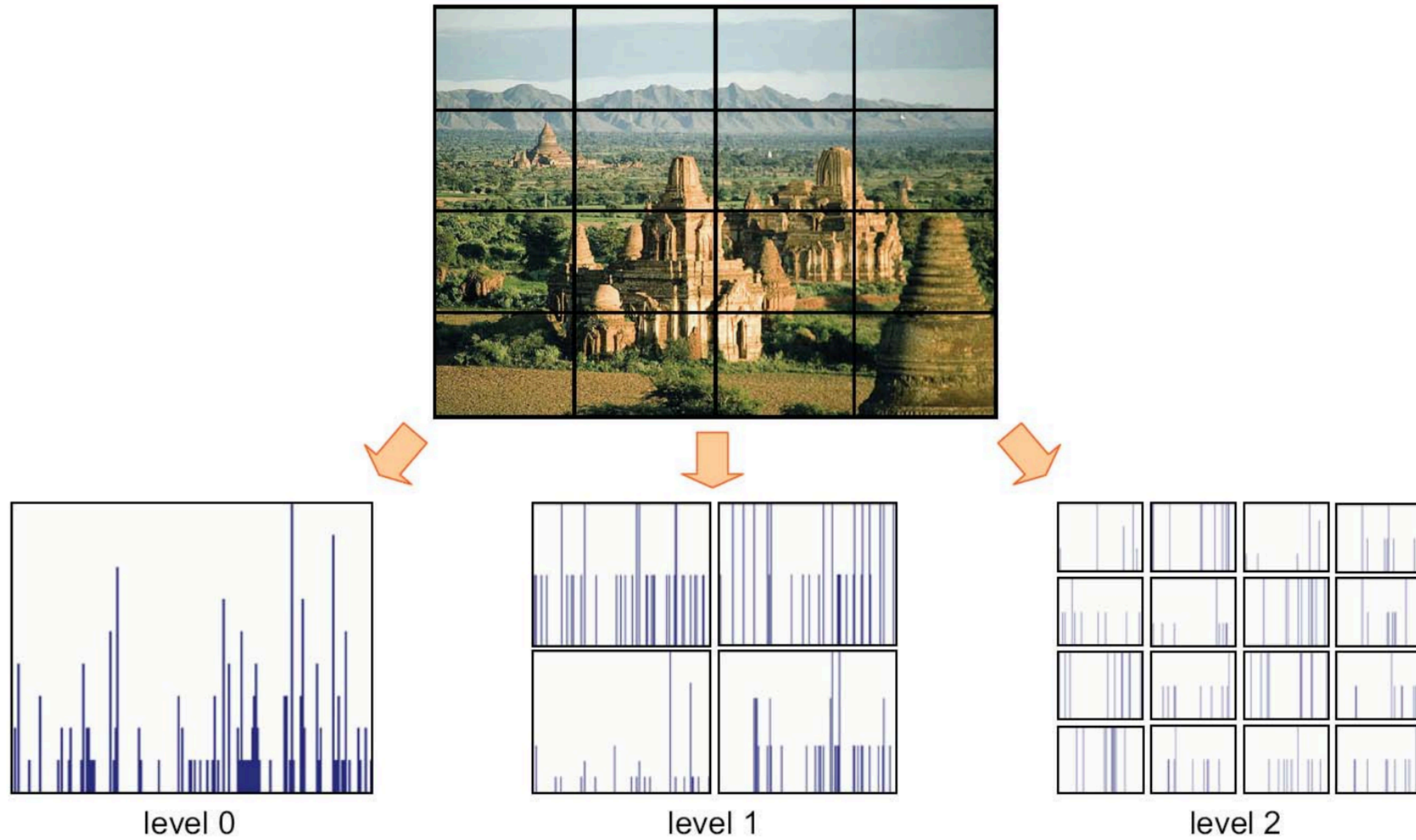
**China, trade,  
surplus, commerce,  
exports, imports, US,  
yuan, bank, domestic,  
foreign, increase,  
trade, value**



# Bag-of-Visual-Words

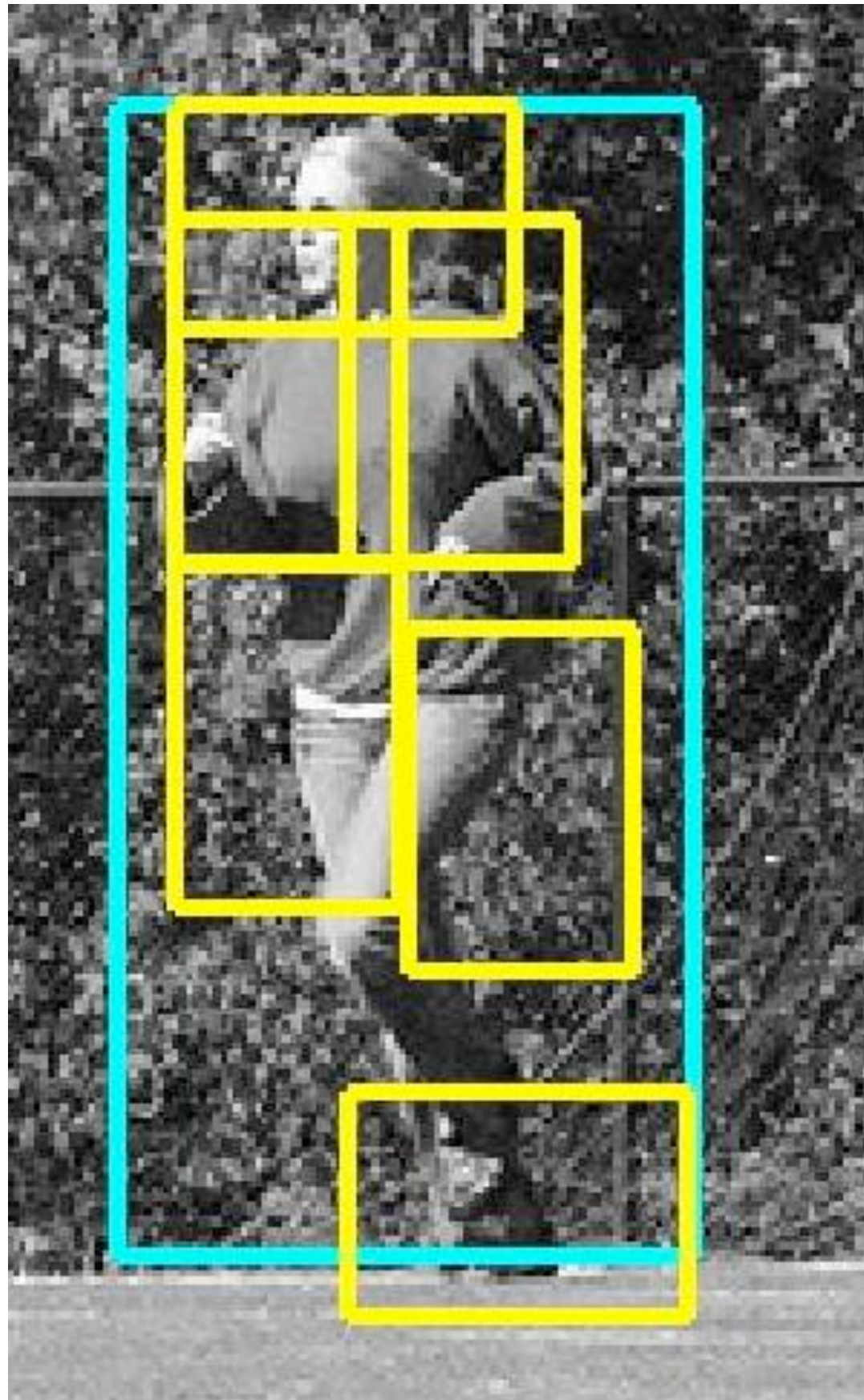


# Beyond Bag of Features

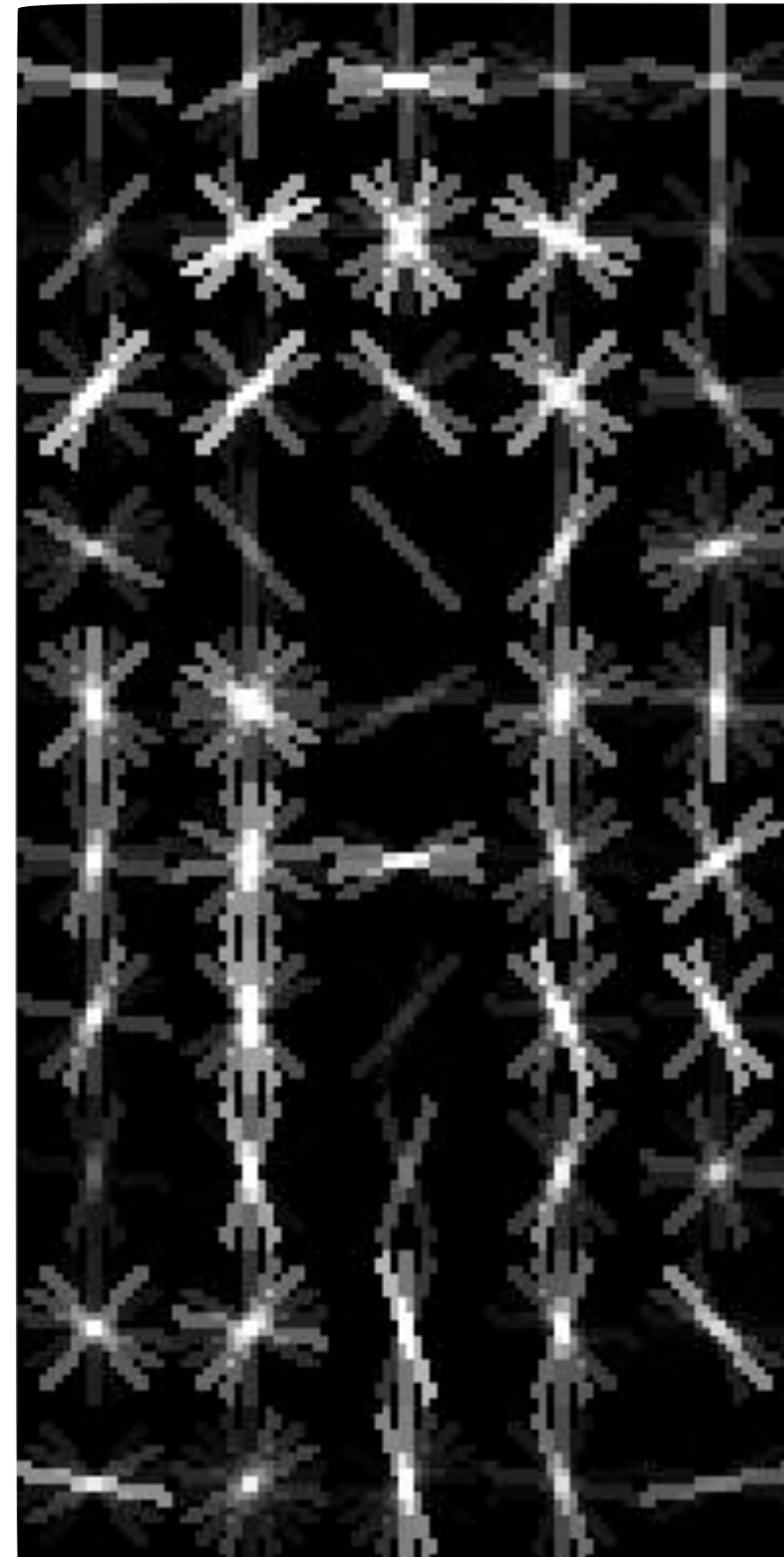


[ Lazebnik, Schmid, Ponce, 2006 ]

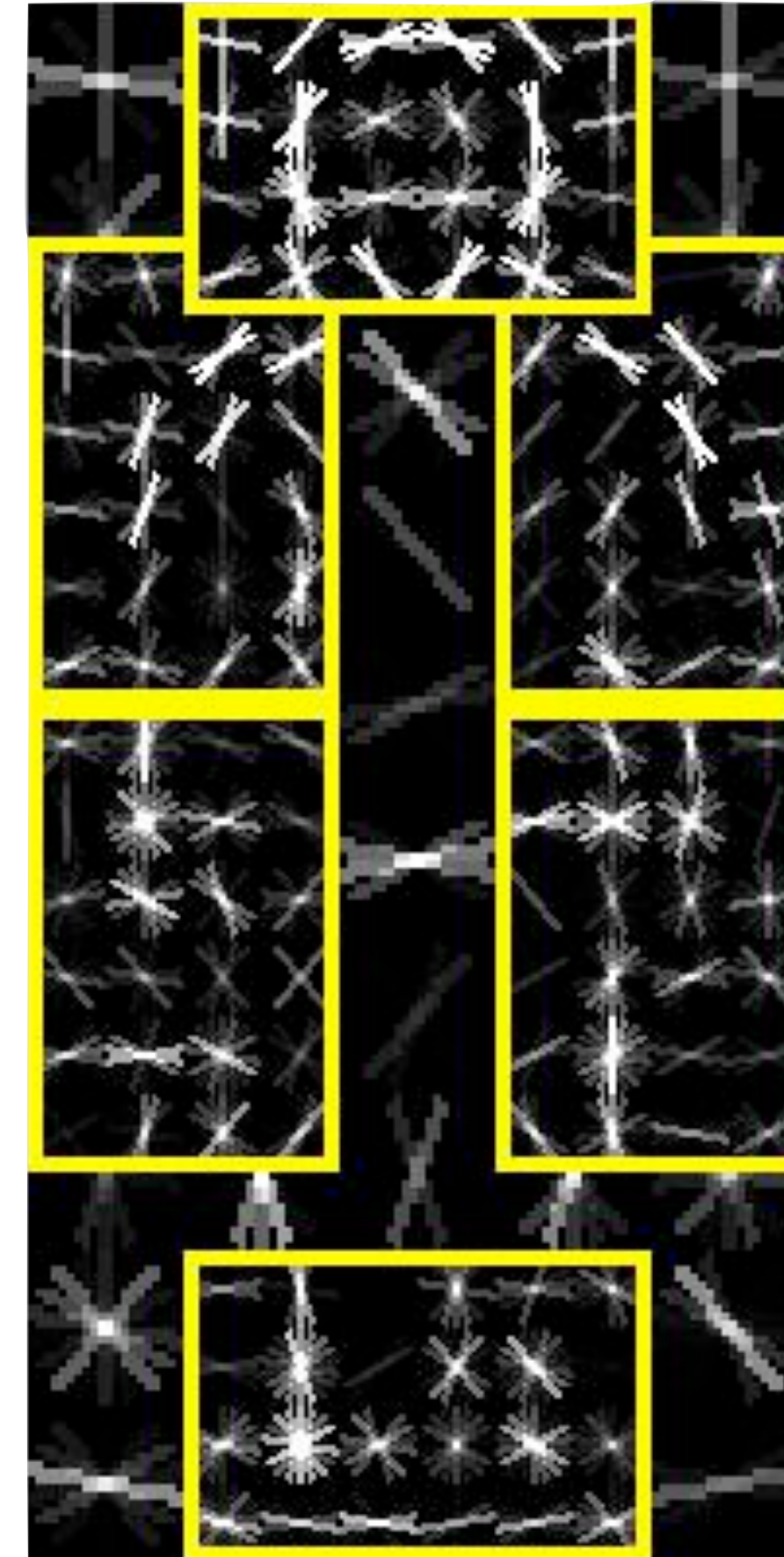
# Deformable Part Models



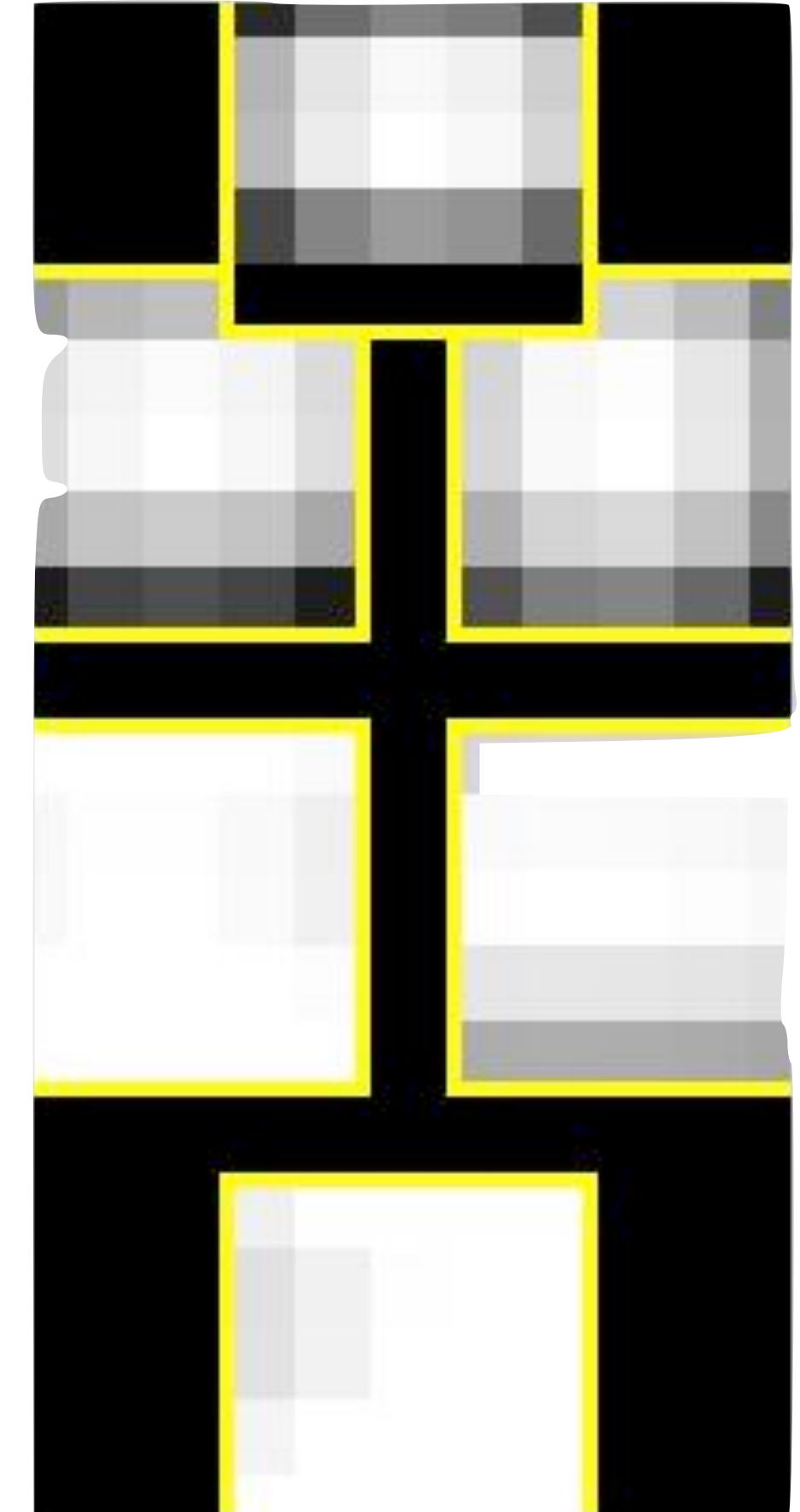
Detection



Root Filter

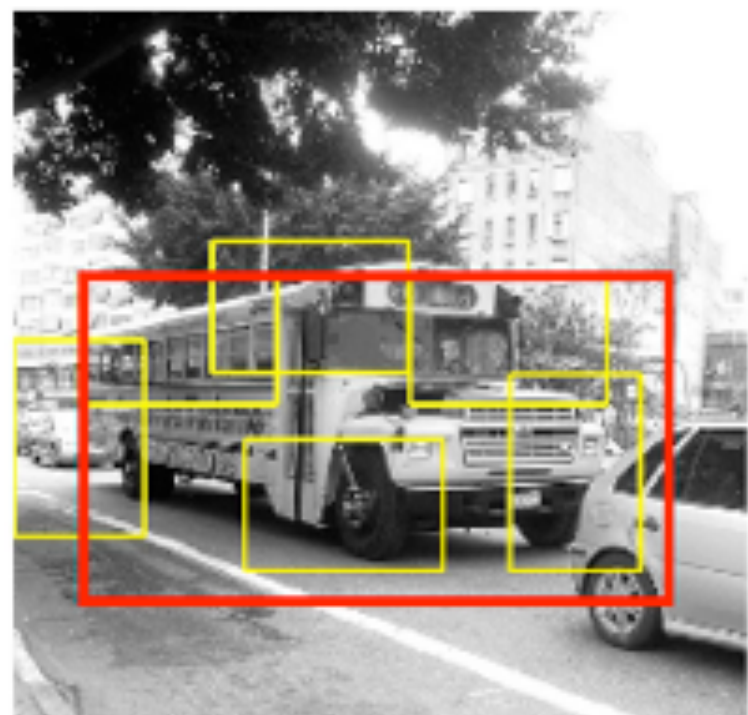
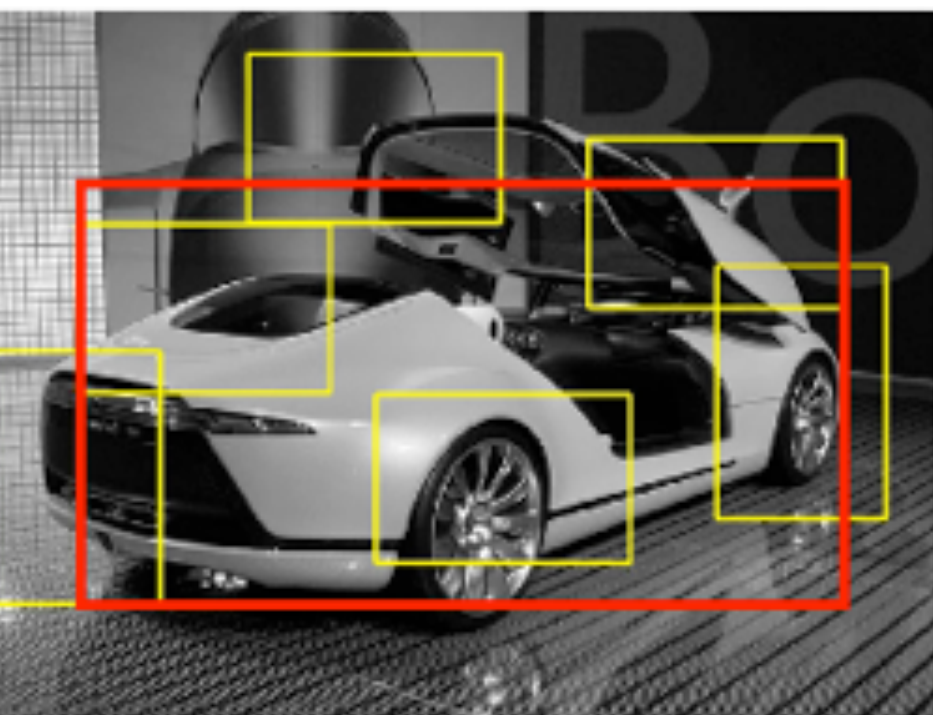
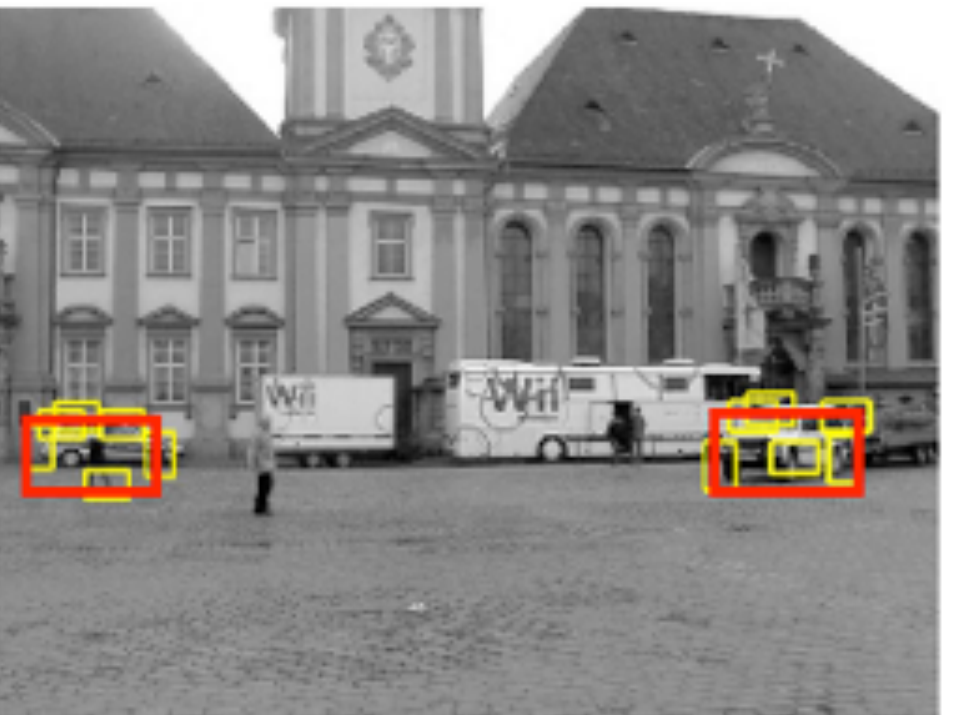
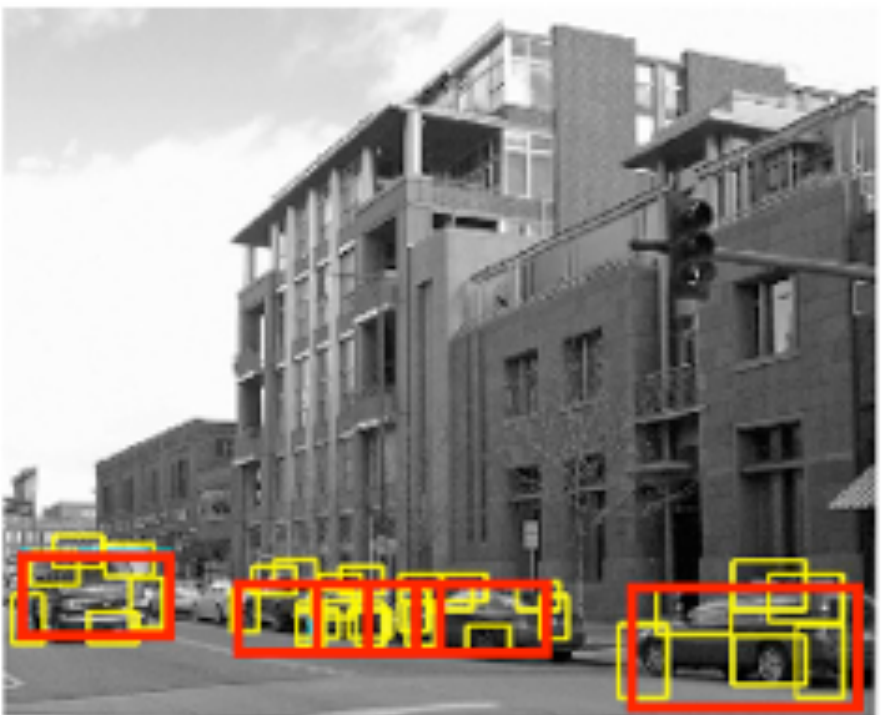
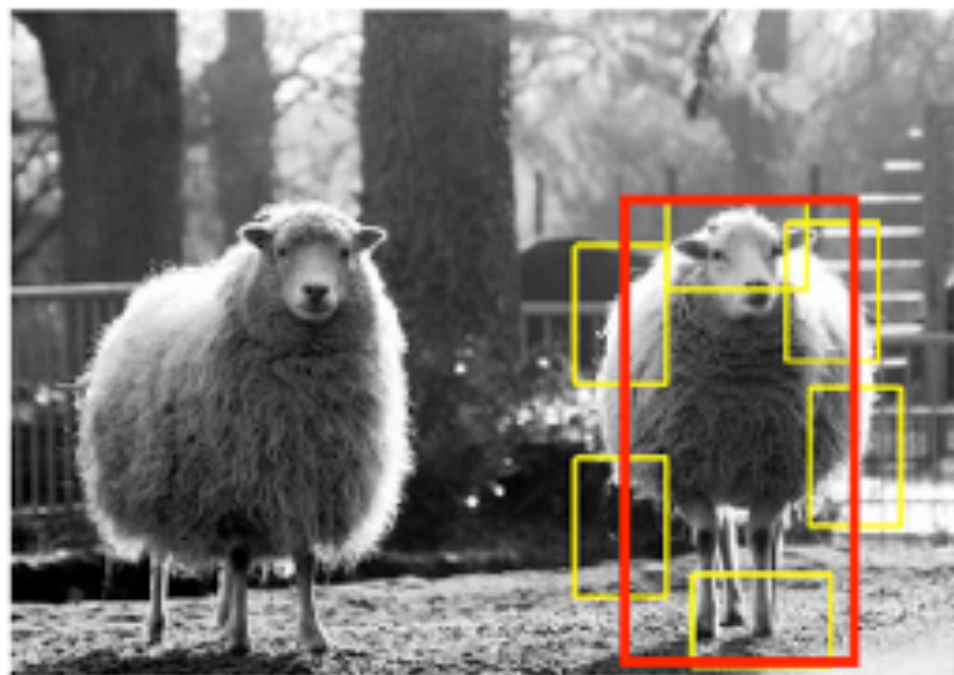
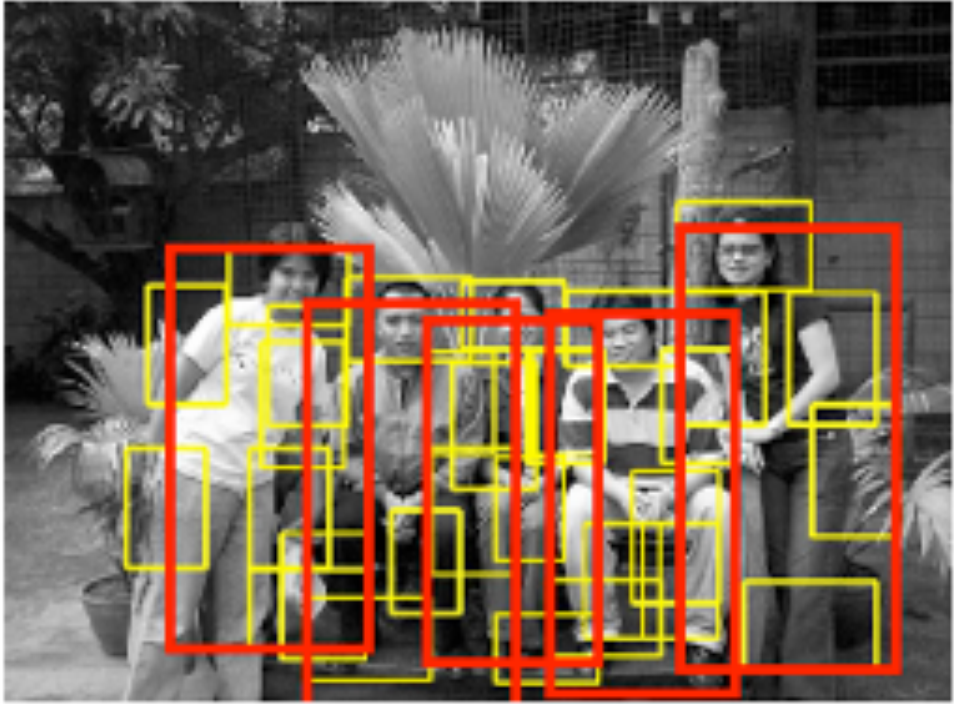


Part Filters

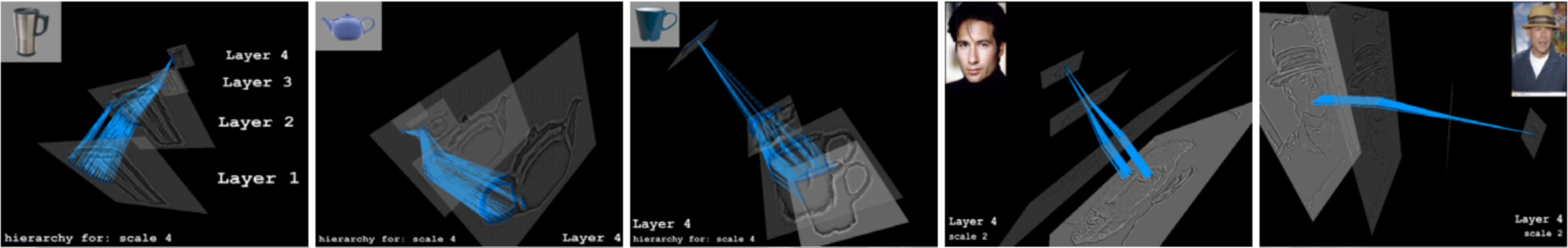
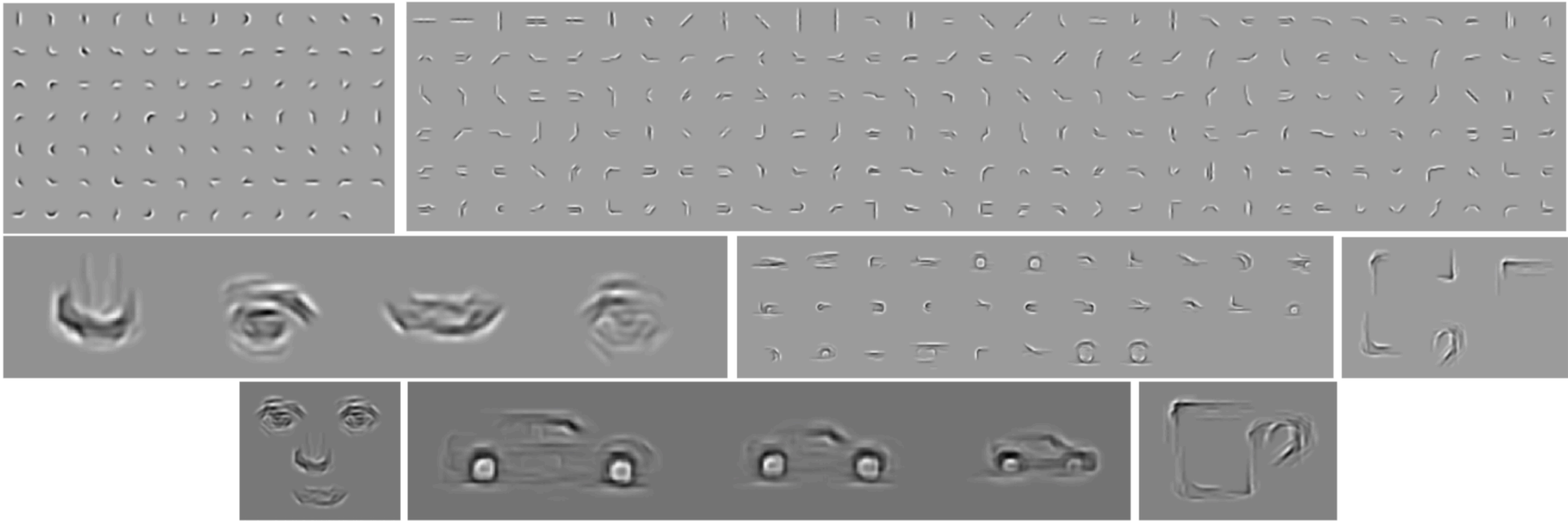


Deformations

# Deformable Part Models



# Hierarchical Models



# PASCAL Visual Object Challenge (VOC)

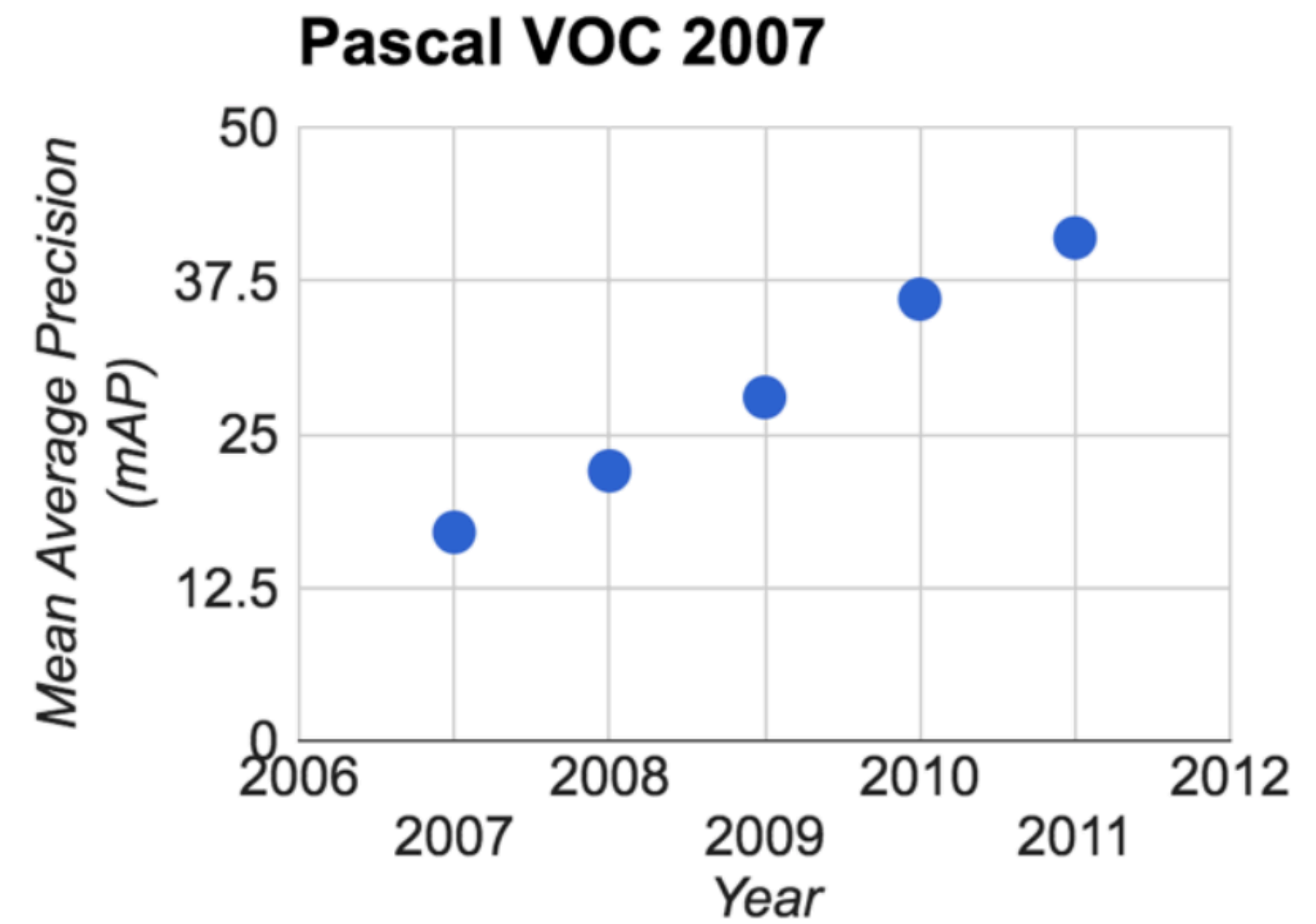
Image is CC BY-SA 3.0



Image is CC0 1.0 public domain

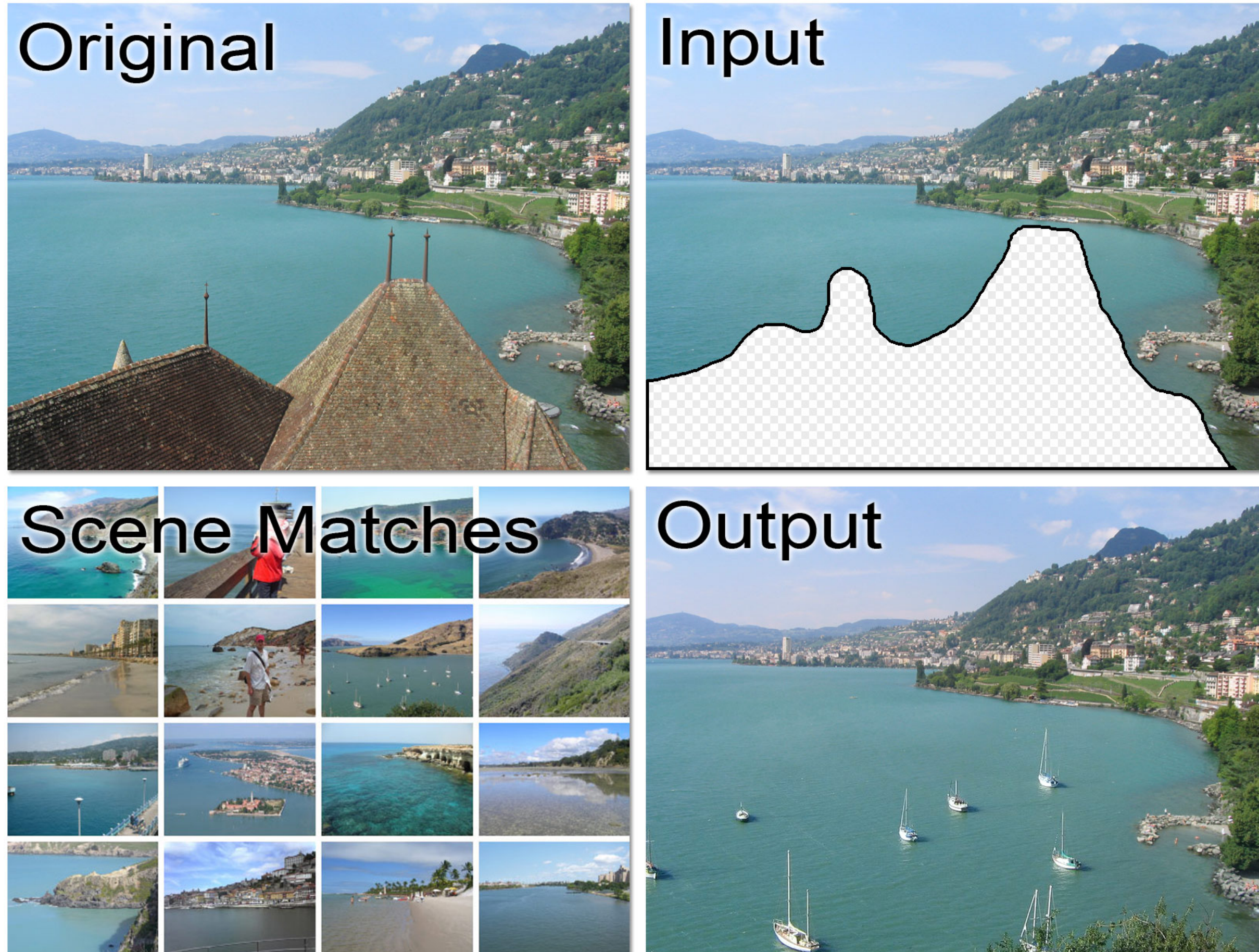


This image is licensed under CC BY-SA 2.0; changes made



[ 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

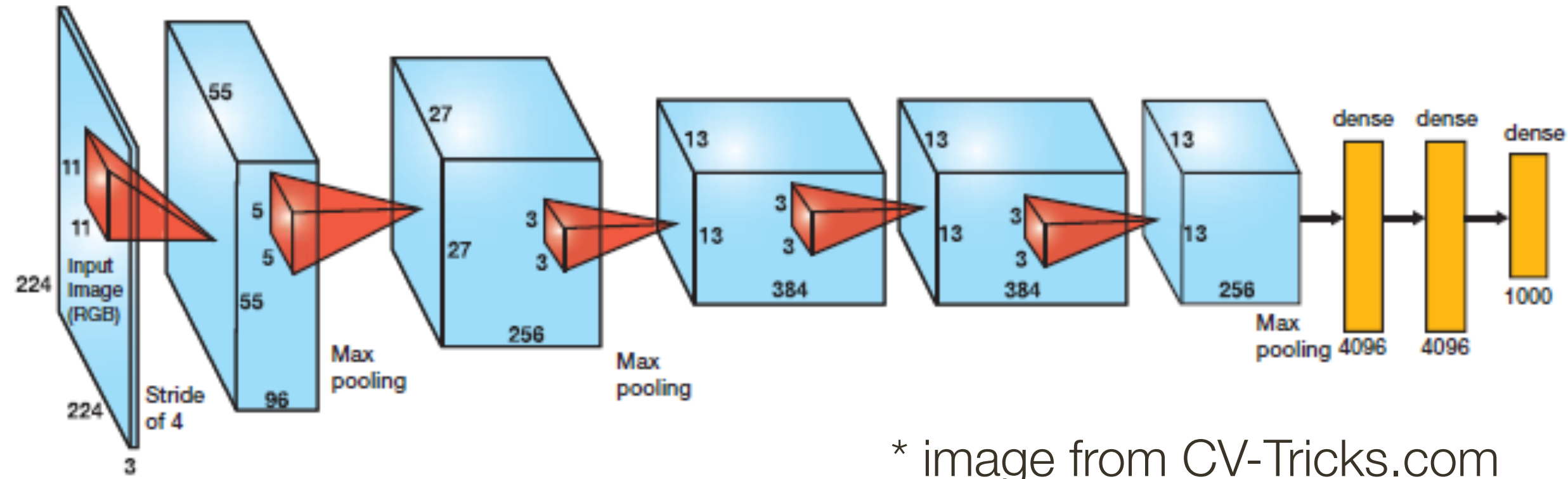
- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plants
  - Tree
  - Flower
  - Food
  - Materials
- Structures
  - Artifact
    - Tools
    - Appliances
    - Structures
- Person
  - Scenes
    - Indoor
    - Geological Formations
  - Sport Activities



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%
<b>CNN</b>	<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

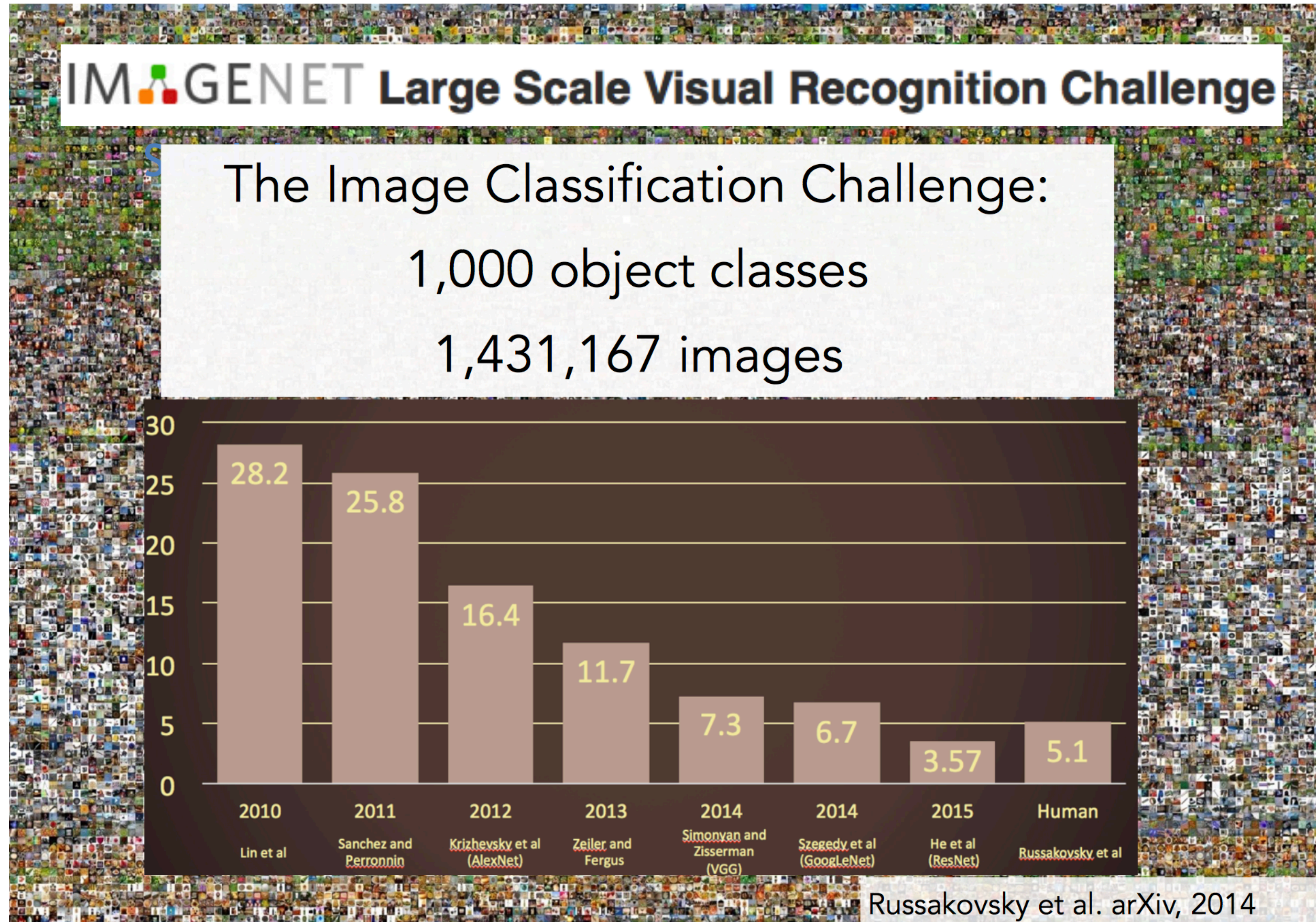
### Abstract

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

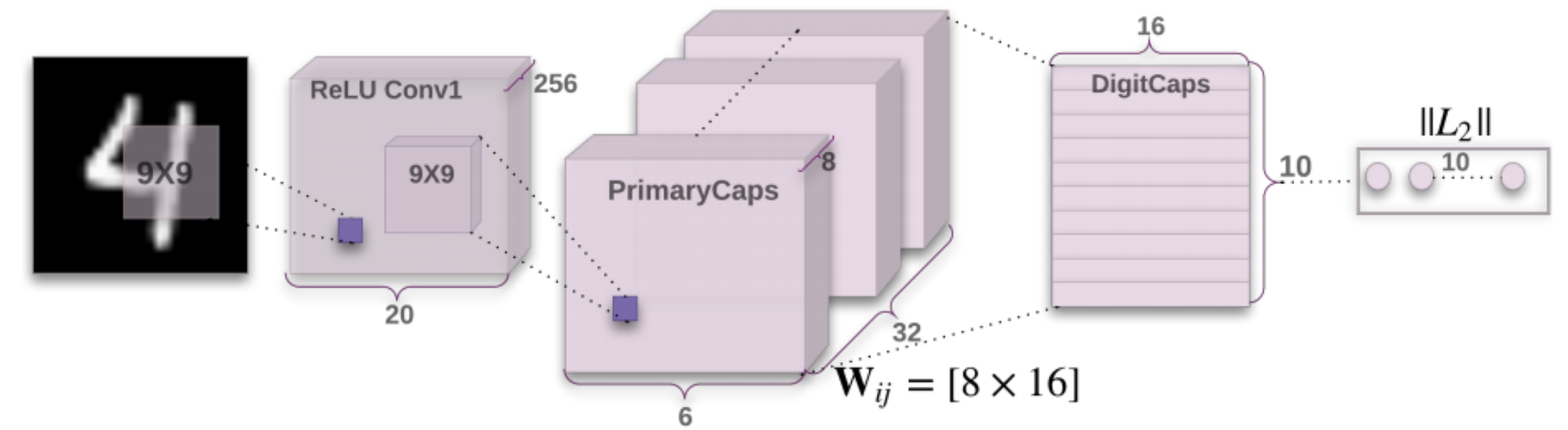


# Final thought ...

- 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 ]



person 0.88

reddish orange color 0.78

light brown color 0.78

starlet 0.66

entertainer 0.66

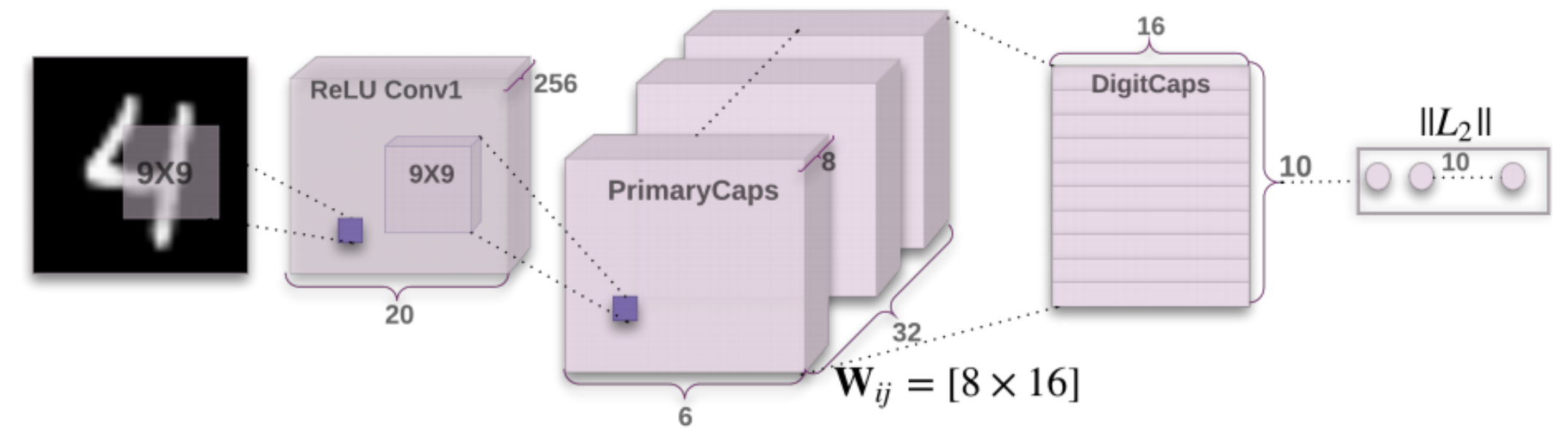
female 0.60

woman 0.59

young lady (heroine) 0.59

# CapsuleNET

Going **back to inverse** graphics



[ Sabour, Frosst, Hinton, NIPS 2017 ]



person 0.88



- reddish orange color 0.78
- light brown color 0.78
- starlet 0.66
- entertainer 0.66
- female 0.60
- woman 0.59
- young lady (heroine) 0.59



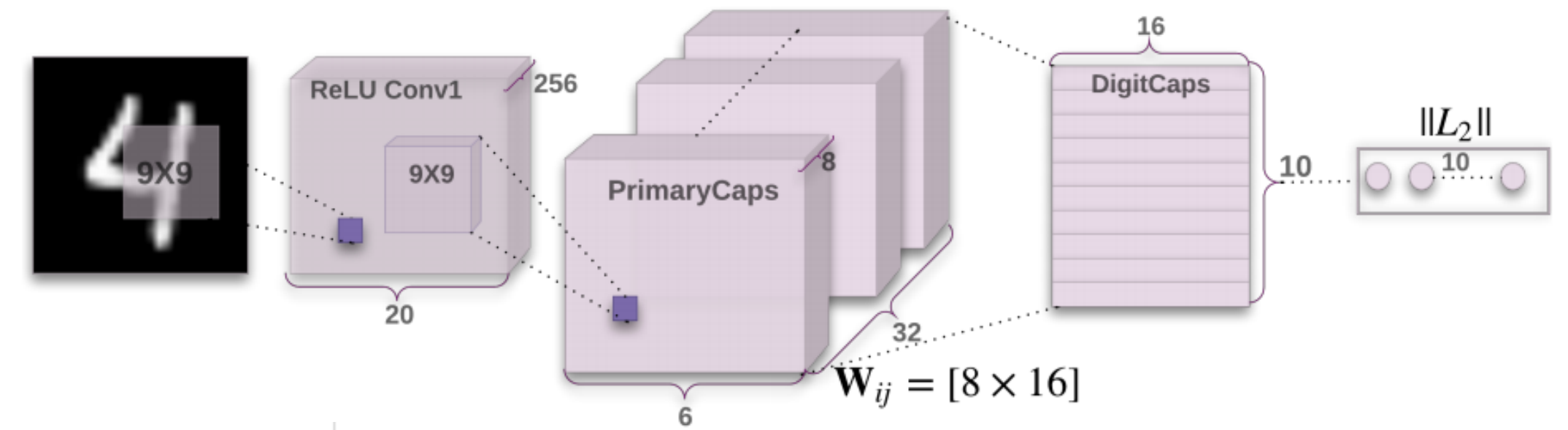
person 0.90



- light brown color 0.84
- starlet 0.77
- entertainer 0.77
- female 0.65
- woman 0.64
- young lady (heroine) 0.64
- reddish orange color 0.64
- newsreader 0.50

# CapsuleNET

Going **back to inverse** graphics



[ Sabour, Frosst, Hinton, NIPS 2017 ]



person 0.88



- reddish orange color 0.78
- light brown color 0.78
- starlet 0.66
- entertainer 0.66
- female 0.60
- woman 0.59
- young lady (heroine) 0.59



person 0.90



- light brown color 0.84
- starlet 0.77
- entertainer 0.77
- female 0.65
- woman 0.64
- young lady (heroine) 0.64
- reddish orange color 0.64
- newsreader 0.50



coal black color 0.79



- hairpiece (hair) 0.71
- dress 0.71
- maroon color 0.71
- person 0.58
- toupee (hairpiece) 0.58
- woman 0.56
- Earrings 0.55
- female 0.50

# Neural Modular Networks

