

Towards Visual Intelligence

Aykut Erdem

The International Symposium on Brain and Cognitive Science
May 6, 2018



HACETTEPE
UNIVERSITY
COMPUTER
VISION LAB



The Purpose of Vision

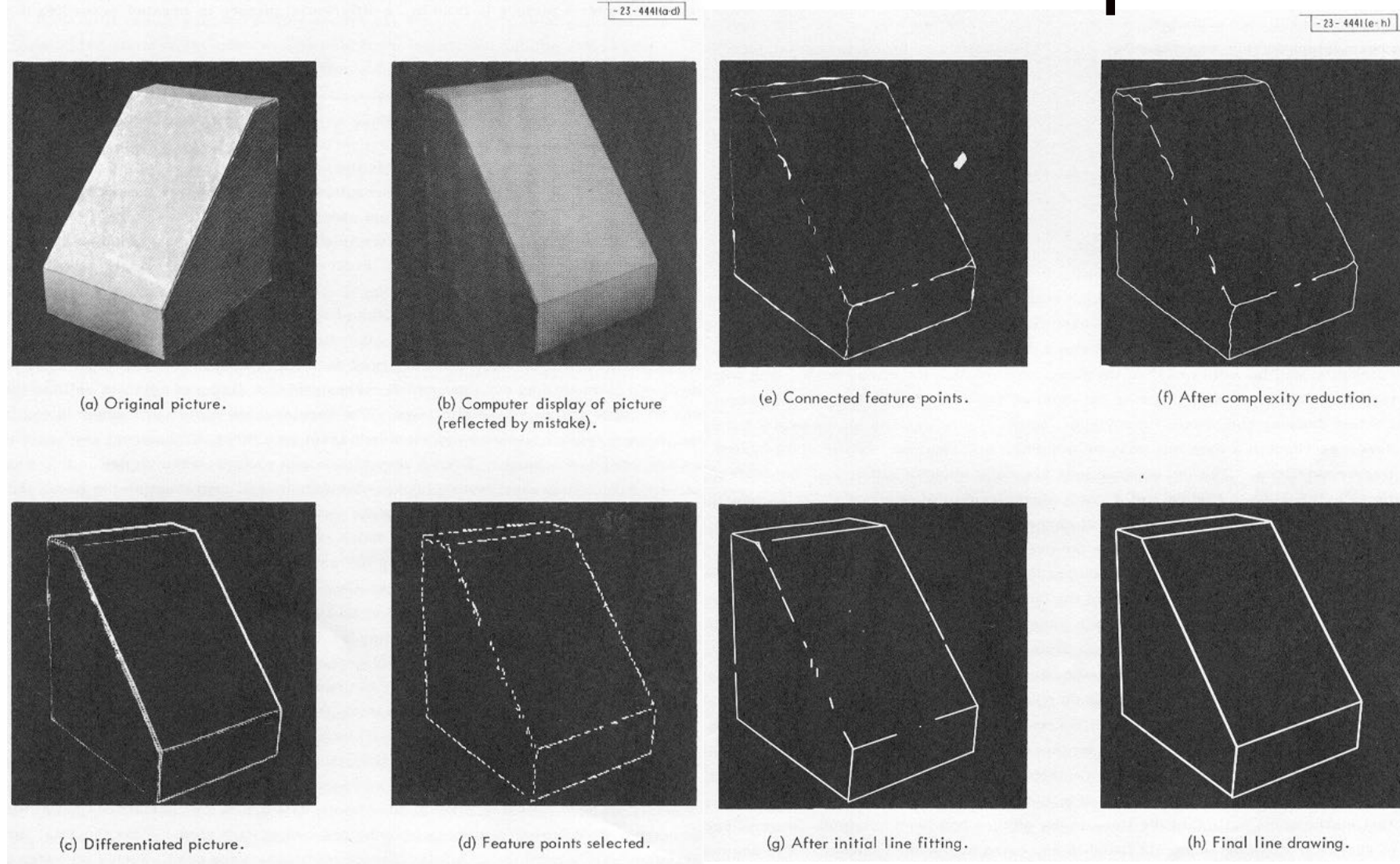


"What does it mean, to see? The plain man's answer (and Aristotle's too) would be, **to know what is where by looking.** In other words, vision is the process of discovering from images what is present in the world, and where it is."

[Marr, 1982]

Image credit: The Sense of Sight (Annie Louisa Swynnerton, 1895) 2

The First PhD Thesis on Computer Vision



- Machine perception of three-dimensional solids [Roberts 1963]

The Summer Vision Project

General goals:

FIGURE-GROUND.

divide a vidisector picture into regions such as likely objects, likely background areas and chaos

REGION DESCRIPTION.

analysis of shape and surface properties

OBJECT IDENTIFICATION.

name objects by matching them with a vocabulary of known objects

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group
Vision Memo. No. 100.

July 7, 1966

THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".

[Papert 1966]

Why does vision appear easy to humans?

- Our brains are specialized to do vision.
- ~50% of the cortex in a human brain is devoted for visual processing (cf. motor control ~20-30%, language ~10-20%)

Visual perception*:

540,000,000 years of data

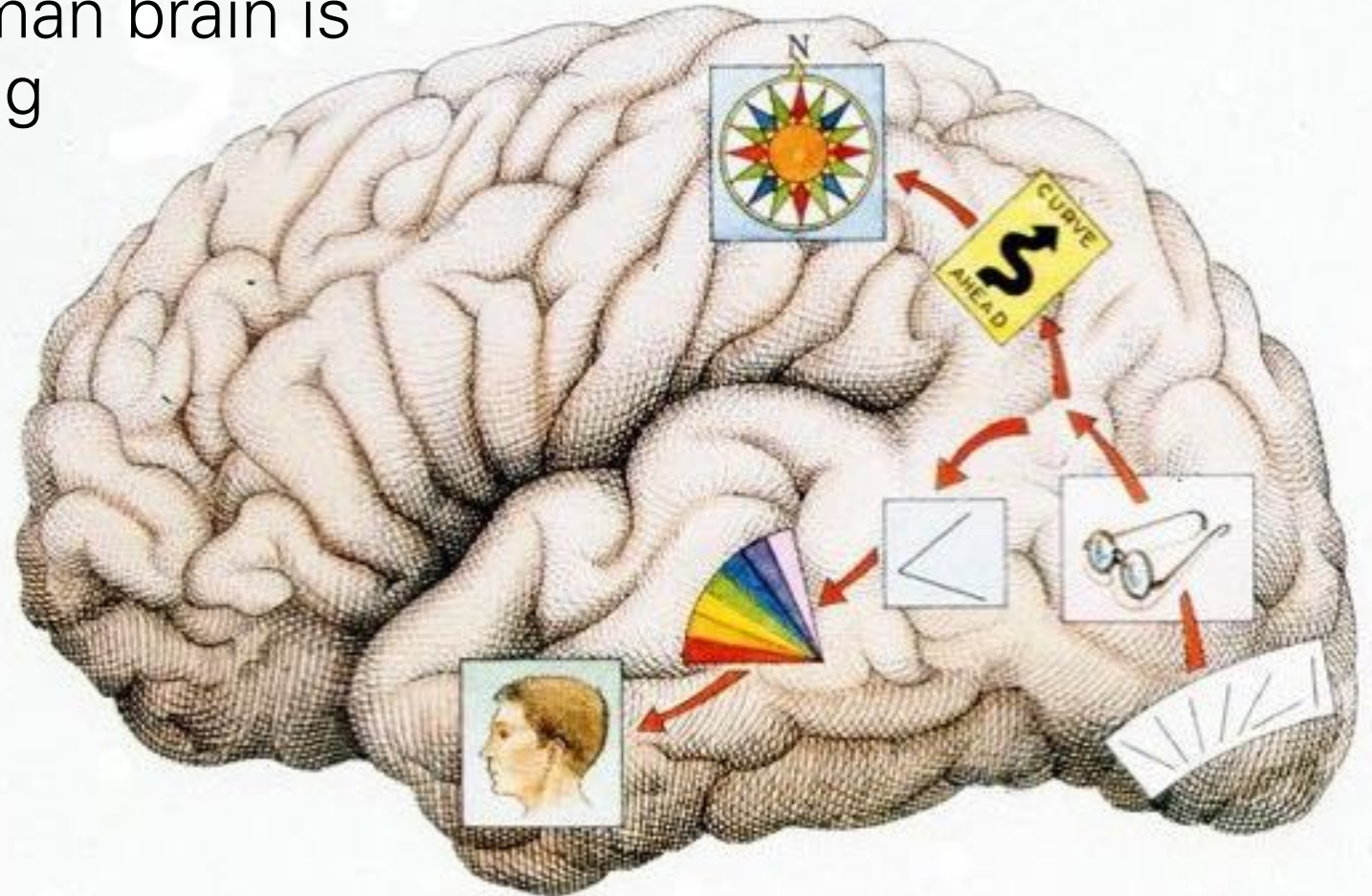
Bipedal movement:

230,000,000 years of data

Abstract thought:

100,000 years of data

*Color vision



Fast Forward to 2012

IMAGENET Large Scale Visual Recognition Challenge (ILSVRC)

- **1.2M** training images, **1K** categories
- Measure top-5 classification error

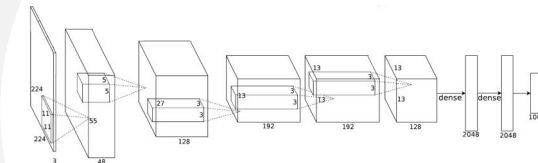
The success of AlexNet, a deep convolutional network (CNN)

- 7 hidden layers (not counting some max pooling layers)
- 60M parameters

2012 Teams	%Error
Supervision (Toronto)	15.3
ISI (Tokyo)	26.1
VGG (Oxford)	26.9
XRCE/INRIA	27.0
UvA (Amsterdam)	29.6
INRIA/LEAR	33.4

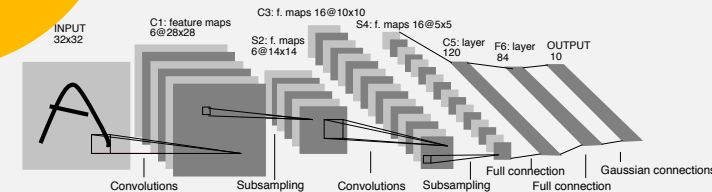
CNN based, non-CNN based

CNNs are biologically inspired by oriented cells in the visual cortex



AlexNet

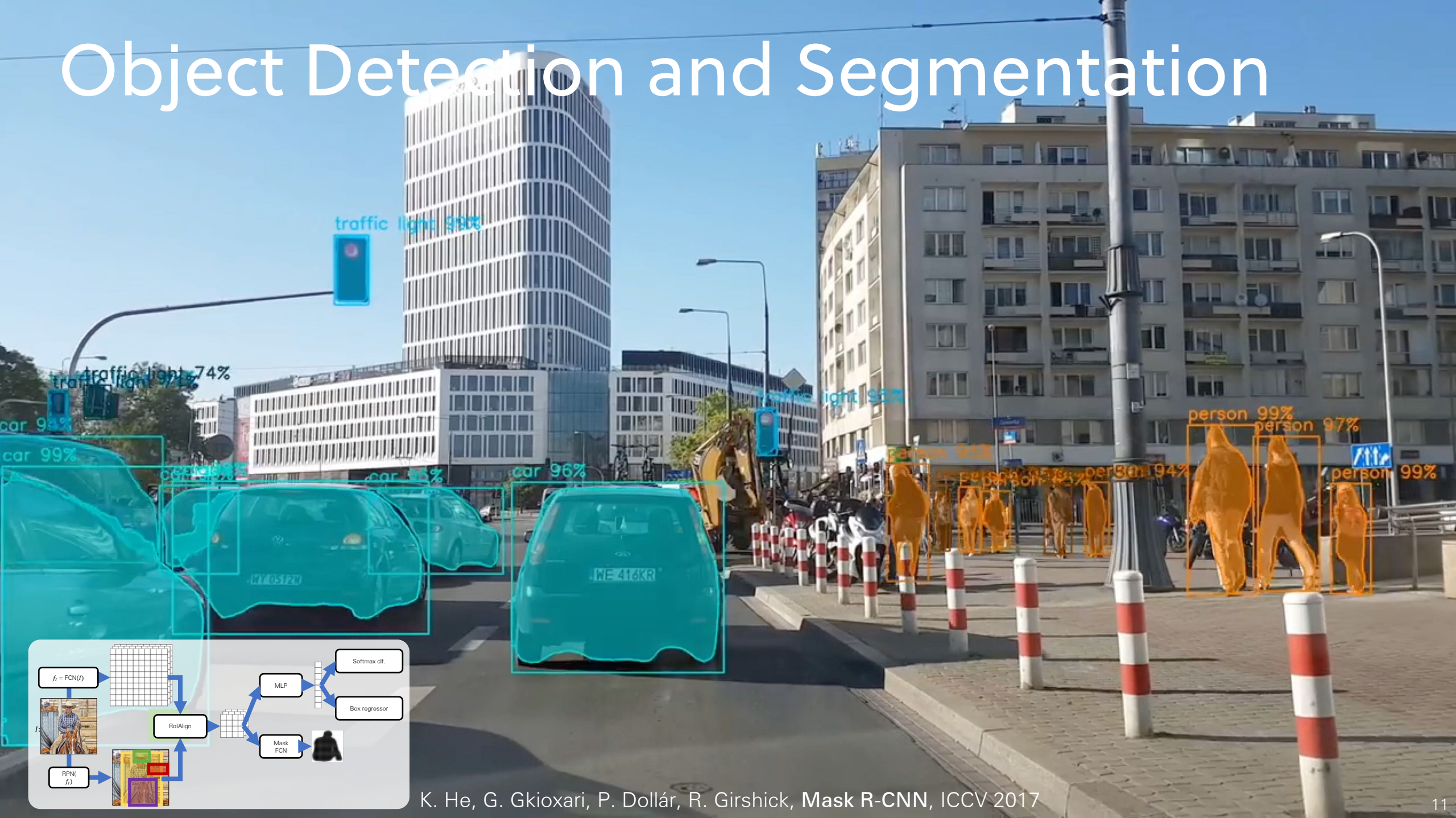
Cat



Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. **Gradient-based learning applied to document recognition**. Proceedings of the IEEE. 86 (11): 2278–2324, 1998.

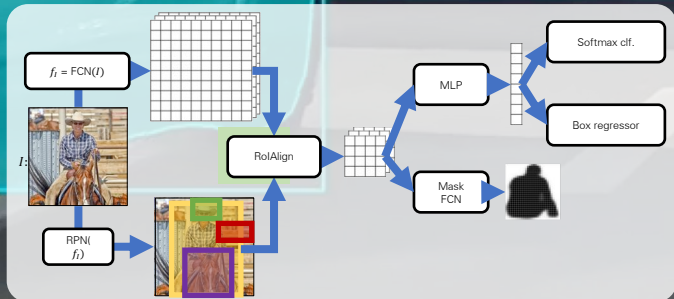
A. Krizhevsky, I. Sutskever, G.E. Hinton **ImageNet Classification with Deep Convolutional Neural Networks**. NIPS 2012

Object Detection and Segmentation



The image displays a street scene with various objects detected and segmented. Cars are highlighted in cyan, and pedestrians in orange. The diagram in the bottom left corner illustrates the Mask R-CNN architecture, showing the flow from input image I through feature extraction $f_i = \text{FCN}(I)$, region proposal network $\text{RPN}(f_i)$, region of interest alignment RoIAlign , and final classification and segmentation heads.

K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask R-CNN, ICCV 2017



K. He, G. Gkioxari, P. Dollár, R. Girshick, **Mask R-CNN**, ICCV 2017

11.4 fps

Pose Estimation

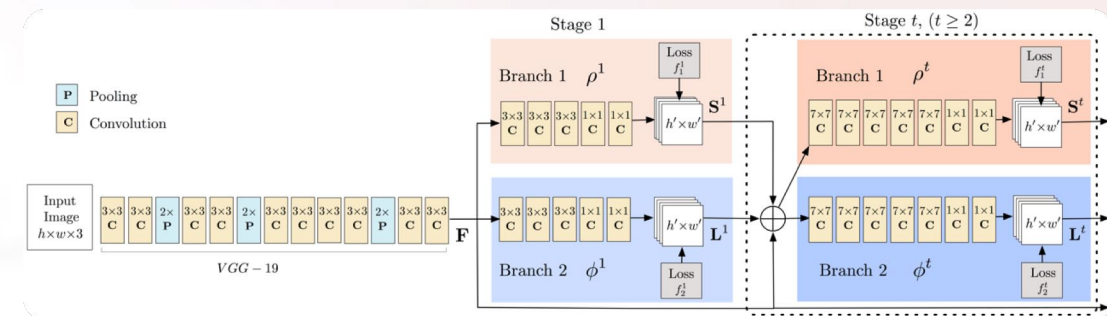


Photo Style Transfer



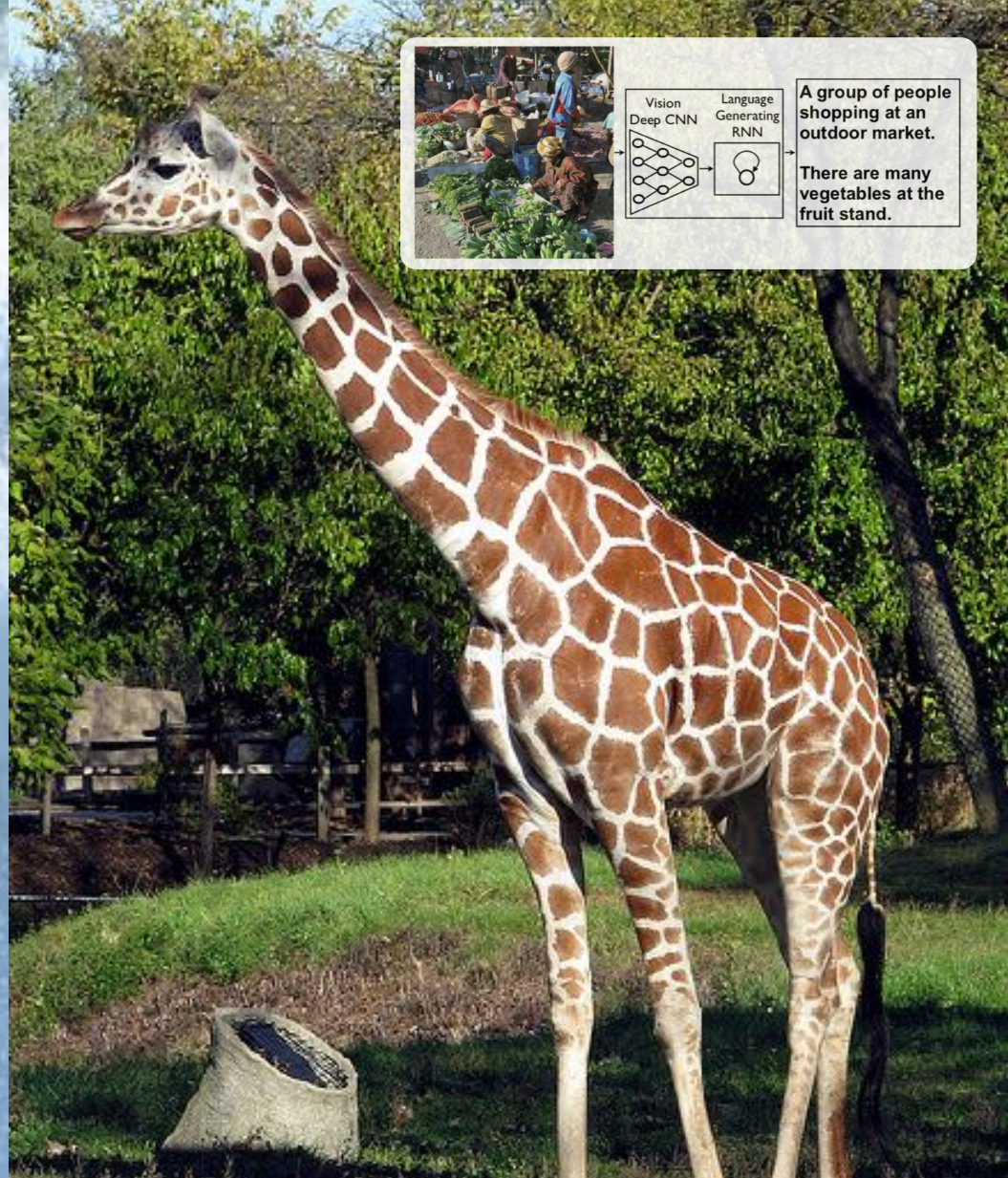
Photo Style Transfer



Image Captioning



A man riding a wave on a surfboard in the water.



A giraffe standing in the grass next to a tree.

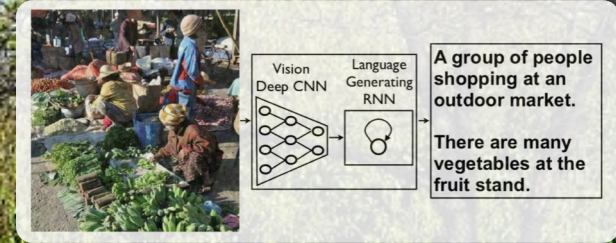
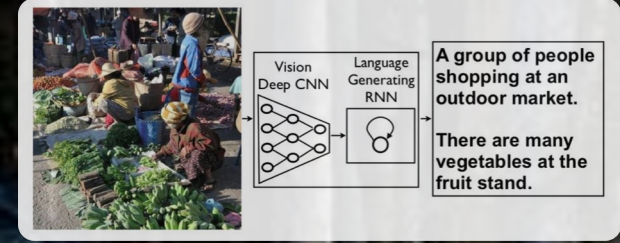
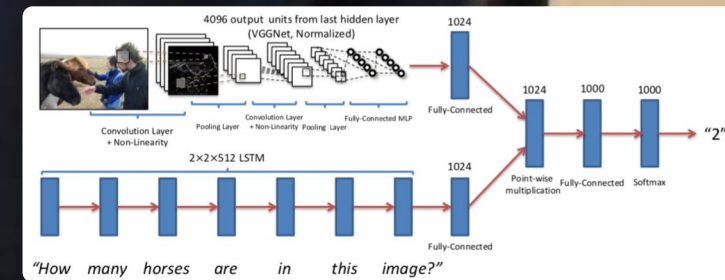


Image Captioning



Yarış pistinde virajı almakta olan bir yarış arabası

Visual Question Answering



Question: What is the girl reaching into?

Answer: apples

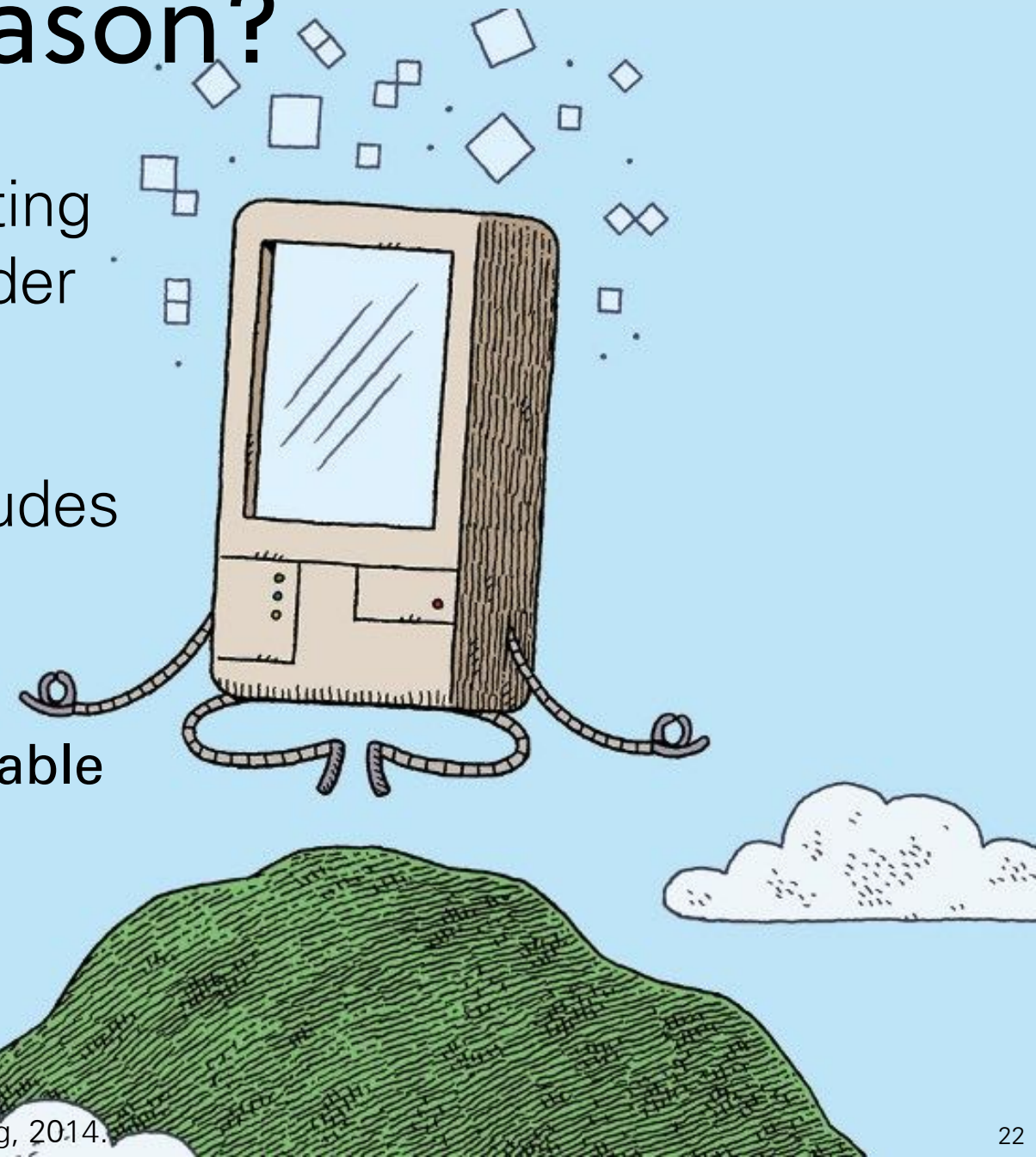
Can Deep Models Reason?

- Reasoning: “Algebraically manipulating previously acquired knowledge in order to answer a new question”
- A very broad definition



Can Deep Models Reason?

- Reasoning: “Algebraically manipulating previously acquired knowledge in order to answer a new question”
- A very broad definition, which includes
 - logical reasoning
 - probabilistic inference
 - composition rules operating on trainable modules

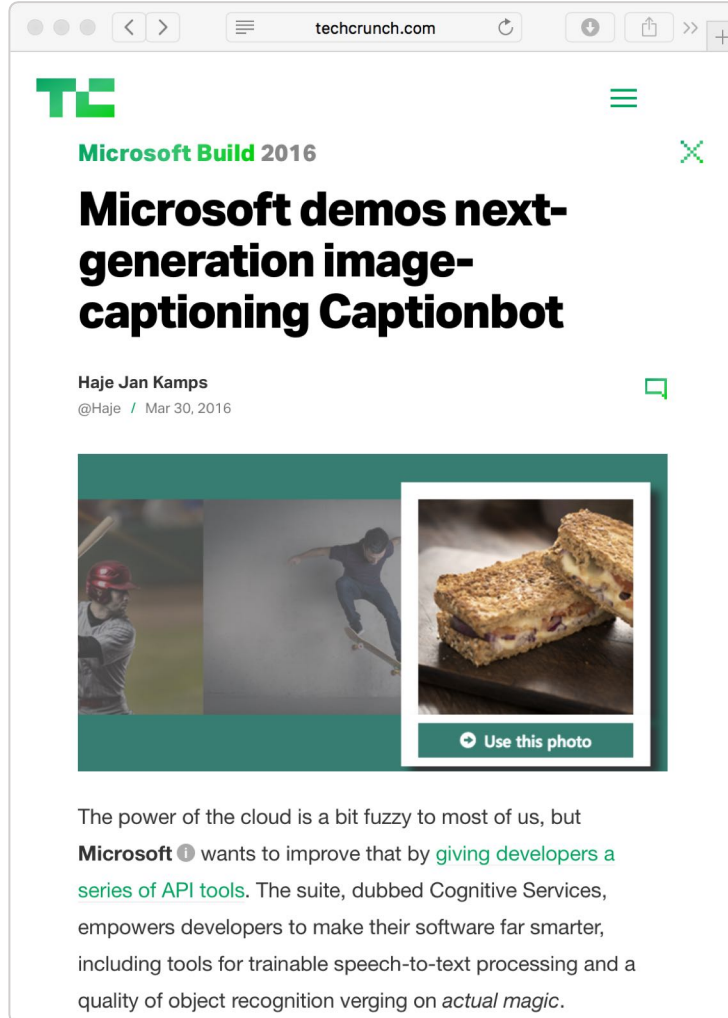


Can Deep Models Reason?

- Deep Learning models are **large correlation engines**
- They use **inductive bias** to learn from training data, which is a double-edged sword
 - Generalize well when **target** and **training** distributions are similar
 - Confuse **correlation** with **causation**



Take 2: Image Captioning



picdescbot @picdescbot · Feb 19
a herd of sheep grazing on a lush green field



6 3 18

Take 2: Image Captioning



picdescbot @picdescbot · Feb 19

a herd of sheep grazing on a lush green field

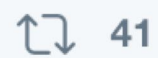




[picdescbot](#) @picdescbot · Mar 8

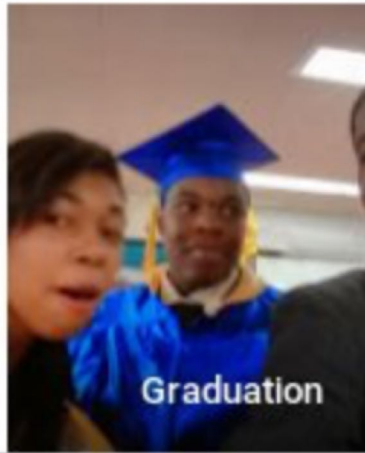
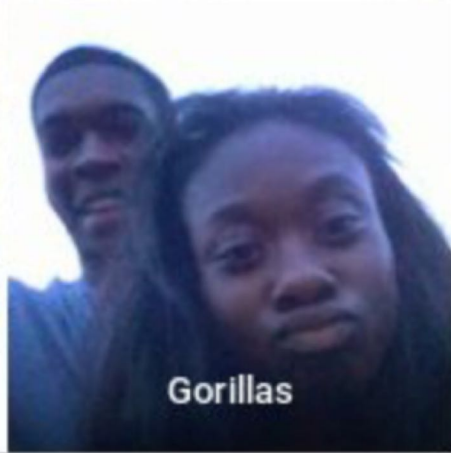
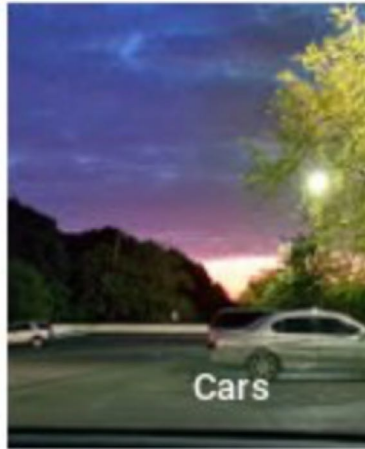
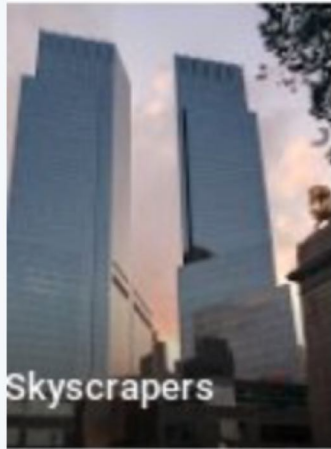


a yellow and orange flowers in a field



Suda yüzmekte olan bir köpek.





jackyalciné is about 40% into the IndieWeb.

@jackyalcine



Google Photos, y'all fucked up. My friend's not a gorilla.

4:22 AM - Jun 29, 2015

♡ 2,280 💬 3,592 people are talking about this

theguardian.com

Support The Guardian

The Guardian

News Opinion Sport Culture Lifestyle

World UK Science Cities Global development

The Google logo, with the letters in their characteristic colors (blue, red, yellow, blue, green, red).

Google says sorry for racist auto-tag in photo app

- Google Photos labelled a picture of two black people as 'gorillas'
- Google Maps and Flickr have also suffered from race-related problems

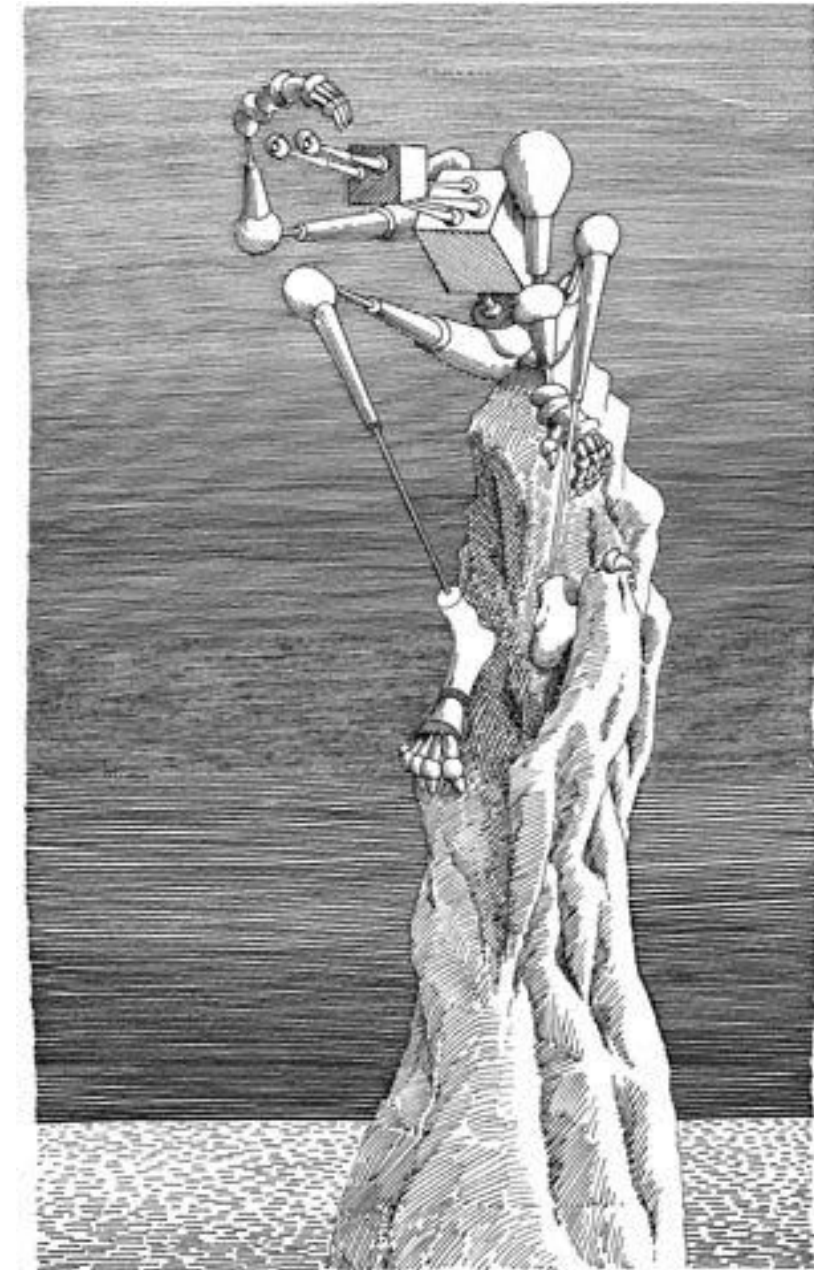
Jana Kasperkevic in New York

@kasperka Email

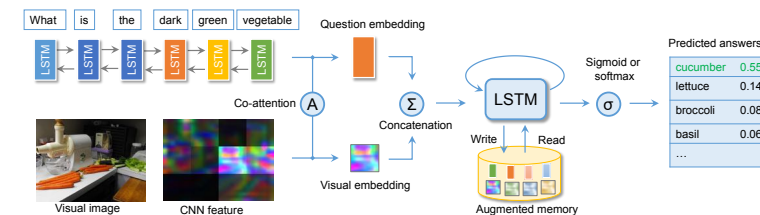
Wed 1 Jul 2015 18.52 BST

Looking Forward

- Intelligence is not just about **Pattern Recognition**
- Learning is the process of **modeling the world...**
 - **explaining** and **understanding** what we see
 - **imagining** things we could see but haven't yet.
 - **problem solving** and **planning** actions to make things real.
 - **building new models** as we learn more about the world.
 - **sharing our models**, communicating to others, understanding their models, and learning from them and with them.



Explaining and understanding what we see



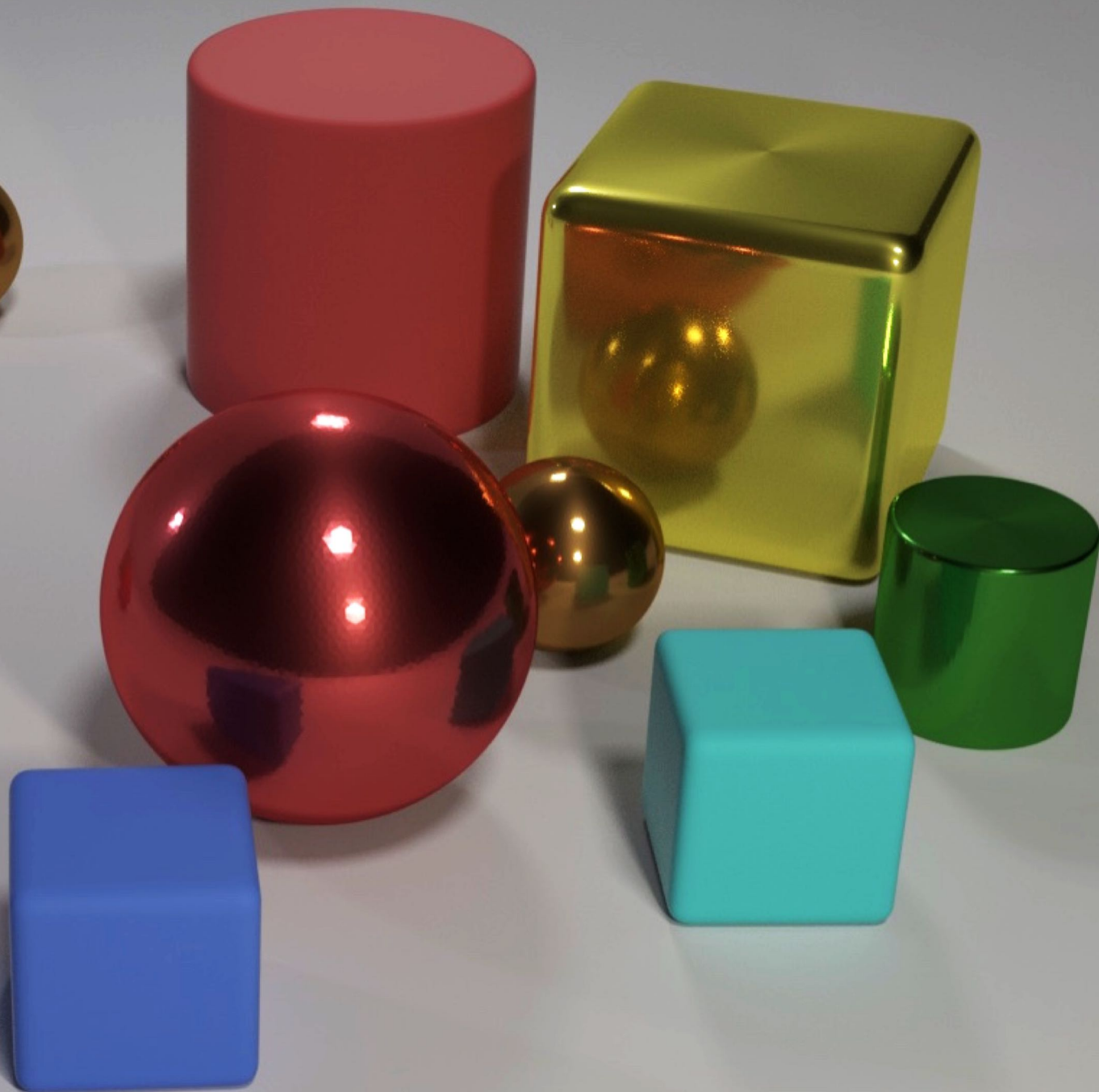
Q: What fruit is showing in this picture?

A: Bananas

Visual Reasoning

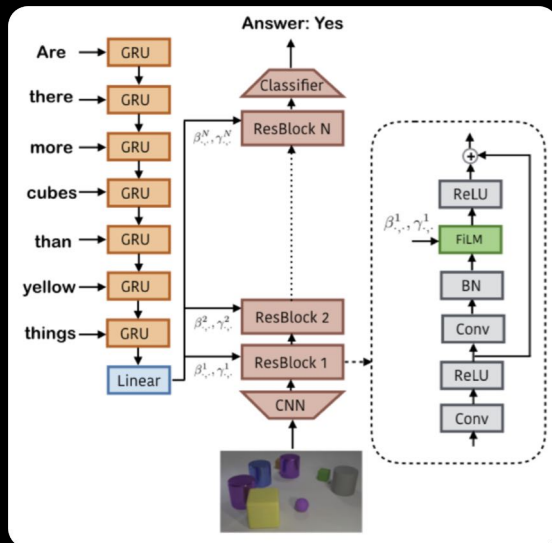
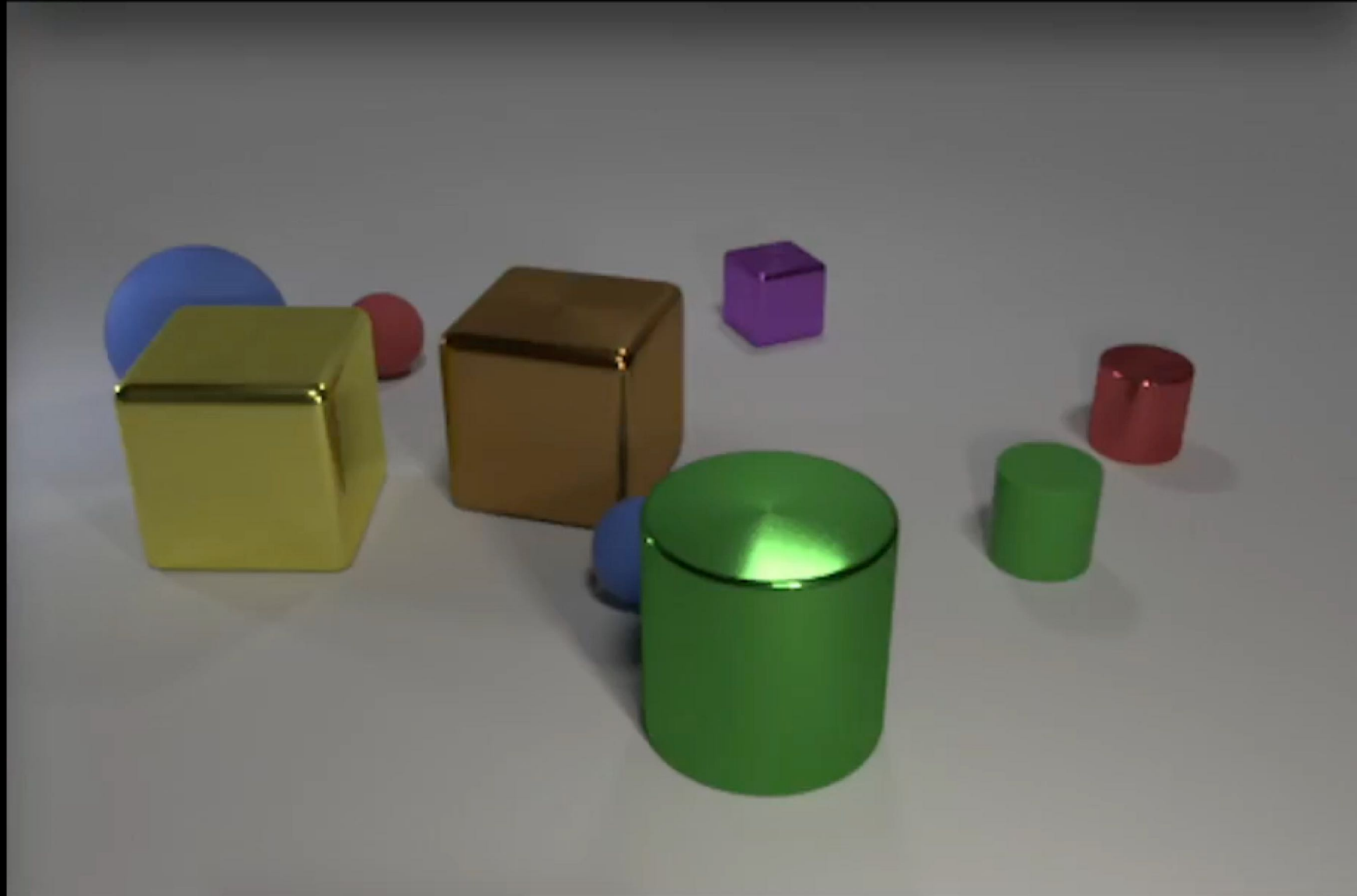
How many objects are
either small cylinders
or red things?

Answer: 5



Visual Reasoning

Ask me something!
>>> H

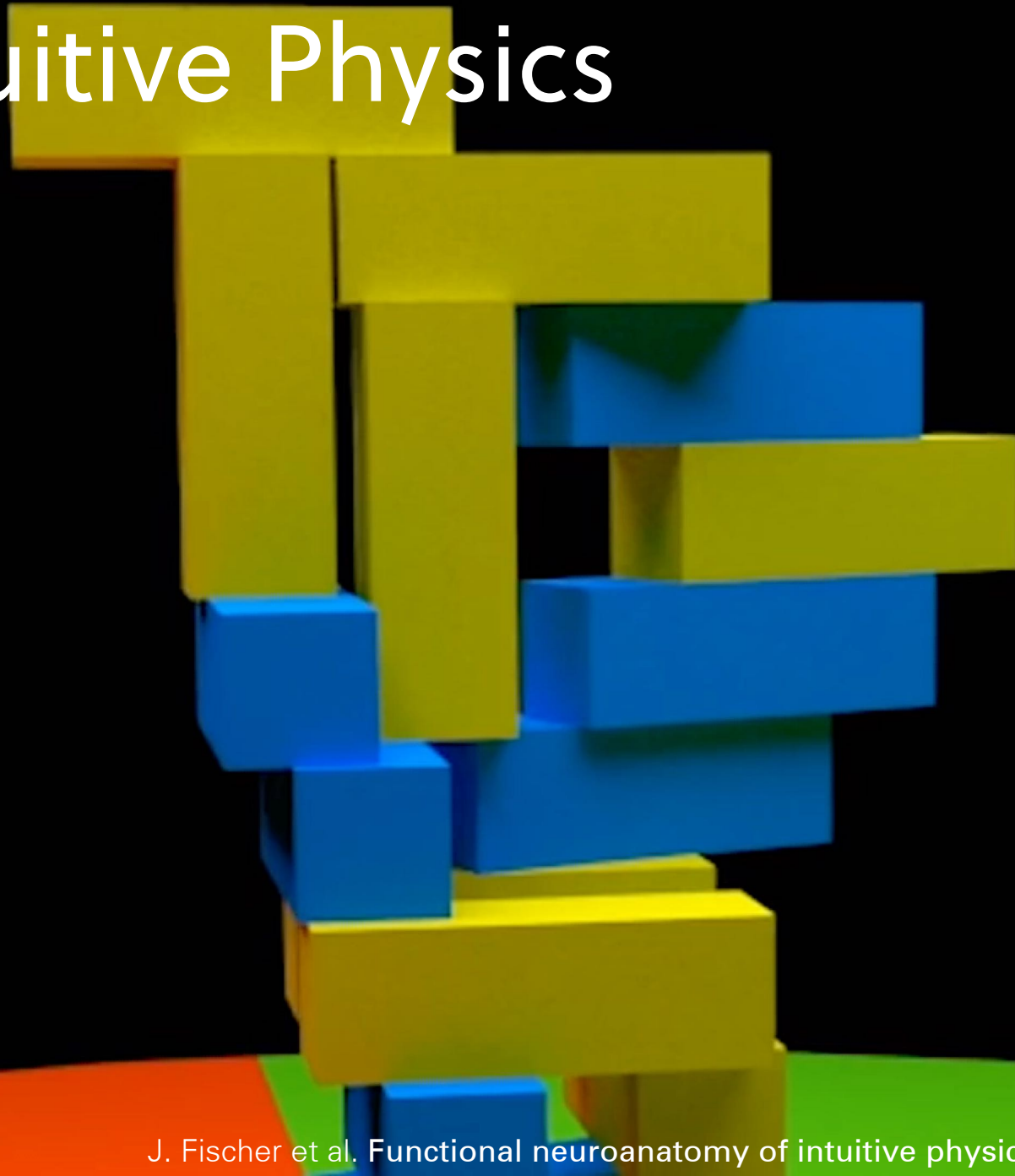


Intuitive Physics

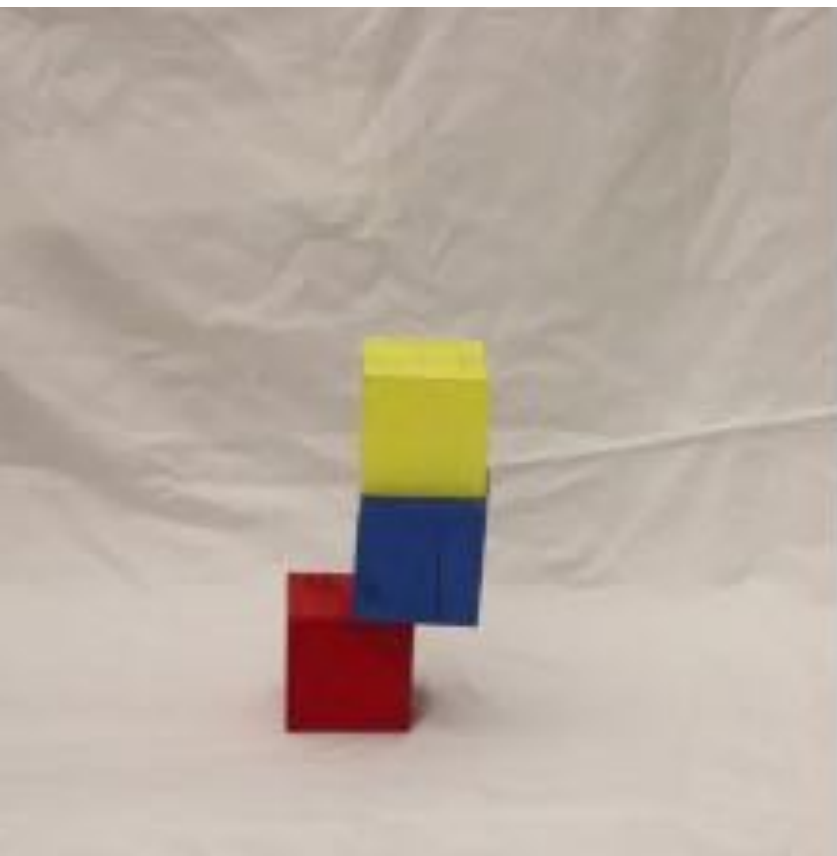
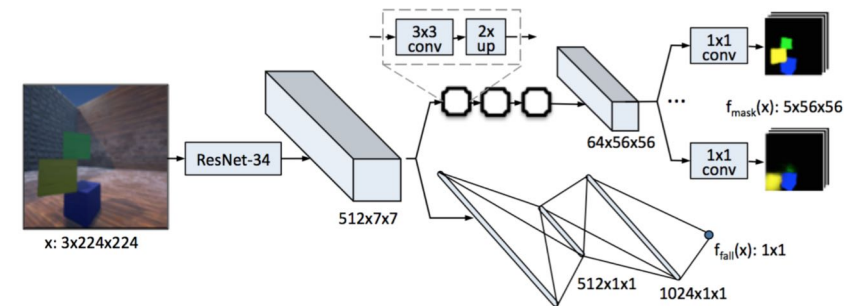
- Common-sense understanding of how the world operates at a physical level
- Helps us to perceive, understand and act with our environment



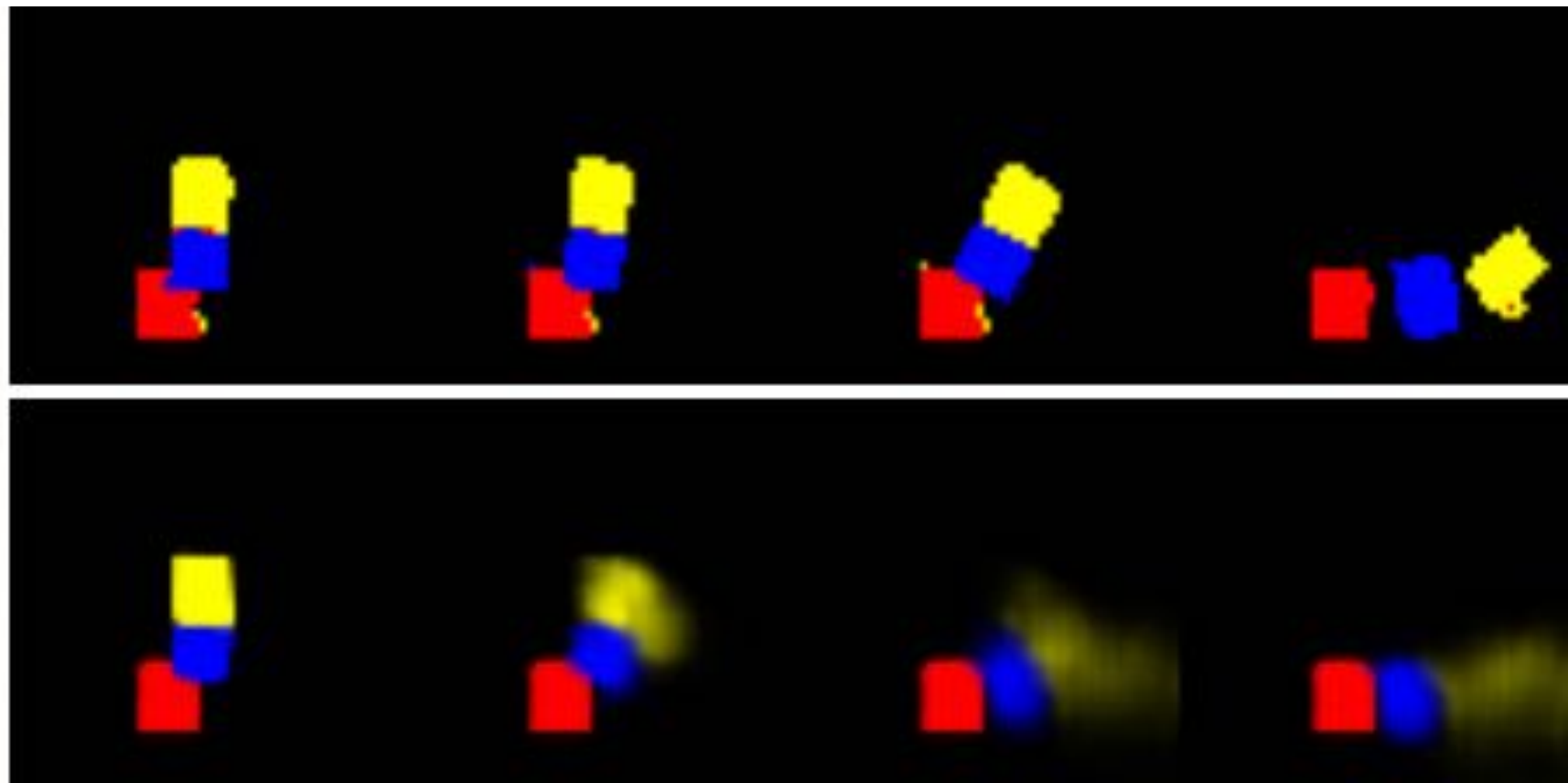
Intuitive Physics



Intuitive Physics



Initial frame



PhysNet predictions of the future

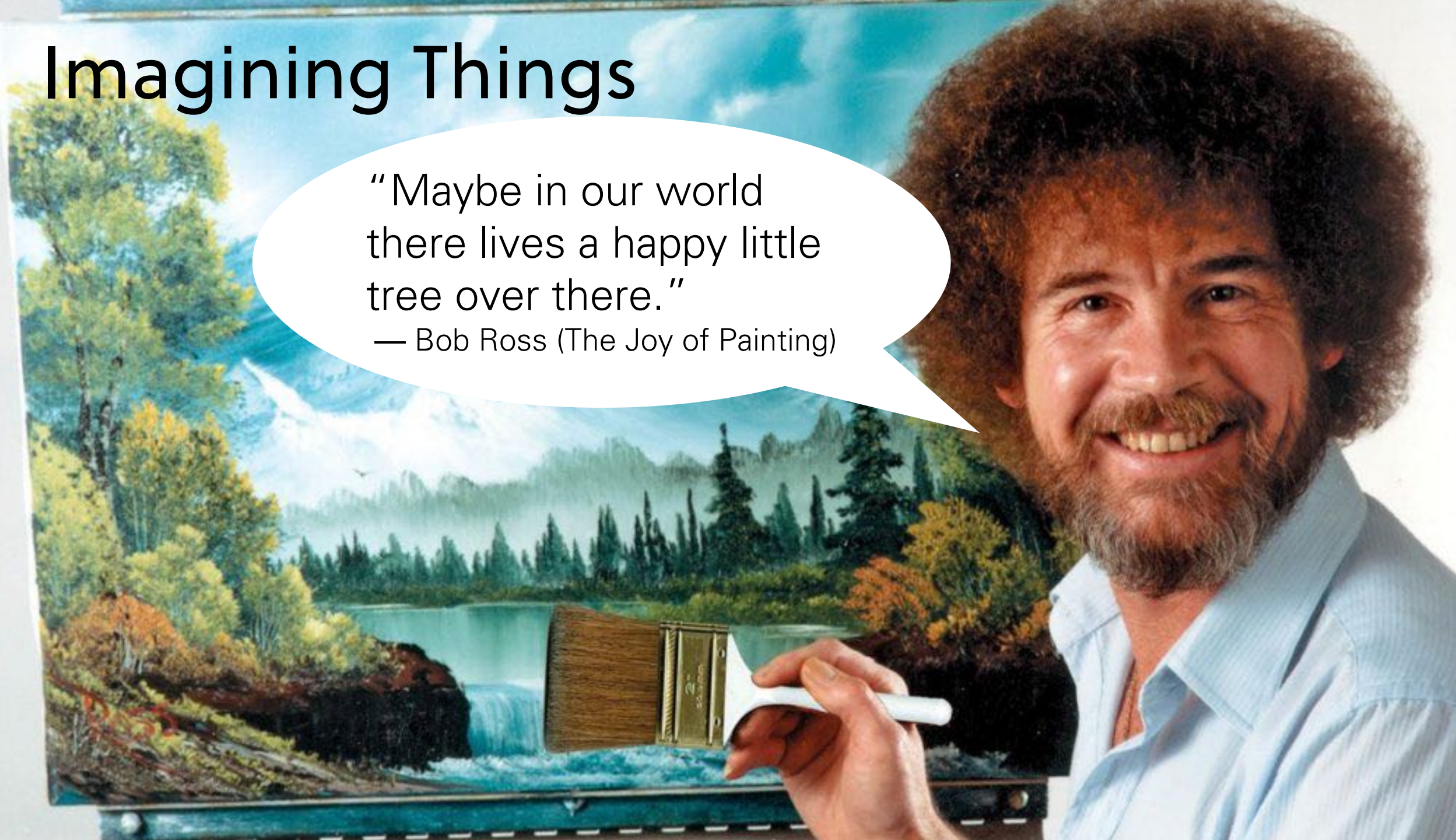


In our world, the future is not deterministic,
there are many possibilities...

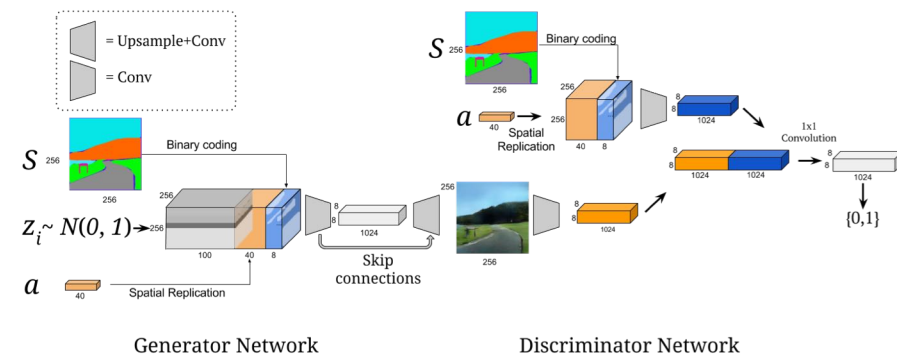
Imagining Things

"Maybe in our world
there lives a happy little
tree over there."

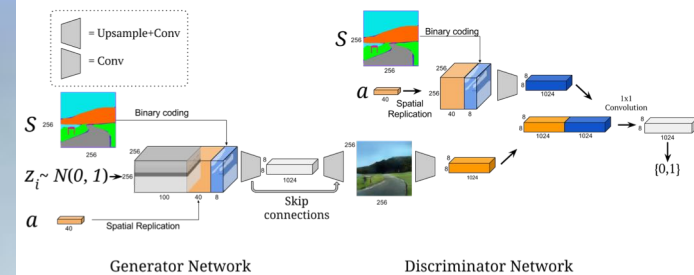
— Bob Ross (The Joy of Painting)



Imagining Things

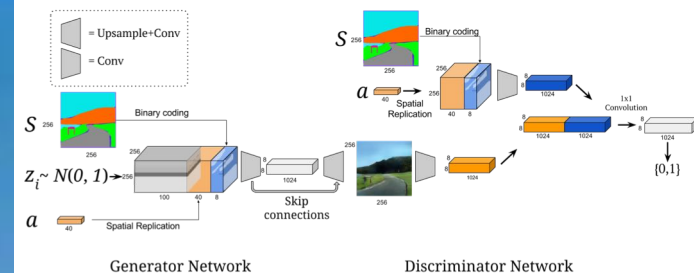


Imagining Things



Semantic Layout

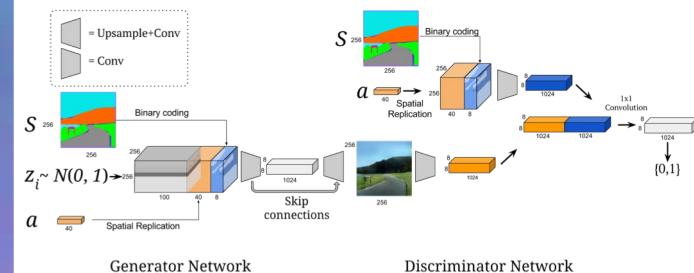
Imagining Things



Semantic Layout

Clear sky + flowers

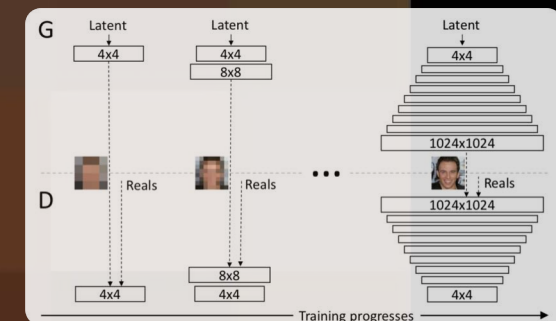
Imagining Things



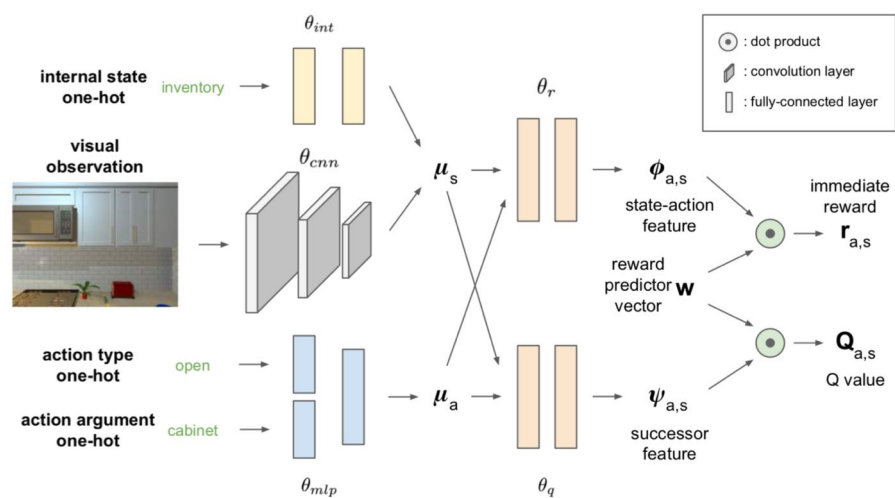
Semantic Layout

Sunset

Imagining Things



Planning



Overhead view of Visual Dynamic Environment



Task: Put bowl in microwave

Initial State

navigate to bowl

1

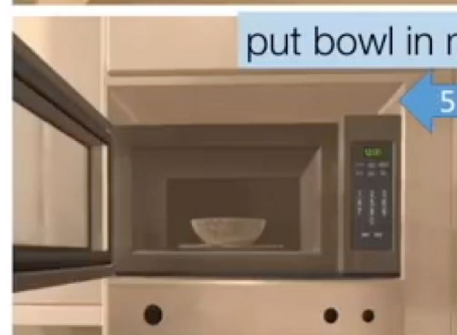


2

pick up bowl

put bowl in microwave

5



navigate to microwave

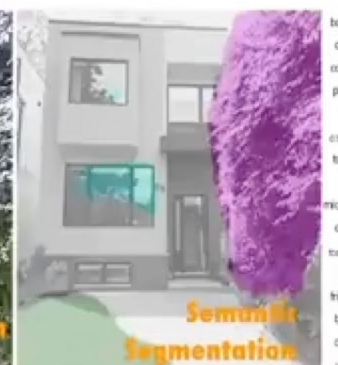
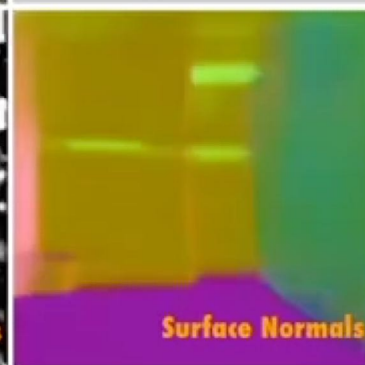
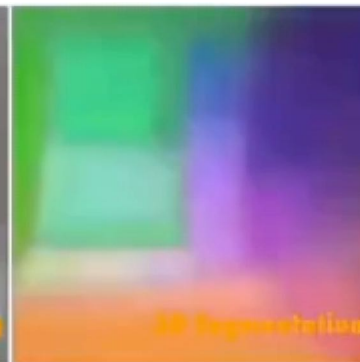
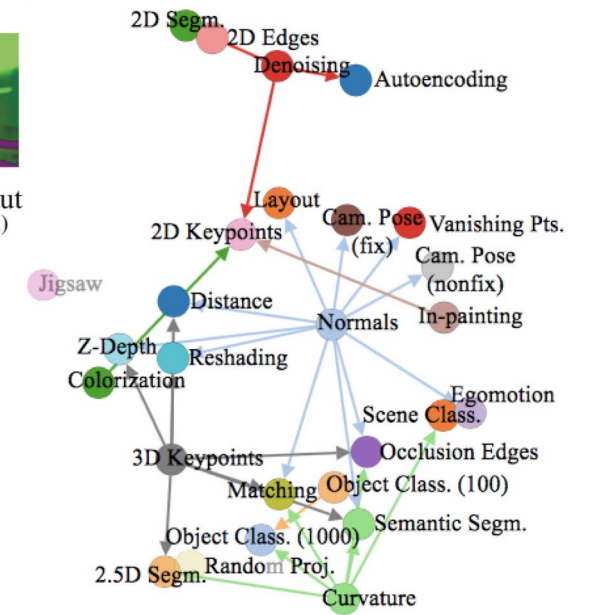
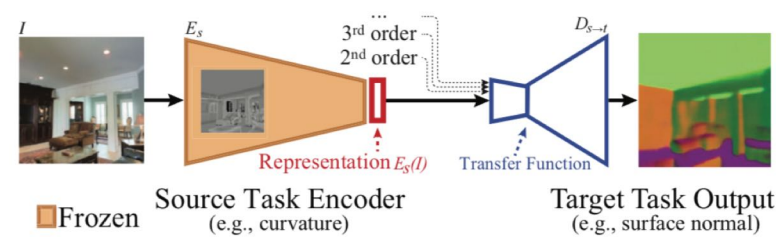
3



open microwave

4

Transfer Learning



Can deep models reason?



- We are not there yet! But we can see real progress soon..