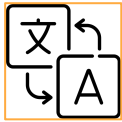# Honouring the 2018 ACM Turing Award Laureates

Geoffrey E. Hinton, Yann LeCun, Yoshua Bengio
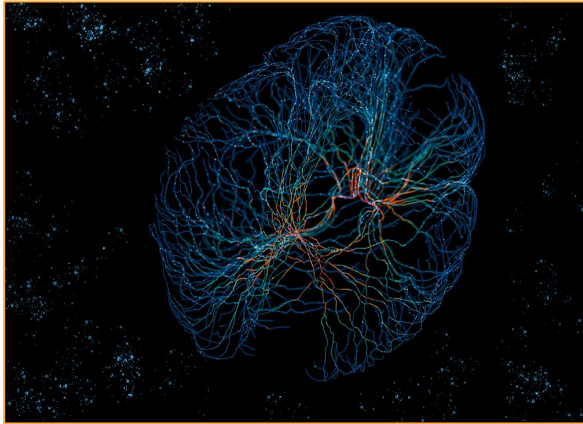
"for conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing"

# Artificial intelligence
# is transforming nearly everything.

AI is the most rapidly growing area in all of science, and certainly one of the most talked-about topics in society.

# The success of AI
# is driven by deep learning.



The incredible advances in AI would not have been possible without some of the foundations, namely, deep learning.

Deep learning is a subset of machine learning that uses **artificial neural networks** to learn from data. Artificial neural networks are inspired by the human brain, and they are able to learn complex patterns from large amounts of data. Deep learning has been used ...

▲ the response from Google Bard

# G. Hinton, Y. LeCun, and Y. Bengio pioneered deep learning.
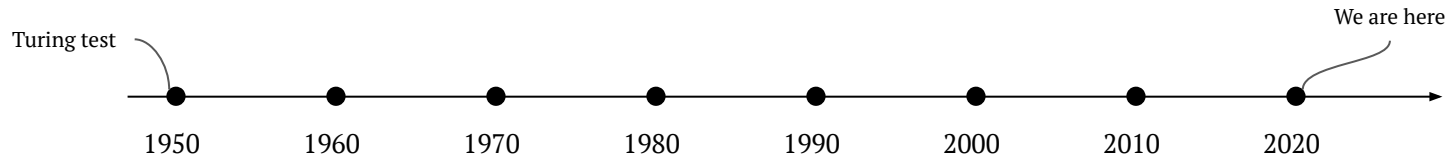


They collectively and independently worked over 30 years.
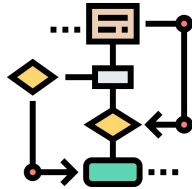
- establish <u>conceptual foundations</u> for deep networks

- identify a lot of <u>interesting phenomena</u>

- develop <u>engineering advances</u> in practice

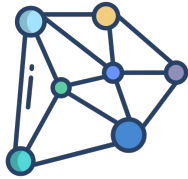Apparently, it wasn't always easy.

# Two paradigms since 1950s

Turing test

We are here

1950 1960 1970 1980 1990 2000 2010 2020

# The two paradigms:
# Symbolic AI and Connectionist AI



## Symbolic AI

- the logic-inspired approach
- use human-readable symbolic rules
- focus on reasoning



## Connectionist AI

- the biological-inspired approach
- learn strengths of the connections in a neural network (vectors)
- focus on learning

# While the learning approach can solve it, the symbolic approach can't.



**What are shown in the photo?**

A man and a chicken.

**What does the man feel and why?**

He is scared of the chicken because it is flying at him.

Image captioning task:
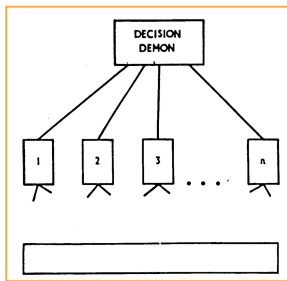
Symbolic AI people tried this for a long time, but <span style="color:red">they failed</span>.
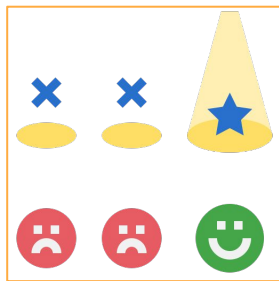
- It's not obvious how to write that program.

It became fairly easy for <span style="color:orange">connectionist AI</span> to solve this task.

- Current methods can capture subtleties.

# The central question:
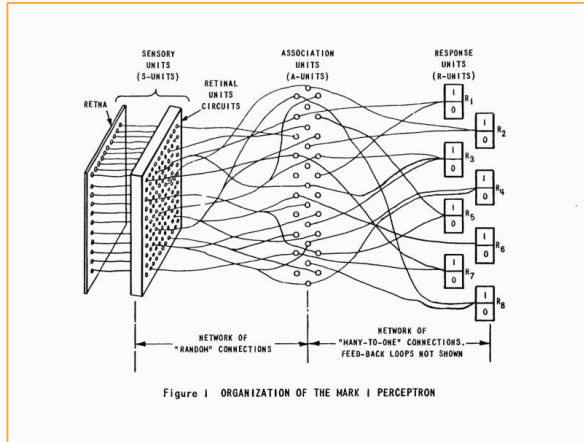# can they learn to do it?



Pandemonium (1959)



trial and error

Large neural networks are very powerful computing devices.

- But can a neural network learn a difficult tasks?

Early researchers like Turing and Selfridge proposed that neural networks with initially random connections could be trained by reinforcement learning.

- This is extremely inefficient.

# Perceptrons:
# A simple learning procedure.
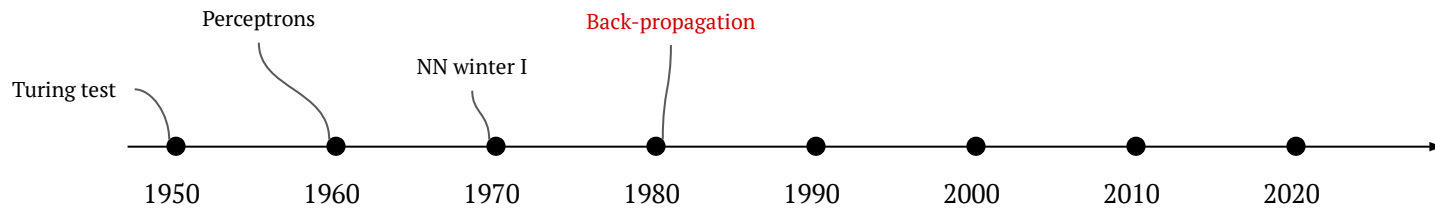


Figure 1 ORGANIZATION OF THE MARK I PERCEPTRON

~1960: Rosenblatt introduced a simple, efficient learning procedure that could figure out how to weight features of the input in order to classify inputs correctly.
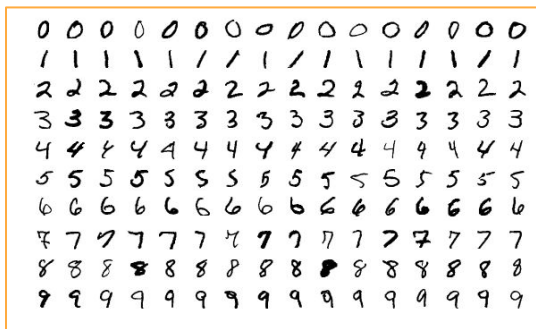
1969: Minsky and Papert showed that perceptrons had some very strong limitations on what they could do.

1970s: The first neural net winter has begun.

# Back-propagation in 1980s

Turing test     Perceptrons     NN winter I     Back-propagation

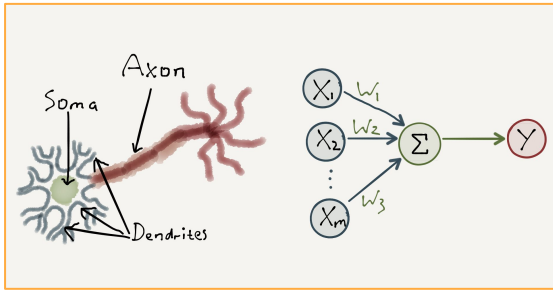1950     1960     1970     1980     1990     2000     2010     2020

# Back-propagation created a lot of excitement.



1980s: The back-propagation procedure allows neural networks to design their own features and have multiple layers of features.
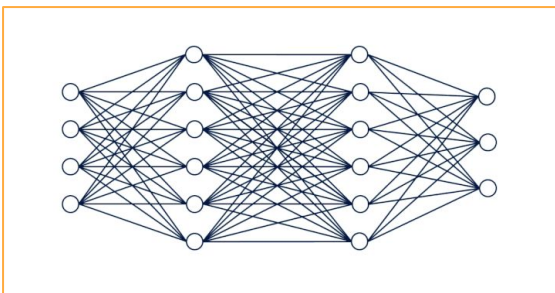
# What is an artificial neuron?



McCulloch-Pitts neuron (1943)

We make a gross idealization of a real neuron so that we can investigate how neurons can collaborate to do computations that are too difficult to program such as:

- Convert the pixel intensity values of an image into a string of words that describe the image.

# What is an artificial neural network?



A feed-forward neural network

You have weights on the incoming weights for each of these neurons, and as you change those incoming weights, you're changing what feature that neuron will respond to. So by learning these weights, you're learning the features.

The question is how are we gonna train it?

- Supervised (or unsupervised) training
- "Mutation" methods?

# How do we train artificial neural networks?
# A "mutation" method that is easy to understand?

Supervised training: Show the network an input vector and tell it to the correct output.

- Adjust the weights to reduce the discrepancy between the correct output and the actual output.

A mutation method:

- Pick one weight. Increase or decrease the weight slightly and measure how well the network now does. Keep it if it helped.
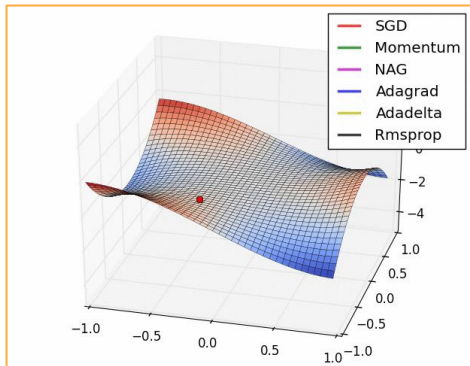
# The backpropagation algorithm



Instead of perturbing the weights one at a time and measuring the effect, use calculus to compute the error gradients for all of the weights at the same time.

- With a million weights, this is more efficient than the mutation method by a factor of a million.
- Once the gradient is computed, perform stochastic gradient descent, which works really well at scale.

# How to learn many layers of features (~1985)



1. Forward pass through the network (for a small batch).

2. Calculate the difference between what you got and what you wanted.

3. Backward pass with the chain rule.

4. Take a stochastic optimisation step (e.g., SGD).

- It works really well at scale.

# A big disappointment,
# and the 2$^{nd}$ neural net winter.

1990s: BP works pretty well, but underperforms the expectations of its proponents on modest-sized datasets.

- Symbolic AI: it is silly to expect to learn difficult tasks in big deep nets that start with random connections and no prior knowledge.

- A series of rejections from NIPS, ICML, and CVPR.

The second neural network winter begins.

REJECTED

# the Deep learning revolution

Turing test

Perceptrons

NN winter I

Back-propagation

NN winter II

Deep Learning Revolution

1950   1960   1970   1980   1990   2000   2010   2020

# The return of backpropagation,
# largely due to a lot of data and compute.

Between 2005 and 2009, several technical advances enabled back-propagation to work better in feed-forward nets.

- Unsupervised pre-training; random dropout of units; rectified linear units.

Back-propagation works amazingly well if you have two things:

- a lot of labeled data
- a lot of convenient compute power (e.g. GPUs)

# Acoustic modeling:
# The killer App (Mohamed, Dahl & Hinton 2009)



Acoustic modeling: for the middle frame of the spectrogram, which piece of which phone in the speaker is trying to express?

Soon after, leading speech groups at MSR, IBM & Google developed them further.

# Object Recognition:
# the 2012 ImageNet object recognition challenge.





The challenge:

- Given a million high-resolution training images of 1000 different classes of object, get the "correct" class in your top 5 bets.

Error rates:

- While all previous / standard ones asymptote at about 25% error, the AlexNet got 16% error.

- By 2015 it was down to 5%. And now it's down to considerably below that.

# A radically new way to do machine translation (Sutskever, Vinayals and Le, 2014)





The sequence-to-sequence model:

- The encoder read reads in the sequence of words, and the decoder expresses the thought in the target language ("thought vector").

A lot of advances since 2014:

- Soft attention, pre-training, and the transformer networks.

The final nail in the coffin of symbolic AI:

- Machine translation is an idea task for symbolic AI because the input is symbols and the output is symbols ("vectors inside").

# Continuing advances

# The future of neural network vision



$x$

"panda"
57.7% confidence

$\mathrm{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"nematode"
8.2% confidence

$x + \epsilon\,\mathrm{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"gibbon"
99.3 % confidence

Adversarial example (2015)



Stacked Capsule Autoencoder (2019)

Convolutional nets get a huge win by wiring in the idea that if a feature is useful in one location it is useful everywhere.

- But they do not recognize objects the same way as us, hence adversarial examples.

People recognize objects by using the viewpoint invariant geometrical relationships between an object and its parts.

- We can make neural networks do this by using transformers.

# The future of neural networks



Forward-forward, Hinton (2022)

Nearly all artificial neural nets use only two time scales: slow adaptation of weights and fast changes in neural activity.
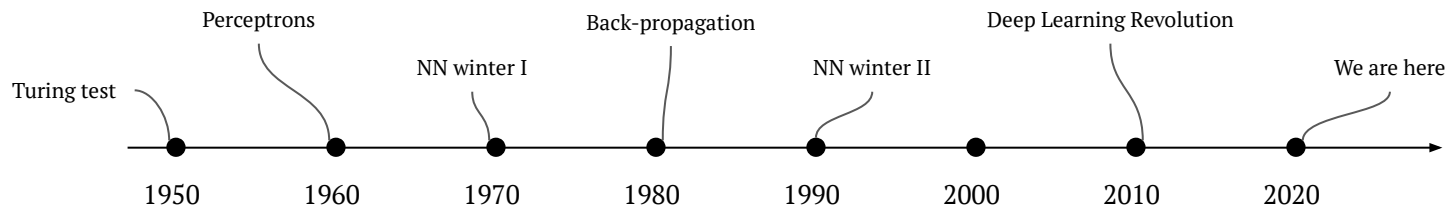
But synapses adapt at multiple different time scales.

- Using fast weights for short-term memory will make neural networks different and better.
- It can improve optimization.
- It allows true recursion (1973, unpublished).

Closing

# Summary &
# The message



We've seen the *history* of deep learning, and its transformative and revolutionary *impact* on society.

Deep learning with large neural networks is becoming *critical* for nearly everything.

# Credits



"<u>a baby is laying on a seat with a remote</u>",
generated using 🤗 Hugging Face API

G. Hinton and Y. LeCun's Turing lectures at 2019 ACM FCRC,

J. Doumont's Trees, maps, and theorems,

various sources on the Internet for visual materials,

the audience who just came to this event (on Friday),

and, of course, my family.

# Thank you

-Namhoon Lee, POSTECH