



Beyond MLPs

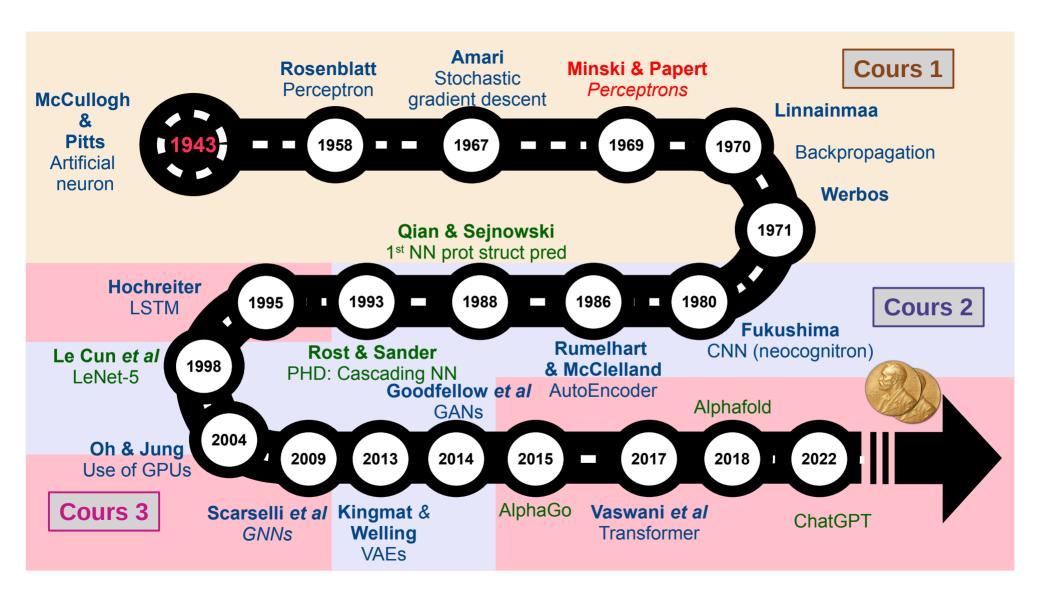
Part 2/2: RNNs, Attention and GNNs

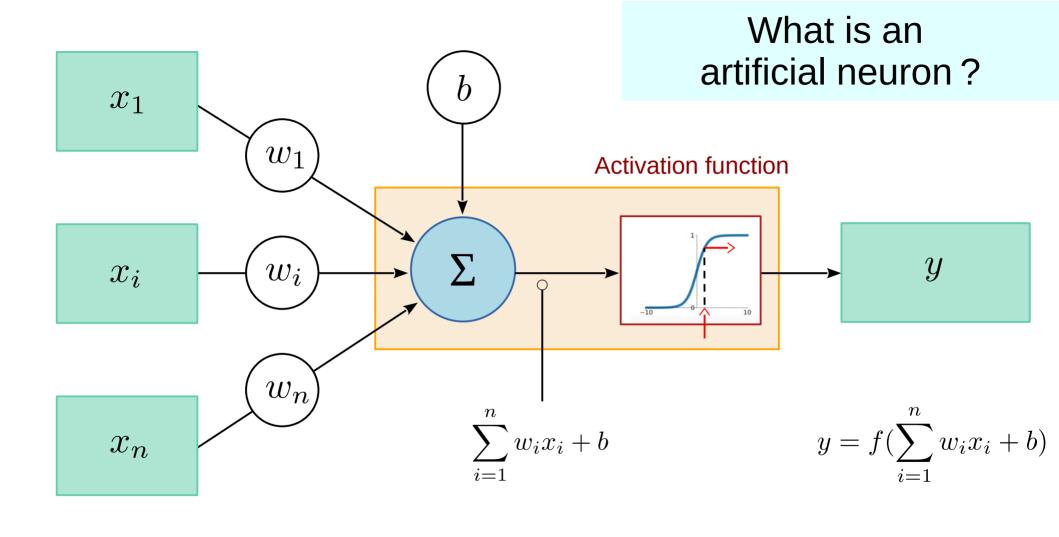
Nicolas Gambardella

nicolas.gambardella@univ-lille.fr

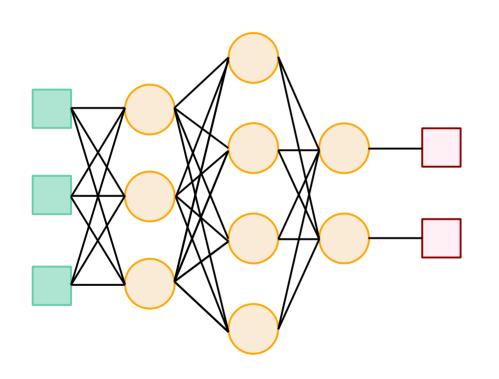




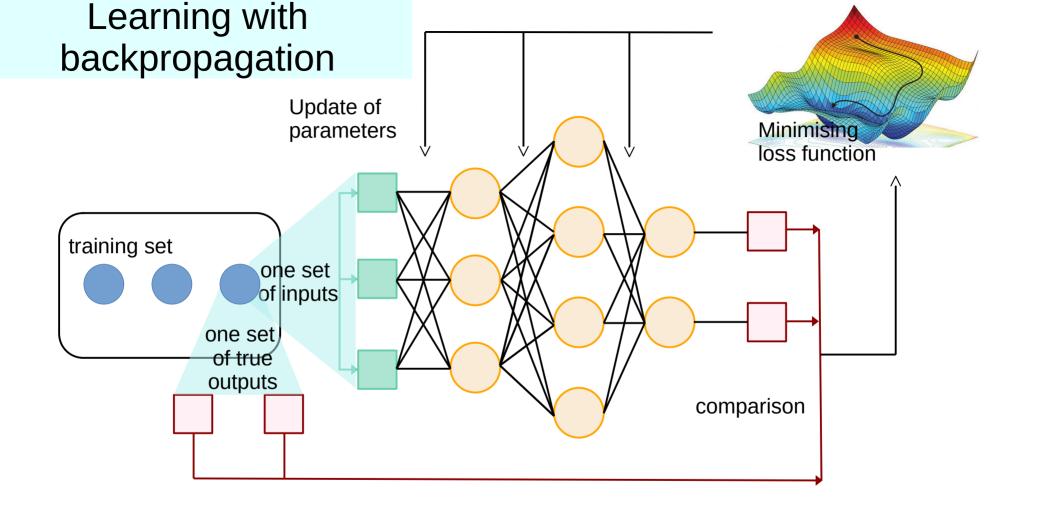




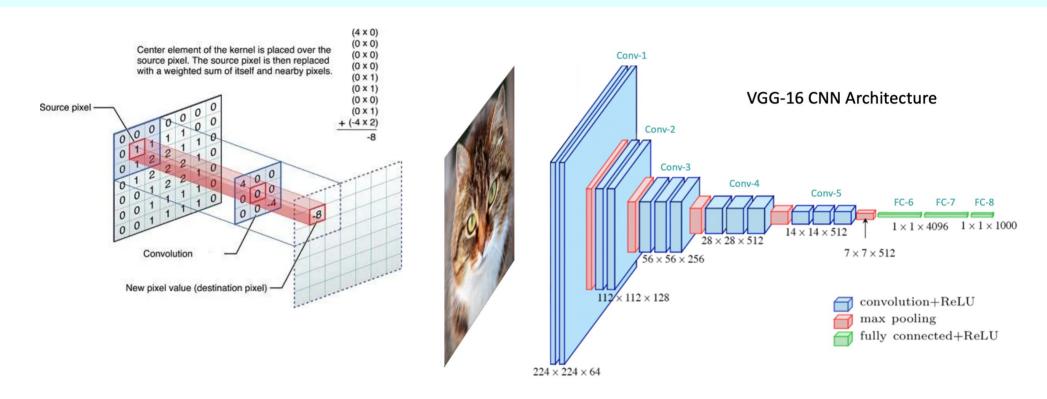
All inputs can be independent and everything connected to everything



Multi-Layer Perceptrons (MLP) or Dense neural networks (DNN) made of Fully Connected layers (FC)

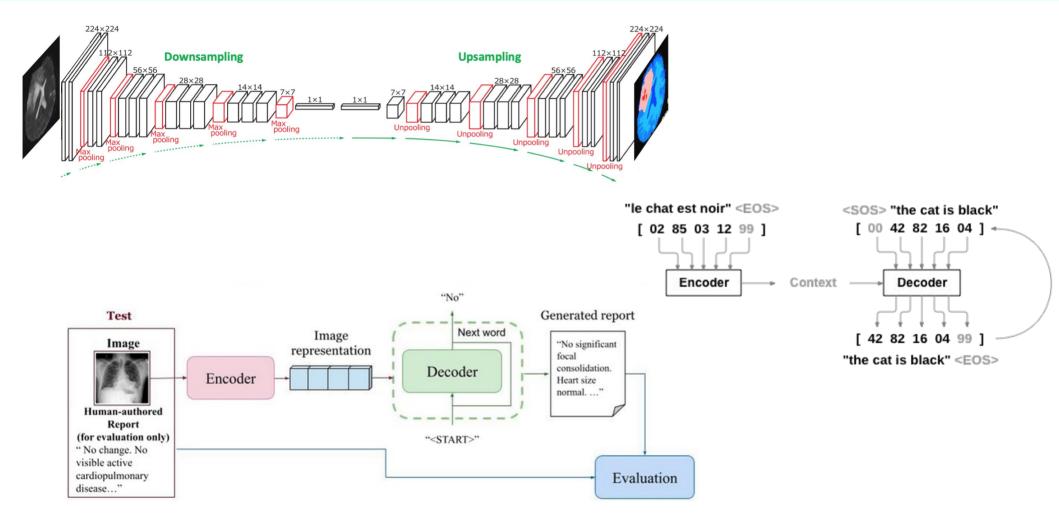


We can detect local features by linking neighbouring inputs

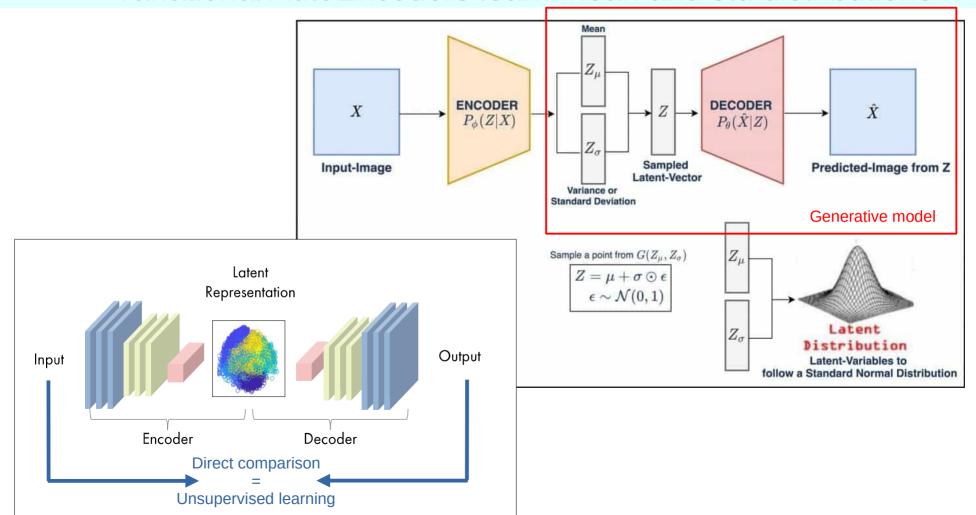


Deep Convolutional Neural Networks (CNN)

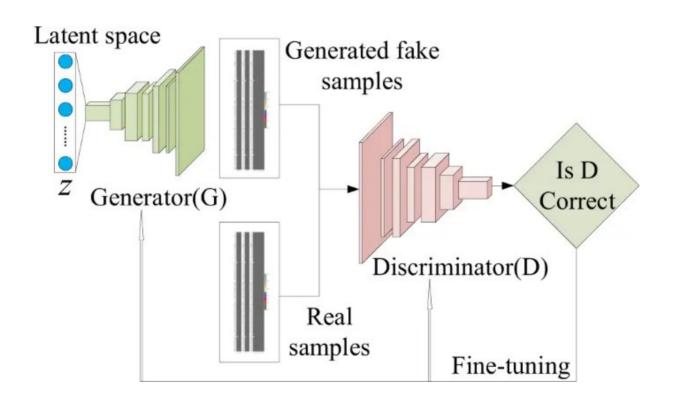
Encoder networks can embed information in a latent space Decoder networks can reconstruct the information from it



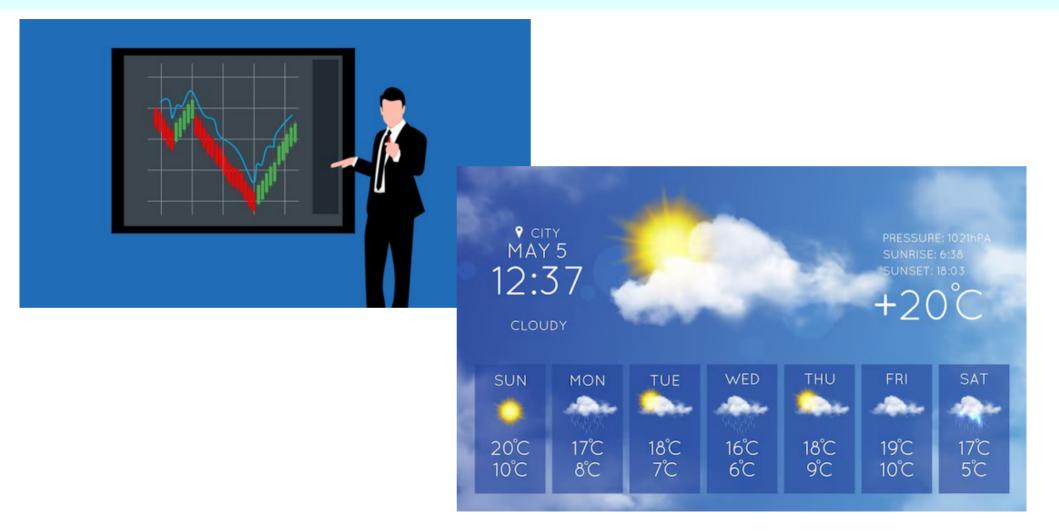
AutoEncoders can train themselves unsupervised Variational AutoEncoders learn mean and std distributions



Generative Adversarial Networks learn by trying to deceive themselves



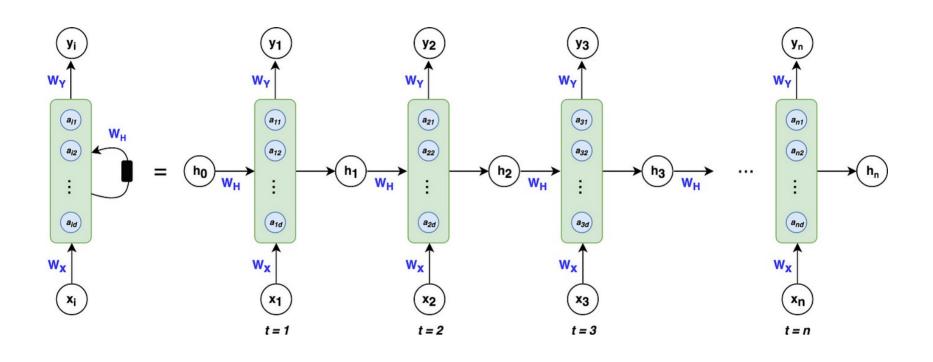
Time series and sequences of variable lengths



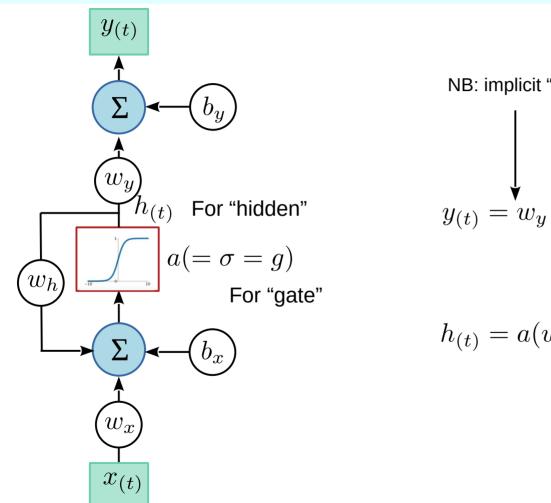
Time series and sequences of variable lengths

"The quick brown fox jumped Speech recognition over the lazy dog." Music generation "There is nothing to like Sentiment classification in this movie." DNA sequence analysis AGCCCCTGTGAGGAACTAG AGCCCCTGTGAGGAACTAG Voulez-vous chanter avec Do you want to sing with Machine translation moi? me? Video activity recognition Running Yesterday, Harry Potter Yesterday, Harry Potter Name entity recognition met Hermione Granger. met Hermione Granger. Andrew Ng

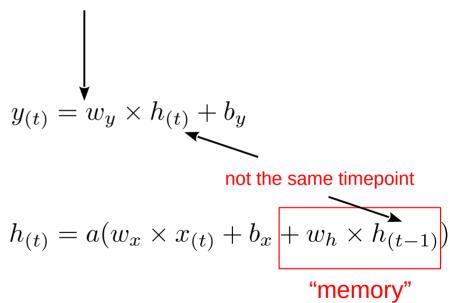
Recurrent Neural Networks: successive inputs are not independent



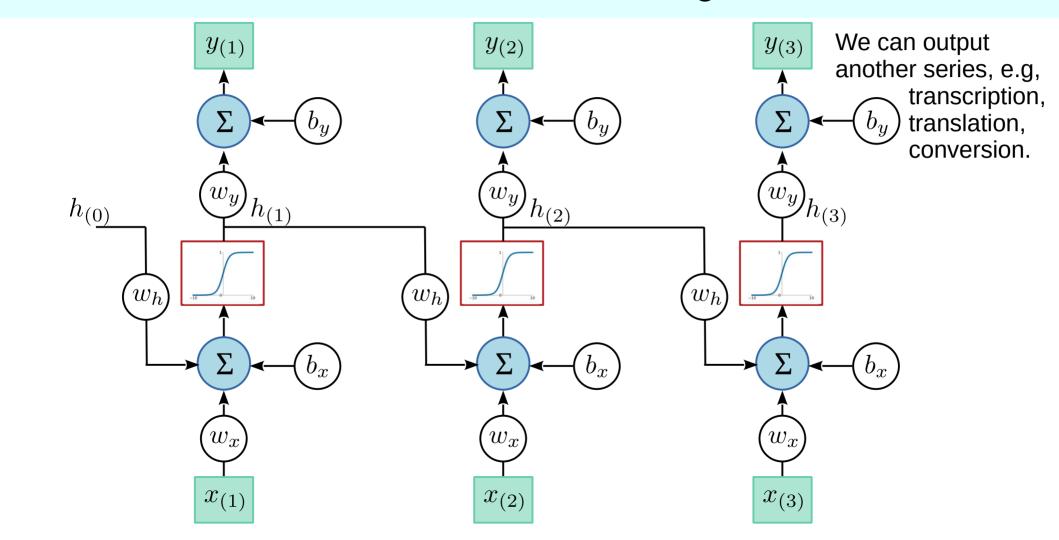
RNN: 1 cell (here, 1 neuron)



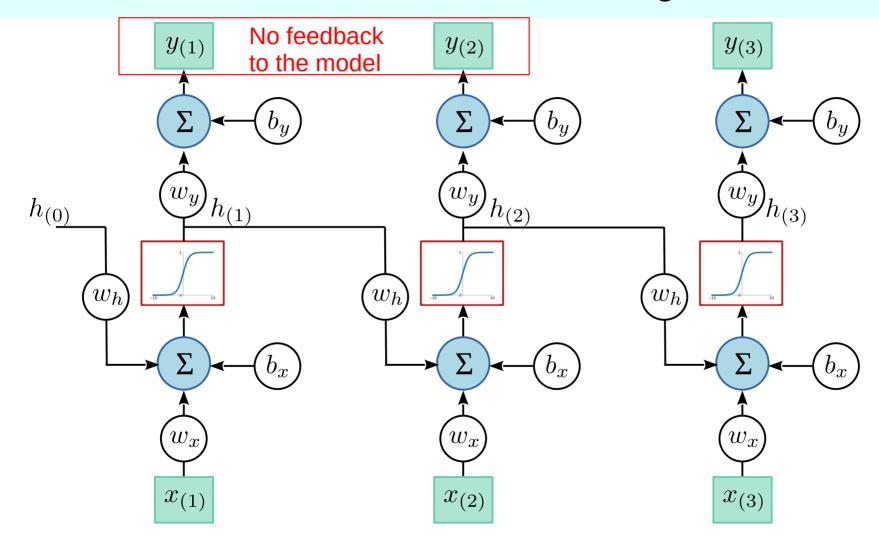
NB: implicit "identity" activation function



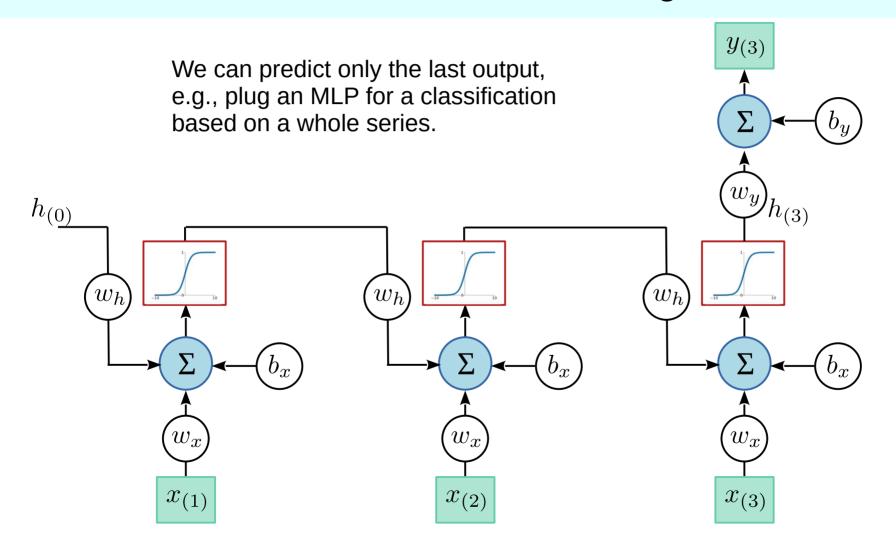
RNN: 1 cell - unfolding



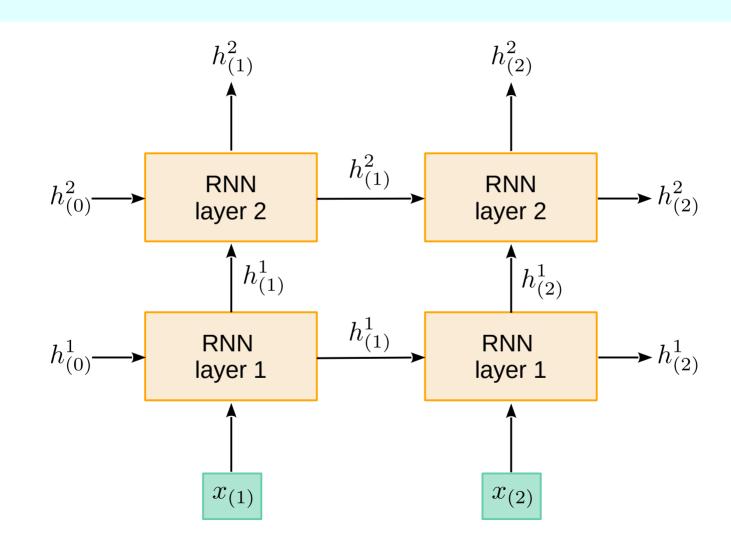
RNN: 1 cell - unfolding



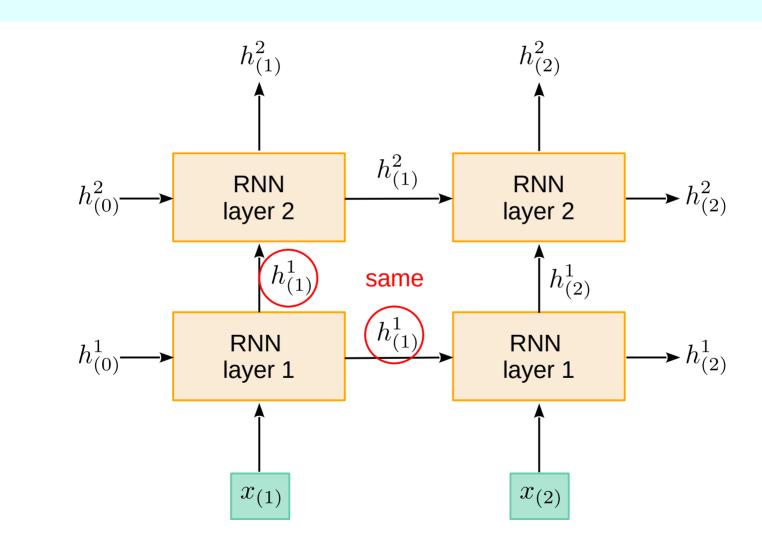
RNN: 1 cell - unfolding



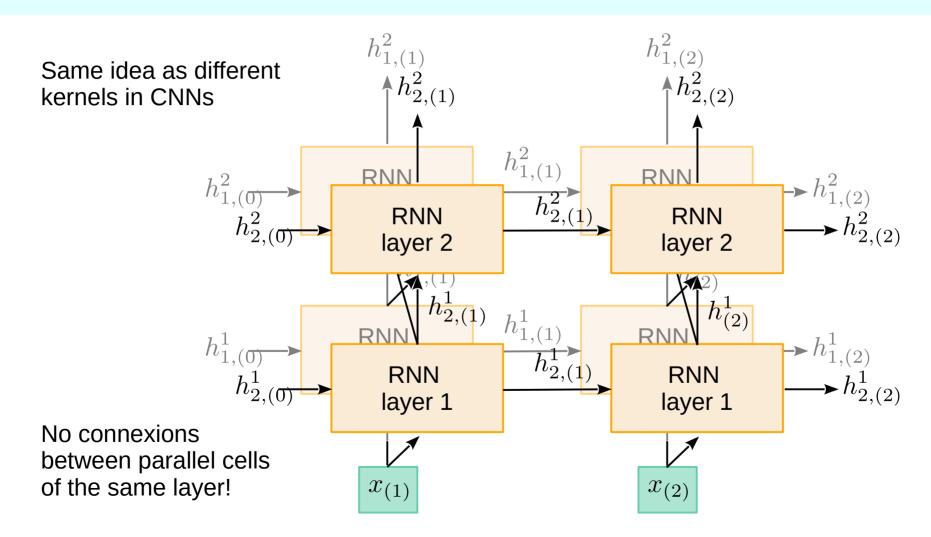
Stacked RNNs



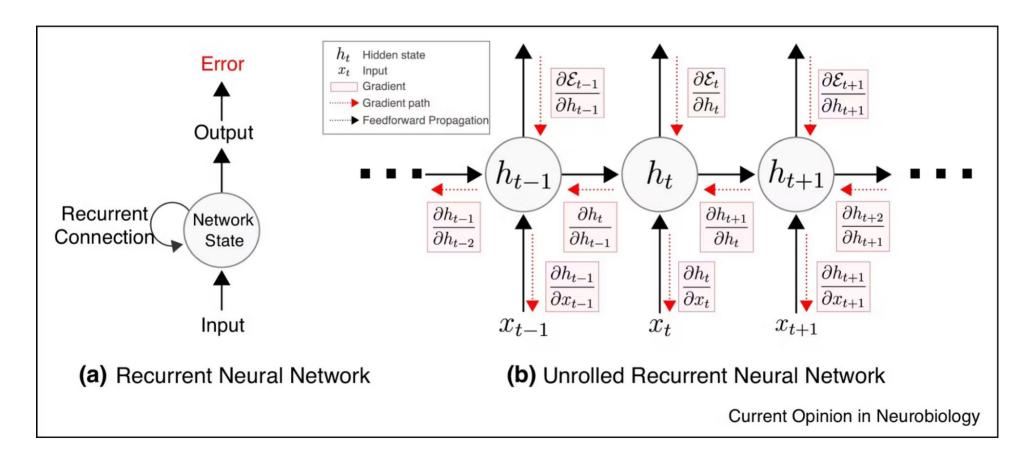
Stacked RNN



Several RNNs may learn different patterns in parallel

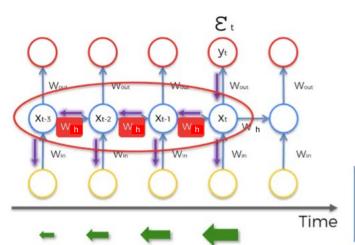


Learning by backpropagation in time



Lillicrap and Santaro (2019) Backpropagation through time and the brain. Curr Op Neurobiol, 55:81-89

Exploding and vanishing gradients



$$\begin{split} \frac{\partial \mathcal{E}}{\partial \theta} &= \sum_{1 \leq t \leq T} \frac{\partial \mathcal{E}_t}{\partial \theta} \\ \frac{\partial \mathcal{E}_t}{\partial \theta} &= \sum_{1 \leq k \leq t} \left(\frac{\partial \mathcal{E}_t}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} \frac{\partial^+ \mathbf{x}_k}{\partial \theta} \right) \\ \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} &= \prod_{t \geq i > k} \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_{i-1}} = \prod_{t \geq i > k} \mathbf{W}_{\mathsf{h}^{T_{cc}}}^T diag(\sigma'(\mathbf{x}_{i-1})) \end{split}$$

 $W_h \sim small \implies Vanishing$ $W_{h} \sim large \implies Exploding$ E.g.: activation function = ReLU

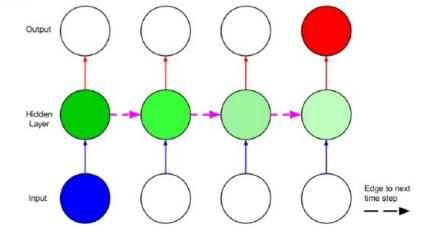
$$\frac{\partial x_i}{\partial x_{i-1}} = w_h$$

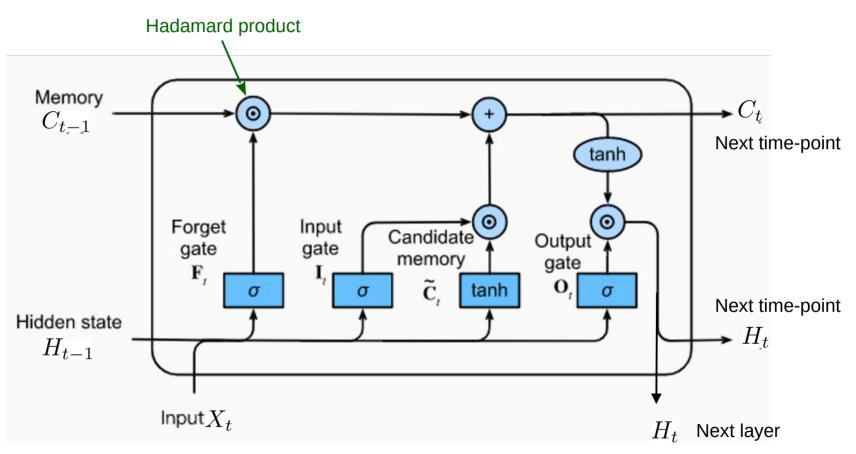
$$w_h = 0.1; \frac{\partial x_{10}}{\partial x_1} = w_h^{10} = 0.0000000001$$

$$w_h = 10; \frac{\partial x_{10}}{\partial x_1} = w_h^{10} = 100000000000$$

Source: SuperDataScience

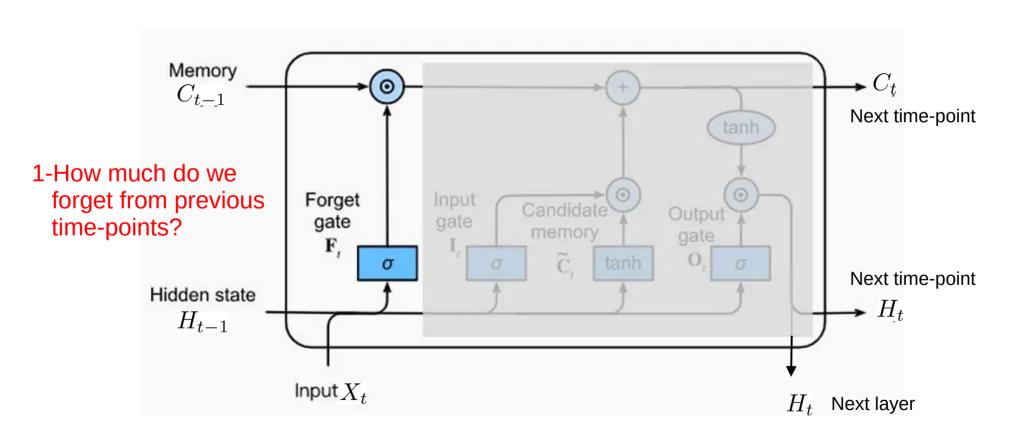
Green = sensitivity of output on input

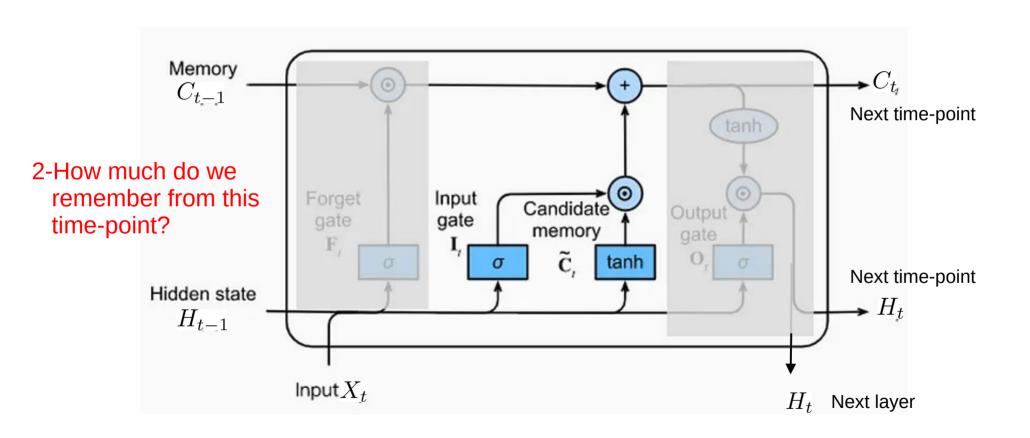


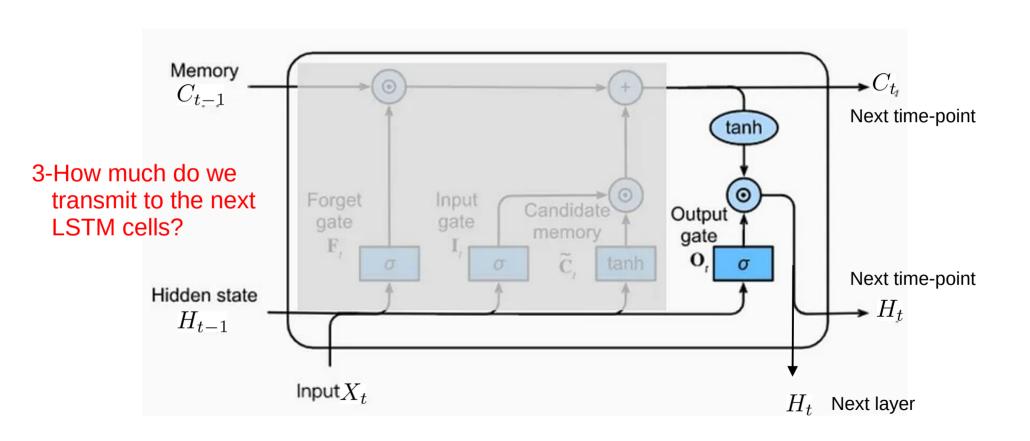


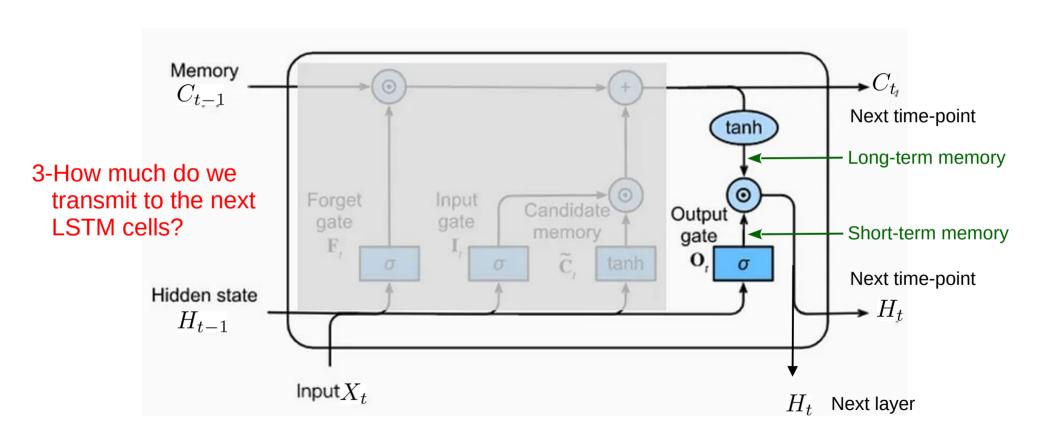
Hochreiter and Schmidhuber (1997) Long short-term memory. Neur Comput, 9(8):1735-1780

Source: Ottavio Calzone (2002) An Intuitive Explanation of LSTM. https://medium.com/@ottaviocalzone

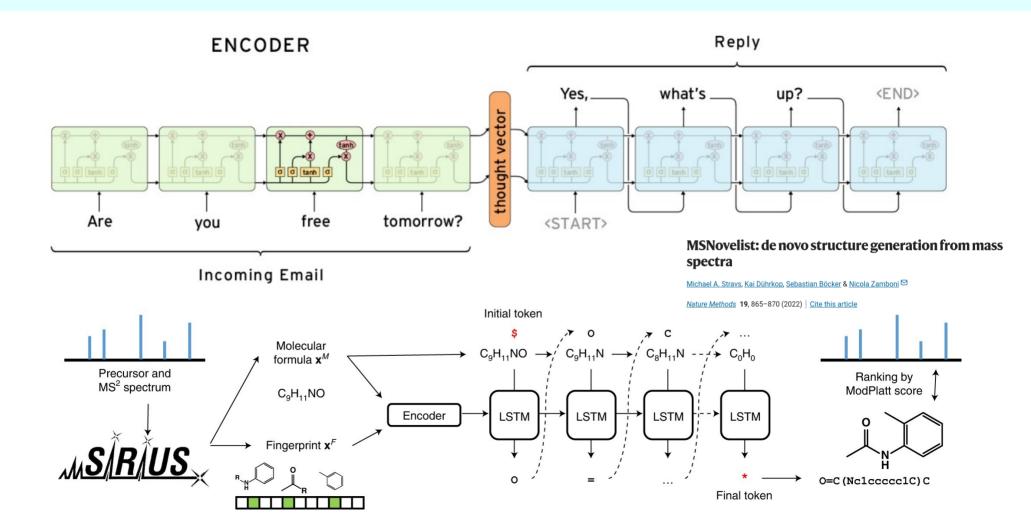








LSTMs for Encoder-Decoder



Example in bioinformatics

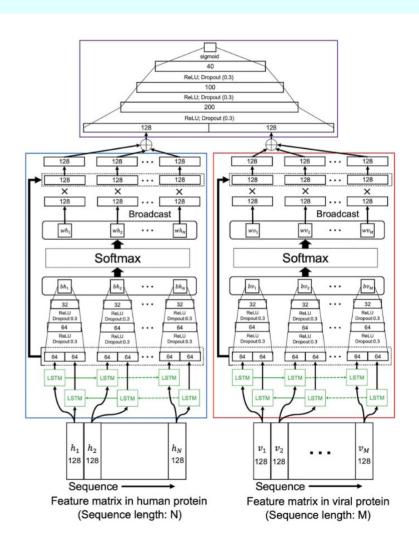


Briefings in Bioinformatics, 22(6), 2021, 1-9

https://doi.org/10.1093/bib/bbab228 Problem Solving Protocol

LSTM-PHV: prediction of human-virus protein-protein interactions by LSTM with word2vec

Sho Tsukiyama, Md Mehedi Hasan, Satoshi Fujii and Hiroyuki Kurata



Example in clinical setting



Contents lists available at ScienceDirect

International Journal of Infectious Diseases



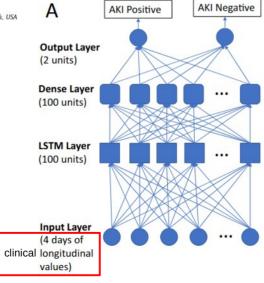
journal homepage: www.elsevier.com/locate/ijid

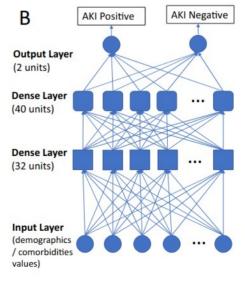
Long-short-term memory machine learning of longitudinal clinical data accurately predicts acute kidney injury onset in COVID-19: a two-center study

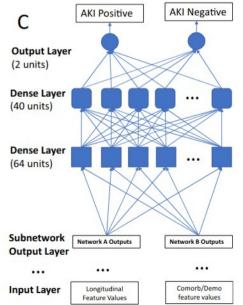


Justin Y. Lu, Joanna Zhu, Jocelyn Zhu, Tim Q Duong*

Department of Radiology, Montefiore Medical Center, Albert Einstein College of Medicine, New York, USA







Attention Is All You Need

Ashish Vaswani' Google Brain avaswani@google.com

Noam Shazeer* Google Brain noam@google.com

Niki Parmar* Google Research nikip@google.com

Jakob Uszkoreit* Google Research usz@google.com

Llion Jones+ Google Research llion@google.com

Aidan N. Gomez* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser* Google Brain lukaszkaiser@google.com

Illia Polosukhin* illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

1 Introduction

Recurrent neural networks, long short-term memory [12] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and transduction problems such as language modeling and machine translation [29] [2]. Numerous efforts have since continued to push the boundaries of recurrent language models and encoder-decoder architectures [31] [21, 13].

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA

The paper that changed everything: the Transfomer

^{*}Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled dot-product attention, multi-head attention and the parameter-free position representation and became the other person involved in nearly every detail. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor2tensor. Llion also experimented with novel model variants, was responsible for our initial codebase, and efficient inference and visualizations. Lukasz and Aidan spent countless long days designing various parts of and implementing tensor2tensor, replacing our earlier codebase, greatly improving results and massively accelerating our research.

Work performed while at Google Brain.

[‡]Work performed while at Google Research.

Attention Is All You Need

Cool title

Ashish Vaswani* Google Brain avaswani@google.com Noam Shazeer* Google Brain noam@google.com Niki Parmar* Google Research Jakob Uszkoreit* Google Research

nikip@google.com usz@google.com

Llion Jones* Google Research llion@google.com Aidan N. Gomez* † University of Toronto aidan@cs.toronto.edu Łukasz Kaiser* Google Brain lukaszkaiser@google.com

Illia Polosukhin* †
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

All authors equal

1 Introduction

Recurrent neural networks, long short-term memory [12] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and

*Equal contribution. Listing order is random.

Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled dot-product attention, multi-head attention and the parameter-free position representation and became the other person involved in nearly every detail. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor/tensor. Llion also experimented with novel model variants, was responsible for our initial codebase, and efficient inference and visualizations. Lukasz and Aidan spent countless long days designing various parts of and implementing tensor/tensor, replacing our earlier codebase, greatly improving results and massively accelerating

†Work performed while at Google Brain. ‡Work performed while at Google Research Never published in a journal

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

The paper that changed everything: the Transfomer

Cited... 138640 times as of 27 October 2024!

Attention Is All You Need

Ashish Vaswani' Google Brain avaswani@google.com

Noam Shazeer* Google Brain noam@google.com

Niki Parmar' Google Research

Jakob Uszkoreit* Google Research nikip@google.com usz@google.com

Llion Jones* Google Research llion@google.com

Aidan N. Gomez* † University of Toronto aidan@cs.toronto.edu Łukasz Kaiser* Google Brain

lukaszkaiser@google.com

Illia Polosukhin*

illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 Englishto-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

1 Introduction

Recurrent neural networks, long short-term memory [12] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and transduction problems such as language modeling and machine translation [29] [2] [5]. Numerous efforts have since continued to push the boundaries of recurrent language models and encoder-decoder architectures [31] [21, 13].

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA

The paper that changed everything: the Transfomer







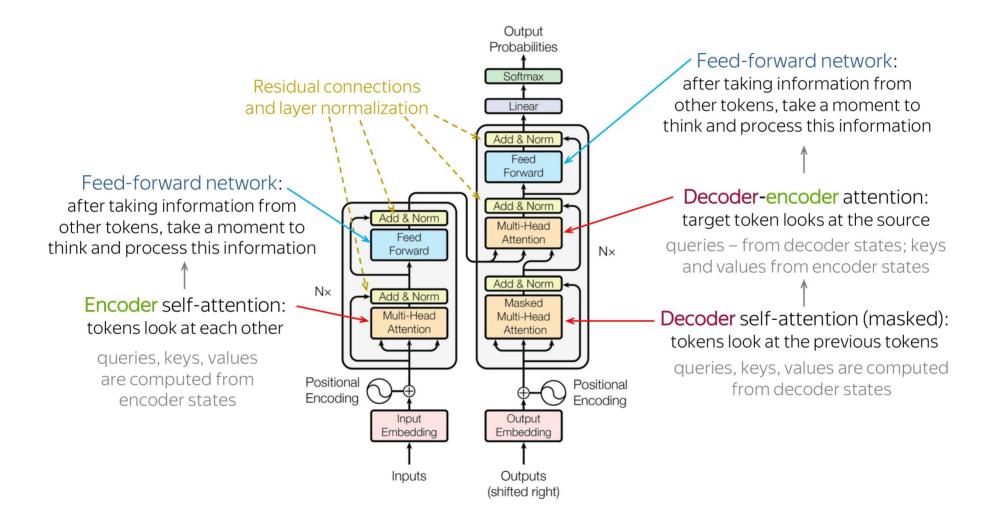


^{*}Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled dot-product attention, multi-head attention and the parameter-free position representation and became the other person involved in nearly every detail. Niki designed, implemented, tuned and evaluated countless model variants in our original codebase and tensor2tensor. Llion also experimented with novel model variants, was responsible for our initial codebase, and efficient inference and visualizations. Lukasz and Aidan spent countless long days designing various parts of and implementing tensor2tensor, replacing our earlier codebase, greatly improving results and massively accelerating our research.

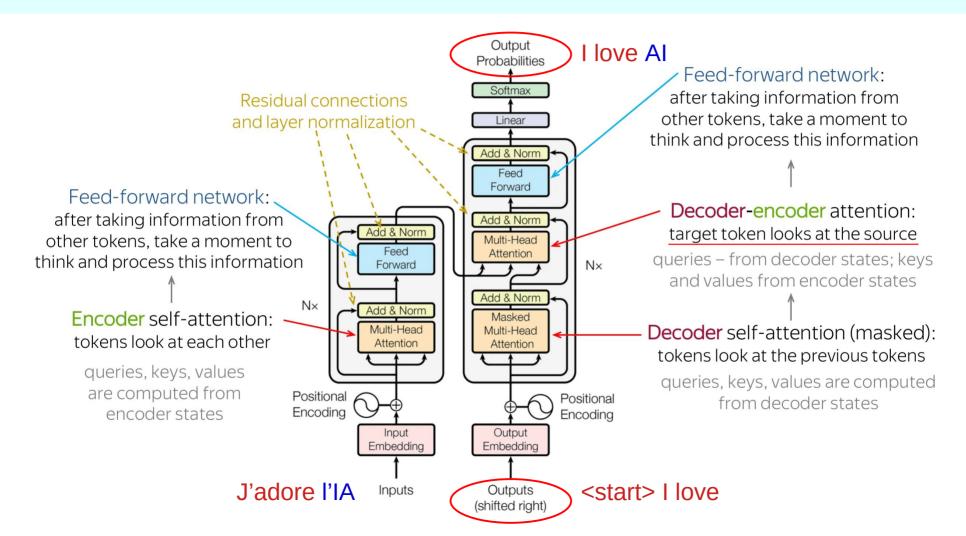
Work performed while at Google Brain.

Work performed while at Google Research.

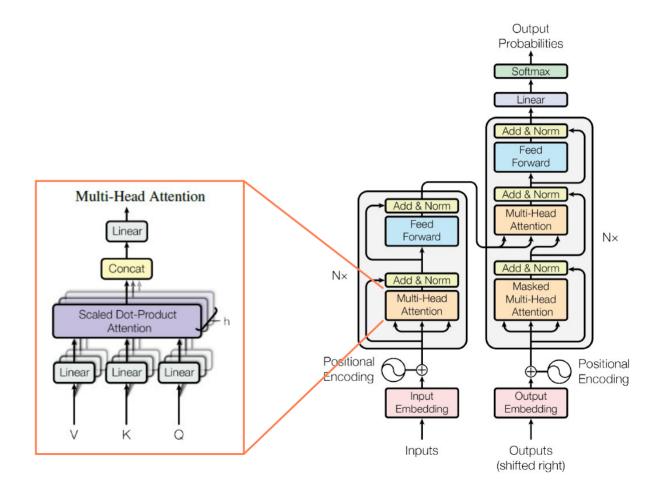
The Transformer: Memory + context = attention



The Transformer: Memory + context = attention

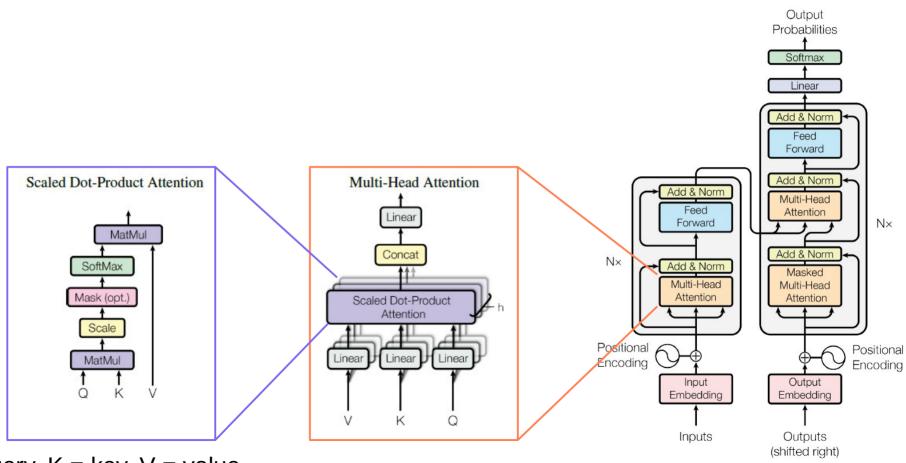


Attention in the Transformer



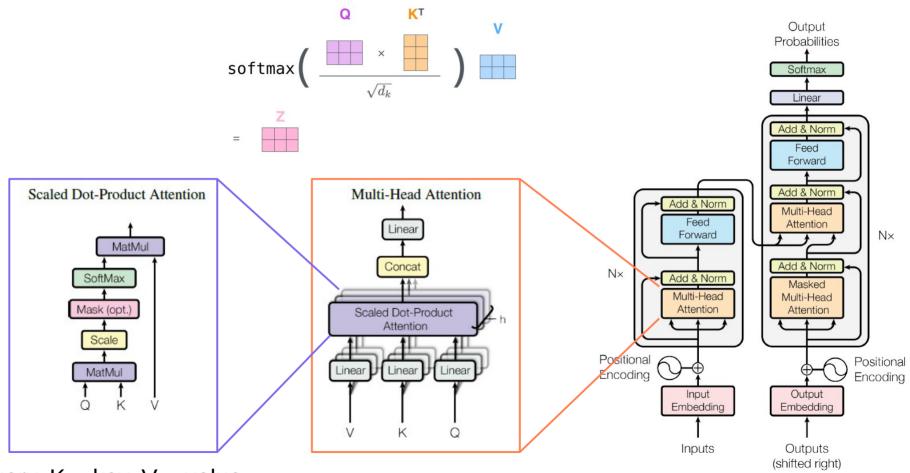
Q = query, K = key, V = value

Attention in the Transformer



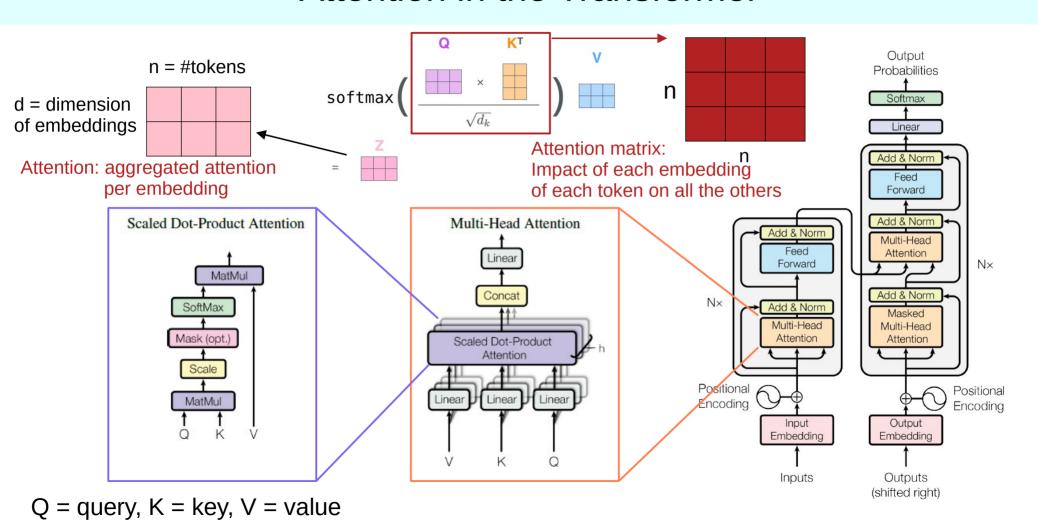
Q = query, K = key, V = value

Attention in the Transformer

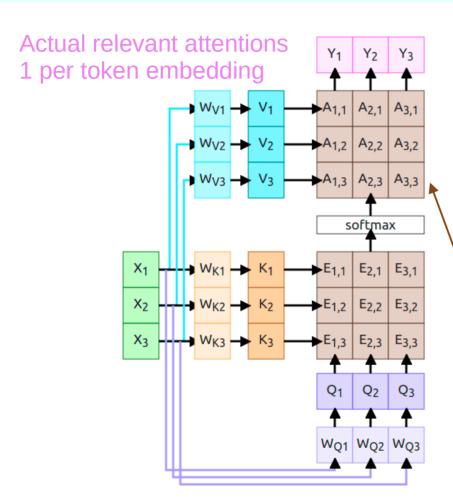


Q = query, K = key, V = value

Attention in the Transformer



Attention in the Transformer



Key matrix: W_K (Shape $D_X \times D_Q$)

Value matrix: W_V (Shape $D_X imes D_V$)

Query matrix: $W_{\mathbb{Q}}$ (Shape $D_X imes D_Q$)

Query vectors: $Q = XW_Q$ (Shape $N_X \times D_Q$)

Key vectors: $oldsymbol{K} = oldsymbol{X} W_K$ (Shape $N_X imes D_Q$)

Value vectors: $V = XW_V$ (Shape $N_X \times D_V$)

Similarity function: scaled dot product

Similarities: $E = Q {m K}^T$ (shape $N_X imes N_X$), $E_{i,j} = Q_i \cdot {m K}_j / \sqrt{D_Q}$

Attention weights: $A = \operatorname{softmax}(E, dim = 1)$ (shape $N_X \times N_X$)

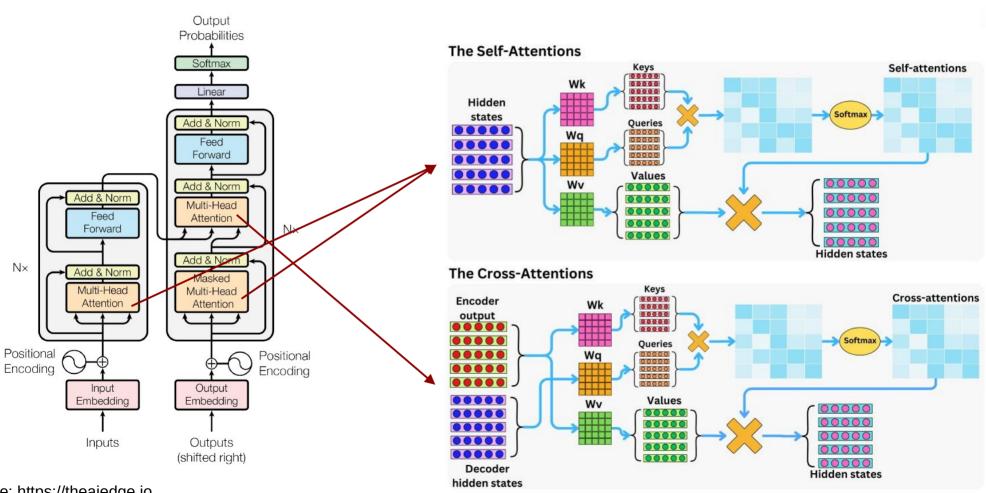
Output: Y = AV (shape $N_X imes D_V$) where $Y_i = \sum_j (A_{i,j}, V_j)$

The attention of all token embeddings on all token embeddings

X is entered W_Q , W_K , and W_V are learned Everything else is computed

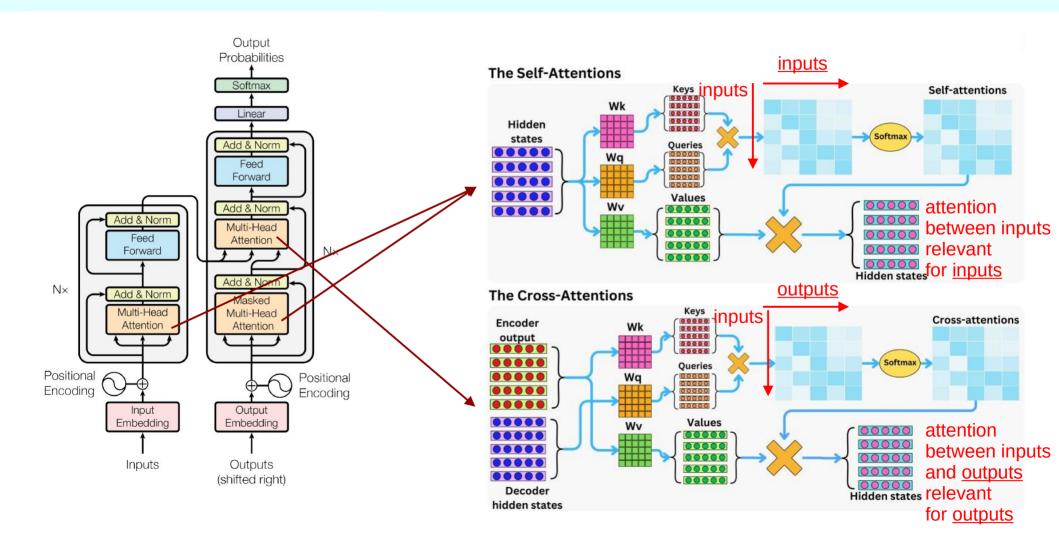
Source: https://erdem.pl/2021/05/introduction-to-attention-mechanism

Self versus cross-attention

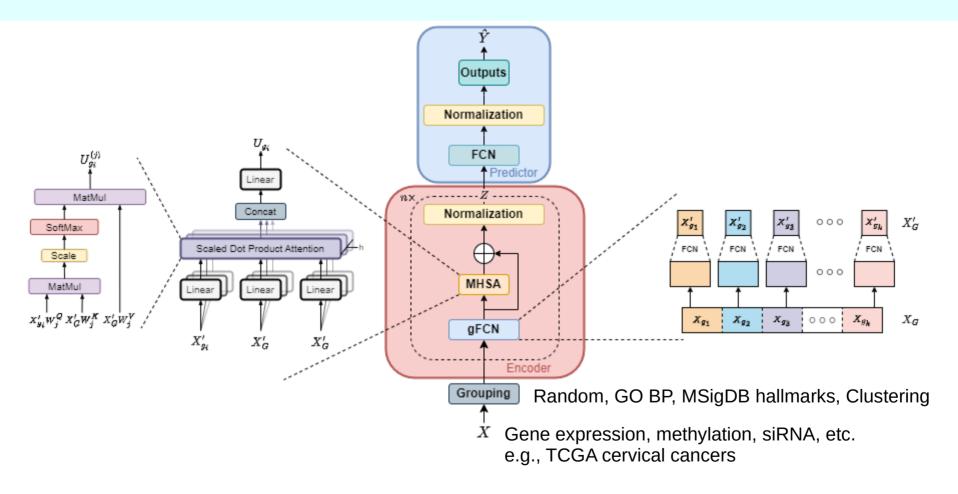


Source: https://theaiedge.io

Self versus cross-attention

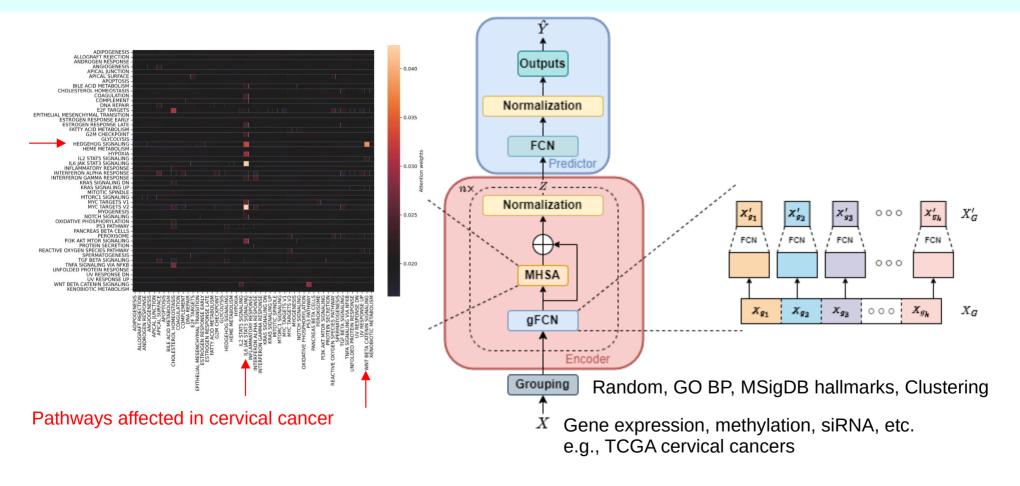


AttOmics



Beaude, A., Rafiee Vahid, M., Augé, F., Zehraoui, F., & Hanczar, B. (2023). AttOmics: attention-based architecture for diagnosis and prognosis from omics data. *Bioinformatics*, 39(Supplement 1), i94-i102.

AttOmics



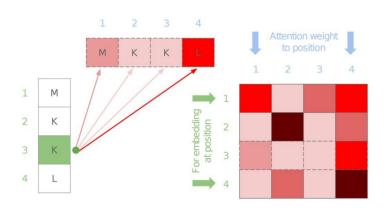
Beaude, A., Rafiee Vahid, M., Augé, F., Zehraoui, F., & Hanczar, B. (2023). AttOmics: attention-based architecture for diagnosis and prognosis from omics data. *Bioinformatics*, 39(Supplement 1), i94-i102.

EnzBERT

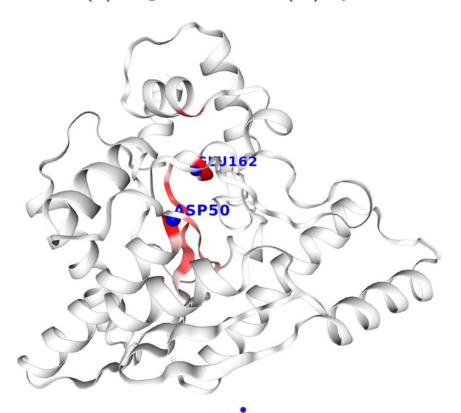
Predicting enzymatic function of protein sequences with attention 8

Nicolas Buton ™, François Coste, Yann Le Cunff

Bioinformatics, Volume 39, Issue 10, October 2023, btad620, https://doi.org/10.1093/bioinformatics/btad620



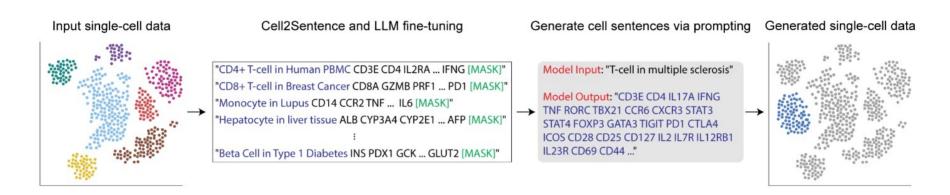
Nh(3)-dependent nad(+) synthetase



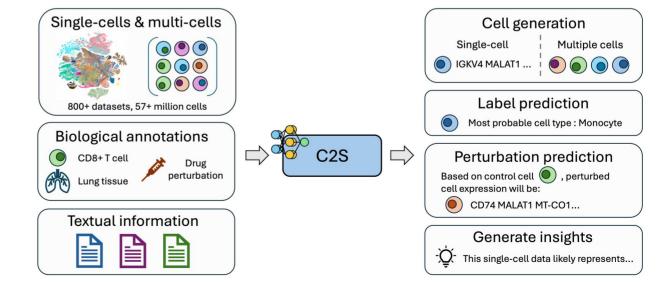
Aggregated attention for each token (amino acid)

- 0 MSMQEKIMRE LHVKPSIDPK QEIEDRVNFL KQYVKKTGAK GFVLGISIMO DSTLAGRLAQ LAVESIREEG GDAQFIAVRL PHGTQQDEDD AQLALKFIKP
- 1 DKSWKFDIKS TVSAFSDQYQ QETGDQLTDF NKGNVKARTR MIAQYAIGGQ EGLLVLO DI ALAVTGFFT KYGDGGADLL PLTGLTKRQG RTLLKELGAP
- 2 ERLYLKEPTA DLLDEKPOOS DETELGISHD EIDDYLEGKE VSAKVSEALE KRYSMTEHKR QVPASMFDDW WK

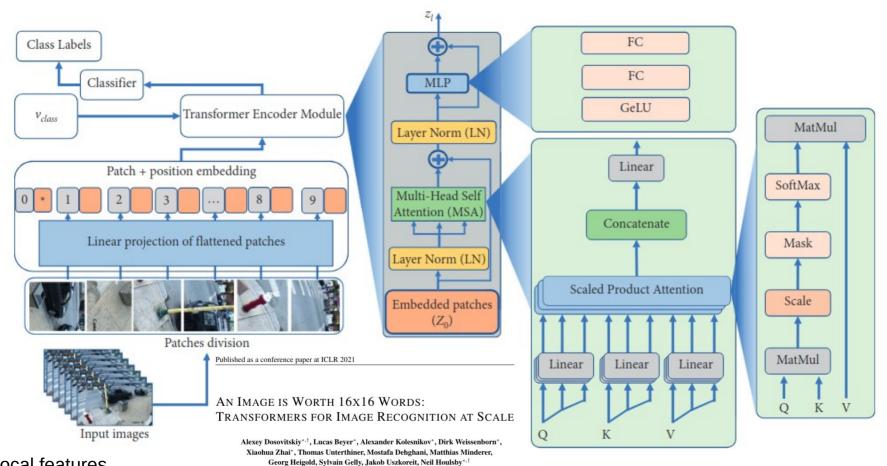
Cell2Sentence



Levine *et al* (2024). Cell2Sentence: Teaching Large Language Models the Language of Biology. *BioRxiv* https://doi.org/10.1101/2023.09.11.557287



Vision Transformer

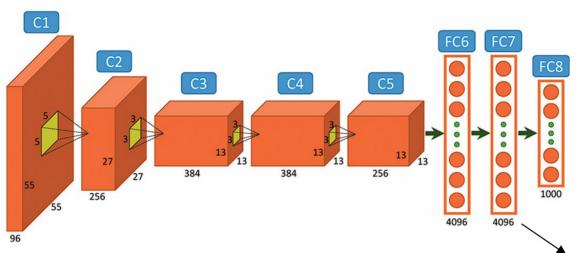


CNN = local features ViT = relations between distant features

*equal technical contribution, *equal advising Google Research, Brain Team {adosovitskiy, neilhoulsby}@google.com

source: https://doi.org/10.1155/2022/3454167

Patches are embedded by CNNs



6 nearest neighbours in the 4096 dimension space

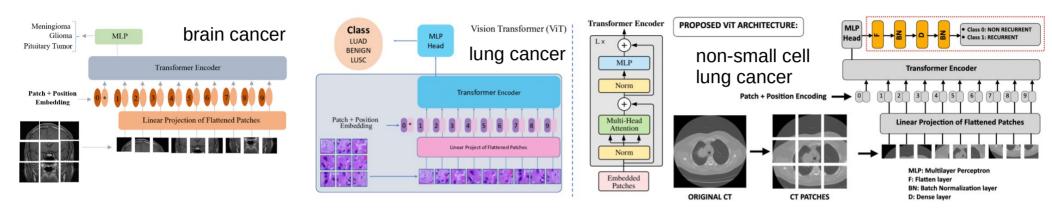
Krizhevsky A, Sutskever I, Hinton GE (2012) ImageNet Classification with Deep Convolutional Neural Networks https://proceedings.neurips.cc/paper/2012/file/ c399862d3b9d6b76c8436e924a68c45b-Paper.pdf

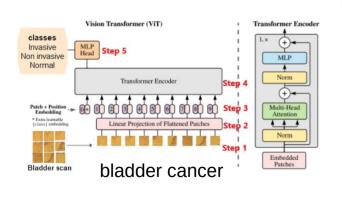
(presenting AlexNet, the first Deep Convolutional Network)

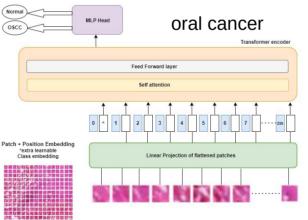
input image

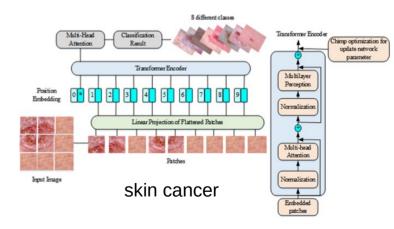


ViTs are replacing vanilla CNNs



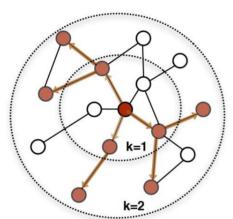




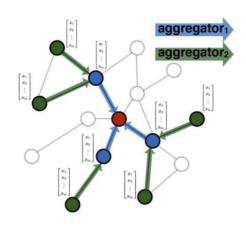


Graph Neural Networks (GNNs)

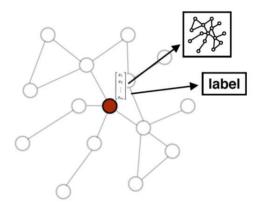
source: https://blogs.nvidia.com/blog/what-are-graph-neural-networks/



1. Sample neighborhood



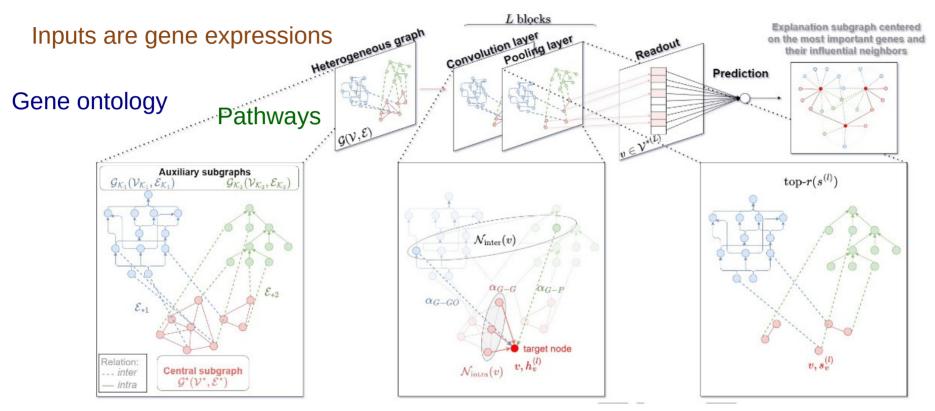
2. Aggregate feature information from neighbors



Predict graph context and label using aggregated information

IEEE TRANSACTIONS ON NEURAL NETWORKS, VOL. 20, NO. 1, JANUARY 2009

GNN can be heterogeneous

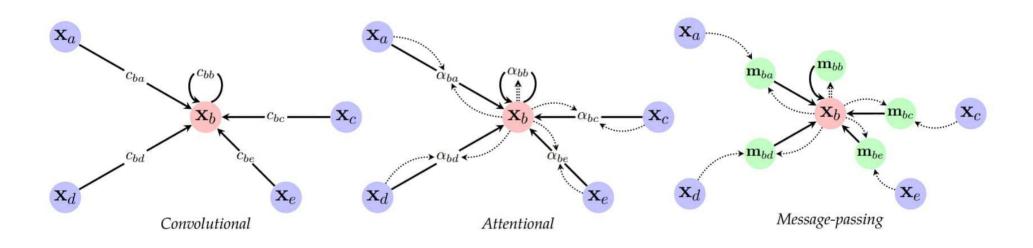


BioHAN: a Knowledge-based Heterogeneous Graph Neural Network for precision medicine on transcriptomic data

https://hal.science/hal-04092210/

Many different ways to update GNNs

source: https://blogs.nvidia.com/blog/what-are-graph-neural-networks/



$$\begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \end{bmatrix} A$$

$$\begin{bmatrix} 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \end{bmatrix}$$

Laplacian matrix

L = D - A

NB: undirected graph → all matrices are symmetric This could be different for a directed graph

Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering

Michaël Defferrard

Xavier Bresson

Pierre Vandergheynst

$$\begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

$$L =$$

$$D = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{array}{c} A \\ B \\ C \\ D \end{array} \quad A = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} \quad L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix}$$

Degree matrix

Adjacency matrix

Laplacian matrix

L = D - A

identity matrix no influence from neighbours

influence from influence from neighbours and immediate neighbours neighbours of neighbours

Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering

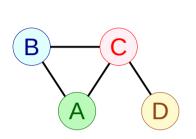
Michaël Defferrard Xavier Bresson Pierre Vanderghevnst

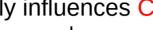
EPFL, Lausanne, Switzerland {michael.defferrard,xavier.bresson,pierre.vandergheynst}@epfl.ch $x' = p_w(L)x$

 $w = |w_0, w_1, w_2, ..., w_k|$

 $p_w(L) = w_0 I + w_1 L + w_2 L^2 + \dots + w_k L^k$

kernel de convolution



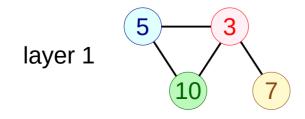


$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix}$$

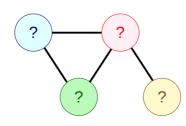
D influences A, B, C

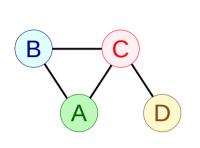
$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix} \qquad L^2 = \begin{bmatrix} 6 & -3 & -4 & 1 \\ -3 & 6 & -4 & 1 \\ -4 & -4 & 12 & -4 \\ 1 & 1 & -4 & 2 \end{bmatrix}$$

$$w = [1, 0.1, 0.01]$$



layer 2





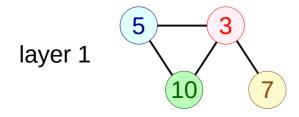


$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix} \qquad L^2 = \begin{bmatrix} 6 & -3 & -4 & 1 \\ -3 & 6 & -4 & 1 \\ -4 & -4 & 12 & -4 \\ 1 & 1 & -4 & 2 \end{bmatrix}$$

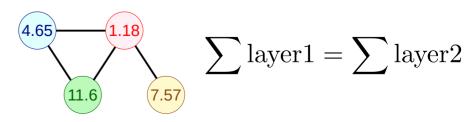
$$-3$$
 -4 1

$$L^{2} = \begin{bmatrix} -3 & 6 & -4 & 1 \\ -4 & -4 & 12 & -4 \\ 1 & 1 & -4 & 2 \end{bmatrix}$$

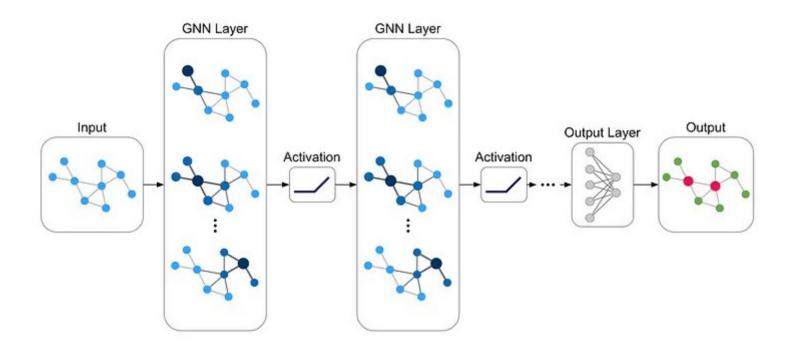
$$w = \begin{bmatrix} 1, 0.1, 0.01 \end{bmatrix} \quad 1 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 10 \\ 5 \\ 3 \\ 7 \end{bmatrix} + 0.1 \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 10 \\ 5 \\ 3 \\ 7 \end{bmatrix} + 0.01 \begin{bmatrix} 6 & -3 & -4 & 1 \\ -3 & 6 & -4 & 1 \\ -4 & -4 & 12 & -4 \\ 1 & 1 & -4 & 2 \end{bmatrix} \begin{bmatrix} 10 \\ 5 \\ 3 \\ 7 \end{bmatrix}$$



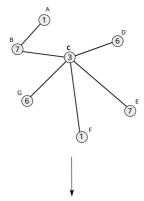
layer 2

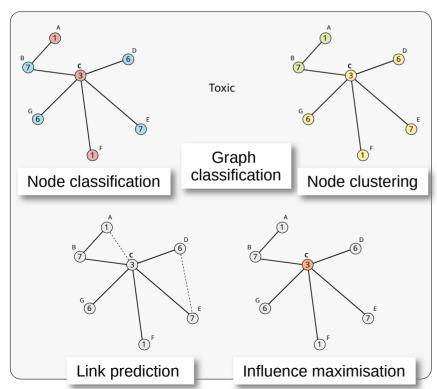


Graph Neural Networks (GNNs)



NB: GNNs generally comprise 3 embeddings that are updated at each iteration, i.e nodes (vertices), edges, and graph





GNNs: What can we do?

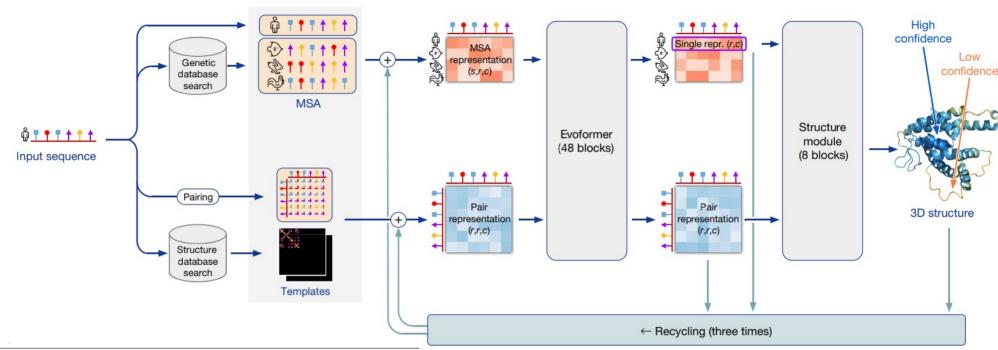
Source: Understanding Convolutions on Graphs https://distill.pub/2021/understanding-gnns/

See also: A Gentle Introduction to Graph Neural Networks https://distill.pub/2021/gnn-intro/

Both by Google Research teams



AlphaFold2



Highly accurate protein structure prediction with AlphaFold

https://doi.org/10.1038/s41586-021-03819-2

Received: 11 May 2021

Accepted: 12 July 2021

Published online: 15 July 2021

Open access

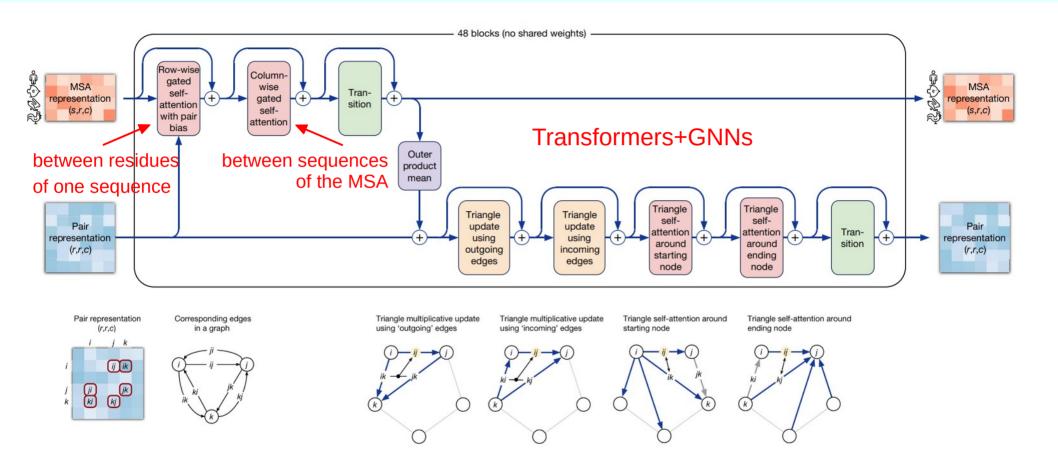
Check for updates

John Jumper¹^{4,83}, Richard Evans¹⁴, Alexander Pritzel¹⁴, Tim Green¹⁴, Michael Figurnov¹⁴, Algus Bates¹⁴, Augustin Židek¹⁴, Anna Potapenko¹⁴, Alexe Bridgland¹⁴, Clemens Meyer¹⁴, Simon A. A. Kohl¹⁴, Rishub Jain¹⁴, Andrew Cowie¹⁴, Bernain Romera-Paredes¹⁴, Stanislav Nikolov¹⁴, Rishub Jain¹⁴, Jonas Adler¹, Trevor Back¹, Stig Petersen¹, David Reiman¹, Ellen Clancy¹, Michal Zielinski¹, Martin Steinegger²³, Michaelina Pacholska³, Tamas Berghammer³, Sebastian Bodenstein¹, David Silver¹, Oriol Vinyals¹, Andrew W. Senior¹, Koray Kavukcuoglu¹, Pushmeet Kohli¹ & Demis Hassabis¹^{4,58}

~93 million parameters (weights+biases)

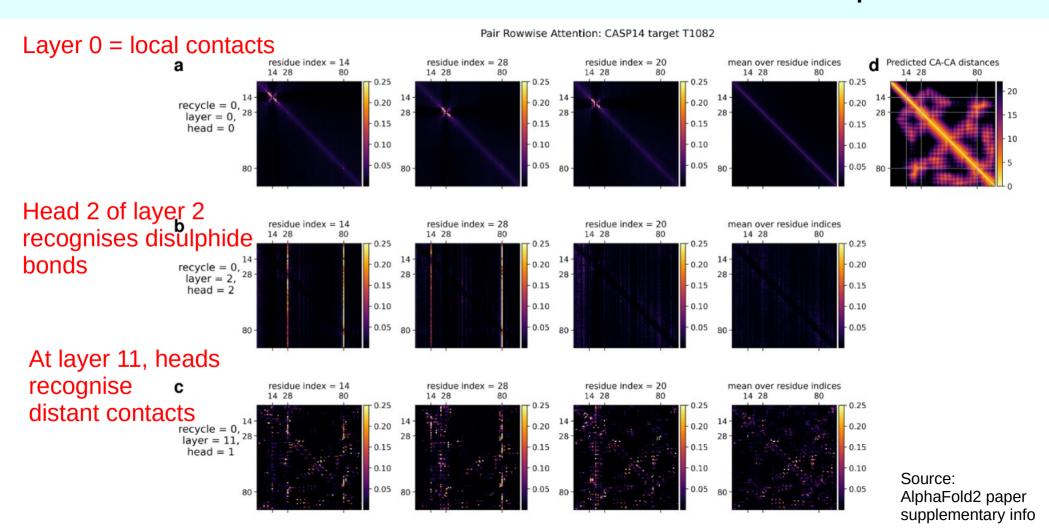
https://github.com/google-deepmind/alphafold

AlphaFold2: evoformer

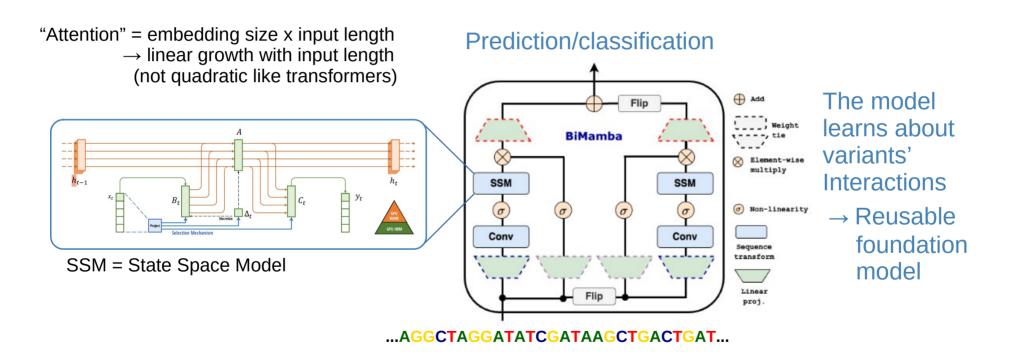


see also: https://www.blopig.com/blog/2021/07/alphafold-2-is-here-whats-behind-the-structure-prediction-miracle/

Row-wise attention: between residues of a sequence



RNNs are back. Rise of the Mamba



Gu and Dao (2023) *arXiv*:2312.00752: Schiff *et al* (2024) *arXiv*:2403.03234







