

Computer Vision Basics

Dr. Sander Ali Khawaja,

Assistant Professor, Department of Telecommunication Engineering
Faculty of Engineering and Technology, University of Sindh, Pakistan

Senior Member, IEEE – Member, ACM

<https://sander-ali.github.io>

Computer Vision & Image Processing



Computer Vision in the News

The Biggest News in Computer Vision in 2022

By DataGen, December 19, 2022

SHARE f t in

2022 has been an incredibly eventful year for the world of computer vision and synthetic data. Here is a roundup of the most exciting developments from this year.

Generative AI models steal the spotlight

These generative models have impressed the world with their ability to generate high-resolution and high-fidelity images from text prompts. 2022 saw the release of multiple powerful text-to-image models that outperformed their predecessors.

Mid-2022, OpenAI released [DALL-E 2](#). It can make realistic edits to existing images, replicate the style of other images, and even expand a single image. [Midjourney](#), another text-to-image model, released its open beta in 2022. Google's Imagen was also released this year.

Stability AI released its brainchild Stable Diffusion 1.0 in 2022, and even updated it to Stable Diffusion 2.0 in November. Based on a novel latent diffusion model, [Stable Diffusion](#) performs competitively with the state-of-the-art for various image generation tasks with a smaller computational footprint. Its developers made the code and model public, a marked departure from its competitors DALL-E 2 and Midjourney.



"A photograph of an astronaut riding a horse" by DALL-E 2 (left) and Stable Diffusion (right)

<https://datagen.tech/blog/2022-computer-vision-year-in-review/>

2

BioTech & Health

Guide dogs will be giving the side-eye to self-driving car tech coming for their jobs

Haje Jan Kamps @Haje / 6:39 PM GMT+5 • January 7, 2023

Comment



Image Credits: Guidi

The visually impaired are getting a helping hand (or a helping belt, as it were) from Korean startup AI Guided. At CES in Las Vegas, the company was showing off some pretty neat tech that incorporates optical and lidar technology along with AI-powered on-device computing to identify obstacles and help with navigation.

The company claims to be able to do advanced object identification to help keep walkers safe, in addition to using gentle haptic feedback to help with wayfinding. The whole system is carried on a belt, leaving the user's hands free.



AI-Guided's upcoming product Guidi, shown off at CES 2023. Image Credit: Haje Kamps/TechCrunch

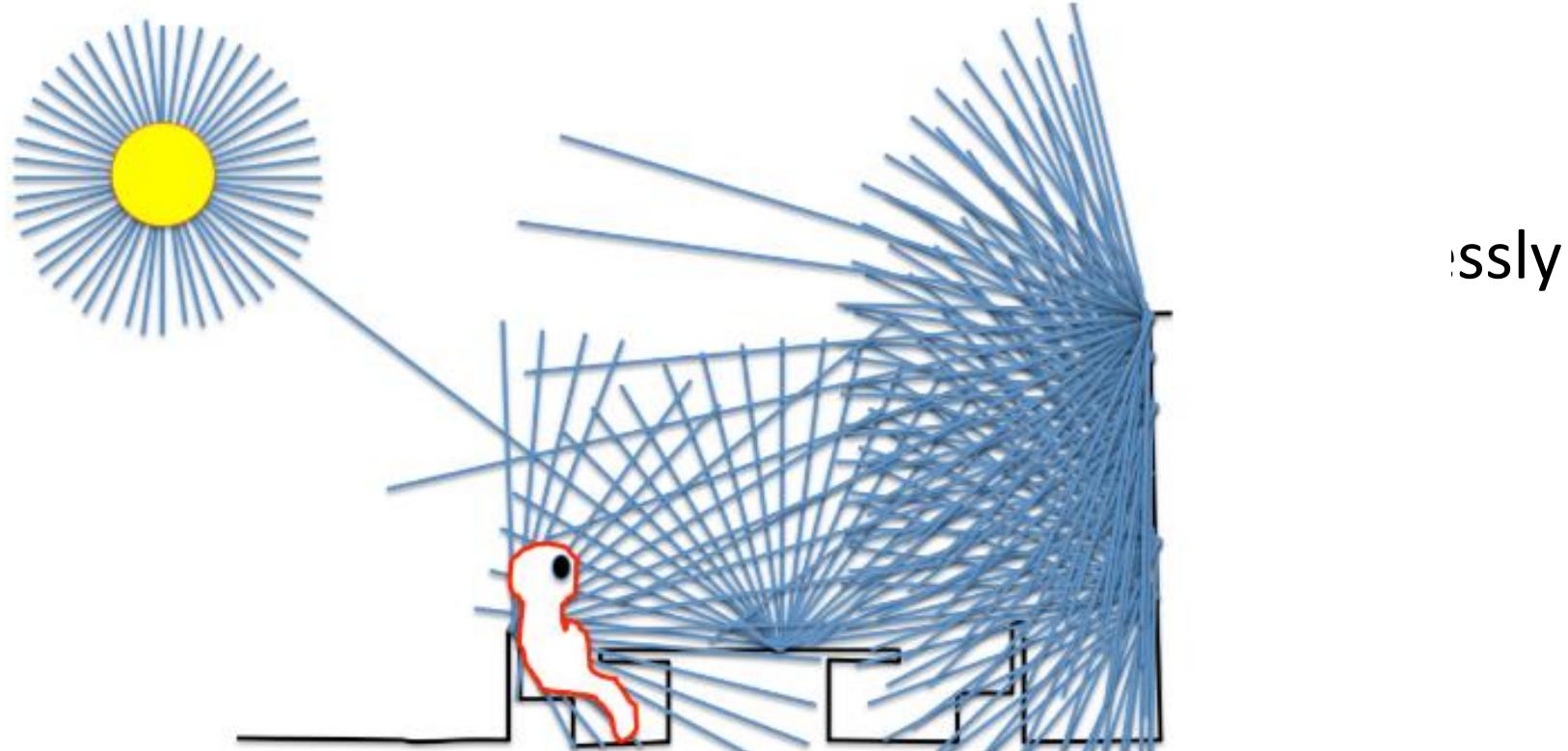
<https://techcrunch.com/2023/01/07/ai-guided/>

Dr. Sander Ali Khowaja



Computer Vision Overview

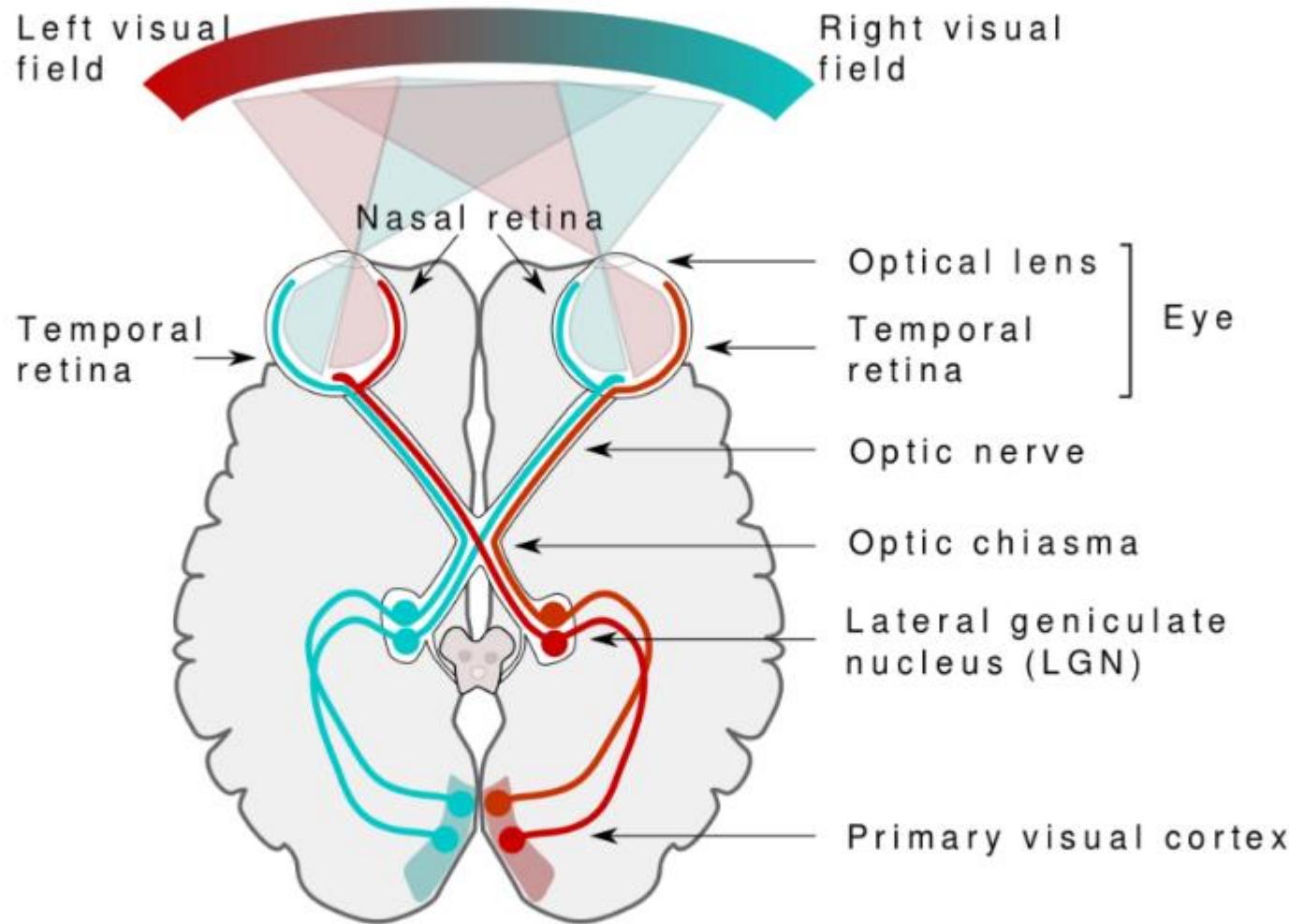
- The objective of the computer vision is to make the computer see and interpret the world like humans and possibly even better



Goal of Computer Vision is to **convert light into meaning** (geometric, semantic, ...)

Image Credits: Antonio Torralba

How it works



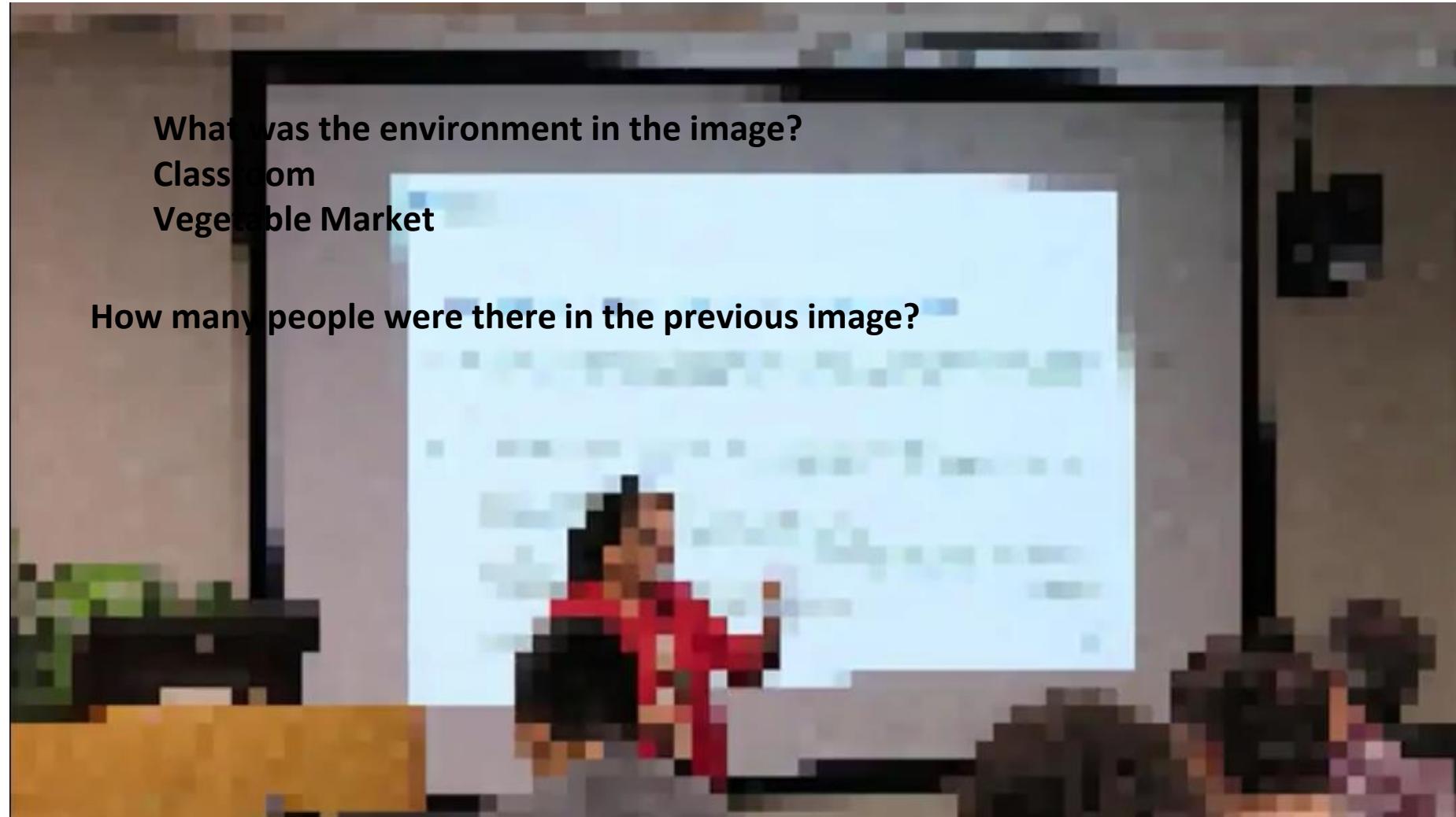
Over **50%** of the processing in the **human brain** is dedicated to **visual information**.

- Theoretical and algorithmic basis to achieve automatic visual understanding.
 - High-level understanding of digital images and videos.
 - Computer Vision aims to come up with computational models of the human visual system.
-
- From engineering point of view:

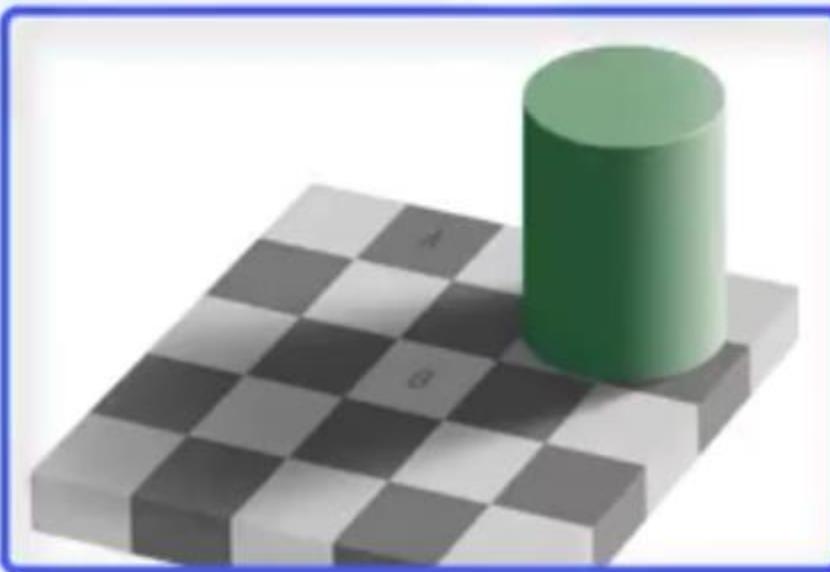
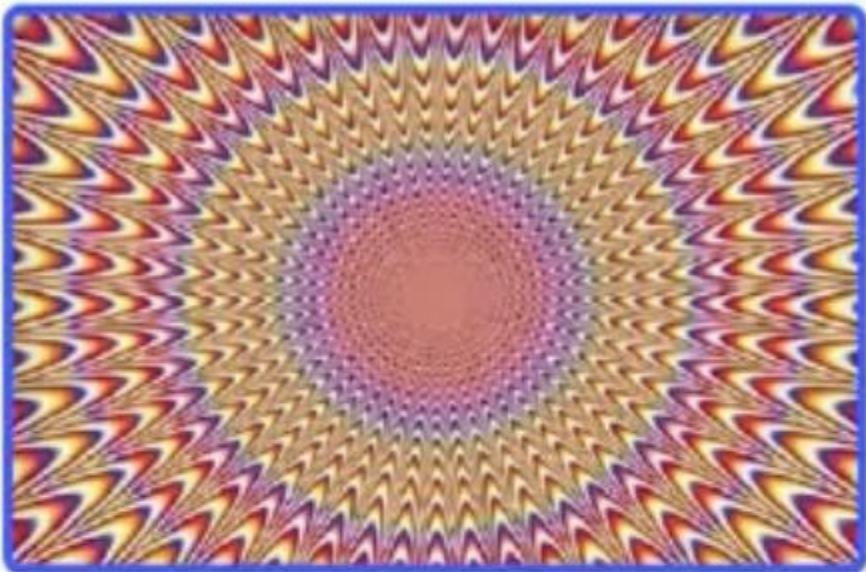
Computer vision aims to build autonomous systems to perform some of the tasks which the human visual system can perform and even surpass it in many cases.



Lets see the power of human vision system



How good is human visual system



More on Illusions



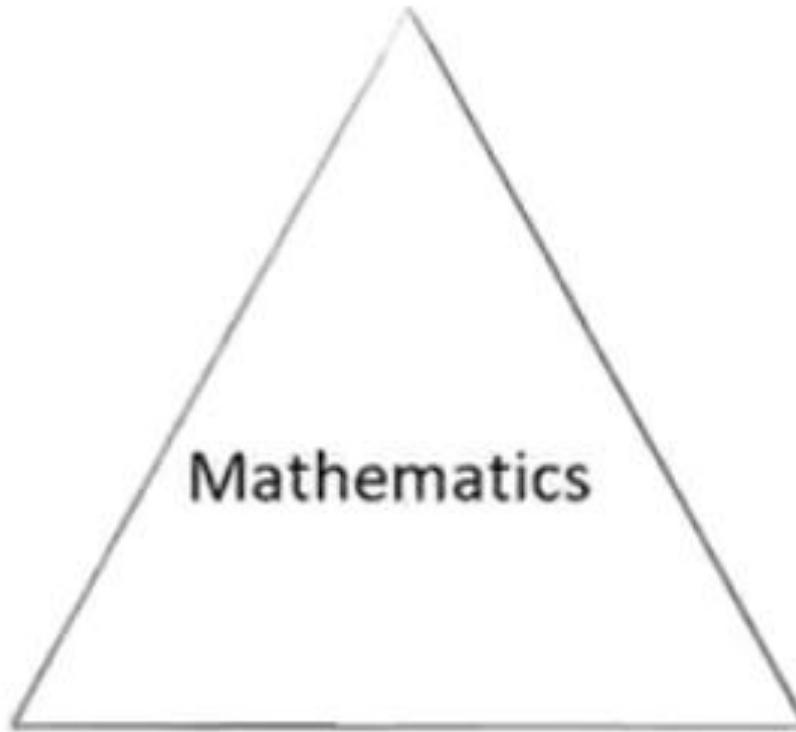
Dr. Sandeep Ankitowaja



Related Fields of Computer Vision



Digital Signal Processing
Computer Vision



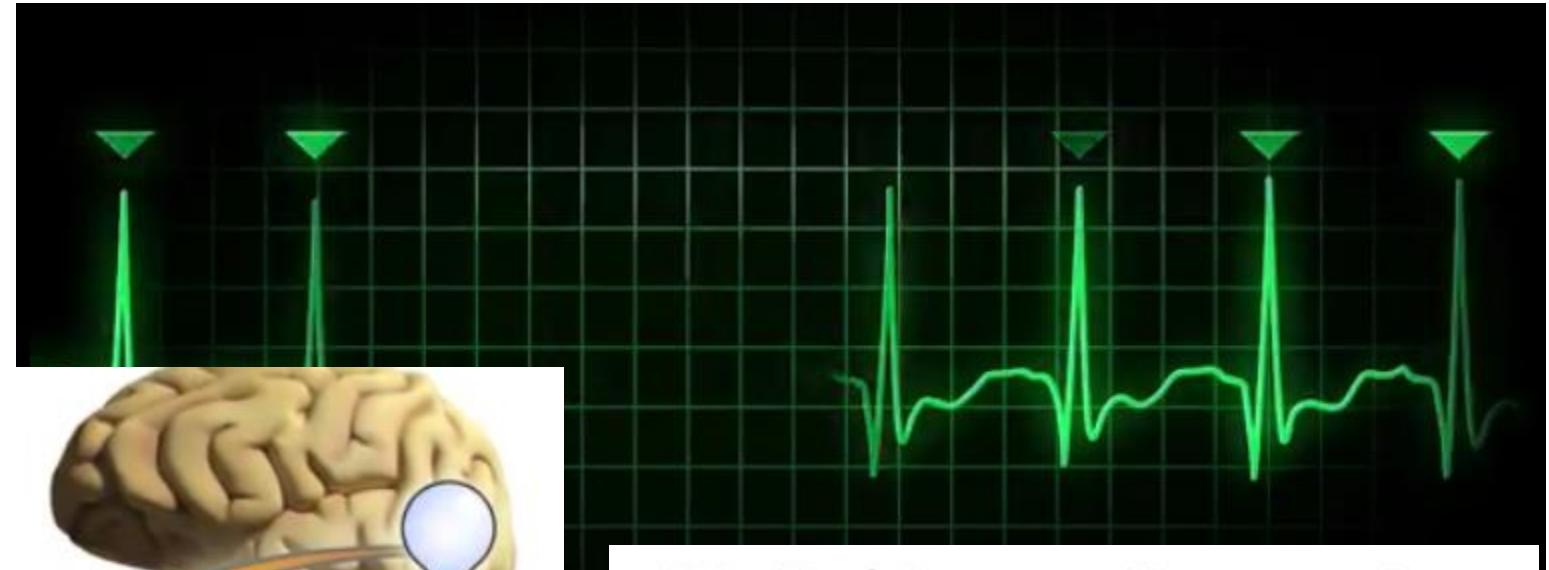
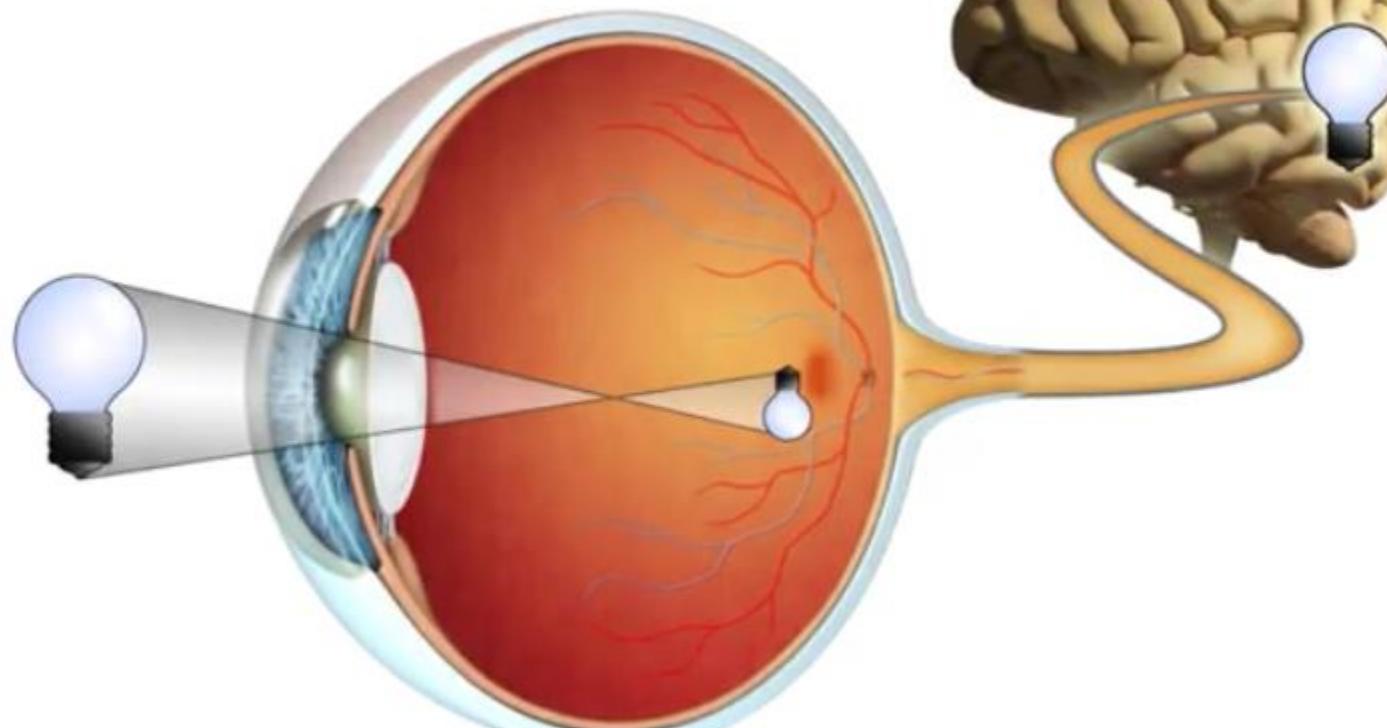
Dr. Sander Ali Khowaja



Computer Vision Pipeline

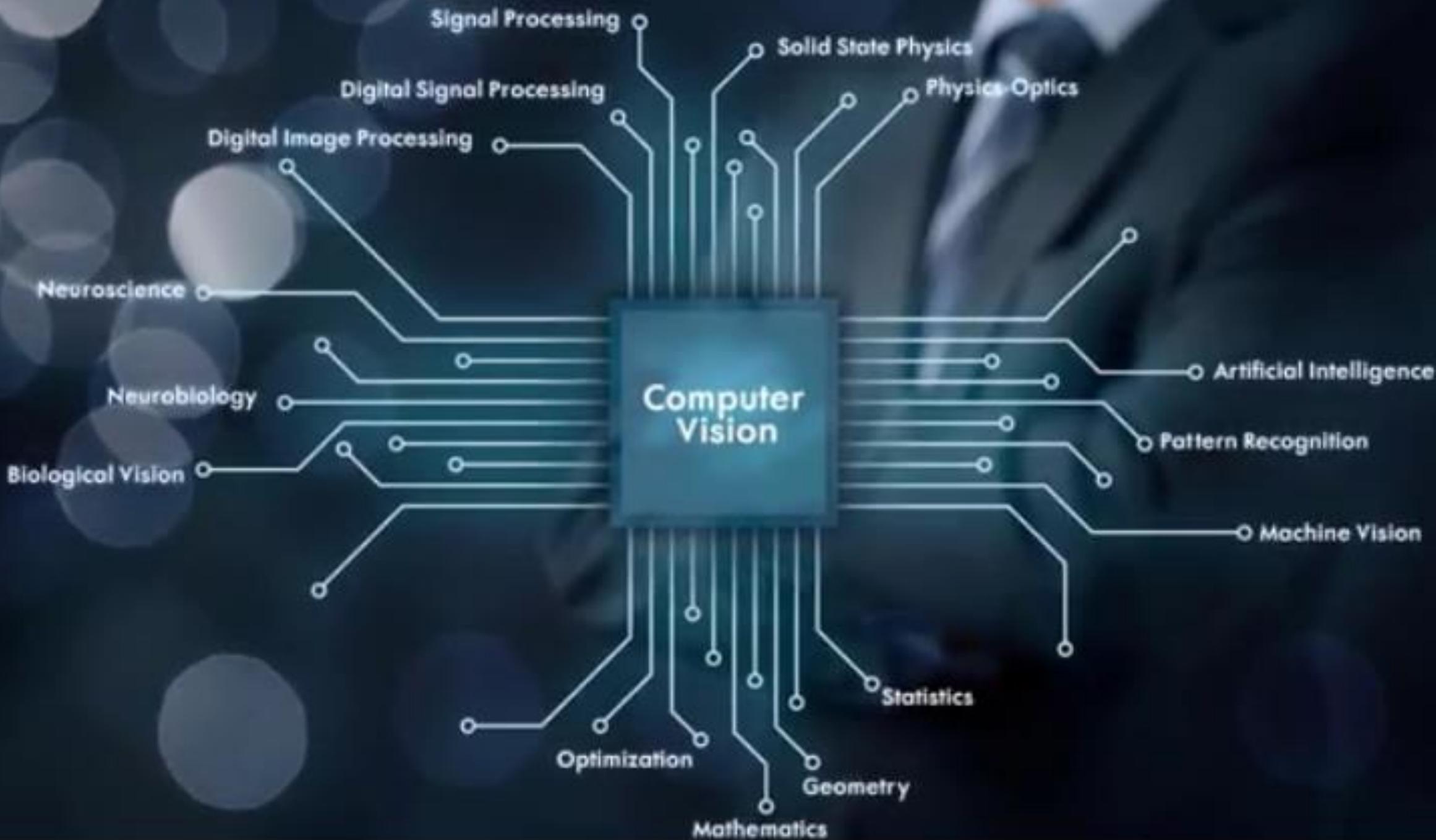
Optics
Electromagnetic Theory
Solid State Physics

Physics

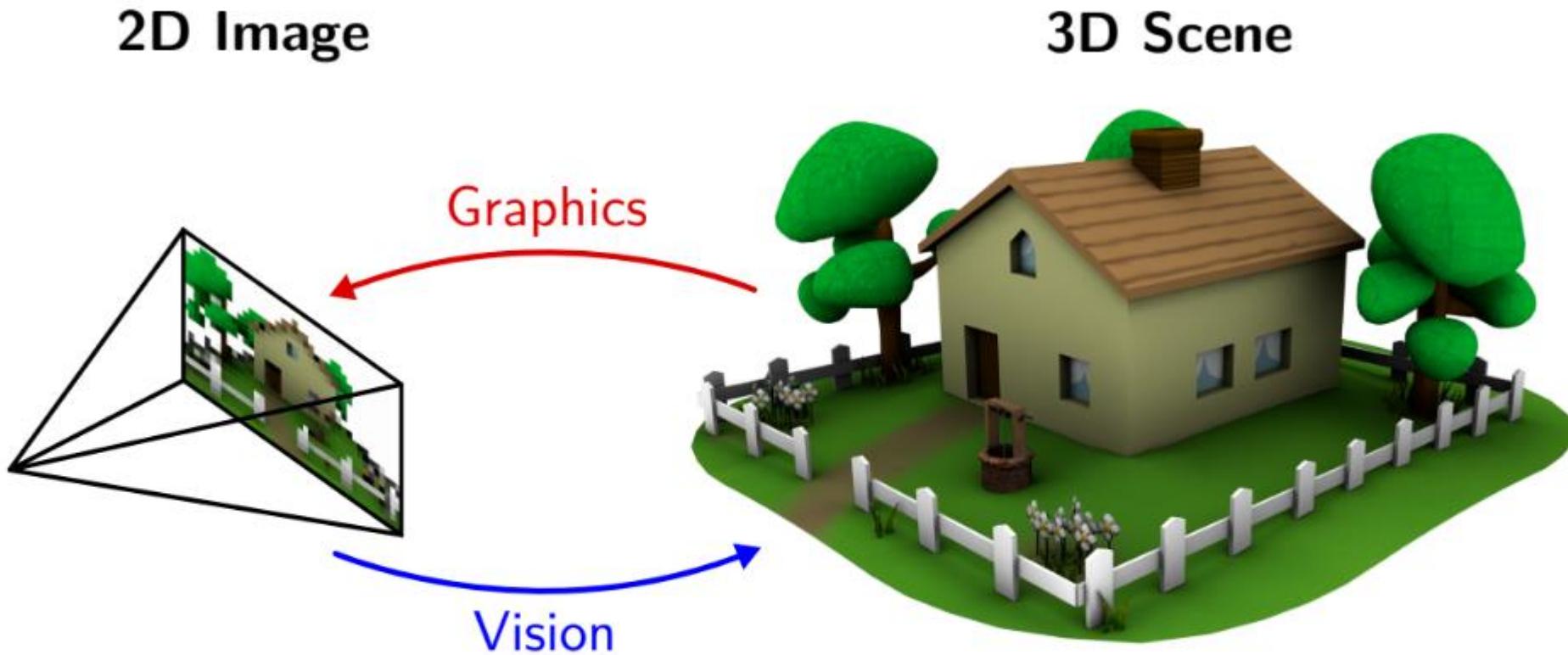


Digital Image Processing





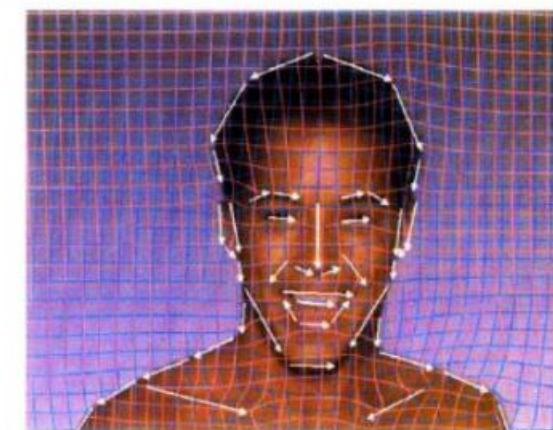
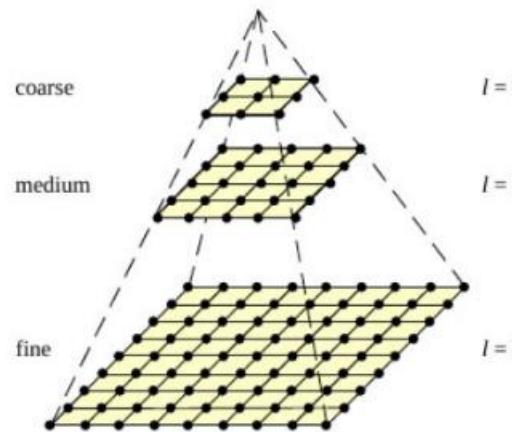
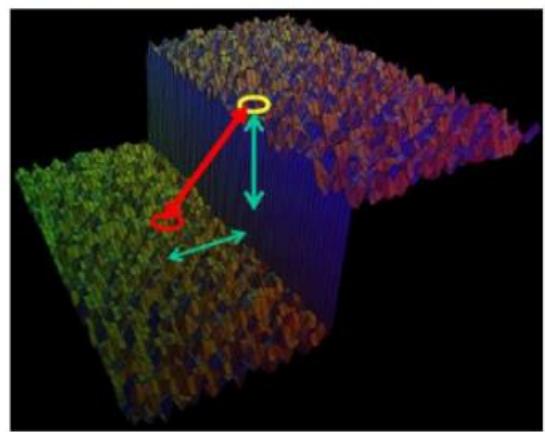
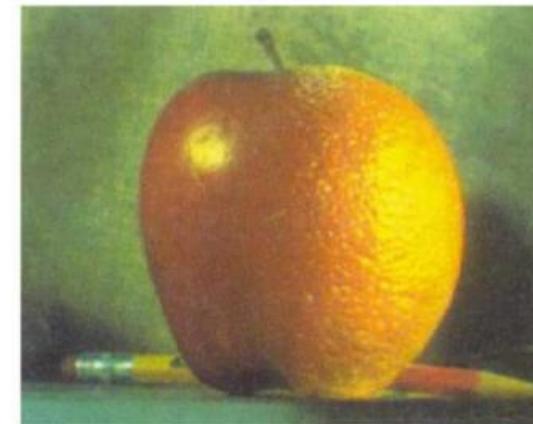
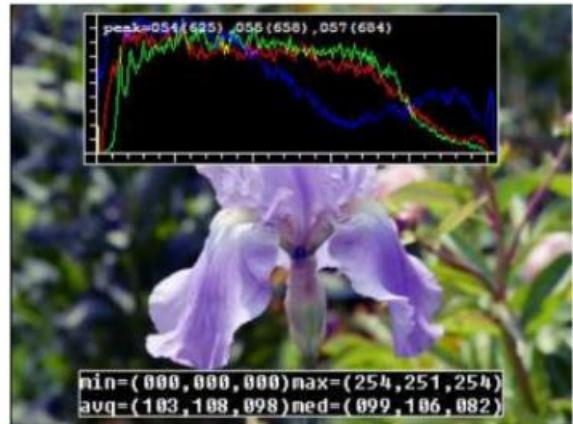
Computer Vision vs Computer Graphics



Computer Vision is an ill-posed inverse problem:

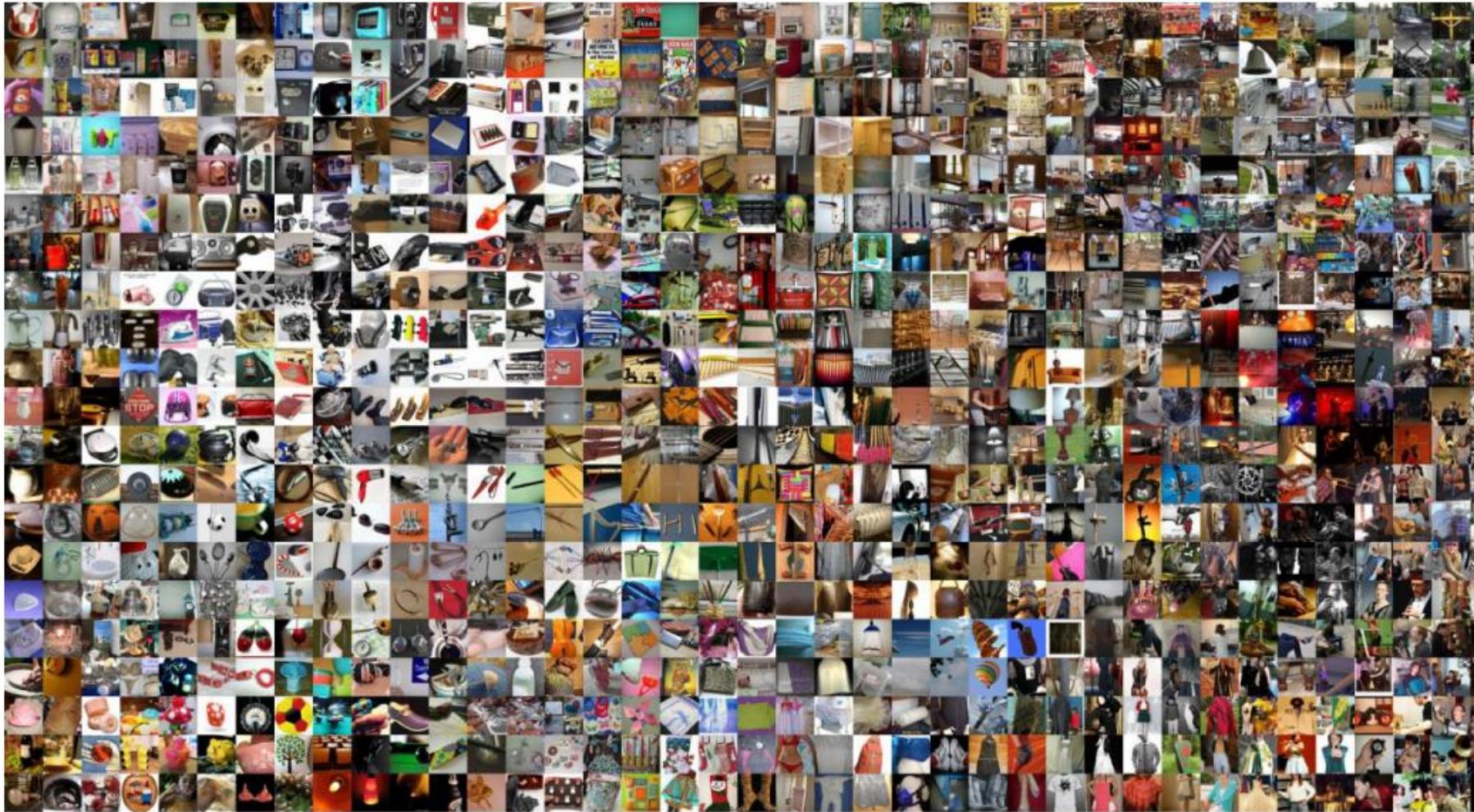
- ▶ Many 3D scenes yield the same 2D image
- ▶ Additional constraints (knowledge about world) required

Computer Vision vs Image Processing



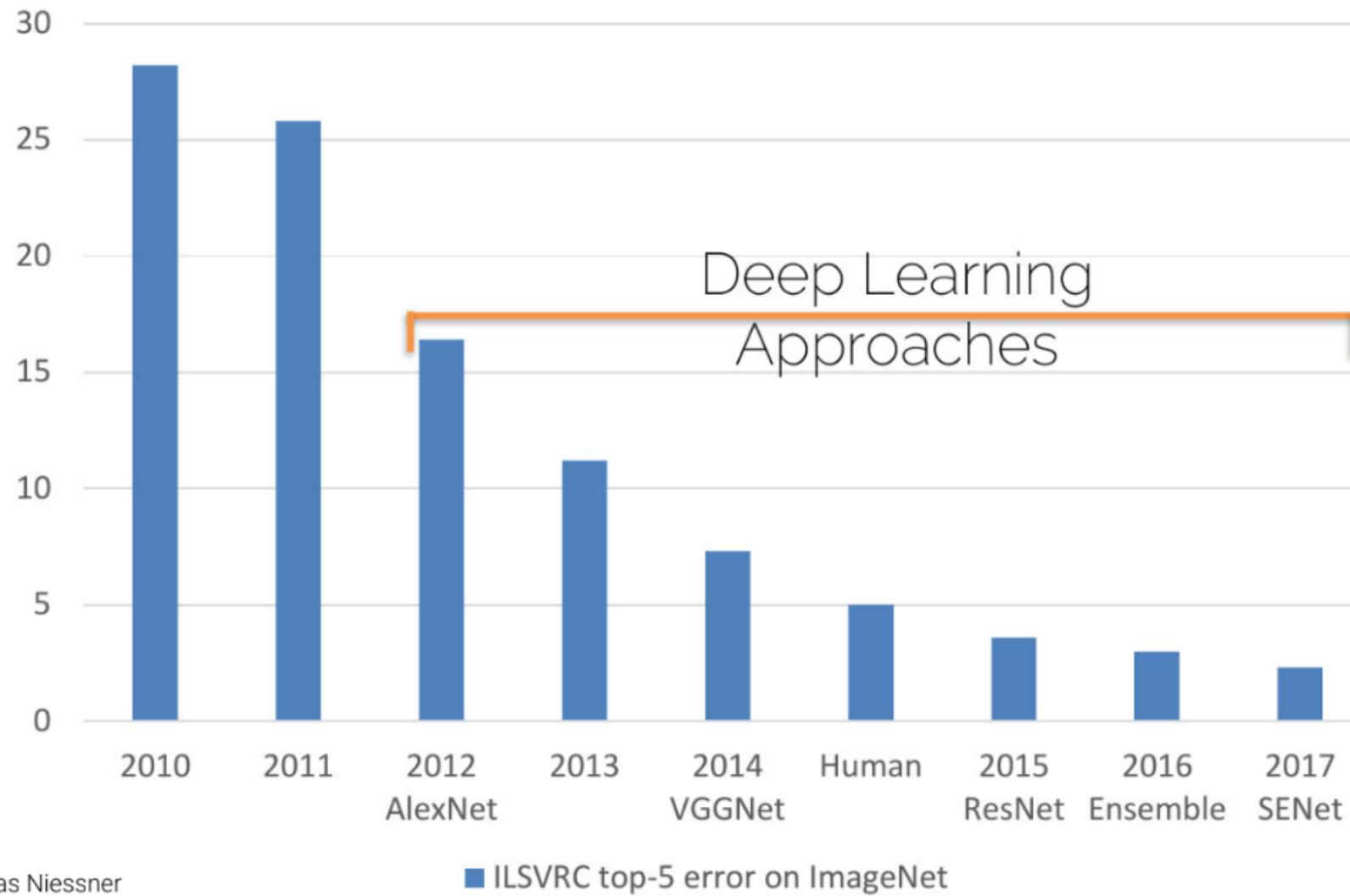
Slide Credits: Rick Szeliski

Computer Vision vs Machine Learning



Deng, Dong, Socher, Li, Li and Li: ImageNet: A large-scale hierarchical image database. CVPR, 2009.

Deep Learning Revolution



Slide Credit: Matthias Niessner



Why is perception hard?

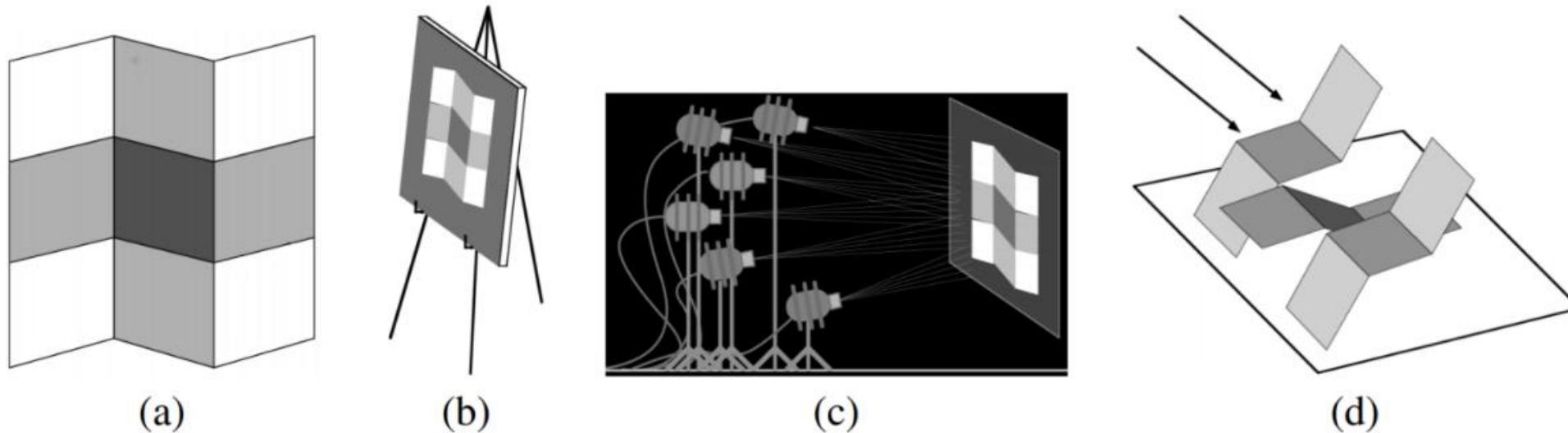


What we see

200	133	110	103	117	90	47	30	32	79	66	65
197	122	123	138	98	100	46	45	22	11	43	55
140	116	165	159	90	56	58	47	26	13	54	102
132	148	119	106	123	57	64	46	21	22	79	94
125	121	80	143	101	55	61	38	20	21	81	65
50	71	74	63	52	39	41	39	32	26	97	66
51	59	62	44	40	40	36	28	27	31	29	44
59	62	70	50	48	35	34	35	26	21	24	32
49	59	65	64	58	34	40	28	26	21	23	124
39	45	47	64	54	34	40	24	19	47	133	207
37	42	39	38	39	50	75	74	105	170	197	167
37	47	33	35	50	108	162	184	184	157	125	112
45	48	35	37	75	148	183	156	83	91	91	116
49	48	54	50	75	158	110	66	74	128	155	149
48	51	57	50	65	91	79	92	101	105	132	132
51	58	66	55	58	52	91	91	88	115	158	174
57	60	61	52	56	61	60	55	92	146	188	190
65	50	54	56	57	51	54	56	80	115	177	187
67	40	40	61	65	48	39	30	36	75	151	181
53	32	36	35	61	43	37	26	29	35	126	189
29	42	107	20	28	41	40	26	30	36	113	200
30	21	32	24	34	37	33	23	25	39	105	171
32	28	19	23	29	36	47	89	132	169	183	128
31	25	62	54	47	44	81	190	227	231	206	155
44	66	99	72	67	63	89	128	127	115	109	157
53	47	47	41	29	32	25	20	41	81	89	175
38	44	61	73	54	48	37	87	90	111	126	189
39	41	83	97	86	91	74	134	131	153	143	185
42	56	98	102	112	111	94	137	121	141	146	181
94	114	114	114	122	113	77	117	117	154	149	169
157	176	116	121	130	139	103	161	148	180	145	125
143	178	182	178	139	153	129	168	175	187	170	152
127	183	203	197	153	164	143	180	195	182	165	211
88	107	127	125	101	107	100	123	149	186	167	215

What the computer sees

Challenges: Images are 2D Projections of the 3D World



Adelson and Pentland's workshop metaphor:

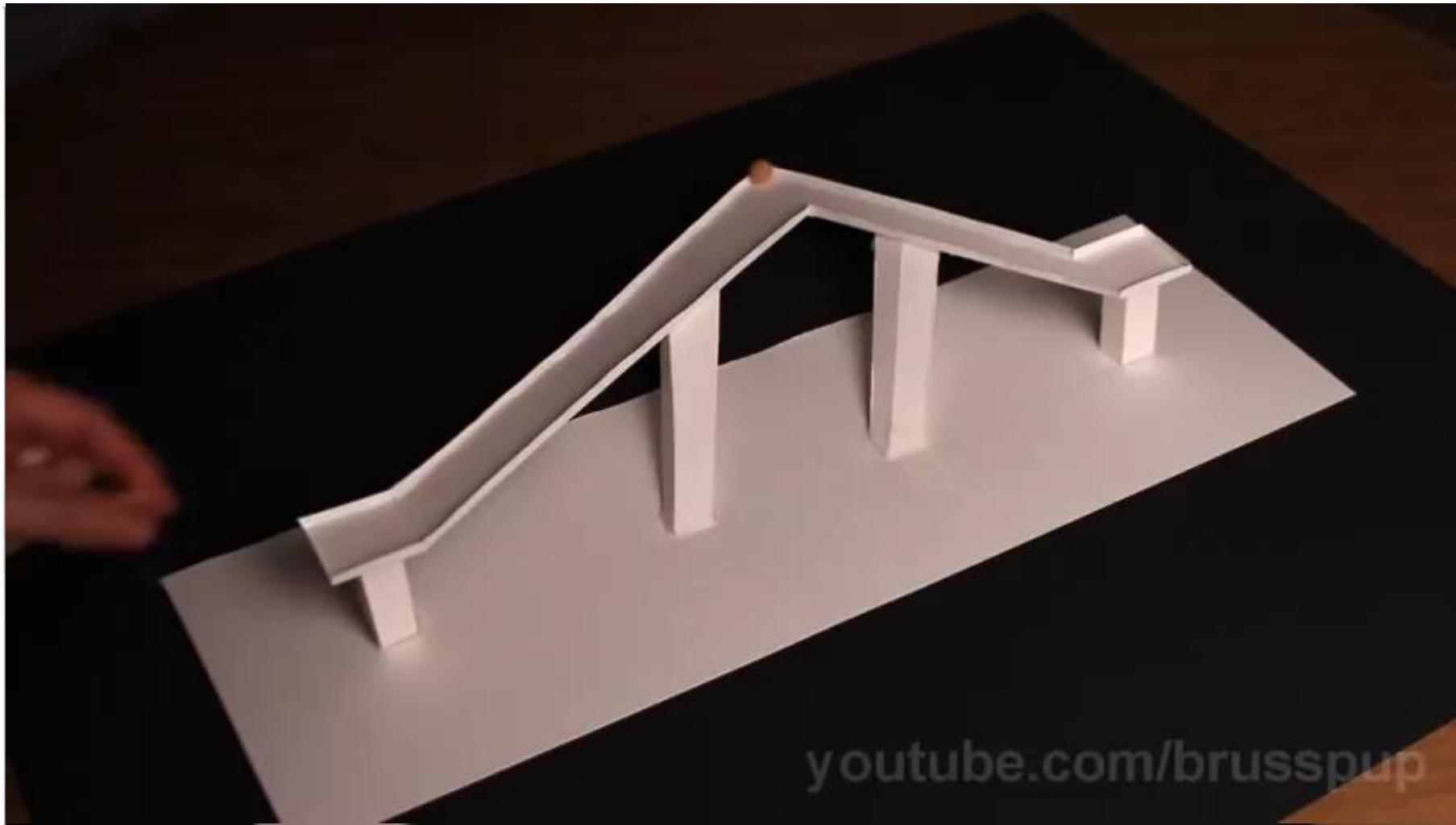
To explain an image (a) in terms of reflectance, lighting and shape, (b) a painter, (c) a light designer and (d) a sculptor will design three different, but plausible, solutions.

Adelson and Pentland: The perception of shading and reflectance. Perception as Bayesian inference, 1996.

Challenges: Ames room Illusion



Challenges: Perspective Illusion



Challenges: Viewpoint Variation



Challenges: Deformation



Xu Beihong (1943)

Challenge: Occlusion



René Magritte (1957)

Challenges: Illumination

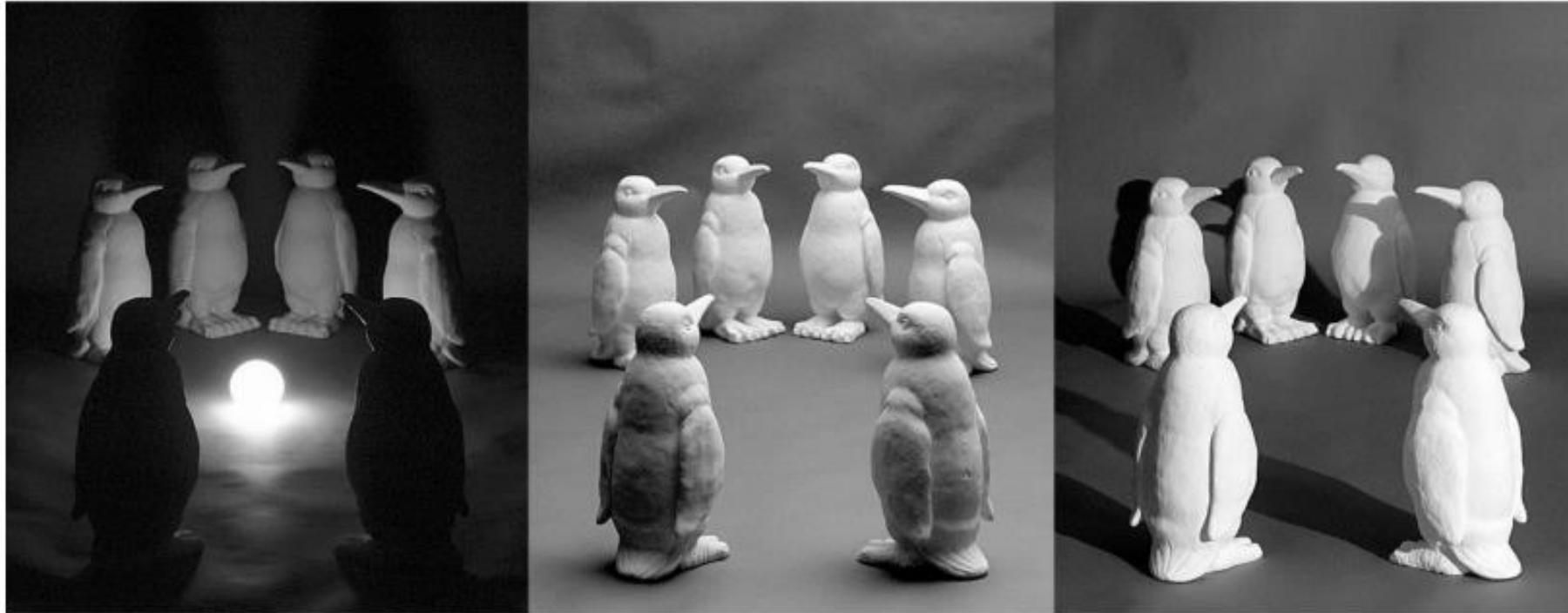


Image Credits: Jan Koenderink

Challenges: Motion



Challenges: Perception vs Measurement



Local Ambiguities



Image Credits: Antonio Torralba

Challenge: Intra class Variation



<http://www.homeworkshop.com/>

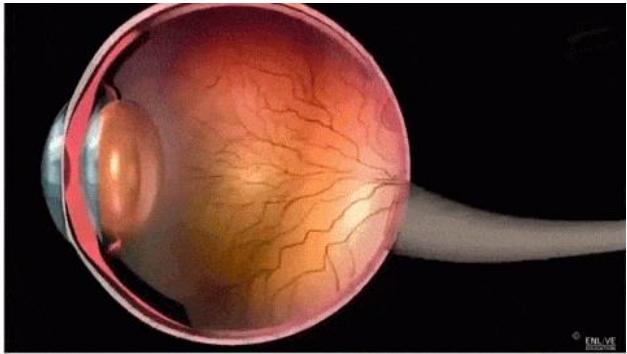
Image Credits: Antonio Torralba

Challenge: Number of Object Categories



Image Credits: Antonio Torralba

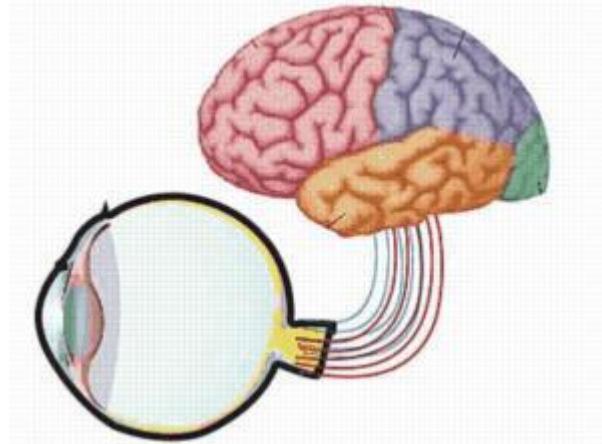
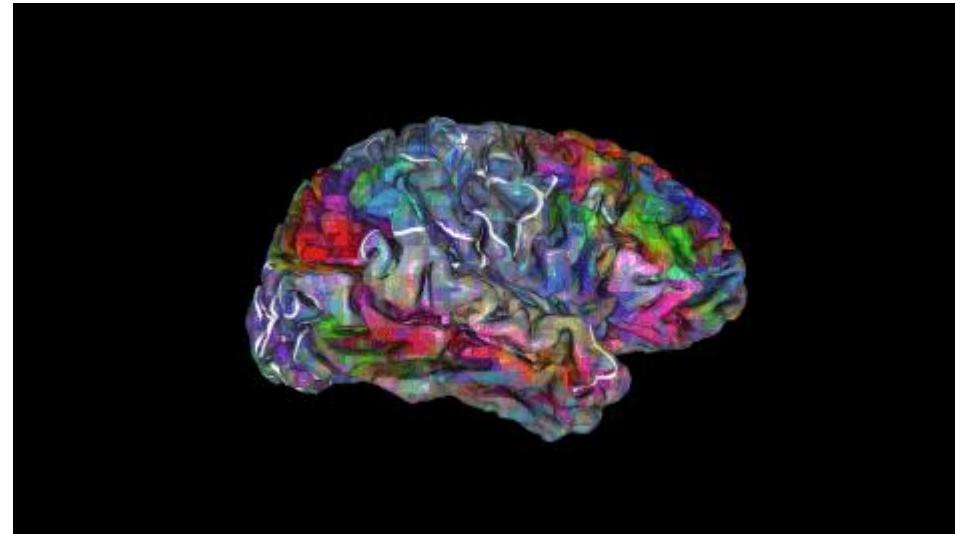
Three distinct lines



Replicating the Eye

Replicating the Visual Cortex

Replicating the rest of the brain



Credits for History of computer Vision

- Svetlana Lazebnik (UIUC): Computer Vision: Looking Back to Look Forward

<https://slazebni.cs.illinois.edu/spring20/>

- Steven Seitz (Univ. of Washington): 3D Computer Vision: Past, Present, and Future

<http://www.youtube.com/watch?v=kylzMr917Rc>

<http://www.cs.washington.edu/homes/seitz/talks/3Dhistory.pdf>



Pre-History



Perspective
Leonardo da Vinci
(1452–1519)



Photometry
Johann Heinrich Lambert
(1728–1777)



Least Squares
Carl Friedrich Gauss
(1777–1855)



Stereopsis
Charles Wheatstone
(1802–1875)

1510: Perspectograph



"Perspective is nothing else than the seeing of an object behind a sheet of glass, smooth and quite transparent, on the surface of which all the things may be marked that are behind this glass. All things transmit their images to the eye by pyramidal lines, and these pyramids are cut by the said glass. The nearer to the eye these are intersected, the smaller the image of their cause will appear."

– Leonardo da Vinci

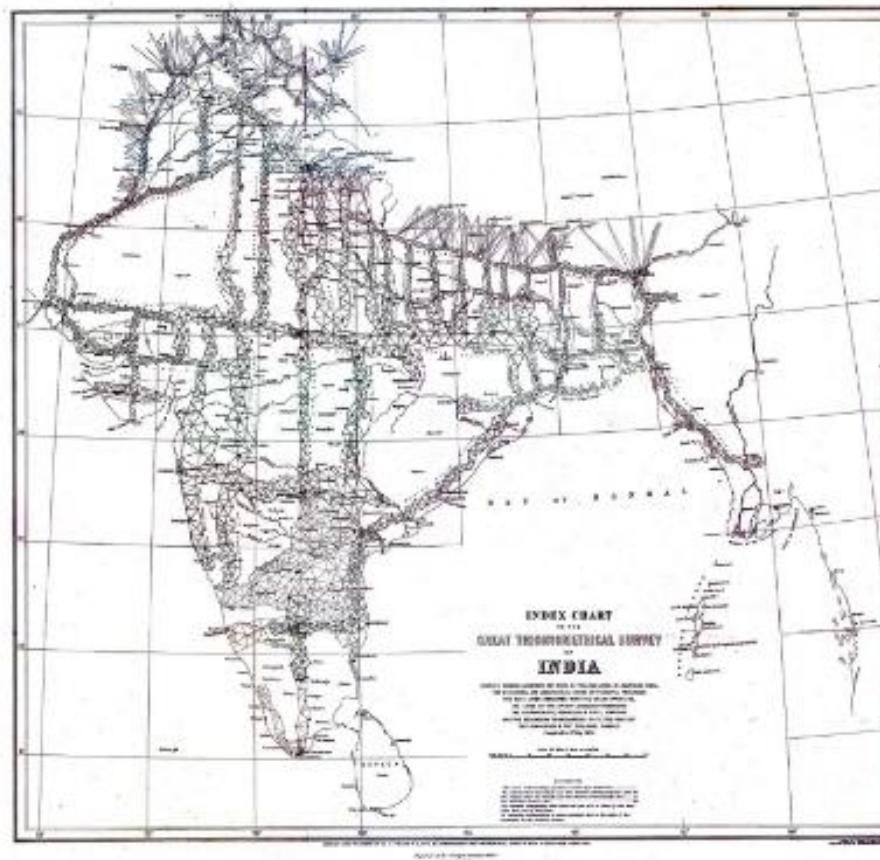
1839: Daguerreotype

- ▶ First publicly available photographic process invented by Louis Daguerre
- ▶ Widely used during the 1840s and 1850s
- ▶ Polish a sheet of silver-plated copper and treat with fumes to make light sensitive
- ▶ Make resulting latent image visible by fuming it with mercury vapor and remove its sensitivity to light by chemical treatment
- ▶ Rinse, dry and seal behind glass



1802-1871: Great Trigonometrical Survey

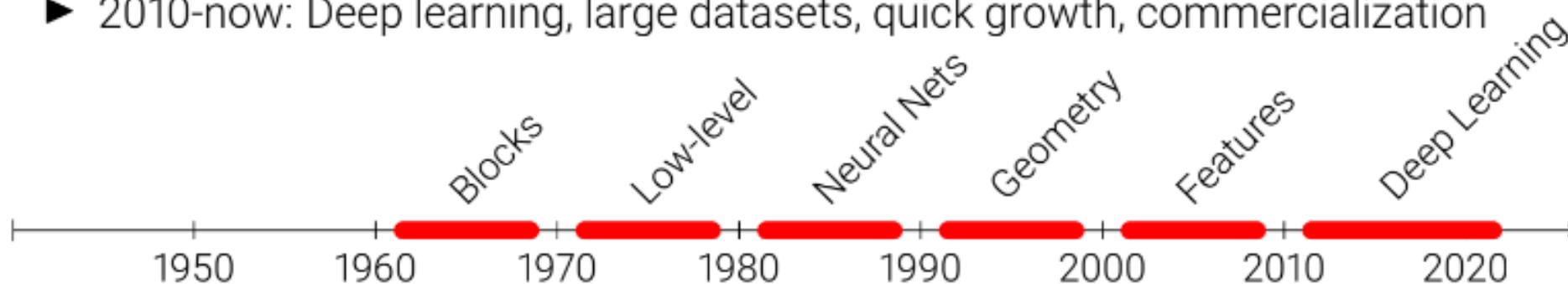
- ▶ Multi-decade project to measure the entire Indian subcontinent with scientific precision
- ▶ Under the leadership of George Everest, the project was made responsible of the Survey of India
- ▶ Manual bundle adjustment proves Mt. Everest highest mountain on earth mountain on earth



Historical Overview

Waves of development:

- ▶ 1960-1970: Blocks Worlds, Edges and Model Fitting
- ▶ 1970-1981: Low-level vision: stereo, flow, shape-from-shading
- ▶ 1985-1988: Neural networks, backpropagation, self-driving
- ▶ 1990-2000: Dense stereo and multi-view stereo, MRFs
- ▶ 2000-2010: Features, descriptors, large-scale structure-from-motion
- ▶ 2010-now: Deep learning, large datasets, quick growth, commercialization



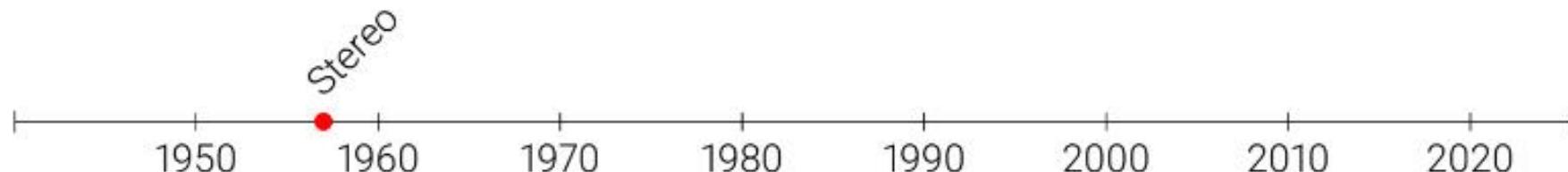
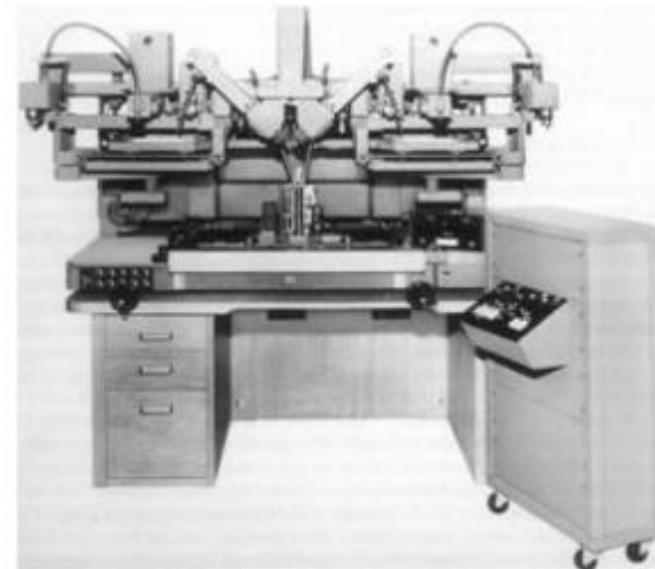
7



A Brief History of Computer Vision

1957: Stereo

- ▶ Gilbert Hobrough demonstrated an analog implementation of stereo image correlation
- ▶ This led to the creation of the Raytheon-Wild B8 Stereomat
- ▶ Used to create Elevation Maps (Photogrammetry, since 1840)



Roberts: Machine perception of 3-d solids. PhD Thesis, 1965.

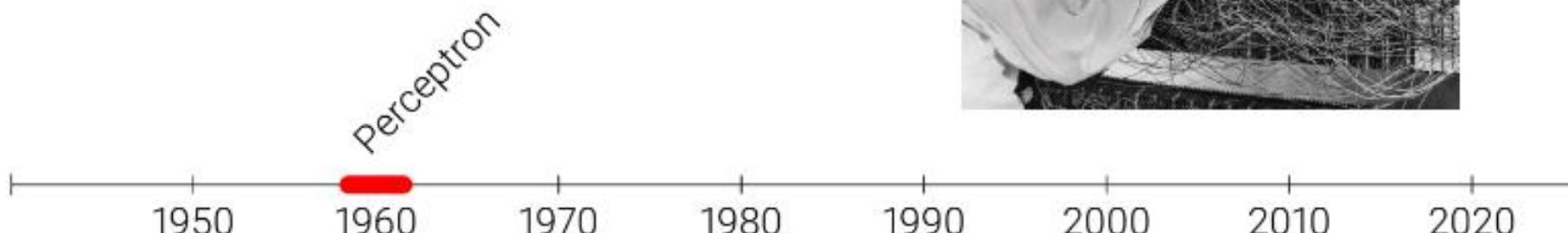
A Brief History of Computer Vision

1958-1962: Rosenblatt's Perceptron

- ▶ First algorithm and implementation to train single linear threshold neuron
- ▶ Optimization of perceptron criterion:

$$\mathcal{L}(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \mathbf{x}_n y_n$$

- ▶ Novikoff proved convergence

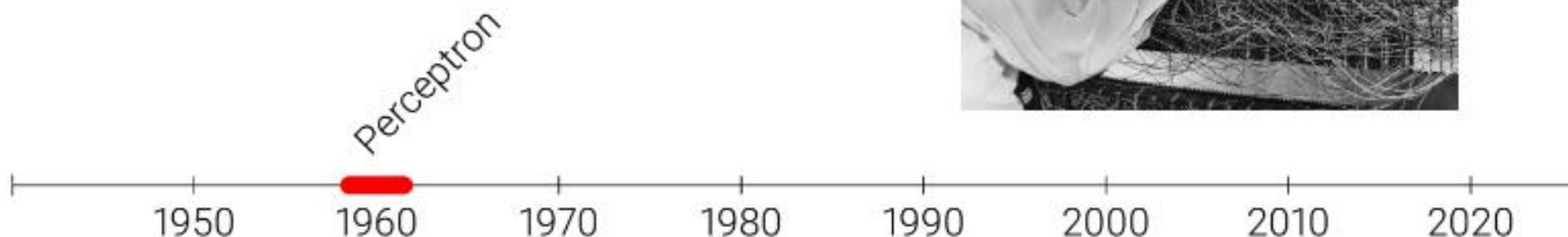


Rosenblatt: The perceptron - a probabilistic model for information storage and organization in the brain. Psychological Review, 1958.

A Brief History of Computer Vision

1958-1962: Rosenblatt's Perceptron

- ▶ First algorithm and implementation to train single linear threshold neuron
- ▶ Overhyped: Rosenblatt claimed that the perceptron will lead to computers that walk, talk, see, write, reproduce and are conscious of their existence



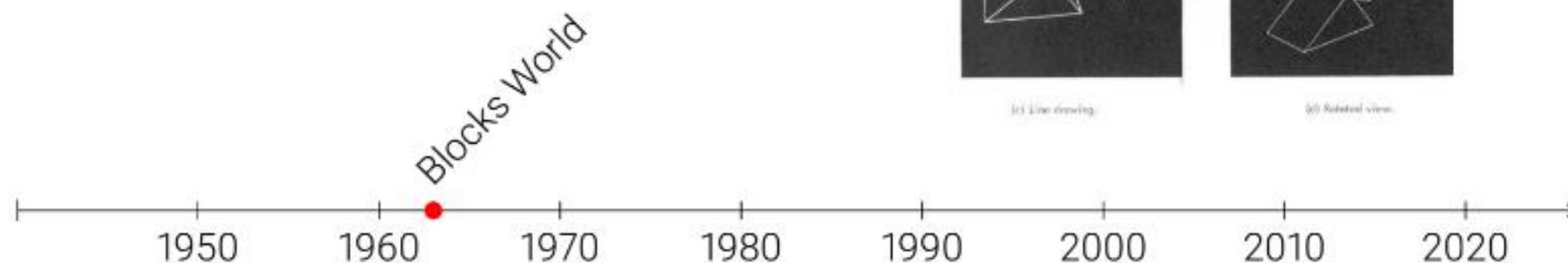
Rosenblatt: The perceptron - a probabilistic model for information storage and organization in the brain. Psychological Review, 1958.



A Brief History of Computer Vision

1963: Larry Robert's Blocks World

- ▶ Scene understanding for robotics
- ▶ Extracts edges as primitives
- ▶ Infers 3D structure of an object from topological structure of the 2D lines
- ▶ Interpret images as projections of 3D scenes, not 2D pattern recognition

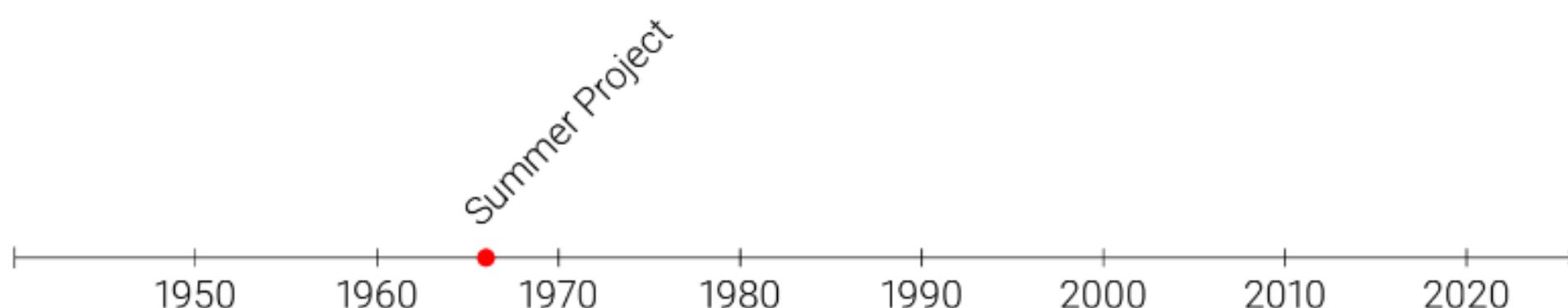


Roberts: Machine Perception of Three-Dimensional Solids. PhD Thesis, 1965.

A Brief History of Computer Vision

1966: MIT Summer Vision Project

- Underestimated the challenge of computer vision, committed to "blocks world"



Paper: The Summer Vision Project. MIT AI Memos, 1966.

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group
Vision Memo, No. 103.

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".

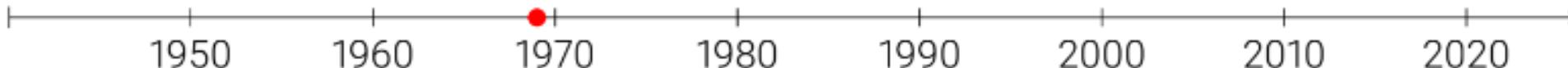


A Brief History of Computer Vision

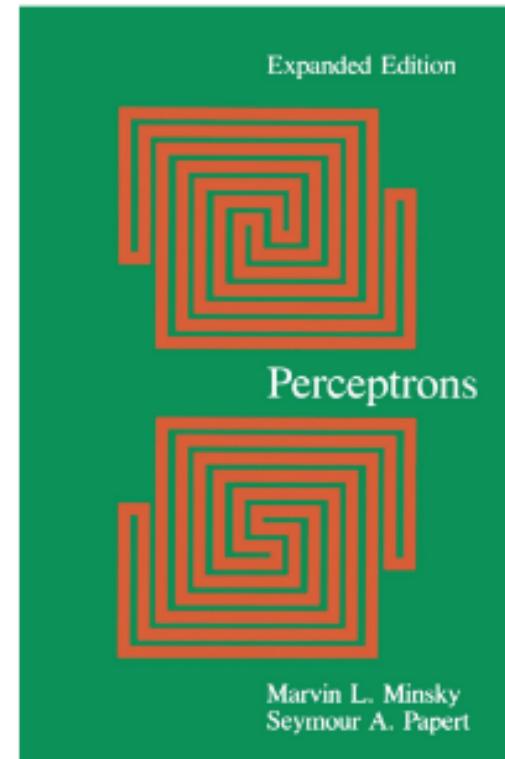
1969: Minsky and Papert publish book

- ▶ Several discouraging results
- ▶ Showed that single-layer perceptrons cannot solve some very simple problems (XOR problem, counting)
- ▶ Symbolic AI research dominates 70s

Minsky/Papert



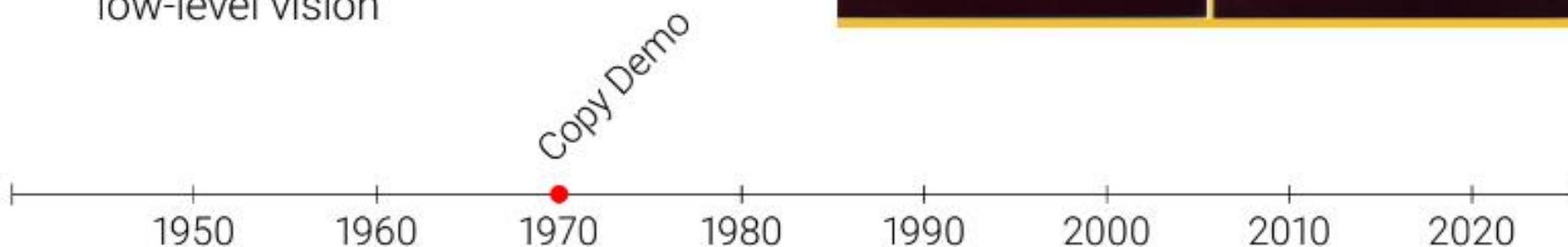
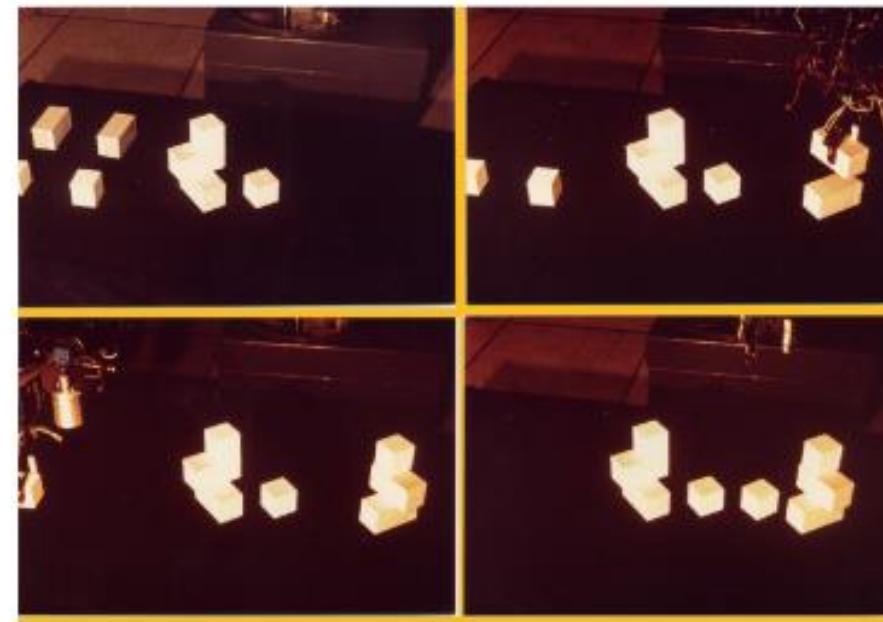
Minsky and Papert: *Perceptrons: An introduction to computational geometry*. MIT Press, 1969.



A Brief History of Computer Vision

1970: MIT Copy Demo

- ▶ Vision system recovers structure of a blocks scene, robot plans and builds copy from another set of blocks
- ▶ Vision, planning and manipulation
- ▶ But low-level edge finding not robust enough for task, led to attention on low-level vision



Paper: The Summer Vision Project. MIT AI Memos, 1966.

A Brief History of Computer Vision

1970: Shape from Shading

- ▶ Recover 3D from single 2D image
- ▶ Assumes Lambertian surface and constant albedo
- ▶ Applies smoothness regularization to constrain the ill-posed problem

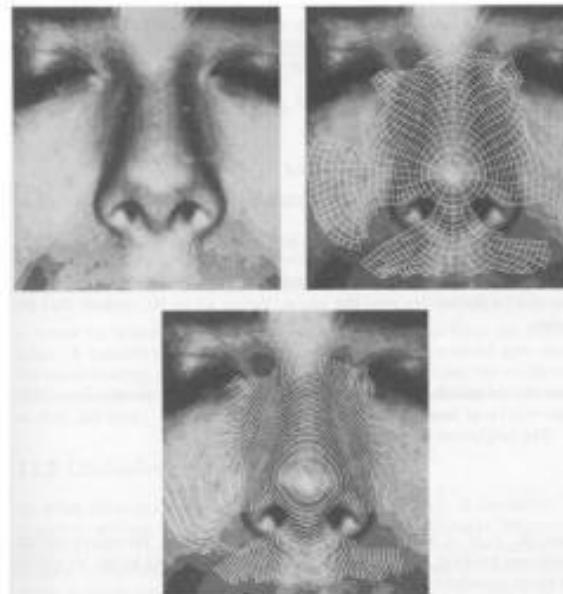
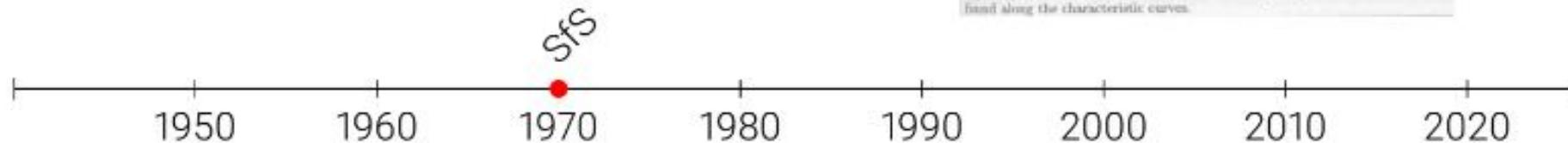


Figure 11-7. The shape-from-shading method is applied here to the recovery of the shape of a nose. The first picture shows the (crudely quantized) gray-level image available to the program. The second picture shows the base characteristics superimposed, while the third shows a contour map computed from the elevations found along the characteristic curves.

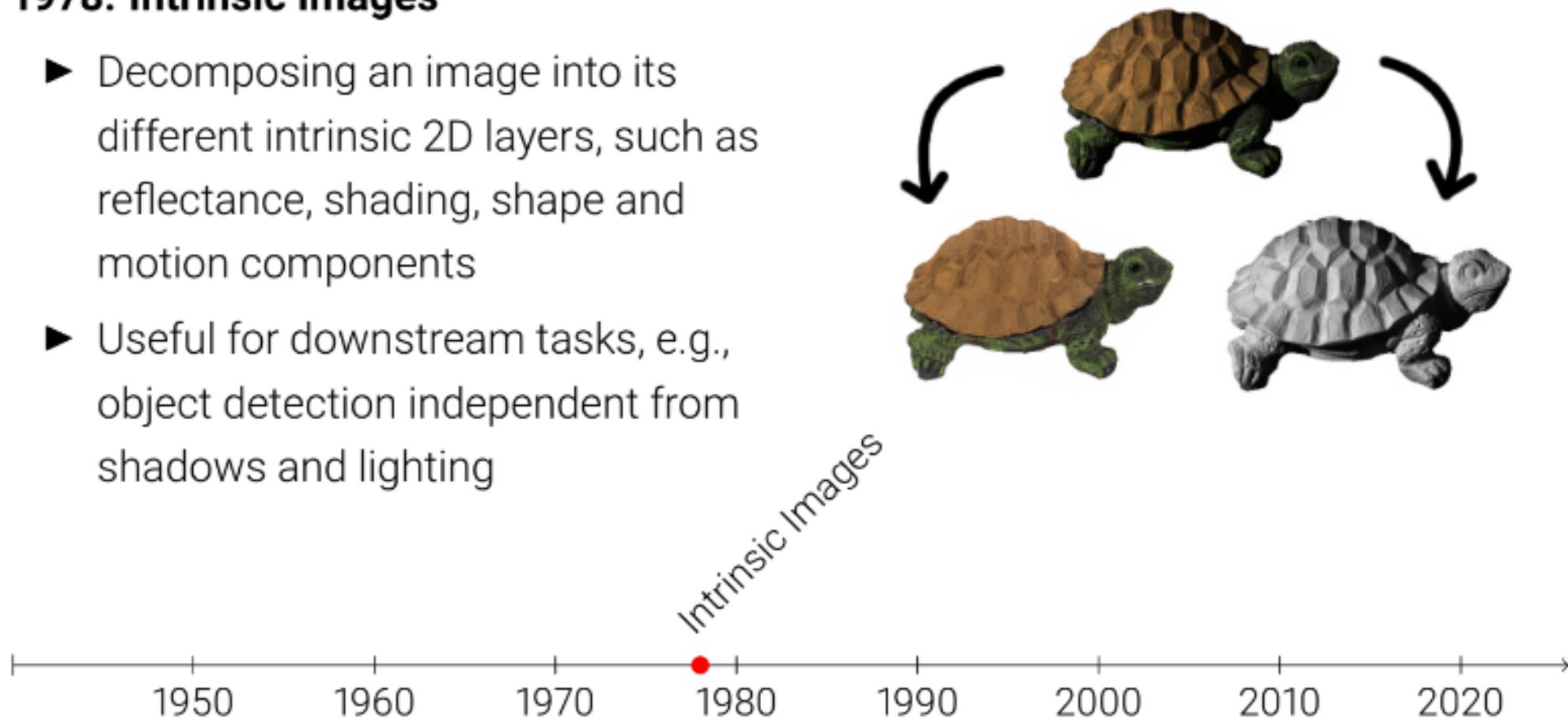


Horn: Shape From Shading: A Method for Obtaining the Shape of a Smooth Opaque Object From One View. MIT TR, 1970.

A Brief History of Computer Vision

1978: Intrinsic Images

- Decomposing an image into its different intrinsic 2D layers, such as reflectance, shading, shape and motion components
- Useful for downstream tasks, e.g., object detection independent from shadows and lighting

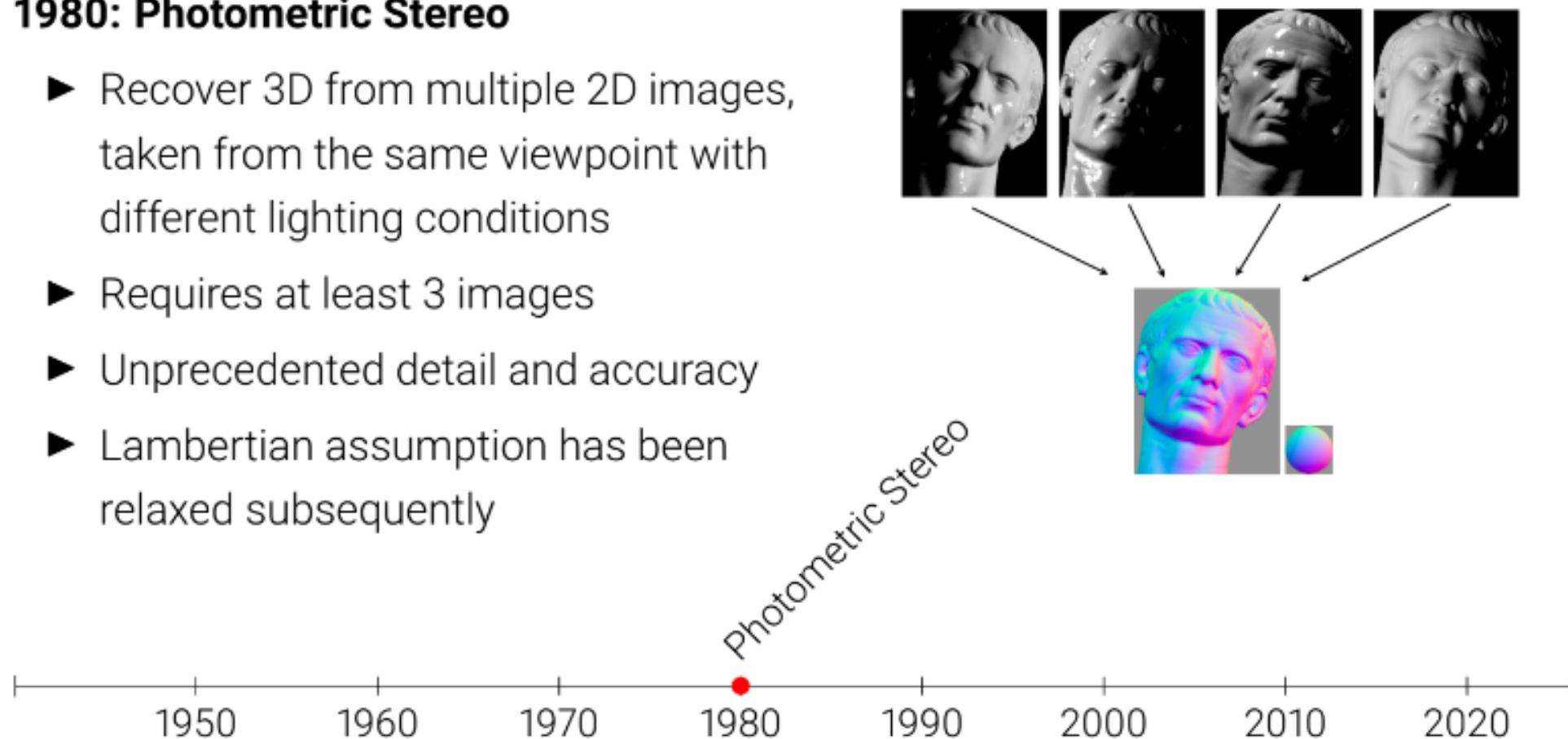


Barrow and Tenenbaum: Recovering intrinsic scene characteristics from images. Computer Vision Systems, 1978.

A Brief History of Computer Vision

1980: Photometric Stereo

- ▶ Recover 3D from multiple 2D images, taken from the same viewpoint with different lighting conditions
- ▶ Requires at least 3 images
- ▶ Unprecedented detail and accuracy
- ▶ Lambertian assumption has been relaxed subsequently

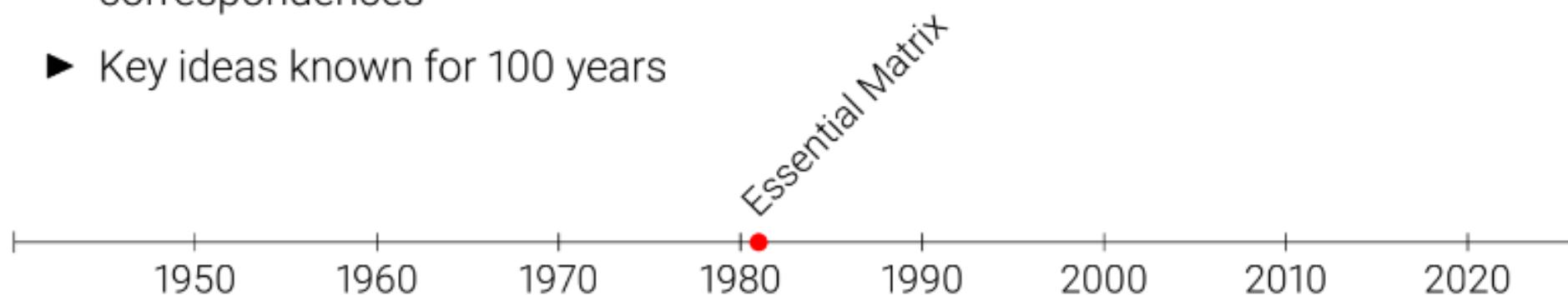
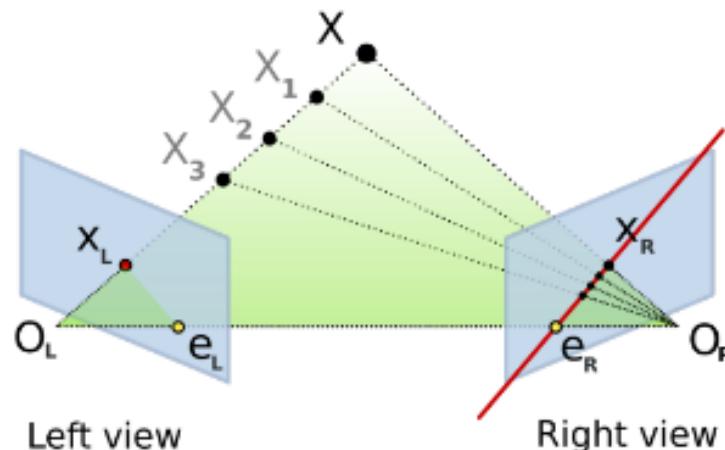


Woodham. Photometric method for determining surface orientation from multiple images. Optical Engineering, 1980.

A Brief History of Computer Vision

1981: Essential Matrix

- ▶ Defines two-view geometry as matrix mapping points to epipolar lines
- ▶ Reduces correspondence search to a 1D problem
- ▶ Can be estimated from a set of 2D correspondences
- ▶ Key ideas known for 100 years

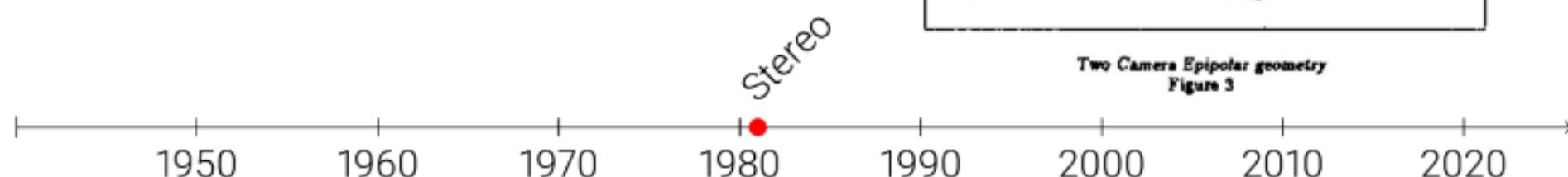


Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. Nature, 1981.

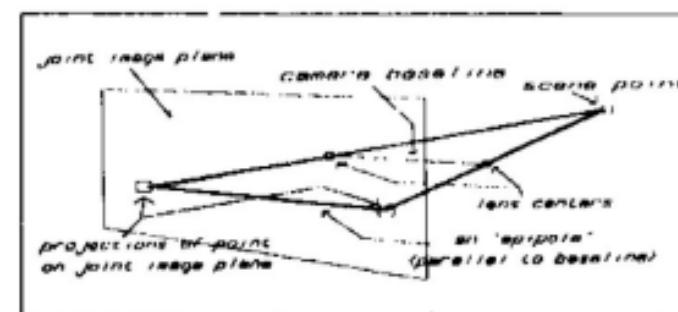
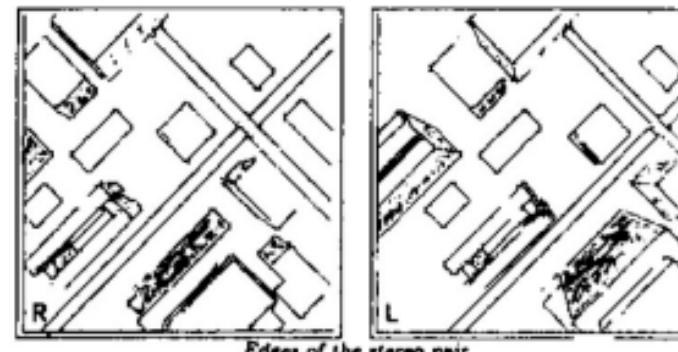
A Brief History of Computer Vision

1981: Binocular Scanline Stereo

- ▶ Correlate points along epipolar lines
- ▶ Use dynamic programming to introduce constraints along scanlines (image rows)
- ▶ Allows for overcoming ambiguities, but streaking artifacts between rows



Baker and Binford: Depth from Edge and Intensity Based Stereo. IJCAI, 1981.

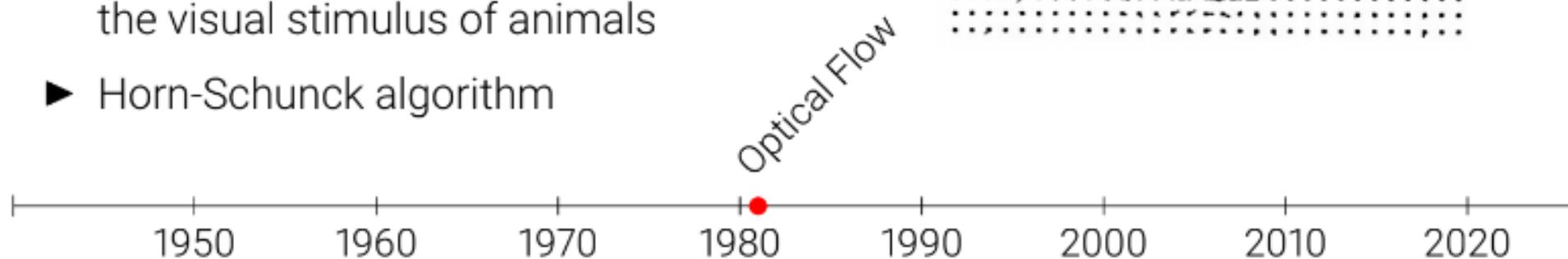


Two Camera Epipolar geometry
Figure 3

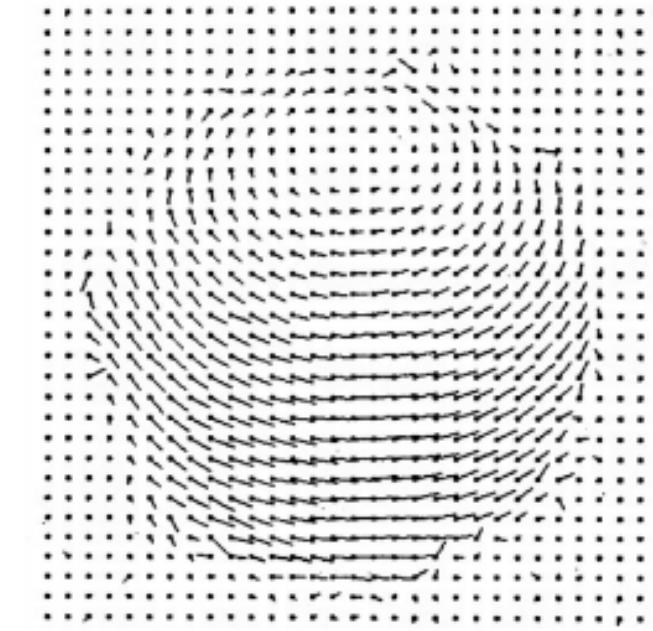
A Brief History of Computer Vision

1981: Dense Optical Flow

- ▶ Pattern of apparent motion of objects, surfaces, and edges in a visual scene
- ▶ Measured by (densely) tracking pixels between two frames
- ▶ Investigated by Gibson to describe the visual stimulus of animals
- ▶ Horn-Schunck algorithm



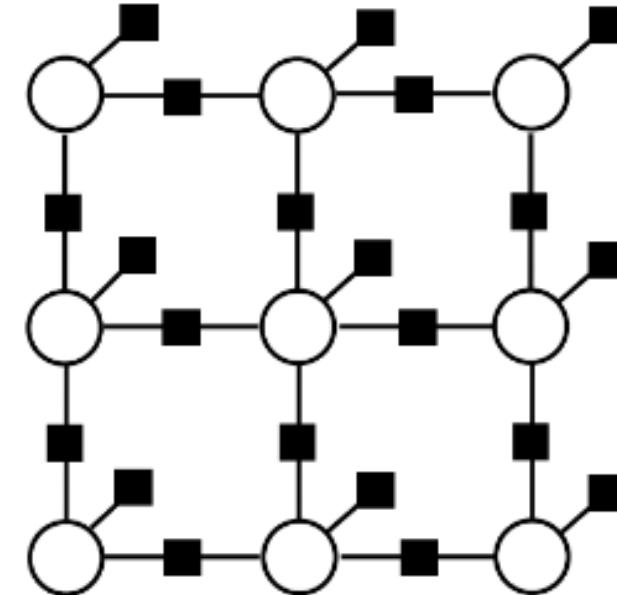
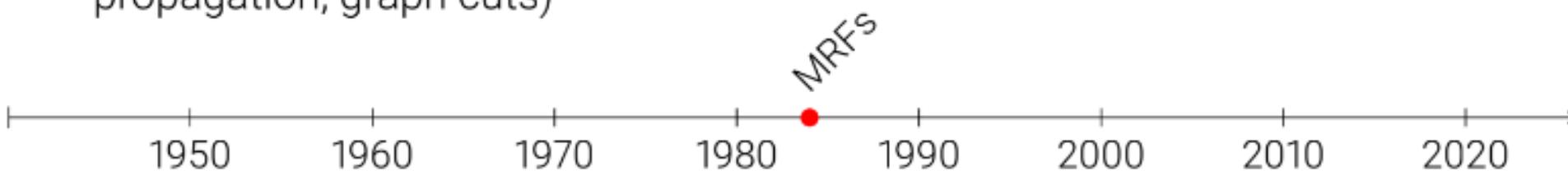
Horn and Schunck: Determining Optical Flow. Artificial Intelligence, 1981.



A Brief History of Computer Vision

1984: Markov Random Fields

- ▶ MRFs for encoding prior knowledge (e.g., about smoothness)
- ▶ Resolves ambiguities in many ill-posed vision problems (e.g., stereo, flow, denoising)
- ▶ Global optimization (e.g., variational inference, sampling, belief propagation, graph cuts)

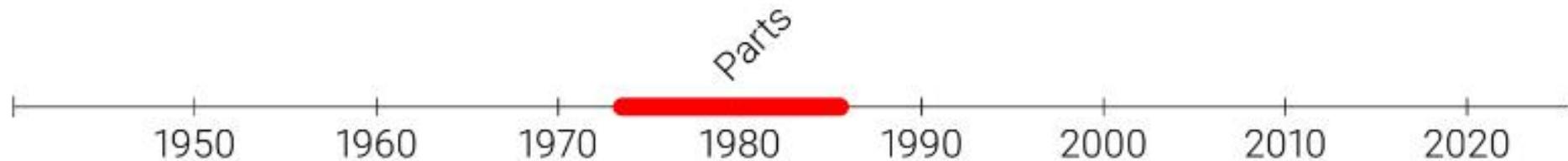
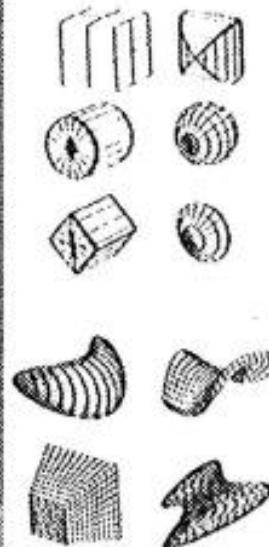
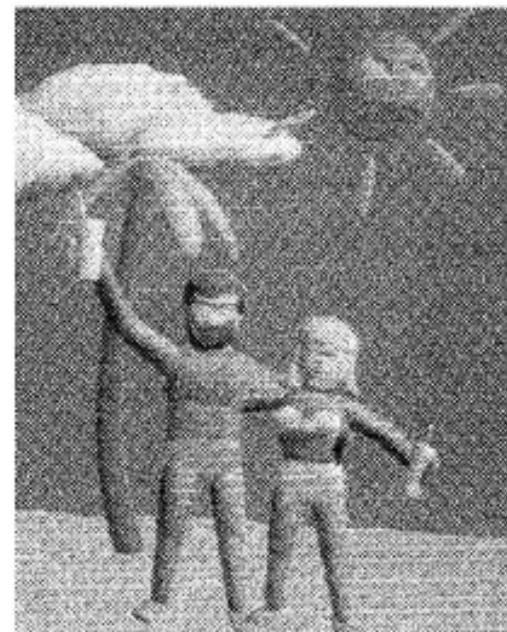


Geman and Geman: Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images. TPAMI, 1984.

A Brief History of Computer Vision

1980s: Part-based Models

- ▶ 1973: Pictorial Structures
- ▶ 1976: Generalized Cylinders
(solids of revolution, swept curves)
- ▶ 1986: Superquadrics
(generalization of quadric surfaces)
- ▶ Express complex relationships
- ▶ Compact representation

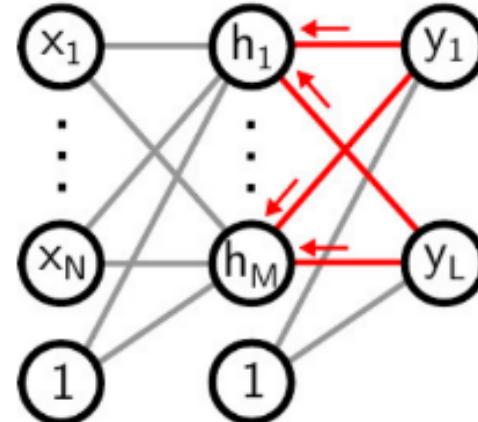
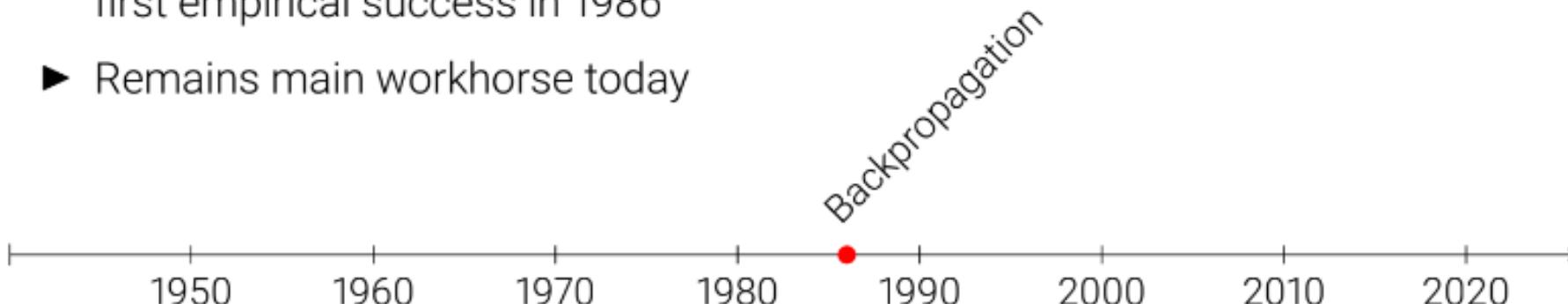


Pentland: Parts: Structured descriptions of shape. AAAI, 1986.

A Brief History of Computer Vision

1986: Backpropagation Algorithm

- ▶ Efficient calculation of gradients in a deep network wrt. network weights
- ▶ Enables application of gradient based learning to deep networks
- ▶ Known since 1961, but first empirical success in 1986
- ▶ Remains main workhorse today

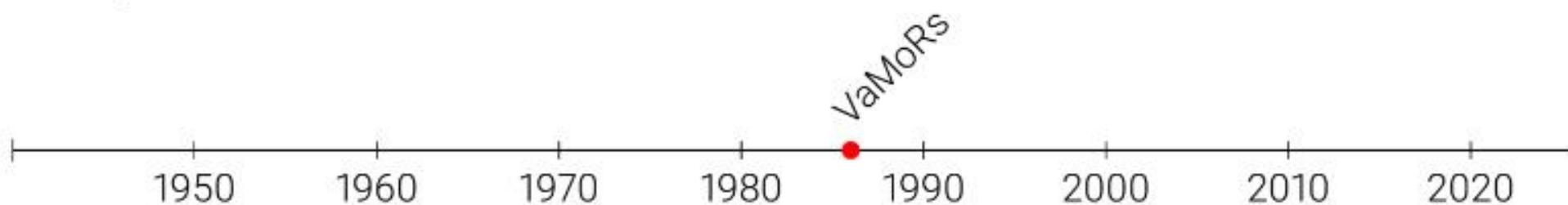


Rumelhart, Hinton and Williams: Learning representations by back-propagating errors. Nature, 1986.

A Brief History of Computer Vision

1986: Self-Driving Car VaMoRs

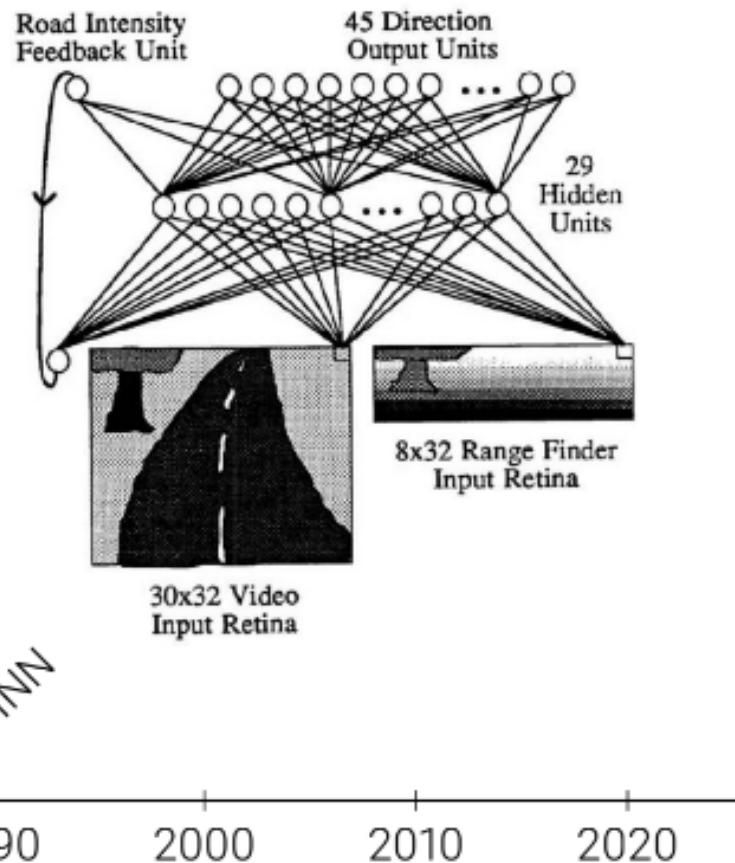
- ▶ Developed by Ernst Dickmanns in context of EUREKA-Prometheus
- ▶ Demonstration to Daimler-Benz Research 1986 in Stuttgart
- ▶ Longitudinal & lateral guidance with lateral acceleration feedback
- ▶ Speed: 0 to 36 km/h



A Brief History of Computer Vision

1988: Self-Driving Car ALVINN

- ▶ Forward-looking, vision based driving
- ▶ Fully connected neural network maps road images to vehicle turn radius
- ▶ Trained on simulated road images
- ▶ Tested on unlined paths, lined city streets and interstate highways
- ▶ 90 consecutive miles at up to 70 mph

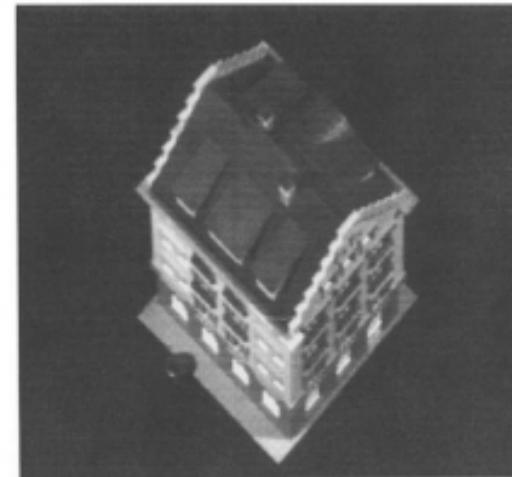


Pomerleau: ALVINN: An Autonomous Land Vehicle in a Neural Network. NIPS, 1988.

A Brief History of Computer Vision

1992: Structure-from-Motion

- ▶ Estimating 3D structures from 2D image sequences of static scenes
- ▶ Requires only a single camera
- ▶ Tomasi-Kanade factorization provides closed-form (SVD-based) solution for orthographic case
- ▶ Today: non-linear least squares

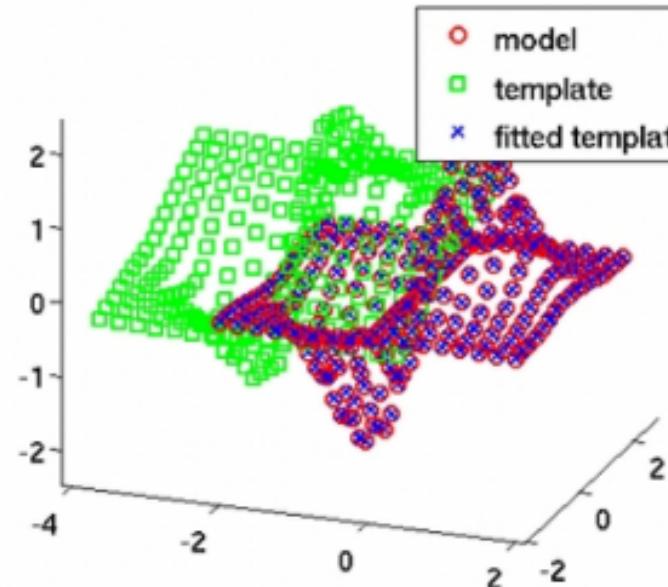


Tomasi and Kanade: Shape and motion from image streams under orthography: a factorization method. IJCV, 1992.

A Brief History of Computer Vision

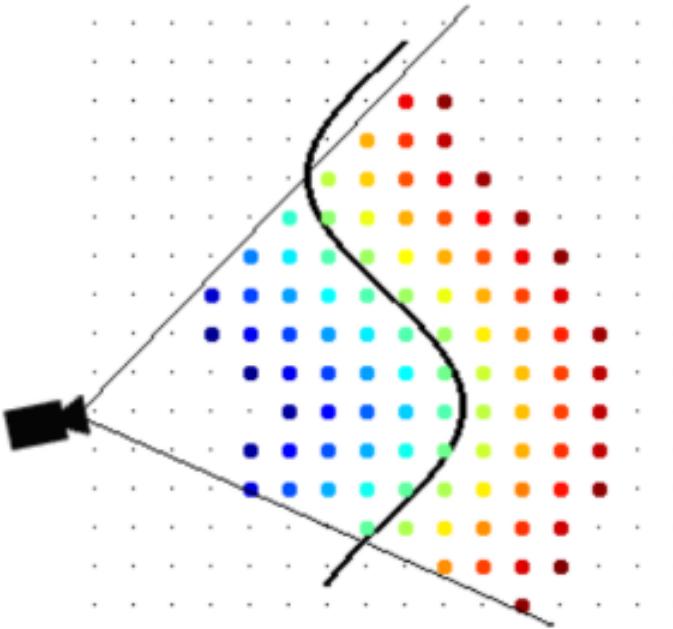
1992: Iterative Closest Points

- ▶ Registering two point clouds by iteratively optimizing a (rigid or non-rigid) transformation
- ▶ Used to aggregate partial 2D or 3D surfaces from different scans, to estimate relative camera poses from point clouds or to localize wrt. a map



Besl and McKay: A Method for Registration of 3-D Shapes. PAMI, 1992.

A Brief History of Computer Vision



1996: Volumetric Fusion

- ▶ Aggregation of multiple implicitly represented surfaces by averaging signed distance values
- ▶ Mesh-extraction as post-processing

Volumetric Fusion



Curless and Levoy: A Volumetric Method for Building Complex Models from Range Images. SIGGRAPH, 1996.



A Brief History of Computer Vision

1998: Multi-View Stereo

- ▶ 3D reconstruction from multiple input images using level-set methods
- ▶ Reconstruction vs. image matching
- ▶ Proper model of visibility
- ▶ Flexible topology
- ▶ Provable convergence
- ▶ Other approaches (dead-ends):
Voxel-coloring, space carving



Fig. 3. Multicamera images of 3D objects. On the left hand side, two crossing synthetic tori (24 images). On the right hand side, real images: two human heads (18 images).

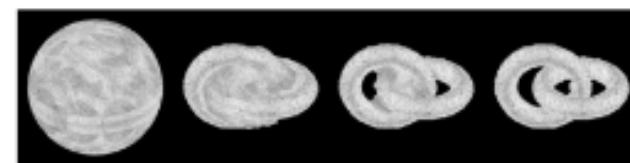
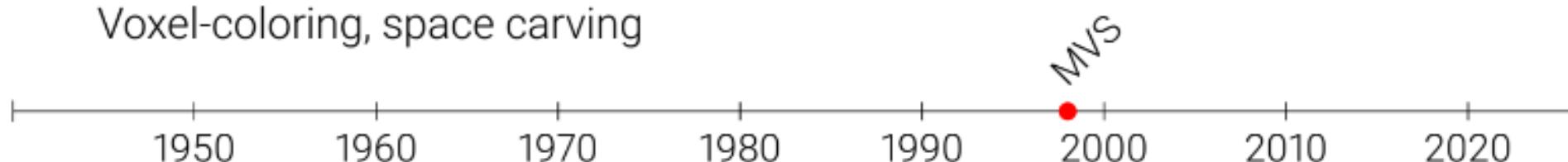


Fig. 4. Resolution of the surface for the two tori.



Faugeras and Keriven: Complete Dense Stereovision Using Level Set Methods. ECCV, 1998.

A Brief History of Computer Vision

1998: Stereo with Graph Cuts

- ▶ Popular discrete MAP inference algorithm for Markov Random Fields
- ▶ First versions included unary and pairwise terms
- ▶ Later versions also included specific forms of higher-order potentials
- ▶ Global reasoning compared to scanline stereo



Graph Cuts

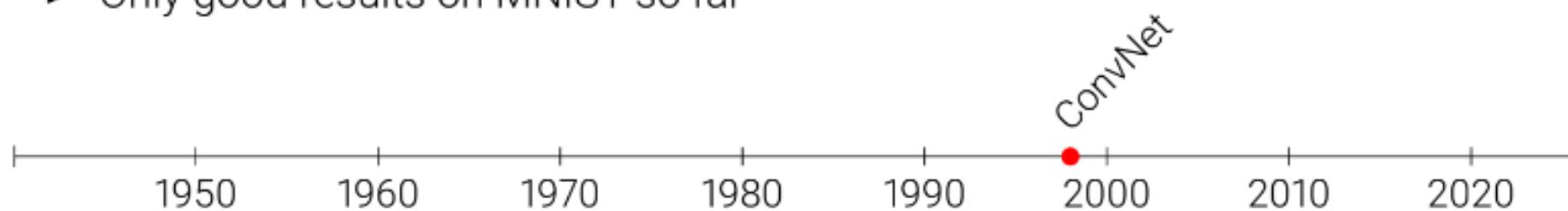
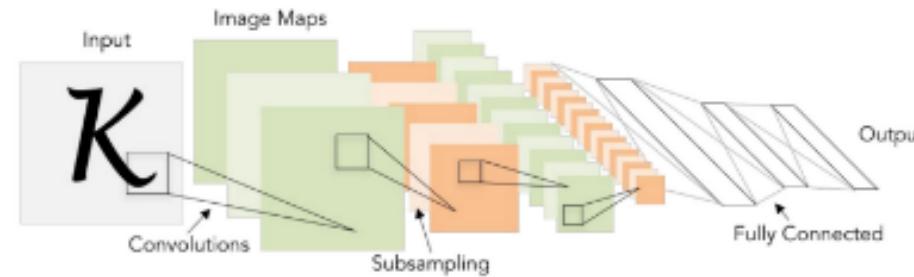


Boykov, Veksler and Zabih: Markov Random Fields with Efficient Approximations. CVPR, 1998.

A Brief History of Computer Vision

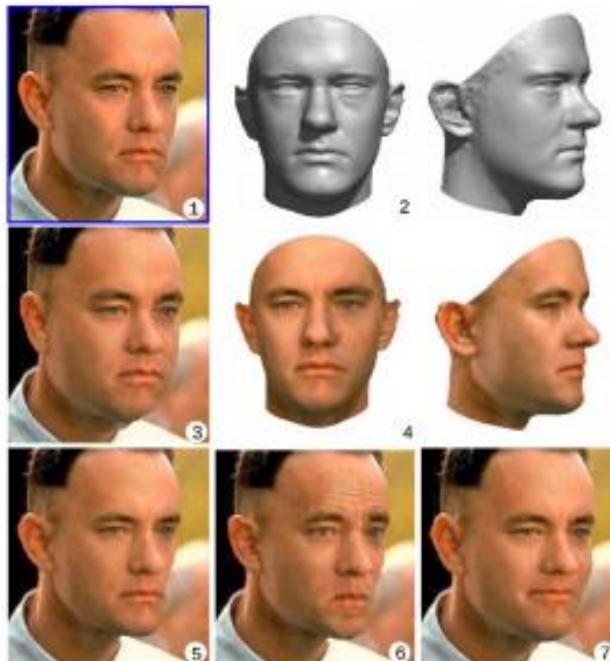
1998: Convolutional Neural Networks

- ▶ Similar to Neocognitron, but trained end-to-end using backpropagation
- ▶ Implements spatial invariance via convolutions and max-pooling
- ▶ Weight sharing reduces parameters
- ▶ Tanh/Softmax activations
- ▶ Only good results on MNIST so far



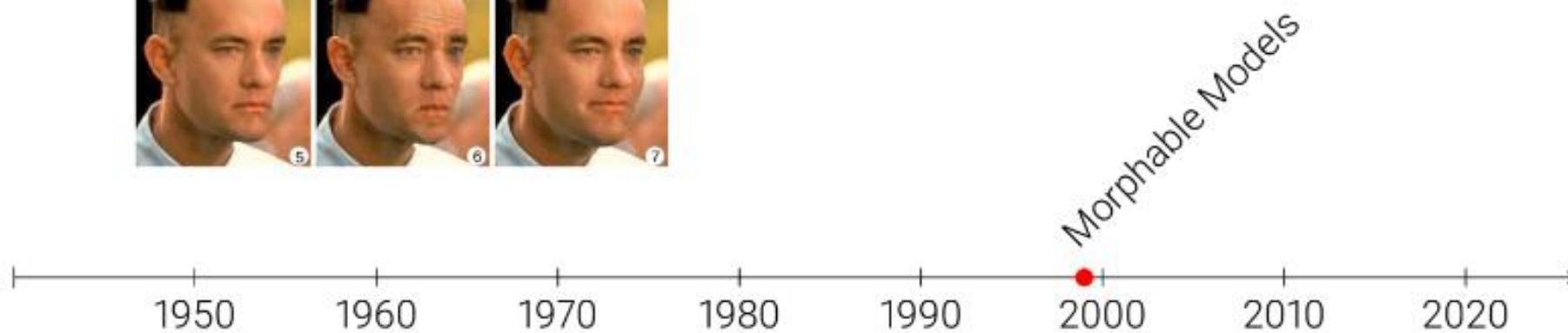
LeCun, Bottou, Bengio and Haffner: Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998.

A Brief History of Computer Vision



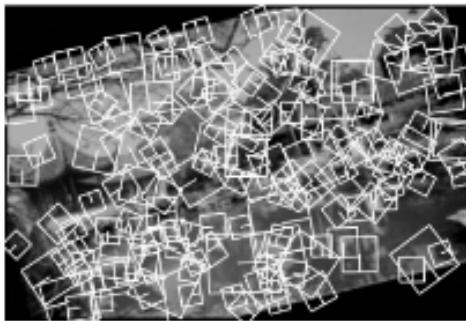
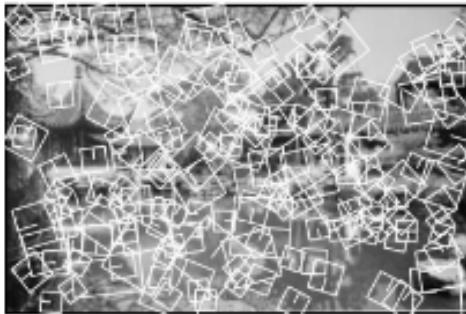
1999: Morphable Models

- ▶ Single-view 3D face reconstruction
- ▶ Linear combination of 200 laser scans of faces
- ▶ Stunning results at the time



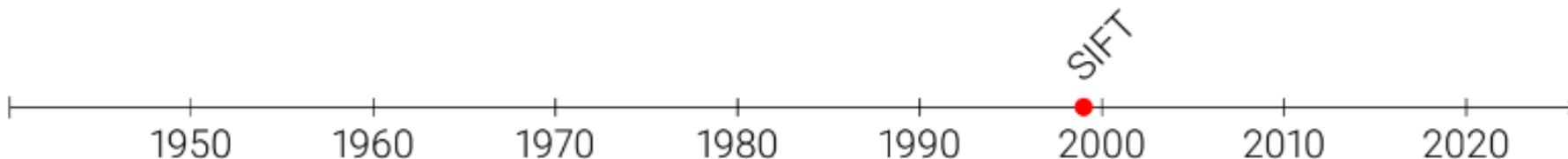
Blanz and Vetter: A Morphable Model for the Synthesis of 3D Faces. SIGGRAPH, 1999.

A Brief History of Computer Vision



1999: SIFT

- ▶ Scale Invariant Feature Transform
- ▶ Detection and description of salient local features in an image
- ▶ Enabled many applications (e.g., image stitching, reconstruction, motion estimation, ...)



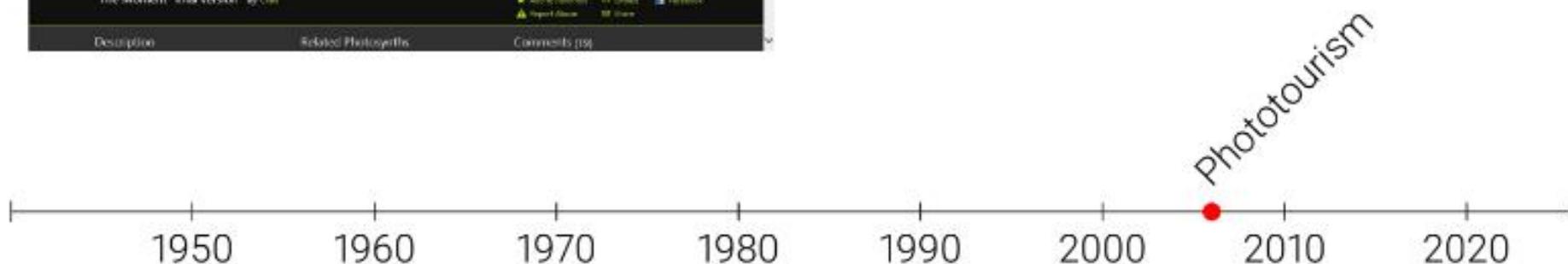
Lowe: Object Recognition from Local Scale-Invariant Features. ICCV, 1999.

A Brief History of Computer Vision



2006: Photo Tourism

- ▶ Large-scale 3D reconstruction from internet photos
- ▶ Key ingredients: SIFT feature matching, bundle adjustment
- ▶ Microsoft Photosynth (discont.)



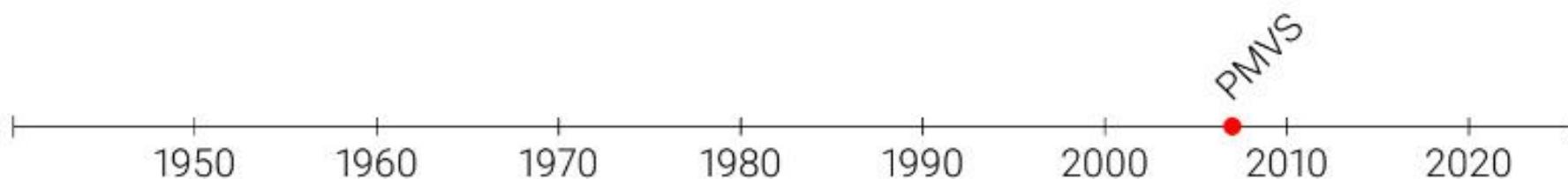
Snavely, Seitz and Szeliski: Photo tourism: exploring photo collections in 3D. SIGGRAPH, 2006.

A Brief History of Computer Vision



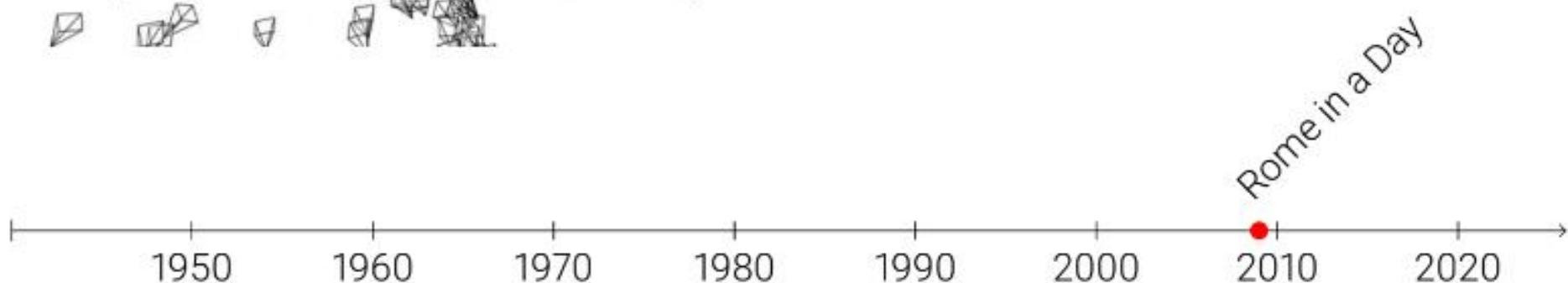
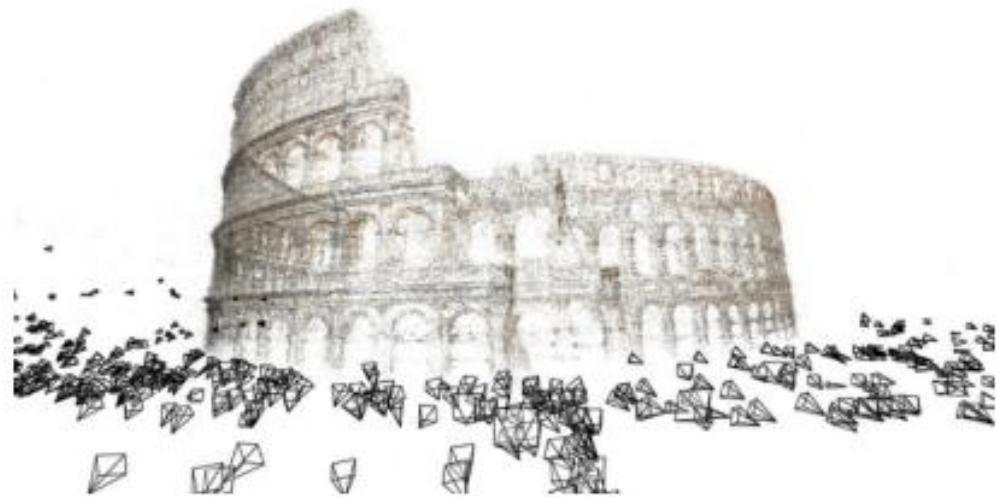
2007: PMVS

- ▶ Patch-based Multi View Stereo
- ▶ Robust reconstruction of various small and large objects
- ▶ Performance of 3D reconstruction techniques continues to increase



Furukawa and Ponce: Accurate, Dense, and Robust Multi-View Stereopsis. CVPR 2007.

A Brief History of Computer Vision

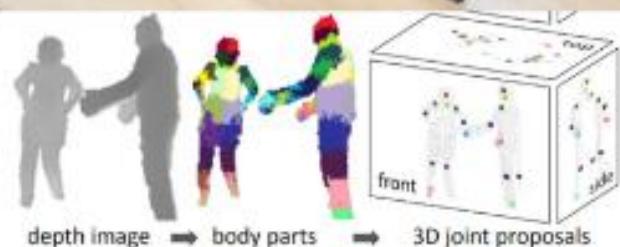


Agarwal, Snavely, Simon, Seitz and Szeliski: Building Rome in a day. ICCV, 2009.

2009: Building Rome in a Day

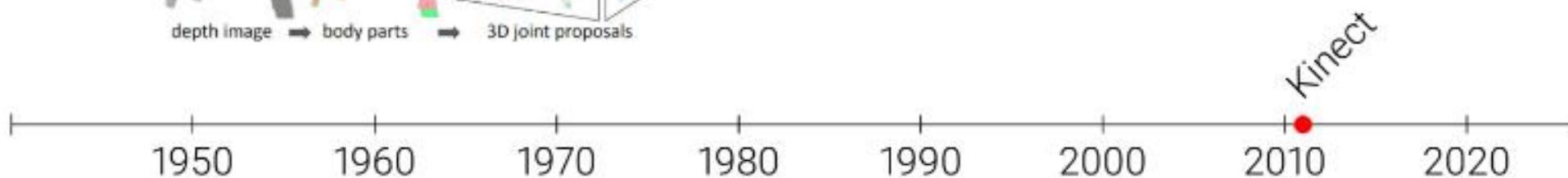
- ▶ 3D reconstruction of landmarks and cities from unstructured Internet photo-collections
- ▶ Follow-up: Rome on a Cloudless Day

A Brief History of Computer Vision



2011: Kinect

- ▶ Active light 3D sensing
- ▶ ML for 3D pose estimation
- ▶ Multiple hardware generations
- ▶ Early versions failed to commercialize but heavily used for robotics and vision research



Shotton et al.: Real-time human pose recognition in parts from single depth images. CVPR, 2011.

A Brief History of Computer Vision

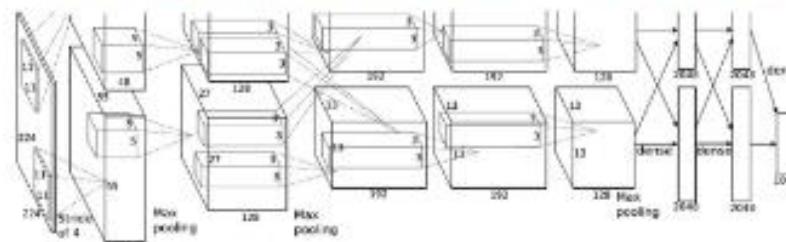
2009-2012: ImageNet and AlexNet

ImageNet

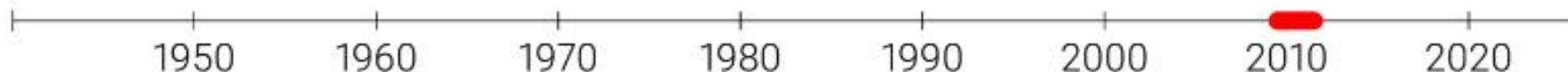
- ▶ Recognition benchmark (ILSVRC)
- ▶ 10 million annotated images
- ▶ 1000 categories

AlexNet

- ▶ First neural network to win ILSVRC via **GPU training, deep models, data**



Image/AlexNet



Krizhevsky, Sutskever and Hinton: ImageNet classification with deep convolutional neural networks. NIPS, 2012.

A Brief History of Computer Vision

2009-2012: ImageNet and AlexNet

ImageNet

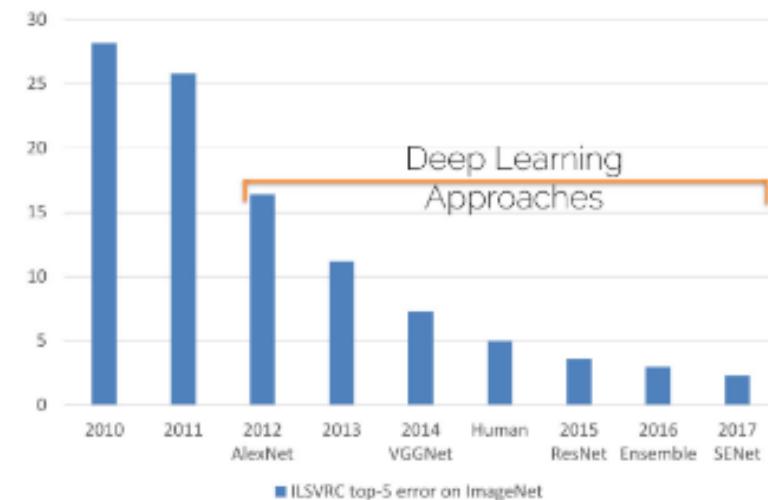
- ▶ Recognition benchmark (ILSVRC)
- ▶ 10 million annotated images
- ▶ 1000 categories

AlexNet

- ▶ First neural network to win ILSVRC via **GPU training, deep models, data**
- ▶ Sparked deep learning revolution



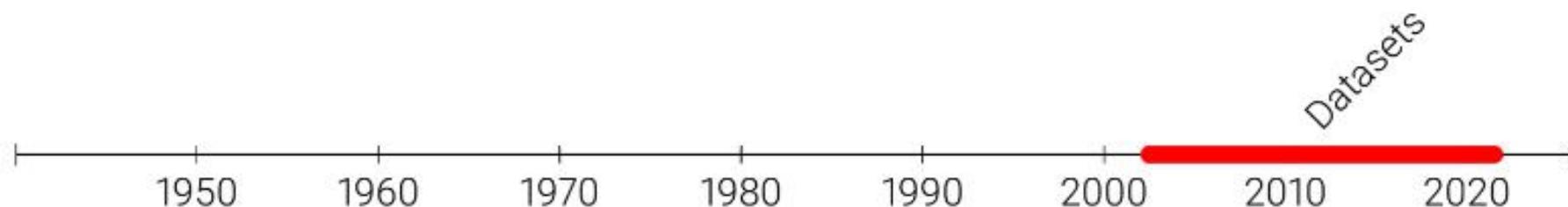
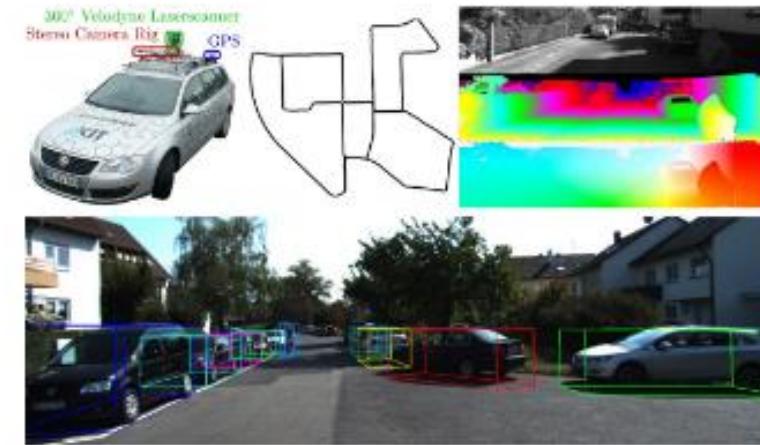
Krizhevsky, Sutskever and Hinton: ImageNet classification with deep convolutional neural networks. NIPS, 2012.



A Brief History of Computer Vision

2002-now: Golden Age of Datasets

- ▶ Middlebury Stereo and Flow
- ▶ KITTI, Cityscapes: Self-driving
- ▶ PASCAL, MS COCO: Recognition
- ▶ ShapeNet, ScanNet: 3D DL
- ▶ Visual Genome: Vision/Language
- ▶ MITOS: Breast cancer

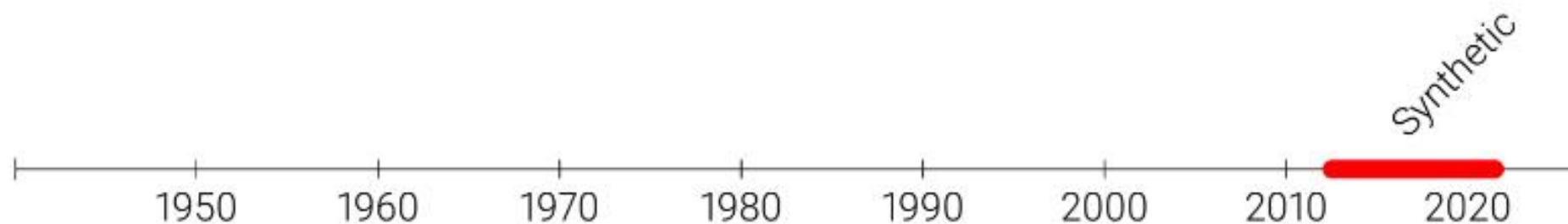


Geiger, Lenz and Urtasun: Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. CVPR, 2012.

A Brief History of Computer Vision

2012-now: Synthetic Data

- ▶ Annotating real data is expensive
- ▶ Led to surge of synthetic datasets
- ▶ Creating 3D assets is also costly
- ▶ But even very simple 3D datasets proved tremendously useful for pre-training (e.g., in optical flow)



Dosovitskiy et al · FlowNet: Learning Optical Flow with Convolutional Networks · ICCV 2015

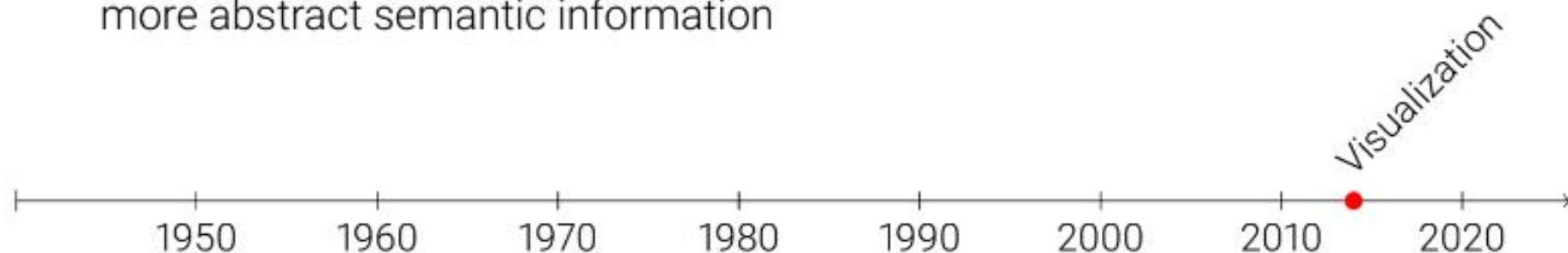
A Brief History of Computer Vision

2014: Visualization

- ▶ Goal: provide insights into what the network (black box) has learned
- ▶ Visualized image regions that most strongly activate various neurons at different layers of the network
- ▶ Found that higher levels capture more abstract semantic information

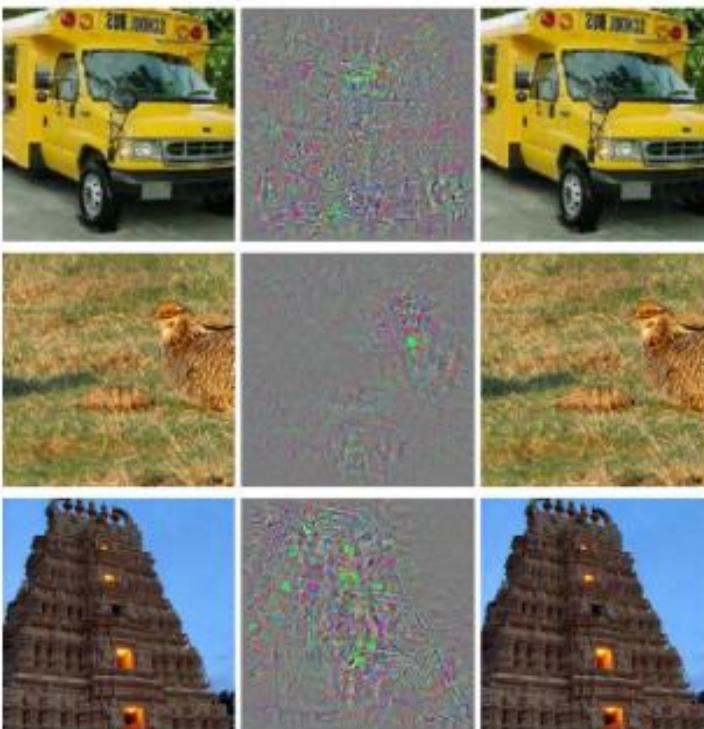


Layer 1 Layer 2 Layer 5



Zeiler and Fergus: CNN Features Off-the-Shelf: An Astounding Baseline for Recognition. CVPR Workshops, 2014.

A Brief History of Computer Vision



2014: Adversarial Examples

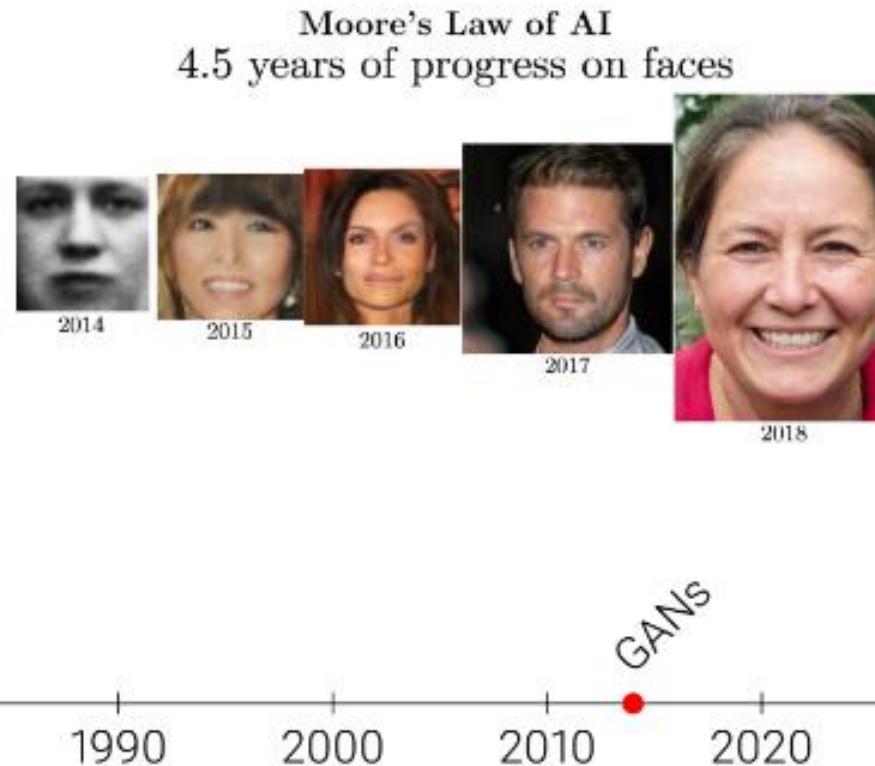
- ▶ Accurate image classifiers can be fooled by imperceptible changes (here magnified for visibility)
- ▶ All images in the right column are classified as “ostrich”

Szegedy et al.: Intriguing properties of neural networks. ICLR, 2014.

A Brief History of Computer Vision

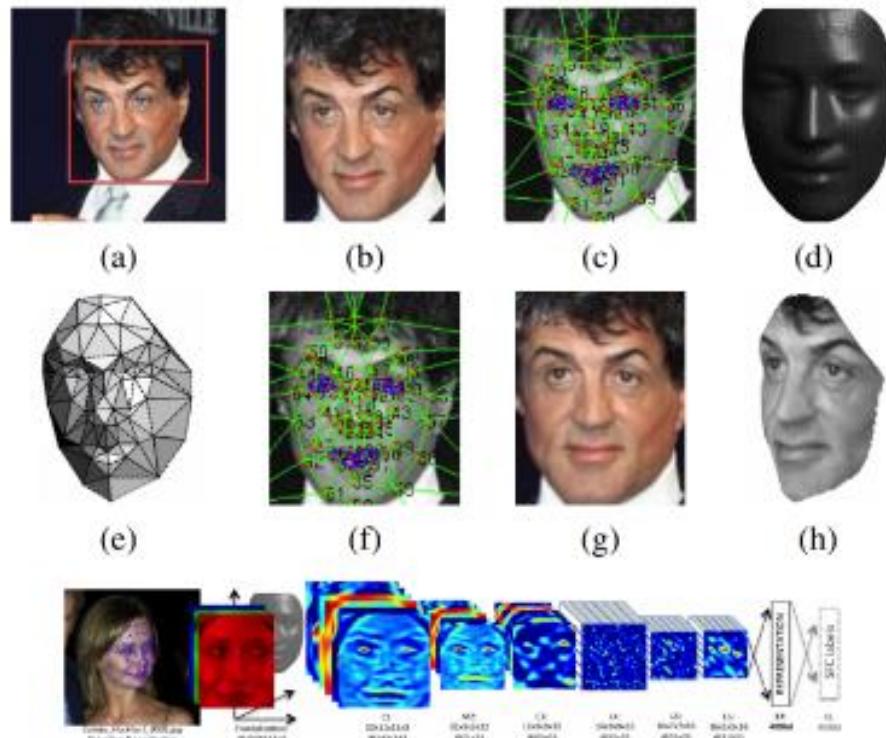
2014: Generative Adversarial Networks

- Deep generative models (VAEs, GANs) produce compelling images
- StyleGAN2 is state-of-the-art
- Results on faces hard to distinguish from real images
- Active research on image translation, domain adaptation, content and scene generation and 3D GANs



Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, Bengio: Generative Adversarial Networks. NIPS, 2014.

A Brief History of Computer Vision

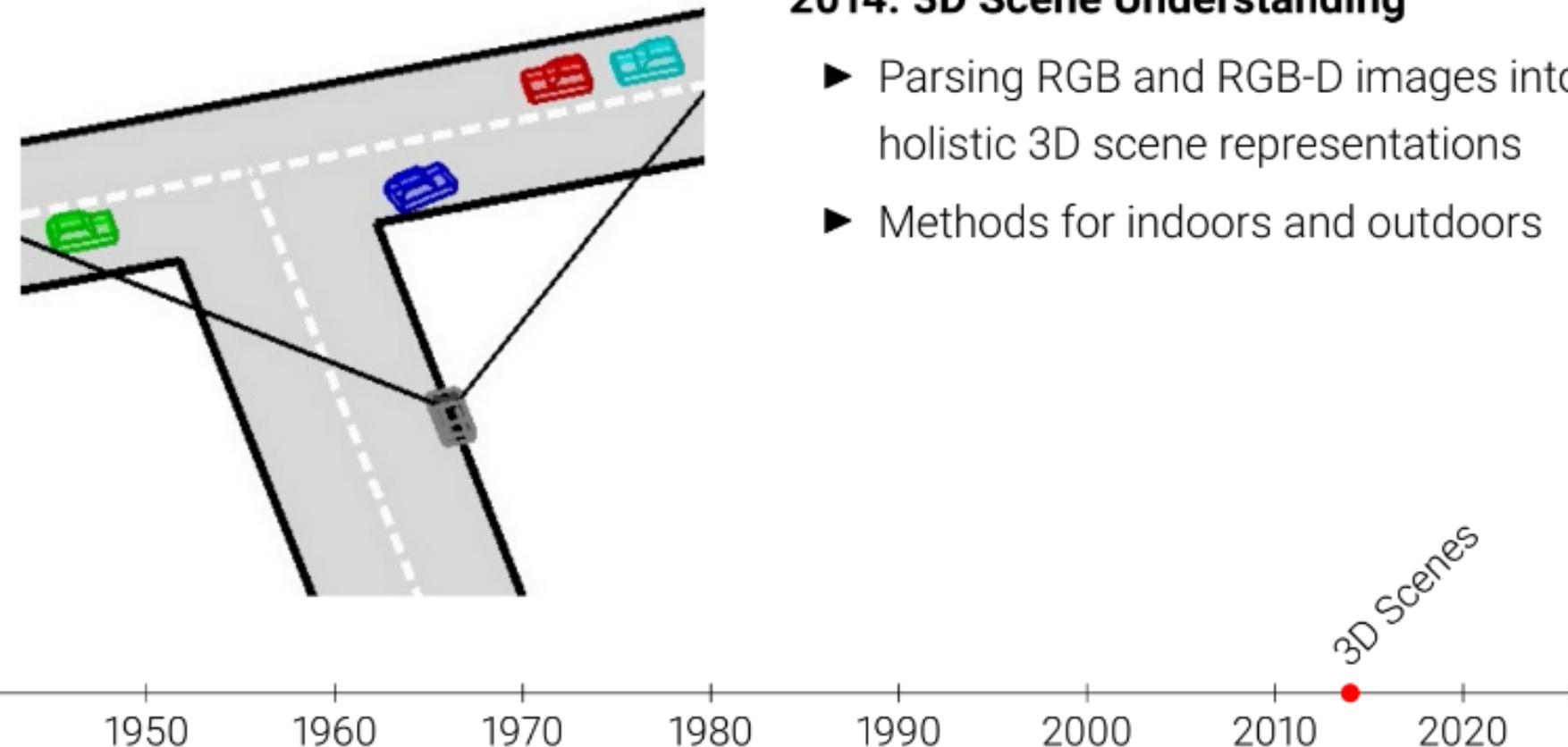


2014: DeepFace

- ▶ Combination of model-based alignment with deep learning for face recognition
- ▶ First model to reach human-level face recognition performance

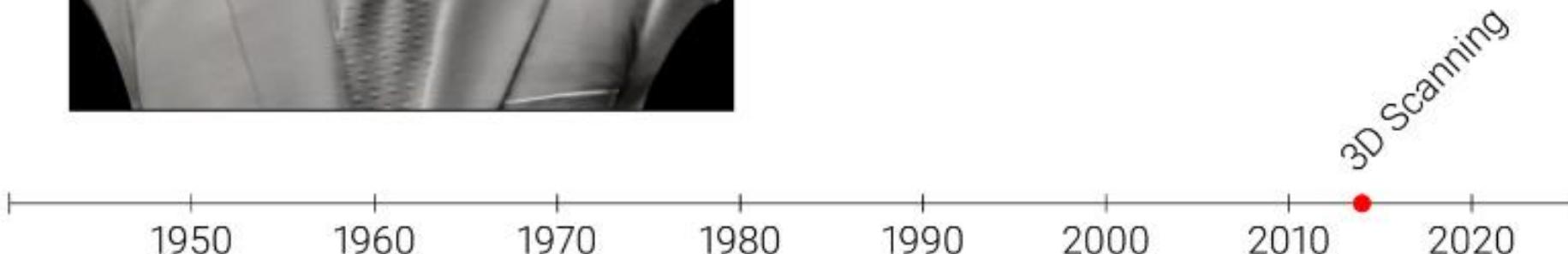
Taigman, Yang, Ranzato and Wolf: DeepFace: Closing the Gap to Human-Level Performance in Face Verification. CVPR, 2014.

A Brief History of Computer Vision



Geiger, Lauer, Wojek, Stiller and Urtasun: 3D Traffic Scene Understanding From Movable Platforms. PAMI, 2014.

A Brief History of Computer Vision

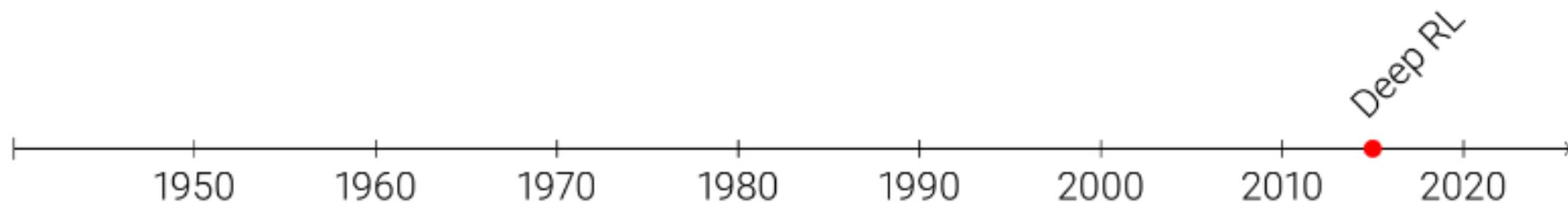
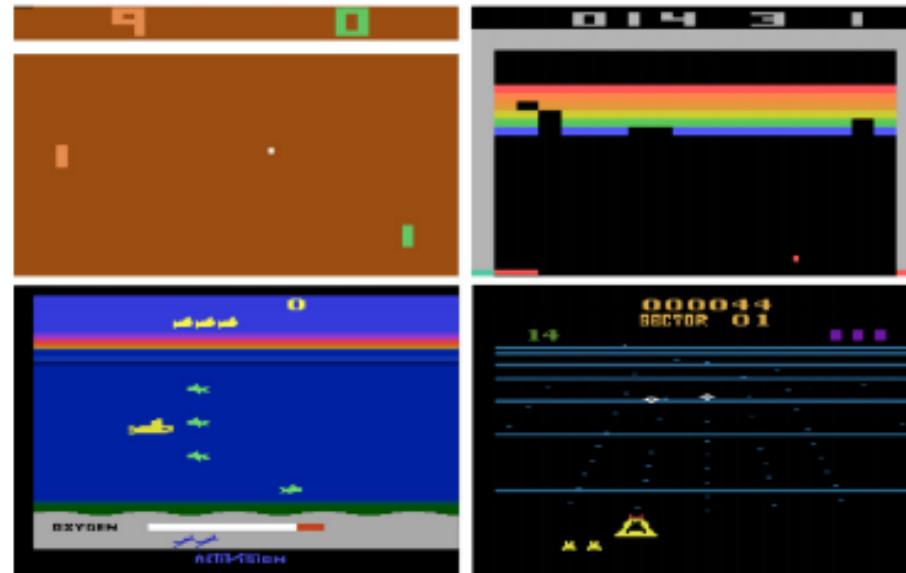


Metallo et al.: Scanning and printing a 3D portrait of president Barack Obama. SIGGRAPH Studio, 2015.

A Brief History of Computer Vision

2015: Deep Reinforcement Learning

- ▶ Learning a policy ($\text{state} \rightarrow \text{action}$) through random exploration and reward signals (e.g., game score)
- ▶ No other supervision
- ▶ Success on many Atari games
- ▶ But some games remain hard



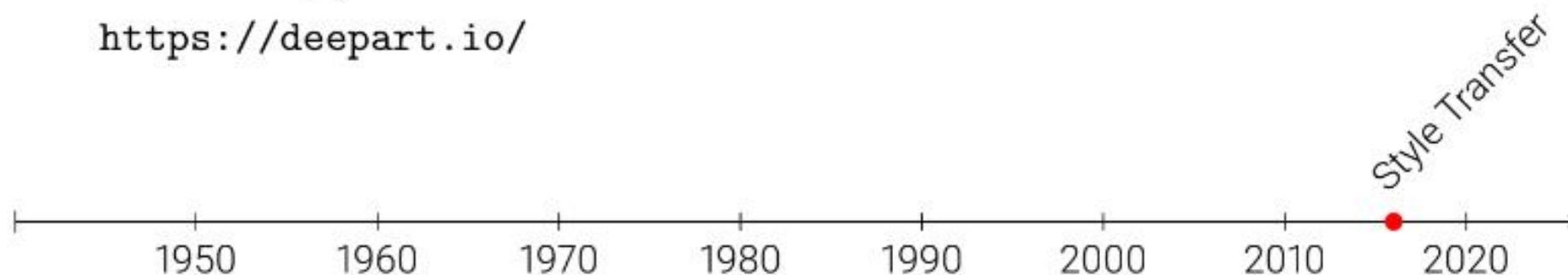
Mnih et al.: Human-level control through deep reinforcement learning. Nature, 2015.

A Brief History of Computer Vision

2016: Style Transfer

- ▶ Manipulate photograph to adopt style of another image (painting)
- ▶ Uses deep network pre-trained on ImageNet for disentangling content from style
- ▶ It is fun! Try yourself:

<https://deepart.io/>

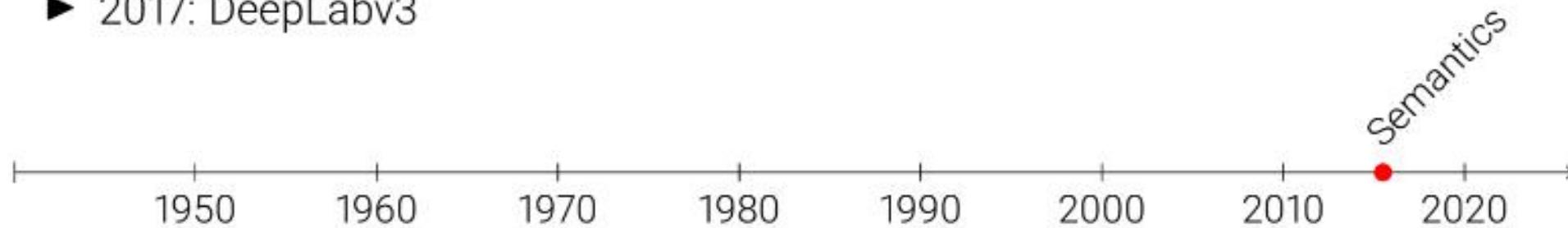
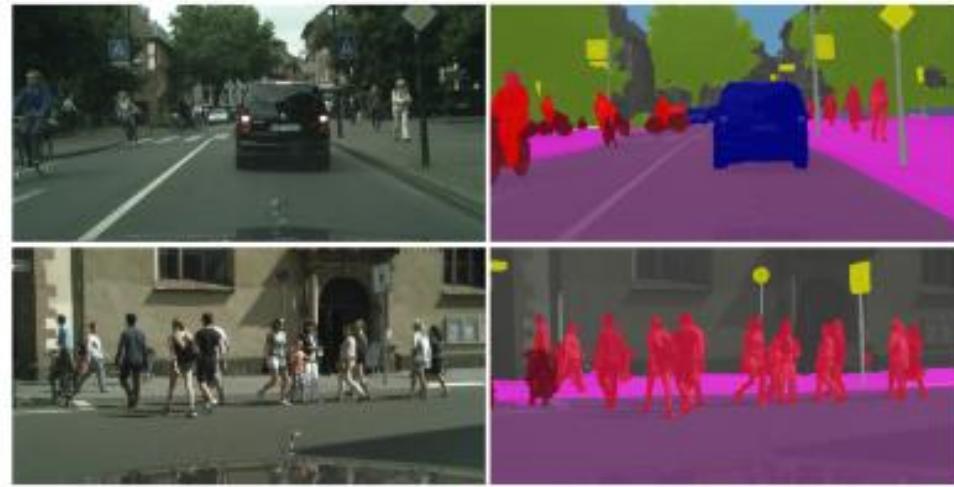


Gatys, Ecker and Bethge: Image Style Transfer Using Convolutional Neural Networks. CVPR, 2016.

A Brief History of Computer vision

2015-2017: Semantic Segmentation

- ▶ Assign semantic class to every pixel
- ▶ Semantic segmentation starts to work on challenging real-world datasets (e.g., CityScapes)
- ▶ 2015: FCN, SegNet
- ▶ 2016: DeepLab, FSO
- ▶ 2017: DeepLabv3

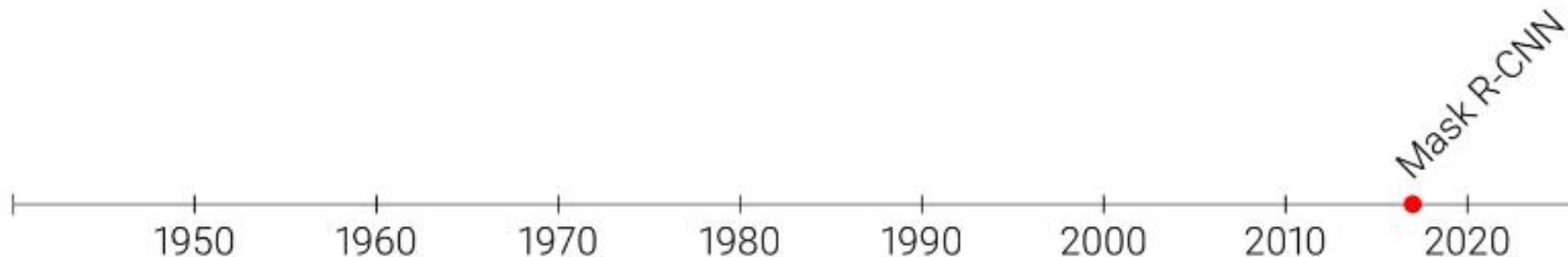
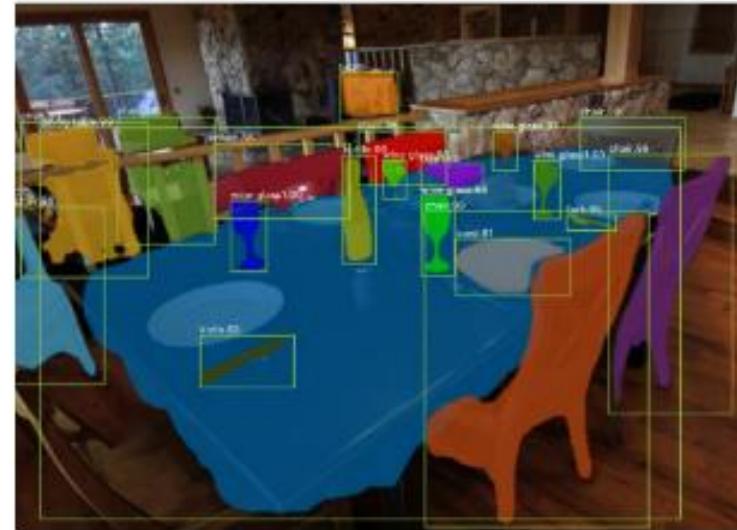


Kundu, Vineet and Koltun: Feature Space Optimization for Semantic Video Segmentation. CVPR, 2016.

A Brief History of Computer Vision

2017: Mask R-CNN

- ▶ Deep neural network for joint object detection and instance segmentation
- ▶ Outputs “structured object”, not only a single number (class label)
- ▶ State-of-the-art on MS-COCO

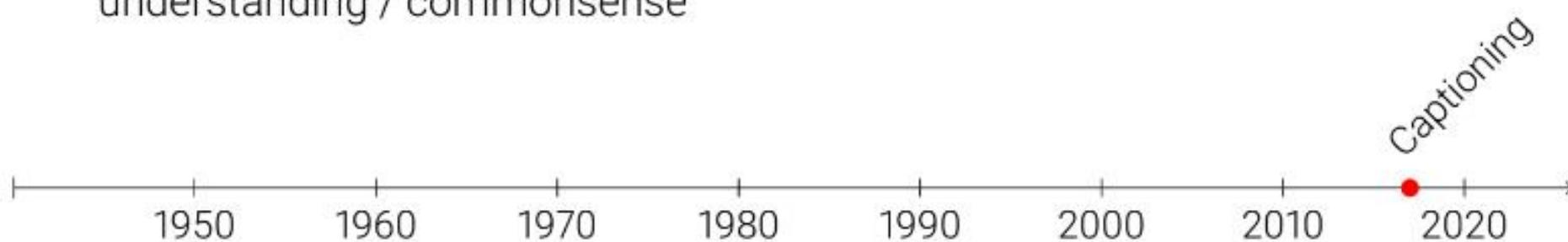


He, Gkioxari, Dollár and Ross Girshick: Mask R-CNN. ICCV, 2017.

A Brief History of Computer Vision

2017: Image Captioning

- Growing interest in combining vision with language
- Several new tasks emerged including image captioning and visual question answering
- However, models still lack understanding / commonsense



Karpathy and Fei-Fei: Deep Visual-Semantic Alignments for Generating Image Descriptions. PAMI, 2017.

A Brief History of Computer Vision



2018: Human Shape and Pose

- ▶ Human pose/shape models mature
- ▶ Rich parametric models (e.g., SMPL, STAR)
- ▶ Regression from RGB images only
- ▶ Models of pose-dependent deformation and clothing

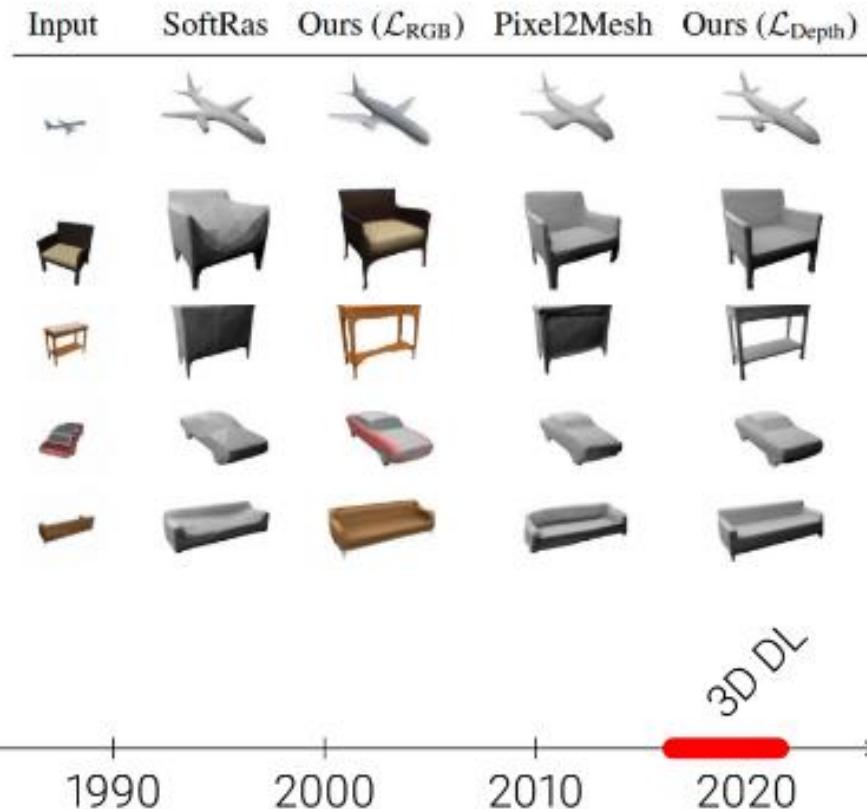


Kanazawa, Black, Jacobs and Malik: End-to-End Recovery of Human Shape and Pose. CVPR, 2018.

A Brief History of Computer Vision

2016-2020: 3D Deep Learning

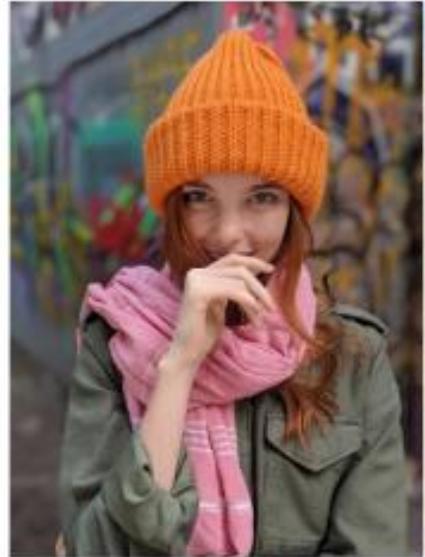
- ▶ First deep models to output 3D representations
- ▶ Voxels, point clouds, meshes, implicit representations
- ▶ Prediction of 3D models even from a single image
- ▶ Geometry, materials, light, motion



Niemeyer, Mescheder, Oechsle, Geiger: Differentiable Volumetric Rendering: Learning Implicit 3D Representations without 3D Supervision. CVPR, 2020.

A Brief History of Computer Vision

Applications and Commercial Products



Google Portrait Mode



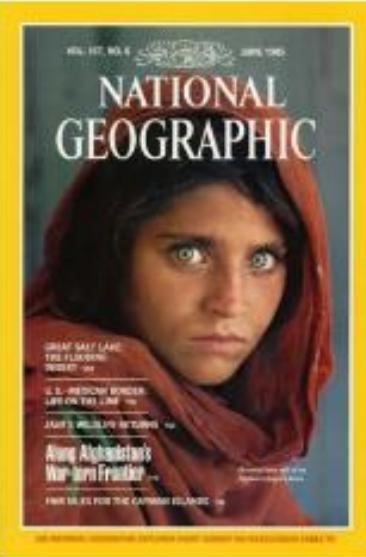
Skydio 2 Drone



Self-Driving Cars



Microsoft Hololens



Iris Recognition



A Brief History of Computer Vision

Current Challenges

- ▶ Un-/Self-Supervised Learning
- ▶ Interactive learning
- ▶ Accuracy (e.g., self-driving)
- ▶ Robustness and generalization
- ▶ Inductive biases
- ▶ Understanding and mathematics
- ▶ Memory and compute
- ▶ Ethics and legal questions

