**CS-503
Visual Intelligence:
Machines and Minds**

Amir Zamir
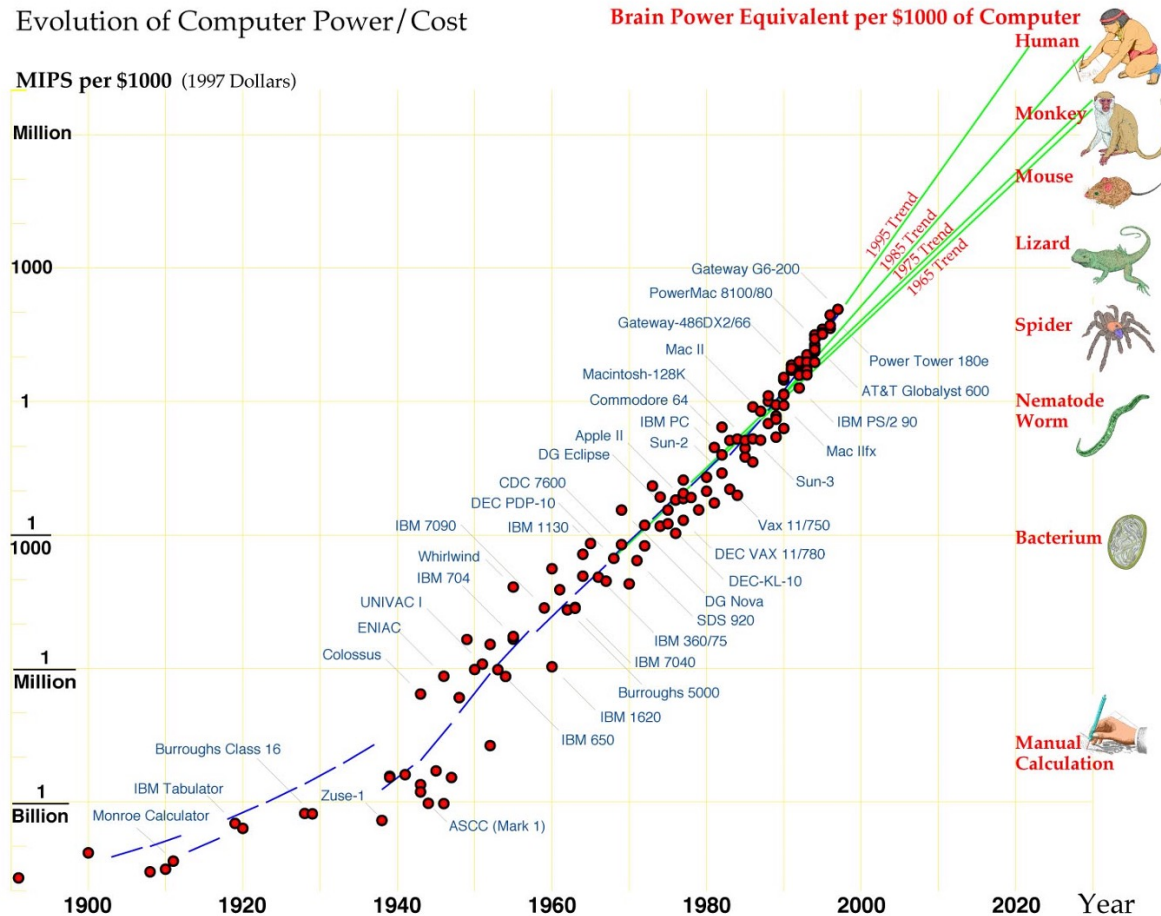
Lecture 2

# Recaps relevant to vision today

- **Data-driven/learning:**
  - e.g., Neural networks, AlexNet, etc.
- **Non-data-driven**
  - e.g., Image formation model, Image transformation, etc.
- "Method" recap, as opposed to historical credit assignment.

Zamir

# Fast Historical Recap

# Historical review

# Bigger picture
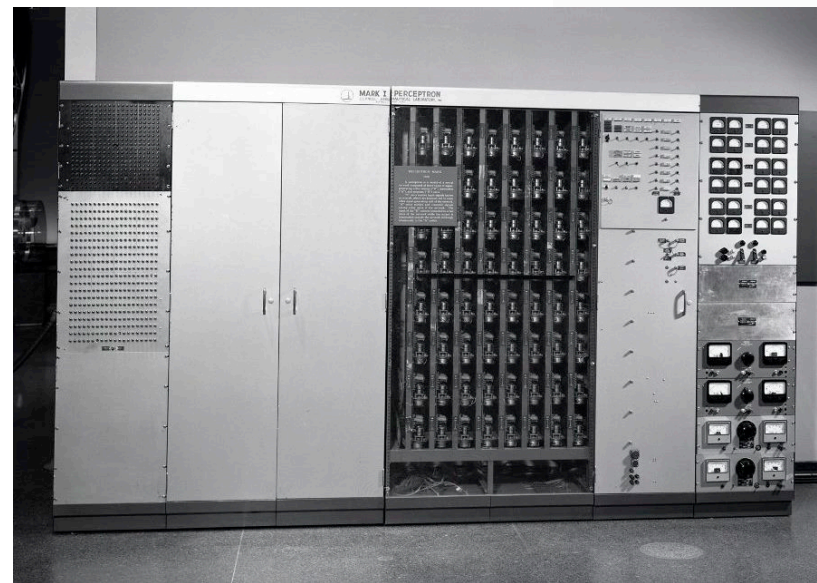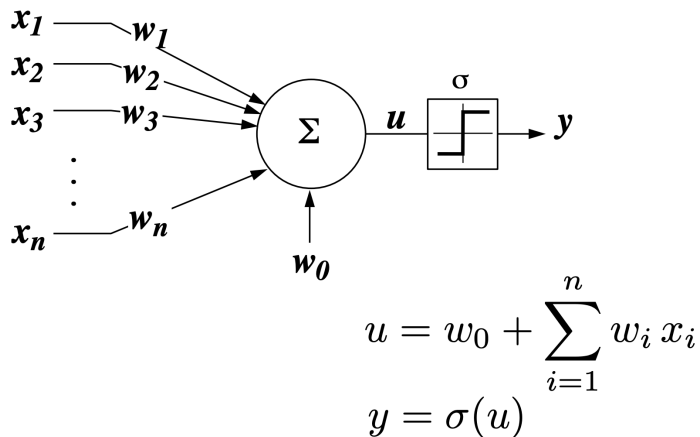


Evolution of Computer Power/Cost

Brain Power Equivalent per $1000 of Computer

Moravec 1998

Zamir

# **Historical review**

- Emission Theory (ca. 400 BC)

Zamir

# **Historical review (neural networks)**

- Perceptron model, Rosenblatt, 1958.





$$u = w_0 + \sum_{i=1}^{n} w_i\, x_i$$

$$y = \sigma(u)$$

Zamir

CS-503: Visual Intelligence: Machines and Minds

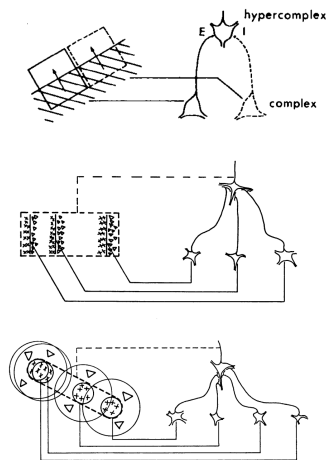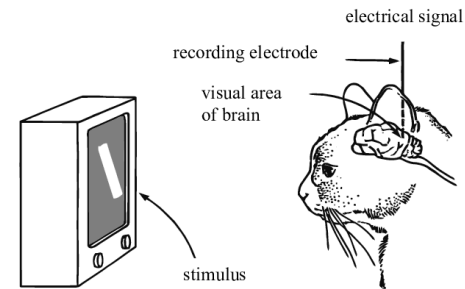# Historical review (neural networks)

Zamir

CS-503: Visual Intelligence: Machines and Minds

- Perceptron model, Rosenblatt, 1958.

# Historical review (neural networks)



- Hubel and Wiesel, ~1962.



Hypercomplex

↑

Complex

↑

Simple

Zamir

# Historical review (neural networks)

Zamir

CS-503: Visual Intelligence: Machines and Minds

- Neocognitron, Fukushima, 1980.
  - Modeled after Hubel and Wiesel.
  - Convolutional.
  - Multi-layer.
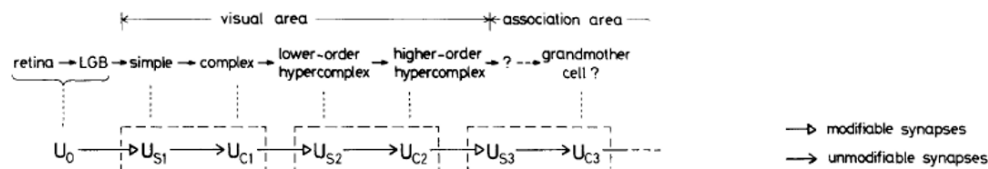  - Hebbian Learning (no backpropagation).



Fig. 1. Correspondence between the hierarchy model by Hubel and Wiesel, and the neural network of the neocognitron
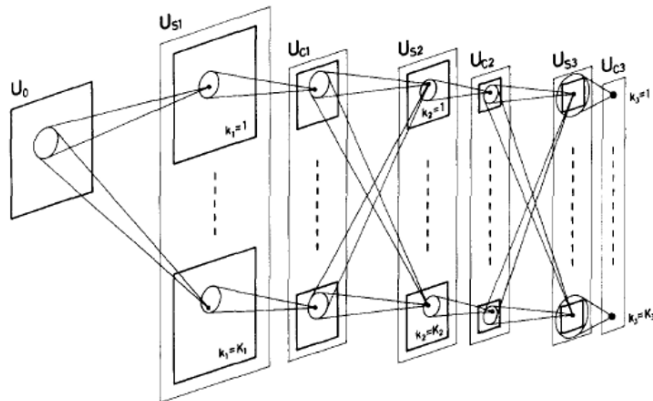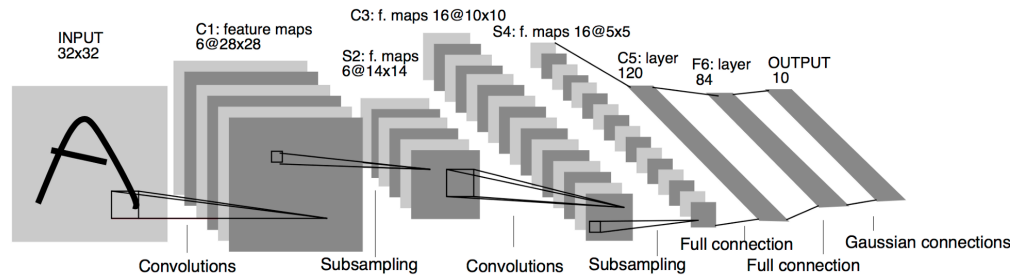
Fig. 2. Schematic diagram illustrating the interconnections between layers in the neocognitron
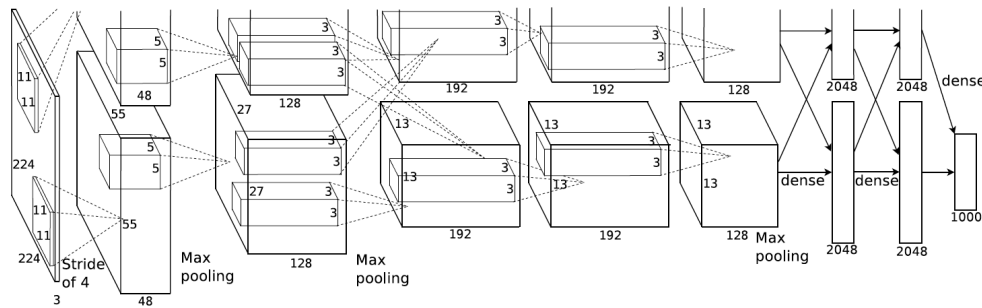
# Historical review (neural networks)

- Backpropagation based learning. LeCun, Hinton, and co ~1989.

  - Backpropagation, as a method, already existed.

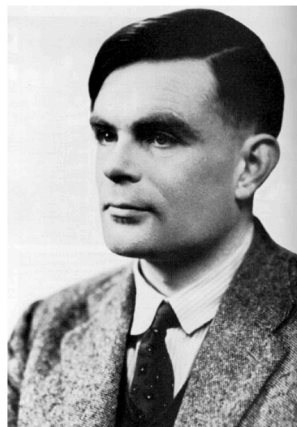  - Most successful for handwritten digit recognition. Didn't work for standard vision tasks.



CS-503: Visual Intelligence: Machines and Minds

Zamir

# Historical review (neural networks)

Zamir

CS-503: Visual Intelligence: Machines and Minds

- AlexNet, Krizhevsky, Sutskever & Hinton. 2012
  - Success at a standard vision task (ImageNet)
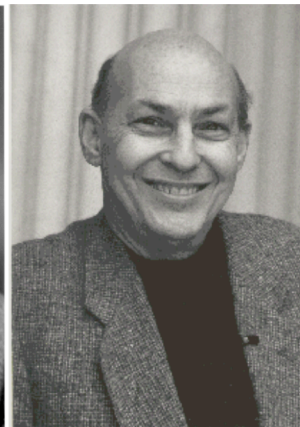  - Deep Learning wave.
  - GPUs

# Historical review

- "AI"



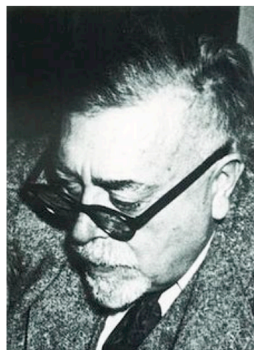Alan Turing    John von Neumann    Marvin Minsky    John McCarthy

Among the most challenging scientific questions of our time are the corresponding analytic and synthetic problems:  How does the brain function?  Can we design a machine which will simulate a brain?
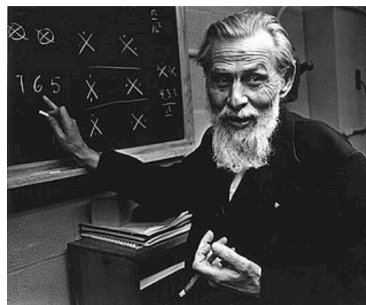-- *Automata Studies*, 1956

B. Olshausen

Zamir

# **Historical review**

- ▪ "Cybernetics"



Norbert Wiener          Warren McCulloch & Walter Pitts          Frank Rosenblatt

"The theory reported here clearly demonstrates the feasibility and fruitfulness of a quantitative statistical approach to the organization of cognitive systems. By the study of systems such as the perceptron, it is hoped that those fundamental laws of organization which are common to all information handling systems, machines and men included, may eventually be understood."  -- Frank Rosenblatt

**The Perceptron:** A Probabilistic Model for Information Storage and Organization in the Brain. In, *Psychological Review*, Vol. 65, No. 6, pp. 386-408, November, 1958.
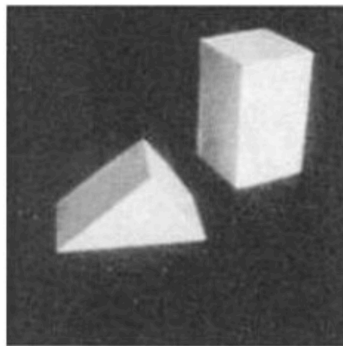
B. Olshausen
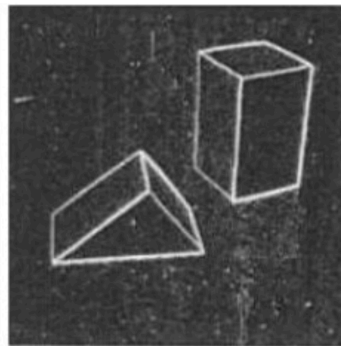
# Historical review (computer vision)

- Larry Roberts Thesis 1963.
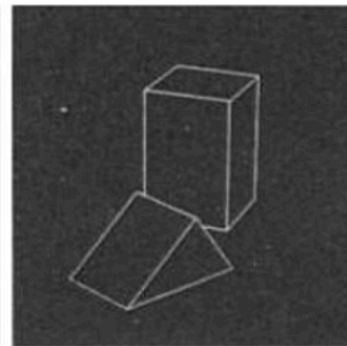  - "Machine Perception of Three-Dimensional Solids"



Larry Roberts
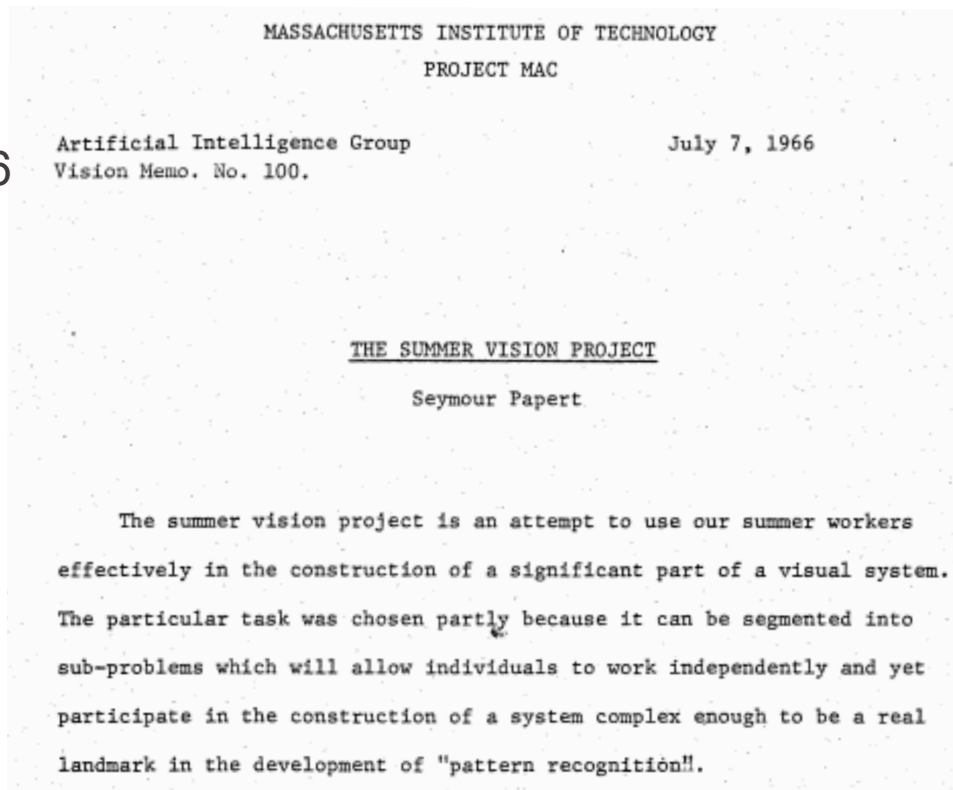"Father of Computer Vision"

Input image

2x2 gradient operator

computed 3D model
rendered from new viewpoint

# Historical review (computer vision)

- Larry Roberts Thesis 1963.
- The Summer Vision Project 1966

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

PROJECT MAC

Artificial Intelligence Group                    July 7, 1966
Vision Memo. No. 100.

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".

# Historical review (computer vision)

- **1960s**: Birth
  - Larry Roberts Thesis 1963.
  - The Summer Vision Project 1966
- **1970s**: Foundational work on image formation
- **1980s**: Applied mathematics: geometry, multi-scale analyses
- **1990s**: Geometric analysis. Vision+graphics. Resurfacing of statistical learning.
- **2000s**: Progress in visual recognition. Pre-deep learning. DPM. PASCAL.
- **2010s**: Deep Learning.

J. Malik

Zamir

CS-503: Visual Intelligence: Machines and Minds

# Historical review AI

A good recap:

https://youtu.be/R3YFxF0n8n8

DOCUMENTARY

THE THINKING MACHINE
A HISTORY OF AI

# Fast Computer Vision Recap

# Vision

- defined as the process of acquiring knowledge about environmental objects and events by extracting information from the light they emit or reflect. [S. Palmer]

    - optical process: image formation

    - cognitive process: phenomena of visual perception

    - physiology of the visual nervous system: biological implementation



S. Palmer

# One useful organization

- Three Rs of Computer Vision.
- Interaction between the Rs.

Recognition

Reconstruction

Reorganization

The three R's of computer vision: Recognition, reconstruction and reorganization. J. Malik et al. 2016.

# Rudimentary Image Formation

# Rudimentary Image Formation

Environmental Object (Distal Stimulus) — Camera — Projected Image (Proximal Stimulus)

A. Pinhole Aperture without Lens --> Sharp Image

B. Large Aperture without Lens --> Fuzzy Image

C. Large Aperture with Lens --> Sharp Image



2-D image | Image plane | Focal point | Virtual image plane | 3-D object

Focal length

# Rudimentary Image Formation

CS-503: Visual Intelligence: Machines and Minds



A. Pinhole Aperture without Lens --> Sharp Image

B. Large Aperture without Lens --> Fuzzy Image

C. Large Aperture with Lens --> Sharp Image

Environmental Object (Distal Stimulus) — Camera — Projected Image (Proximal Stimulus)



(a) The large pinhole eye of the cephalopod mollusc Nautilus. (b) Corneal eyes of a jumping spider Platycryptus. (c) Concave mirror eyes of the scallop Pecten. (d) Primitive compound eye of the ark clam Barbatia. (e) Compound eye of a male robberfly Holocephala. (f) Mirror compound eye of the shrimp Palaemonetes.



A. Pigment cup eye

B. Lens eye

C. Corneal eye

D. Concave mirror eye

A. Basic compound eye

B. Apposition eye

D. Refracting superposition eye

D. Reflecting superposition eye

# Image Transformation

Zamir



Image registration.

# Image Transformation

Zamir

| Transformation | Matrix | # DoF | Preserves | Icon |
|---|---|---|---|---|
| translation | $\left[\, I \mid t \,\right]_{2\times3}$ | 2 | orientation | |
| rigid (Euclidean) | $\left[\, R \mid t \,\right]_{2\times3}$ | 3 | lengths | |
| similarity | $\left[\, sR \mid t \,\right]_{2\times3}$ | 4 | angles | |
| affine | $\left[\, A \,\right]_{2\times3}$ | 6 | parallelism | |
| projective | $\left[\, \tilde{H} \,\right]_{3\times3}$ | 8 | straight lines | |



| Transformation | Before | After |
|---|---|---|
| Projective | | |
| Affine | | |
| Similarity | | |
| Euclidean | | |

# Image Transformation

CS-503: Visual Intelligence: Machines and Minds

| Transformation | Matrix | # DoF | Preserves | Icon |
|---|---|---|---|---|
| translation | $\left[\begin{array}{c|c} I & t \end{array}\right]_{2\times3}$ | 2 | orientation | |
| rigid (Euclidean) | $\left[\begin{array}{c|c} R & t \end{array}\right]_{2\times3}$ | 3 | lengths | |
| similarity | $\left[\begin{array}{c|c} sR & t \end{array}\right]_{2\times3}$ | 4 | angles | |
| affine | $\left[\begin{array}{c} A \end{array}\right]_{2\times3}$ | 6 | parallelism | |
| projective | $\left[\begin{array}{c} \tilde{H} \end{array}\right]_{3\times3}$ | 8 | straight lines | |

projective

affine

Euclidean

Identity

A Square

Rotation  Shearing  Translation  Scaling

# Correspondences

# Correspondences

CS-503: Visual Intelligence: Machines and Minds



Figure 2. The DAISY descriptor. Each circle represents a region where the radius is proportional to the standard deviations of the Gaussian kernels and the '+' sign represents the locations where we sample the convolved orientation maps center being a pixel location where we compute the descriptor. By overlapping the regions we achieve smooth transitions between the regions and a degree of rotational robustness. The radii of the outer regions are increased to have an equal sampling of the rotational axis which is necessary for robustness against rotation.



Image gradients

Keypoint descriptor

SIFT, Lowe et al. 2003
DAISY, Tola et al. 2008

# Correspondences

Zamir

CS-503: Visual Intelligence: Machines and Minds

- Now gone deep.



e.g. Zagoruyko, S., Komodakis, N.: Learning to compare image patches via convolutional neural networks (2015).

# Dynamic Perspective

- Motion. Video.

# **Dynamic Perspective**

Zamir

- Motion. Video.

Eadweard Muybridge

# Dynamic Perspective

Zamir

- Motion. Video.

Eadweard Muybridge



CS-503: Visual Intelligence: Machines and Minds

# Dynamic Perspective

- Optical Flow.
- J J Gibson's examples:



(a)

(b)

Zamir

# Dynamic Perspective

- Optical Flow

# Dynamic Perspective

Zamir

CS-503: Visual Intelligence: Machines and Minds

- Recognition in Videos
  - E.g. actions

Images

Video

Image classification

"Is there a dog in the image?"

Object detection

"Is there a dog and where is it in the image?"

Action classification

"Is there a person diving in the video?"

Action detection

"Is there a person diving and where is it in the video?"

# Dynamic Perspective→3D

- Moving (or multi) Camera

## Binocular Stereopsis



Perceived Object

Left Image

Right Image

Left Eye

Right Eye

# 3D

Zamir

CS-503: Visual Intelligence: Machines and Minds

- Moving (or multi) Camera

## Depth from Triangulation

Passive Stereopsis

Camera 1    Camera 2

Active Stereopsis

Camera    Projector

# 3D

Zamir

CS-503: Visual Intelligence: Machines and Minds

## Depth from Triangulation



Camera 1          Camera 2          Camera          Projector

Passive Stereopsis          Active Stereopsis

IR speckle pattern
projector

RGB

Monochrome

XBOX 360

Projected speckle pattern

# 3D

Zamir

CS-503: Visual Intelligence: Machines and Minds

- Moving (or multi) Camera



Building Rome in a Day. Snavely et al.

Zamir

# Other 3D estimation cues

- e.g. Pictorial cues



CS-503: Visual Intelligence: Machines and Minds

# List of 3D cues
# (Vision Science, Palmer)

| INFORMATION SOURCE | Ocular/ Optical | Binocular/ Monocular | Static/ Dynamic | Relative/ Absolute | Qualitative/ Quantitative |
|---|---|---|---|---|---|
| Accommodation | ocular | monocular | static | absolute | quantitative |
| Convergence | ocular | binocular | static | absolute | quantitative |
| Binocular Disparity | optical | binocular | static | relative | quantitative |
| Motion Parallax | optical | monocular | dynamic | relative | quantitative |
| Texture Accretion/Deletion | optical | monocular | dyanmic | relative | qualitative |
| Convergence of Parallels | optical | monocular | static | relative | quantitative |
| Position relative to Horizon | optical | monocular | static | relative | quantitative |
| Relative Size | optical | monocular | static | relative | quantitative |
| Familiar Size | optical | monocular | static | absolute | quantitative |
| Texture Gradients | optical | monocular | static | relative | quantitative |
| Edge Interpretation | optical | monocular | static | relative | qualitative |
| Shading and Shadows | optical | monocular | static | relative | qualitative |
| Aerial Perspective | optical | monocular | static | relative | qualitative |

J. Malik

# Binocular Stereopsis

J. Malik

# Accommodation

Thick Lens → Close

Thin Lens → Far

J. Malik

# Convergence

Large Angle => Close

Small Angle => Far

$$d = \frac{c}{2 \tan (a/2)}$$

Convergence Angle in degrees

Distance in meters (d)

J. Malik

# Shading

J. Malik

# Recognition

Zamir



| **Semantic Segmentation** | **Classification + Localization** | **Object Detection** | **Instance Segmentation** |
|---|---|---|---|
| GRASS, CAT, TREE, SKY | CAT | DOG, DOG, CAT | DOG, DOG, CAT |
| No objects, just pixels | Single Object | Multiple Object | |

This image is CC0 public domain

# Recognition

Zamir

warped region

aeroplane? no.
⋮
person? yes.
⋮
tvmonitor? no.

CNN

figure credit: R. Girshick et al.

input image | region proposals ~2,000 | **1 CNN for each region** | classify regions

## R-CNN pipeline

R. Girshick, J. Donahue, T. Darrell, & J. Malik. "Rich feature hierarchies for accurate object detection and semantic segmentation". CVPR 2014

# Embodied/Active Vision

CS-503: Visual Intelligence: Machines and Minds

# Embodied/Active Vision

# "Foundation" models
# Vision-Language Models
# Multimodal Models



Flamingo,, Alayrac et al, 2022.
LLaVA, Liu et al, 2023
4M, Mizrahi & Bachmann et al. 2023.

# "Foundation" models
# Vision-Language Models
# Multimodal Models



Figure 3: **Flamingo architecture overview.** Flamingo is a family of visual language models (VLMs) that take as input visual data interleaved with text and produce free-form text as output.

Flamingo,, Alayrac et al, 2022.
LLaVA, Liu et al, 2023
4M, Mizrahi & Bachmann et al. 2023.

CS-503: Visual Intelligence: Machines and Minds

Zamir

# "Foundation" models
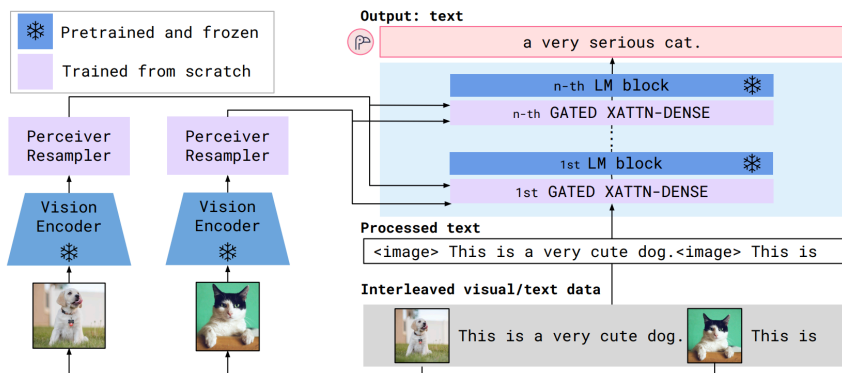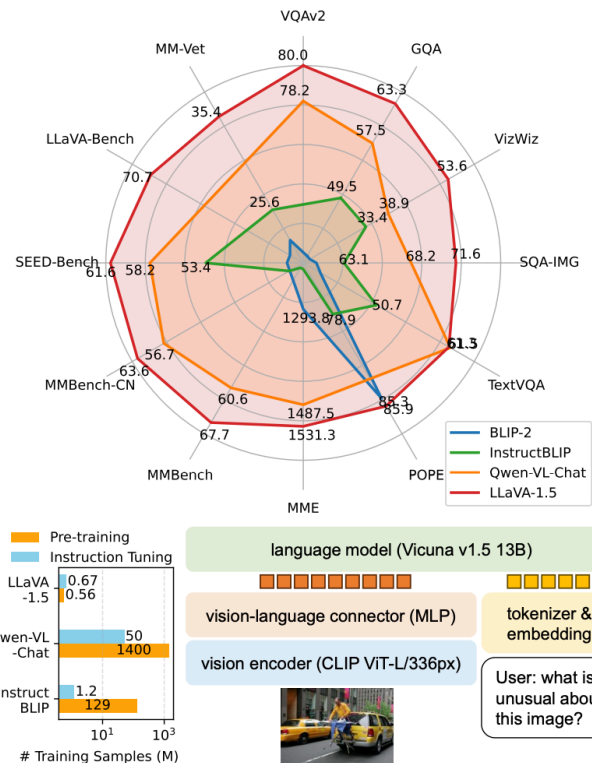# Vision-Language Models
# Multimodal Models



Figure 3: **Flamingo architecture overview.** Flamingo is a family of visual language models (VLMs) that take as input visual data interleaved with text and produce free-form text as output.

Flamingo,, Alayrac et al, 2022.
LLaVA, Liu et al, 2023
4M, Mizrahi & Bachmann et al. 2023.

CS-503: Visual Intelligence: Machines and Minds

Zamir

# "Foundation" models
# Vision-Language Models
# Multimodal Models



Flamingo,, Alayrac et al, 2022.
LLaVA, Liu et al, 2023
4M, Mizrahi & Bachmann et al. 2023.

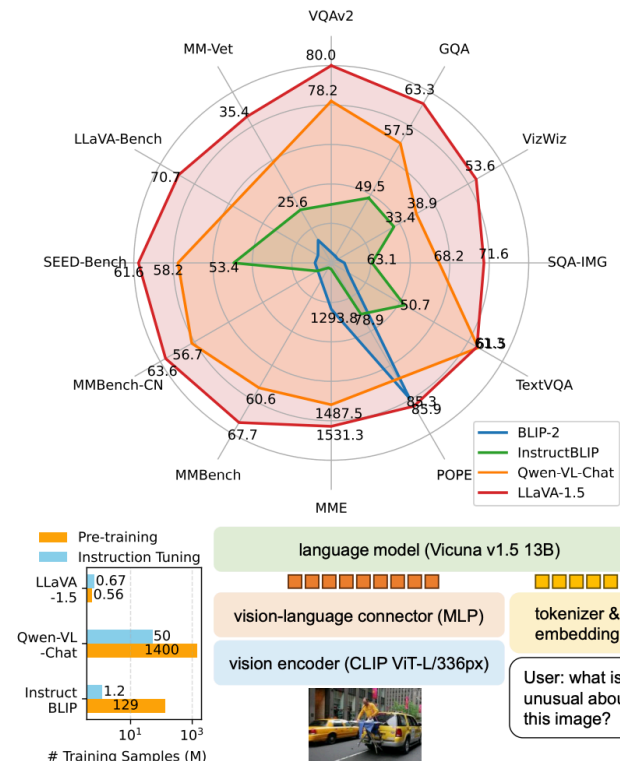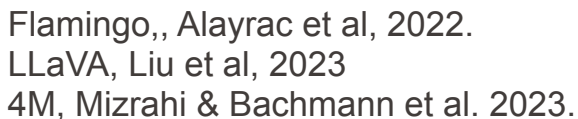CS-503: Visual Intelligence: Machines and Minds

Zamir

# Image Processing Perspective

- e.g. White Balancing





- Manual
  - Choose color-neutral object in the photos and normalize
- Automatic (AWB)
  - Grey World:  force average color of scene to grey
  - White World: force brightest object to white

# Image Processing Perspective

- Point Processing



a b
c d

**FIGURE 3.9**
(a) Aerial image.
(b)–(d) Results of
applying the
transformation in
Eq. (3.2-3) with
$c = 1$ and
$\gamma = 3.0, 4.0,$ and
$5.0$, respectively.
(Original image
for this example
courtesy of
NASA.)

The simplest kind of range transformations are these independent of position x,y:

$$g = T(f)$$

This is called point processing.

e.g. Gain and Bias transform:

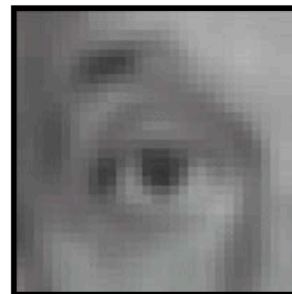g(x,y) = a*f(x,y) + b

A. Efros.

# Image Processing Perspective

Original

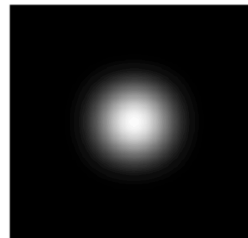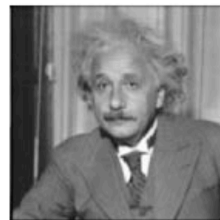$$\frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} =$$
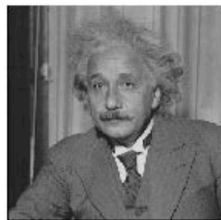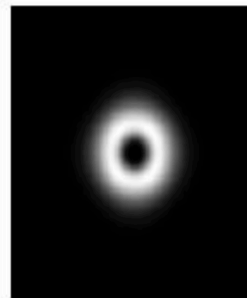
Blur (with a mean filter)

A. Efros.

# Image Processing Perspective

Zamir

CS-503: Visual Intelligence: Machines and Minds
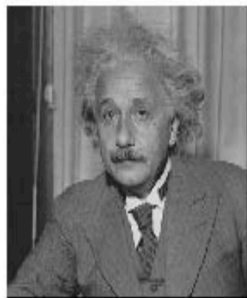
- e.g. Frequency Analyses

low-pass:



High-pass / band-pass:

Zamir

# Not in the recap (In SOTA lectures)

- Neural Network architecture
- Embodied vision simulators and active agent training
- Foundation models: language, multimodal, generative, etc. models.
- Generalization and Robustness

CS-503: Visual Intelligence: Machines and Minds

# Some left out elephants

- Plenoptic function ->
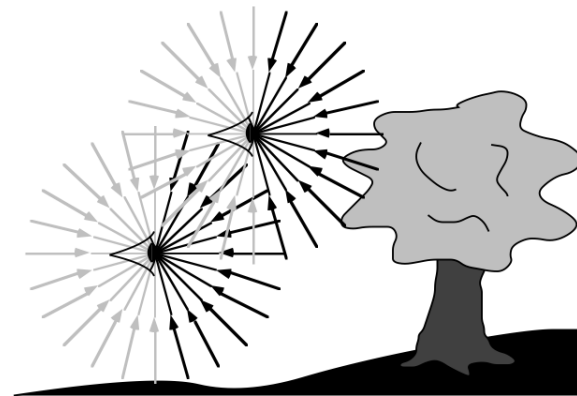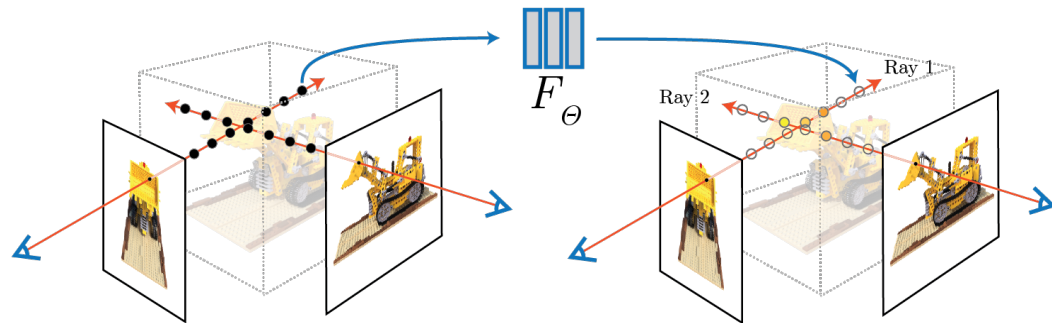  Radiance Fields ->
  Neural Radiance Fields



**Fig.1.3**
The plenoptic function describes the information available to an observer at any point in space and time. Shown here are two schematic eyes-which one should consider to have punctate pupils-gathering pencils of light rays. A real observer cannot see the light rays coming from behind, but the plenoptic function does include these rays.



$F_\Theta$

Ray 1

Ray 2

Zamir

CS-503: Visual Intelligence: Machines and Minds

# Some left out elephants

- Plenoptic function -> Radiance Fields -> Neural Radiance Fields

- 3D Gaussain Splatting



Fig.1.3
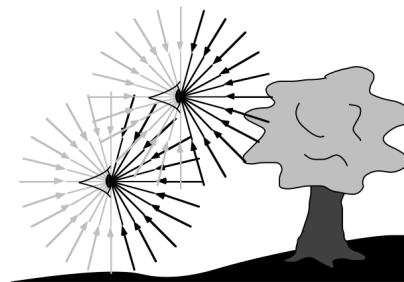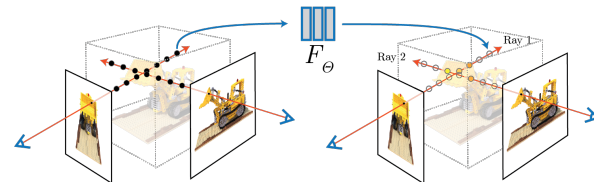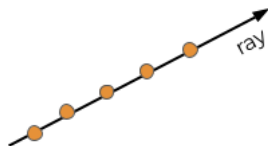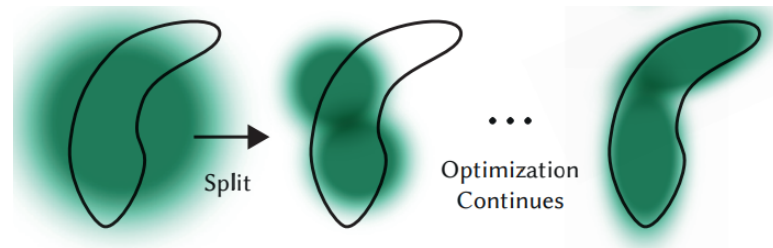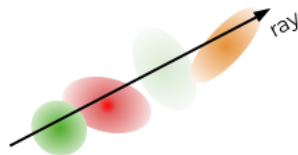The plenoptic function describes the information available to an observer at any point in space and time. Shown here are two schematic eyes-which one should consider to have punctate pupils-gathering pencils of light rays. A real observer cannot see the light rays coming from behind, but the plenoptic function does include these rays.
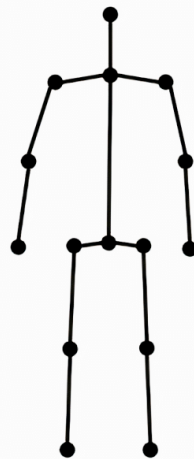




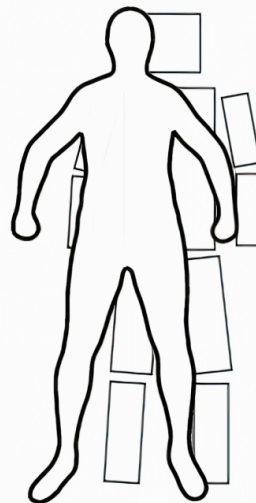Split     Optimization Continues

**NeRF**     **Gaussian Splatting**



Zamir

CS-503: Visual Intelligence: Machines and Minds
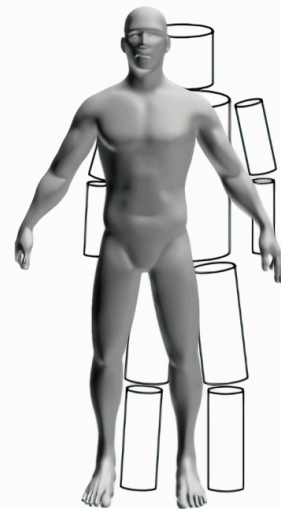
# other (more specific) problems

**HUMAN BODY MODELS**



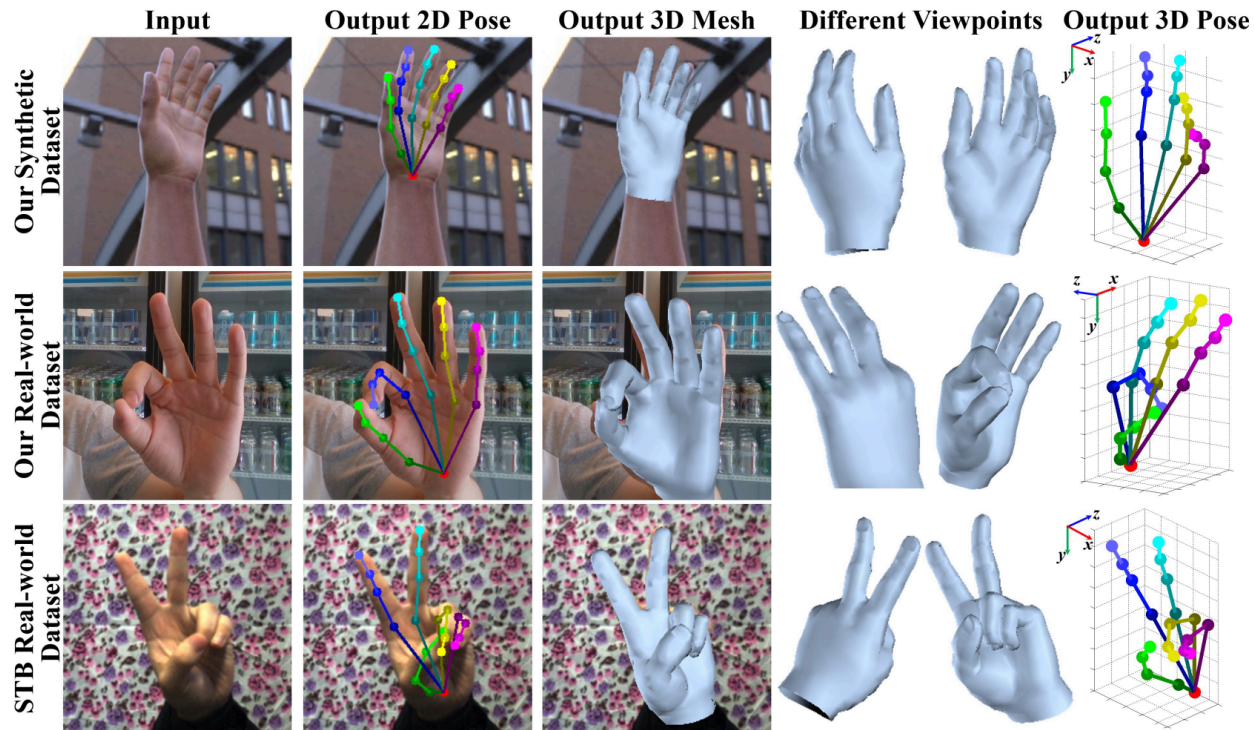skeleton-based
model

contour-based
model

volume-based
model

# other (more specific) problems

Zamir

Viola & Jones

# other (more specific) problems

Zamir



|  | Input | Output 2D Pose | Output 3D Mesh | Different Viewpoints | Output 3D Pose |
|---|---|---|---|---|---|
| Our Synthetic Dataset | | | | | |
| Our Real-world Dataset | | | | | |
| STB Real-world Dataset | | | | | |

Ge et al, 2019

Zamir

# Questions?

https://vilab.epfl.ch/