



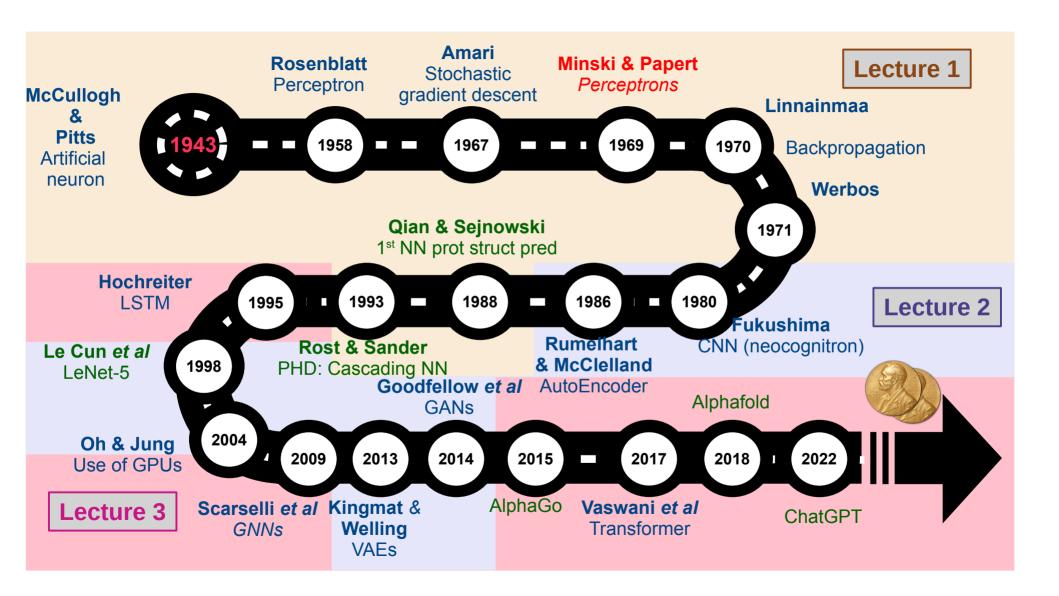
Beyond MLPs

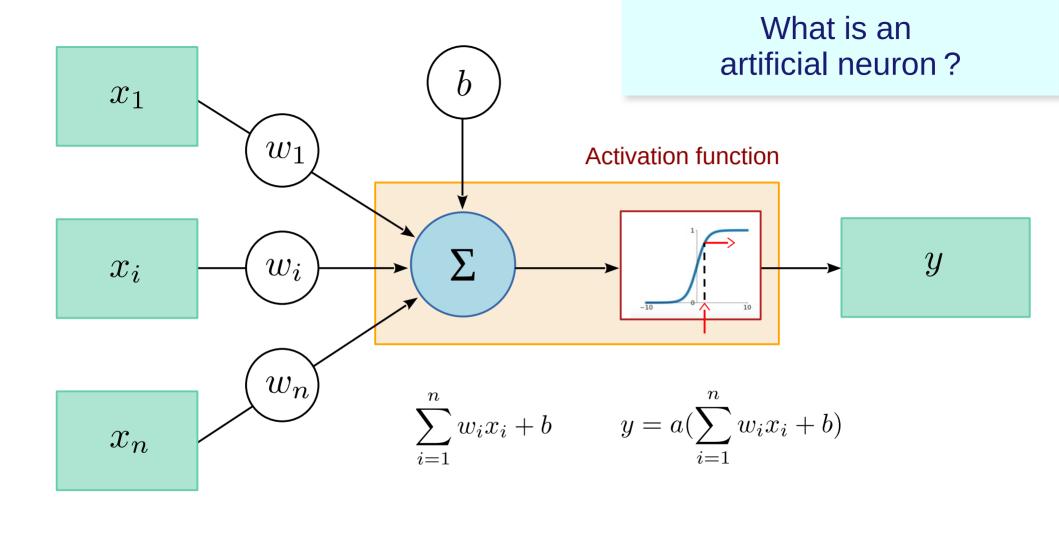
Part 2/2: RNNs, Attention and GNNs

nicolas.gambardella@univ-lille.fr

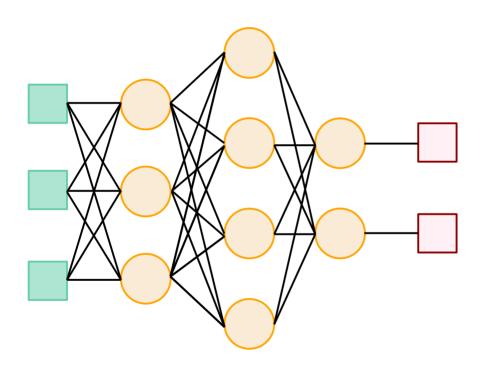




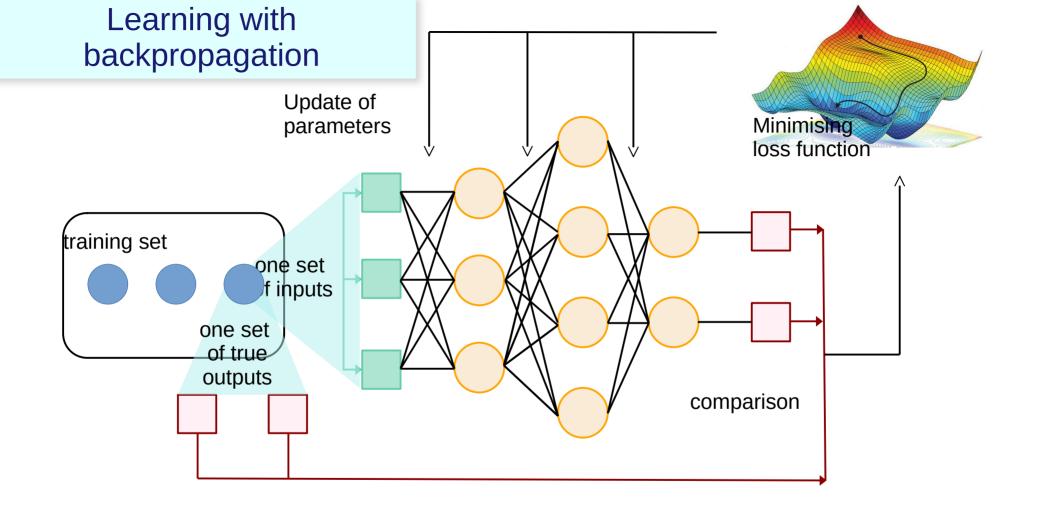




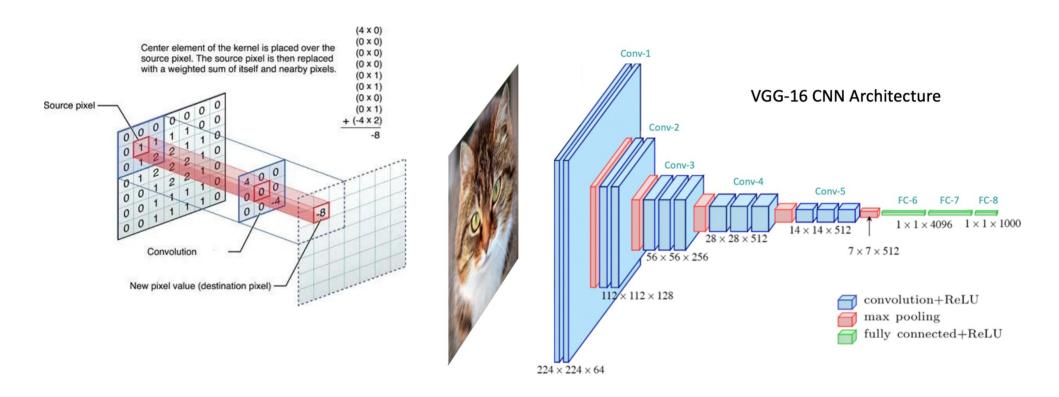
All inputs can be independent and everything connected to everything



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

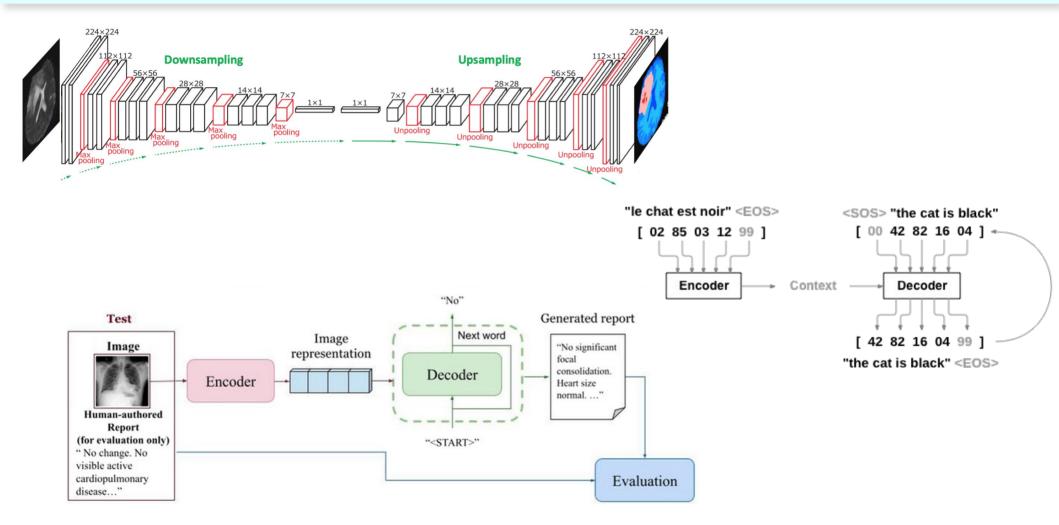


We can detect local features by linking neighbouring inputs

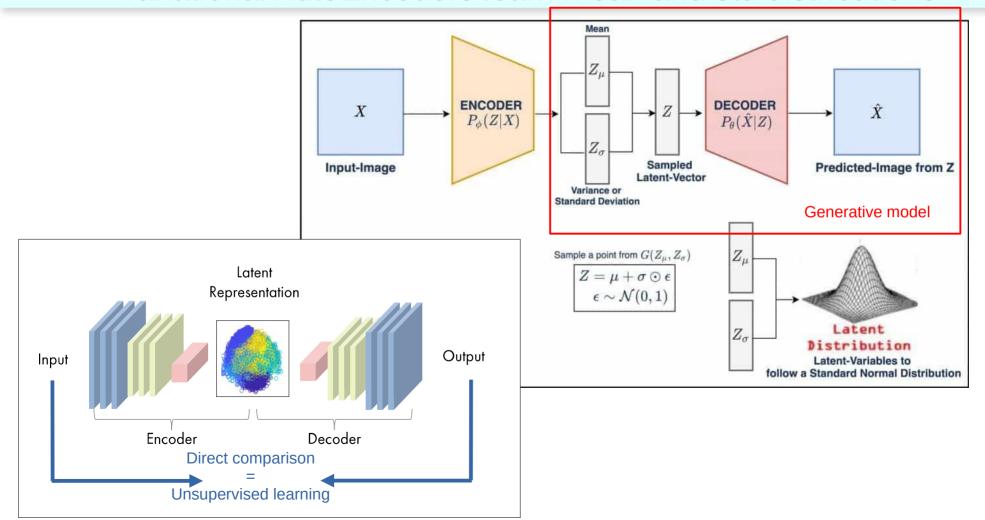


Deep Convolutional Neural Networks (CNN)

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



Blogs/Media



Search	Q
	Advanced Search

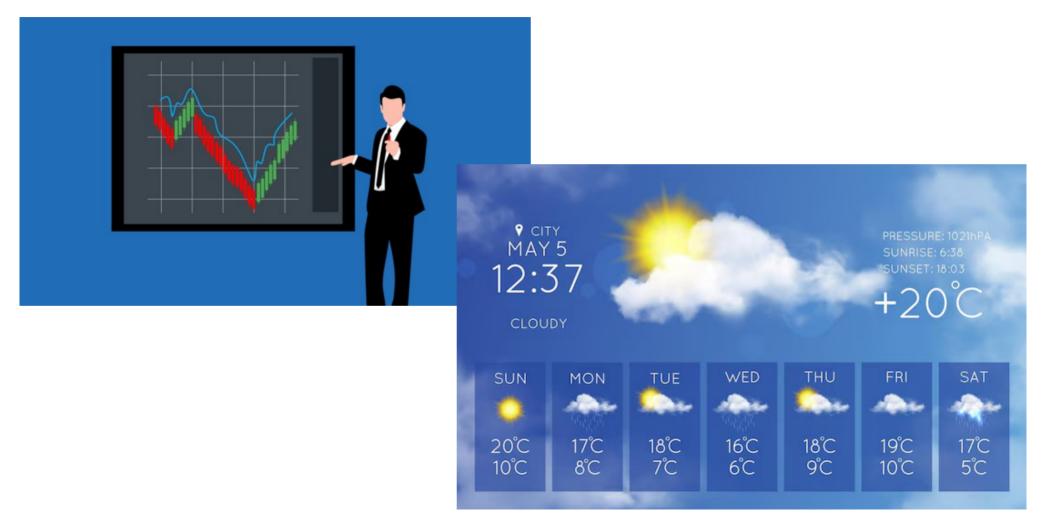
New Results ▲ Follow this preprint <a> ○ ○ Previous Next (Deep learning models reading clinical data and liver omics strongly distinguish Posted October 10, 2025. NASH from steatosis and suggest new genes involved in liver disease severity Download PDF ✓ Email Nicolas Gambardella, Smaïn Fettem, Mathilde Boissel, Lijiao Ning, Violeta Raverdy, A Share ▼ Print/Save Options Marwa Afnouch, Souhila Amanzougarene, Mehdi Derhourhi, Bénédicte Toussaint, Emmanuel Vaillant, Supplementary Material Citation Tools Amna Khamis, Philippe Lefebvre, @ Bart Staels, Francois Pattou, Philippe Froquel, Amelie Bonnefond R Get QR code doi: https://doi.org/10.1101/2025.10.10.681581 This article is a preprint and has not been certified by peer review [what does this mean?]. Subject Area Abstract Info/History Metrics Preview PDF Molecular Biology Abstract **Reviews and Context** Background & Aims Metabolic dysfunction-associated steatotic liver disease (MASLD) is a frequent co-morbidity of obesity and diabetes, with prevalence increasing worldwide. Comment 0 Recognising liver disease stages and elucidating the molecular underpinning of their TRIP Peer Reviews **Y** progression are thus medically important. Methods Using data gathered from 300 patients with Community Reviews ÷ obesity of the ABOS cohort, we selected non-redundant clinical variables, gene expressions and CpGs methylation levels most associated with severity using unsupervised approaches to **Automated Services** O₀

train a multi-module, multi-layer perceptron predicting patients liver status. Results The

combination of five models trained on the three modalities reached an AUC of 0.945 on a



Time series and sequences of variable lengths



Time series and sequences of variable lengths

Speech recognition

Music generation

Sentiment classification

DNA sequence analysis

Machine translation

Video activity recognition

Name entity recognition



"There is nothing to like in this movie."

AGCCCCTGTGAGGAACTAG ---

Voulez-vous chanter avec moi?



Yesterday, Harry Potter met Hermione Granger. "The quick brown fox jumped over the lazy dog."



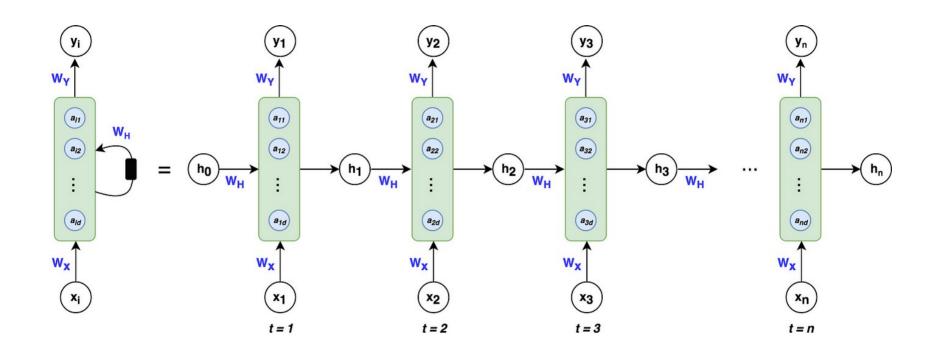
AGCCCCTGTGAGGAACTAG

Do you want to sing with me?

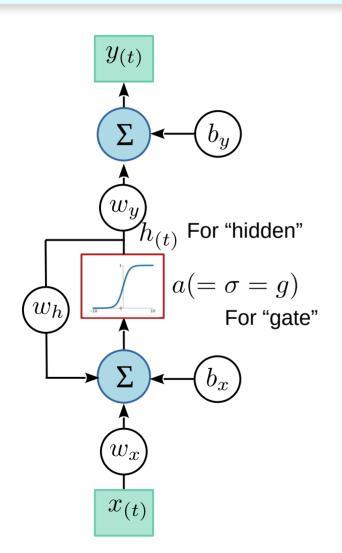
Running

Yesterday, Harry Potter met Hermione Granger. Andrew Ng

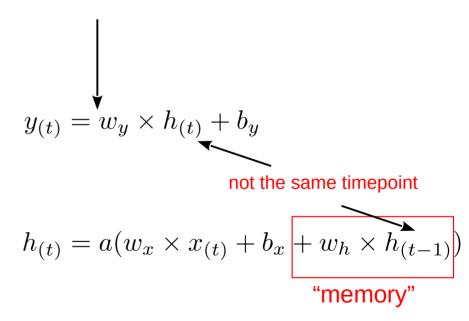
Recurrent Neural Networks: successive inputs are not independent



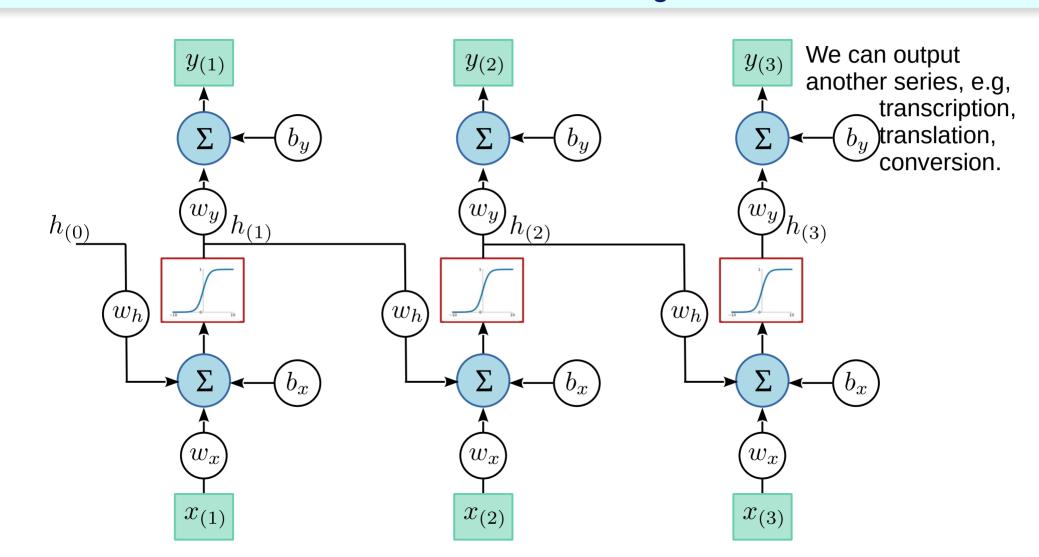
RNN: 1 cell (here, 1 neuron)



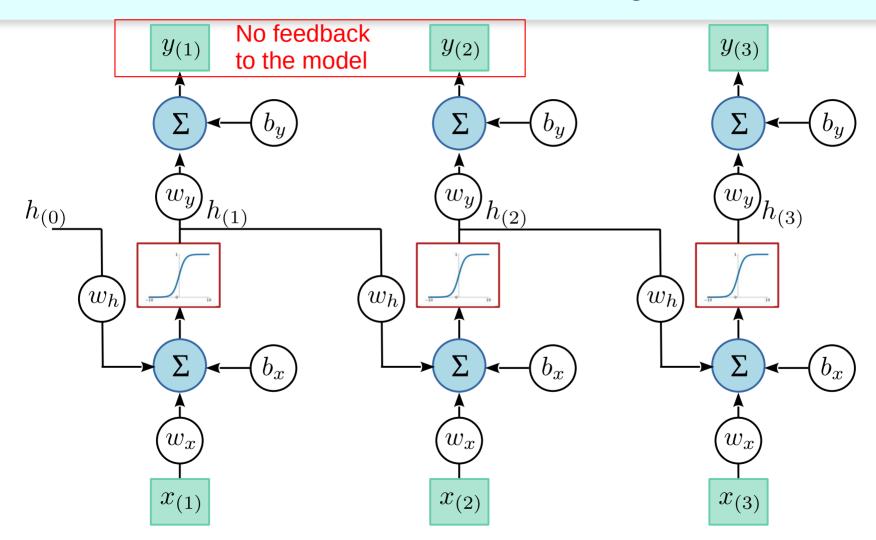
NB: implicit "identity" activation function



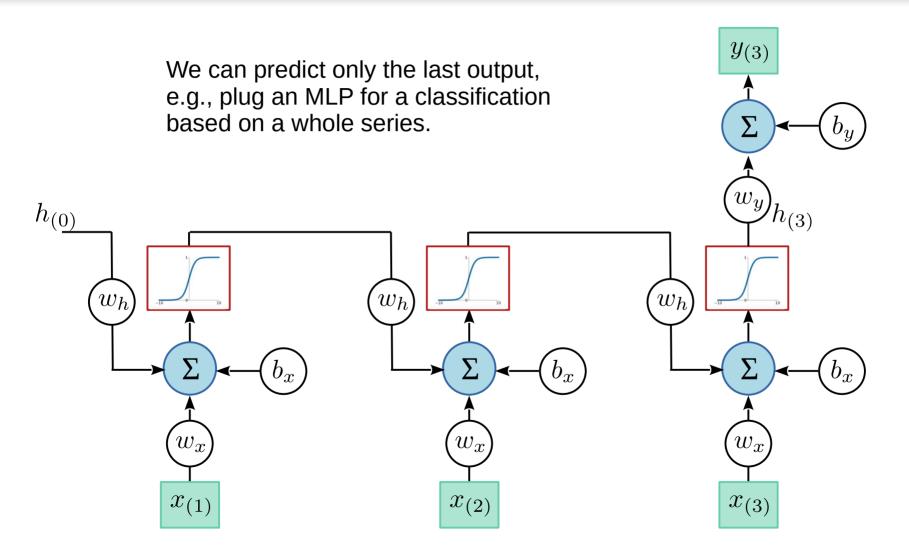
RNN: 1 cell - unfolding



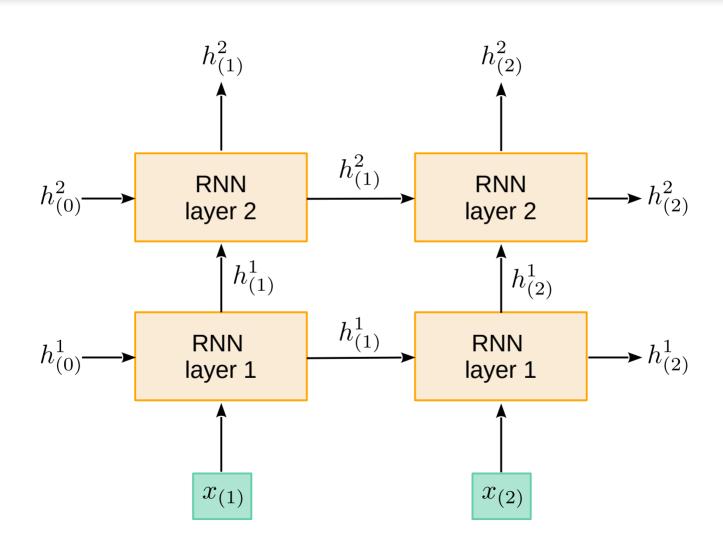
RNN: 1 cell - unfolding



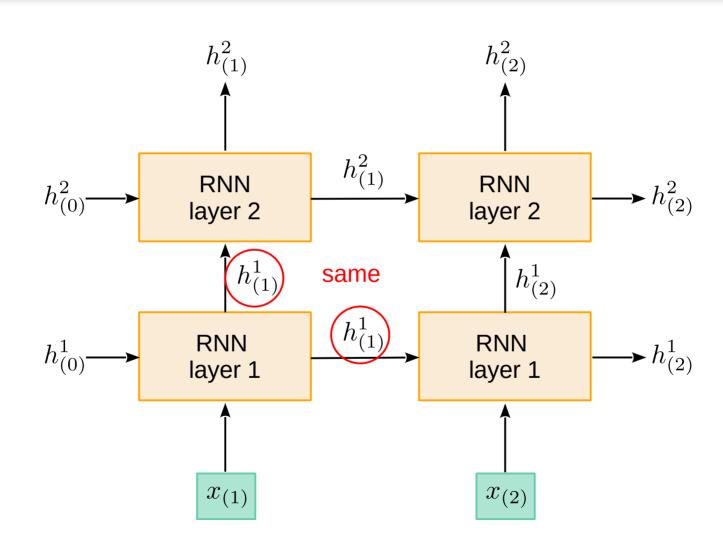
RNN: 1 cell - unfolding



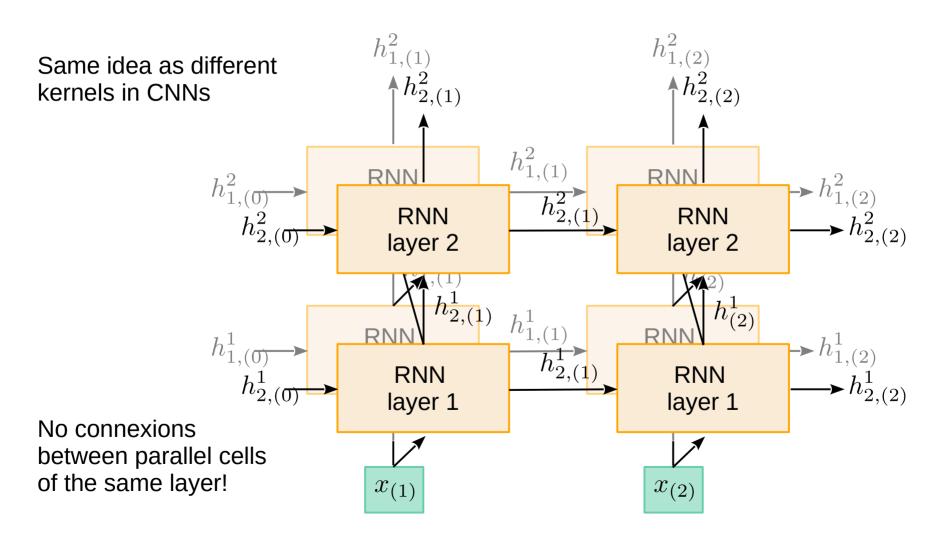
Stacked RNNs



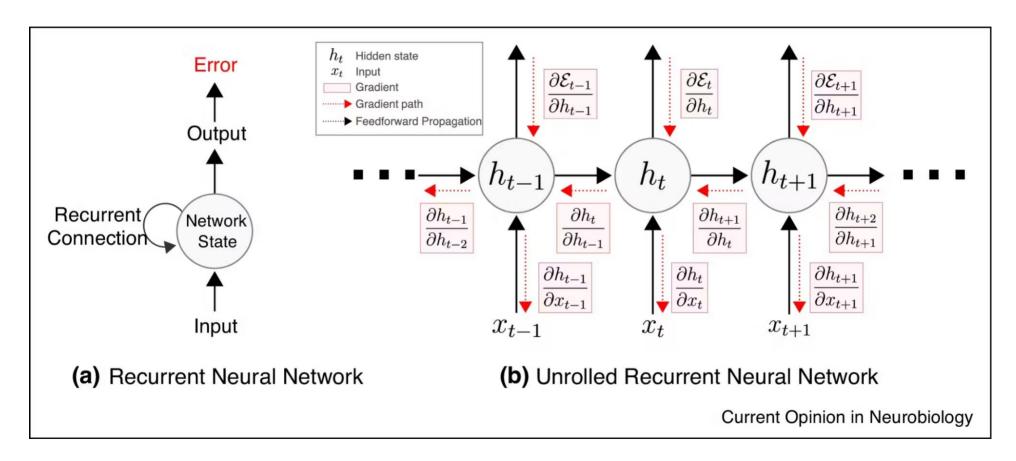
Stacked RNN



Several RNNs may learn different patterns in parallel

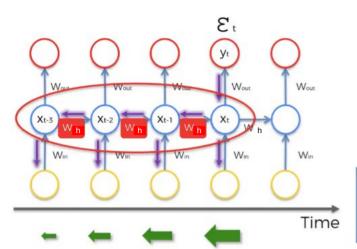


Training RNNs by backpropagation in time



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

Exploding and vanishing gradients



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

 $W_h \sim small \implies Vanishing$ $W_{h} \sim large \implies Exploding$

E.g.: activation function = ReLU

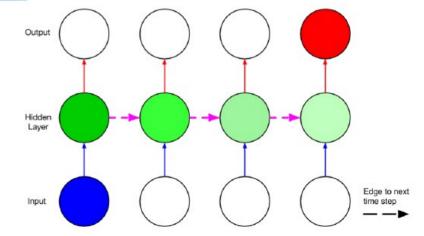
$$\frac{\partial x_i}{\partial x_{i-1}} = w_h$$

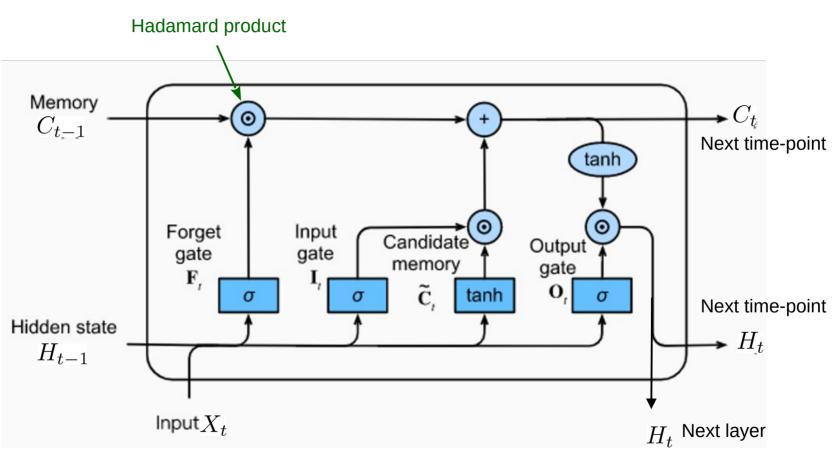
$$w_h = 0.1; \frac{\partial x_{10}}{\partial x_1} = w_h^{10} = 0.0000000001$$

$$w_h = 10; \frac{\partial x_{10}}{\partial x_1} = w_h^{10} = 10000000000$$

Source: SuperDataScience

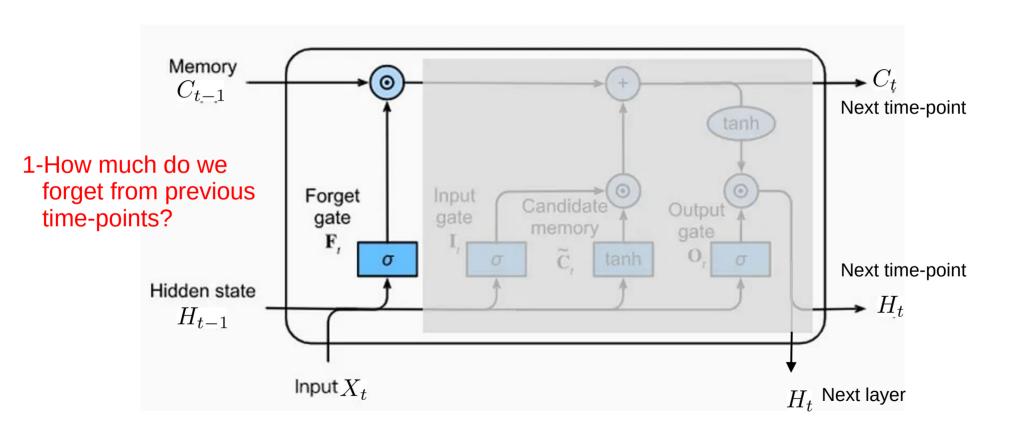
Green = sensitivity of output on input

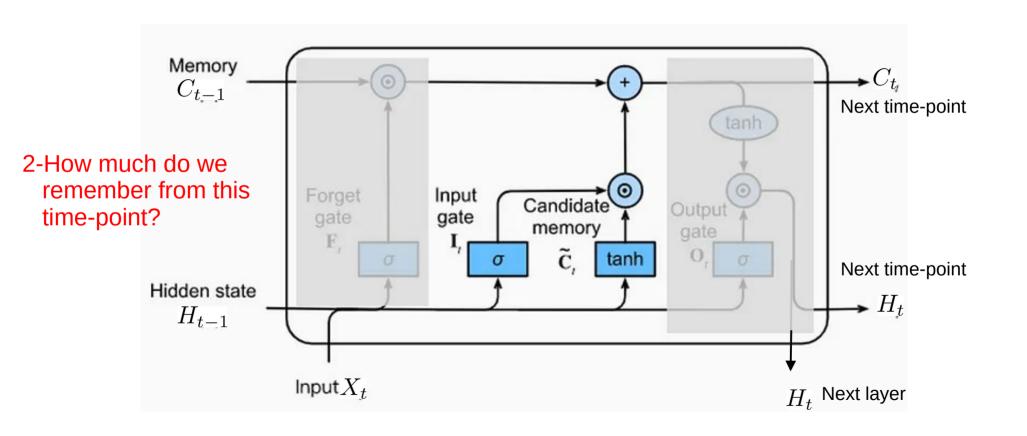


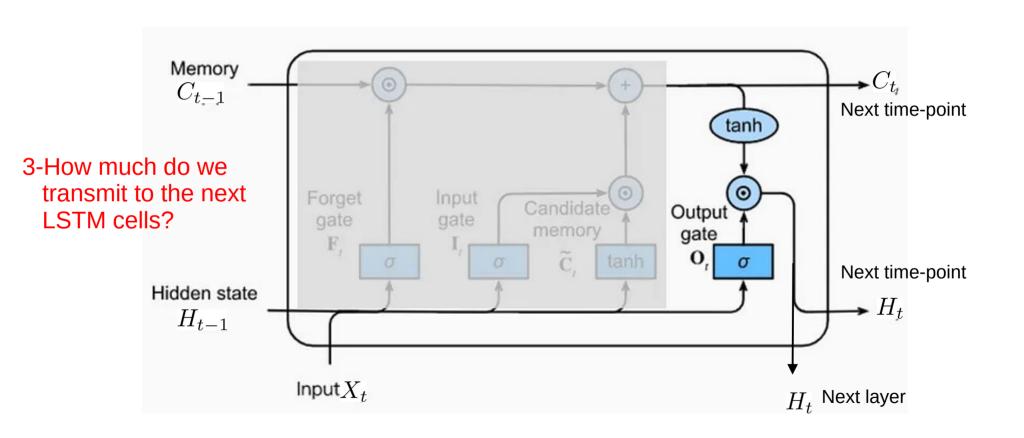


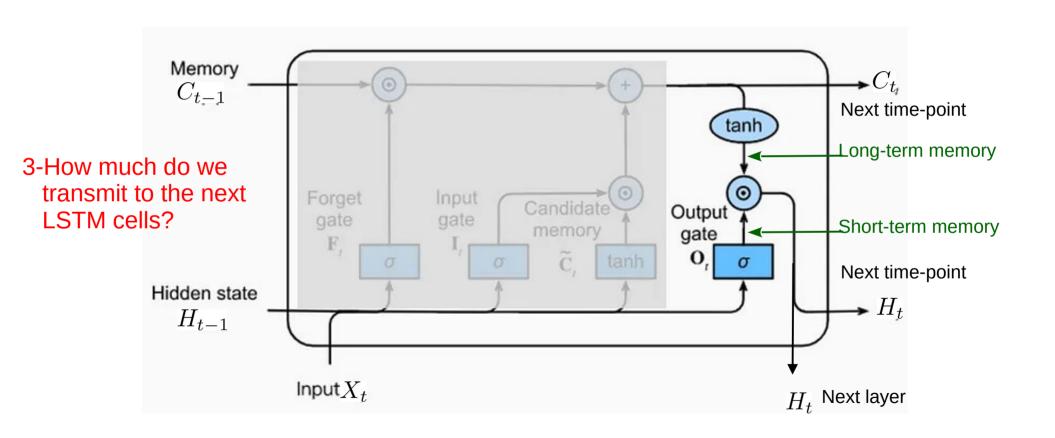
Hochreiter and Schmidhuber (1997) Long short-term memory. Neur Comput, 9(8):1735-1780

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

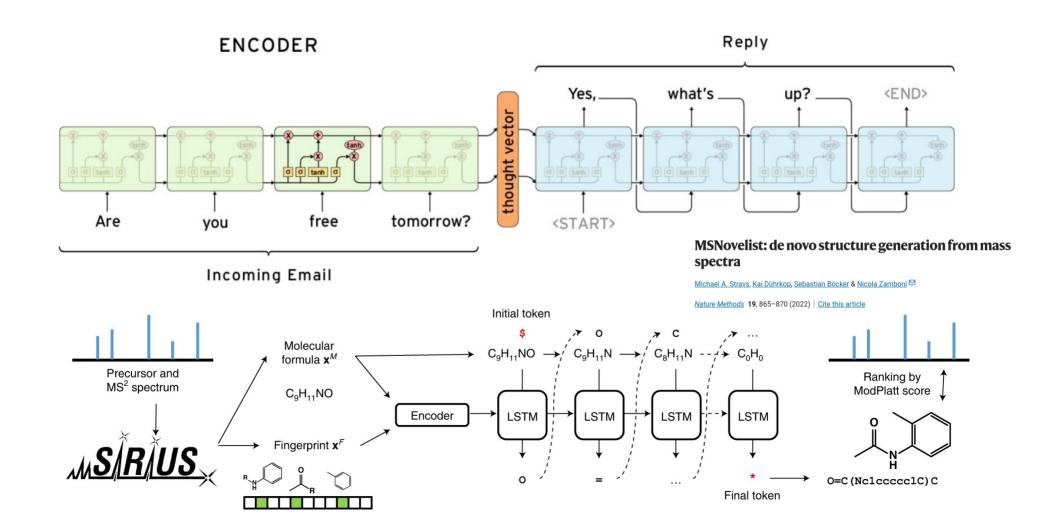




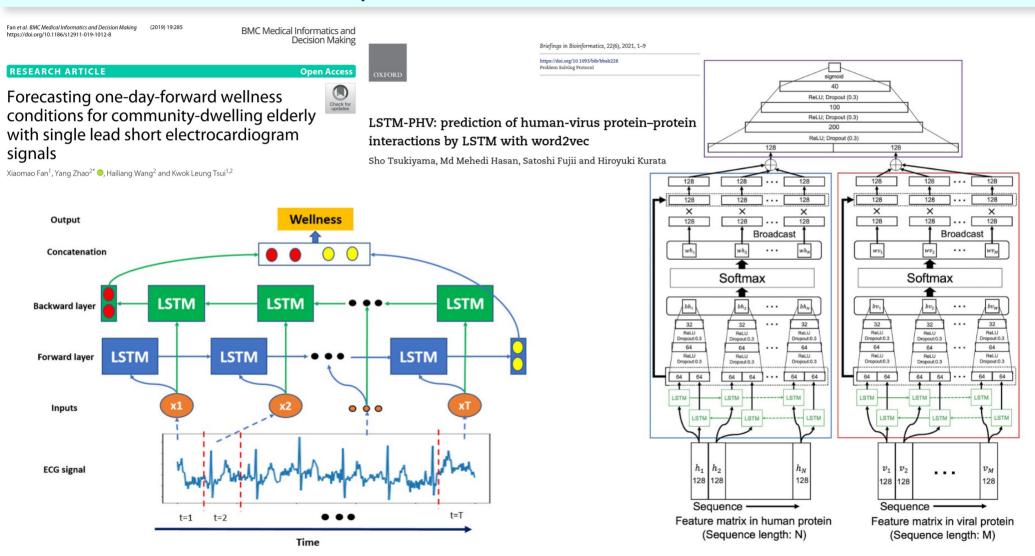




LSTMs for Encoder-Decoder



Examples in the biomedical domain



Examples in the biomedical domain



Contents lists available at ScienceDirect

International Journal of Infectious Diseases



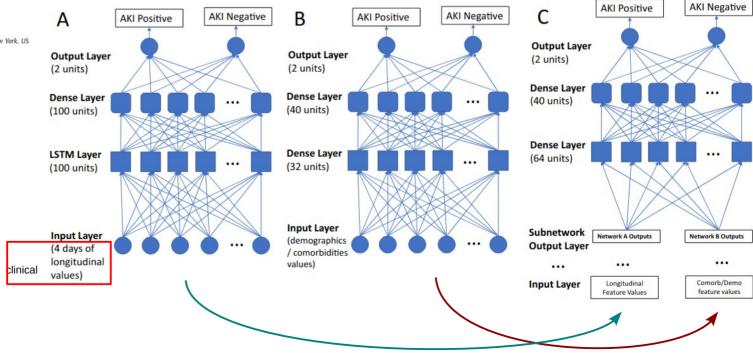
journal homepage: www.elsevier.com/locate/ijid

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, US



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 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.

All authors equal

1 Introduction

Recurrent neural networks, long short-term memory [12] and gated recurrent [7] neural networks

*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 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

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

Work performed while at Google Research

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

The paper that changed everything: the Transfomer

Cited... 198836 times as of 14 October 2025!

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] [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







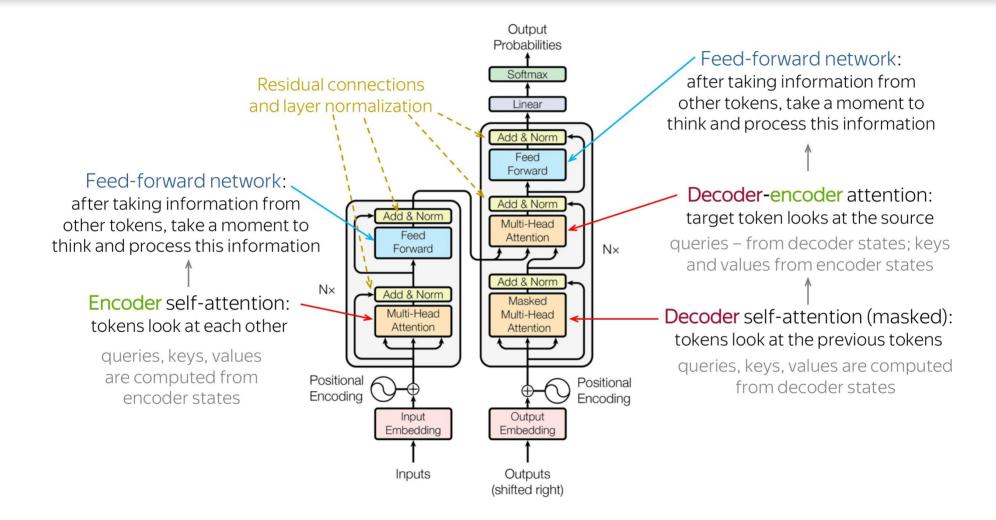


^{*}Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled 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.

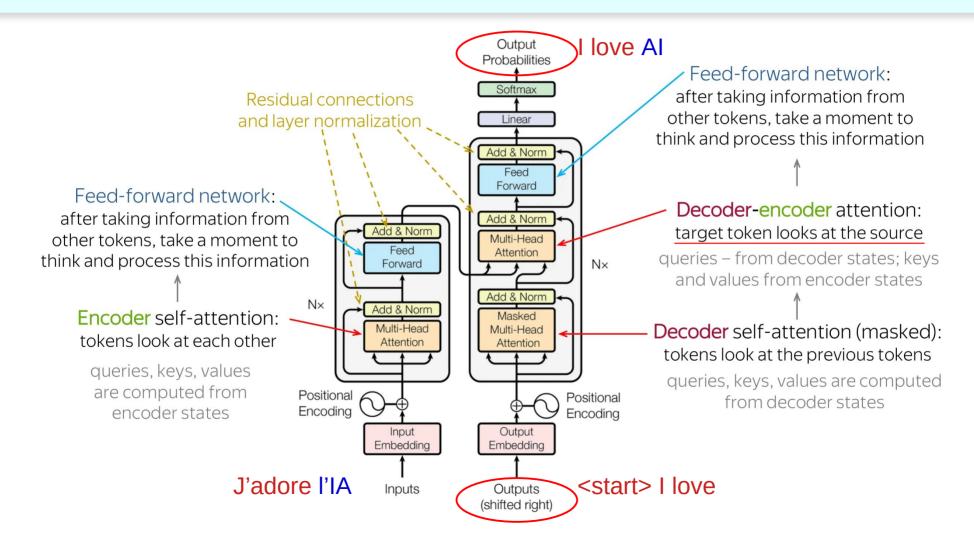
Work performed while at Google Brain.

Work performed while at Google Research.

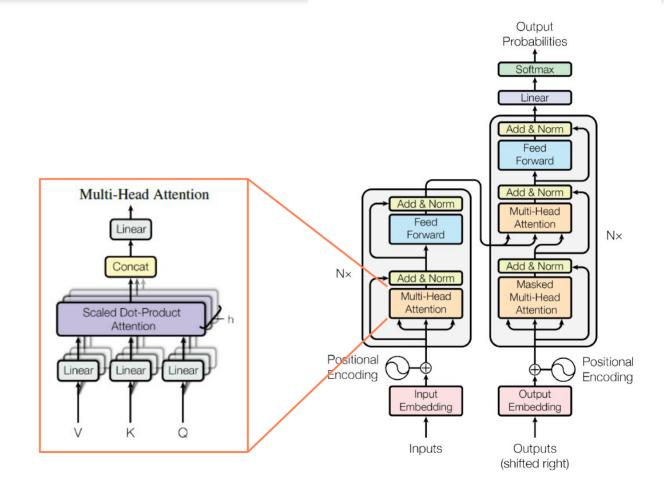
The Transformer: Memory + context = attention



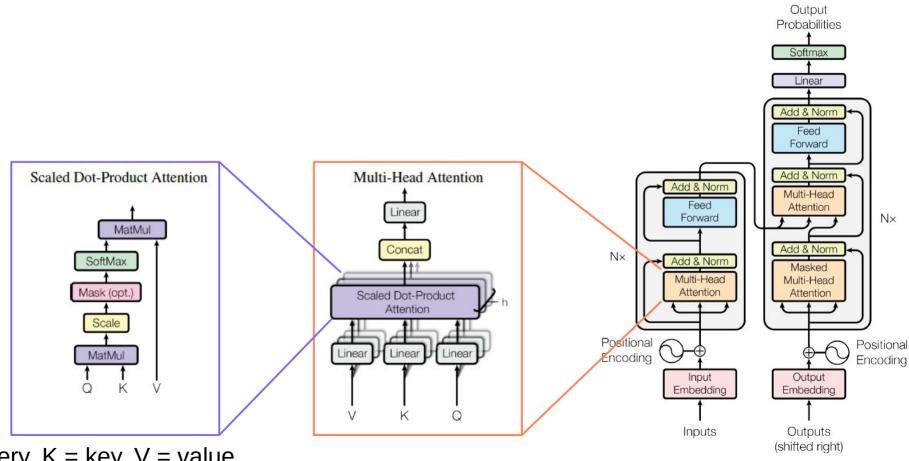
The Transformer: Memory + context = attention



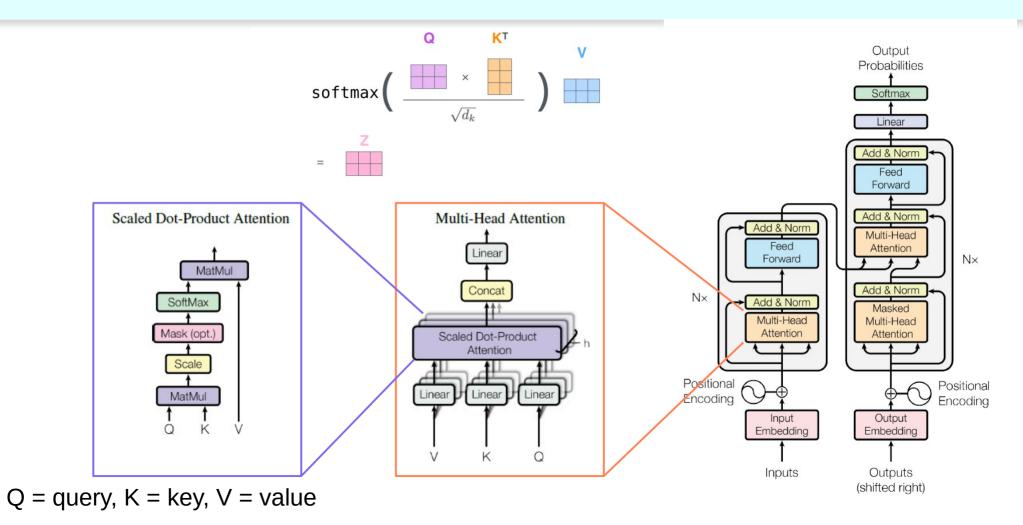
Attention in the Transformer

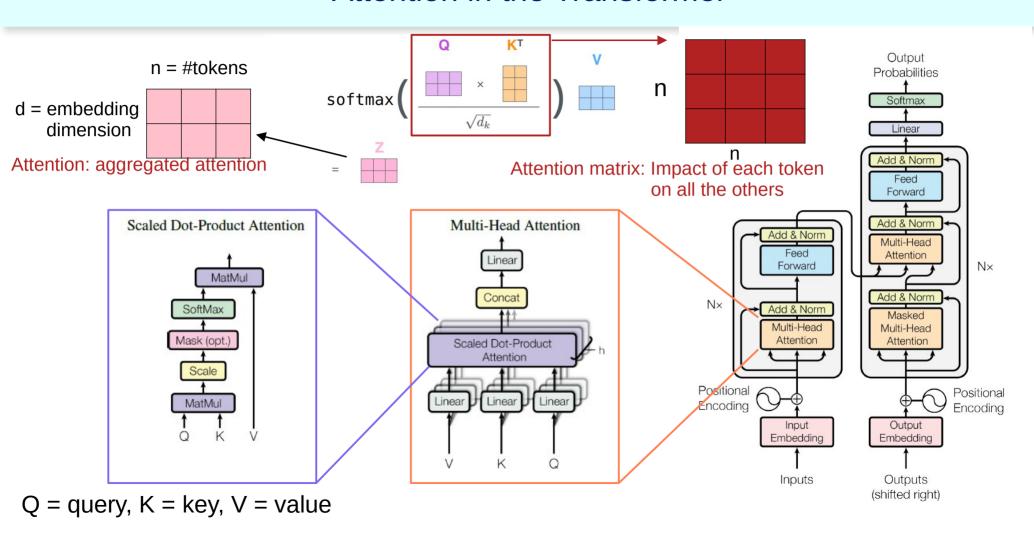


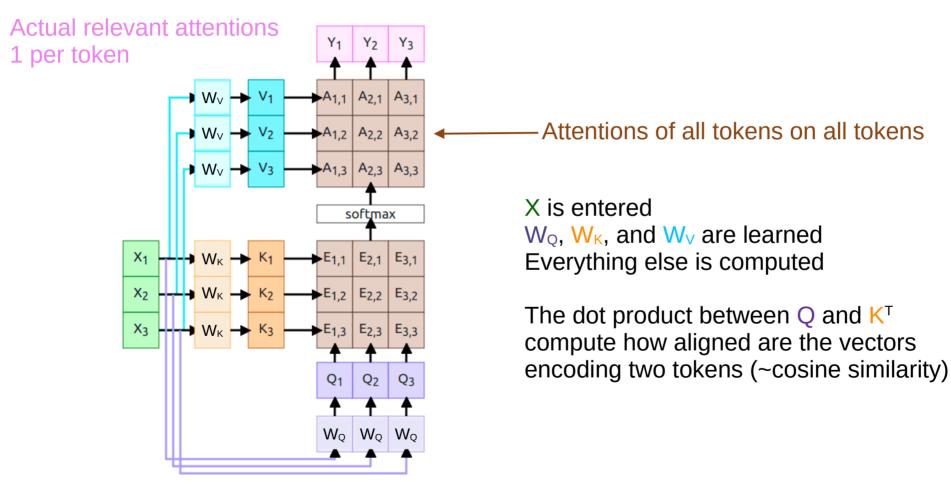
Q = query, K = key, V = value



Q = query, K = key, V = value

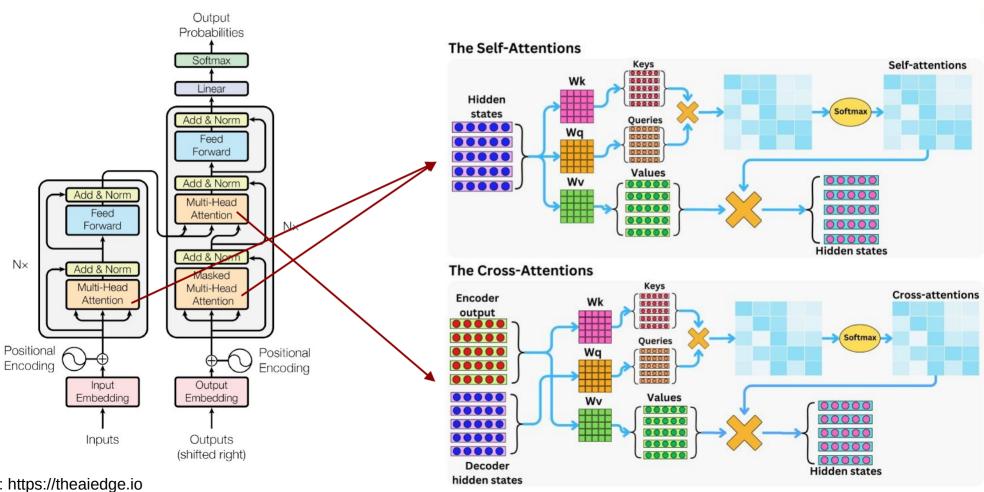






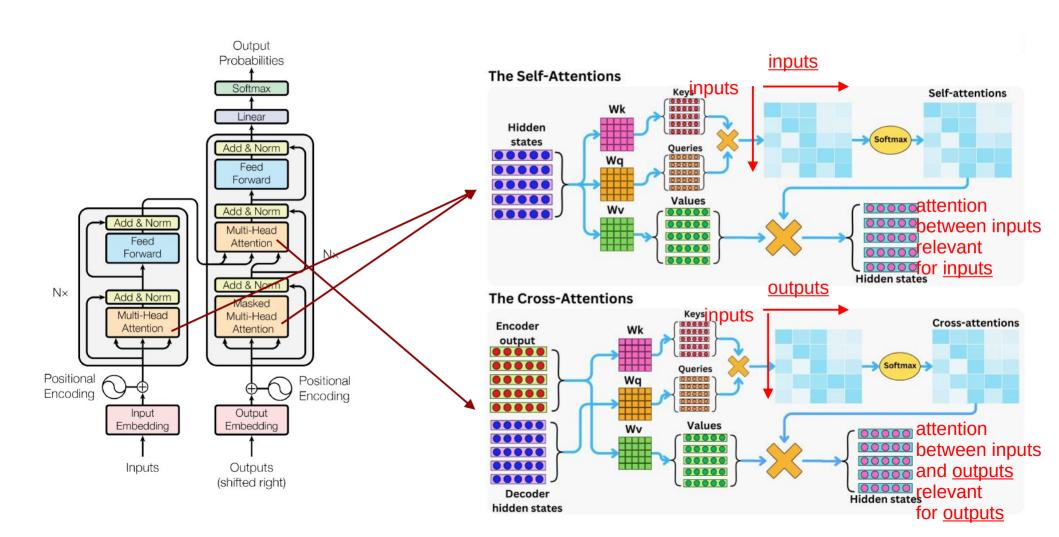
Source: https://erdem.pl/2021/05/introduction-to-attention-mechanism

Self versus cross-attention

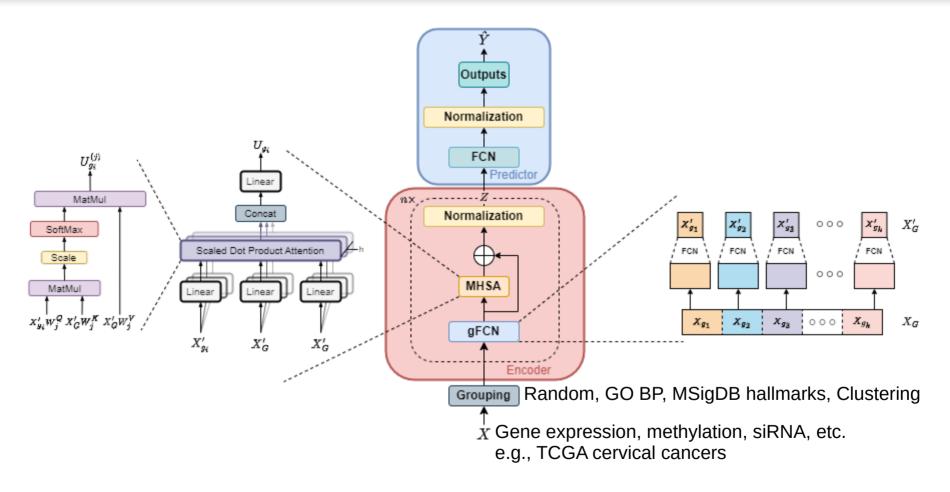


ource: https://theaiedge.io

Self versus cross-attention

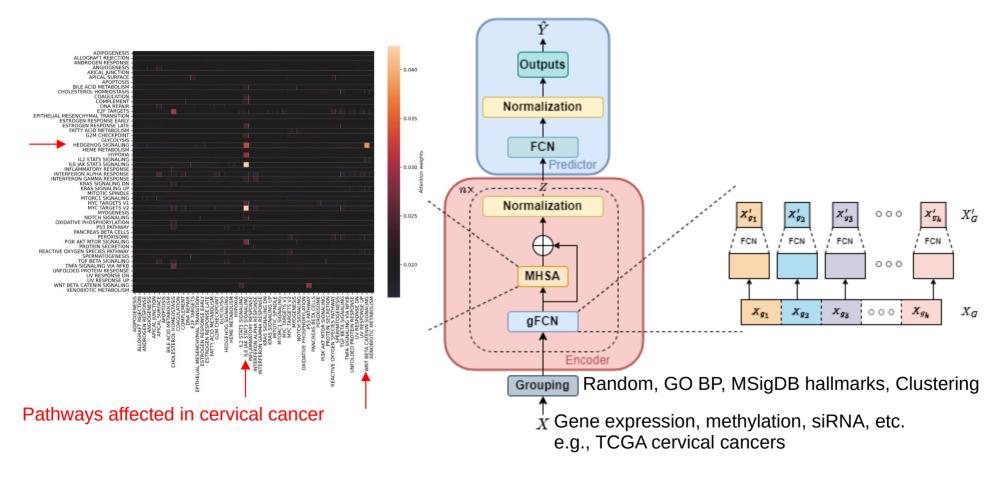


AttOmics: Omics values as tokens



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: Omics values as tokens

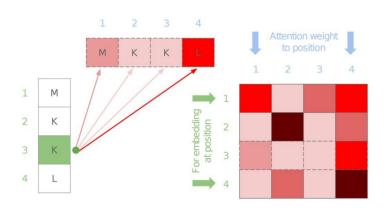


Beaude, A., Rafiee Vahid, M., Augé, F., Zehraoui, F., & Hanczar, B. (2023). AttOmics: attention-based architecture for diagnosis and prognosis from omics data. *Bioinformatics*, 39(Supplement_1), i94-i102.

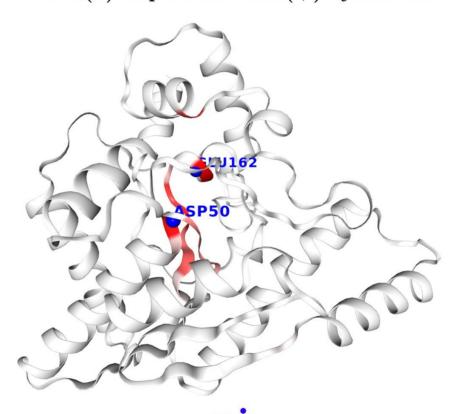
EnzBERT: amino-acids as tokens

Predicting enzymatic function of protein sequences with attention 3

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



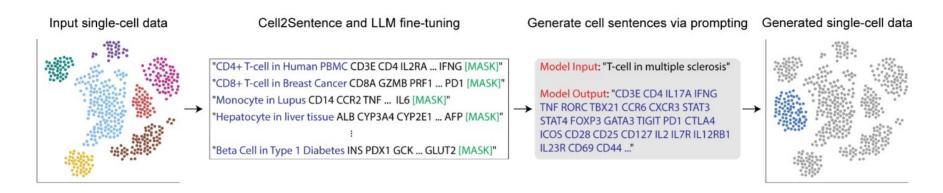
Nh(3)-dependent nad(+) synthetase



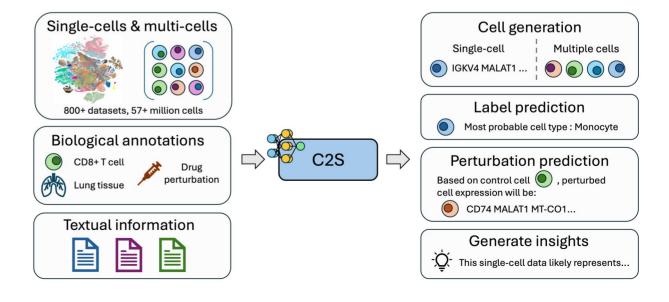
Aggregated attention for each token (amino acid)

- 0 MSMQEKIMRE LHVKPSIDPK QEIEDRVNFL KQYVKKTGAK GFVLGI
- 1 DKSWKFDIKS TVSAFSDQYQ QETGDQLTDF NKGNVKARTR MIAQYAIGGQ EGLLVLSIDH #AEAVTGFFT KYGDGGADLL PLTGLTKRQG RTLLKELGAF
- 2 ERLYLKEPTA DLLDEKPOOS DETELGISHD EIDDYLEGKE VSAKVSEALE KRYSMTEHKR QVPASMFDDW WK

Cell2Sentence: gene names as token



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

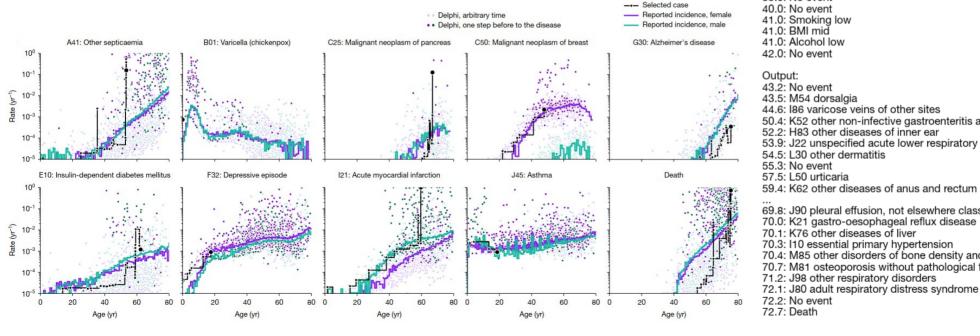


Delphi-2M: Life events as tokens

Article

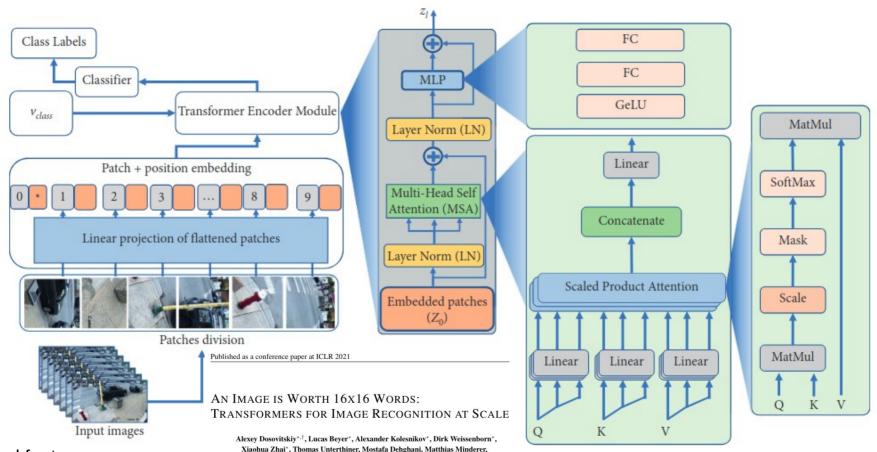
Learning the natural history of human disease with generative transformers

Artem Shmatko^{1,2,3,13}, Alexander Wolfgang Jung^{2,4,5,6,13}, Kumar Gaurav^{2,13}, Søren Brunak^{4,7}, https://doi.org/10.1038/s41586-025-09529-3 Laust Hvas Mortensen^{5,7,8}, Ewan Birney^{2,∞}, Tom Fitzgerald^{2,∞} & Moritz Gerstung^{1,2,9,10,11,12,∞} Received: 18 May 2024 Accepted: 13 August 2025 Decision-making in healthcare relies on understanding patients' past and current Published online: 17 September 2025 health states to predict and ultimately change their future course¹⁻³. Artificial



Input: Age: Token 0.0: Male 2.0: B01 varicella (chickenpox) 3.0: L20 atopic dermatitis 5.0: No event 10.0: No event 15.0: No event 20.0: No event 20.0: G43 migraine 21.0: E73 lactose intolerance 22.0: B27 infectious mononucleosis 25.0: No event 28.0: J11 influenza, virus not identified 30.0: No event 35.0: No event 40.0: No event 41.0: Smoking low 41.0: BMI mid 41.0: Alcohol low 42.0: No event Output: 43.2: No event 43.5: M54 dorsalgia 44.6: I86 varicose veins of other sites 50.4: K52 other non-infective gastroenteritis and colitis 52.2: H83 other diseases of inner ear 53.9: J22 unspecified acute lower respiratory infection 54.5: L30 other dermatitis 55.3: No event 57.5: L50 urticaria 59.4: K62 other diseases of anus and rectum 69.8: J90 pleural effusion, not elsewhere classified 70.0: K21 gastro-oesophageal reflux disease 70.1: K76 other diseases of liver 70.3: I10 essential primary hypertension 70.4: M85 other disorders of bone density and structure 70.7: M81 osteoporosis without pathological fracture 71.2: J98 other respiratory disorders

Vision Transformer



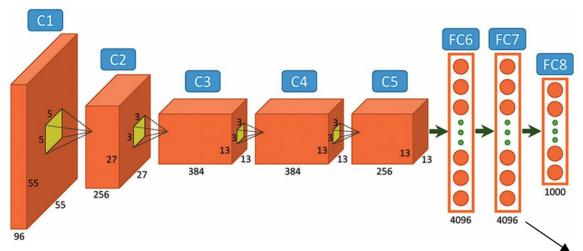
CNN = local features
ViT = relations between distant features

ua Zhai*, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby*.† *equal technical contribution, †equal advising

"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

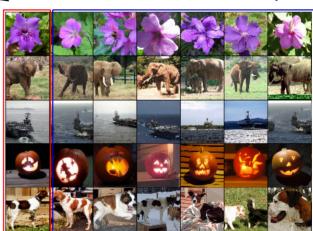


6 nearest neighbours in the 4096 dimension space

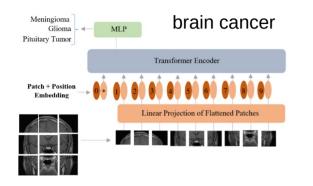
Krizhevsky A, Sutskever I, Hinton GE (2012) ImageNet Classification with Deep Convolutional Neural Networks https://proceedings.neurips.cc/paper/2012/file/ c399862d3b9d6b76c8436e924a68c45b-Paper.pdf

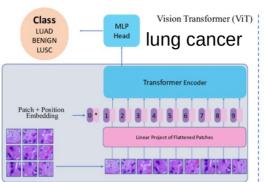
(presenting AlexNet, the first Deep Convolutional Network)

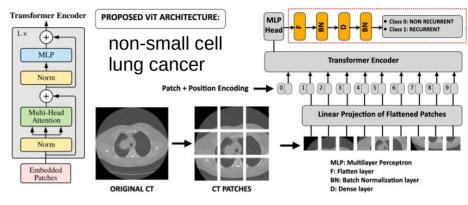
input image

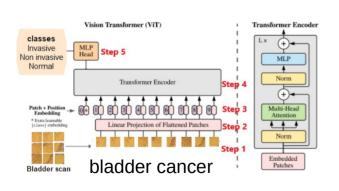


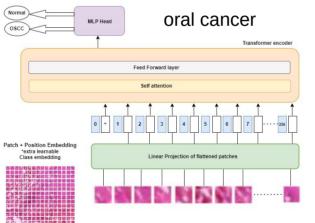
ViTs are replacing vanilla CNNs

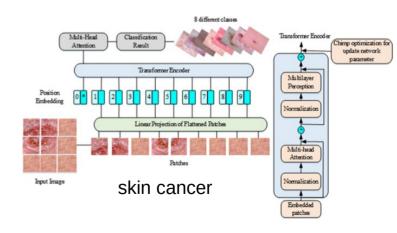




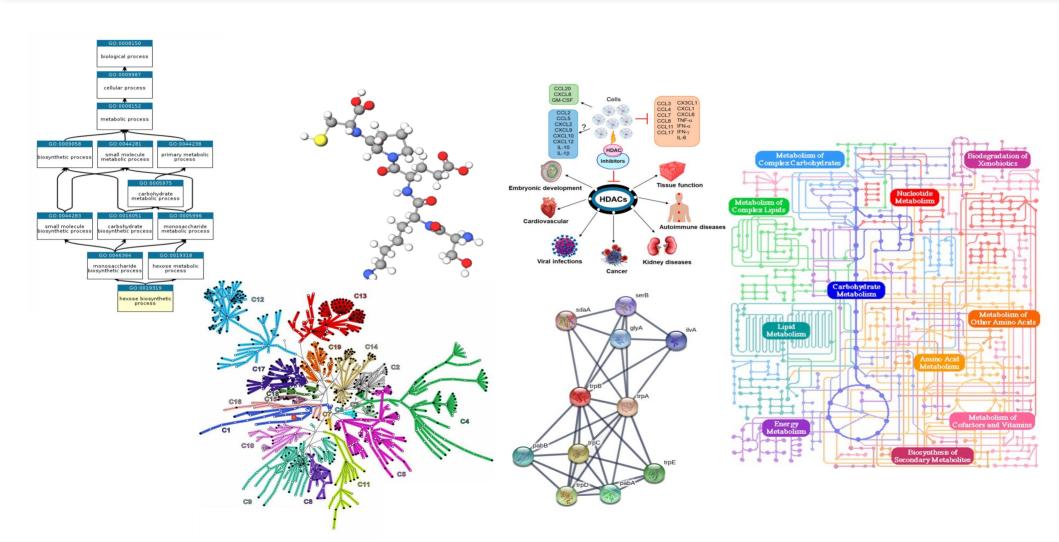






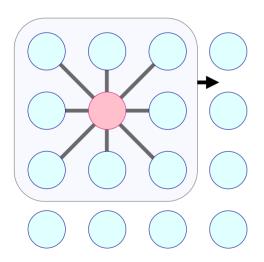


Most biological knowledge comes as graphs



61

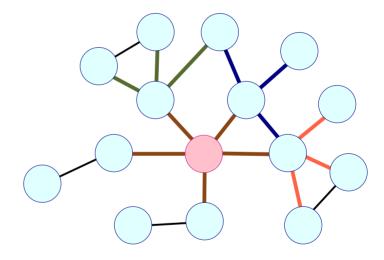
CNNs



Regular grid (same number of neighbours) Homogeneous kernels

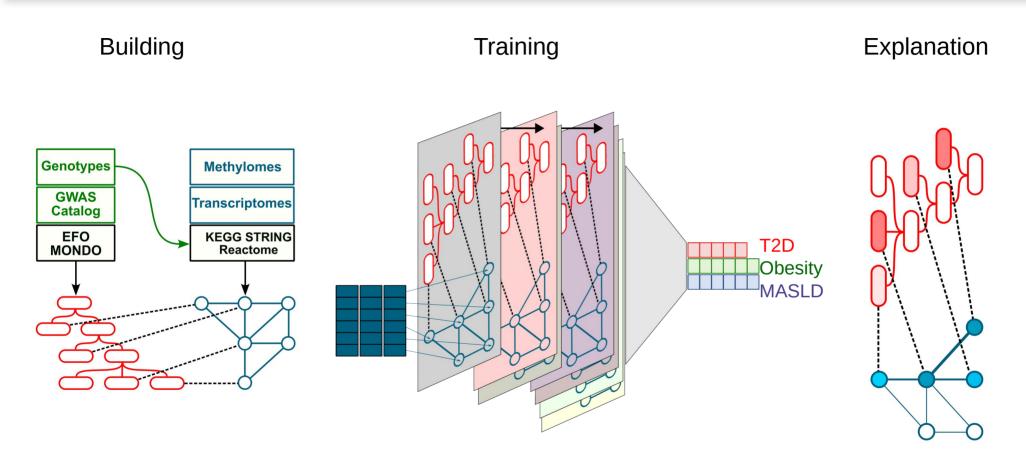
The Graph Neural Network Model

Franco Scarselli, Marco Gori, *Fellow, IEEE*, Ah Chung Tsoi, Markus Hagenbuchner, *Member, IEEE*, and Gabriele Monfardini



Any number of neighbours Information passed from neighbours depends on contexts and positions.

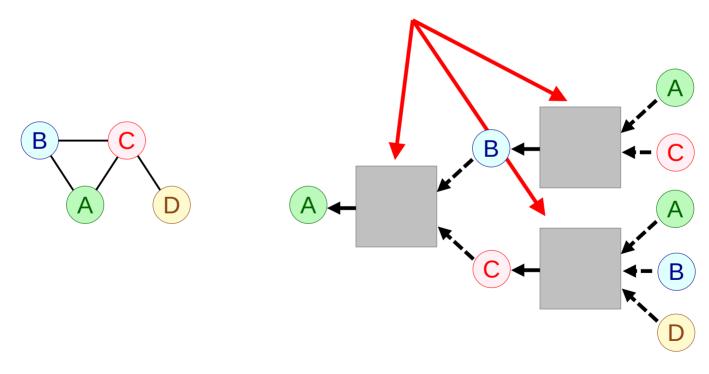
GNN can be heterogeneous



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

Many different ways to update GNNs

Can be message passing (MLP), convolutions, attention-based, or KANs



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

L = D - A

Laplacian matrix

Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering

Michaël Defferrard

Xavier Bresson

Pierre Vandergheynst

$$\mathbf{D} = egin{bmatrix} \mathbf{2} & 0 & 0 & 0 \ 0 & \mathbf{2} & 0 & 0 \ 0 & 0 & 0 & 1 \ \end{bmatrix} egin{bmatrix} \mathsf{A} \ \mathsf{B} \ \mathsf{C} \ \mathsf{C} \ \mathsf{D} \ \mathsf{C} \ \mathsf{C}$$

$$L = D - A$$

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

Degree matrix

identity matrix no influence from neighbours Adjacency matrix

 $y = p_w(L)x$

Laplacian matrix

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

EPFL, Lausanne, Switzerland

Pierre Vandergheynst

kernel de convolution $w = [w_0, w_1, w_2, ..., w_k]$

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

 ${\tt \{michael.defferrard, xavier.bresson, pierre.vandergheynst\}@epfl.ch}$ 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

D only influences C

D influences A, B, C

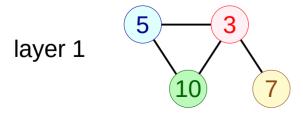




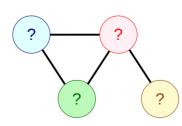
$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix}$$

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

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



layer 2

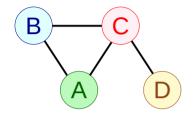




D influences A. B. C



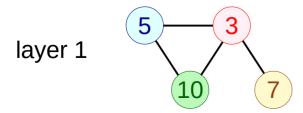




$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix} \qquad L^2 = \begin{bmatrix} 6 & -3 & -4 & 1 \\ -3 & 6 & -4 & 1 \\ -4 & -4 & 12 & -4 \\ 1 & 1 & -4 & 2 \end{bmatrix}$$

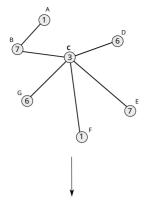
$$L^{2} = \begin{vmatrix} 6 & -3 & -4 & 1 \\ -3 & 6 & -4 & 1 \\ -4 & -4 & 12 & -4 \\ 1 & 1 & -4 & 2 \end{vmatrix}$$

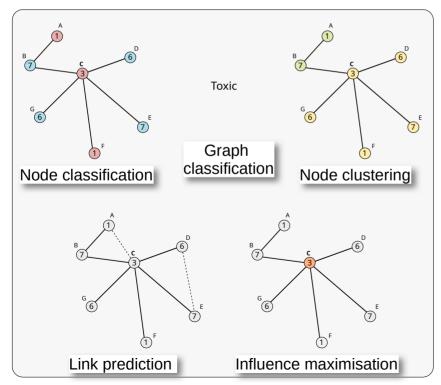
$$w = \begin{bmatrix} 1, 0.1, 0.01 \end{bmatrix} \quad 1 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 10 \\ 5 \\ 3 \\ 7 \end{bmatrix} + 0.1 \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 2 & -1 & 0 \\ -1 & -1 & 3 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix}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



layer 2







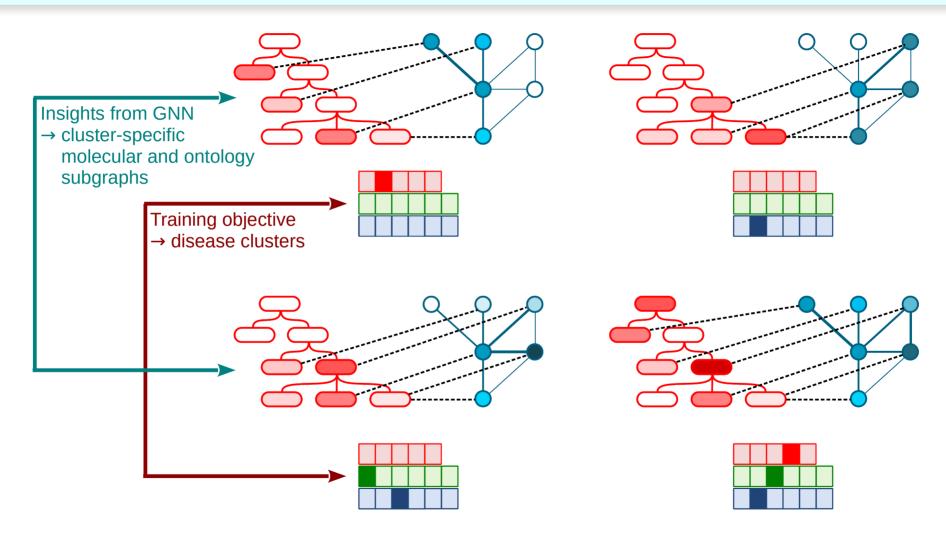
What can we do with GNN?

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/

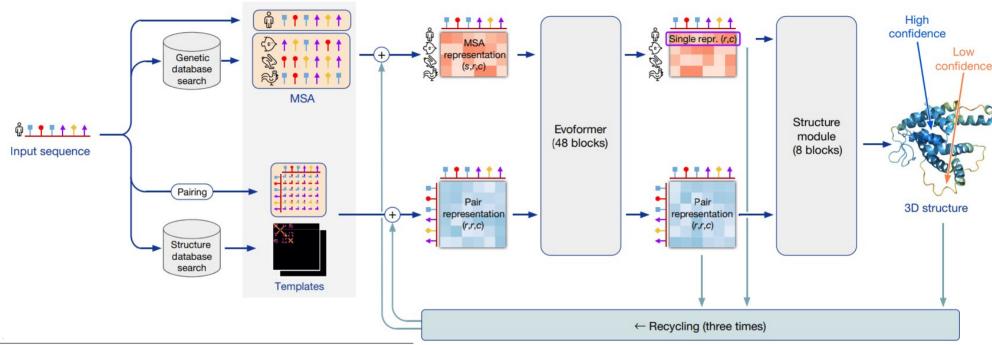
Both by Google Research teams

GNN insights can be subgraphs





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

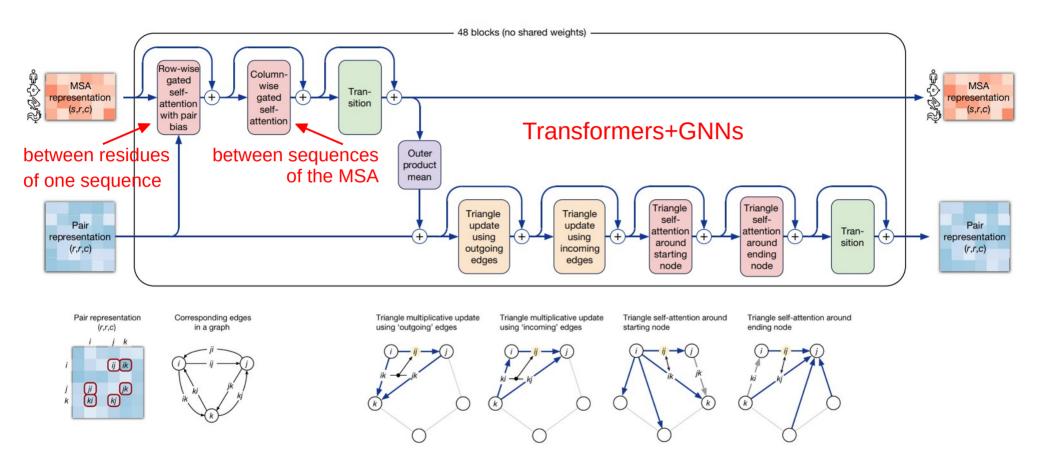
Check for updates

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)

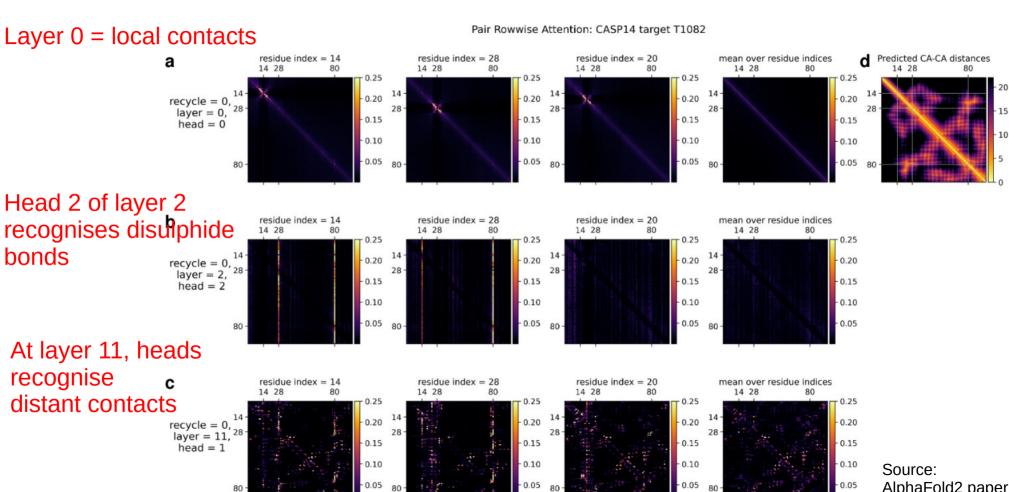
https://github.com/google-deepmind/alphafold

AlphaFold2: evoformer



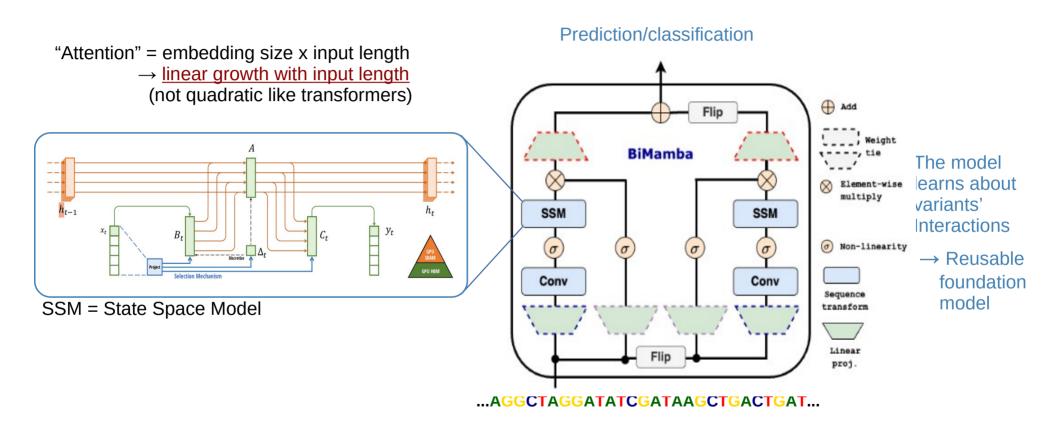
see also: https://www.blopig.com/blog/2021/07/alphafold-2-is-here-whats-behind-the-structure-prediction-miracle/

Row-wise attention: between residues of a sequence



AlphaFold2 paper supplementary info

RNNs are back. Rise of the Mamba



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

Mamba everywhere

Computer Science > Computer Vision and Pattern Recognition

[Submitted on 17 Jan 2024 (v1), last revised 14 Nov 2024 (this version, v3)]

Vision Mamba: Efficient Visual Representation Learning with Bidirectional State Space Model

Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, Xinggang Wang

MambaVision: A Hybrid Mamba-Transformer Vision Backbone

JOURNAL ARTICLE

MambaCpG: an accurate model for single-cell DNA methylation status imputation using mamba 3

Qi Zhao, Ze Li, Qian Mao, Tingwei Chen, Yiran Zhang, Bingle Li, Zheng Zhao ☎, Xiaoya Fan ☎

Briefings in Bioinformatics, Volume 26, Issue 4, July 2025, bbaf360, https://doi.org/10.1093/bib/bbaf360

Published: 28 July 2025 Article history ▼

Ali Hatamizadeh, Jan Kautz NVIDIA

{ahatamizadeh, jkautz}@nvidia.com

Computer Science > Machine Learning

[Submitted on 15 Feb 2025 (v1), last revised 18 Feb 2025 (this version, v2)]

HybriDNA: A Hybrid Transformer-Mamba2 Long-Range DNA Language Model

Mingqian Ma, Guoqing Liu, Chuan Cao, Pan Deng, Tri Dao, Albert Gu, Peiran Jin, Zhao Yang, Yingce Xia, Renqian Luo, Pipi H

Computer Science > Machine Learning

[Submitted on 13 Feb 2024 (v1), last revised 19 Feb 2024 (this version, v2)]

Graph Mamba: Towards Learning on Graphs with State Space Models

Ali Behrouz, Farnoosh Hashemi

Computer Science > Machine Learning

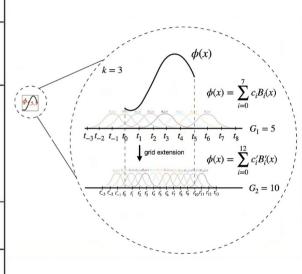
[Submitted on 1 Feb 2024]

Graph-Mamba: Towards Long-Range Graph Sequence Modeling with Selective State Spaces

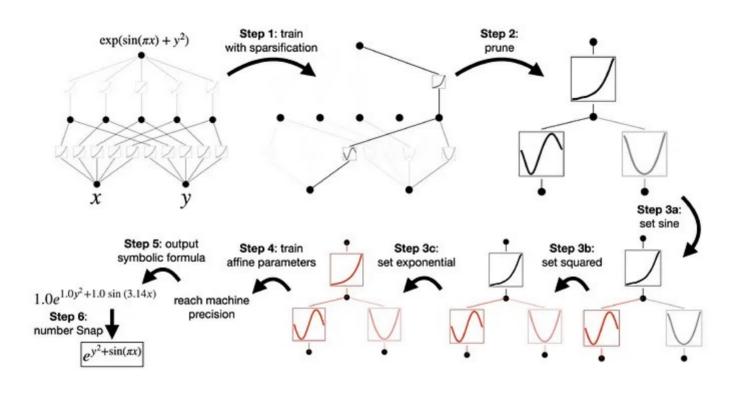
Chloe Wang, Oleksii Tsepa, Jun Ma, Bo Wang

Kolmogorov Arnold Networks

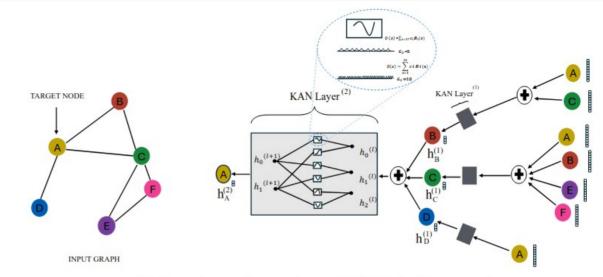
Model	Multi-Layer Perceptron (MLP)	Kolmogorov-Arnold Network (KAN)
Theorem	Universal Approximation Theorem	Kolmogorov-Arnold Representation Theorem
Formula (Shallow)	$f(\mathbf{x}) \approx \sum_{i=1}^{N(\epsilon)} a_i \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i)$	$f(\mathbf{x}) = \sum_{q=1}^{2n+1} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right)$
Model (Shallow)	fixed activation functions on nodes learnable weights on edges	learnable activation functions on edges sum operation on nodes
Formula (Deep)	$\mathrm{MLP}(\mathbf{x}) = (\mathbf{W}_3 \circ \sigma_2 \circ \mathbf{W}_2 \circ \sigma_1 \circ \mathbf{W}_1)(\mathbf{x})$	$KAN(\mathbf{x}) = (\mathbf{\Phi}_3 \circ \mathbf{\Phi}_2 \circ \mathbf{\Phi}_1)(\mathbf{x})$
Model (Deep)	(c) $\begin{array}{c c} W_3 \\ \hline \\ W_2 \\ \hline \\ W_1 \\ \hline \\ W_2 \\ \hline \\ Uinear, \\ learnable \\ X \\ \hline \\ X \\ \end{array}$	(d) Φ_3 Φ_2 nonlinear, learnable Φ_1 X



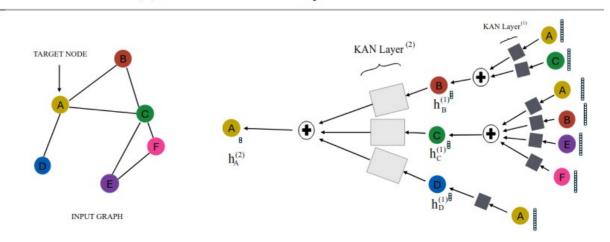
KANs to extract symbolic formulas



KANs in GNNs



(b) Overview of a two-layer GKAN Architecture 1.



Kiamari *et al.*(2024) GKAN: Graph Kolmogorov-Arnold Networks. https://doi.org/10.48550/arXiv.2406.06470

How many architectures can you cram into one model?

nature > scientific reports > articles > article

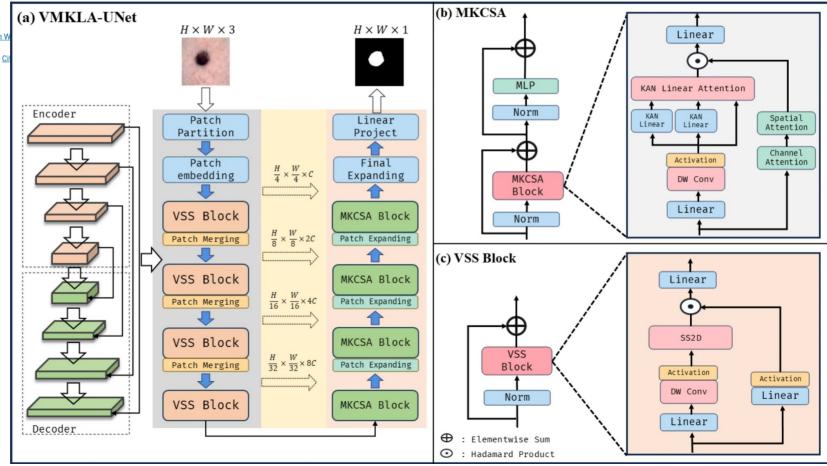
Article Open access Published: 17 April 2025

VMKLA-UNet: vision Mamba with KAN linear attention U-

Net

Chenhong Su, Xuegang Luo, Shiqing Li, Li Chen & Juan W

Scientific Reports 15, Article number: 13258 (2025)

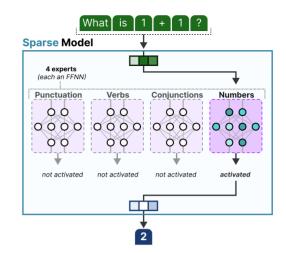


Topics to explore

Mixture Of Experts

https://huggingface.co/blog/moe

https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-mixture-of-experts



Reinforcement learning

https://fr.mathworks.com/content/dam/mathworks/ebook/gated/reinforcement-learning-ebook-all-chapters.pdf

https://web.stanford.edu/class/psych209/Readings/SuttonBartoIPRLBook2ndEd.pdf

https://arxiv.org/pdf/2412.05265



RAG

https://en.wikipedia.org/wiki/Retrieval-augmented_generation https://blogs.nvidia.com/blog/what-is-retrieval-augmented-generation/

