

CS7.601

Deep Learning : Theory and Practices

Monsoon 2023

Naresh Manwani
Machine Learning Lab
IIIT Hyderabad

What Is It?



WIKIPEDIA
The Free Encyclopedia

Deep learning

From Wikipedia, the free encyclopedia

Deep learning (*deep machine learning*, or *deep structured learning*, or *hierarchical learning*, or sometimes *DL*) is a branch of [machine learning](#) based on a set of [algorithms](#) that attempt to model high-level abstractions in data by using model architectures, with complex structures or otherwise, composed of [multiple non-linear transformations](#).[\[1\]](#)(p198)[\[2\]](#)[\[3\]](#)[\[4\]](#)

Learning about Deep Neural Networks

Yann Lecun quote: DNNs require: “an interplay between intuitive insights, theoretical modeling, practical implementations, empirical studies, and scientific analyses”

i.e. there isn’t a framework or core set of principles to explain everything (c.f. graphical models for machine learning).

We try to cover the ground in Lecun’s quote.

This Course

Goals:

- Introduce deep learning to a broad audience.
- Review principles and techniques for understanding deep networks.
- Develop skill at designing networks for applications.

Prerequisites

- Knowledge of calculus and linear algebra, probability theory and optimization methods
- Statistical Methods in AI (Must have done)
- Programming: assignments will mostly use Python.

Course Content

CO-1: Perceptron, convergence proof. Feedforward neural network, Representation power of feedforward neural network (Universal Approximation Theorem, limitations of shallow networks. [4 Lectures]

CO-2: Back propagation, loss surfaces, learning rates, optimization for deep networks: gradient descent (GD), momentum based GD, Nesterov accelerated GD, stochastic GD, AdaGrad, RMSProp, Adam. [3 Lectures]

CO-3: Greedy layerwise pre-training, better activation functions, better weight initialization methods, batch normalization [2 Lecture]

CO-4: Bias variance tradeoff: overfitting and under-fitting. L2 regularization, early stopping, dataset augmentation, parameter sharing and tying, injecting noise at input, ensemble methods, dropout. [3 Lectures]

CO-5: Auto-encoders and relation to PCA, regularization in auto-encoders, denoising auto-encoders, sparse auto-encoders, contractive auto-encoders, variational auto-encoders (VAEs), mutual information and the information bottleneck. Word2vec and its relationship to latent semantic indexing (LSI). Unsupervised Learning, Restricted Boltzmann Machines (RBMs), Contrastive divergence for RBMs. Generative Adversarial Networks (GANs), Diffusion Models. [5 Lectures]

Course Content

CO-6: Convolutional neural networks (CNNs), backpropagation in CNNs, LeNet, AlexNet, Inception, VGG, GoogLeNet, ResNet. [2 Lectures]

CO-7: Recurrent neural networks, backpropagation through time (BPTT), vanishing and exploding gradients, truncated BPTT, stability, bidirectional RNNs, gated recurrent units (GRUs), long short term memory (LSTM), solving the vanishing gradient problem with LSTMs. [3 Lectures]

CO-8: Encoder Decoder Models, Attention Mechanism, Hierarchical Attention, Transformers, Graph Neural Networks. [4 Lectures]

.

This Course

Books :

1. Simon Haykin. 1998. Neural Networks: A Comprehensive Foundation (2nd ed.), Prentice Hall PTR, Upper Saddle River, NJ, USA.
2. Ian Goodfellow and Yoshua Bengio and Aaron Courville, Deep Learning, MIT Press, 2016.
3. R. Rojas: Neural Networks, Springer-Verlag, Berlin, 1996.
4. Zhang, A., Lipton, Z.-C., Li, M., Smola, A.-J.\ 2021.\ Dive into Deep Learning.

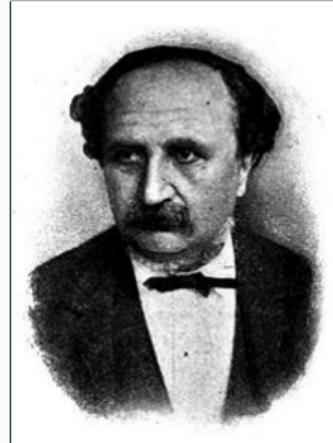
Grading Policy

| Type of Evaluation | Weightage (in %) |
|-----------------------------------|------------------|
| Quiz 1 | 7.5% |
| Quiz 2 | 7.5% |
| Mid Set | 15% |
| End Sem | 30% |
| Assignments-4 | 32% |
| Scribing (2 People for 1 Lecture) | 8% |

Chapter 1: Biological Neurons

Reticular Theory

Joseph von Gerlach proposed that the nervous system is a single continuous network as opposed to a network of many discrete cells!



1871-1873



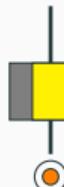
Reticular theory

Staining Technique

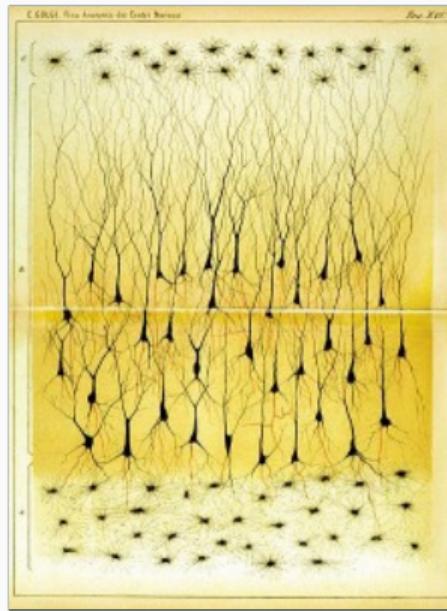
Camillo Golgi discovered a chemical reaction that allowed him to examine nervous tissue in much greater detail than ever before

He was a proponent of Reticular theory.

1871-1873

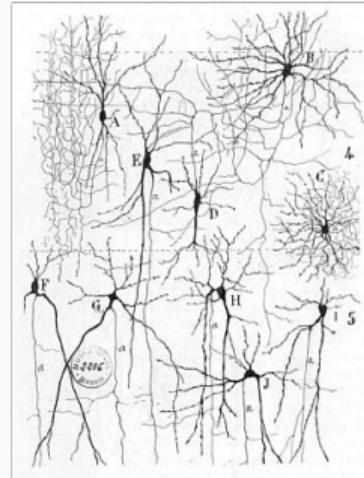


Reticular theory



Neuron Doctrine

Santiago Ramón y Cajal used Golgi's technique to study the nervous system and proposed that it is actually made up of discrete individual cells forming a network (as opposed to a single continuous network)

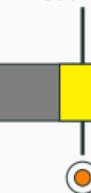


1871-1873



Reticular theory

1888-1891

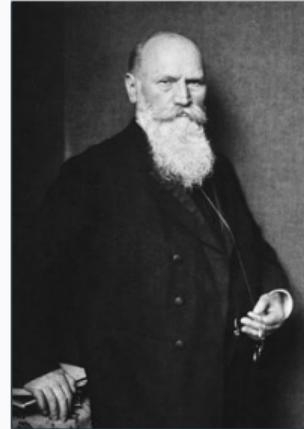


Neuron Doctrine

The Term Neuron

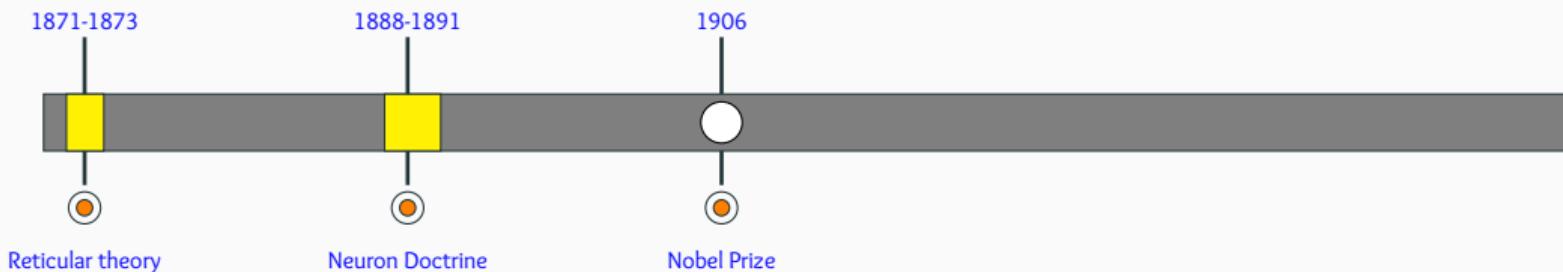
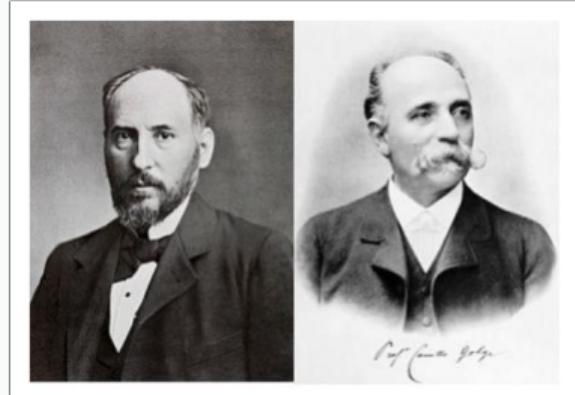
The term neuron was coined by Heinrich Wilhelm Gottfried von Waldeyer-Hartz around 1891.

He further consolidated the Neuron Doctrine.



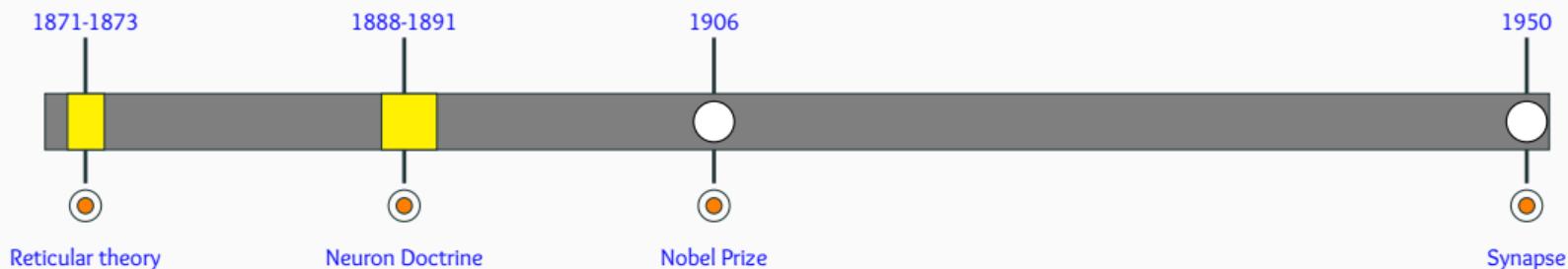
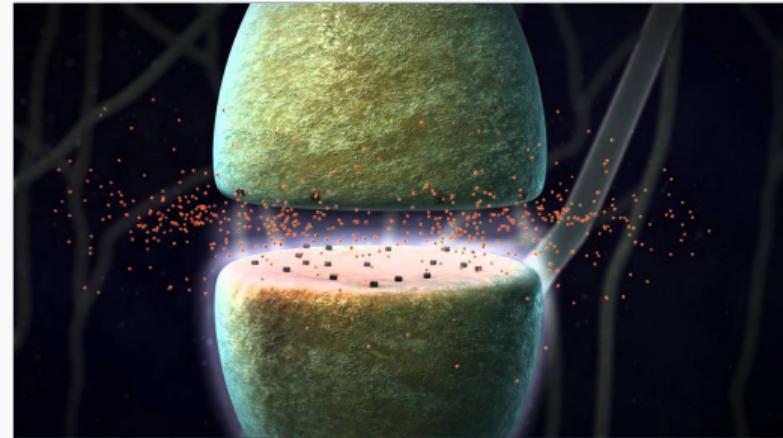
Nobel Prize

Both Golgi (reticular theory) and Cajal (neuron doctrine) were jointly awarded the 1906 Nobel Prize for Physiology or Medicine, that resulted in lasting conflicting ideas and controversies between the two scientists.



The Final Word

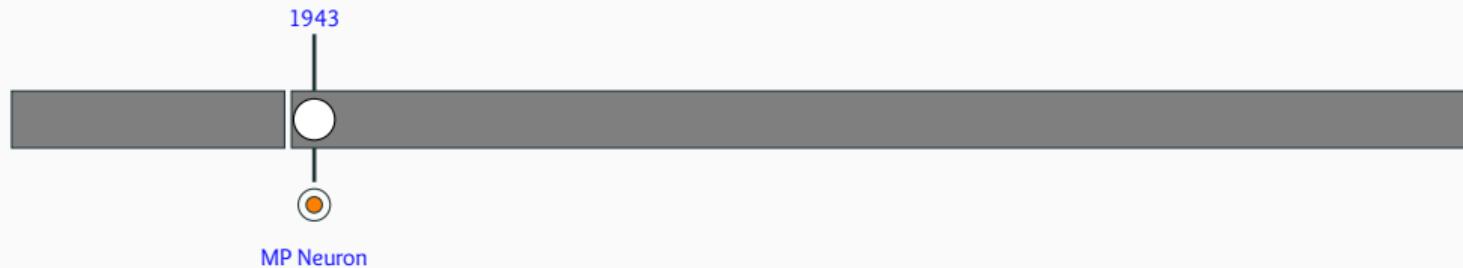
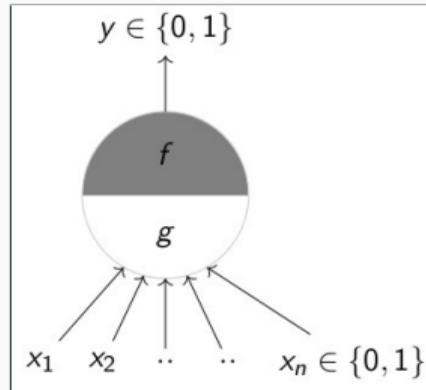
In 1950s electron microscopy finally confirmed the neuron doctrine by unambiguously demonstrating that nerve cells were individual cells interconnected through synapses (a network of many individual neurons).



Chapter 2: From Spring to Winter of AI

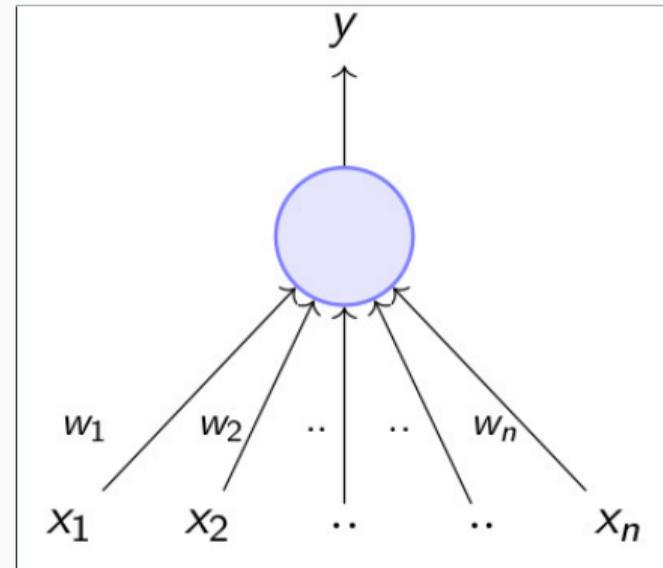
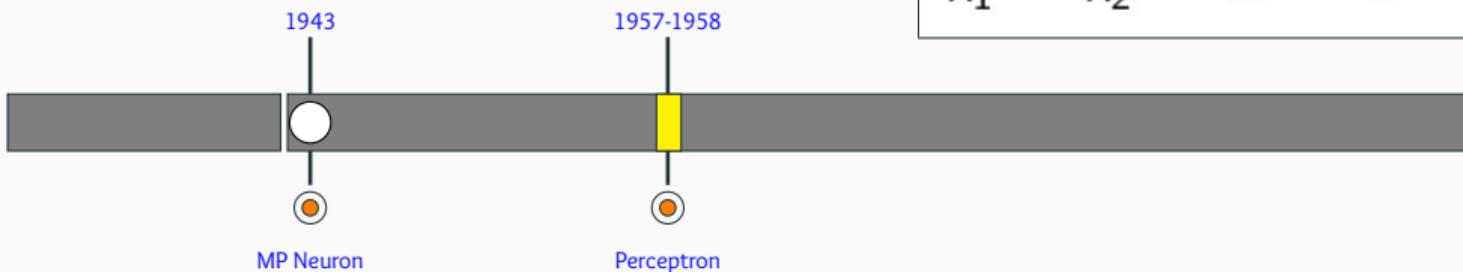
McCulloch Pitts Neuron

McCulloch (neuroscientist) and Pitts (logician) proposed a highly simplified model of the neuron (1943)^[2]



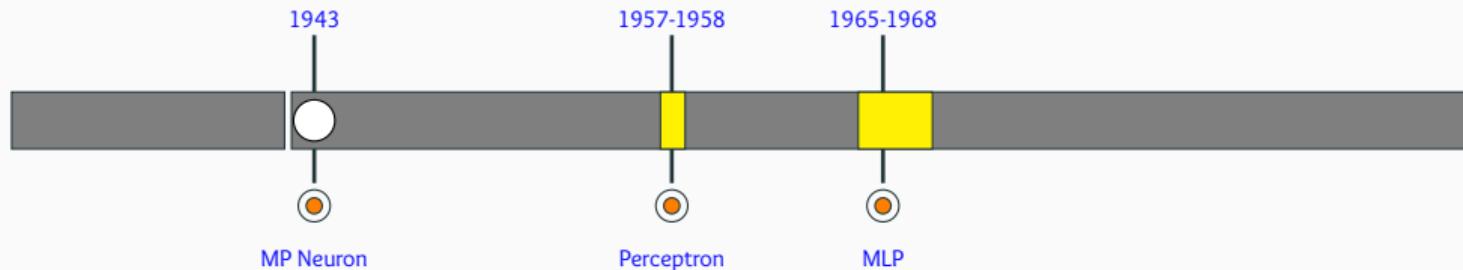
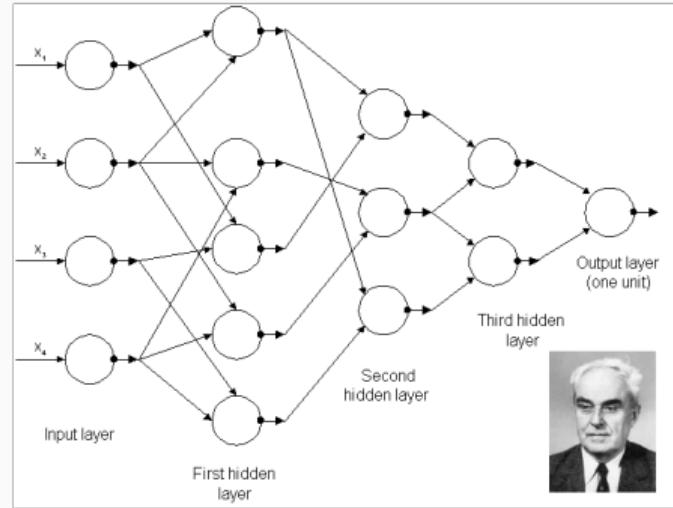
Perceptron

“the perceptron may eventually be able to learn, make decisions, and translate languages” -Frank Rosenblatt



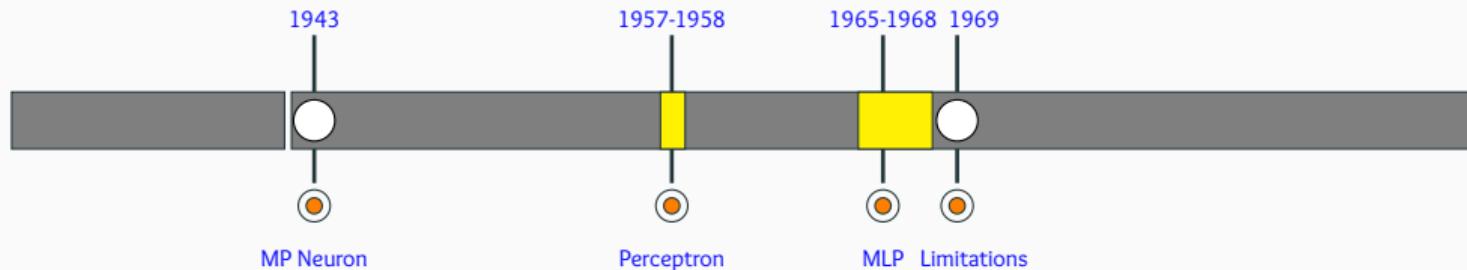
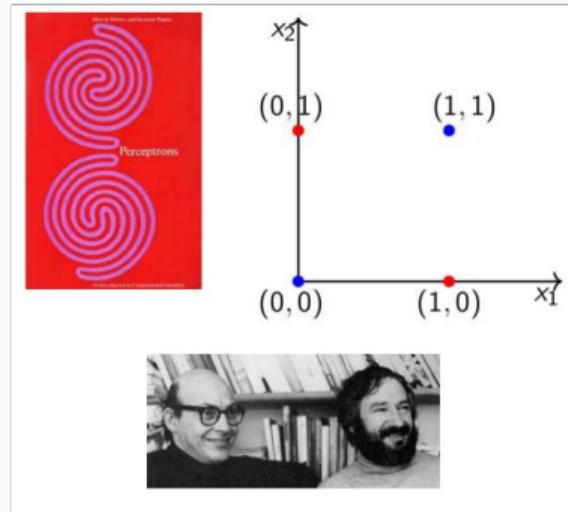
First generation Multilayer Perceptrons

Ivakhnenko et. al. [3]



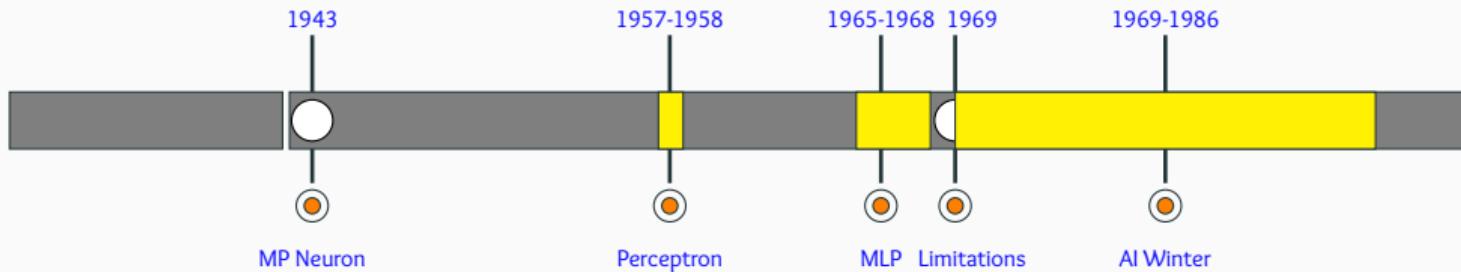
Perceptron Limitations

In their now famous book “Perceptrons”, Minsky and Papert outlined the limits of what perceptrons could do^[4]



AI Winter of connectionism

Almost lead to the abandonment of
connectionist AI

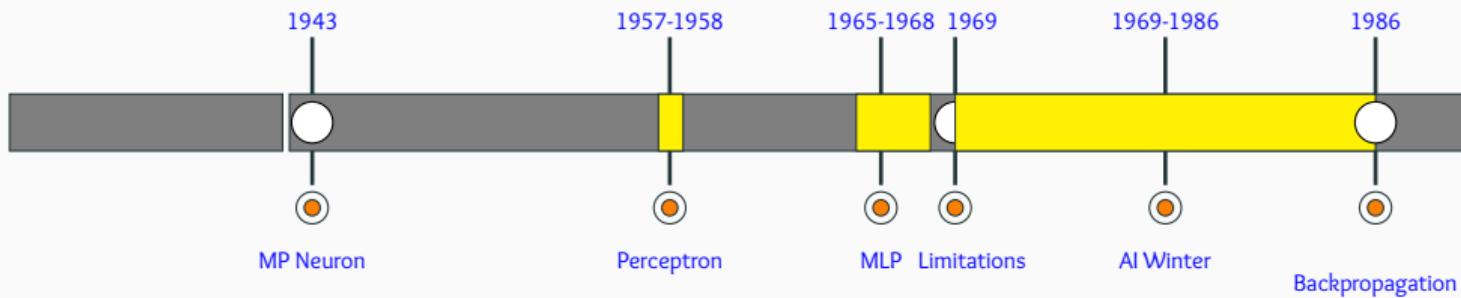
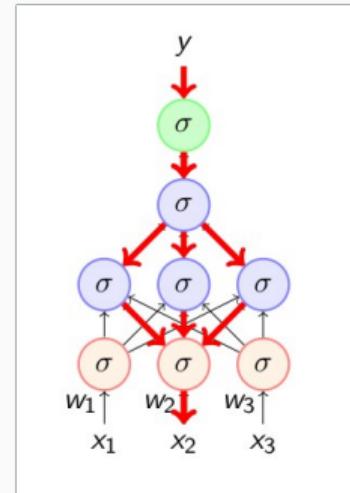


Backpropagation

Discovered and rediscovered several times throughout 1960's and 1970's

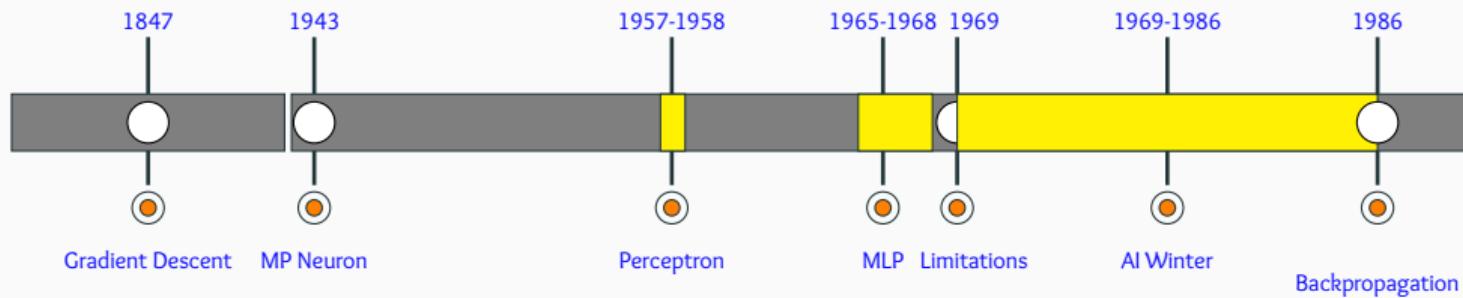
Werbos(1982)^[5] first used it in the context of artificial neural networks

Eventually popularized by the work of Rumelhart et. al. in 1986^[6]



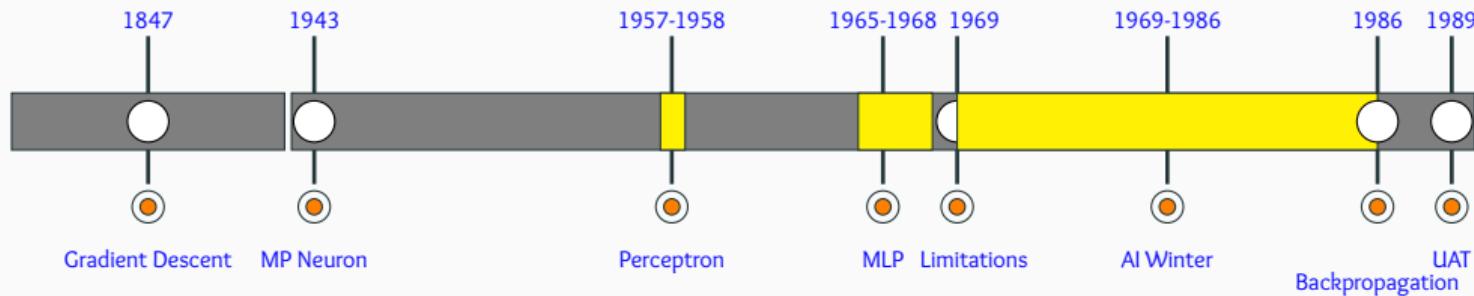
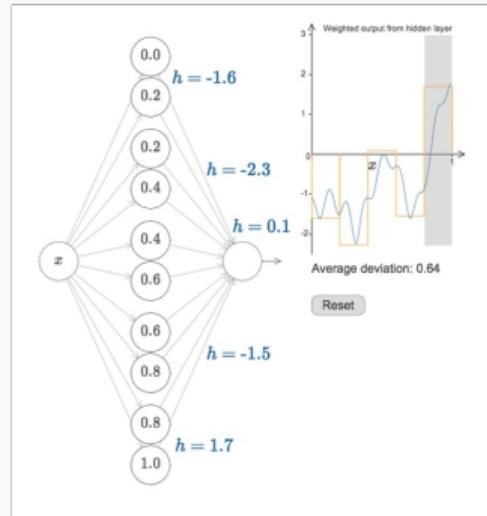
Gradient Descent

Cauchy discovered Gradient Descent
motivated by the need to compute the orbit
of heavenly bodies



Universal Approximation Theorem

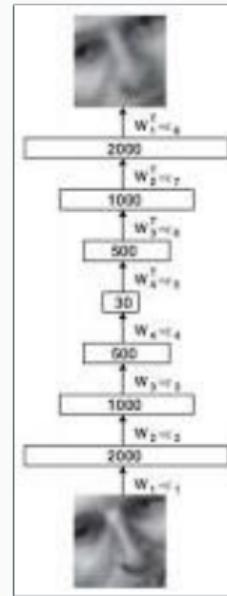
A multilayered network of neurons with a single hidden layer can be used to approximate any continuous function to any desired precision [7]



Chapter 3: The Deep Revival

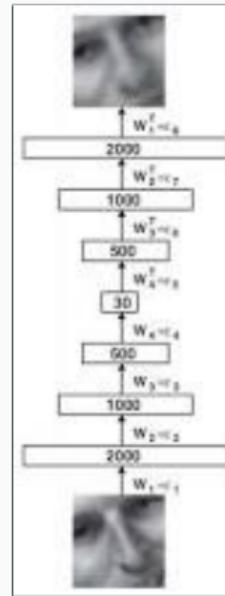
Unsupervised Pre-Training

Hinton and Salakhutdinov described an effective way of initializing the weights that allows deep autoencoder networks to learn a low-dimensional representation of data. [8]



Unsupervised Pre-Training

The idea of unsupervised pre-training actually dates back to 1991-1993 (J. Schmidhuber) when it was used to train a “Very Deep Learner”



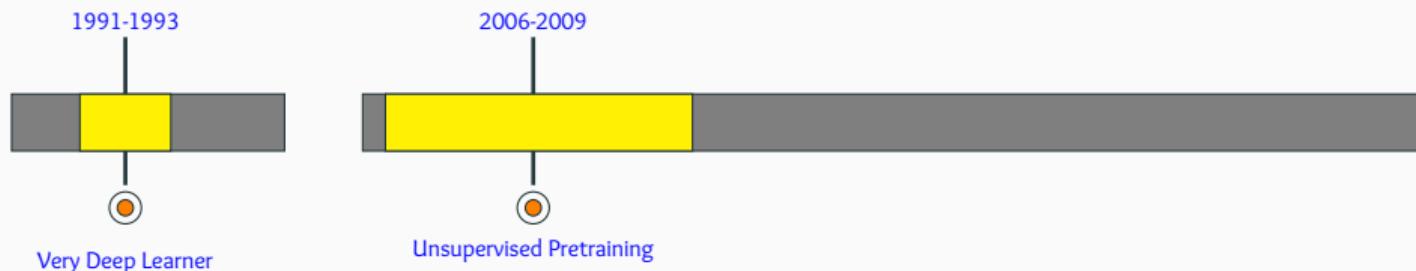
More insights (2007-2009)

Further Investigations into the effectiveness
of Unsupervised Pre-training

Greedy Layer-Wise Training of Deep Networks

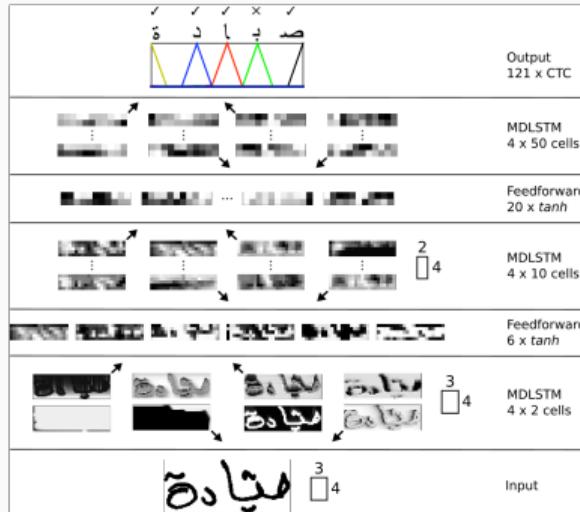
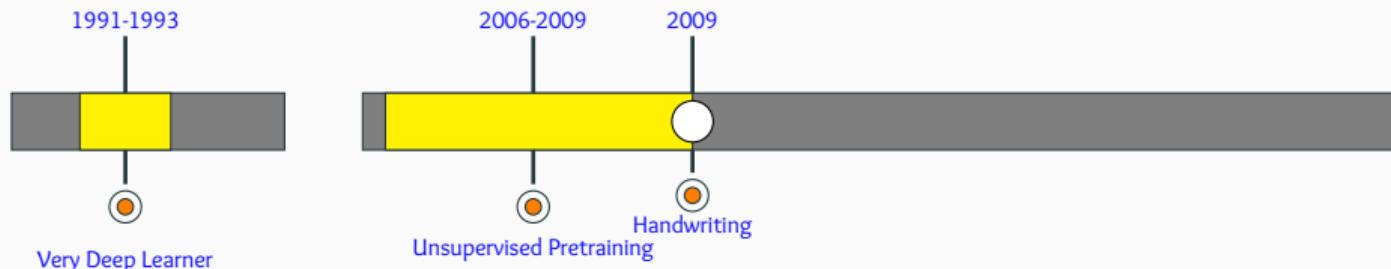
Why Does Unsupervised Pre-training Help Deep Learning?

Exploring Strategies for Training Deep Neural Networks



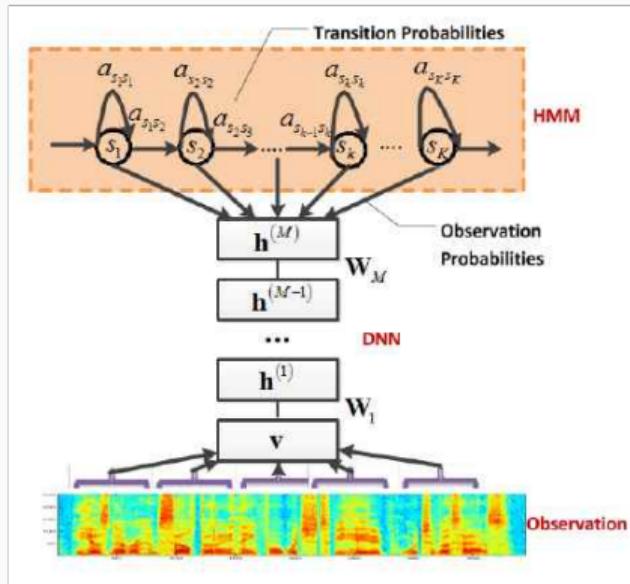
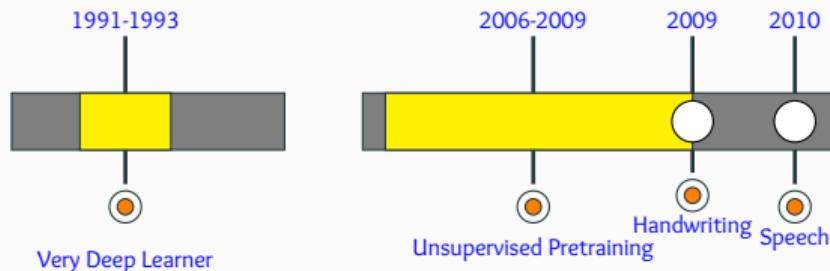
Success in Handwriting Recognition

Graves et. al. outperformed all entries in an international Arabic handwriting recognition competition [9]



Success in Speech Recognition

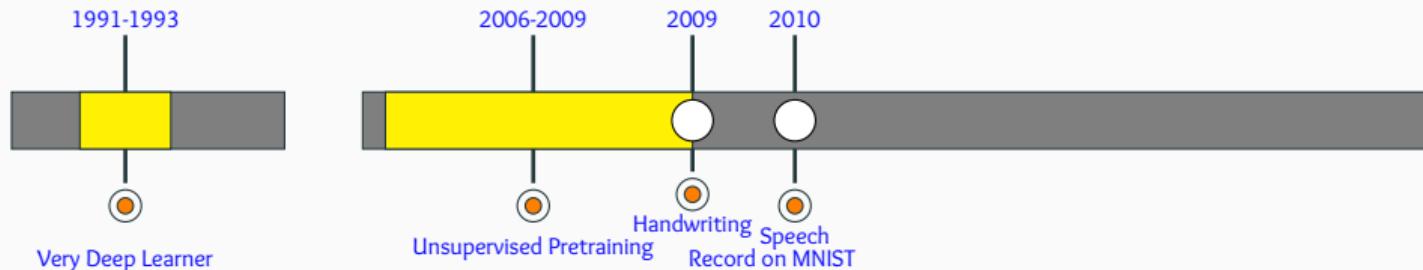
Dahl et. al. showed relative error reduction of 16.0% and 23.2% over a state of the art system [10]



New record on MNIST

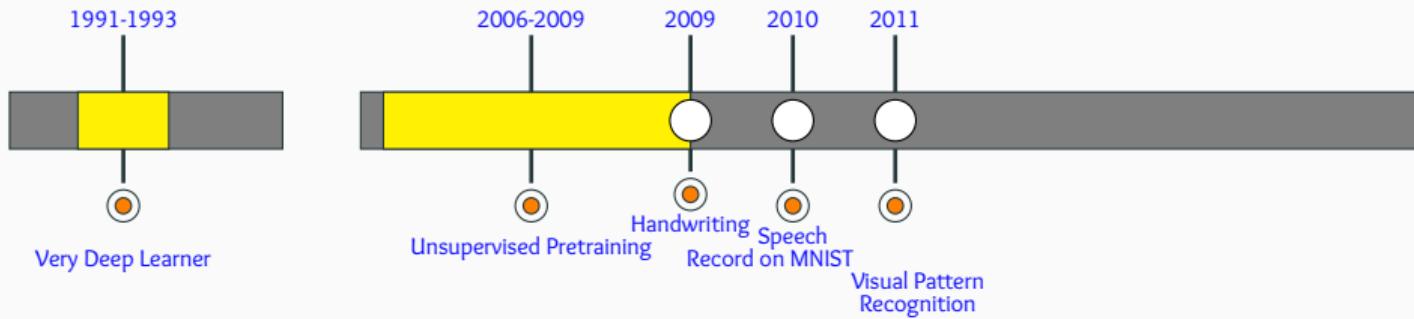
Ciresan et. al. set a new record on the MNIST dataset using good old backpropagation on GPUs (GPUs enter the scene) [11]

| | | | | | | |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| 1 2 17 | 7 1 71 | 9 8 98 | 9 9 59 | 9 9 79 | 5 5 35 | 3 8 23 |
| 4 9 49 | 3 5 35 | 9 4 97 | 4 9 49 | 4 4 94 | 2 2 02 | 5 5 35 |
| 1 6 16 | 9 4 94 | 0 0 60 | 6 6 06 | 6 6 86 | 1 1 79 | 1 1 71 |
| 4 9 49 | 0 0 50 | 3 5 35 | 8 8 98 | 3 9 79 | 7 7 17 | 1 1 61 |
| 2 7 27 | 8 8 58 | 2 2 78 | 1 6 16 | 6 5 65 | 4 4 94 | 0 0 60 |
| | | | | | | |



First Superhuman Visual Pattern Recognition

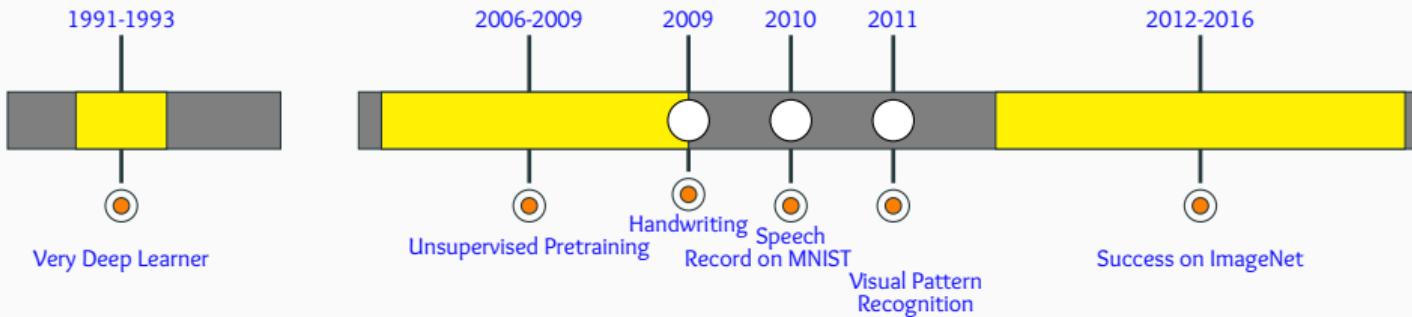
D. C. Ciresan et. al. achieved 0.56% error rate in the IJCNN Traffic Sign Recognition Competition^[12]



Winning more visual recognition challenges



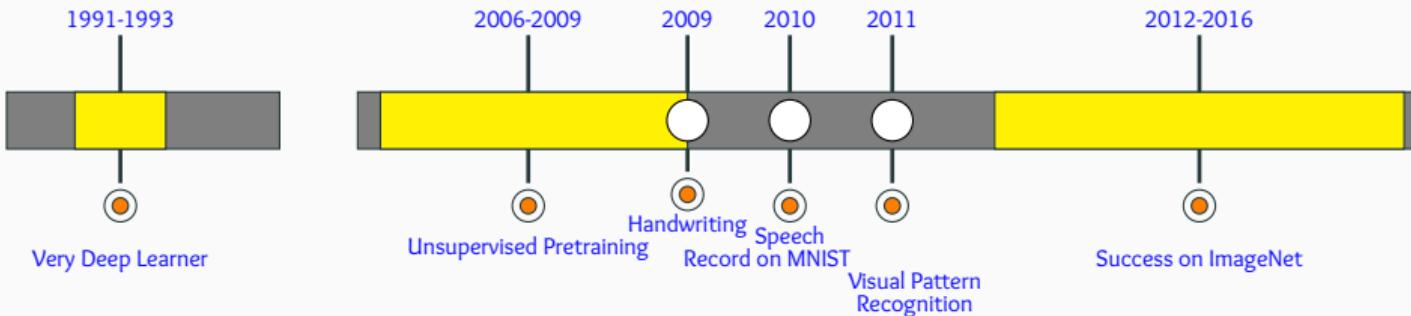
| Network | Error | Layers |
|-------------------------|-------|--------|
| AlexNet ^[13] | 16.0% | 8 |



Winning more visual recognition challenges



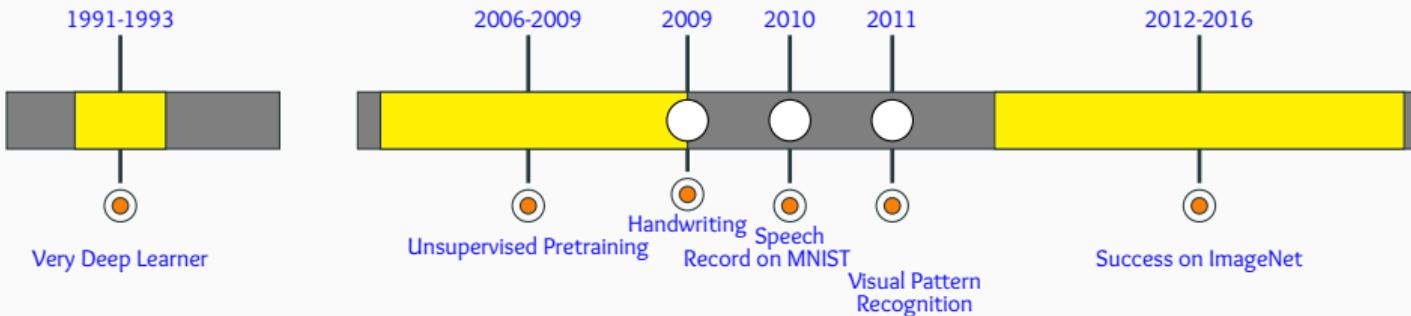
| Network | Error | Layers |
|-------------------------|-------|--------|
| AlexNet ^[13] | 16.0% | 8 |
| ZFNet ^[14] | 11.2% | 8 |



Winning more visual recognition challenges



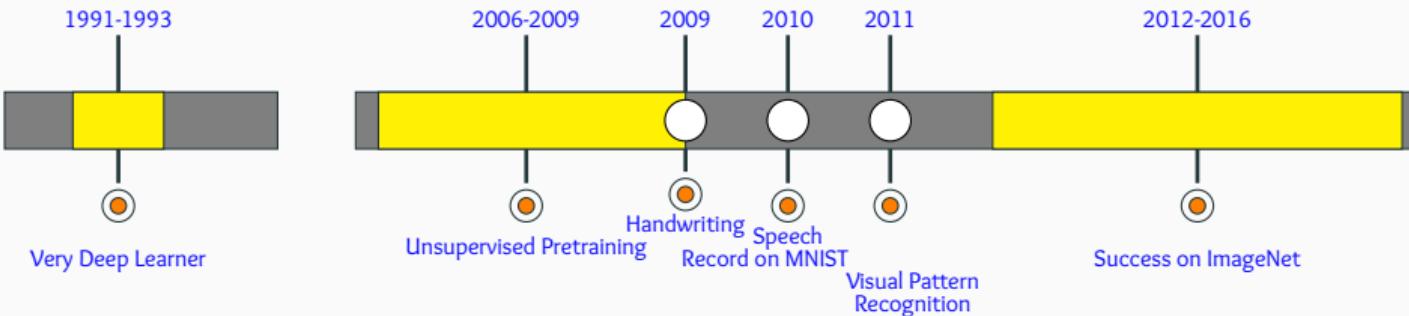
| Network | Error | Layers |
|-------------------------|-------|--------|
| AlexNet ^[13] | 16.0% | 8 |
| ZFNet ^[14] | 11.2% | 8 |
| VGGNet ^[15] | 7.3% | 19 |



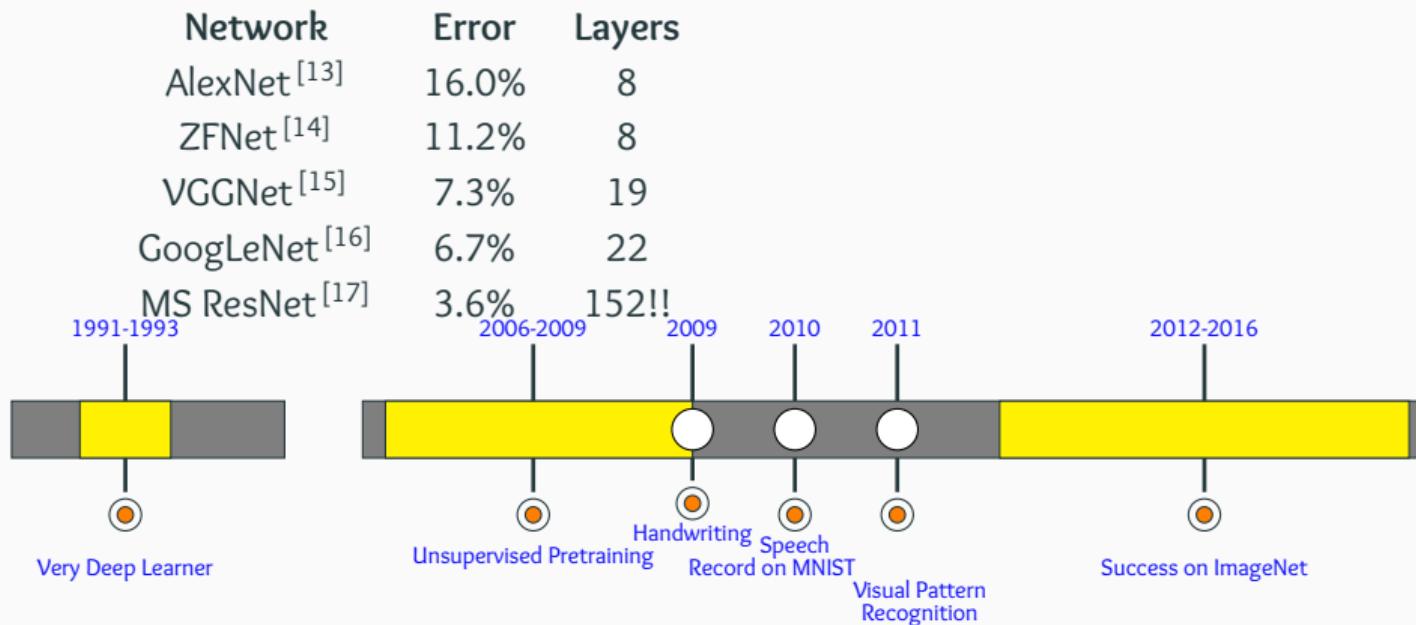
Winning more visual recognition challenges



| Network | Error | Layers |
|---------------------------|-------|--------|
| AlexNet ^[13] | 16.0% | 8 |
| ZFNet ^[14] | 11.2% | 8 |
| VGGNet ^[15] | 7.3% | 19 |
| GoogLeNet ^[16] | 6.7% | 22 |



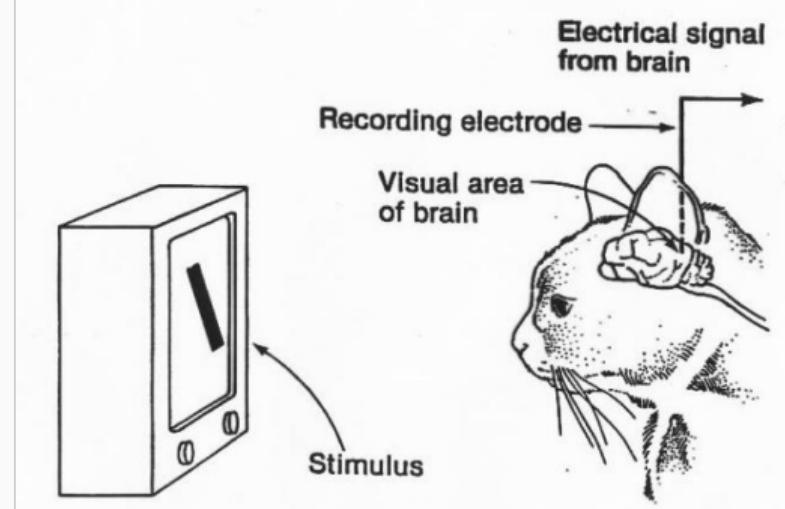
Winning more visual recognition challenges



Chapter 4: From Cats to Convolutional Neural Networks

Hubel and Wiesel Experiment

Experimentally showed that each neuron has a fixed receptive field - i.e. a neuron will fire only in response to a visual stimuli in a specific region in the visual space^[18]



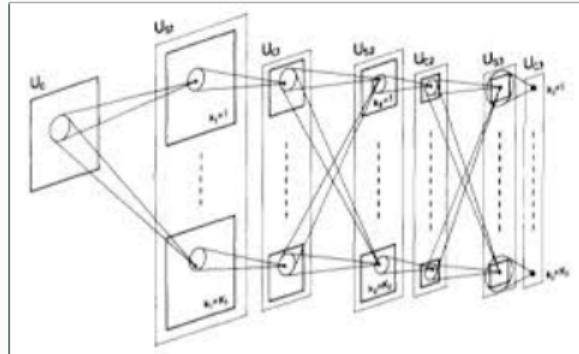
1959



H and W experiment

Neocognitron

Used for Handwritten character recognition
and pattern recognition (Fukushima et.
al.) [19]



1959



H and W experiment

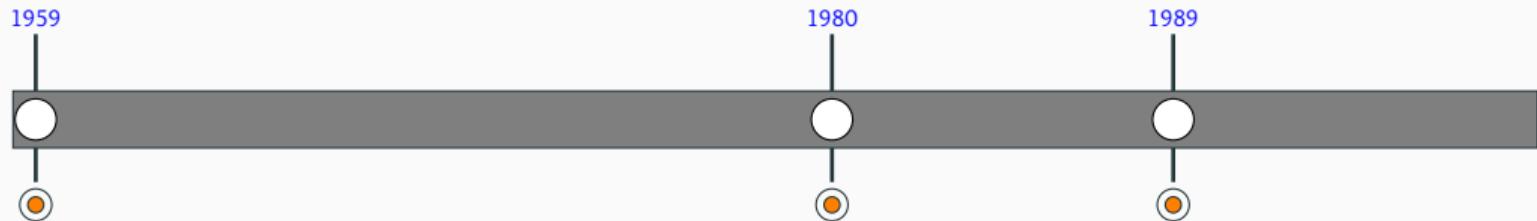
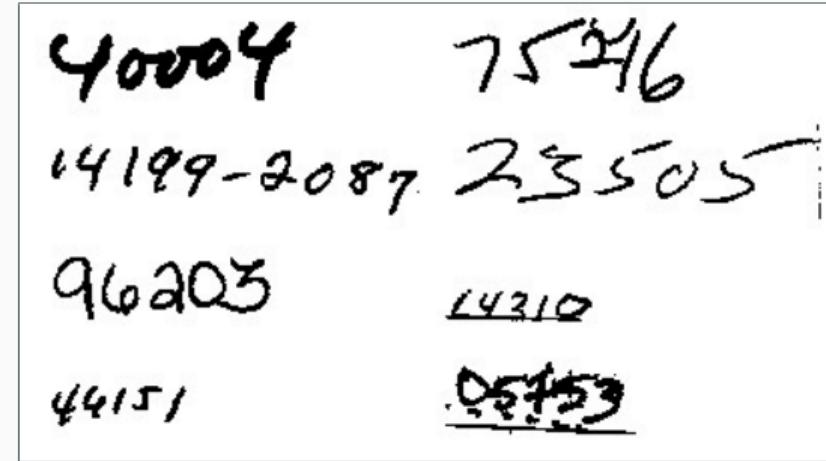
1980



Neocognitron

Convolutional Neural Network

Handwriting digit recognition using backpropagation over a Convolutional Neural Network (LeCun et. al.) [20]



H and W experiment

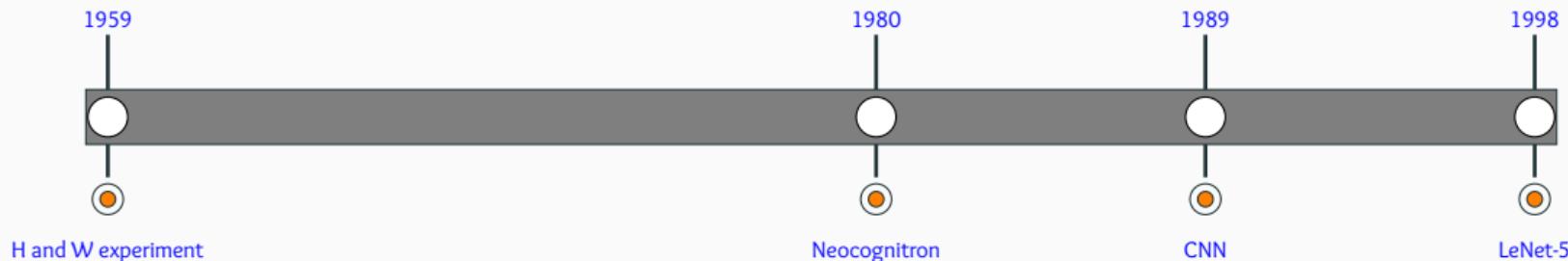
Neocognitron

CNN

LeNet-5

Introduced the (now famous) MNIST dataset
(LeCun et. al.)^[21]

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 3 | 6 | 8 | 1 | 7 | 9 | 6 | 6 | 9 | 1 |
| 6 | 7 | 5 | 7 | 8 | 6 | 3 | 4 | 8 | 5 |
| 2 | 1 | 7 | 9 | 7 | 1 | 2 | 8 | 4 | 5 |
| 4 | 8 | 1 | 9 | 0 | 1 | 8 | 8 | 9 | 4 |
| 7 | 6 | 1 | 8 | 6 | 4 | 1 | 5 | 6 | 0 |
| 7 | 5 | 9 | 2 | 6 | 5 | 8 | 1 | 9 | 7 |
| 2 | 2 | 2 | 2 | 3 | 4 | 4 | 8 | 0 | |
| 0 | 2 | 3 | 8 | 0 | 7 | 3 | 8 | 5 | 7 |
| 0 | 1 | 4 | 6 | 4 | 6 | 0 | 2 | 4 | 3 |
| 7 | 1 | 2 | 8 | 7 | 6 | 9 | 8 | 6 | 1 |



H and W experiment

Neocognitron

CNN

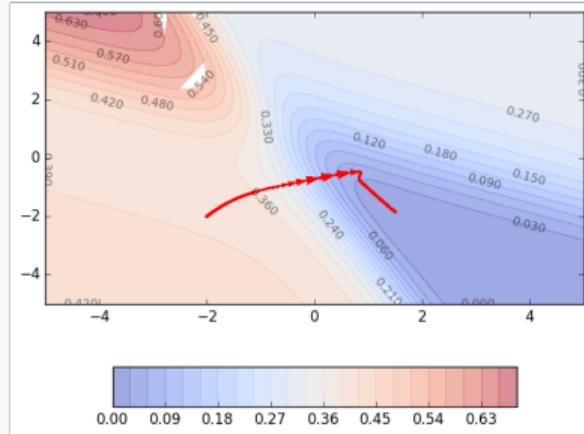
LeNet-5

An algorithm inspired by an experiment on cats is today used to detect cats in videos :-)

Chapter 5: Faster, higher, stronger

Better Optimization Methods

Faster convergence, better accuracies



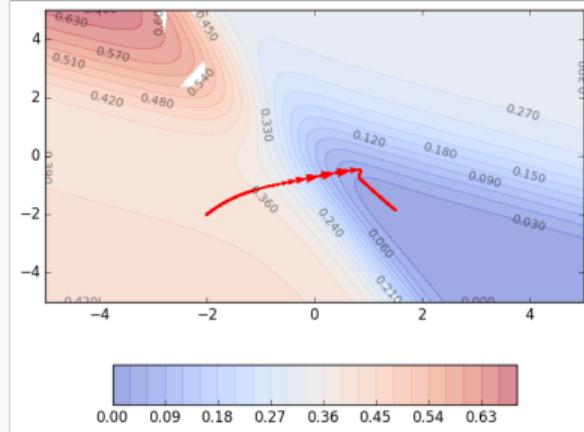
1983



Nesterov

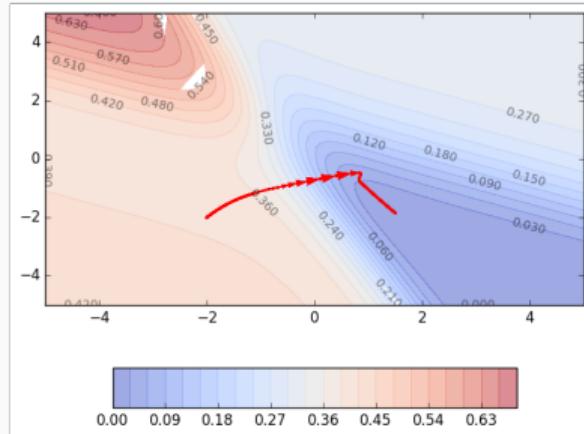
Better Optimization Methods

Faster convergence, better accuracies



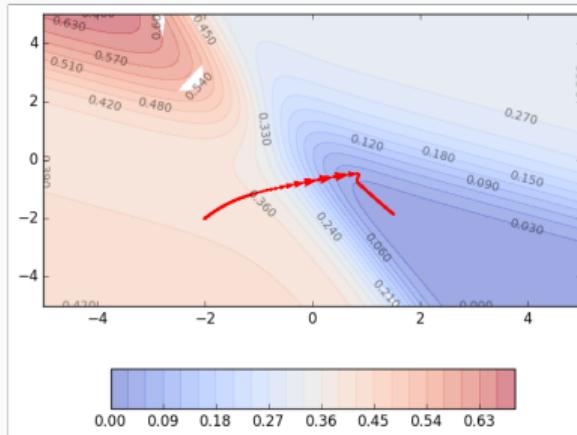
Better Optimization Methods

Faster convergence, better accuracies



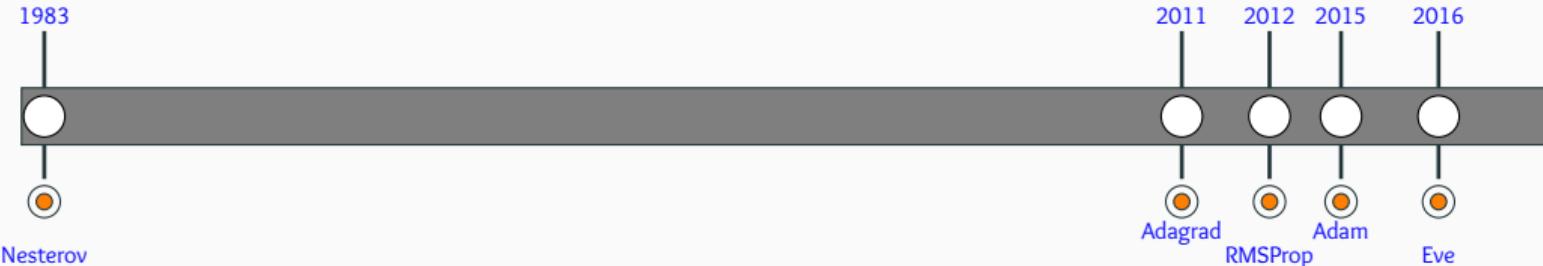
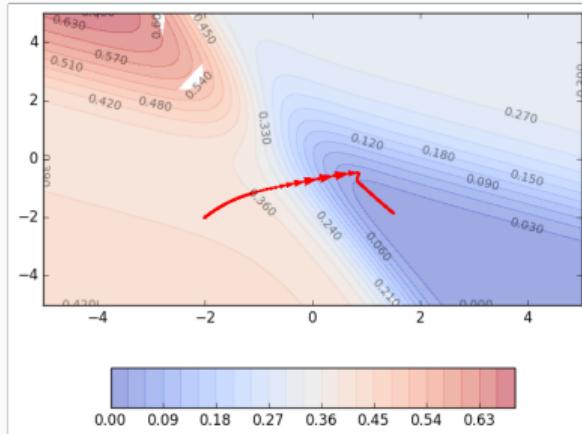
Better Optimization Methods

Faster convergence, better accuracies



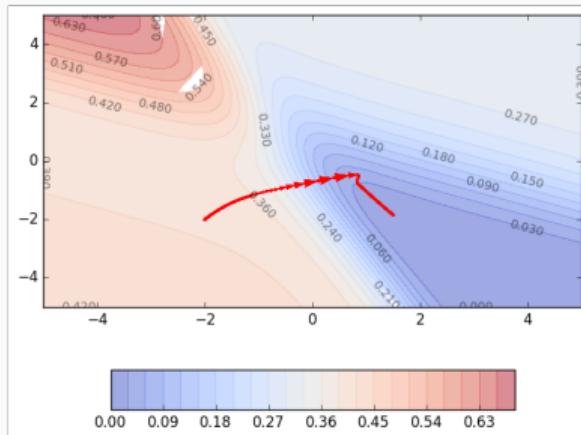
Better Optimization Methods

Faster convergence, better accuracies



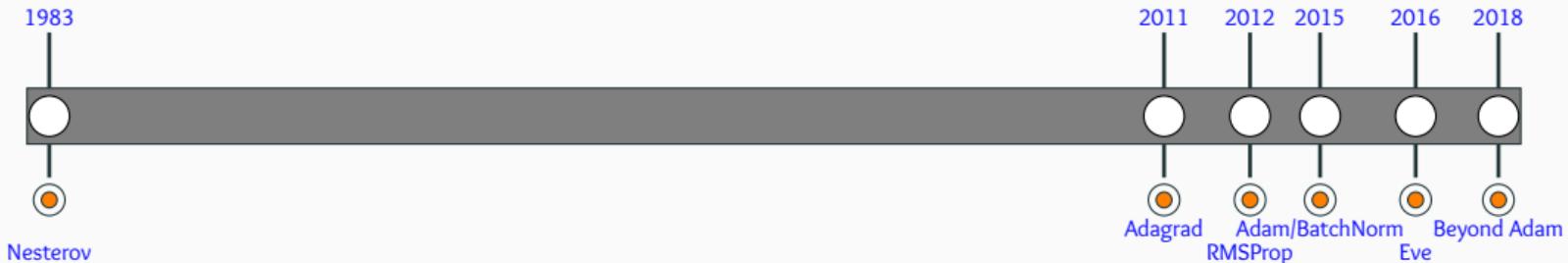
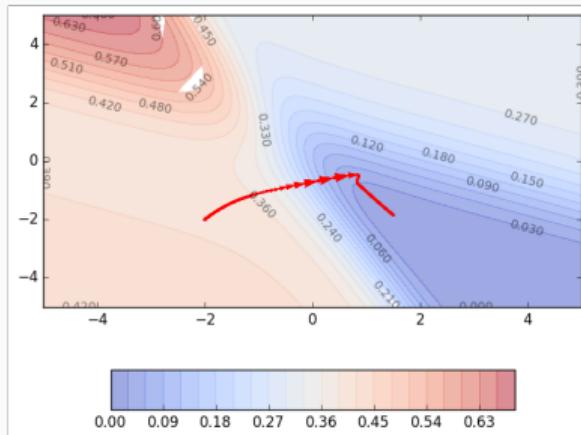
Better Optimization Methods

Faster convergence, better accuracies



Better Optimization Methods

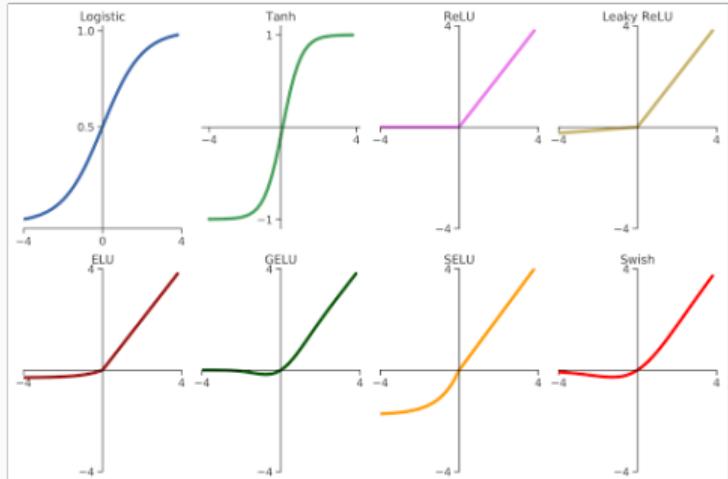
Faster convergence, better accuracies



Better Activation Functions

We have come a long way from the initial days when the logistic function was the default activation function in NNs!

Over the past few years many new functions have been proposed leading to better convergence and/or performance!



Chapter 6: The Curious Case of Sequences

Sequences

They are everywhere

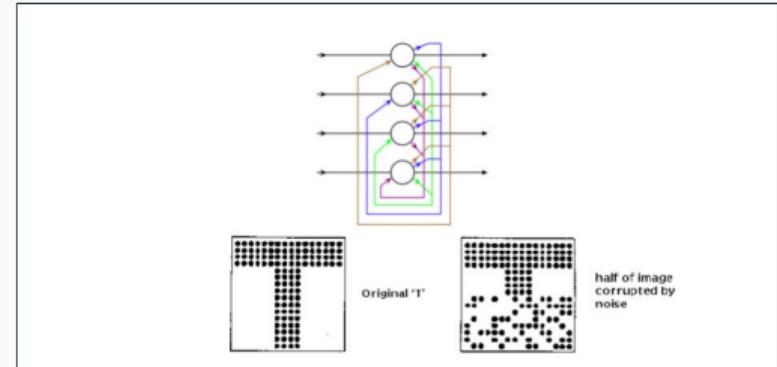
Time series, speech, music, text, video

Each unit in the sequence interacts with
other units

Need models to capture this interaction

Hopfield Network

Content-addressable memory systems for
storing and retrieving patterns^[22]

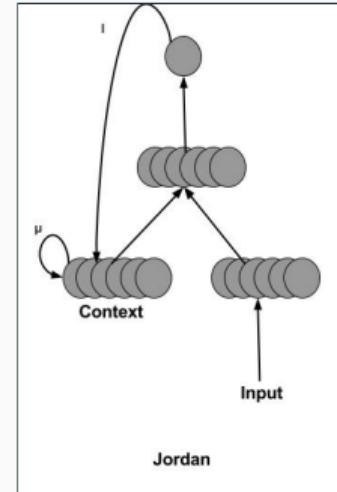


1982



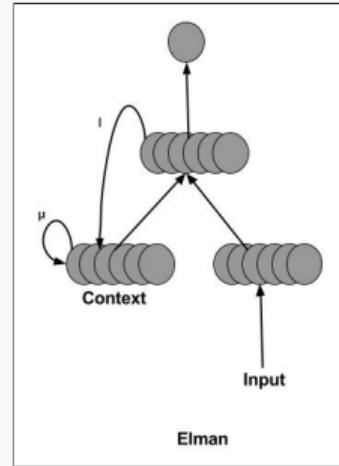
Jordan Network

The output state of each time step is fed to the next time step thereby allowing interactions between time steps in the sequence



Elman Network

The hidden state of each time step is fed to the next time step thereby allowing interactions between time steps in the sequence



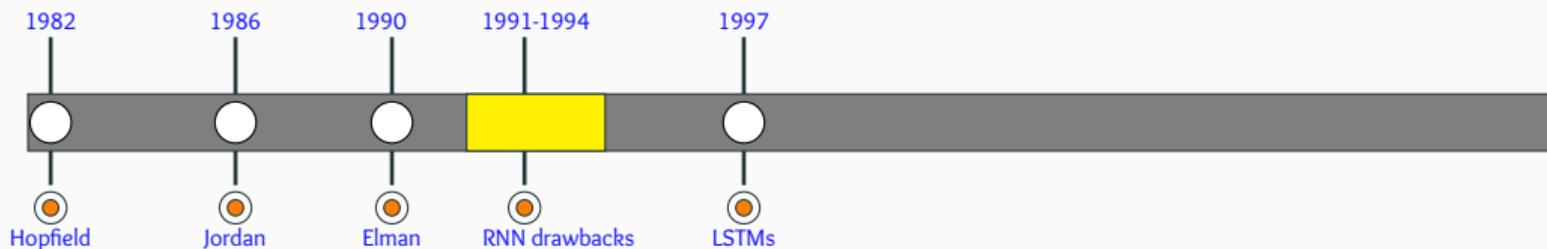
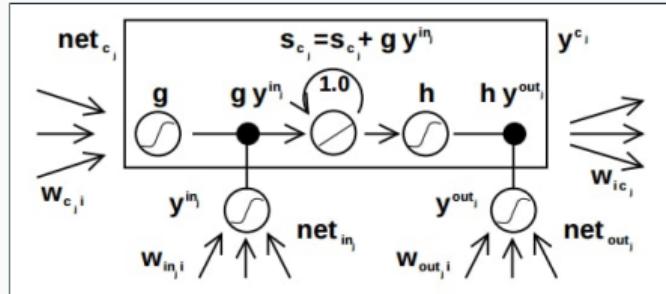
Drawbacks of RNNs

Hochreiter et. al. and Bengio et. al. showed the difficulty in training RNNs (the problem of exploding and vanishing gradients)



Long Short Term Memory

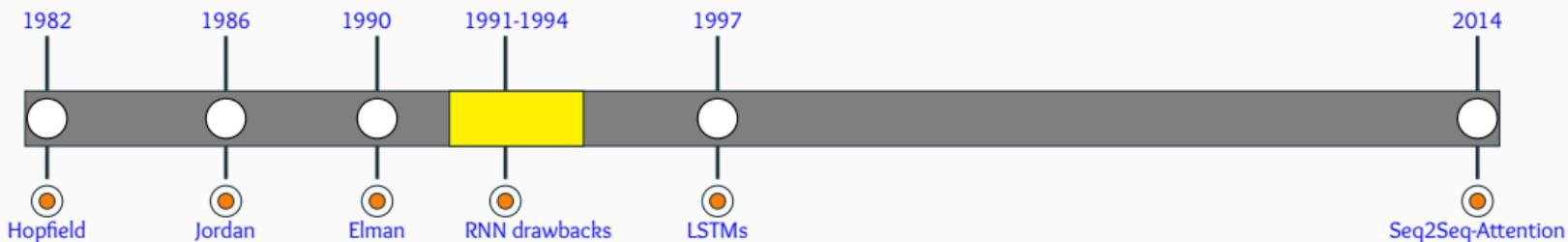
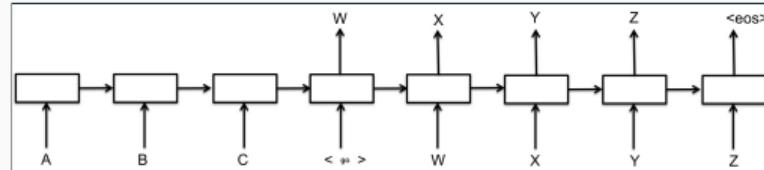
Showed that LSTMs can solve complex long time lag tasks that could never be solved before



Sequence To Sequence Models

Initial success in using RNNs/LSTMs for
large scale Sequence To Sequence
Learning Problems

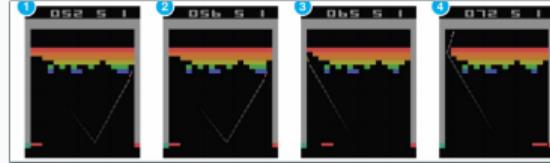
Introduction of Attention which is
perhaps the idea of the decade!



Chapter 7: Beating humans at their own game (literally)

Playing Atari Games

Human-level control through deep reinforcement learning for playing Atari Games [23]



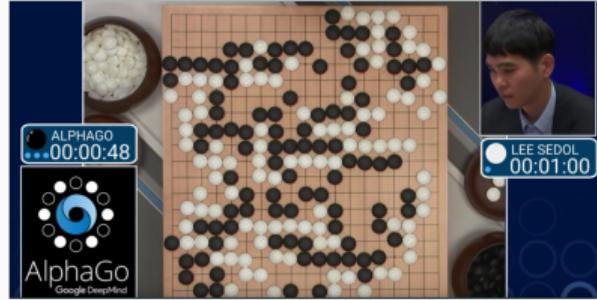
Let's GO

Alpha Go Zero - Best Go player ever,
surpassing human players^[24]

GO is more complex than chess because
of number of possible moves

No brute force backtracking unlike
previous chess agents

2015



Taking a shot at Poker

DeepStack defeated 11 professional poker players with only one outside the margin of statistical significance^[25]



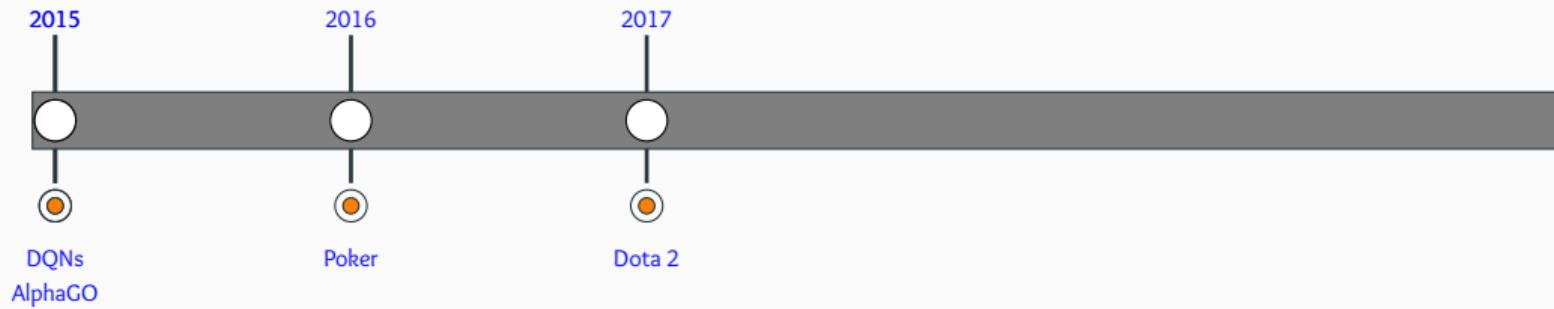
UNIVERSITY OF ALBERTA

48



Defense of the Ancients

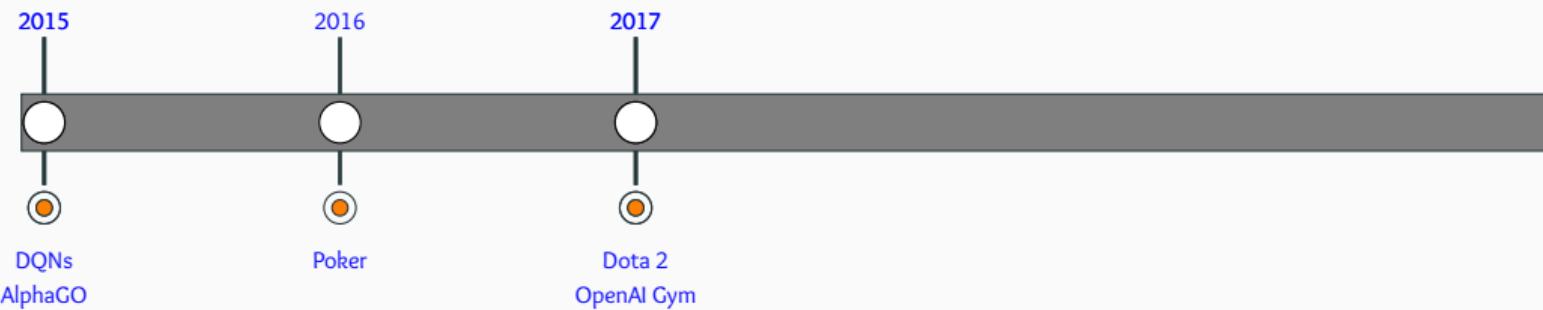
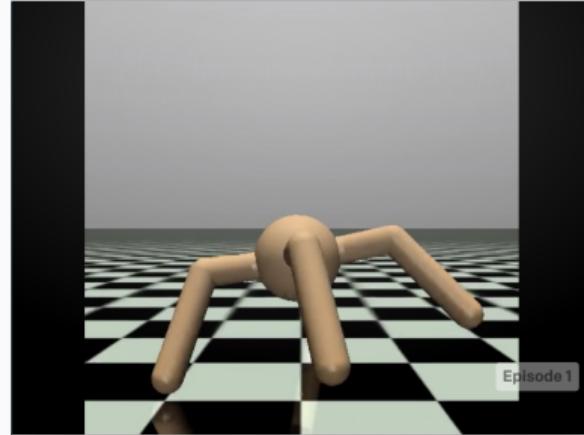
“Our Dota 2 AI, called OpenAI Five, learned by playing over 10,000 years of games against itself. It demonstrated the ability to achieve expert-level performance, learn human–AI cooperation, and operate at internet scale.” – OpenAI



A toolkit for RL

OpenAI Gym^a is a toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Pinball.

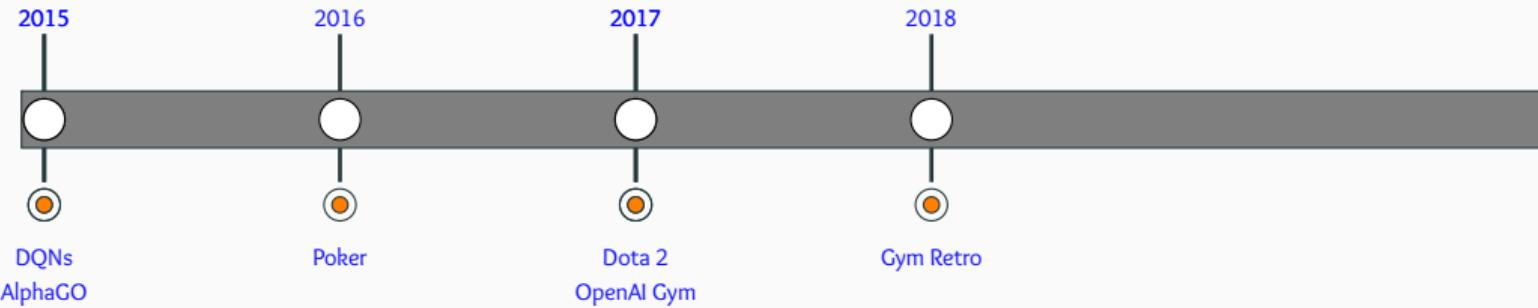
^a<https://gym.openai.com/>



RL for a 1000 games!

Open AI Gym Retro^a: a platform for reinforcement learning research on games which contains 1,000 games across a variety of backing emulators.

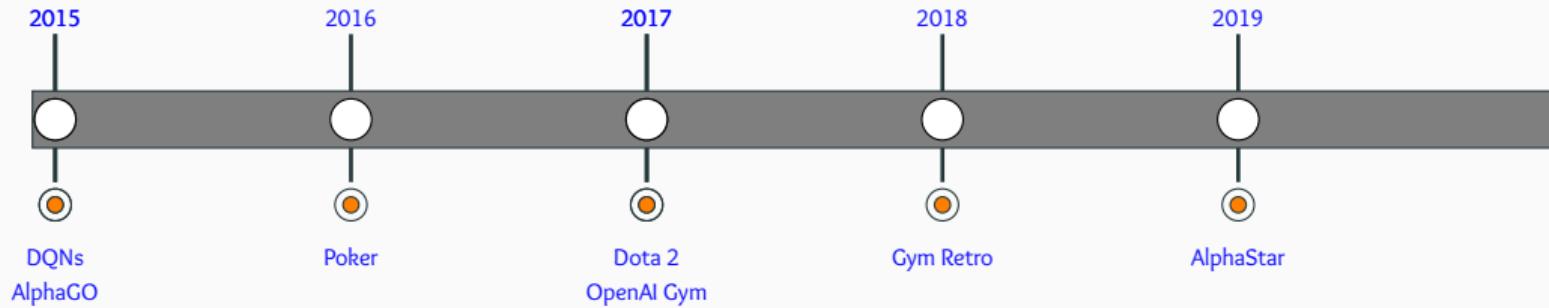
^a<https://openai.com/blog/gym-retro/>



Complex Strategy Games

AlphaStar^a learned to balance short and long-term goals and adapt to unexpected situations while playing using the same maps and conditions as humans

^a<https://deepmind.com/>

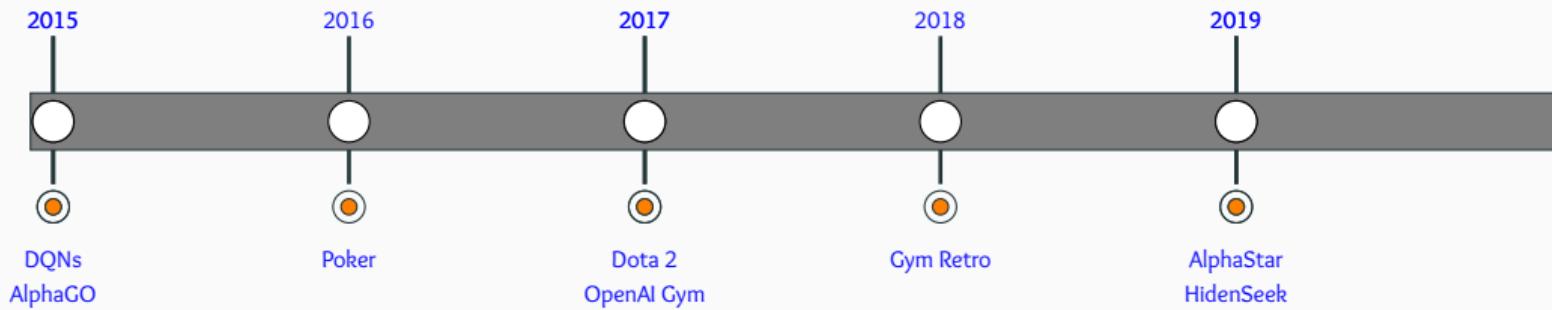


Learning to Hide

OpenAI demonstrated agents which can learn complex strategies such as chase and hide, build a defensive shelter, break a shelter, use a ramp to search inside a shelter and so on!

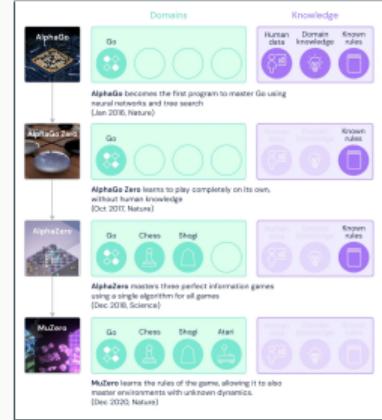


<https://openai.com/blog/emergent-tool-use/>

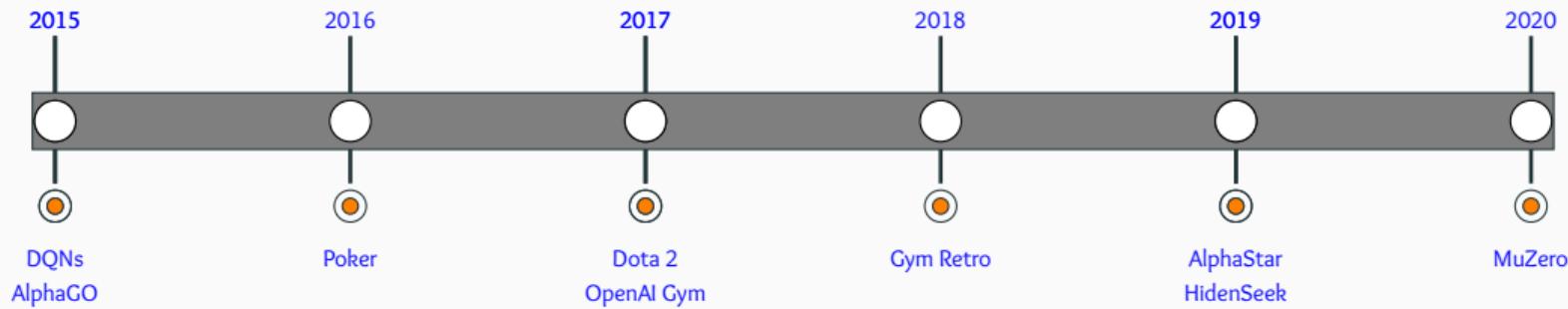


Jack of all, Master of all!

MuZero masters Go, chess, shogi and Atari without needing to be told the rules, thanks to its ability to plan winning strategies in unknown environments.



<https://deepmind.com/blog>



Chapter 8: The Madness (2013-)

He sat on a chair.

Language Modeling

Mikolov et al. (2010)^[26]

Kiros et al. (2015)^[27]

Kim et al. (2015)^[28]



Speech Recognition

Hinton et al. (2012)^[29]

Graves et al. (2013)^[30]

Chorowski et al. (2015)^[31]

Sak et al. (2015)^[32]

MACHINE TRANSLATION



Machine Translation

Kalchbrenner et al. (2013)^[33]

Cho et al. (2014)^[34]

Bahdanau et al. (2015)^[35]

Jean et al. (2015)^[36]

Gulcehre et al. (2015)^[37]

Sutskever et al. (2014)^[38]

Luong et al. (2015)^[39]

Zheng et al. (2017)^[40]

Cheng et al. (2016)^[41]

Chen et al. (2017)^[42]

Firat et al. (2016)^[43]

| Time | User | Utterance |
|----------|-----------|--|
| 03:44 | Old | I dont run graphical ubuntu, I run ubuntu server. |
| 03:45 | kuja | Taru: Haha sucker. |
| 03:45 | Taru | Kuja: ? |
| 03:45 | bur[n]er | Old: you can use "ps ax" and "kill (PID#)" |
| 03:45 | kuja | Taru: Anyways, you made the changes right? |
| 03:45 | Taru | Kuja: Yes. |
| 03:45 | LiveCD | or killall speedlink |
| 03:45 | kuja | Taru: Then from the terminal type: sudo apt-get update |
| 03:46 | _pm | if i install the beta version, how can i update it when the final version comes out? |
| 03:46 | Taru | Kuja: I did. |
| Sender | Recipient | Utterance |
| Old | | I dont run graphical ubuntu, I run ubuntu server. |
| bur[n]er | Old | you can use "ps ax" and "kill (PID#)" |
| kuja | Taru | Haha sucker. |
| Taru | Kuja | ? |
| kuja | Taru | Anyways, you made the changes right? |
| Taru | Kuja | Yes. |
| kuja | Taru | Then from the terminal type: sudo apt-get update |
| Taru | Kuja | I did. |

Conversation Modeling

Shang et al. (2015)^[44]

Vinyals et al. (2015)^[45]

Lowe et al. (2015)^[46]

Dodge et al. (2015)^[47]

Weston et al. (2016)^[48]

Serban et al. (2016)^[49]

Bordes et al. (2017)^[50]

Serban et al. (2017)^[51]

Task 1: Single Supporting Fact

Mary went to the bathroom.
John moved to the hallway.
Mary travelled to the office.
Where is Mary? A:office

Task 2: Two Supporting Facts

John is in the playground.
John picked up the football.
Bob went to the kitchen.
Where is the football? A:playground

Task 3: Three Supporting Facts

John picked up the apple.
John went to the office.
John went to the kitchen.
John dropped the apple.
Where was the apple before the kitchen? A:office

Task 4: Two Argument Relations

The office is north of the bedroom.
The bedroom is north of the bathroom.
The kitchen is west of the garden.
What is north of the bedroom? A: office
What is the bedroom north of? A: bathroom

Question Answering

Hermann et al. (2015)^[52]

Chen et al. (2016)^[53]

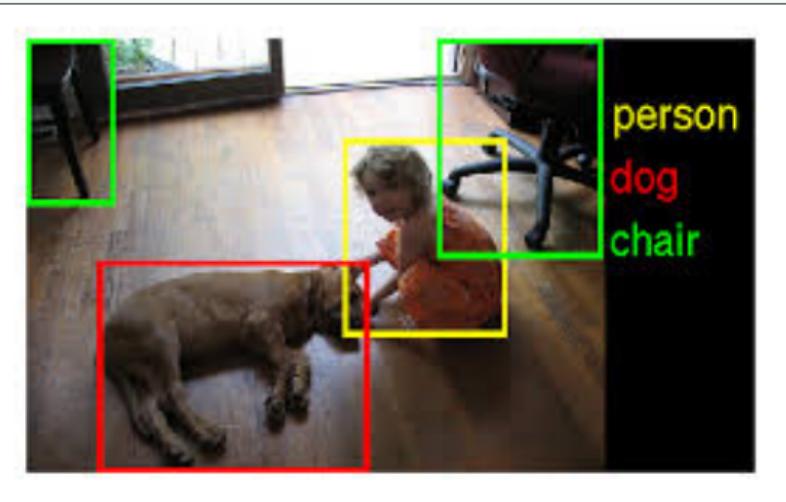
Xiong et al. (2016)^[54]

Seo et al. (2016)^[55]

Dhingra et al. (2017)^[56]

Wang et al. (2017)^[57]

Hu et al. (2017)^[58]



Object Detection/Recognition

Semantic Segmentation (Long et al, 2015) [59]

Recurrent CNNs (Liang et al., 2015) [60]

Faster RCNN (Ren et al., 2015) [61]

Inside-Outside Net (Bell et al., 2015) [62]

YOLO9000 (Redmon et al., 2016) [63]

R-FCN (Dai et al., 2016) [64]

Mask R-CNN (He at al., 2017) [65]

Video Object segmentation (Caelles et al., 2017) [66]



Visual Tracking

Choi et al. (2017)^[67]

Yun et al. (2017)^[68]

Alahi et al. (2017)^[69]

Retr.
Gen.



1. Top view of the lights of a city at night, with a well-illuminated square in front of a church in the foreground;
2. People on the stairs in front of an illuminated cathedral with two towers at night;

A square with burning street lamps and a street in the foreground;



1. Tourists are sitting at a long table with beer bottles on it in a rather dark restaurant and are raising their bierglaeser;
2. Tourists are sitting at a long table with a white table-cloth in a somewhat dark restaurant;

Tourists are sitting at a long table with a white table cloth and are eating;

Image Captioning

Mao et al. (2014)^[70]

Mao et al. (2015)^[71]

Kiros et al. (2015)^[72]

Donahue et al. (2015)^[73]

Vinyals et al. (2015)^[74]

Karpathy et al. (2015)^[75]

Fang et al. (2015)^[76]

Chen et al. (2015)^[77]



A group of young men playing a game of soccer



A man riding a wave on top of a surfboard.

Video Captioning

Donahue et al. (2014)^[78]

Venugopalan et al. (2014)^[79]

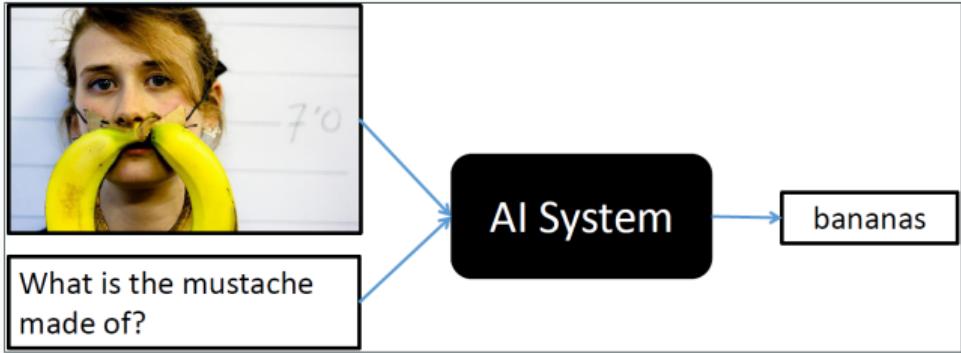
Pan et al. (2015)^[80]

Yao et al. (2015)^[81]

Rohrbach et al. (2015)^[82]

Zhu et al. (2015)^[83]

Cho et al. (2015)^[34]



Visual Question Answering

Santoro et al. (2017)^[84]

Hu et al. (2017)^[85]

Johnson et al. (2017)^[86]

Ben-younes et al. (2017)^[87]

Malinowski et al. (2017)^[88]

Kazemi et al. (2016)^[89]

She _____.



She opens the _____.



Question: What is the cat doing? Answer: playing with a tablet

Video Question Answering

Tapaswi et. al. 2016^[90]

Zeng et. al. 2016^[91]

Maharaj et. al. 2017^[92]

Zhao et. al. 2017^[93]

Yu Youngjae et. al. 2017^[94]

Xue Hongyang et. al. 2017^[95]

Mazaheri et. al. 2017^[96]



Video Summarization

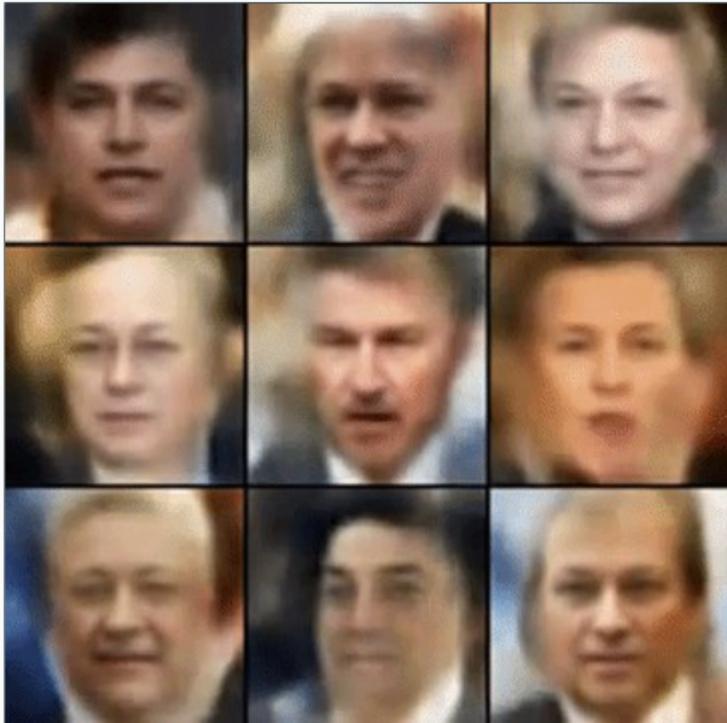
Chheng 2007^[97]

Ajmal 2012^[98]

Zhang Ke 2016^[99]

Zhong Ji 2017^[100]

Panda 2017^[101]



Generating Authentic Photos

Variational Autoencoders

(Kingma et. al., 2013) [102]

Generative Adversarial

Networks (Goodfellow et. al.,
2014) [103]

Plug & Play generative nets

(Nguyen et al., 2016) [104]

Progressive Growing of GANs

(Karras et al., 2017) [105]



Generating Raw Audio

Wavenets (Oord et. al.,
2016)^[106]



Pixel RNNs

(Oord et al., 2016)^[107]

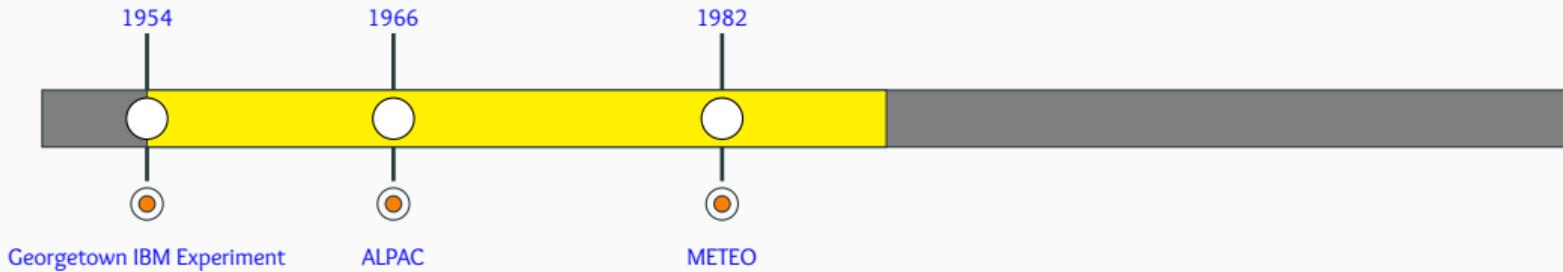
(Oord et al., 2016)^[108]

(Salimans et al., 2017)^[109]

Chapter 9: The Rise of the Transformers

Rule Based Systems

Initial Machine Translation Systems used hand crafted rules and dictionaries to translate sentences between few politically important language pairs (e.g., English -Russian). They could not live up to the initial hype and were panned by the ALPAC report (1966)



Statistical MT

The IBM Models for Machine Translation gave a boost to the idea of data driven statistical NLP which then ruled NLP for the next 2 decades till Deep Learning took over!

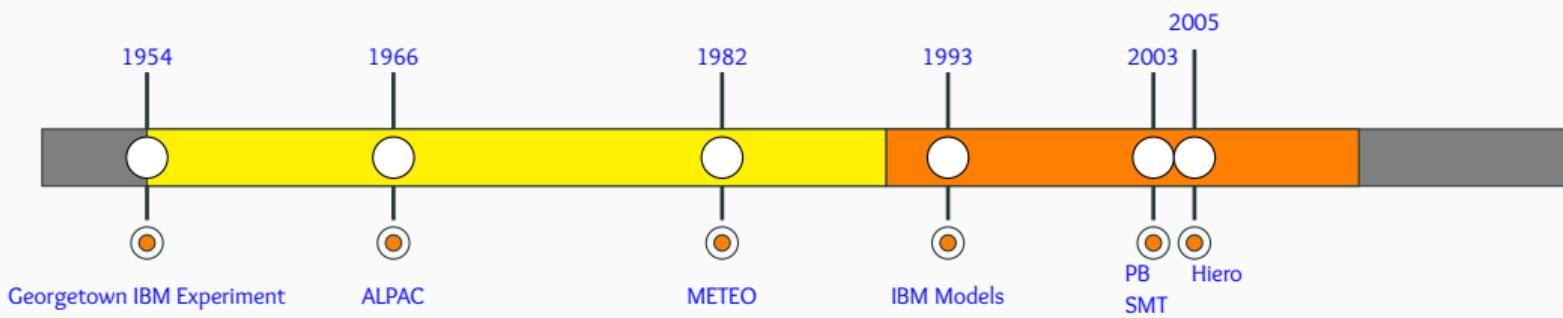
The Mathematics of Statistical Machine Translation: Parameter Estimation

Peter F. Brown*
IBM T.J. Watson Research Center

Stephen A. Della Pietra*
IBM T.J. Watson Research Center

Vincent J. Della Pietra*
IBM T.J. Watson Research Center

Robert L. Mercer*
IBM T.J. Watson Research Center



Neural MT

The introduction of seq2seq models and attention^[35] (perhaps, the idea of the decade!) lead to a paradigm shift in NLP ushering the era of bigger, hungrier (more data), better models!

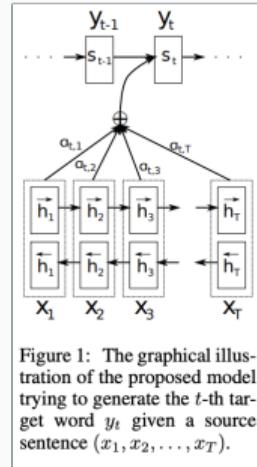
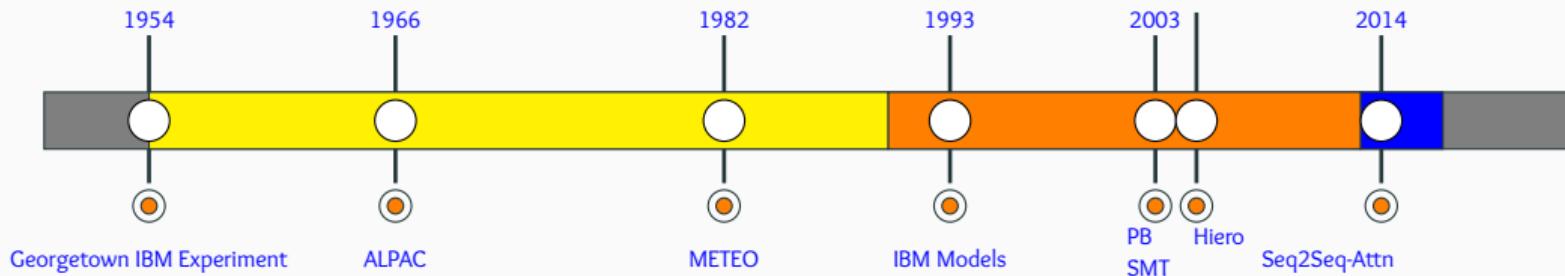


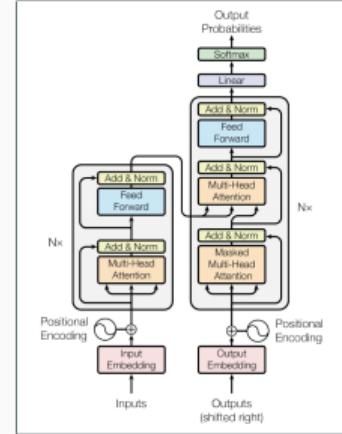
Figure 1: The graphical illustration of the proposed model trying to generate the t -th target word y_t given a source sentence (x_1, x_2, \dots, x_T) .

Source: Bahdanau et. al. [35]
2005

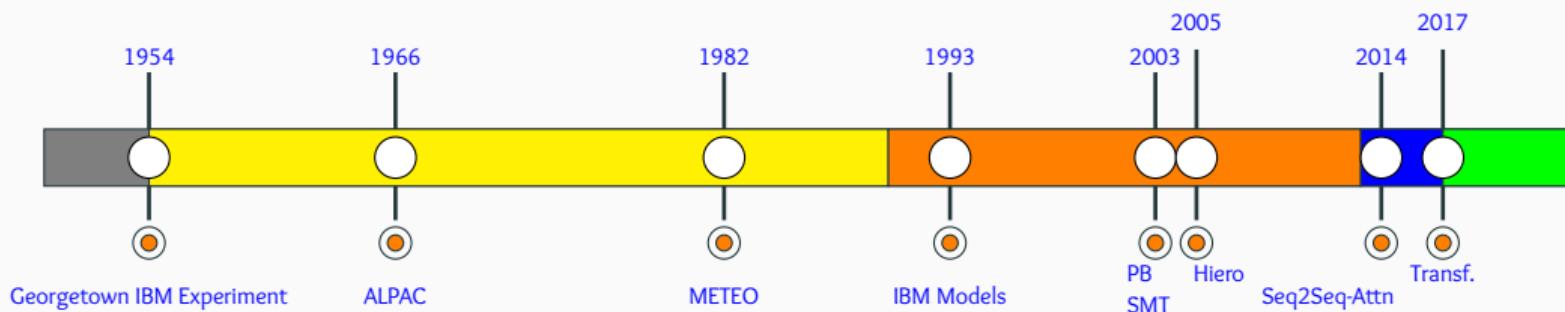


The Transformer Revolution

It is rare for a field to see two dramatic paradigm shifts in a short span of 4 years!
Since their inception transformers have taken the NLP world by storm leading to the development of insanely big models trained on obscene amounts of data!

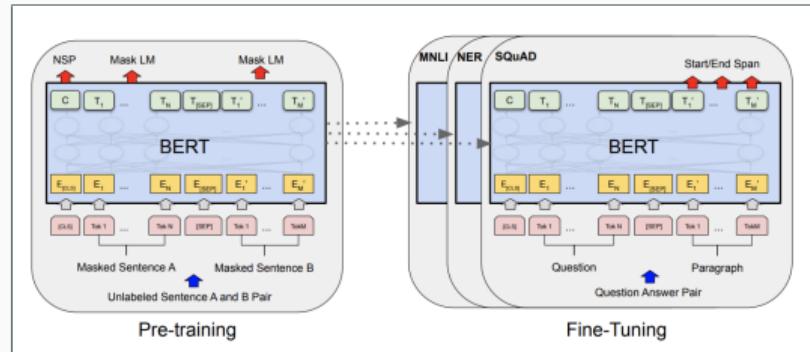


Source: Vaswani et. al. [110]

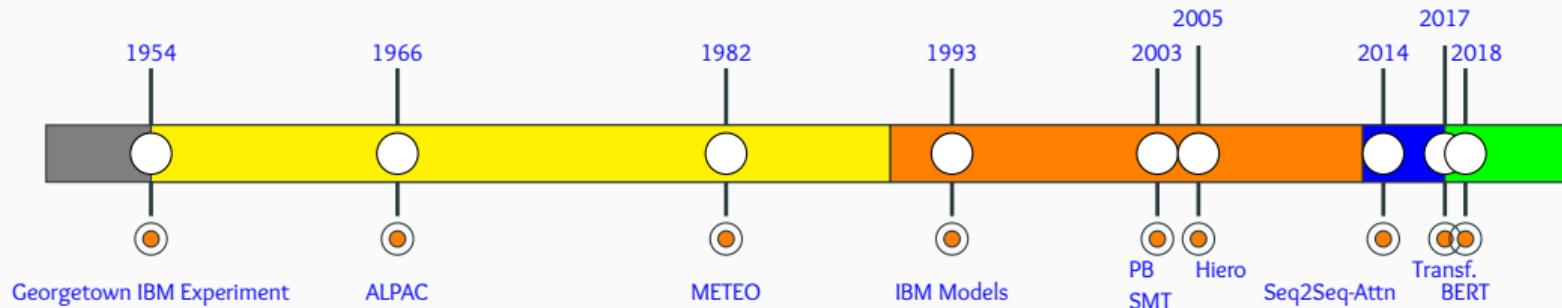


The Transformer Revolution

Most NLP applications today are driven by BERT and its variants. The key idea here was to learn general language characteristics using large amounts of unlabeled corpora and then fine-tune the model for specific downstream tasks.



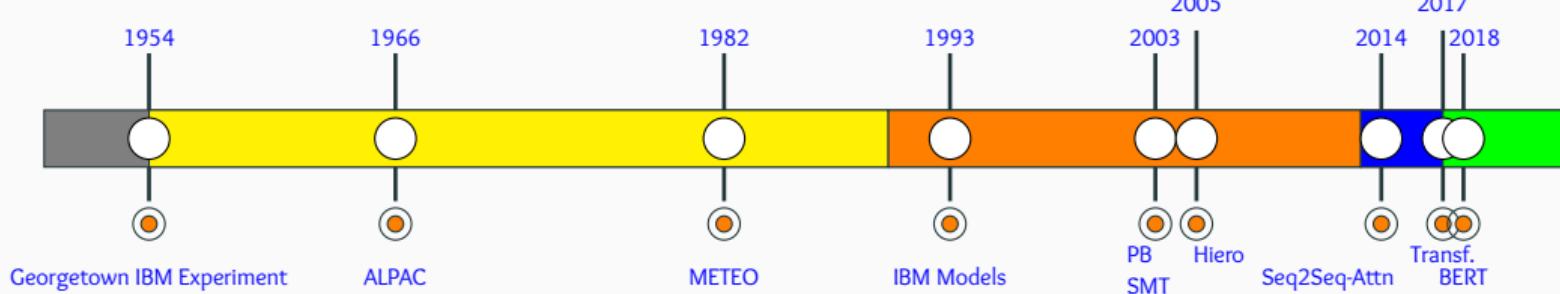
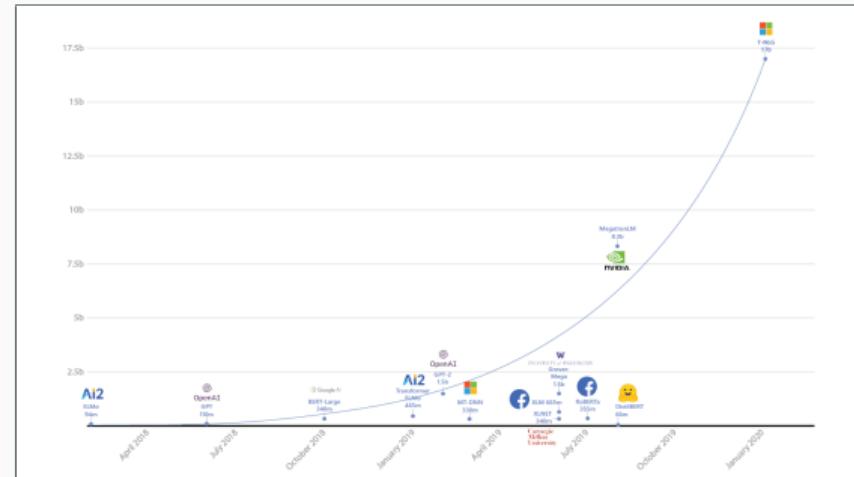
Source: Devlin et. al. [111]



The Billion Parameter Club

The models are becoming bigger and bigger and bigger!

Source: <https://msturing.org/>



The Trillion Parameter Club

Trained on 100 languages, with a total of 13B examples, 1 Trillion Parameters on 2048 TPUs!

This is insane!

GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding

Dmitry Lepikhin
lepiikhin@google.com

HyoukJoong Lee
hyouklee@google.com

Yuanzhong Xu
yuanzx@google.com

Dehao Chen
dehao@google.com

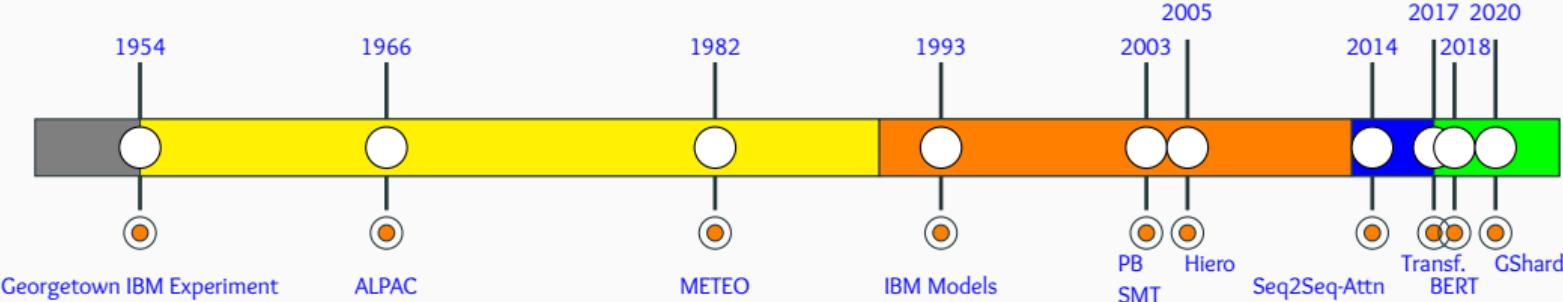
Orhan Firat
orhanf@google.com

Yanping Huang
huangyp@google.com

Maxim Krikun
krikun@google.com

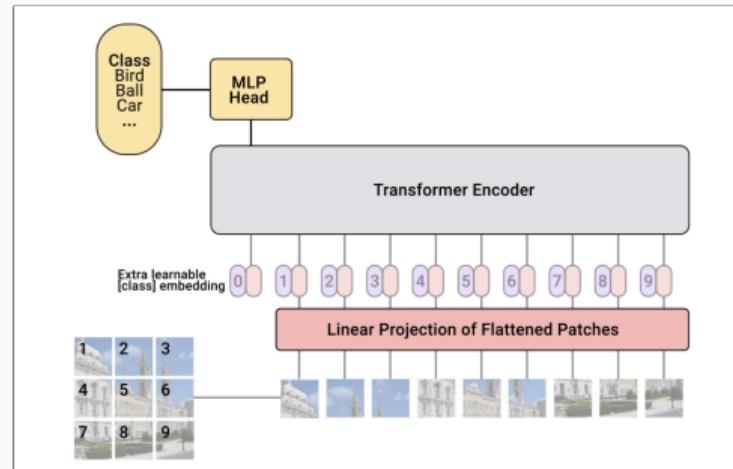
Noam Shazeer
noam@google.com

Zhifeng Chen
zhifengc@google.com

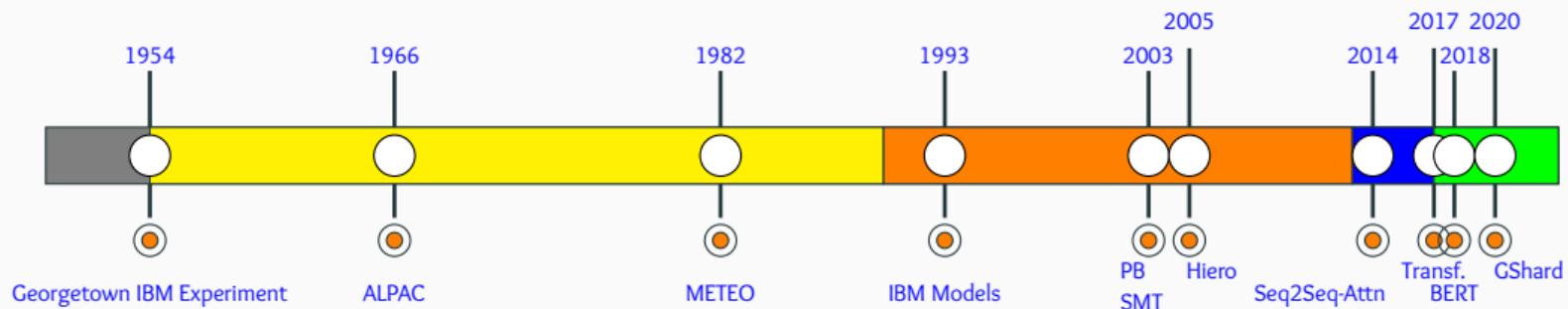


From Language To Vision

A vision model^a based as closely as possible on the Transformer architecture originally designed for text-based tasks (another paradigm shift from CNNs which have been around since 1980s!)



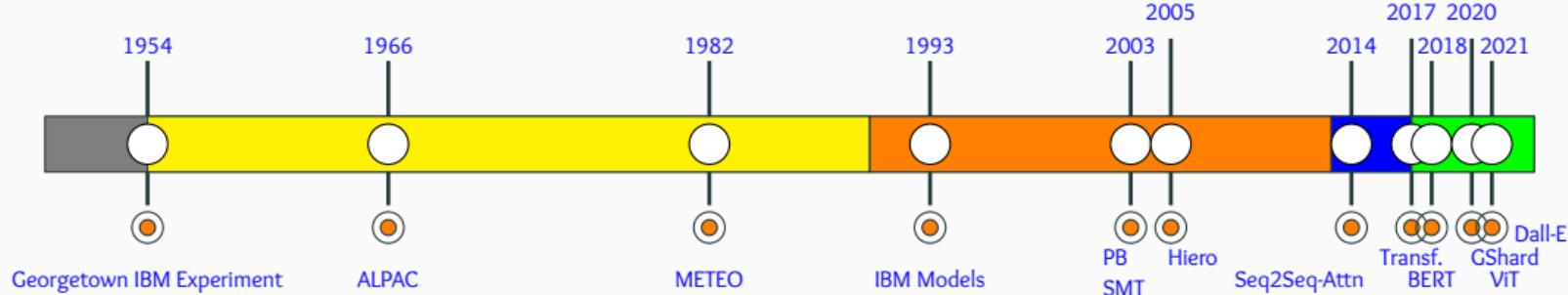
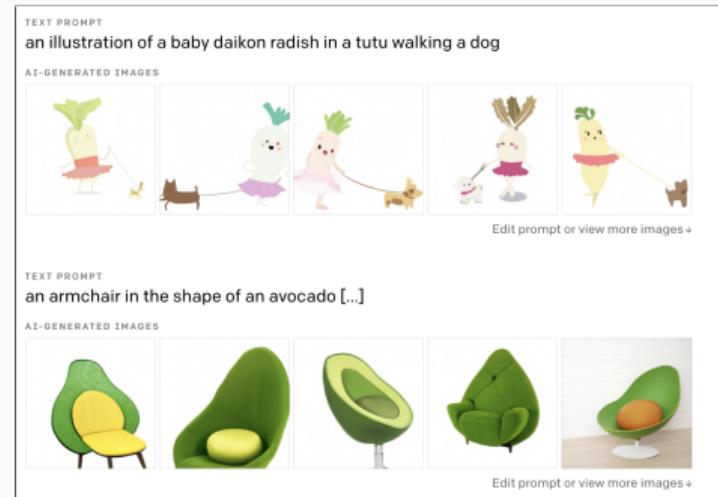
^aSource:<https://ai.googleblog.com/2020/12/transformers-for-image-recognition-at.html>



From Language To Vision

DALL·E^a is a 12-billion parameter version of GPT-3 trained to generate images from text descriptions, using a dataset of text–image pairs.

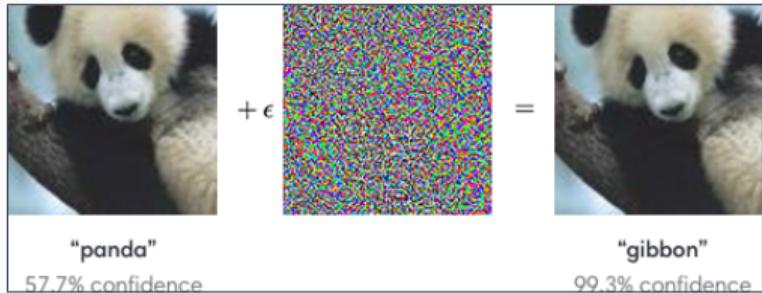
^a<https://openai.com/blog/dall-e/>



Chapter 10: Calls for Sanity (Interpretable, Fair, Responsible, Green AI)

The Paradox of Deep Learning

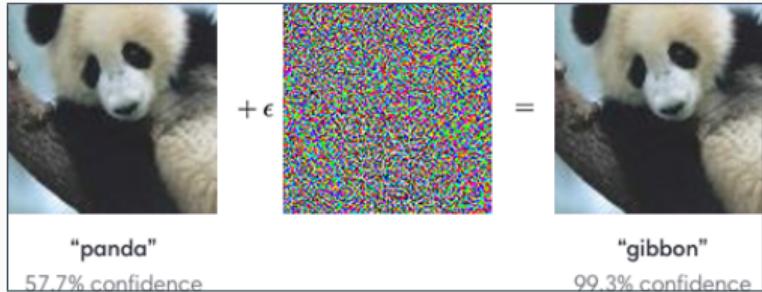
Why does deep learning work so well despite



*<https://arxiv.org/pdf/1710.05468.pdf>

The Paradox of Deep Learning

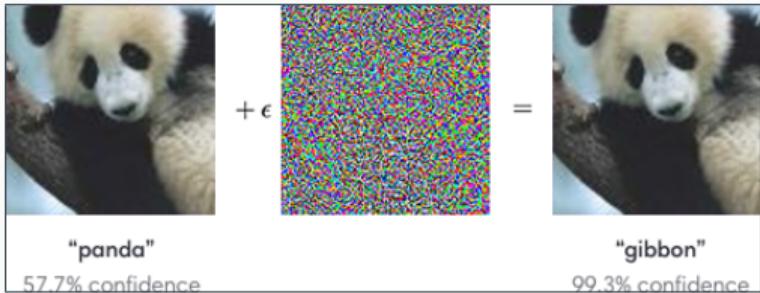
Why does deep learning work so well despite
high capacity (susceptible to overfitting)



*<https://arxiv.org/pdf/1710.05468.pdf>

The Paradox of Deep Learning

Why does deep learning work so well despite
high capacity (susceptible to overfitting)
numerical instability (vanishing/exploding gradients)

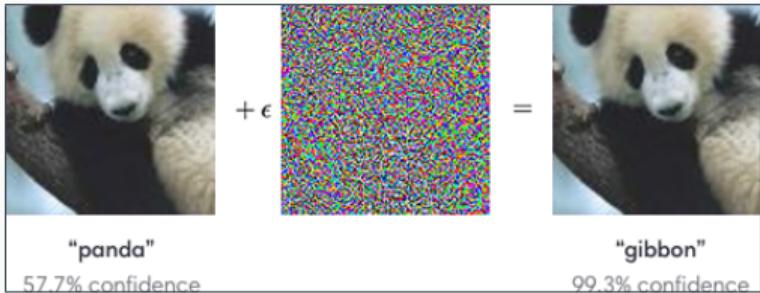


*<https://arxiv.org/pdf/1710.05468.pdf>

The Paradox of Deep Learning

Why does deep learning work so well despite

- high capacity (susceptible to overfitting)
- numerical instability (vanishing/exploding gradients)
- sharp minima (leading to overfitting)

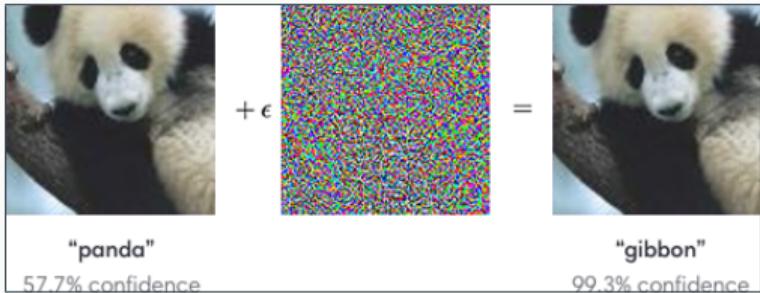


*<https://arxiv.org/pdf/1710.05468.pdf>

The Paradox of Deep Learning

Why does deep learning work so well despite

- high capacity (susceptible to overfitting)
- numerical instability (vanishing/exploding gradients)
- sharp minima (leading to overfitting)
- non-robustness (see figure)

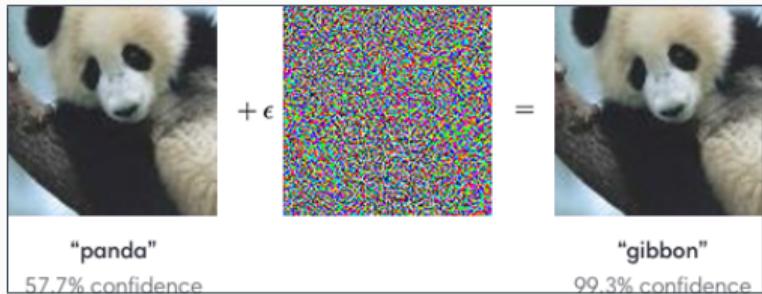


*<https://arxiv.org/pdf/1710.05468.pdf>

The Paradox of Deep Learning

Why does deep learning work so well despite

- high capacity (susceptible to overfitting)
- numerical instability (vanishing/exploding gradients)
- sharp minima (leading to overfitting)
- non-robustness (see figure)



No clear answers yet but ...

*<https://arxiv.org/pdf/1710.05468.pdf>

The Paradox of Deep Learning

Why does deep learning work so well despite

- high capacity (susceptible to overfitting)

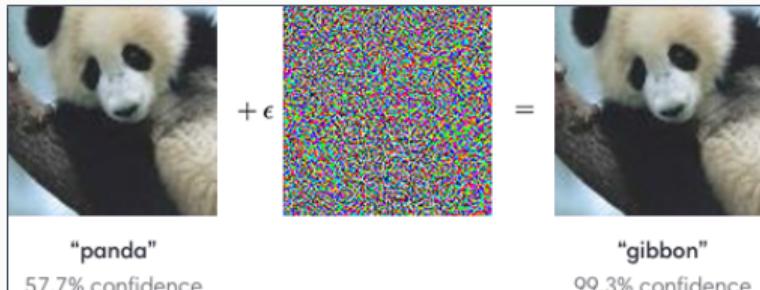
- numerical instability (vanishing/exploding gradients)

- sharp minima (leading to overfitting)

- non-robustness (see figure)

No clear answers yet but ...

Slowly but steadily there is increasing emphasis on
explainability and theoretical justifications!*



*<https://arxiv.org/pdf/1710.05468.pdf>

The Paradox of Deep Learning

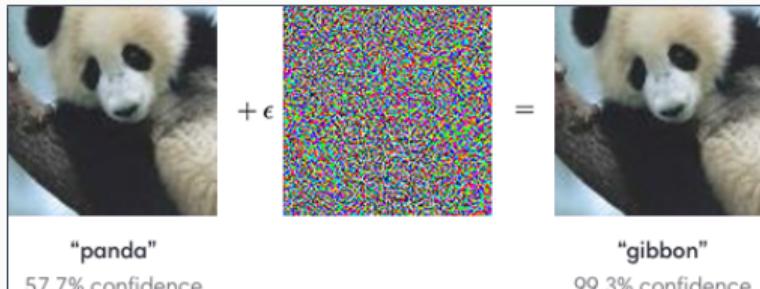
Why does deep learning work so well despite

- high capacity (susceptible to overfitting)

- numerical instability (vanishing/exploding gradients)

- sharp minima (leading to overfitting)

- non-robustness (see figure)



No clear answers yet but ...

Slowly but steadily there is increasing emphasis on
explainability and theoretical justifications!*

Hopefully this will bring sanity to the proceedings !

*<https://arxiv.org/pdf/1710.05468.pdf>

Tell me why!

Workshop on Human Interpretability in
Machine Learning

*We still do not know much about why DL models
do what they do!*

2016



WHI

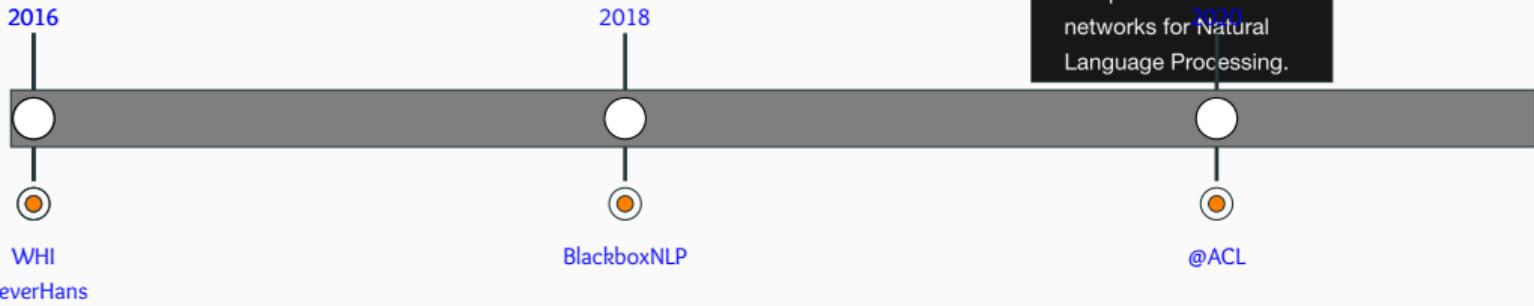


Tell me why!

Push for analyzing and interpreting neural networks for NLP

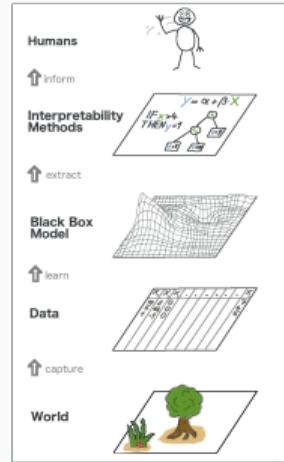
Analyzing
and
interpreting
neural
networks
for NLP

Revealing the content
of the neural black box:
workshop on the
analysis and
interpretation of neural
networks for Natural
Language Processing.

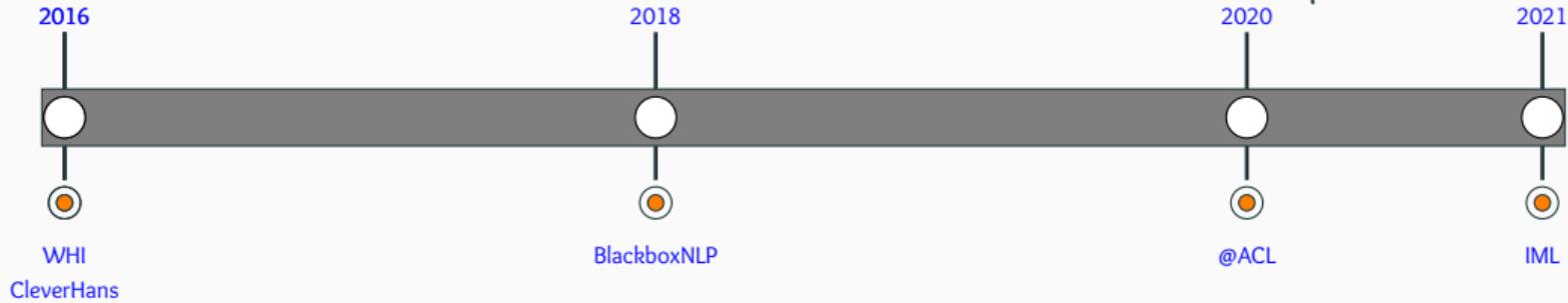


Tell me why!

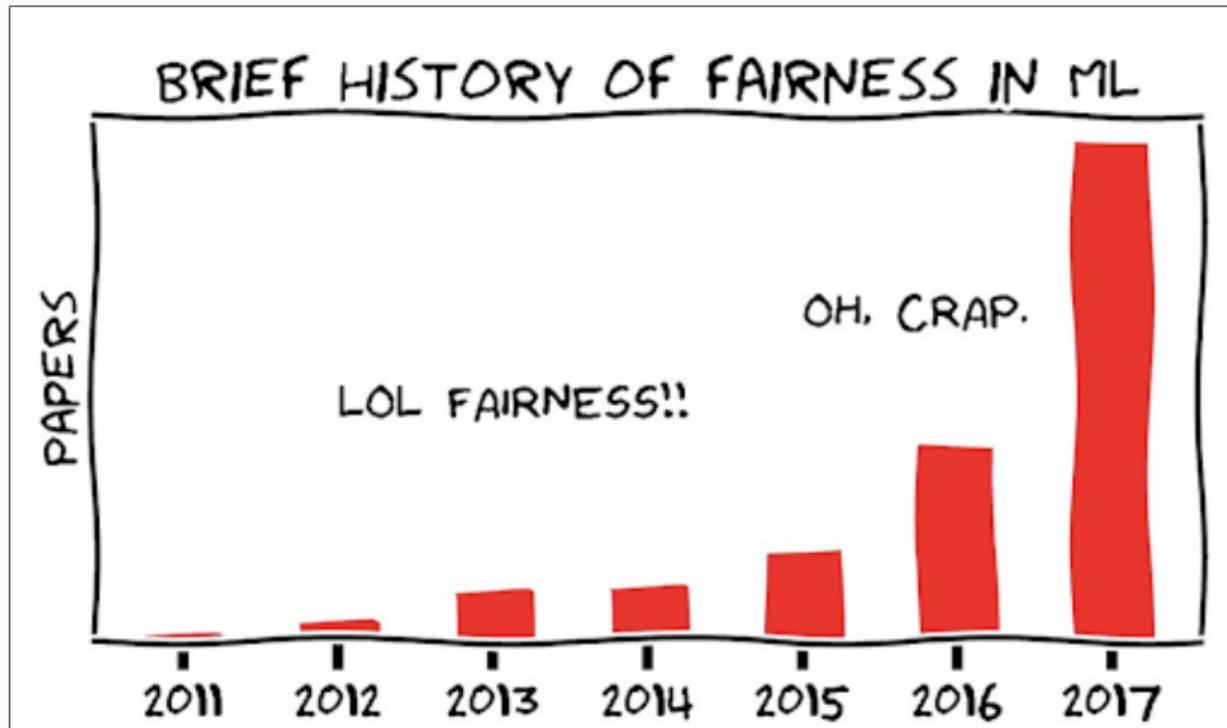
Interpretable Machine Learning: A Guide for Making Black Box Models Explainable. –
Christoph Molnar



Source: IML: Christoph Molnar



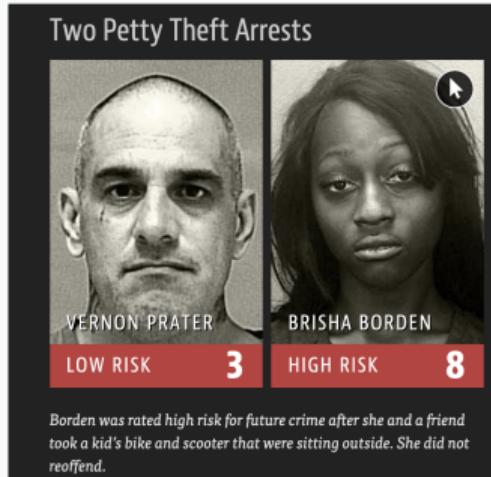
Be Fair and Responsible!



Source: <https://fairmlclass.github.io/> (Moritz Hardt)

Be Fair and Responsible!

“There’s software used across the country to predict future criminals. And it’s biased against blacks.” - Propublica



2016



Machine Bias

Source:

<https://www.propublica.org/article/machine-bias>

Be Fair and Responsible!

“Facial Recognition Is Accurate, if You’re a White Guy” - MIT Media

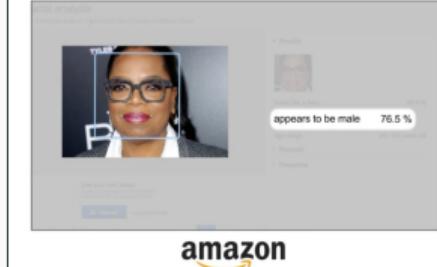
Gender Shades audit, 2018

Accuracy in gender classification

| | Darker Male | Darker Female | Lighter Male | Lighter Female | Largest Gap |
|-----------|-------------|---------------|--------------|----------------|-------------|
| IBM | 88.0% | 65.3% | 99.7% | 92.9% | 34.4% |
| Megvii | 99.3% | 65.5% | 99.2% | 94.0% | 33.8% |
| Microsoft | 94.0% | 79.2% | 100.0% | 98.3% | 20.8% |

Chart: MIT Technology Review • Source: Joy Buolamwini & Timnit Gebru • Created with Datawrapper

Oprah Winfrey

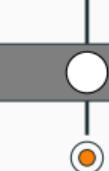


2016



Machine Bias

2018

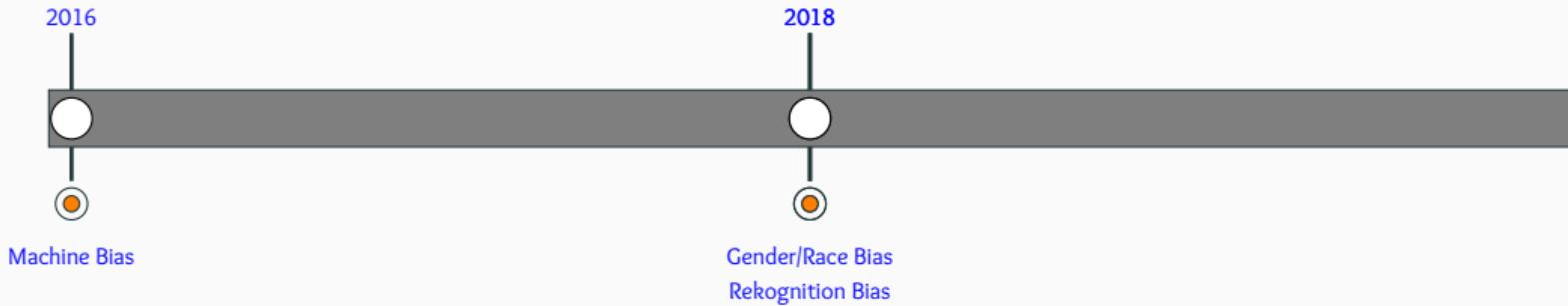
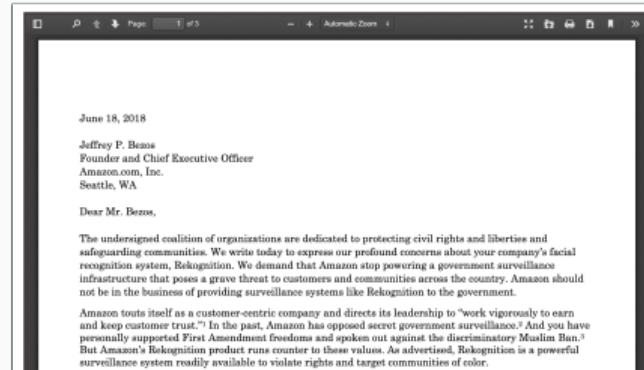


Gender/Race Bias

Source: Joy Buolamwini (Youtube)

Be Fair and Responsible!

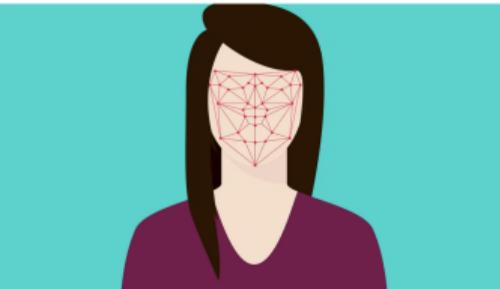
In 2018, nearly 70 civil rights and research organizations wrote a letter to Jeff Bezos demanding that Amazon stop providing face recognition technology to governments.



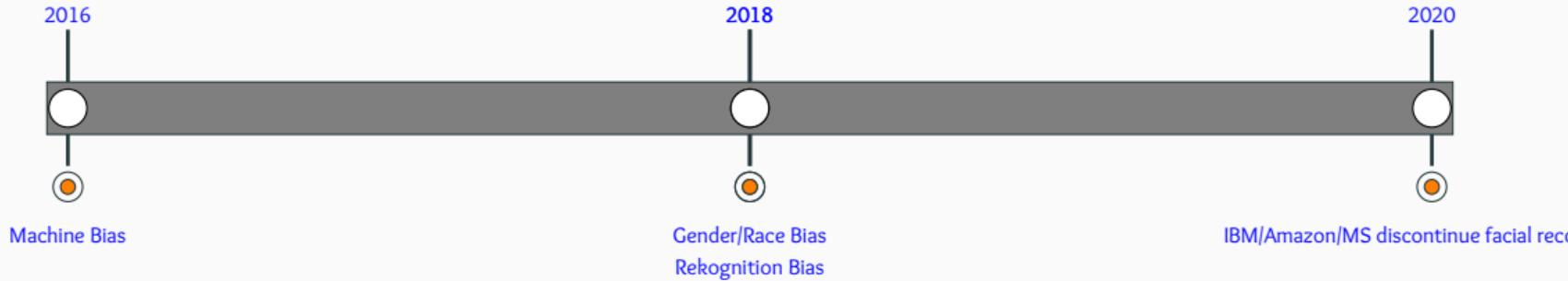
Be Fair and Responsible!

Microsoft refuses to sell police its facial-recognition technology, following similar moves by Amazon and IBM

IBM says it is no longer working on face recognition because it's used for racial profiling

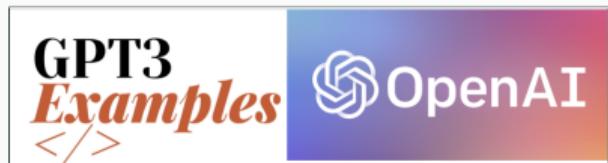
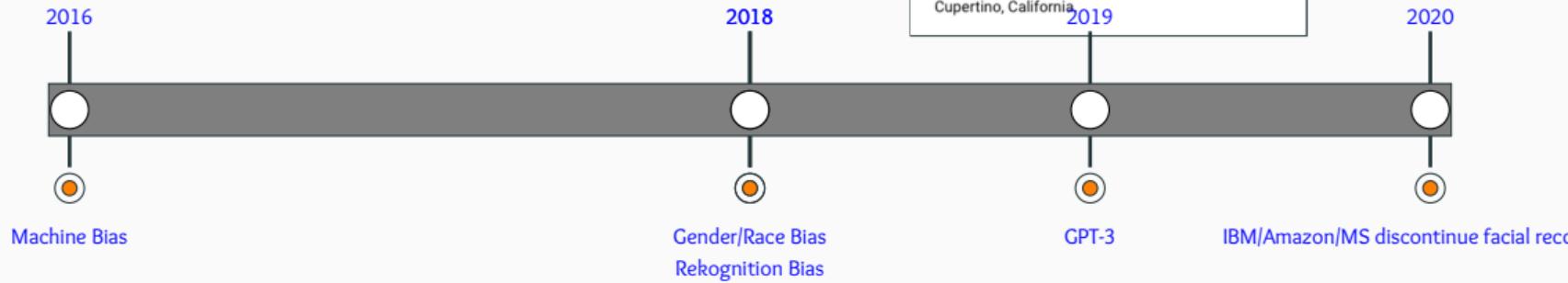


Technology
Microsoft won't sell police its facial-recognition technology, following similar moves by Amazon and IBM



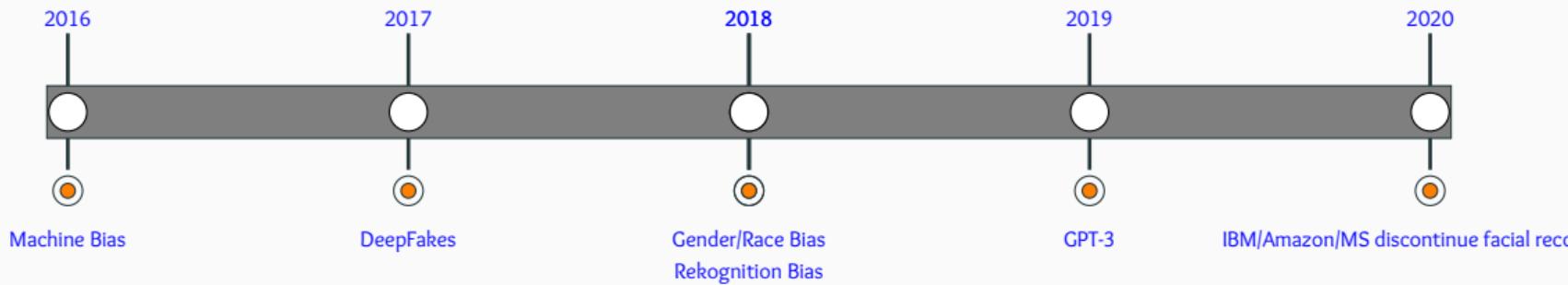
Be Fair and Responsible!

“Due to our concerns about malicious applications of the technology, we are not releasing the trained model.” — OpenAI



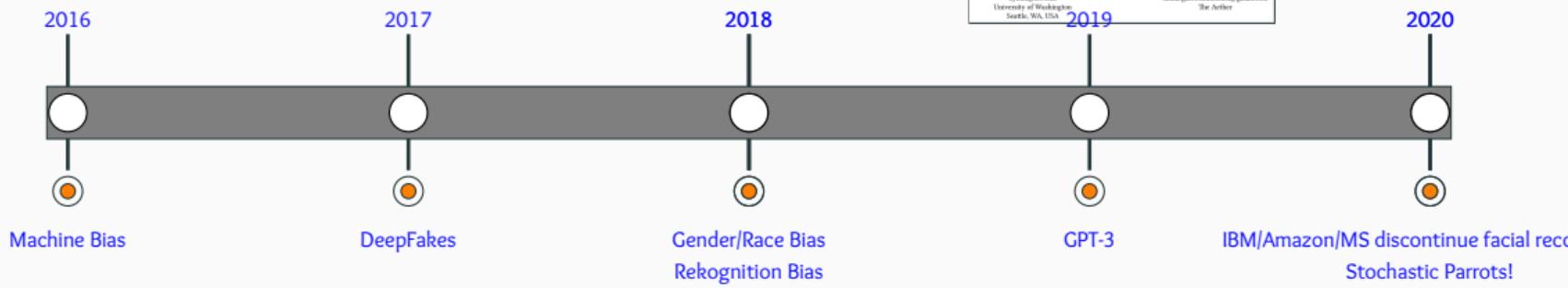
Be Fair and Responsible!

What started off as an innocuous project for mimicking facial expressions has since lead to many apps and creation of fake videos for blackmailing, pronography and swaying elections!



Be Fair and Responsible!

“Models are only as good as the data. Be responsible while curating data.” – *Bender et al.*



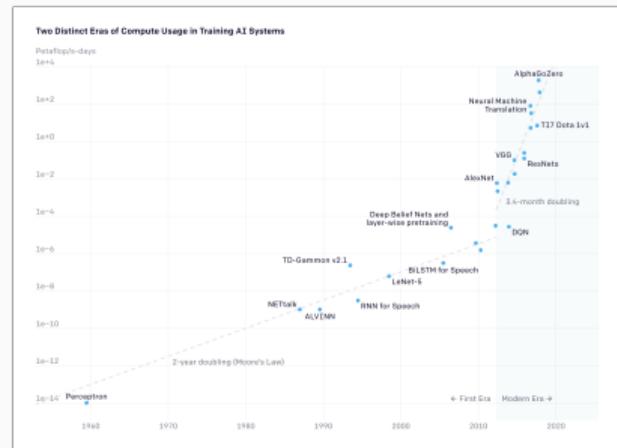
Push for Green AI

The computations required for deep learning research have been doubling every few months, resulting in an estimated 300,000x increase from 2012 to 2018 – AllenAI

Ironically, deep learning was inspired by the human brain, which is remarkably energy efficient.



GreenAI



<https://openai.com/blog/ai-and-compute/>

Push for Green AI

Call for energy and policy considerations for Deep Learning

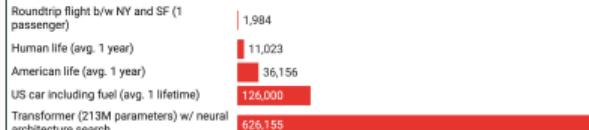
| | Date of original paper | Energy consumption (kWh) | Carbon footprint (lbs of CO2e) | Cloud compute cost (USD) |
|---|------------------------|--------------------------|--------------------------------|--------------------------|
| Transformer (65M parameters) | Jun, 2017 | 27 | 26 | \$41-\$140 |
| Transformer (213M parameters) | Jun, 2017 | 201 | 192 | \$289-\$981 |
| ELMo | Feb, 2018 | 275 | 262 | \$433-\$1,472 |
| BERT (110M parameters) | Oct, 2018 | 1,507 | 1,438 | \$3,751-\$12,571 |
| Transformer (213M parameters) w/ neural architecture search | Jan, 2019 | 656,547 | 626,155 | \$942,973-\$3,201,722 |
| GPT2 | Feb, 2019 | - | - | \$12,902-\$43,008 |

Note: because of a lack of power draw data on GPT2's training hardware, the researchers weren't able to calculate its carbon footprint.

Table: MIT Technology Review • Source: Strubell et al. • Created with Datawrapper

Common carbon footprint benchmarks

in lbs of CO2 equivalent



2019



GreenAI
Energy-Aware NLP

Push for Green AI

“Is it fair that the residents of the Maldives (likely to be underwater by 2100) or the 800,000 people in Sudan affected by drastic floods pay the environmental price of training and deploying ever larger English LMs, when similar large-scale models aren’t being produced for Dhivehi or Sudanese Arabic?” – Bender et. al.



GreenAI
Energy-Aware NLP

On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?

Emily M. Bender*
ebender@uw.edu
University of Washington
Seattle, WA, USA

Angelina McMillan-Major
aymm@uw.edu
University of Washington
Seattle, WA, USA

Timnit Gebru*
timnit@blackinai.org
Black in AI
Palo Alto, CA, USA

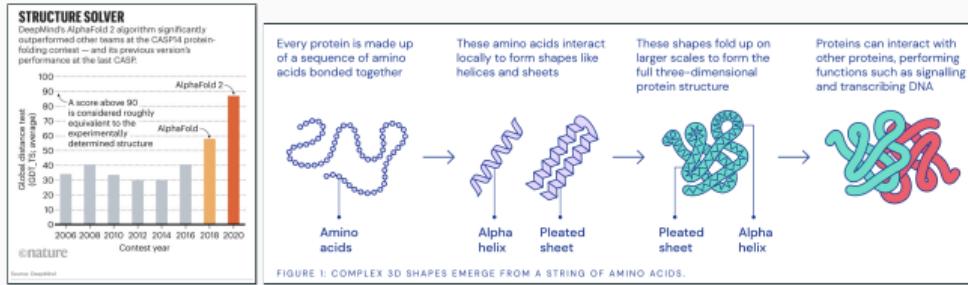
Shmargaret Shmitchell
shmargaret.shmitchell@gmail.com
The Aether



Stochastic Parrots

Chapter 11: The AI revolution in Scientific Research (exciting times ahead!)

Accelerating Scientific Discovery^a



^a<https://deepmind.com/blog/article/AlphaFold-Using-AI-for-scientific-discovery>

<https://ocean.org/stories/spotting-seals-from-space>

<https://www.quantamagazine.org/how-artificial-intelligence-is-changing-science-20190311/>

Accelerating Scientific Discovery^a

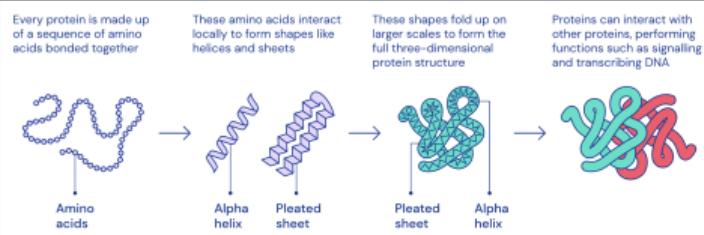
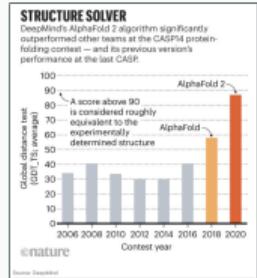
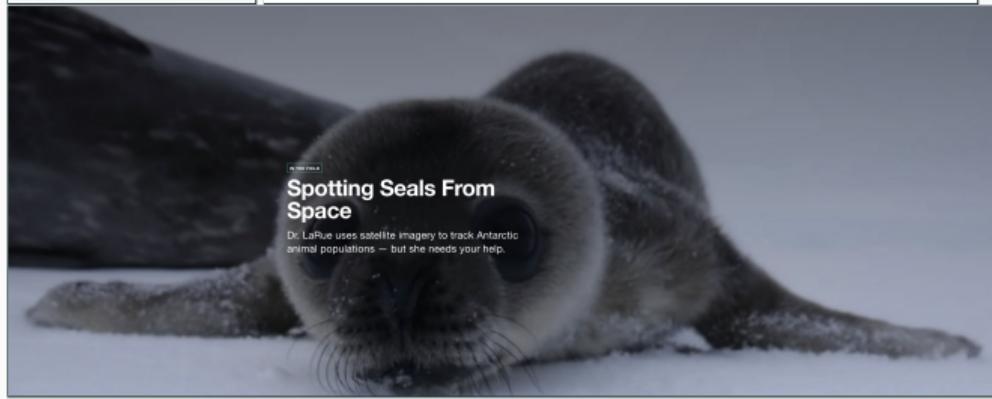


FIGURE 1: COMPLEX 3D SHAPES EMERGE FROM A STRING OF AMINO ACIDS.

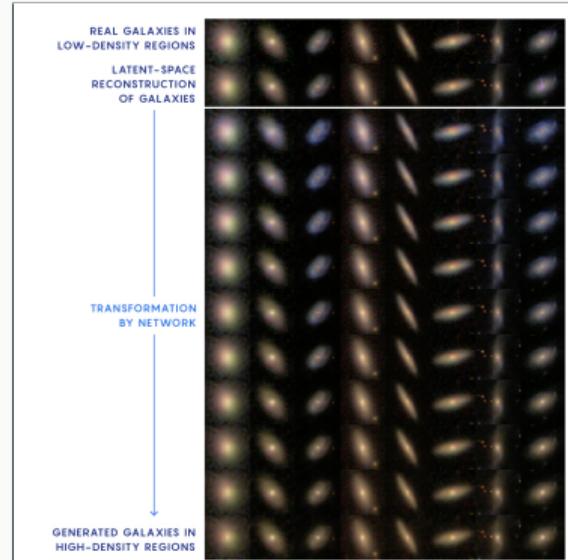
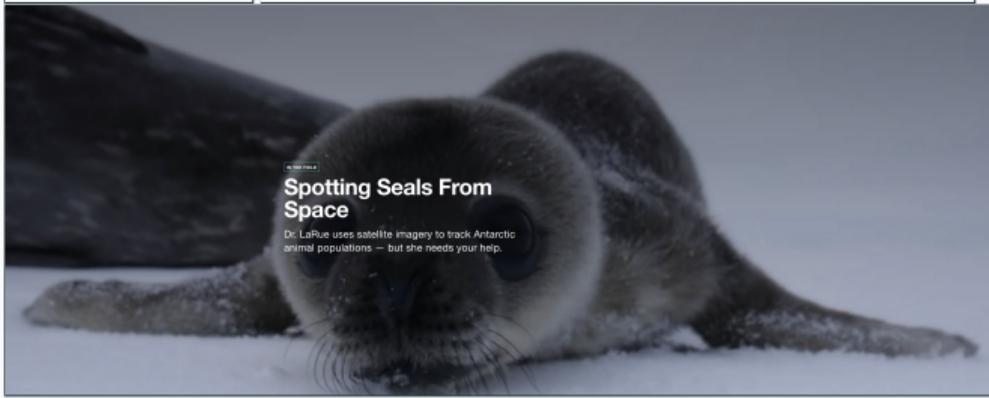
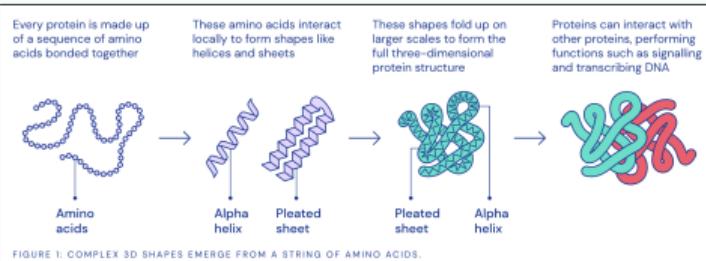
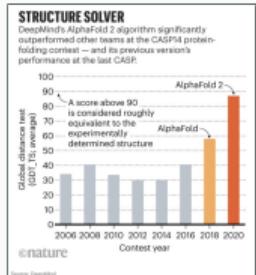


^a<https://deepmind.com/blog/article/AlphaFold-Using-AI-for-scientific-discovery>

<https://ocean.org/stories/spotting-seals-from-space>

<https://www.quantamagazine.org/how-artificial-intelligence-is-changing-science-20190311/>

Accelerating Scientific Discovery^a



Using generative modeling, astrophysicists could investigate how galaxies change when they go from low-density regions of the cosmos to high-density regions, and what physical processes are responsible for these changes.

Adapted from K. Schawinski et al.; Source doi: 10.1086/0004-6365/201833800

^a<https://deepmind.com/blog/article/AlphaFold-Using-AI-for-scientific-discovery>

<https://ocean.org/stories/spotting-seals-from-space>

<https://www.quantamagazine.org/how-artificial-intelligence-is-changing-science-20190311/>

<https://github.com/ChristosChristofidis/awesome-deep-learning>

2020

Healthcare



Finance & Insurance



Transportation



Construction



Retail & Warehousing



Govt. & City Planning



Legal



Mining



Food & Agriculture



Media & Entertainment



Energy



Education



Telecom



Manufacturing



Real Estate



CROSS-INDUSTRY TECH

AI Processors



NLP, NLG, & Computer Vision



Sales & CRM



AI Model Development



Cybersecurity



BI & Ops Intel



DevOps & Model Monitoring



Other R&D



Source: <https://www.cbinsights.com/research/artificial-intelligence-top-startups/>

ⁱSource: <https://www.cbinsights.com/research/artificial-intelligence-top-startups/>

References

- [1] Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, 2015.
- [2] W.S.McCulloch and W.Pitts. A logival calculus of the ideas imminent in nervous activity. 1943.
- [3] A.G. Ivakhnenko and V.G. Lapa. Cybernetic predicting devices. 1965.
- [4] M.Minsky and S.Papert. Perceptrons. 1969.
- [5] P. J. Werbos. Applications of advances in nonlinear sensitivity analysis. In *Proceedings of the 10th IFIP Conference, 31.8 - 4.9, NYC*, pages 762–770, 1981.
- [6] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing*, volume 1, pages 318–362. MIT Press, 1986.
- [7] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- [8] Ruslan Salakhutdinov and Geoffrey Hinton. An efficient learning procedure for deep boltzmann machines. *Neural Comput.*, 24(8):1967–2006, August 2012.
- [9] Alex Graves and Jürgen Schmidhuber. Offline handwriting recognition with multidimensional recurrent neural networks. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21*, pages 545–552. Curran Associates, Inc., 2009.
- [10] G. E. Dahl, Dong Yu, Li Deng, and A. Acero. Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition. *Trans. Audio, Speech and Lang. Proc.*, 20(1):30–42, January 2012.

References ii

- [11] Dan Claudio Ciresan, Ueli Meier, Luca Maria Gambardella, and Jürgen Schmidhuber. Deep big simple neural nets excel on handwritten digit recognition. *CoRR*, abs/1003.0358, 2010.
- [12] Dan C. Ciresan, Ueli Meier, and Jürgen Schmidhuber. Multi-column deep neural networks for image classification. *CoRR*, abs/1202.2745, 2012.
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [14] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. *CoRR*, abs/1311.2901, 2013.
- [15] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [16] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [18] D. H. Wiesel and T. N. Hubel. Receptive fields of single neurones in the cat's striate cortex. *J. Physiol.*, 148:574–591, 1959.
- [19] K. Fukushima. Neocognitron: A self-organizing neural network for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36(4):193–202, 1980.
- [20] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Back-propagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):541–551, 1989.
- [21] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998.

References iii

- [22] J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proc. of the National Academy of Sciences*, 79:2554–2558, 1982.
- [23] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [24] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [25] Matej Moravcík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael H. Bowling. Deepstack: Expert-level artificial intelligence in no-limit poker. *CoRR*, abs/1701.01724, 2017.
- [26] Tomas Mikolov, Martin Karafiát, Lukáš Burget, Jan Černocký, and Sanjeev Khudanpur. Recurrent neural network based language model. In *INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association, Makuhari, Chiba, Japan, September 26-30, 2010*, pages 1045–1048, 2010.
- [27] Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 3294–3302, 2015.
- [28] Yoon Kim, Yacine Jernite, David Sontag, and Alexander M. Rush. Character-aware neural language models. *CoRR*, abs/1508.06615, 2015.
- [29] Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Process. Mag.*, 29(6):82–97, 2012.
- [30] Alex Graves, Abdel-rahman Mohamed, and Geoffrey E. Hinton. Speech recognition with deep recurrent neural networks. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2013, Vancouver, BC, Canada, May 26-31, 2013*, pages 6645–6649, 2013.

References iv

- [31] Jan Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. Attention-based models for speech recognition. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 577–585, 2015.
- [32] Hasim Sak, Andrew W. Senior, Kanishka Rao, and Fran oise Beaufays. Fast and accurate recurrent neural network acoustic models for speech recognition. In *INTERSPEECH 2015, 16th Annual Conference of the International Speech Communication Association, Dresden, Germany, September 6-10, 2015*, pages 1468–1472, 2015.
- [33] Nal Kalchbrenner and Phil Blunsom. Recurrent continuous translation models. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, 18-21 October 2013, Grand Hyatt Seattle, Seattle, Washington, USA, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1700–1709, 2013.
- [34] Kyunghyun Cho, Bart van Merrienoer,  aglar G l ehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1724–1734, 2014.
- [35] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2014.
- [36] S bastien Jean, KyungHyun Cho, Roland Memisevic, and Yoshua Bengio. On using very large target vocabulary for neural machine translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, 2015, Beijing, China, Volume 1: Long Papers*, pages 1–10, 2015.
- [37]  aglar G l ehre, Orhan Firat, Kelvin Xu, Kyunghyun Cho, Lo  Barrault, Huei-Chi Lin, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. On using monolingual corpora in neural machine translation. *CoRR*, abs/1503.03535, 2015.

References v

- [38] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 3104–3112, 2014.
- [39] Thang Luong, Hieu Pham, and Christopher D. Manning. Effective approaches to attention-based neural machine translation. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 1412–1421, 2015.
- [40] Hao Zheng, Yong Cheng, and Yang Liu. Maximum expected likelihood estimation for zero-resource neural machine translation. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, pages 4251–4257, 2017.
- [41] Yong Cheng, Qian Yang, Yang Liu, Maosong Sun, and Wei Xu. Joint training for pivot-based neural machine translation. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, pages 3974–3980, 2017.
- [42] Yun Chen, Yang Liu, Yong Cheng, and Victor O. K. Li. A teacher-student framework for zero-resource neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1925–1935, 2017.
- [43] Orhan Firat, Baskaran Sankaran, Yaser Al-Onaizan, Fatos T. Yarman-Vural, and Kyunghyun Cho. Zero-resource translation with multi-lingual neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 268–277, 2016.
- [44] Lifeng Shang, Zhengdong Lu, and Hang Li. Neural responding machine for short-text conversation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, 2015, Beijing, China, Volume 1: Long Papers*, pages 1577–1586, 2015.
- [45] Oriol Vinyals and Quoc V. Le. A neural conversational model. *CoRR*, abs/1506.05869, 2015.

References vi

- [46] Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. The ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems. In *Proceedings of the SIGDIAL 2015 Conference, The 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue, 2-4 September 2015, Prague, Czech Republic*, pages 285–294, 2015.
- [47] Jesse Dodge, Andreea Gane, Xiang Zhang, Antoine Bordes, Sumit Chopra, Alexander H. Miller, Arthur Szlam, and Jason Weston. Evaluating prerequisite qualities for learning end-to-end dialog systems. *CoRR*, abs/1511.06931, 2015.
- [48] Jason Weston, Antoine Bordes, Sumit Chopra, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. *CoRR*, abs/1502.05698, 2015.
- [49] Iulian Vlad Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron C. Courville, and Yoshua Bengio. A hierarchical latent variable encoder-decoder model for generating dialogues. *CoRR*, abs/1605.06069, 2016.
- [50] Antoine Bordes and Jason Weston. Learning end-to-end goal-oriented dialog. *CoRR*, abs/1605.07683, 2016.
- [51] Iulian Vlad Serban, Chinnadhurai Sankar, Mathieu Germain, Saizheng Zhang, Zhouhan Lin, Sandeep Subramanian, Taesup Kim, Michael Pieper, Sarath Chandar, Nan Rosemary Ke, Sai Mudumba, Alexandre de Brébisson, Jose Sotelo, Dendi Suhubdy, Vincent Michalski, Alexandre Nguyen, Joelle Pineau, and Yoshua Bengio. A deep reinforcement learning chatbot. *CoRR*, abs/1709.02349, 2017.
- [52] Karl Moritz Hermann, Tomás Kocišký, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 1693–1701, 2015.
- [53] Danqi Chen, Jason Bolton, and Christopher D. Manning. A thorough examination of the cnn/daily mail reading comprehension task. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*, 2016.
- [54] Caiming Xiong, Victor Zhong, and Richard Socher. Dynamic coattention networks for question answering. *CoRR*, abs/1611.01604, 2016.

References vii

- [55] Min Joon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. Bidirectional attention flow for machine comprehension. *CoRR*, abs/1611.01603, 2016.
- [56] Bhuvan Dhingra, Hanxiao Liu, Zhilin Yang, William W. Cohen, and Ruslan Salakhutdinov. Gated-attention readers for text comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 1832–1846, 2017.
- [57] Wenhui Wang, Nan Yang, Furu Wei, Baobao Chang, and Ming Zhou. Gated self-matching networks for reading comprehension and question answering. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, pages 189–198, 2017.
- [58] Minghao Hu, Yuxing Peng, and Xipeng Qiu. Mnemonic reader for machine comprehension. *CoRR*, abs/1705.02798, 2017.
- [59] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 3431–3440, 2015.
- [60] Ming Liang and Xiaolin Hu. Recurrent convolutional neural network for object recognition. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 3367–3375, 2015.
- [61] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(6):1137–1149, 2017.
- [62] Sean Bell, C. Lawrence Zitnick, Kavita Bala, and Ross B. Girshick. Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks. *CoRR*, abs/1512.04143, 2015.
- [63] Joseph Redmon and Ali Farhadi. YOLO9000: better, faster, stronger. *CoRR*, abs/1612.08242, 2016.

References viii

- [64] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-FCN: object detection via region-based fully convolutional networks. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 379–387, 2016.
- [65] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2980–2988, 2017.
- [66] Sergi Caelles, Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. One-shot video object segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 5320–5329, 2017.
- [67] Janghoon Choi, Junseok Kwon, and Kyoung Mu Lee. Visual tracking by reinforced decision making. *CoRR*, abs/1702.06291, 2017.
- [68] Sangdoo Yun, Jongwon Choi, Youngjoon Yoo, Kimin Yun, and Jin Young Choi. Action-decision networks for visual tracking with deep reinforcement learning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1349–1358, 2017.
- [69] Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Tracking the untrackable: Learning to track multiple cues with long-term dependencies. *CoRR*, abs/1701.01909, 2017.
- [70] Junhua Mao, Wei Xu, Yi Yang, Jiang Wang, and Alan L. Yuille. Deep captioning with multimodal recurrent neural networks (m-rnn). *CoRR*, abs/1412.6632, 2014.
- [71] Junhua Mao, Xu Wei, Yi Yang, Jiang Wang, Zhiheng Huang, and Alan L. Yuille. Learning like a child: Fast novel visual concept learning from sentence descriptions of images. In *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [72] Ryan Kiros, Ruslan Salakhutdinov, and Richard S. Zemel. Unifying visual-semantic embeddings with multimodal neural language models. *CoRR*, abs/1411.2539, 2014.

References ix

- [73] Jeff Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Trevor Darrell, and Kate Saenko. Long-term recurrent convolutional networks for visual recognition and description. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 2625–2634, 2015.
- [74] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 3156–3164, 2015.
- [75] Andrej Karpathy and Fei-Fei Li. Deep visual-semantic alignments for generating image descriptions. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 3128–3137, 2015.
- [76] Hao Fang, Saurabh Gupta, Forrest N. Iandola, Rupesh Kumar Srivastava, Li Deng, Piotr Dollár, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C. Platt, C. Lawrence Zitnick, and Geoffrey Zweig. From captions to visual concepts and back. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1473–1482, 2015.
- [77] Kan Chen, Jiang Wang, Liang-Chieh Chen, Haoyuan Gao, Wei Xu, and Ram Nevatia. ABC-CNN: an attention based convolutional neural network for visual question answering. *CoRR*, abs/1511.05960, 2015.
- [78] Jeff Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. *CoRR*, abs/1411.4389, 2014.
- [79] Subhashini Venugopalan, Huijuan Xu, Jeff Donahue, Marcus Rohrbach, Raymond J. Mooney, and Kate Saenko. Translating videos to natural language using deep recurrent neural networks. In *NAACL HLT 2015, The 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA, May 31 - June 5, 2015*, pages 1494–1504, 2015.
- [80] Yingwei Pan, Tao Mei, Ting Yao, Houqiang Li, and Yong Rui. Jointly modeling embedding and translation to bridge video and language. *CoRR*, abs/1505.01861, 2015.

References x

- [81] Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher J. Pal, Hugo Larochelle, and Aaron C. Courville. Describing videos by exploiting temporal structure. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 4507–4515, 2015.
- [82] Anna Rohrbach, Marcus Rohrbach, Wei Qiu, Annemarie Friedrich, Manfred Pinkal, and Bernt Schiele. Coherent multi-sentence video description with variable level of detail. In *Pattern Recognition - 36th German Conference, GCPR 2014, Münster, Germany, September 2-5, 2014, Proceedings*, pages 184–195, 2014.
- [83] Linchao Zhu, Zhongwen Xu, Yi Yang, and Alexander G. Hauptmann. Uncovering temporal context for video question and answering. *CoRR*, abs/1511.04670, 2015.
- [84] Adam Santoro, David Raposo, David G. T. Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Tim Lillicrap. A simple neural network module for relational reasoning. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 4974–4983, 2017.
- [85] Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. Learning to reason: End-to-end module networks for visual question answering. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 804–813, 2017.
- [86] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1988–1997, 2017.
- [87] Hedi Ben-younes, Rémi Cadène, Matthieu Cord, and Nicolas Thome. MUTAN: multimodal tucker fusion for visual question answering. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2631–2639, 2017.
- [88] Mateusz Malinowski, Marcus Rohrbach, and Mario Fritz. Ask your neurons: A neural-based approach to answering questions about images. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 1–9, 2015.
- [89] Vahid Kazemi and Ali Elqursh. Show, ask, attend, and answer: A strong baseline for visual question answering. *CoRR*, abs/1704.03162, 2017.

References xi

- [90] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. Movieqa: Understanding stories in movies through question-answering. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 4631–4640, 2016.
- [91] Kuo-Hao Zeng, Tseng-Hung Chen, Ching-Yao Chuang, Yuan-Hong Liao, Juan Carlos Niebles, and Min Sun. Leveraging video descriptions to learn video question answering. *CoRR*, abs/1611.04021, 2016.
- [92] Tegan Maharaj, Nicolas Ballas, Anna Rohrbach, Aaron C. Courville, and Christopher Joseph Pal. A dataset and exploration of models for understanding video data through fill-in-the-blank question-answering. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 7359–7368, 2017.
- [93] Zhou Zhao, Qifan Yang, Deng Cai, Xiaofei He, and Yueteng Zhuang. Video question answering via hierarchical spatio-temporal attention networks. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, pages 3518–3524, 2017.
- [94] Youngjae Yu, Hyungjin Ko, Jongwook Choi, and Gunhee Kim. End-to-end concept word detection for video captioning, retrieval, and question answering. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 3261–3269, 2017.
- [95] Hongyang Xue, Zhou Zhao, and Deng Cai. The forgettable-watcher model for video question answering. *CoRR*, abs/1705.01253, 2017.
- [96] Amir Mazaheri, Dong Zhang, and Mubarak Shah. Video fill in the blank with merging lstms. *CoRR*, abs/1610.04062, 2016.
- [97] Tommy Chheng. Video summarization using clustering.
- [98] Muhammad Ajmal, Muhammad Husnain Ashraf, Muhammad Shakir, Yasir Abbas, and Faiz Ali Shah. Video summarization: Techniques and classification. In *Computer Vision and Graphics - International Conference, ICCUG 2012, Warsaw, Poland, September 24-26, 2012. Proceedings*, pages 1–13, 2012.

References xii

- [99] Ke Zhang, Wei-Lun Chao, Fei Sha, and Kristen Grauman. Video summarization with long short-term memory. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VII*, pages 766–782, 2016.
- [100] Zhong Ji, Kailin Xiong, Yanwei Pang, and Xuelong Li. Video summarization with attention-based encoder-decoder networks. *CoRR*, abs/1708.09545, 2017.
- [101] Rameswar Panda, Niluthpol Chowdhury Mithun, and Amit K. Roy-Chowdhury. Diversity-aware multi-video summarization. *IEEE Trans. Image Processing*, 26(10):4712–4724, 2017.
- [102] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013.
- [103] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680, 2014.
- [104] Anh Nguyen, Jason Yosinski, Yoshua Bengio, Alexey Dosovitskiy, and Jeff Clune. Plug & play generative networks: Conditional iterative generation of images in latent space. *CoRR*, abs/1612.00005, 2016.
- [105] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *CoRR*, abs/1710.10196, 2017.
- [106] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alexander Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. In *Arxiv*, 2016.
- [107] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. *arXiv preprint arXiv:1601.06759*, 2016.

References xiii

- [108] Aaron van den Oord, Nal Kalchbrenner, Lasse Espeholt, koray kavukcuoglu, Oriol Vinyals, and Alex Graves. Conditional image generation with pixelcnn decoders. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4790–4798. Curran Associates, Inc., 2016.
- [109] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint arXiv:1701.05517*, 2017.
- [110] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017.
- [111] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics, 2019.