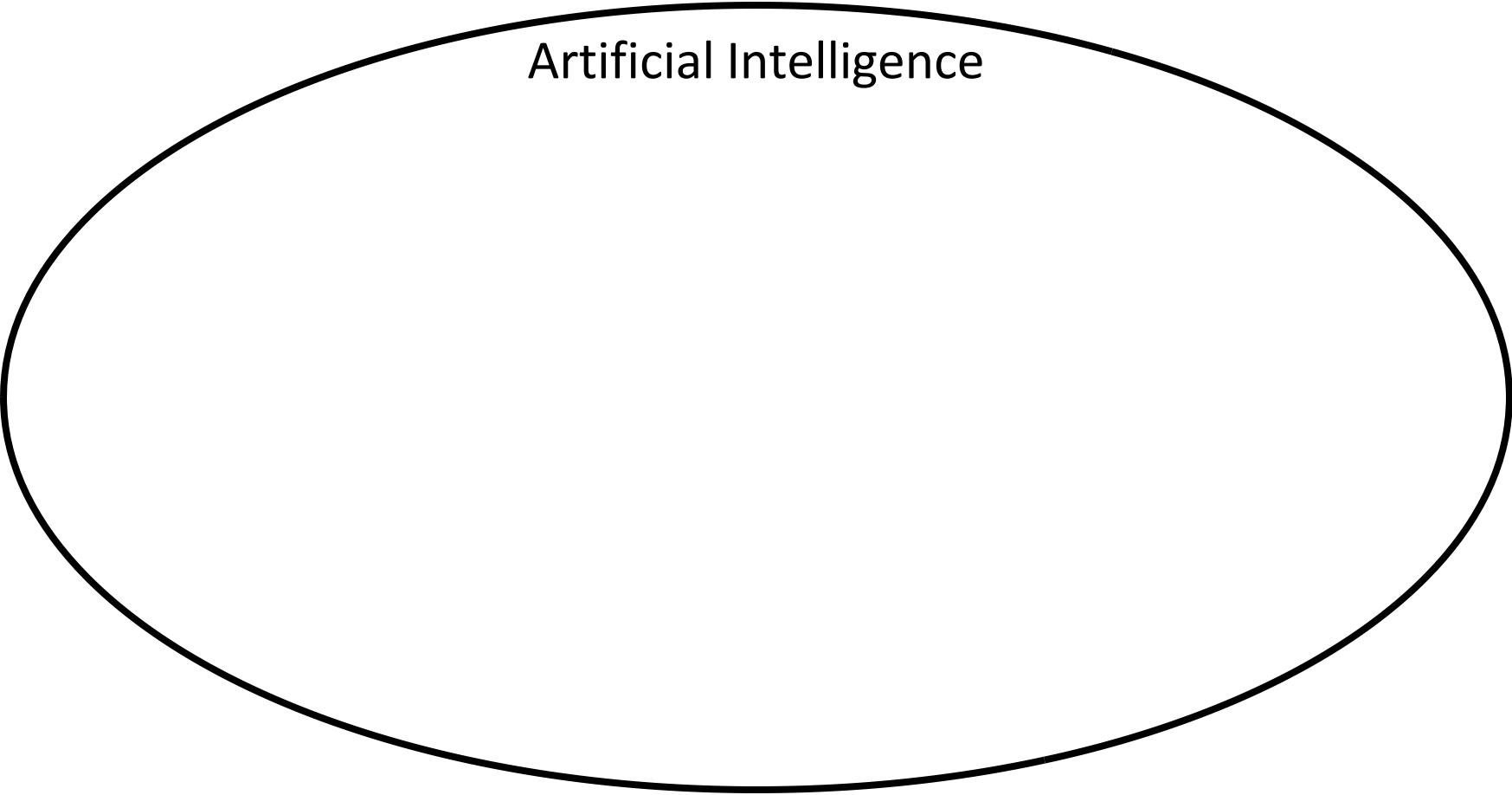




CS231n: Deep Learning for Computer Vision

Lecture 1: Introduction



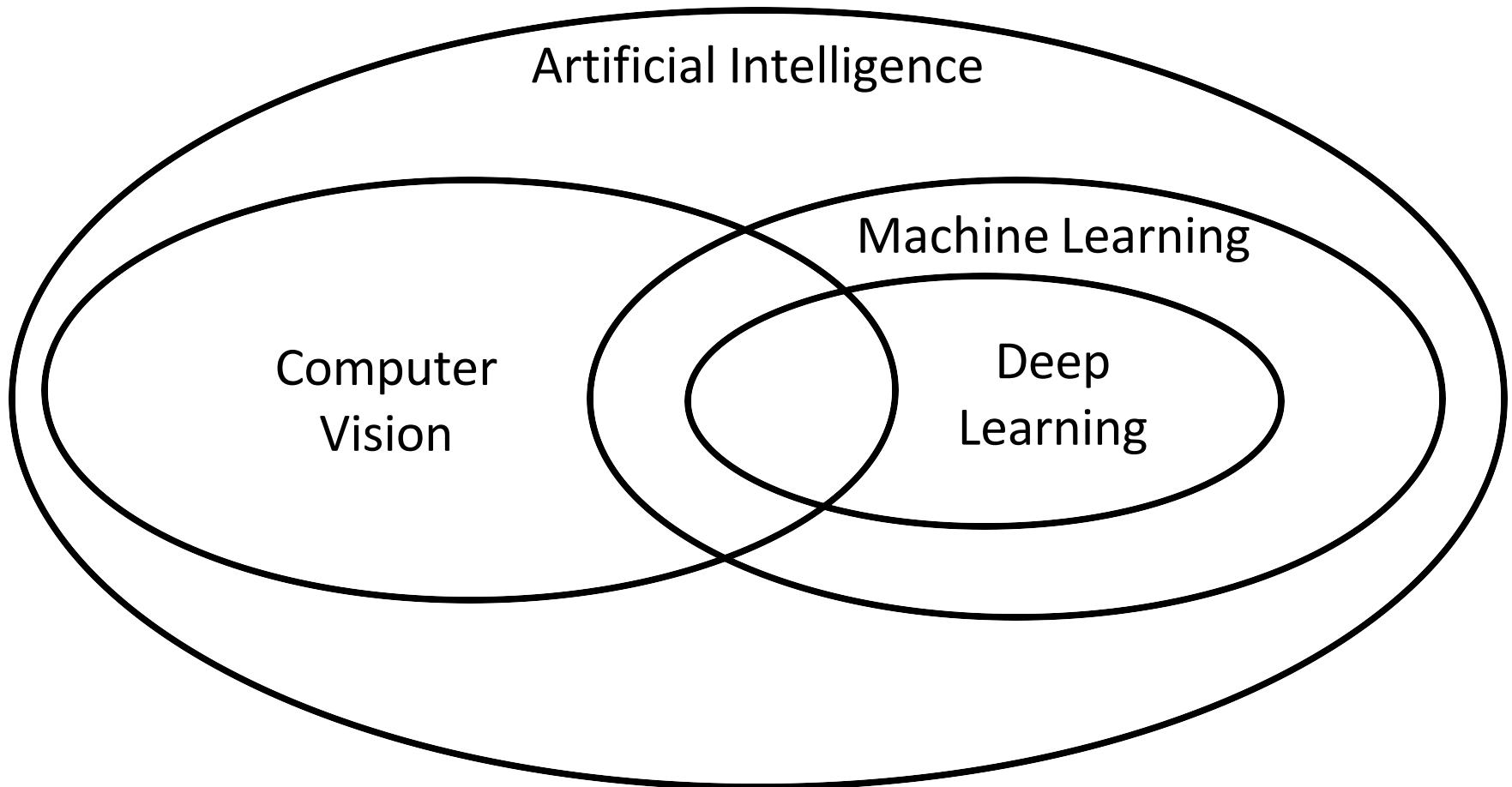
Artificial Intelligence

Artificial Intelligence

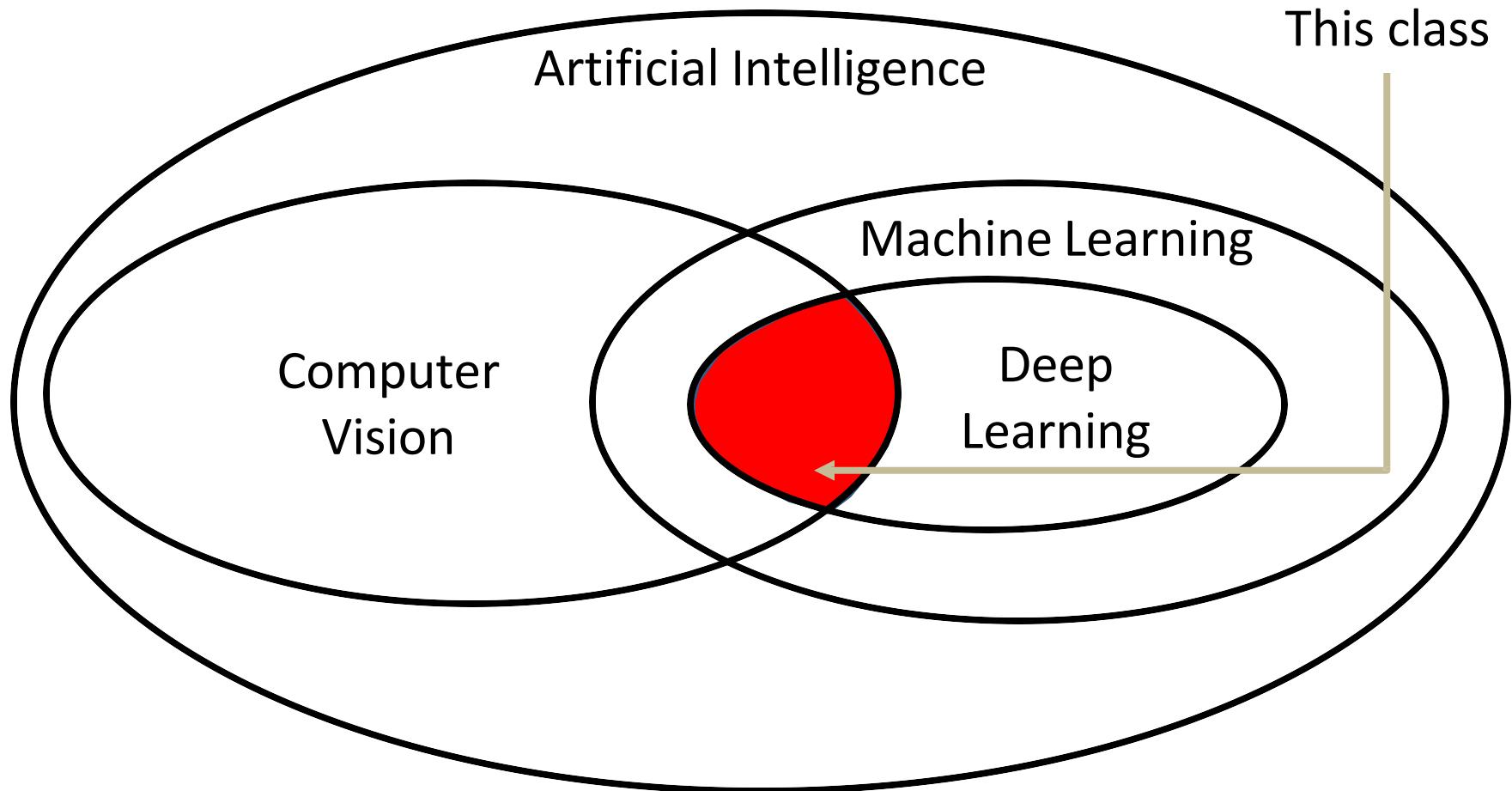
Machine Learning

Computer
Vision

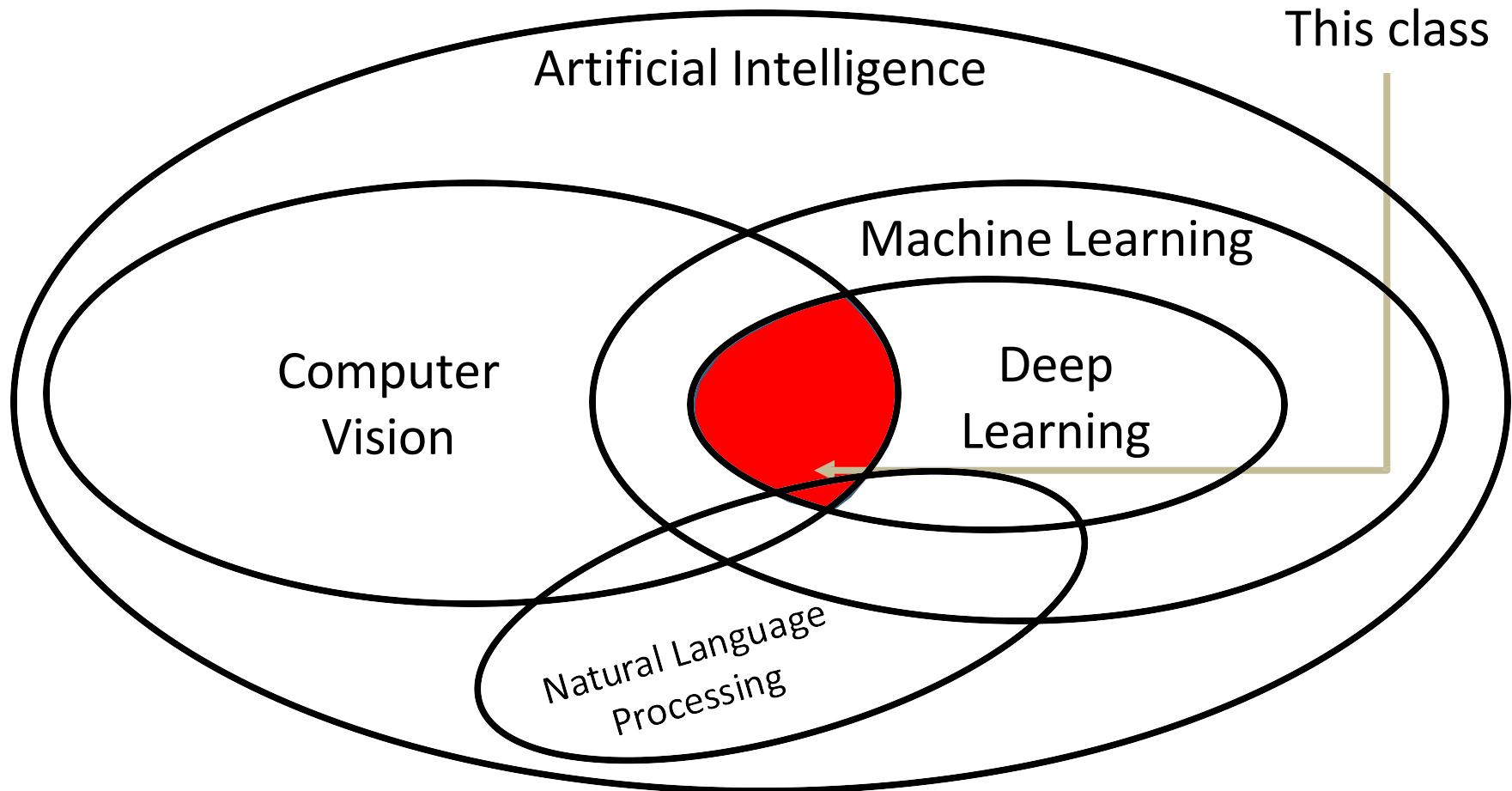
Artificial Intelligence



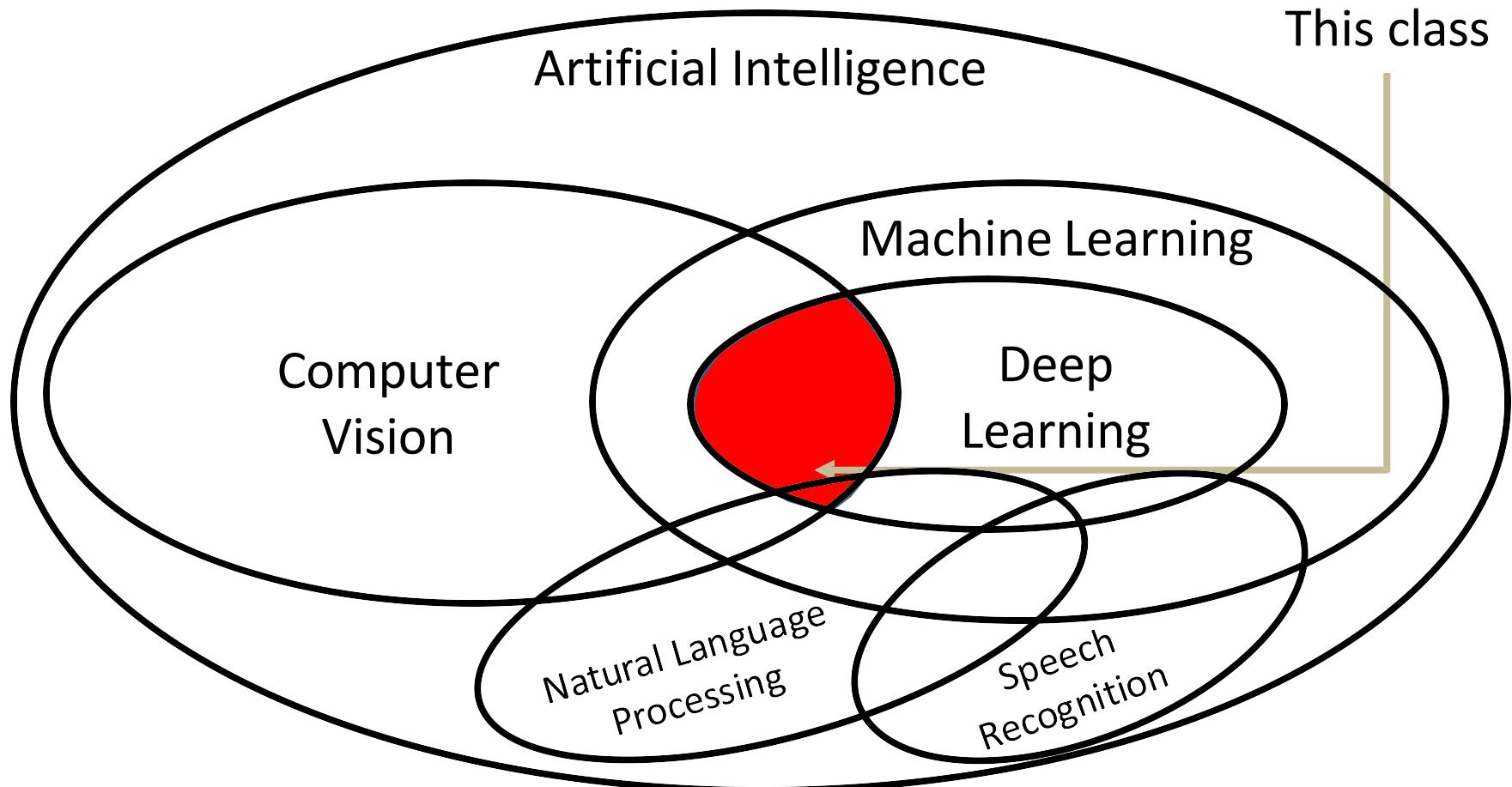
Slide inspiration: Justin Johnson



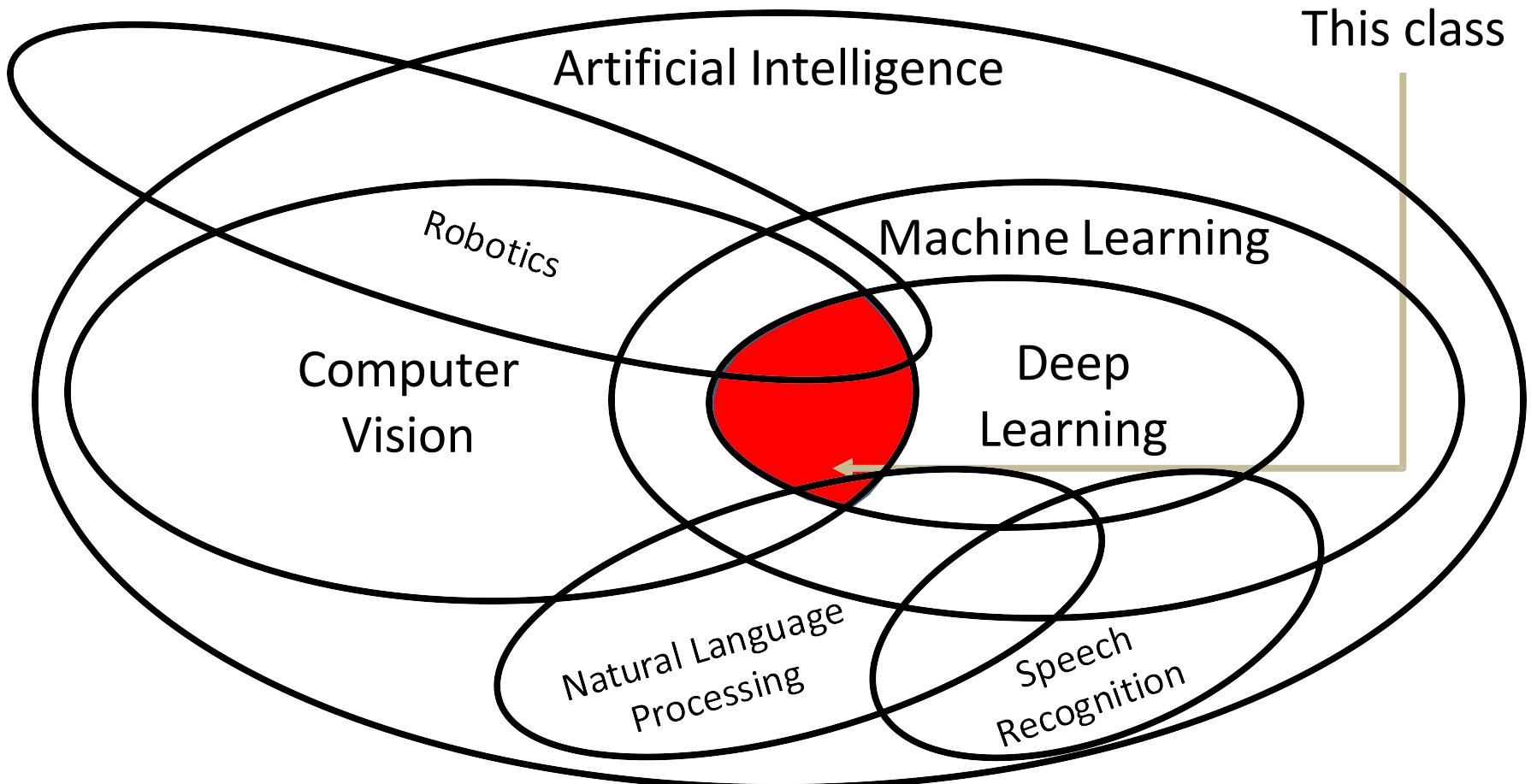
Slide inspiration: Justin Johnson



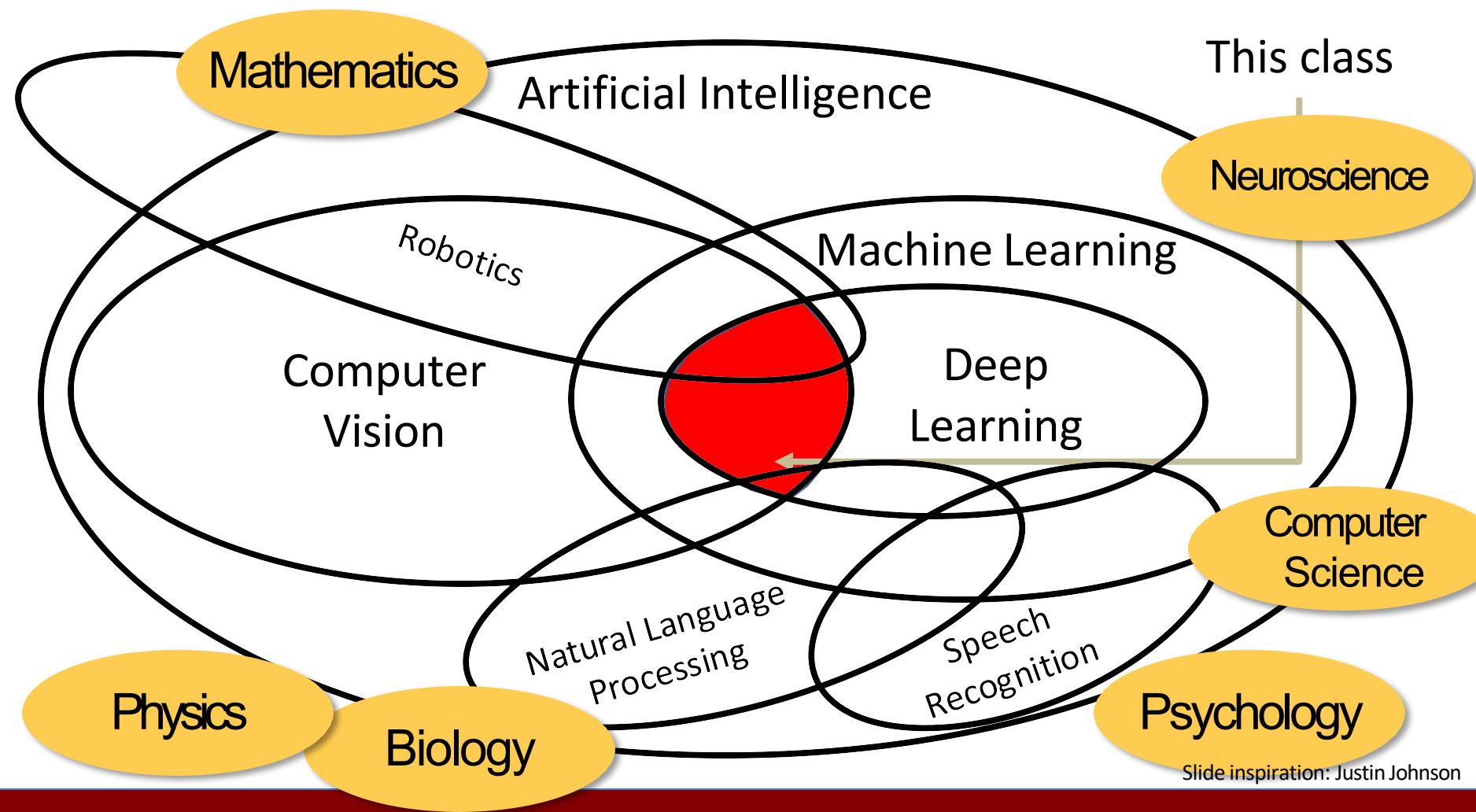
Slide inspiration: Justin Johnson



Slide inspiration: Justin Johnson



Slide inspiration: Justin Johnson



Slide inspiration: Justin Johnson

Today's agenda

- A brief history of computer vision and deep learning
- CS231n overview

Evolution's Big Bang: Cambrian Explosion, 530-540million years, B.C.



[This image](#) is licensed under [CC-BY 2.5](#)



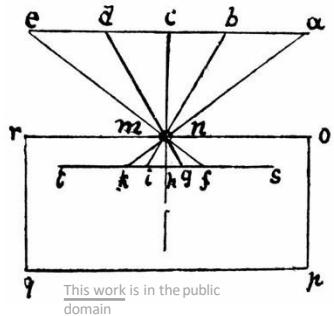
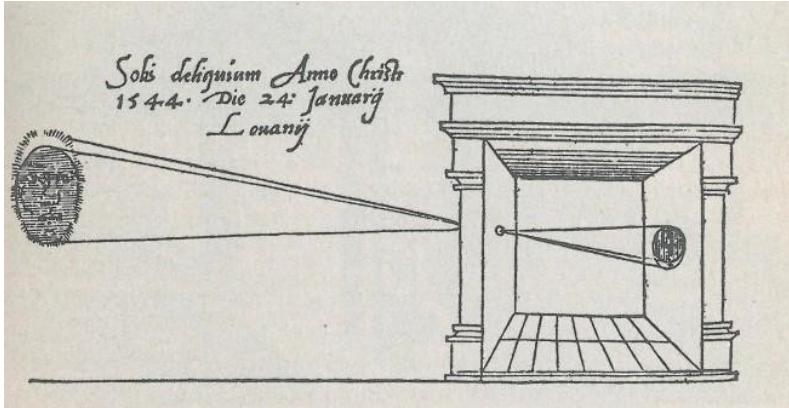
[This image](#) is licensed under [CC-BY 2.5](#)



[This image](#) is licensed under [CC-BY 3.0](#)

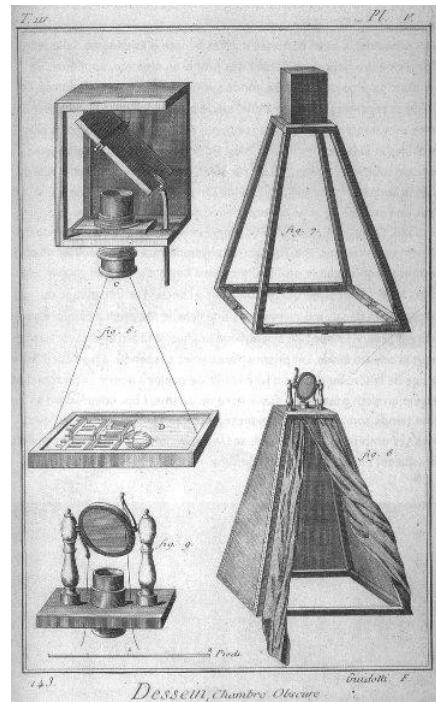
Camera Obscura

Gemma Frisius, 1545



Leonardo da Vinci,
16th Century AD

Encyclopedia, 18th Century



Computer Vision is everywhere!



Left to right:
[Image by Roger H Goun](#) is licensed under [CC BY 2.0](#)
[Image](#) is [CC0 1.0](#) public domain
[Image](#) is [CC0 1.0](#) public domain
[Image](#) is [CC0 1.0](#) public domain



Left to right:
[Image](#) is [free to use](#)
[Image](#) is [CC0 1.0](#) public domain
[Image](#) by [NASA](#) is licensed under [CC BY 2.0](#)
[Image](#) is [CC0 1.0](#) public domain

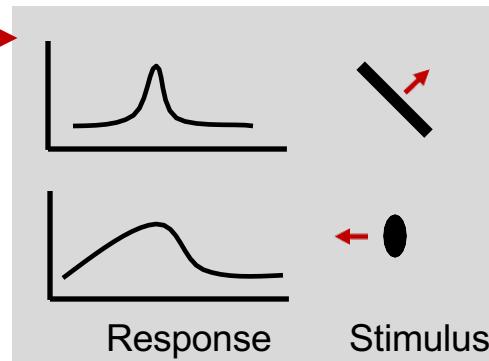
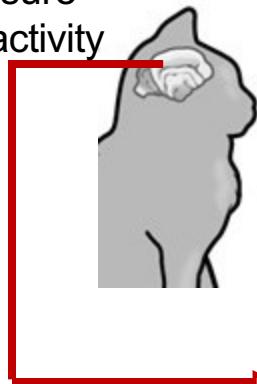


Bottom row, left to right
[Image](#) is [CC0 1.0](#) public domain
[Image](#) by [Derek Keats](#) is licensed under [CC BY 2.0](#)
changes made
[Image](#) is public domain
[Image](#) is licensed under [CC-BY 2.0](#); changes made

Where did we come from?

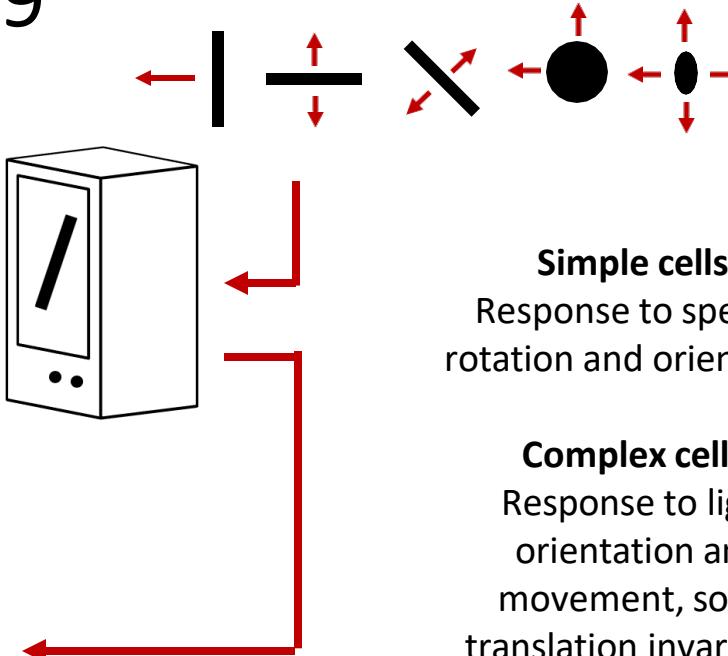
Hubel and Wiesel, 1959

Measure
brain activity



Cat image by [CNX OpenStax](#) is licensed under [CC BY 4.0](#) changes made

1959
Hubel & Wiesel

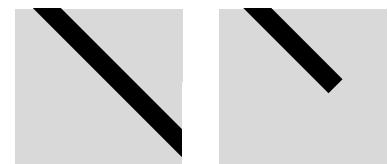


Simple cells:

Response to specific rotation and orientation

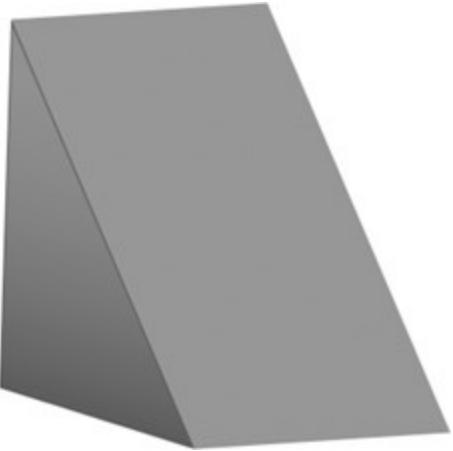
Complex cells:

Response to light orientation and movement, some translation invariance

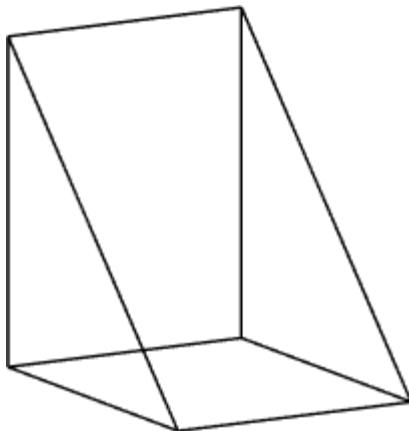


Slide inspiration: Justin Johnson

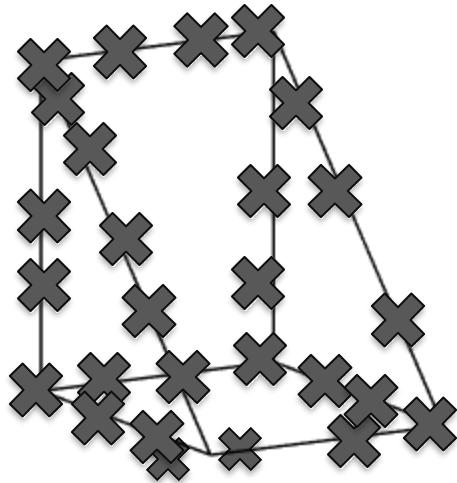
Larry Roberts, 1963



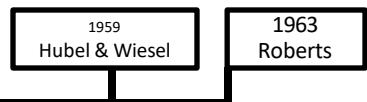
(a) Original picture



(b) Differentiated picture



(c) Feature points selected

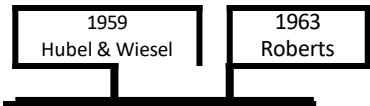


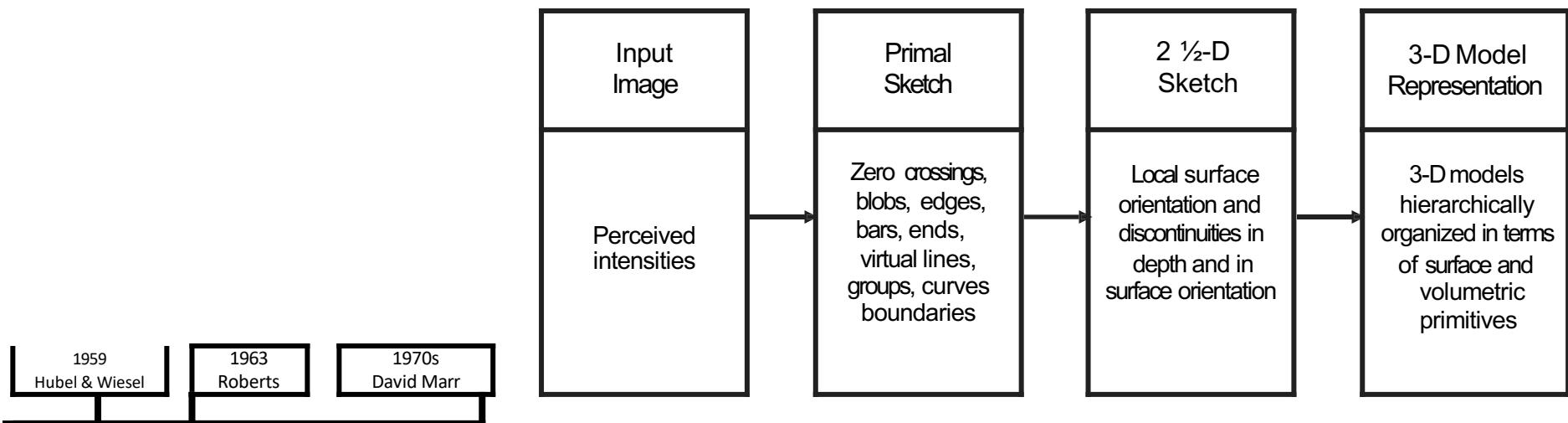
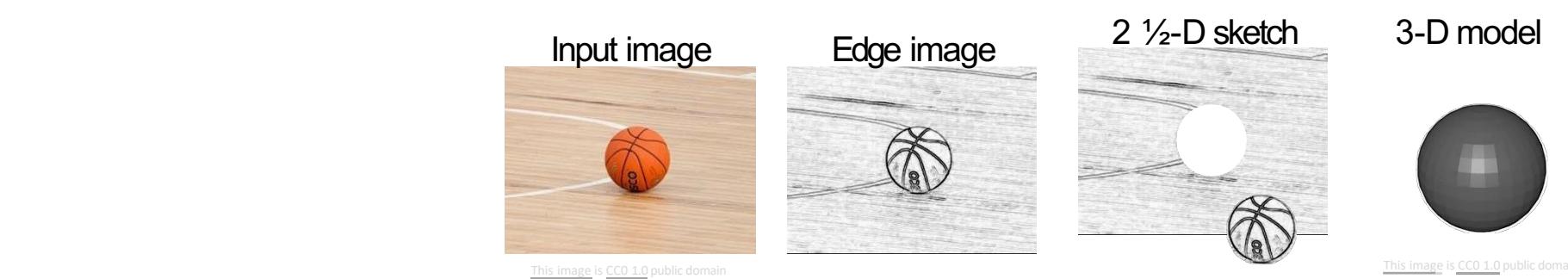
Lawrence Gilman Roberts, "Machine Perception of Three-Dimensional Solids", 1963

Slide inspiration: Justin Johnson

The Summer Vision Project- Seymour Paper

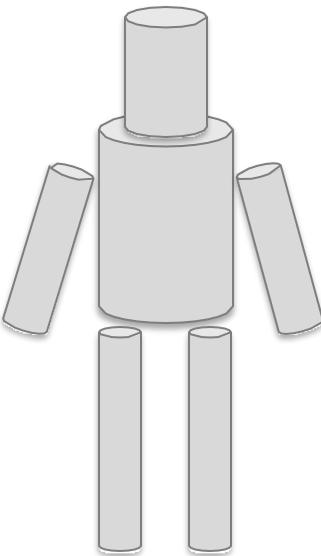
<https://dspace.mit.edu/handle/1721.1/6125>



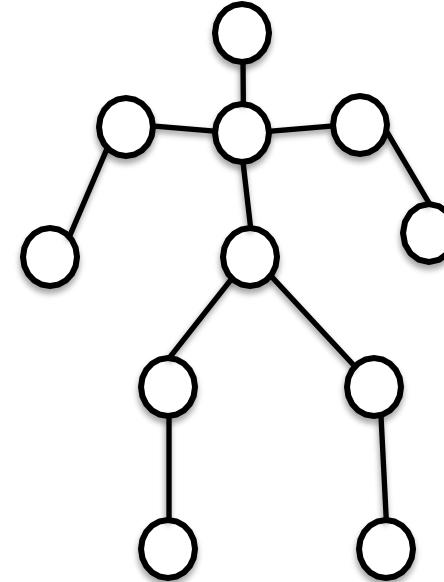


Stages of Visual Representation, David Marr, 1970s

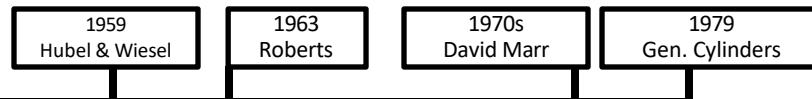
Recognition via Parts (1970s)



Generalized Cylinders,
Brooks and Binford,
1979



Pictorial Structures,
Fischler and Elshlager, 1973



Recognition via Edge Detection (1980s)



1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

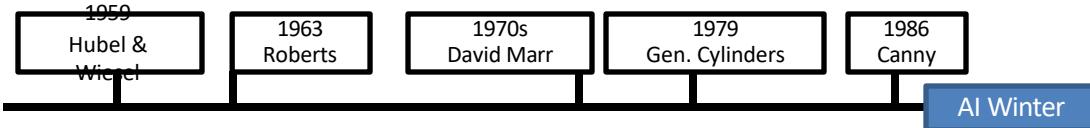
1979
Gen. Cylinders

1986
Canny

John Canny, 1986
David Lowe, 1987

Arriving at an “AI winter”

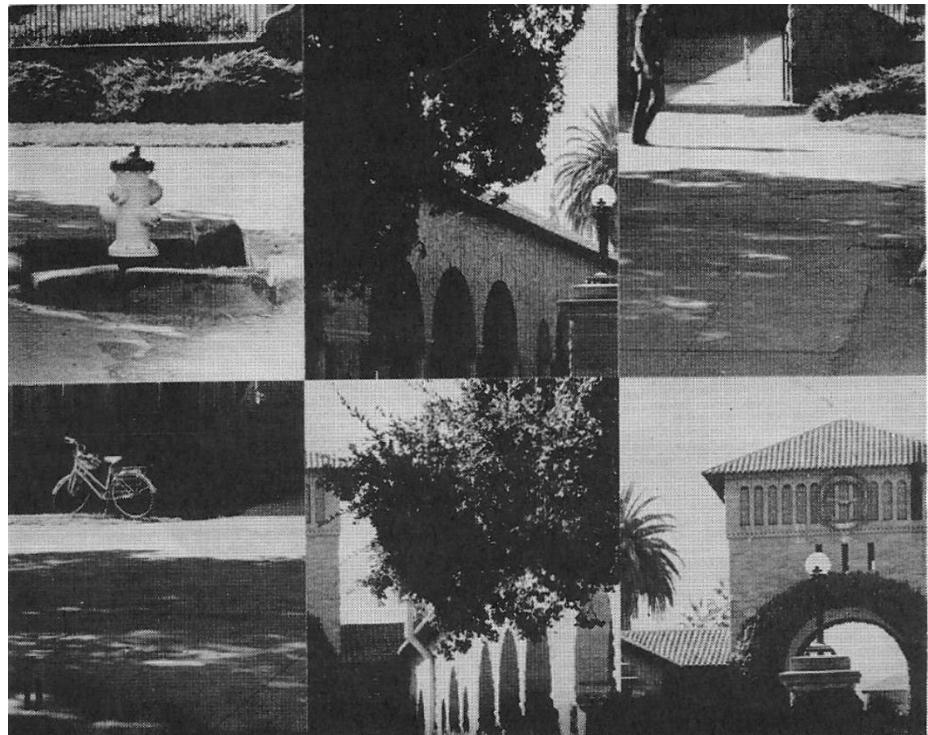
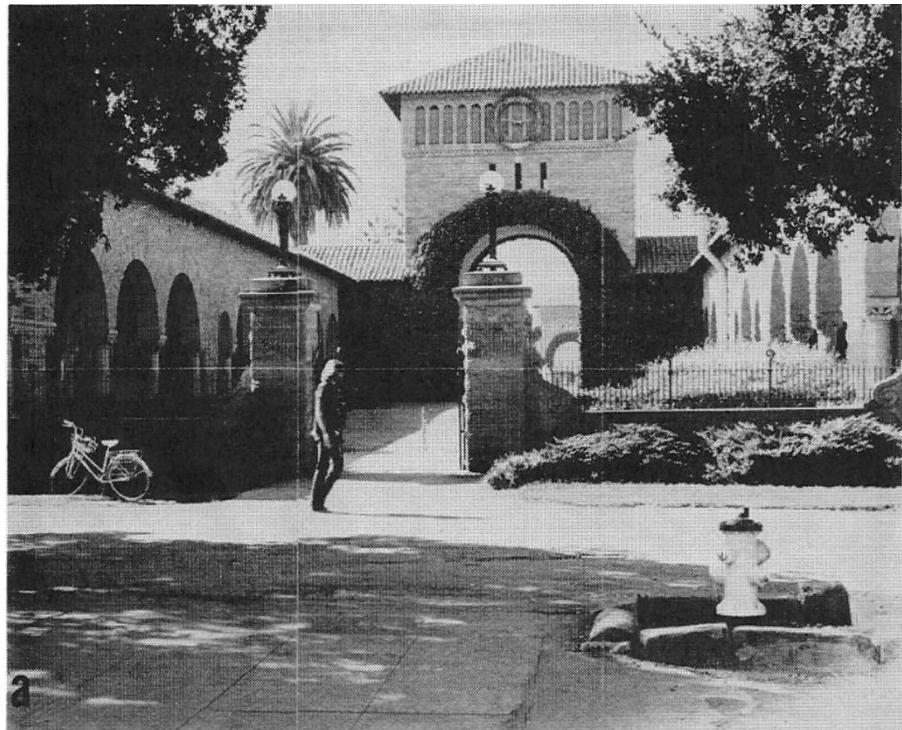
- Enthusiasm (and funding!) for AI research dwindled
- “Expert Systems” failed to deliver on their promises
- But subfields of AI continues to grow
 - Computer vision, NLP, robotics, compbio, etc.



In the meantime...seminal work in
cognitive and neuroscience

Perceiving Real-World Scenes

Irving Biederman



AAAS (aaas.org) Material: Copyright permission: Figure 1a,b from Irving Biederman,
Perceiving Real-World Scenes. Science 177,77- 80(1972). DOI:10.1126/science.177.4043.77

I. Biederman, *Science*, 1972

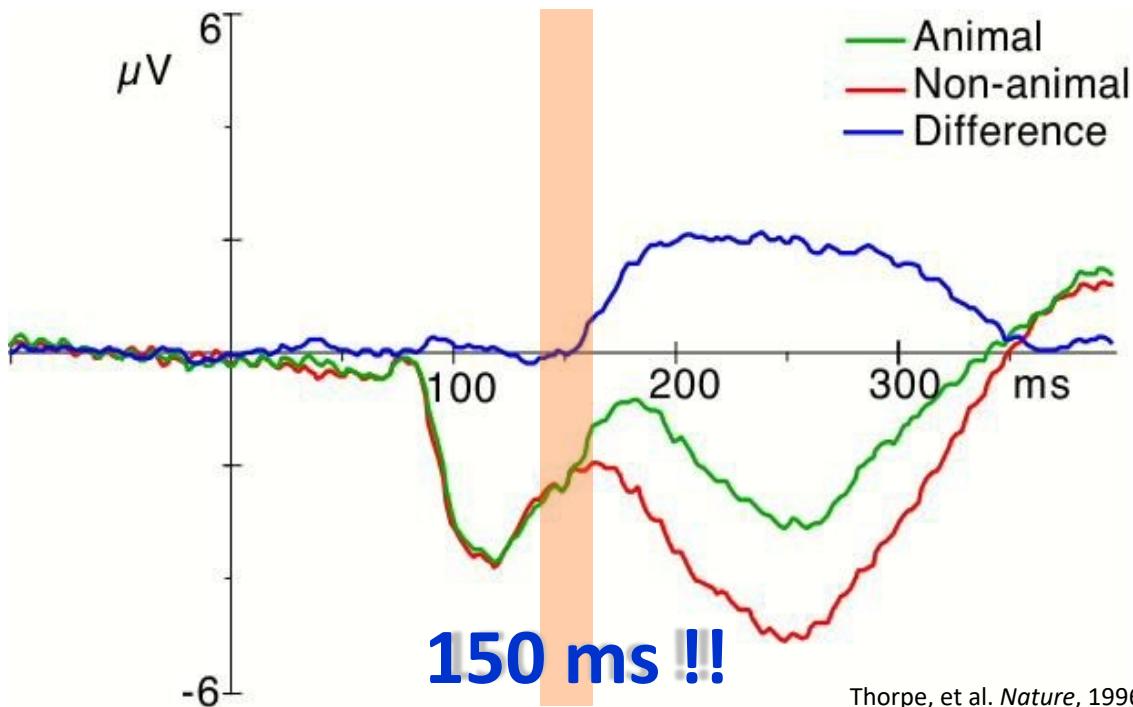
Rapid Serial Visual Perception (RSVP)



Potter, etc. 1970s

Speed of processing in the human visual system

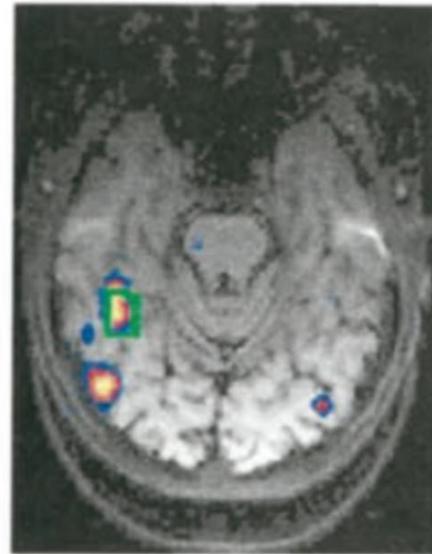
Simon Thorpe, Denis Fize & Catherine Marlot



Thorpe, et al. *Nature*, 1996

Neural correlates of object & scene recognition

Faces > Houses

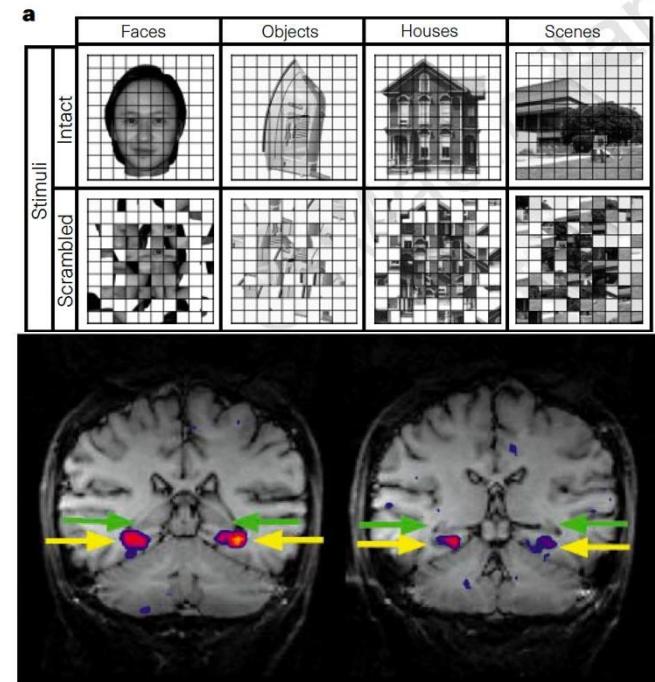


% signal change

Kanwisher et al. J. Neuro. 1997

Top image on left source: CC-BY-NC-SA 4.0

Bottom image on left source: Wikimedia public domain



Epstein & Kanwisher, Nature, 1998

Visual recognition is a fundamental task
for visual intelligence

Recognition via Grouping (1990s)



1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

1979
Gen. Cylinders

1986
Canny

1997
Norm. Cuts

AI Winter

Normalized Cuts, Shi and Malik, 1997

Recognition via Matching (2000s)



[Image](#) is public domain



[Image](#) is public domain

1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

1979
Gen. Cylinders

1986
Canny

1997
Norm. Cuts

1999
SIFT

AI Winter

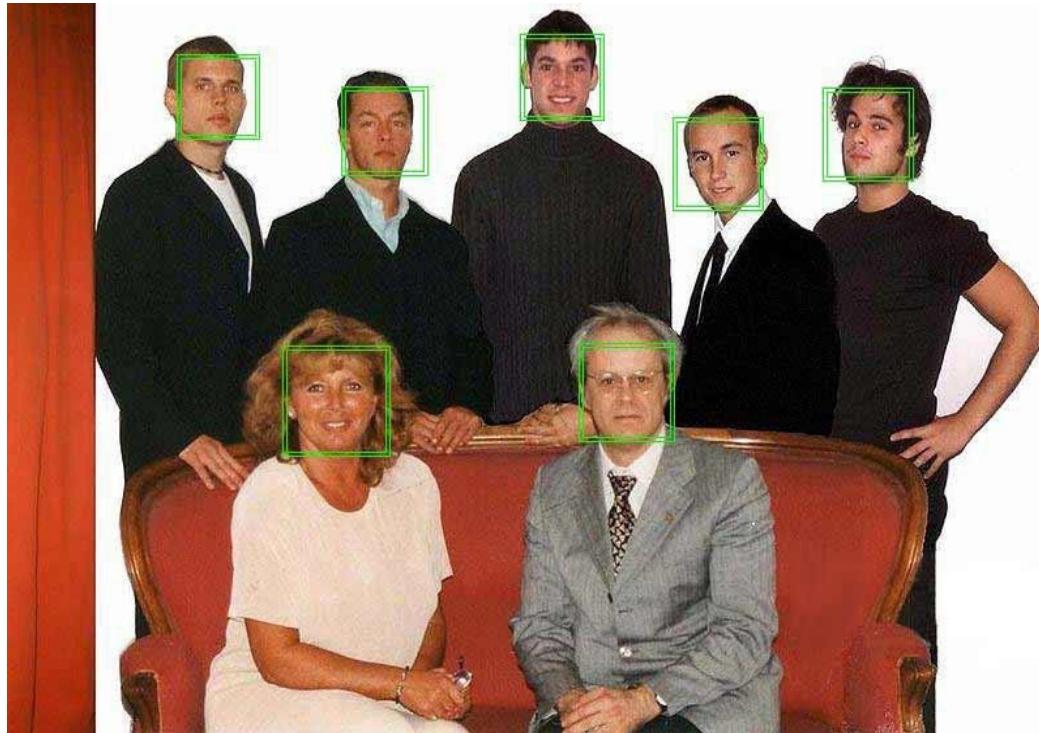
SIFT, David
Lowe, 1999

Slide inspiration: Justin Johnson

Face Detection

Viola and Jones, 2001

One of the first successful applications of machine learning to vision



1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

1979
Gen. Cylinders

1986
Canny

1997
Norm. Cuts

1999
SIFT

2001
V&J

AI Winter

Image source: Wikimedia public domain

Caltech 101 images



1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

1979
Gen. Cylinders

1986
Canny

1997
Norm. Cuts

1999
SIFT

2001
V&J

2004, 2007
Caltech101;
PASCAL

AI Winter

PASCAL Visual Object Challenge

[Image](#) is CC0 1.0 public domain



[Image](#) is CC0 1.0 public domain

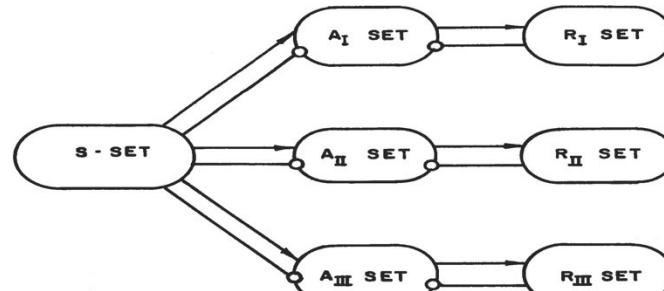
Perceptron

Learning representations by back-propagating errors

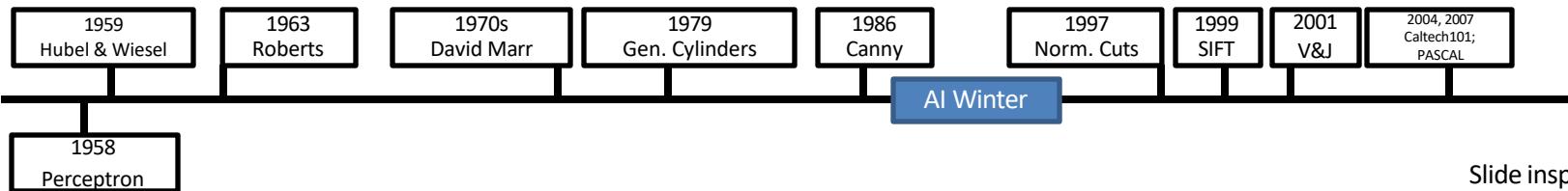
David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California,
San Diego, La Jolla, California 92093, USA

† Department of Computer Science, Carnegie-Mellon University,
Pittsburgh, Philadelphia 15213, USA



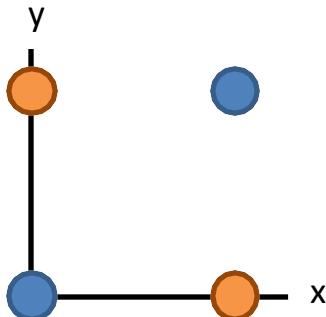
Frank Rosenblatt, ~1957



Slide inspiration: Justin Johnson

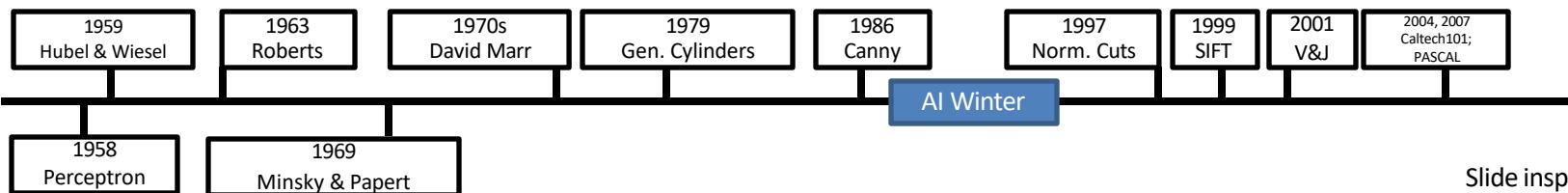
Minsky and Papert, 1969

X	Y	F(x,y)
0	0	0
0	1	1
1	0	1
1	1	0



Showed that Perceptrons could not learn the XOR function

Caused a lot of disillusionment in the field



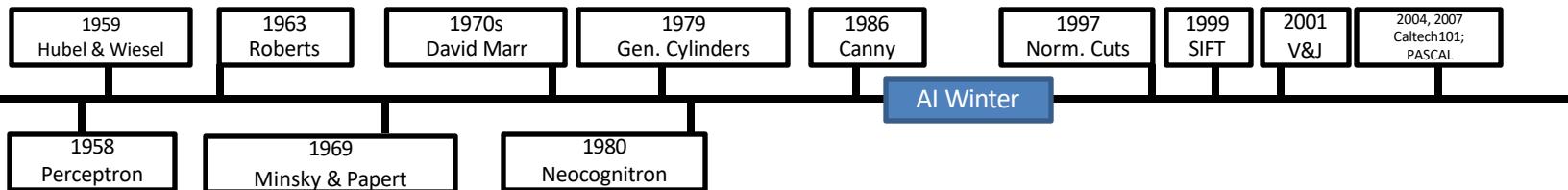
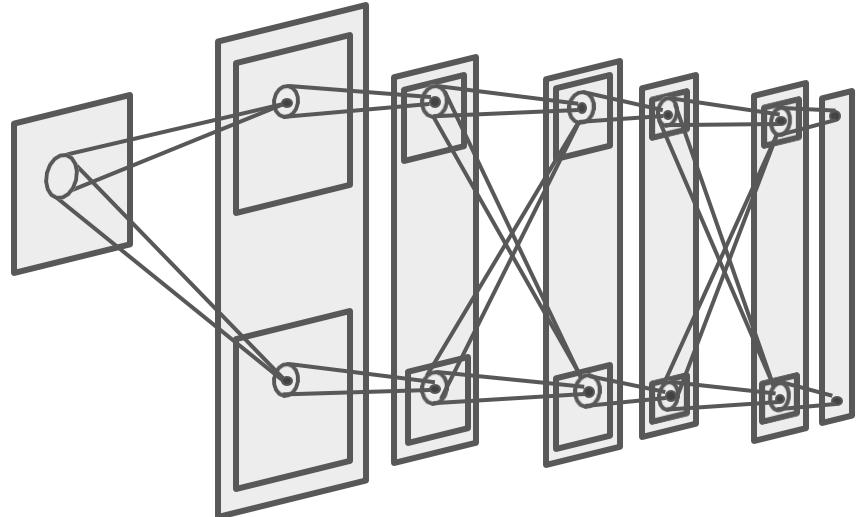
Slide inspiration: Justin Johnson

Neocognitron: Fukushima, 1980

Computational model the visual system,
directly inspired by Hubel and Wiesel's
hierarchy of complex and simple cells

Interleaved simple cells (convolution)
and complex cells (pooling)

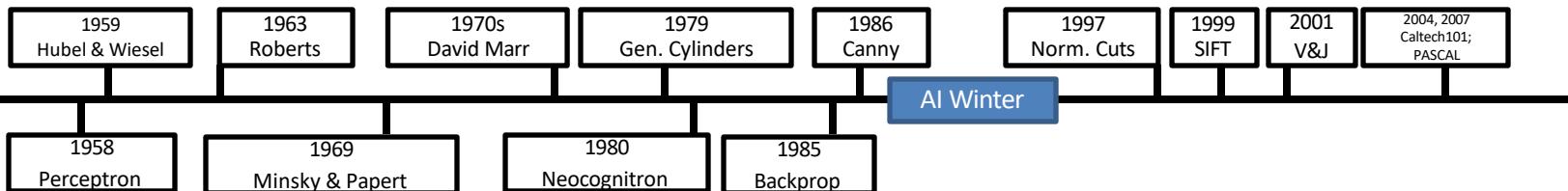
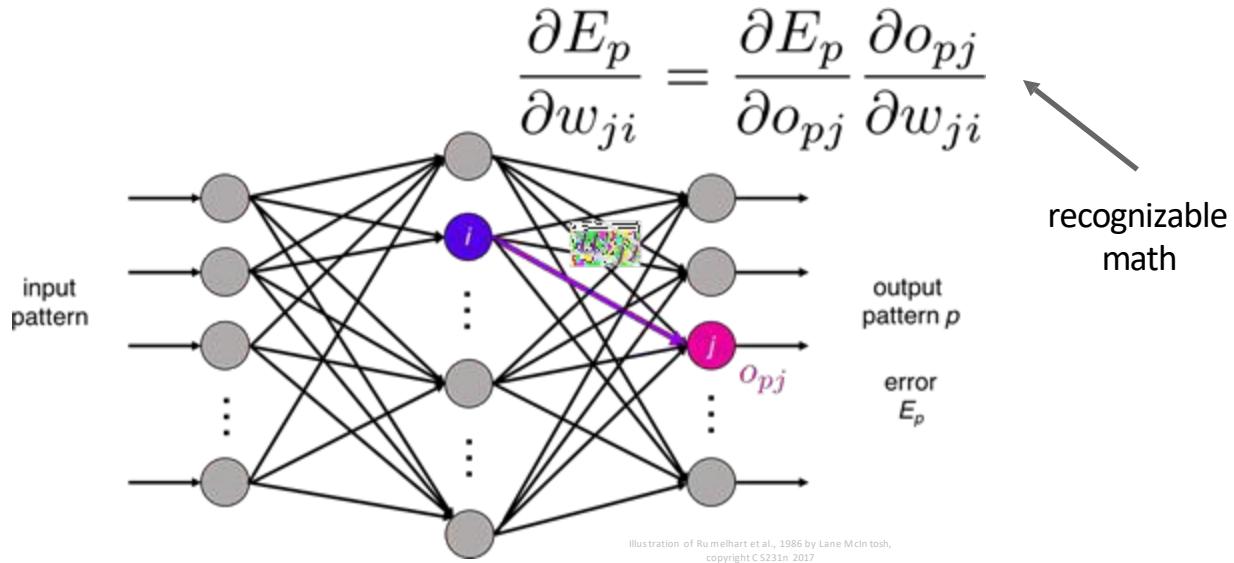
No practical training algorithm



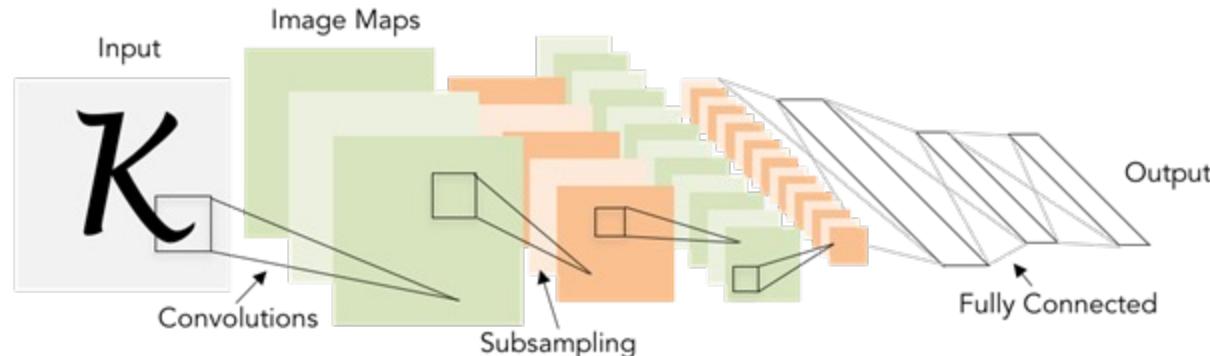
Backprop: Rumelhart, Hinton, and Williams, 1986

Introduced backpropagation for computing gradients in neural networks

Successfully trained perceptrons with multiple layers



Convolutional Networks: LeCun et al, 1998

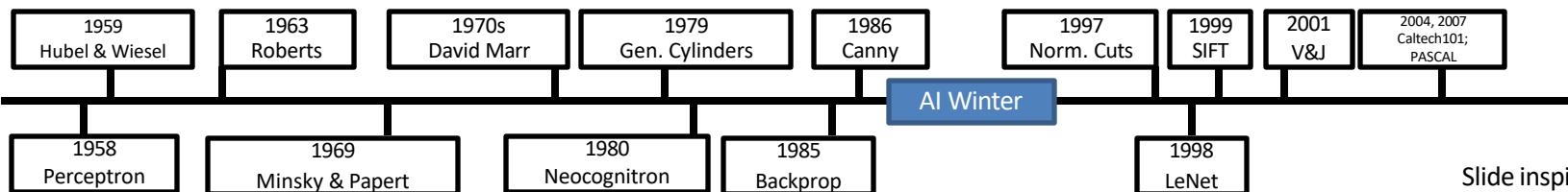


Applied backprop algorithm to a Neocognitron-like architecture

Learned to recognize handwritten digits

Was deployed in a commercial system by NEC, processed handwritten checks

Very similar to our modern convolutional networks!



Slide inspiration: Justin Johnson

2000s: “Deep Learning”

People tried to train neural networks that were deeper and deeper

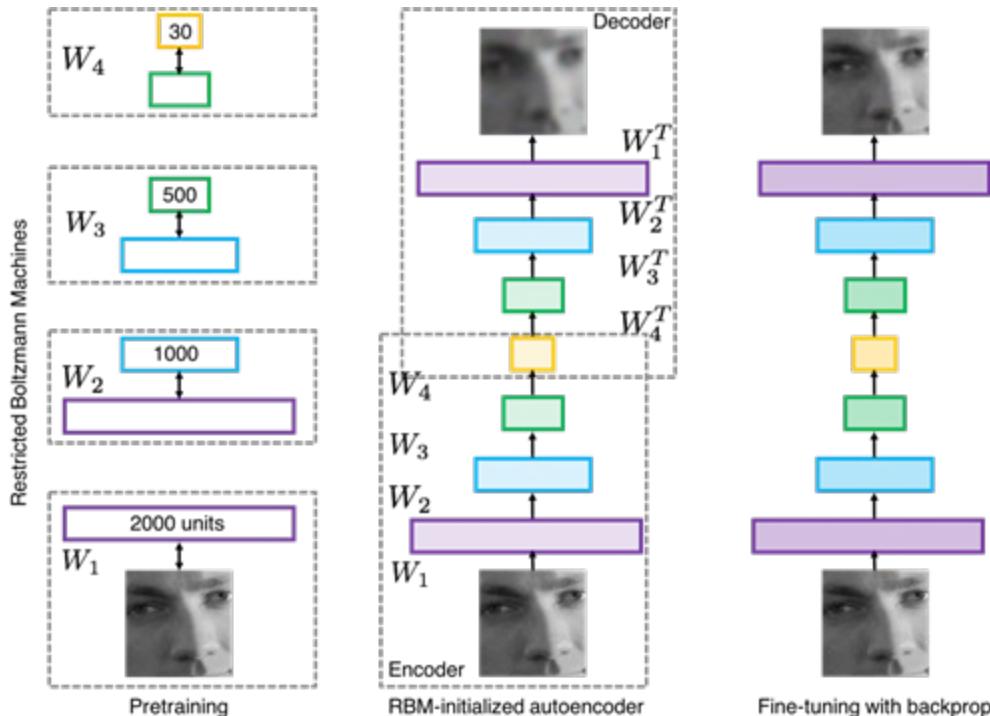
Not a mainstream research topic at this time

Hinton and Salakhutdinov, 2006

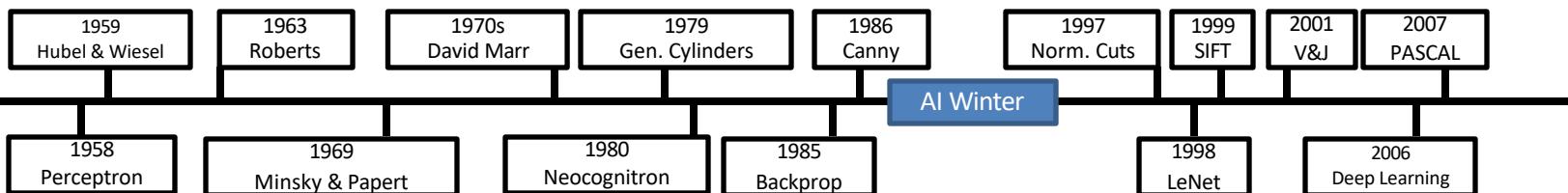
Bengio et al, 2007

Lee et al, 2009

Glorot and Bengio, 2010



Fine-tuning with backprop



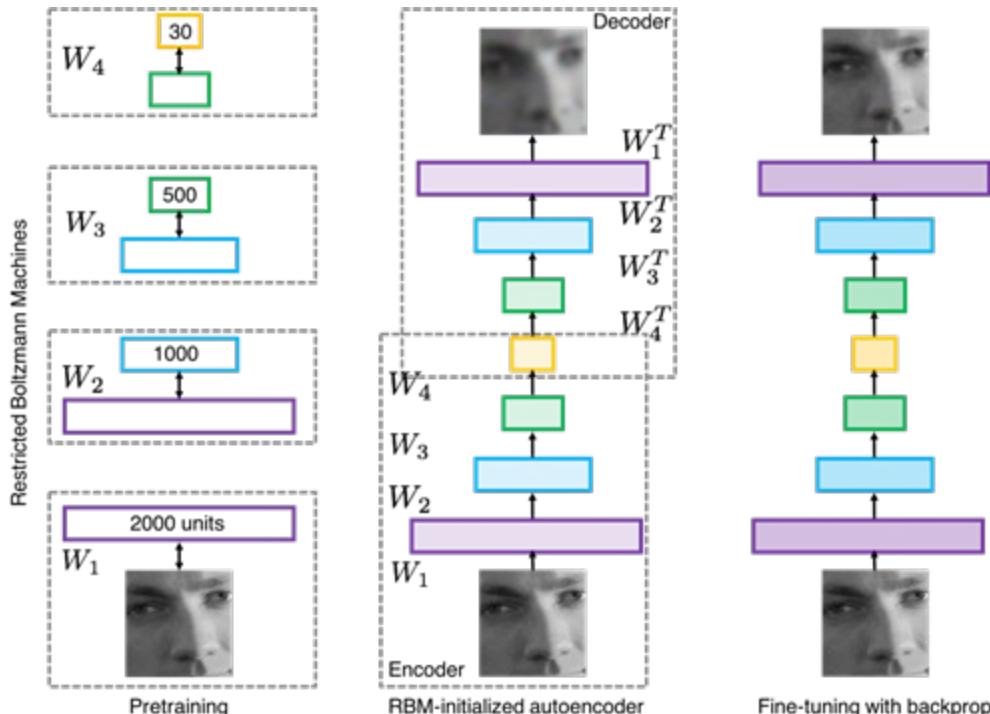
2000s: “Deep Learning”

People tried to train neural networks that were deeper and deeper

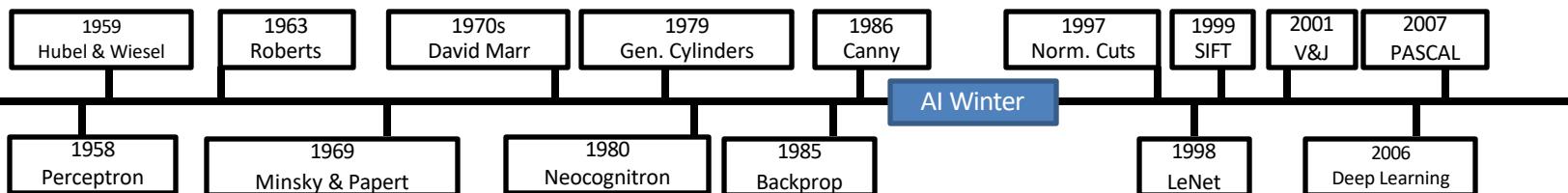
Not a mainstream research topic at this time

No good dataset to work on

Hinton and Salakhutdinov, 2006
Bengio et al, 2007
Lee et al, 2009
Glorot and Bengio, 2010



Fine-tuning with backprop



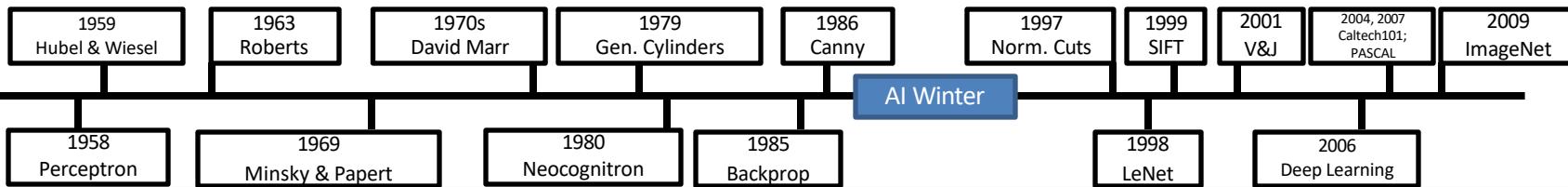
IMAGENET Large Scale Visual Recognition Challenge

The Image Classification Challenge:
1,000 object classes
1,431,167 images

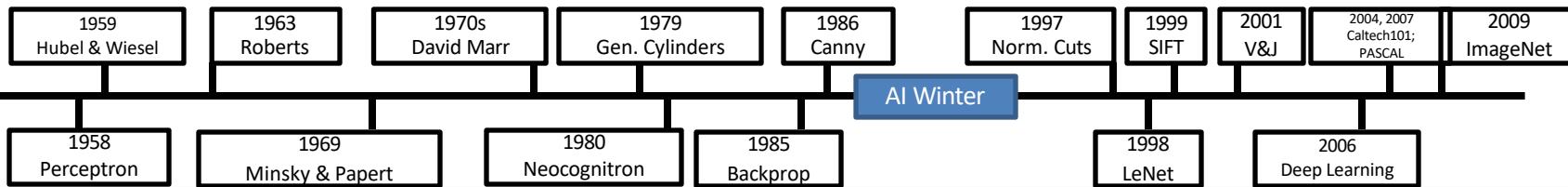
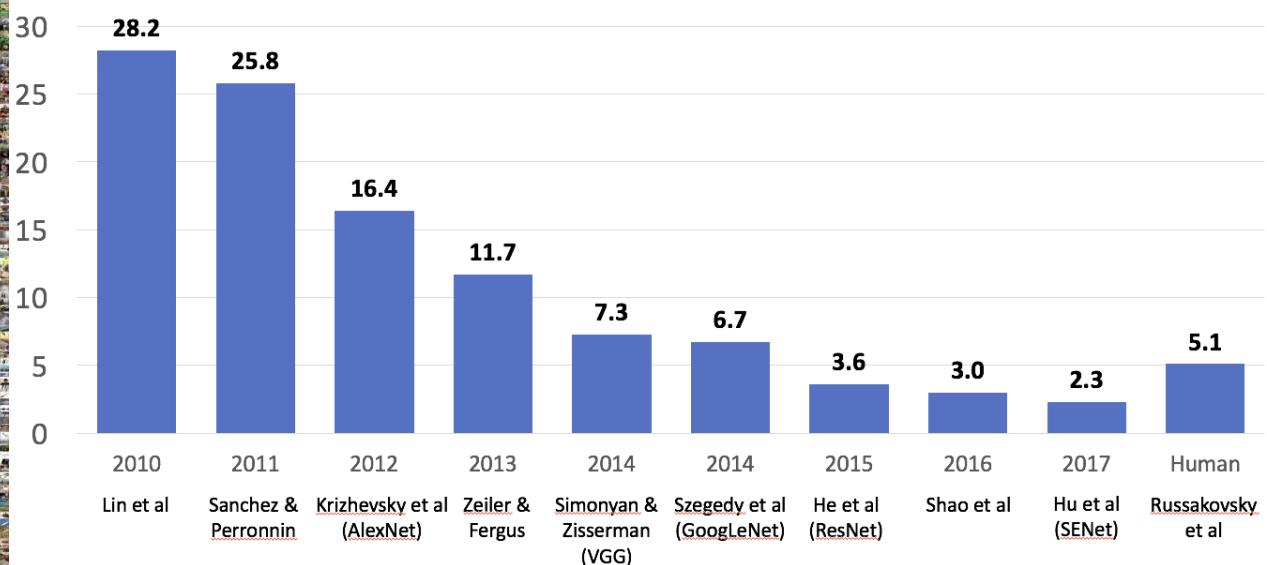


Output:
Scale
T-shirt
Steel drum
Drumstick
Mud turtle

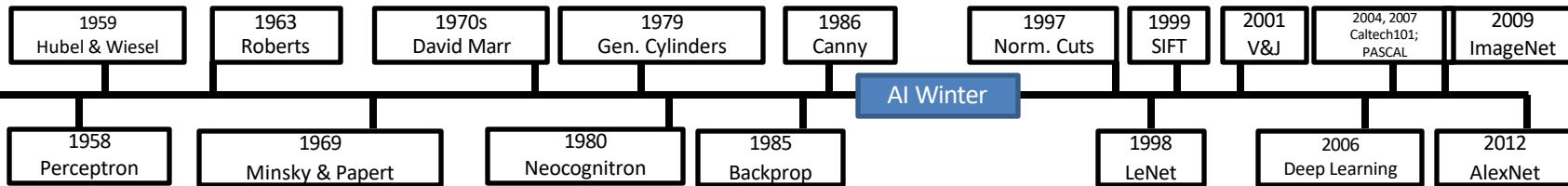
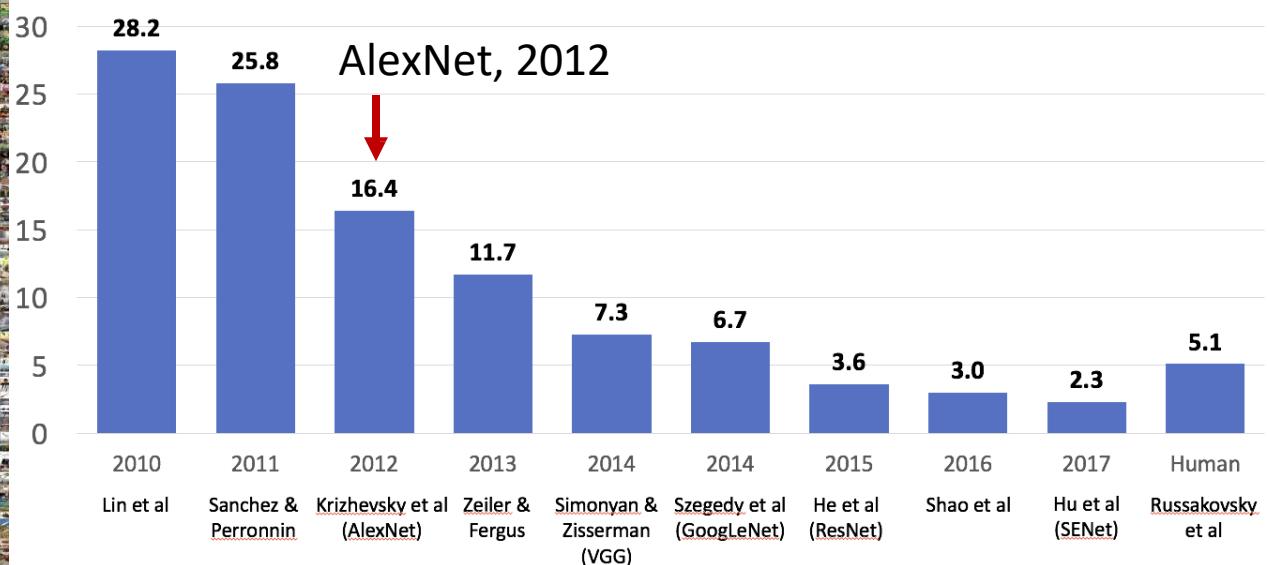
Deng et al, 2009
Russakovsky et al. IJCV 2015



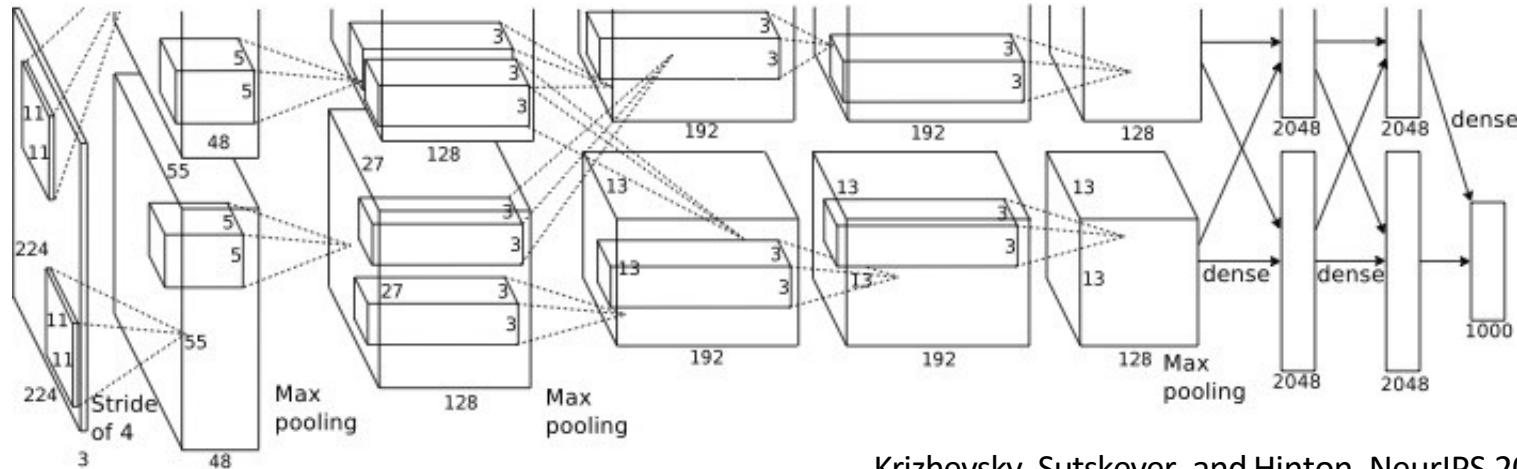
IMAGENET Large Scale Visual Recognition Challenge



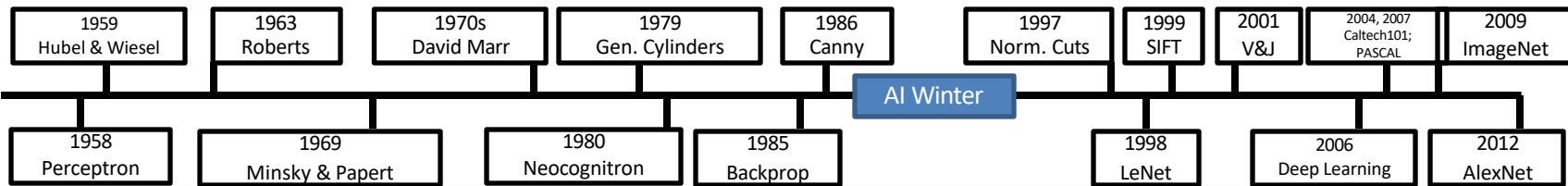
IMAGENET Large Scale Visual Recognition Challenge



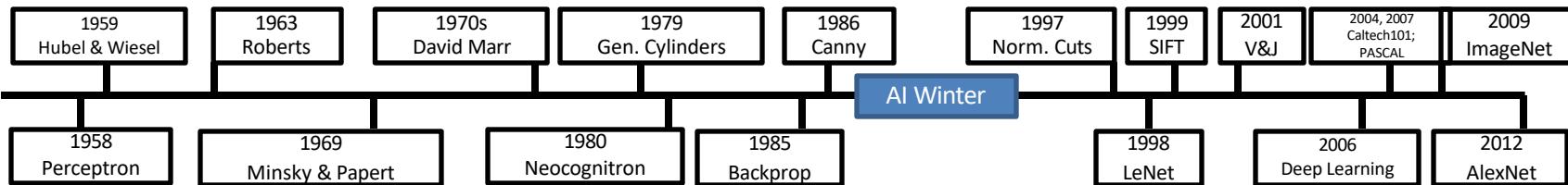
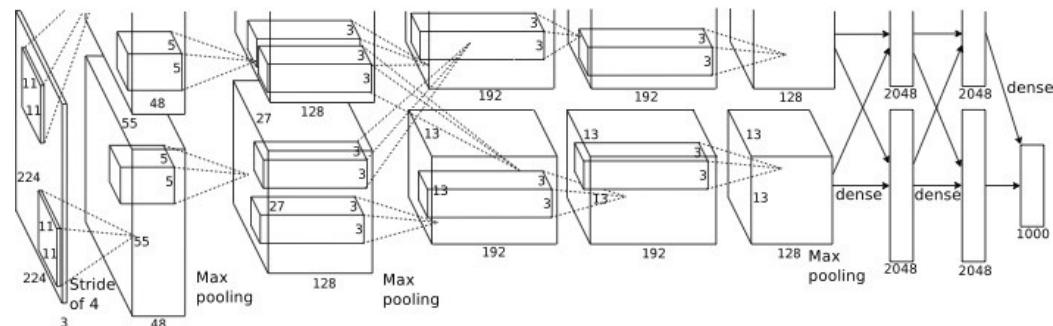
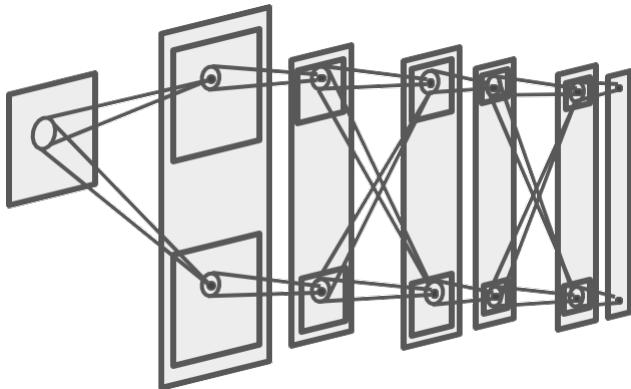
AlexNet: Deep Learning Goes Mainstream



Krizhevsky, Sutskever, and Hinton, NeurIPS 2012

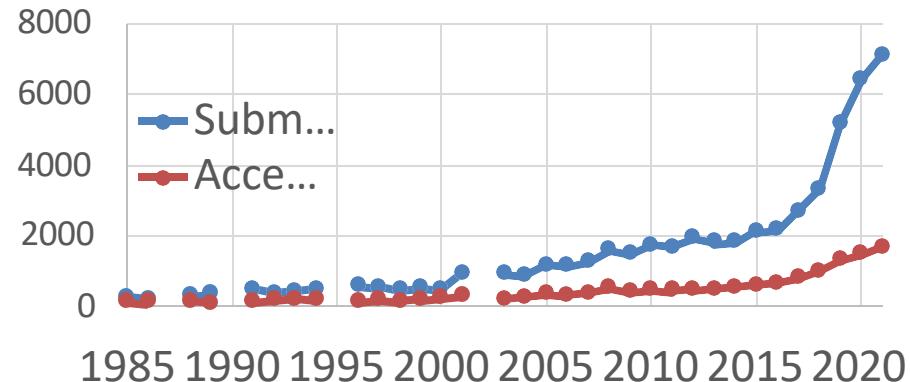


AlexNet vs. Neocognitron: 32 years apart

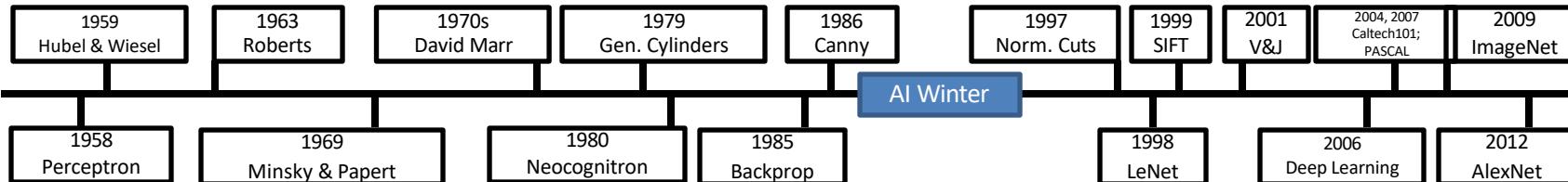


2012 to Present: Deep Learning Explosion

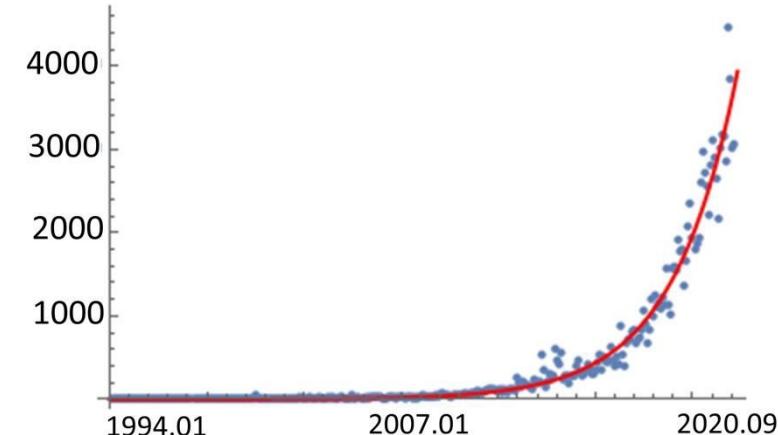
CVPR Papers



Publications at top Computer Vision conference



ML+AI arXiv papers per month

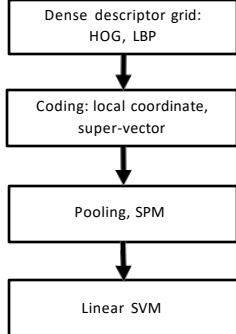


arXiv papers per month [\(source\)](#)

2012 to Present: Deep Learning is Everywhere

Year 2010

NEC-UUC

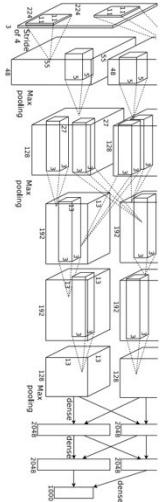


[Lin CVPR 2011]

Lion image by Swissfrog
is licensed under CC BY 3.0

Year 2012

SuperVision

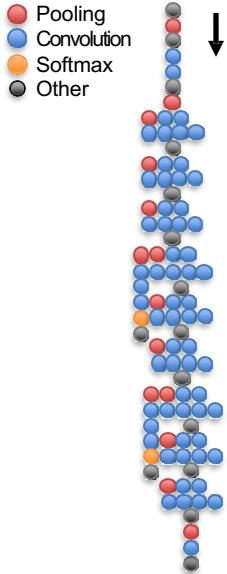


[Krizhevsky NIPS 2012]

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.
Reproduced with permission.

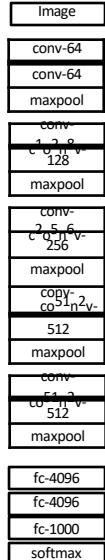
Year 2014

GoogLeNet



[Szegedy arxiv 2014]

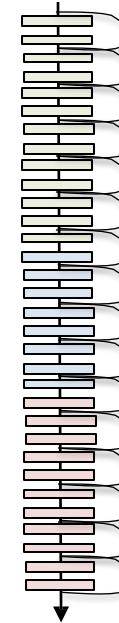
VGG



[Simonyan arxiv 2014]

Year 2015

MSRA



[He ICCV 2015]

2012 to Present: Deep Learning is Everywhere

Image Classification



Image Retrieval



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

2012 to Present: Deep Learning is Everywhere

Object Detection

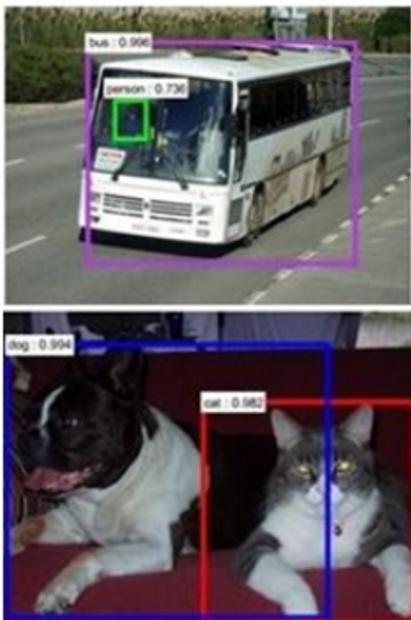


Image Segmentation

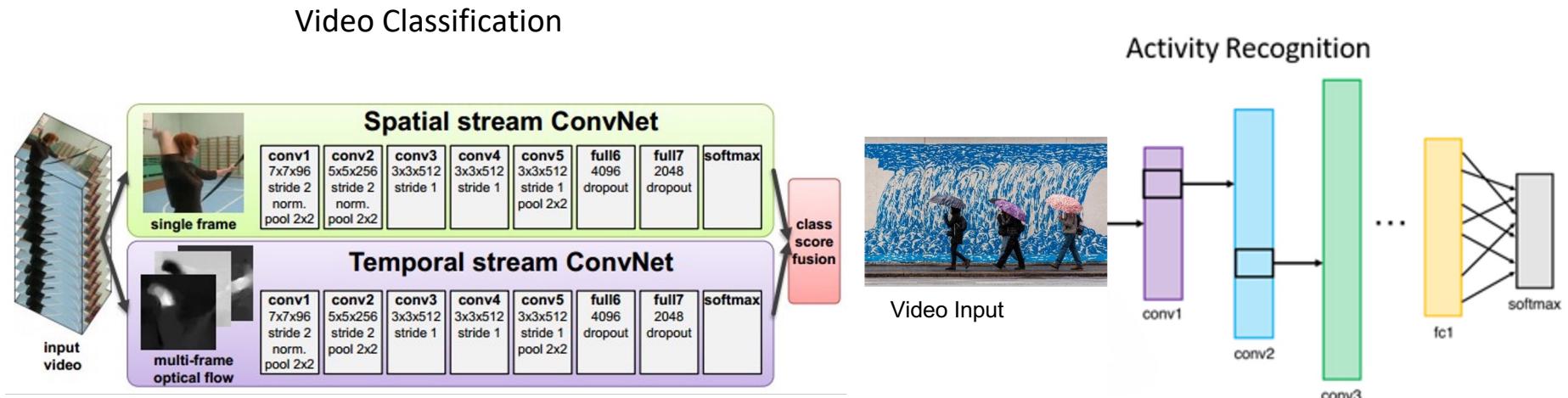


Figures copyright Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun, 2015. Reproduced with permission.
Images adapted from above source

Figures copyright Clement Farabet, 2012. Reproduced with permission.

[Farabet et al., 2012]

2012 to Present: Deep Learning is Everywhere



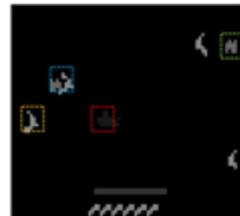
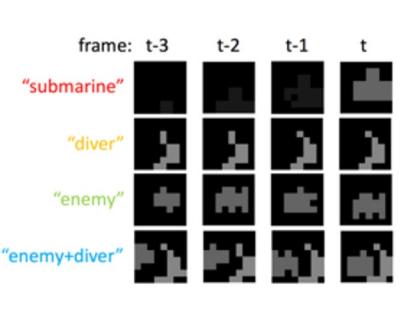
Simonyan et al,
2014

2012 to Present: Deep Learning is Everywhere

Pose Recognition (Toshev and Szegedy, 2014)



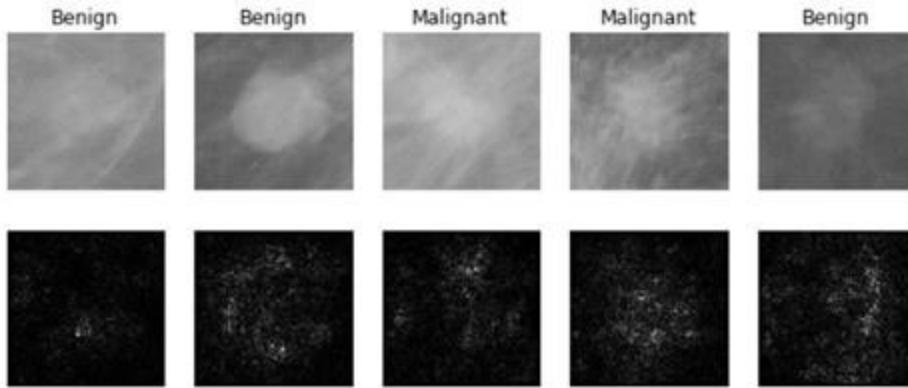
Playing Atari games (Guo et al, 2014)



Logo source: Wikimedia Commons Public Domain

2012 to Present: Deep Learning is Everywhere

Medical Imaging



Levy et al, 2016

Figure reproduced with permission

Whale recognition



Galaxy Classification



Dieleman et al, 2014

From left to right: public domain by NASA, usage permitted by
ESA/Hubble, public domain by NASA, and public domain

Kaggle Challenge

This image by Chirstin Khan is in the public domain and originally came from the U.S. NOAA.

2012 to Present: Deep Learning is Everywhere



*A white teddy bear
sitting in the grass*



A man in a baseball uniform throwing a ball



*A woman is holding
a cat in her hand*



*A man riding a wave
on top of a surfboard*



*A cat sitting on a
suitcase on the floor*



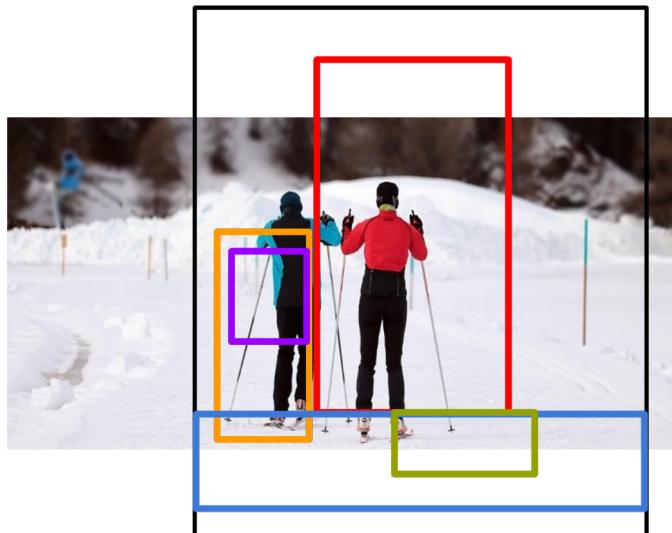
A woman standing on a beach holding a surfboard

Image Captioning

Vinyals et al, 2015
Karpathy and , 2015

All images are CC0 Public domain :
<https://pixabay.com/en/leisure-activity-recreation/>
<https://pixabay.com/en/teddy-dolph-bear-cute-teddy-hear-1623434/>
<https://pixabay.com/en/surf-wave-sun-supper-sun-surf-1668716/>
<https://pixabay.com/en/woman-female-model-nude-trait-adult-983967/>
<https://pixabay.com/en/hands-stand-land-meditation-496008/>
<https://pixabay.com/en/baseball-player-shorts-top-in-field-1045263/>

2012 to Present: Deep Learning is Everywhere



Results:

spatial, comparative, asymmetrical, verb,
prepositional

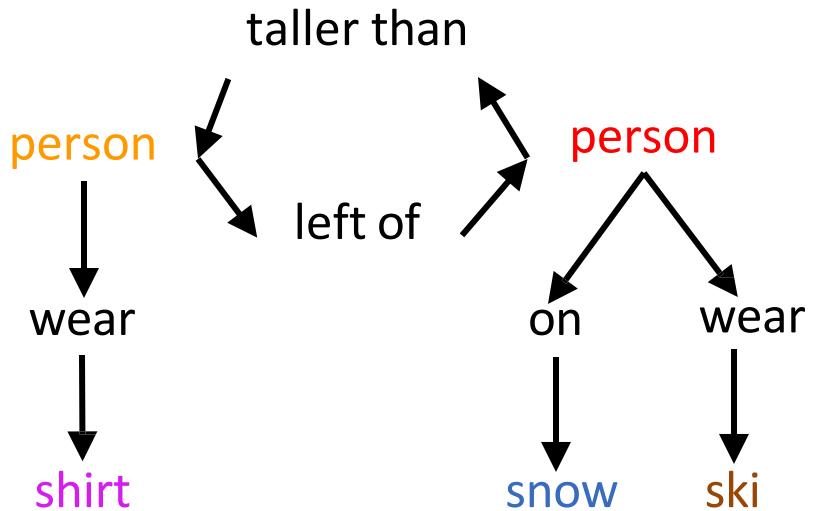
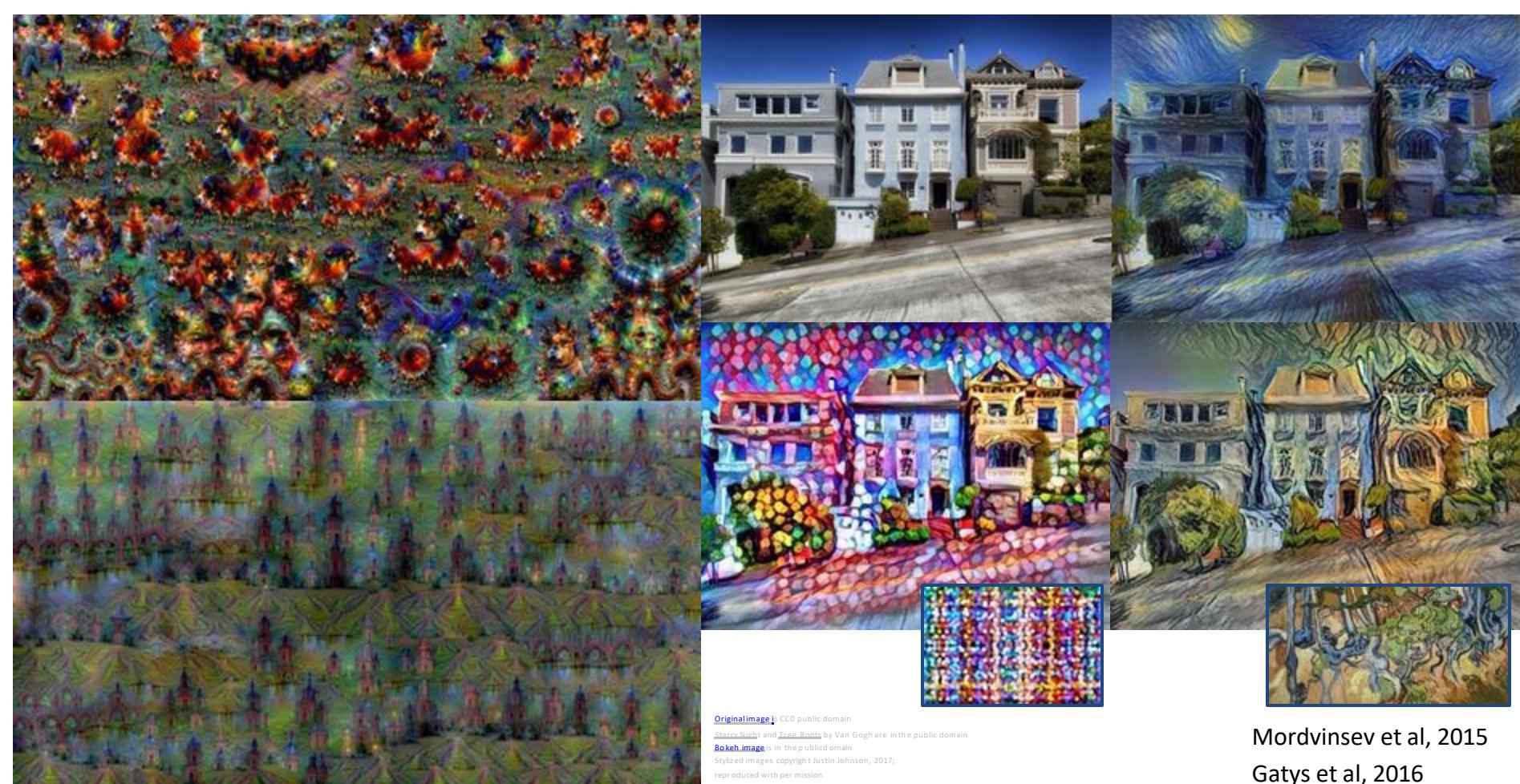


Image has been adapted

Krishna*, Lu*, Bernstein, , ECCV 2016



[Original image](#) is CC0 public domain

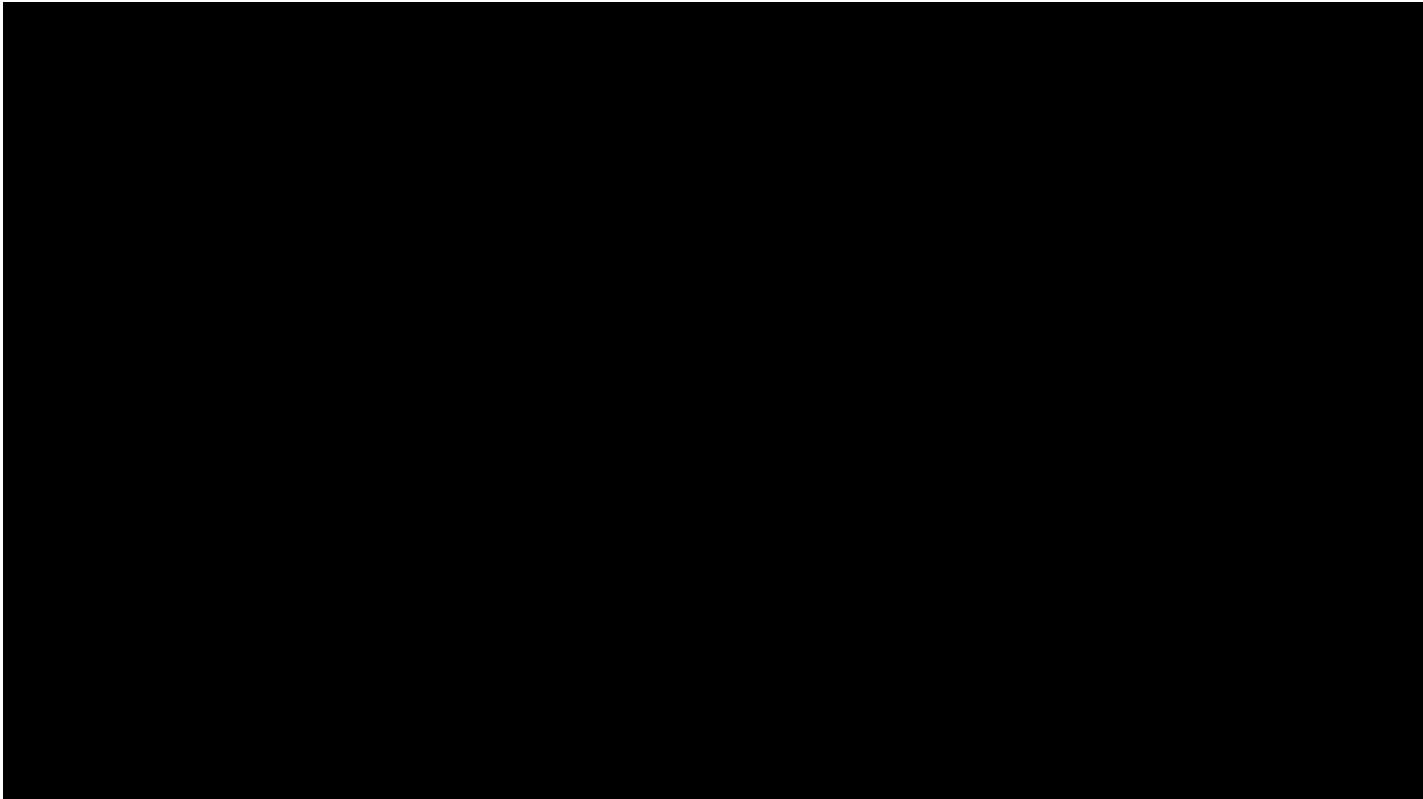
[Starry Night](#) and [Tree Roots](#) by Van Gogh are in the public domain

[Bokeh image](#) is in the public domain

Stylized images copyright Justin Johnson, 2017;
reproduced with permission

Mordvinsev et al, 2015
Gatys et al, 2016

2012 to Present: Deep Learning is Everywhere



Karras et al, "Progressive Growing of GANs for Improved Quality, Stability, and Variation", ICLR 2018

2012 to Present: Deep Learning is Everywhere

TEXT PROMPT

an armchair in the shape of an avocado. an armchair imitating an avocado.

AI-GENERATED IMAGES



2012 to Present: Deep Learning is Everywhere

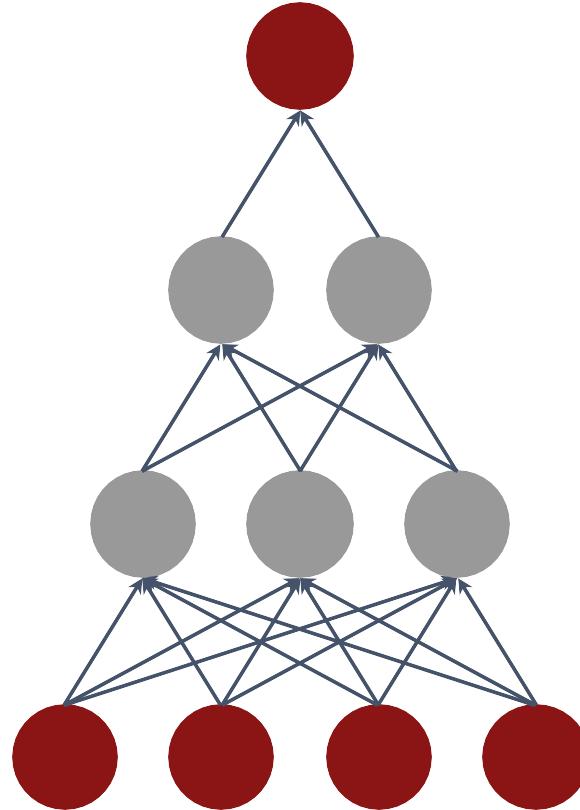
TEXT PROMPT

an armchair in the shape of a peach. an armchair imitating a peach.

AI-GENERATED IMAGES



Computation



Data

Algorithms

GFLOP per Dollar

● CPU ● GPU (FP32)

RTX 3080 →

RTX 3090 →

Deep Learning Explosion

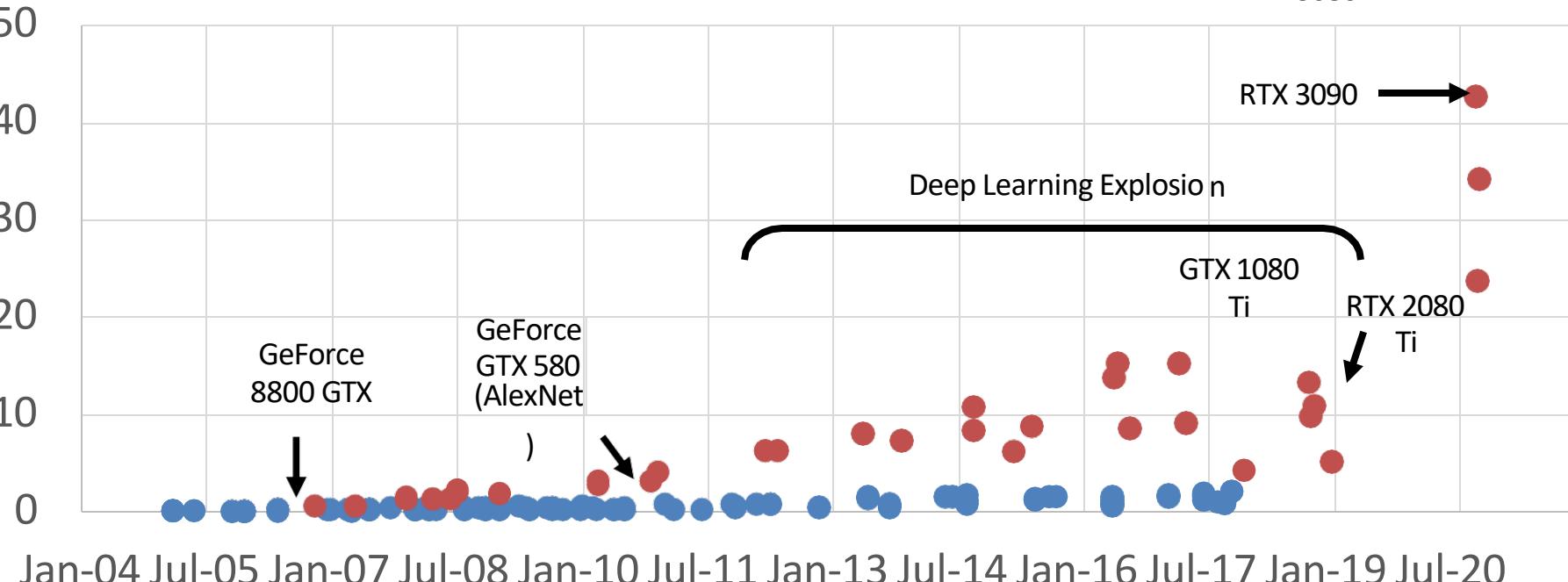
GTX 1080

Ti

RTX 2080
Ti

GeForce
8800 GTX

GeForce
GTX 580
(AlexNet)

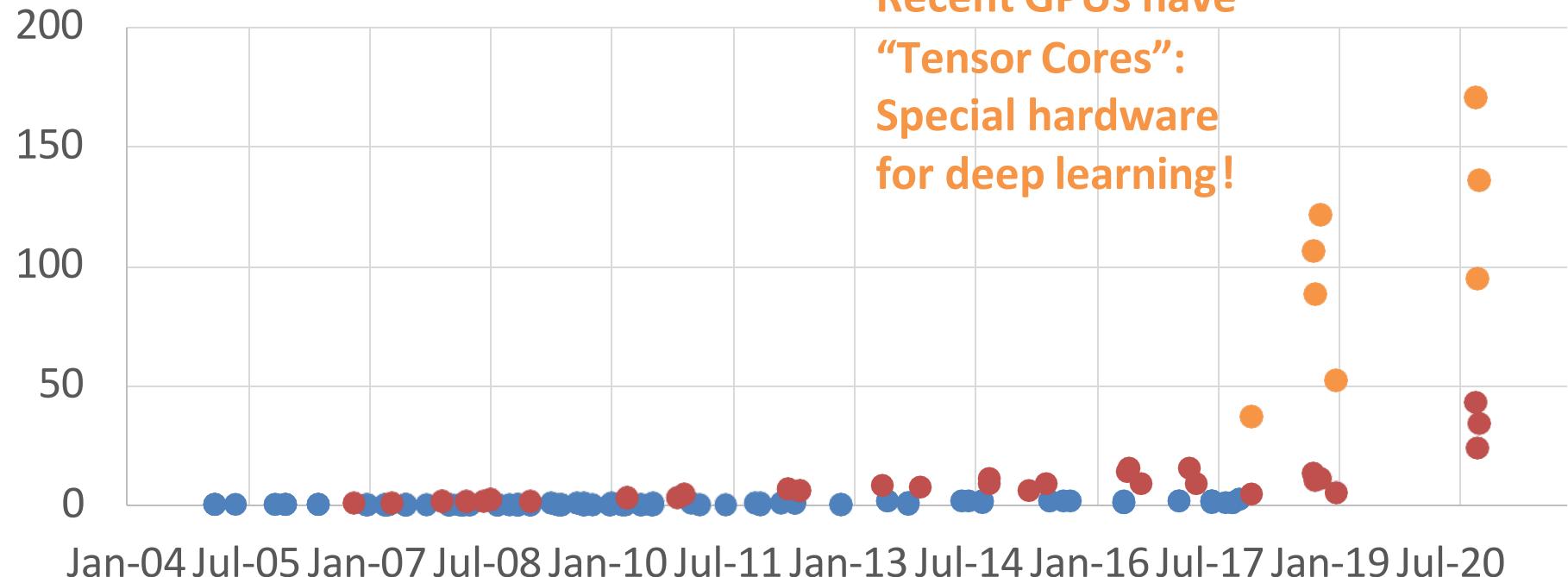


GFLOP per Dollar

● CPU ● GPU (FP32) ● GPU (Tensor Core)

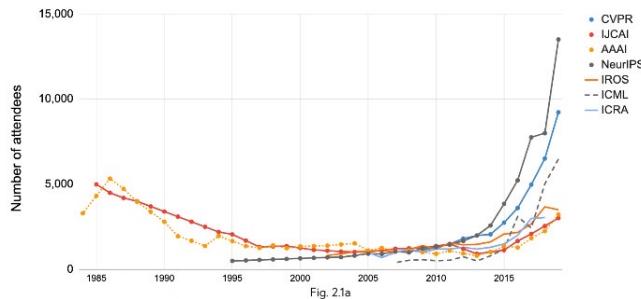
Recent GPUs have

“Tensor Cores”:
Special hardware
for deep learning!



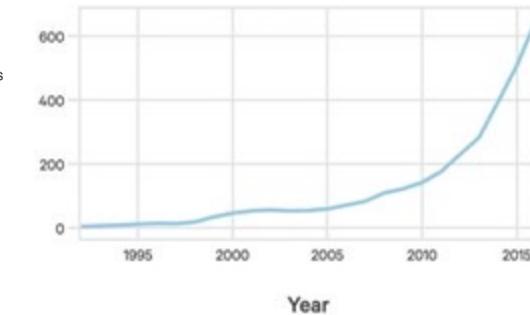
AI's Explosive Growth & Impact

Attendance at large conferences (1984-2019)
Source: Conference provided data.



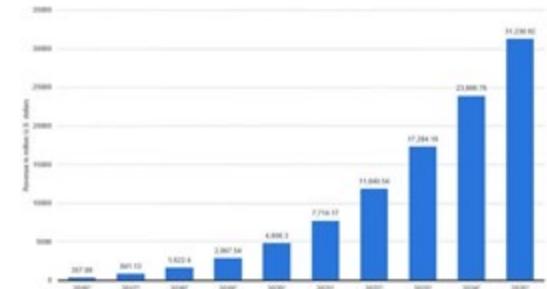
Number of attendance
At AI conferences

Source: The Gradient



Startups Developing AI
Systems

Source: Crunchbase, VentureSource, Sand
Hill Econometrics



Enterprise Application AI
Revenue

Source: Statista

Despite the successes, computer
vision still has a long way to go

Computer Vision Can Cause Harm

Harmful Stereotypes

Affect people's lives

https://www.youtube.com/watch?app=desktop&v=fMym_BKWQzk

The Trouble with Bias

Barocas et al, "The Problem With Bias: Allocative Versus
Representational Harms in Machine Learning", SIGCIS 2017
Kate Crawford, "The Trouble with Bias", NeurIPS 2017 Keynote

Technology

A face-scanning algorithm increasingly decides whether you deserve the job

HireVue claims it uses artificial intelligence to decide who's best for a job. Outside experts call it 'profoundly disturbing.'

Source: <https://www.washingtonpost.com/technology/2019/10/22/ai-hiring-face-scanning-algorithm-increasingly-decides-whether-you-deserve-job/>

<https://www.hirevue.com/platform/online-video-interviewing-software>

Example Credit: Timnit Gebru

Computer Vision Can Save Lives

How to take care of seniors while keeping them safe?



Early Symptom Detection
of COVID-19



Monitor Patients with
Mild Symptoms



Manage Chronic Conditions

Versatile



Mobility



Infection



Sleep



Diet

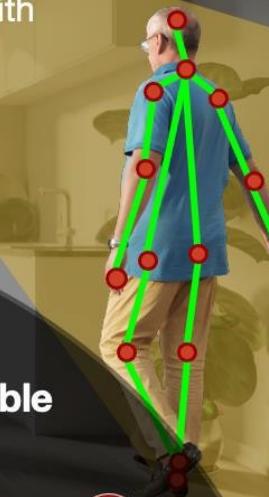
Scalable



Low-cost



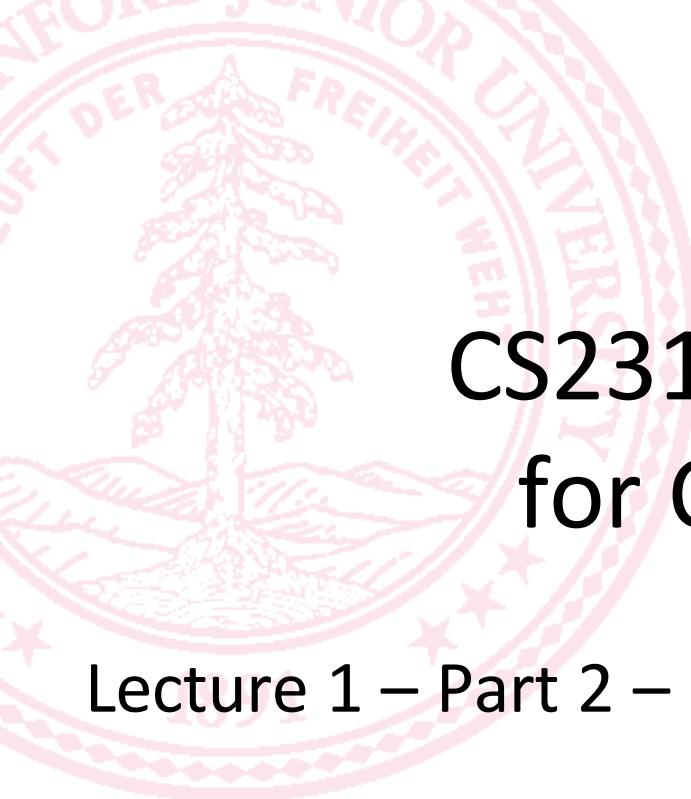
Burden-free



And there is a lot we don't know how to do



This image is
copyright-free United
States government
[work](#)



CS231n: Deep Learning for Computer Vision

Lecture 1 – Part 2 – Overview

CS231n overview

- Deep Learning Basics
- Perceiving and Understanding the Visual World
- Generative and Interactive Visual Intelligence
- Human-Centered Applications and Implications

Deep Learning Basics

- Image Classification: A core task in Computer Vision



→ cat

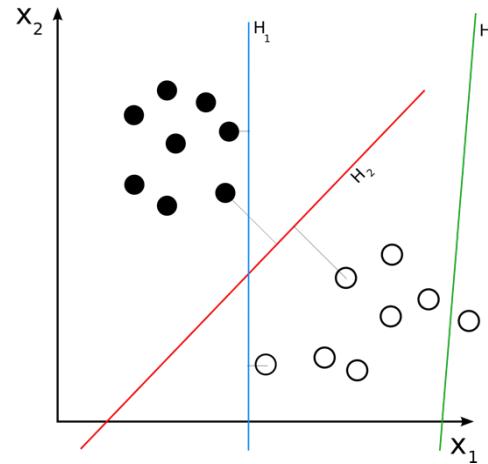
This image by [Nikita](#) is
licensed under [CC-BY 2.0](#)

Deep Learning Basics

- Image Classification: A core task in Computer Vision



cat



Linear Classifier

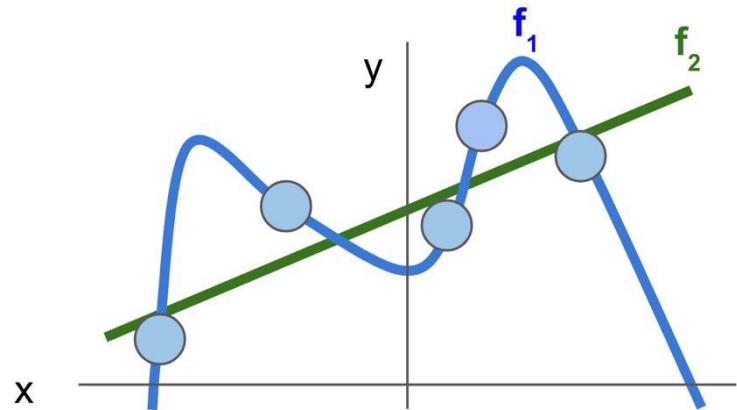
This image by [Nikita](#) is
licensed under [CC-BY 2.0](#)

Deep Learning Basics

- Image Classification: A core task in Computer Vision



This image by [Nikita](#) is
licensed under [CC-BY 2.0](#)



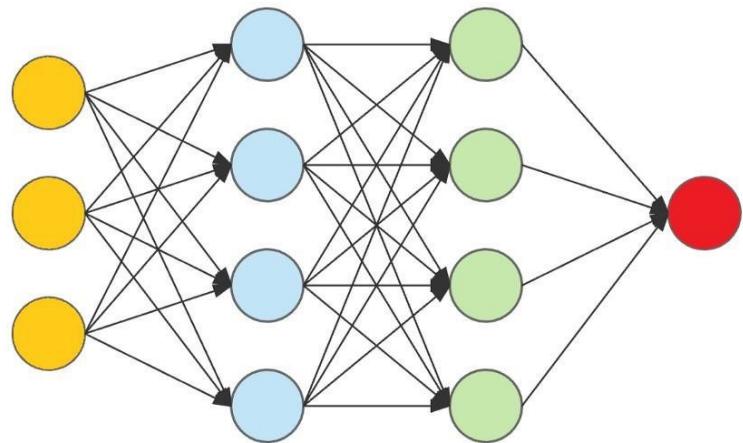
Regularization & Optimization

Deep Learning Basics

- Image Classification: A core task in Computer Vision



→ cat



Neural Networks

This image by [Nikita](#) is
licensed under [CC-BY 2.0](#)

CS231n overview

- Deep Learning Basics
- Perceiving and Understanding the Visual World
- Generative and Interactive Visual Intelligence
- Human-Centered Applications and Implications

CS231n overview

- Deep Learning Basics
- Perceiving and Understanding the Visual World
- Generative and Interactive Visual Intelligence
- Human-Centered Applications and Implications

Perceiving and Understanding the Visual World



Tasks

Models

Tasks Beyond Image Classification

Classification



CAT

No spatial extent

Semantic Segmentation



GRASS, CAT, TREE,
SKY

No objects, just pixels

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



DOG, DOG, CAT

[This image](#) is [CC0 public domain](#)

Tasks Beyond Image Classification

Video
Classification

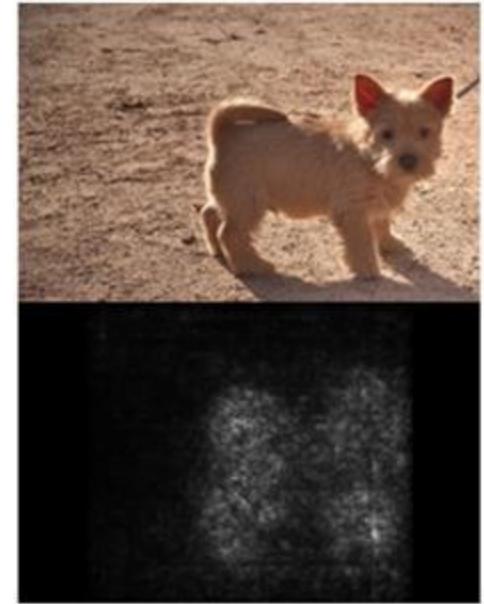


Running? Jumping?

Multimodal Video
Understanding



Visualization &
Understanding



Models Beyond Multi-Layer Perceptron

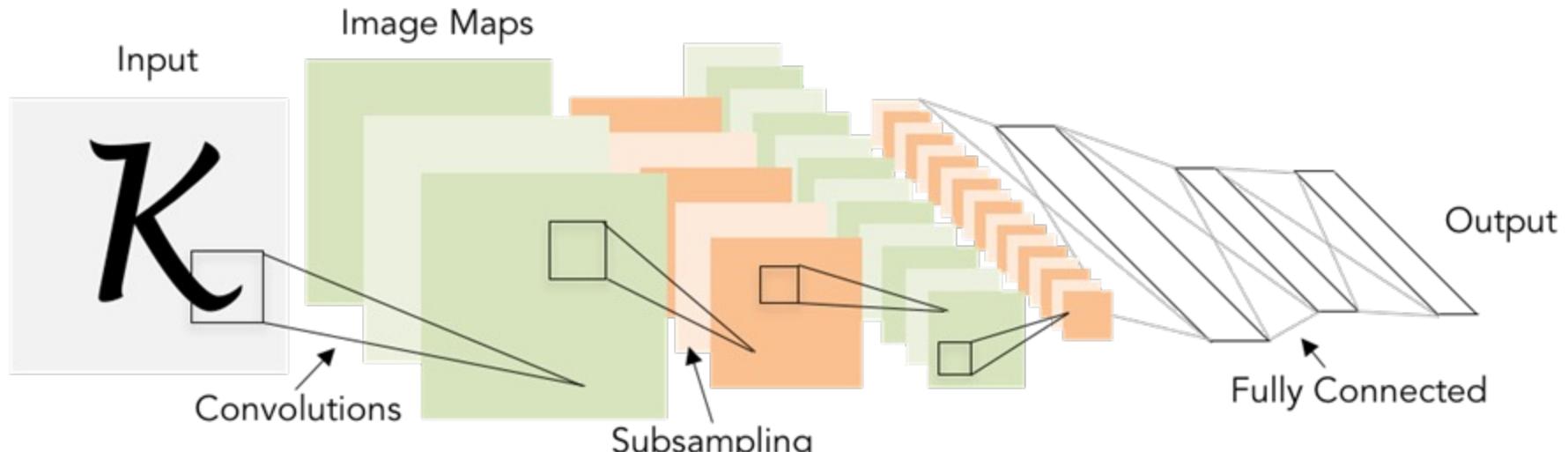
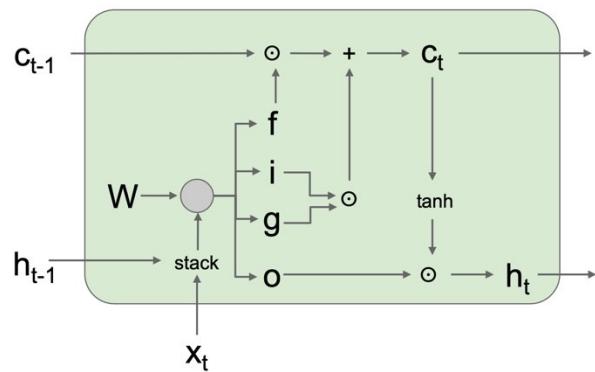
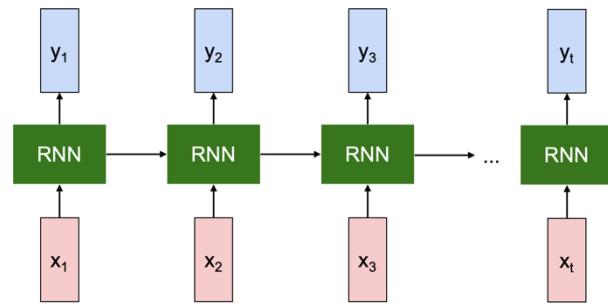


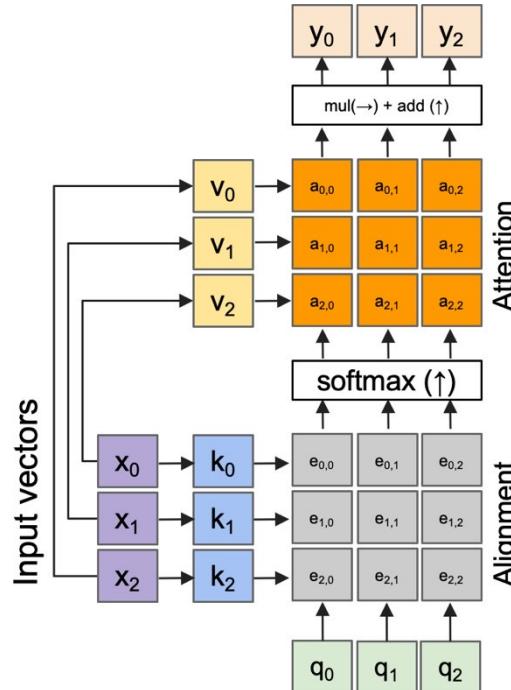
Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

Convolutional neural network

Models Beyond Multi-Layer Perceptron



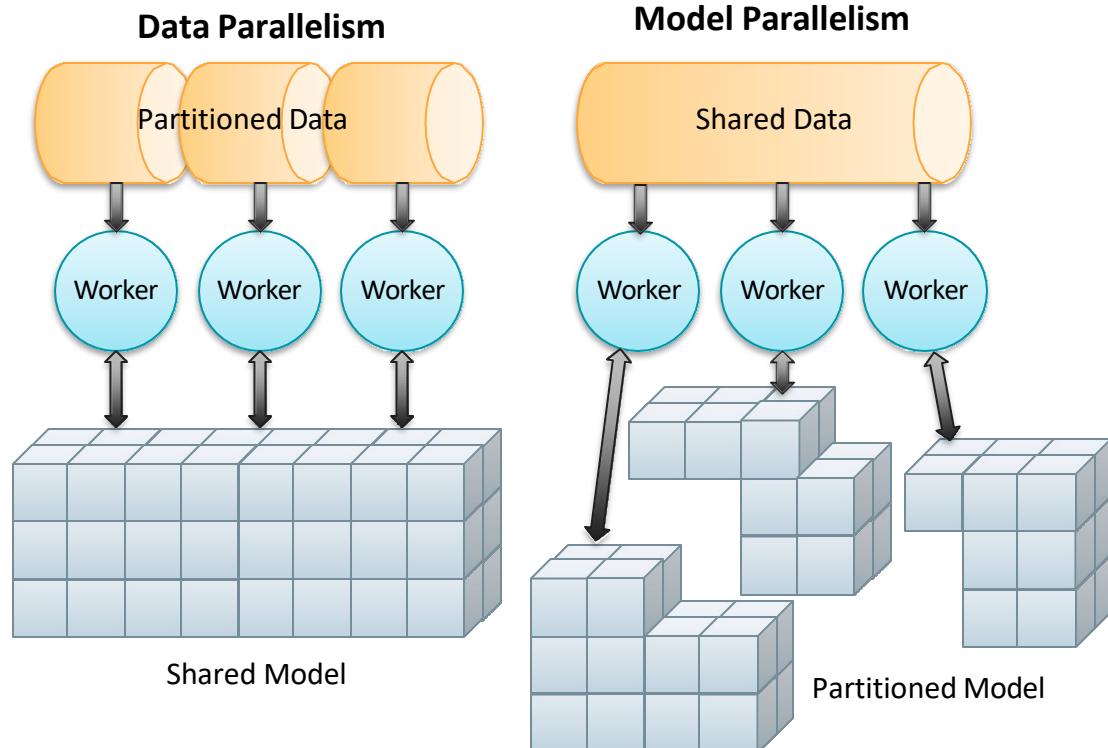
Recurrent neural network



Attention mechanism / Transformers

Large Scale Distributed Training

- Train Large Models on big datasets faster
- Scale beyond single GPU/machine limitations
- How?
 - Data Parallelism: Copy the model to all workers, split the data
 - Model Parallelism: Split model across devices
 - Synchronous vs. Asynchronous gradient updates



CS231n overview

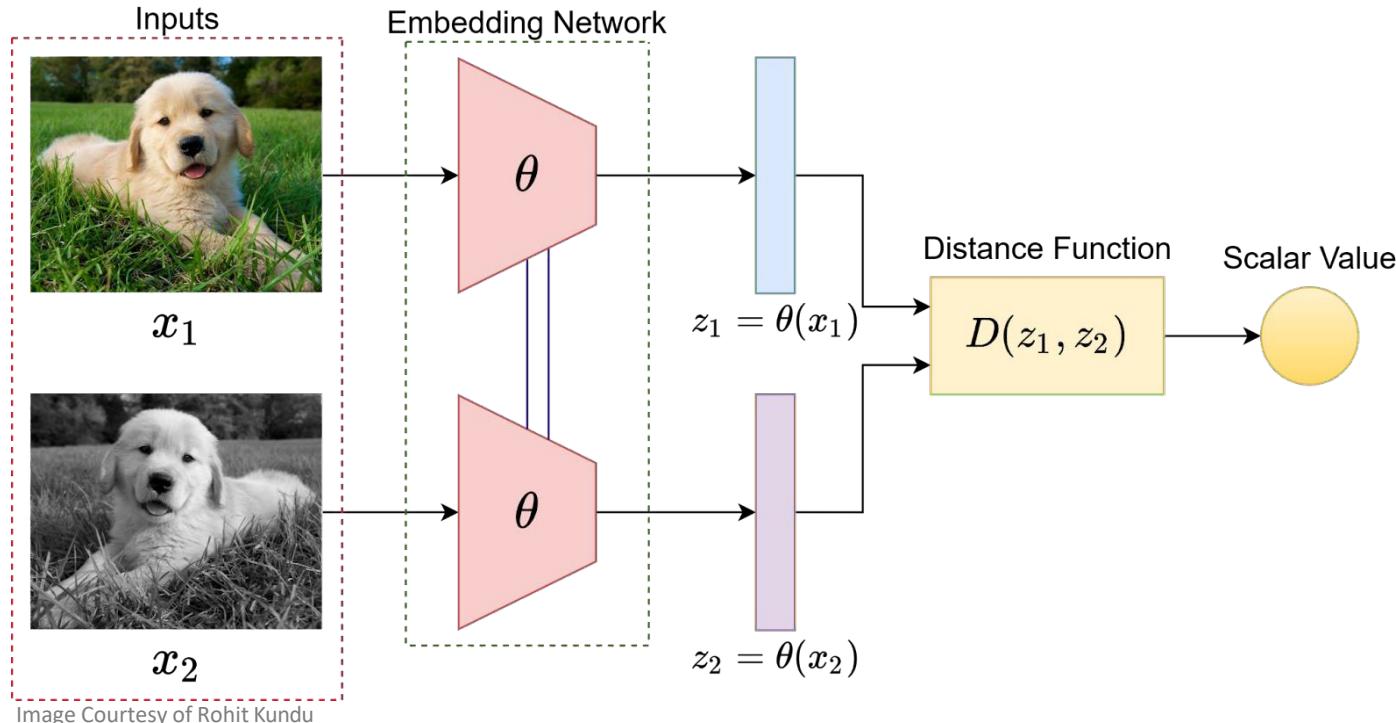
- Deep Learning Basics
- Perceiving and Understanding the Visual World
- Generative and Interactive Visual Intelligence
- Human-Centered Applications and Implications

CS231n overview

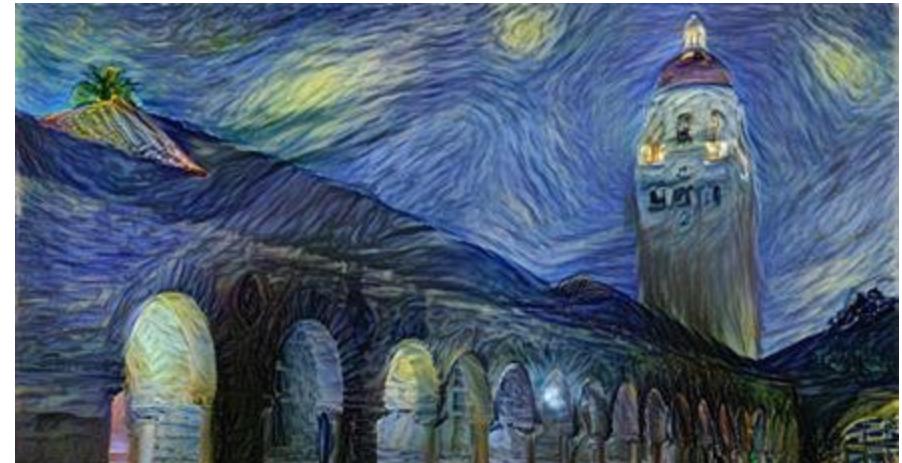
- Deep Learning Basics
- Perceiving and Understanding the Visual World
- **Generative and Interactive Visual Intelligence**
- Human-Centered Applications and Implications

Beyond 2D Recognition

Beyond 2D Recognition: Self-supervised Learning



Beyond 2D Recognition: Generative Modeling



Style Transfer

Beyond 2D Recognition: Generative Modeling

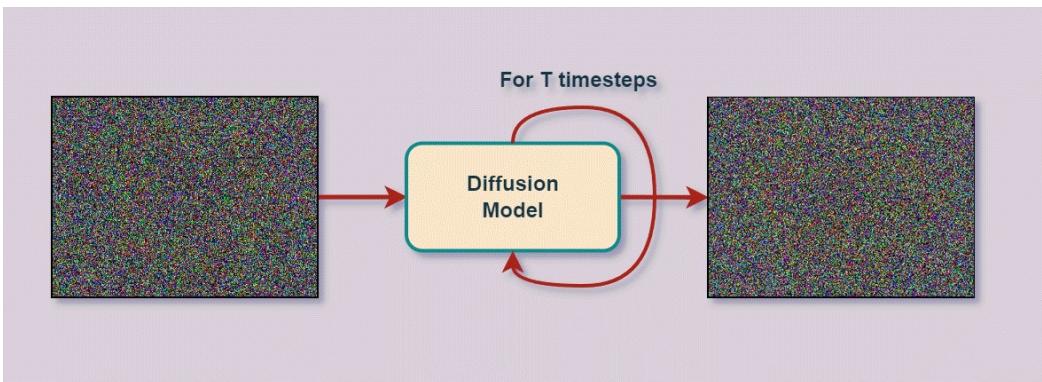


“Teddy bears working on new
AI research underwater with
1990s technology”

DALL-E 2

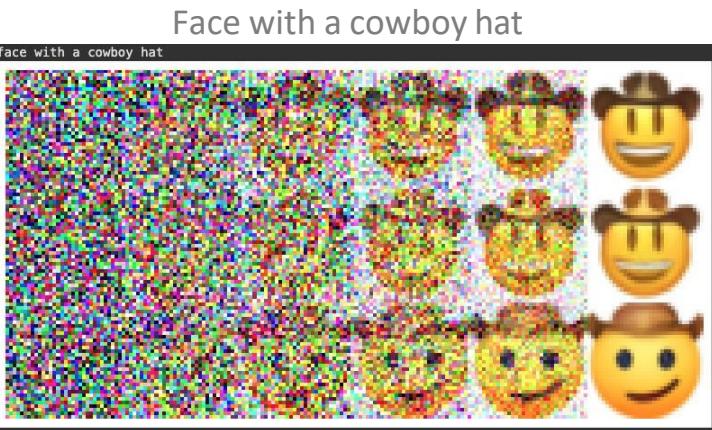
Beyond 2D Recognition: Generative Modeling

Image Generation using Diffusion Models

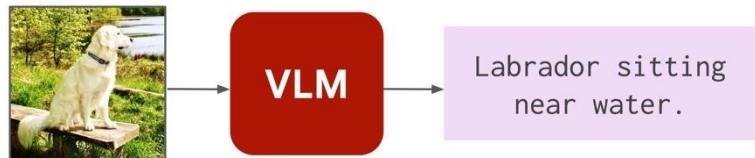
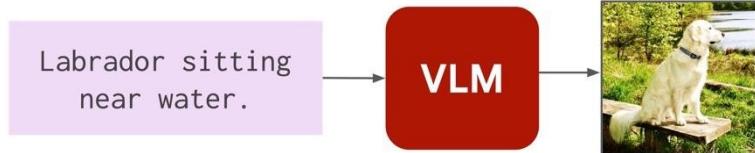


<https://learnopencv.com/image-generation-using-diffusion-models/>

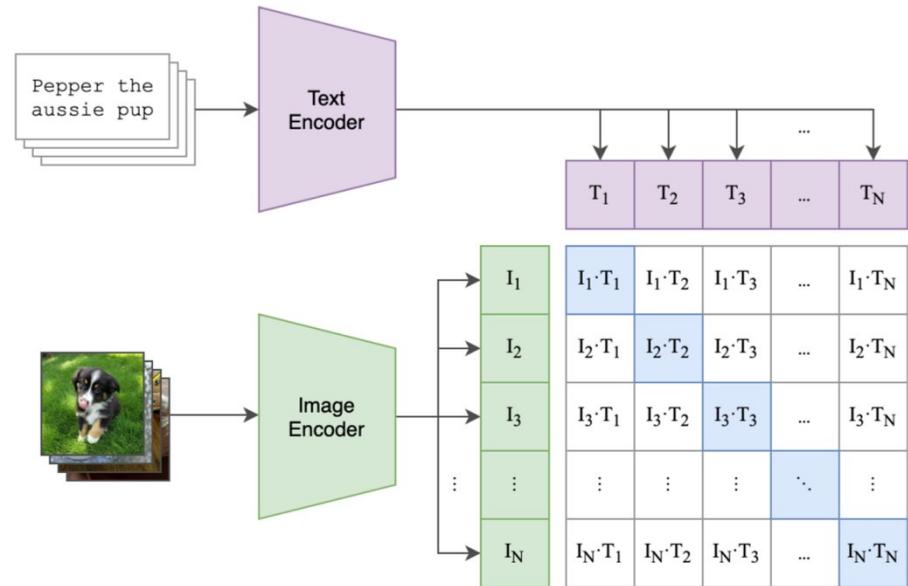
You will learn and implement a generative model in Assignment 3 that generates emojis from text inputs



Beyond 2D Recognition: Vision Language Models



Yasunaga, Michihiro, et al. "Retrieval-augmented multimodal language modeling." arXiv preprint arXiv:2211.12561 (2022).



Contrastive pre-training in CLIP. The blue squares are the pairs for which we want to optimize the similarity. Image derived from <https://github.com/openai/CLIP>

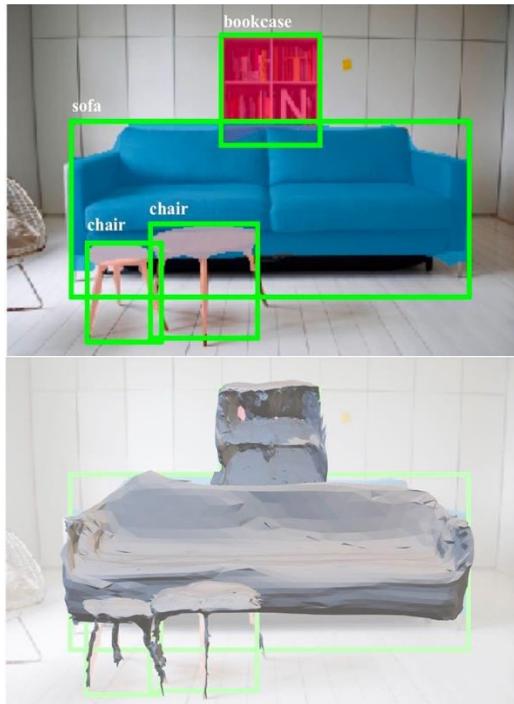
Beyond 2D Recognition: 3D Vision



Choy et al., 3D-R2N2: Recurrent Reconstruction Neural Network (2016)



Zhou et al., 3D Shape Generation and Completion through Point-Voxel Diffusion (2021)



Gkioxari et al., "Mesh R-CNN", ICCV 2019

Beyond 2D Recognition: Embodied Intelligence



Clean Your House After a Wild Party

BEHAVIOR Task #1

Li et al., BEHAVIOR-1K: A Benchmark for Embodied AI with 1,000 Everyday Activities and Realistic Simulation (2022)



Mandlekar and Xu et al., Learning to Generalize Across Long-Horizon Tasks from Human Demonstrations (2020)

CS231n overview

- Deep Learning Basics
- Perceiving and Understanding the Visual World
- **Generative and Interactive Visual Intelligence**
- Human-Centered Applications and Implications

CS231n overview

- Deep Learning Basics
- Perceiving and Understanding the Visual World
- Generative and Interactive Visual Intelligence
- Human-Centered Applications and Implications

2018 Turing Award for deep learning

most prestigious technical award, is given for major contributions of lasting importance to computing.



Geoffrey Hinton



Yoshua Bengio



Yann LeCun

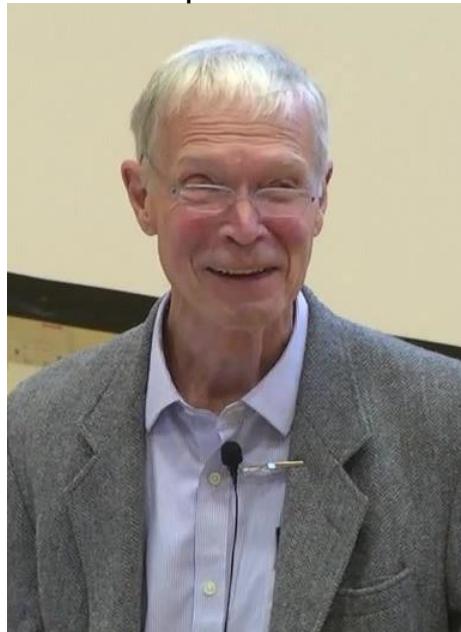
2024 Nobel Prize in Physics

Hinton speaking at the Nobel Prize
Lectures in Stockholm in 2024



[This image](#) is [CC0 public domain](#)

John Hopfield in 2016



[This image](#) is [CC0 public domain](#)

In 2024, he was jointly awarded the [Nobel Prize in Physics](#) with [John Hopfield](#) “for foundational discoveries and inventions that enable machine learning with artificial neural networks.”

Learning objectives

Formalize computer vision applications into tasks

- Formalize inputs and outputs for vision-related problems
- Understand what data and computational requirements you need to train a model

Develop and train vision models

- Learn to code, debug, and train convolutional neural networks.
- Learn how to use software frameworks like PyTorch and TensorFlow

Gain an understanding of where the field is and where it is headed

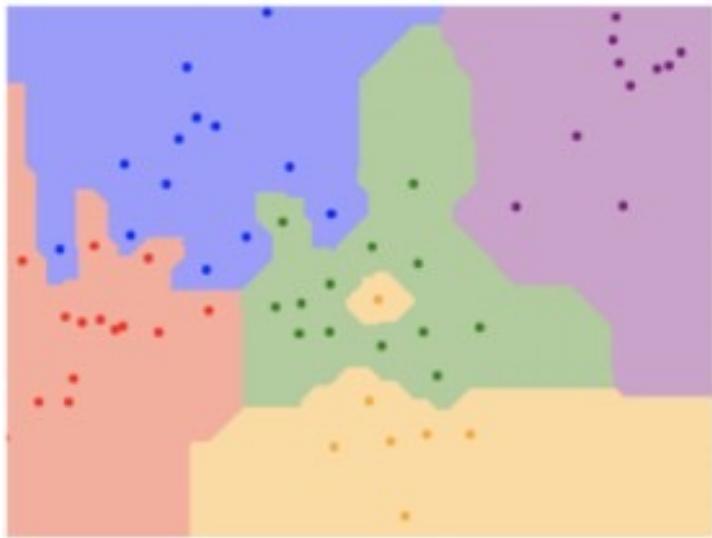
- What new research has come out in the last 0-5 years?
- What are open research challenges?
- What ethical and societal considerations should we consider before deployment?

CS231n: Deep Learning for Computer Vision

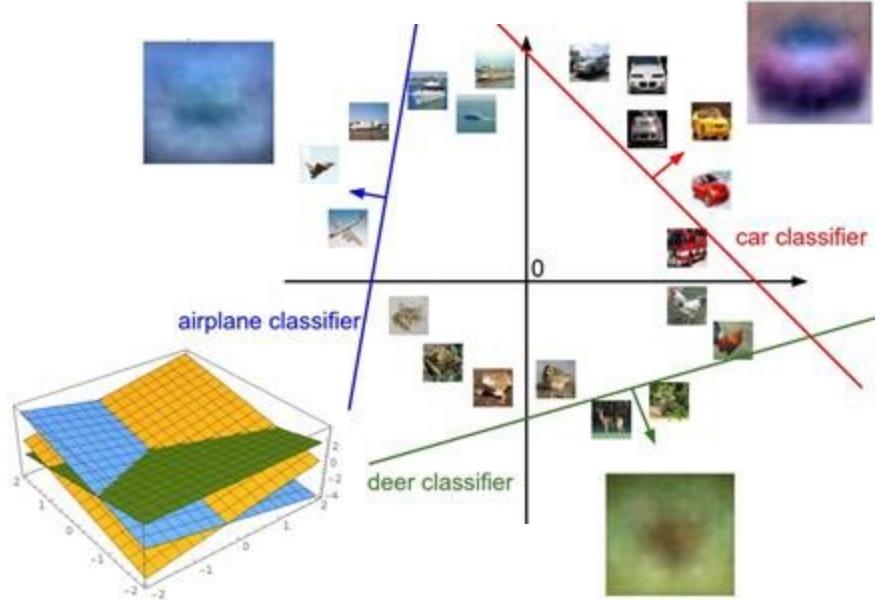
- Deep Learning Basics (Lecture 2 – 4)
- Perceiving and Understanding the Visual World (Lecture 5 – 12)
- Reconstructing and Interacting with the Visual World (Lecture 13 – 17)
- Human-Centered Artificial Intelligence (Lecture 18)

Next time: Image classification with Linear Classifiers

k- nearest neighbor



Linear classification



Plot created using [Wolfram Cloud](#)