# Beyond MLPs

# Part 2/2: RNNs, Attention and GNNs
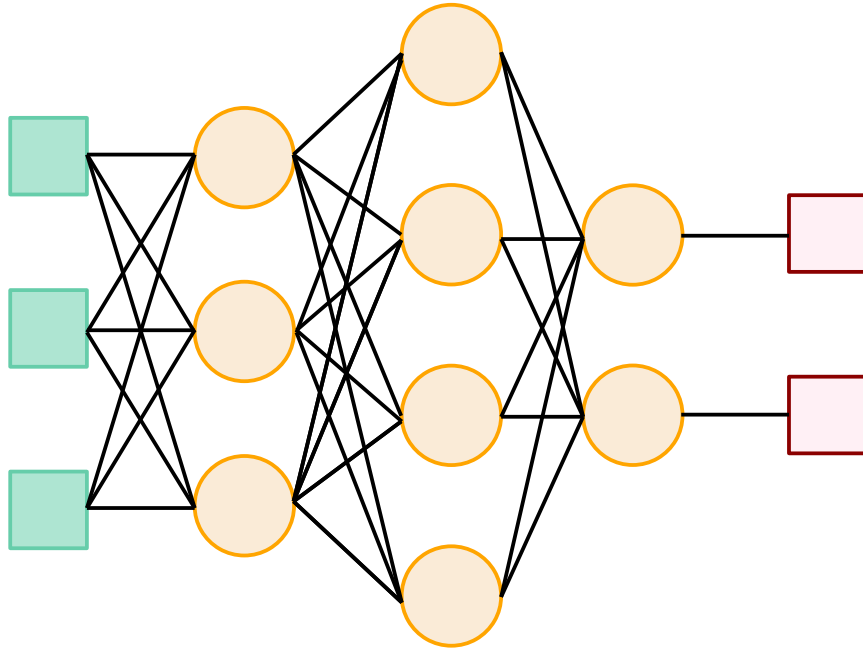
*nicolas.gambardella@univ-lille.fr*

**McCullogh & Pitts** Artificial neuron

**Rosenblatt** Perceptron

**Amari** Stochastic gradient descent

**Minski & Papert** *Perceptrons*

**Linnainmaa** Backpropagation

**Werbos**

**Lecture 1**

1943 — 1958 — 1967 — 1969 — 1970 — 1971

**Qian & Sejnowski** 1st NN prot struct pred

**Hochreiter** LSTM

**Le Cun** *et al* LeNet-5

**Oh & Jung** Use of GPUs

1995 — 1993 — 1988 — 1986 — 1980

1998 — 2004

**Lecture 2**

**Fukushima** CNN (neocognitron)

**Rumelhart & McClelland** AutoEncoder

**Rost & Sander** PHD: Cascading NN

**Goodfellow** *et al* GANs

Alphafold

**Lecture 3**

**Scarselli** *et al* *GNNs*

**Kingmat & Welling** VAEs

2009 — 2013 — 2014 — 2015 — 2017 — 2018 — 2022

AlphaGo

**Vaswani** *et al* Transformer

ChatGPT

What is an artificial neuron?

$x_1$

$x_i$

$x_n$

$w_1$

$w_i$

$w_n$

$b$

$\Sigma$

Activation function

$y$

$$\sum_{i=1}^{n} w_i x_i + b \qquad y = a(\sum_{i=1}^{n} w_i x_i + b)$$
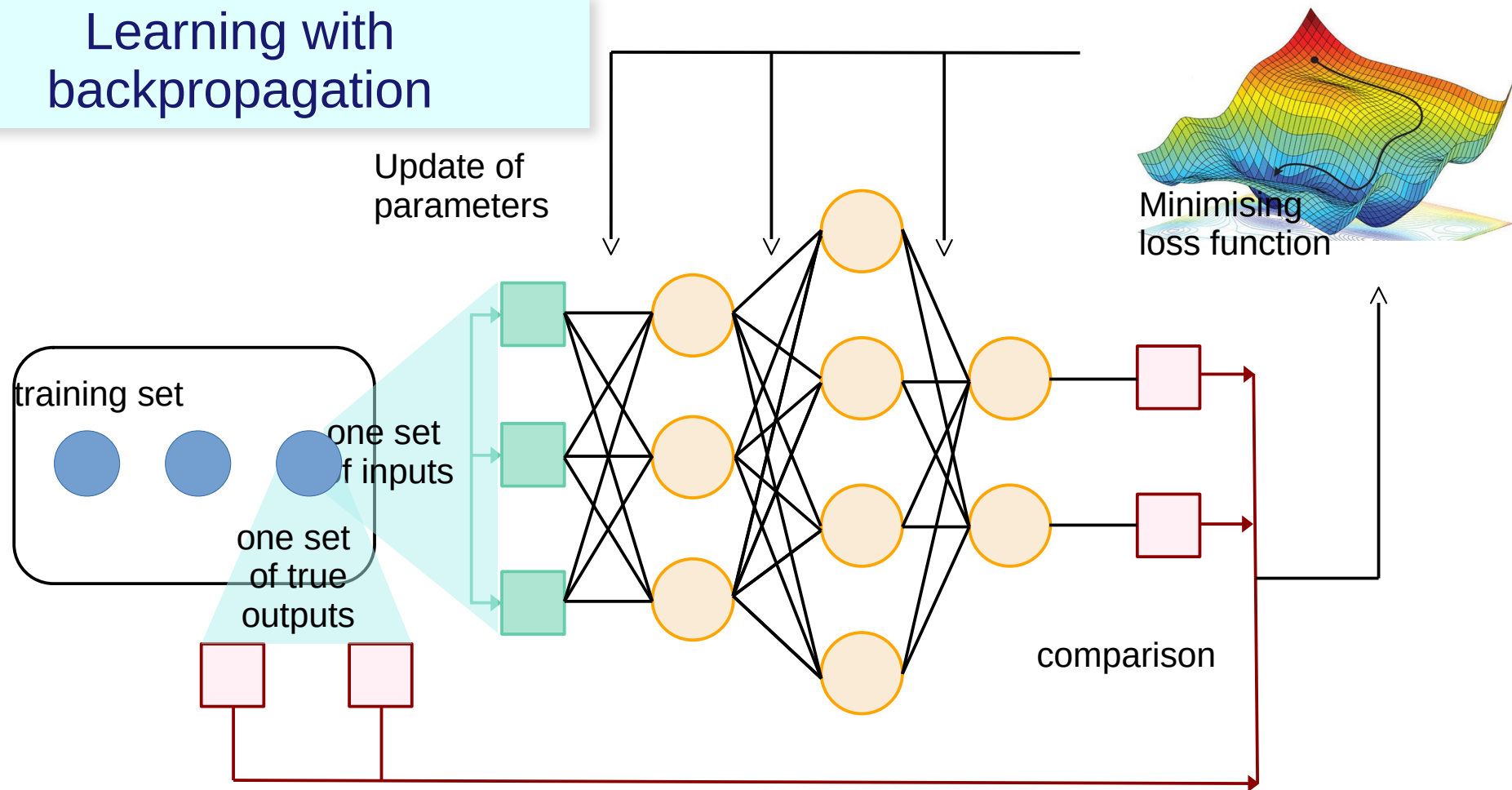
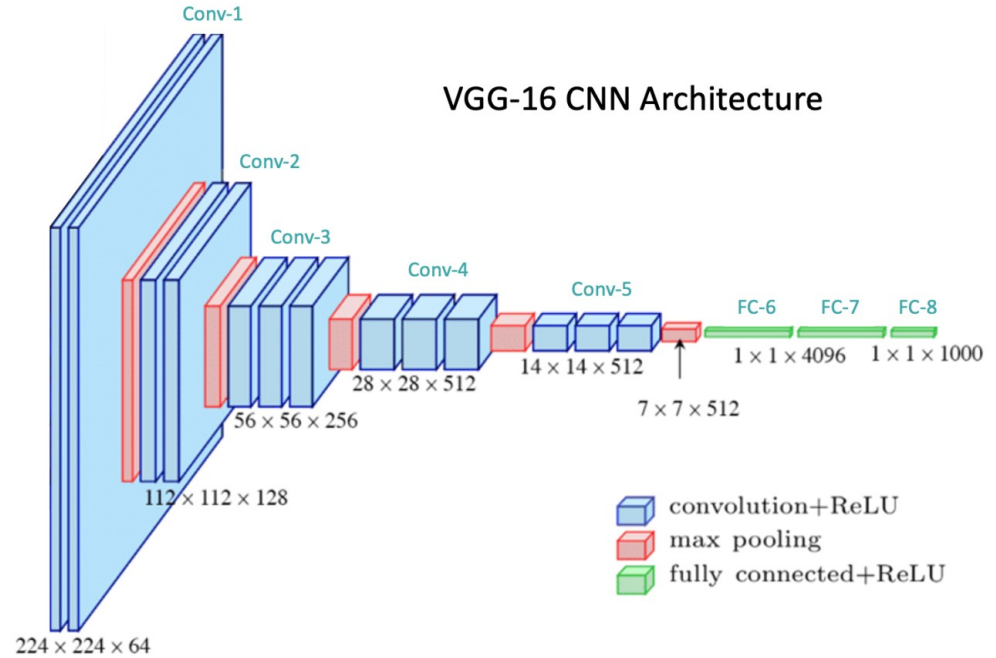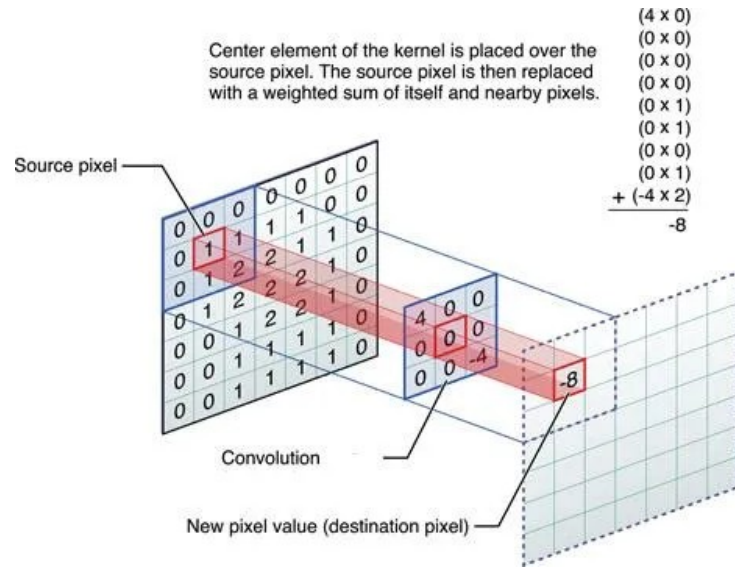# All inputs can be independent and everything connected to everything



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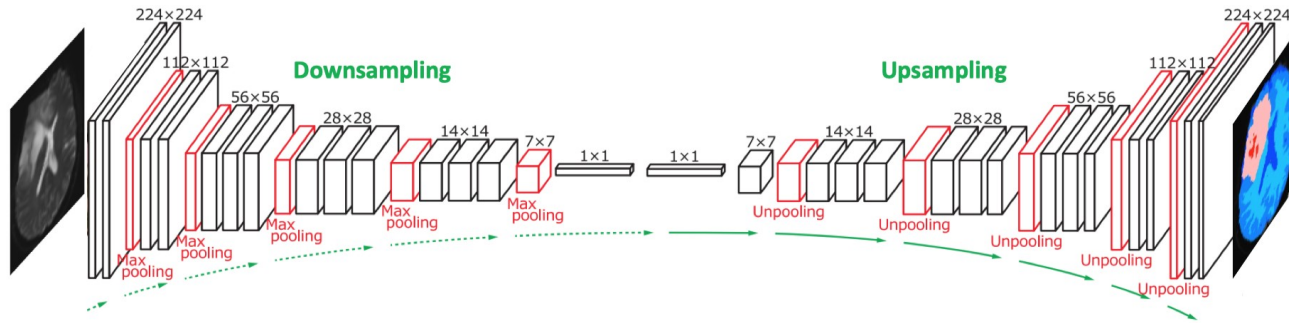made of Fully Connected layers (FC)

# Learning with backpropagation

Update of parameters

Minimising loss function

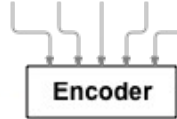training set

one set of inputs

one set of true outputs

comparison

Deep Convolutional Neural Networks (CNN)

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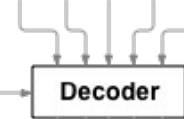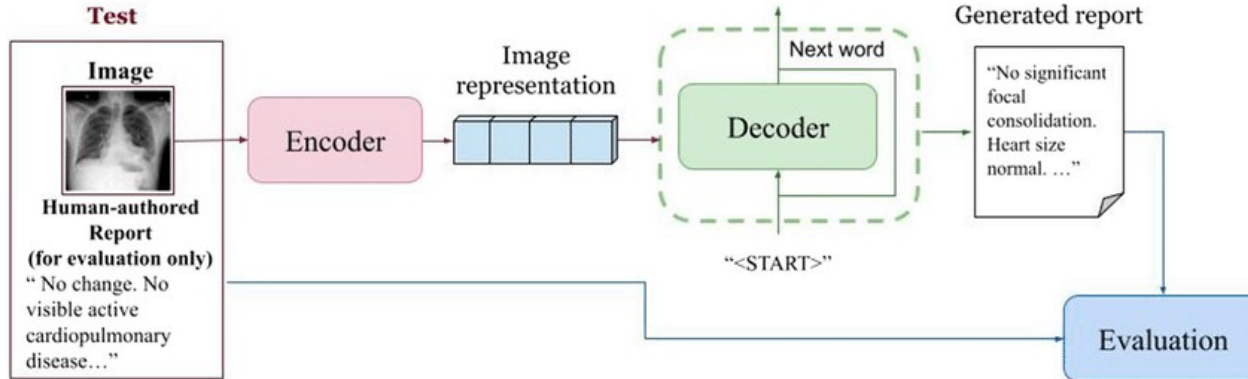# Decoder networks can reconstruct the information from it

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



Mean

$Z_\mu$

**ENCODER**
$P_\phi(Z|X)$

$X$

Input-Image

$Z_\sigma$

Variance or
Standard Deviation

$Z$

Sampled
Latent-Vector

**DECODER**
$P_\theta(\hat{X}|Z)$

$\hat{X}$

Predicted-Image from Z

Generative model

Sample a point from $G(Z_\mu, Z_\sigma)$

$$Z = \mu + \sigma \odot \epsilon$$
$$\epsilon \sim \mathcal{N}(0,1)$$

$Z_\mu$

$Z_\sigma$

**Latent
Distribution**
Latent-Variables to
follow a Standard Normal Distribution

Latent
Representation

Input

Encoder

Decoder

Output

Direct comparison
=
Unsupervised learning

Search

Advanced Search

New Results

🔔 Follow this preprint   ◀ Previous                                    Next ▶

## Deep learning models reading clinical data and liver omics strongly distinguish NASH from steatosis and suggest new genes involved in liver disease severity

🔘 Nicolas Gambardella, Smaïn Fettem, Mathilde Boissel, Lijiao Ning, Violeta Raverdy, Marwa Afnouch, Souhila Amanzougarene, Mehdi Derhourhi, Bénédicte Toussaint, Emmanuel Vaillant, Amna Khamis, Philippe Lefebvre, 🔘 Bart Staels, Francois Pattou, Philippe Froguel, Amelie Bonnefond

doi: https://doi.org/10.1101/2025.10.10.681581 ⓒⓡ

Posted October 10, 2025.

🔗 **Download PDF**          ✉ Email
🗎 **Print/Save Options**     ➦ Share
🗎 Supplementary Material    🌐 Citation Tools
                            ▦ Get QR code

**Subject Area**

Molecular Biology ▶

| Abstract | Info/History | Metrics |            🗎 Preview PDF

## Abstract

Background & Aims Metabolic dysfunction-associated steatotic liver disease (MASLD) is a frequent co-morbidity of obesity and diabetes, with prevalence increasing worldwide. Recognising liver disease stages and elucidating the molecular underpinning of their progression are thus medically important. Methods Using data gathered from 300 patients with obesity of the ABOS cohort, we selected non-redundant clinical variables, gene expression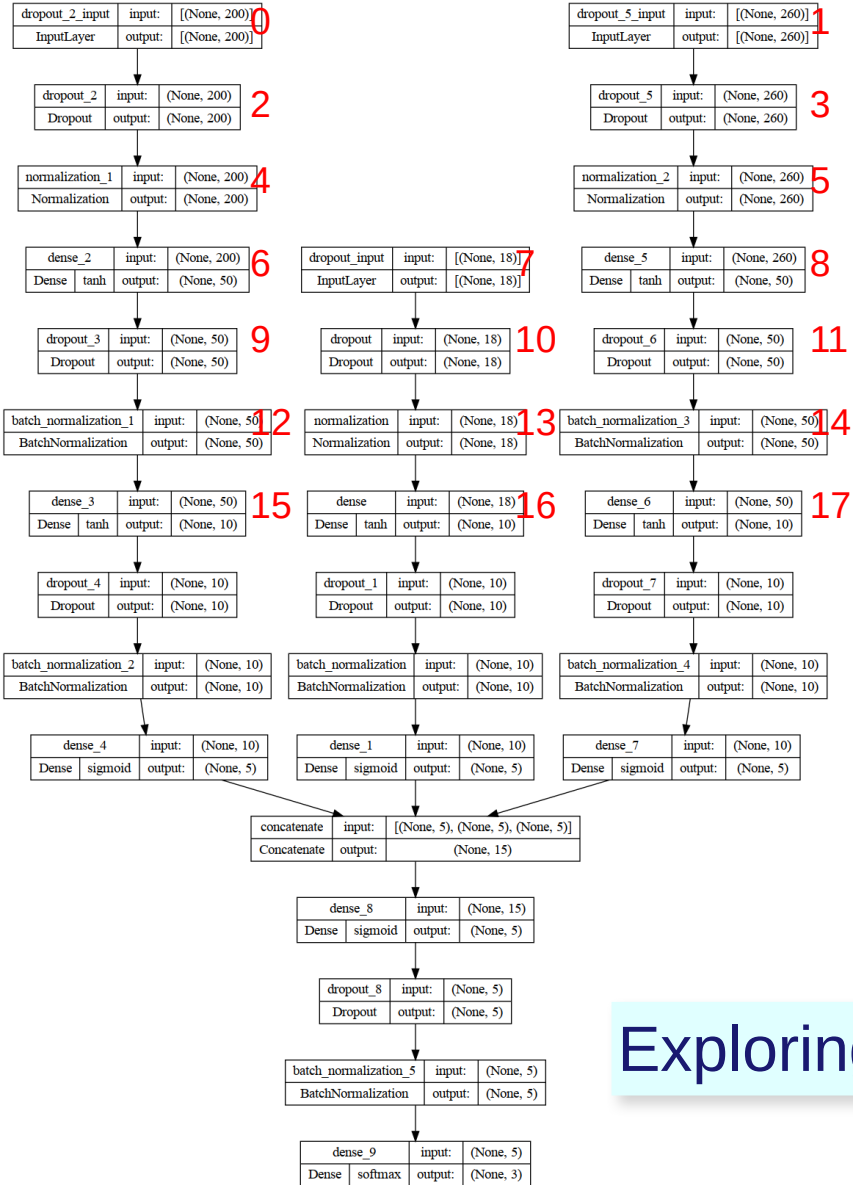s and CpGs methylation levels most associated with severity using unsupervised approaches to train a multi-module, multi-layer perceptron predicting patients liver status. Results The combination of five models trained on the three modalities reached an AUC of 0.945 on a

**Reviews and Context**

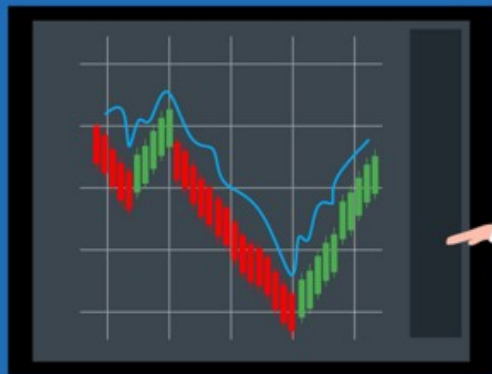| 0 | Comment | 💬 |
| 0 | TRIP Peer Reviews | ☑ |
| 0 | Community Reviews | 👥 |
| 1 | Automated Services | ⚙ |
| 0 | Blogs/Media | 🖥 |
| 0 | Author Videos | 🎬 |

```
218
219   # Before concatenation - ClinDat, layer 15, "dense_3" from model1
220   latCoordClinDatlist = {}
221   for x in range(1,nbmodels+1):
222       model = modlist["model{0}".format(x)]
223       intermediate_model = keras.Model(inputs  = model.input,
224                                        outputs = model.layers[15].output)
225       activations = intermediate_model.predict([ClinDat_X, RNAseq_X, Methylo_X])
226       act_df = pd.DataFrame(activations)
227       latCoordClinDatlist["latCoord{0}".format(x)] = pd.concat([init_cols, act_df], axis = 1)
228       latCoordClinDatlist["latCoord{0}".format(x)].to_csv(path.join(ResultDir,rootname,
      ↳      latCoordfilename +str(x)+"-Whole_ClinDatBefConcat.csv"), index = False)
229
230   # Before concatenation - RNAseq, layer 16, "dense_3" from model1
231   latCoordRNAseqlist = {}
232   for x in range(1,nbmodels+1):
233       model = modlist["model{0}".format(x)]
234       intermediate_model = keras.Model(inputs  = model.input,
235                                        outputs = model.layers[16].output)
236       activations = intermediate_model.predict([ClinDat_X, RNAseq_X, Methylo_X])
237       act_df = pd.DataFrame(activations)
238       latCoordRNAseqlist["latCoord{0}".format(x)] = pd.concat([init_cols, act_df], axis = 1)
239       latCoordRNAseqlist["latCoord{0}".format(x)].to_csv(path.join(ResultDir,rootname,
      ↳      latCoordfilename +str(x)+"-Whole_RNAseqBefConcat.csv"), index = False)
240
241   # Before concatenation - Methylo, layer 17, "dense_6" from model1
242   latCoordMethylolist = {}
243   for x in range(1,nbmodels+1):
244       model = modlist["model{0}".format(x)]
245       intermediate_model = keras.Model(inputs  = model.input,
246                                        outputs = model.layers[17].output)
247       activations = intermediate_model.predict([ClinDat_X, RNAseq_X, Methylo_X])
248       act_df = pd.DataFrame(activations)
249       latCoordMethylolist["latCoord{0}".format(x)] = pd.concat([init_cols, act_df], axis = 1)
250       latCoordMethylolist["latCoord{0}".format(x)].to_csv(path.join(ResultDir,rootname,
      ↳      latCoordfilename +str(x)+"-Whole_MethyloBefConcat.csv"), index = False)
251
252
```

Build a partial model from the top to the latent space we want

get the activation values

Exploring latent spaces

# Time series and sequences of variable lengths

# Time series and sequences of variable lengths

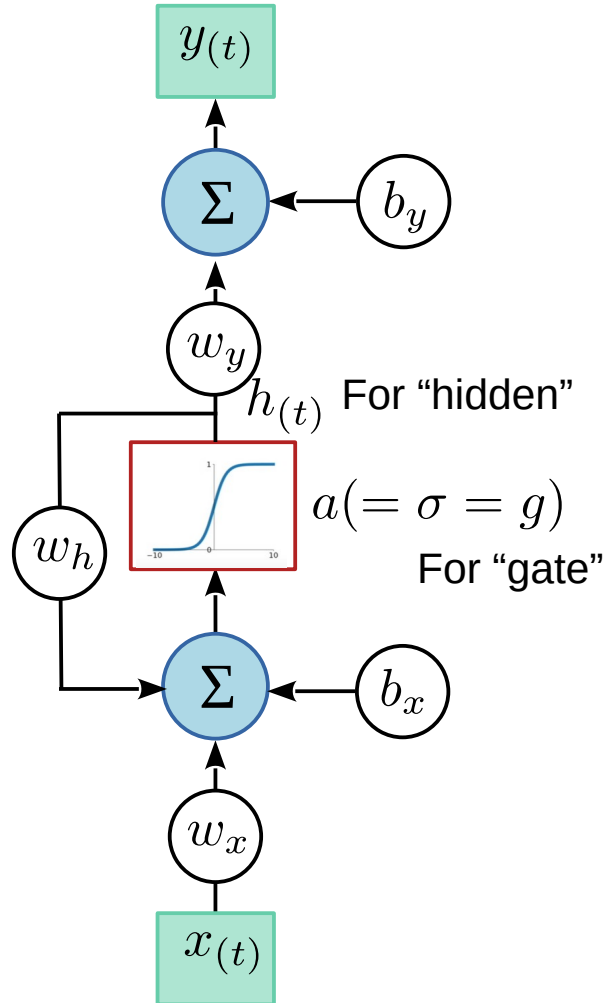| | | |
|---|---|---|
| Speech recognition | [audio waveform] → | "The quick brown fox jumped over the lazy dog." |
| Music generation | Ø → | [musical notes] |
| Sentiment classification | "There is nothing to like in this movie." → | ★☆☆☆☆ |
| DNA sequence analysis | AGCCCCTGTGAGGAACTAG → | AGCCCCTGTGAGGAACTAG |
| Machine translation | Voulez-vous chanter avec moi? → | Do you want to sing with me? |
| Video activity recognition | [video frames] → | Running |
| Name entity recognition | Yesterday, Harry Potter met Hermione Granger. → | Yesterday, Harry Potter met Hermione Granger. |

Andrew Ng

# Recurrent Neural Networks: successive inputs are not independent

$y_{(t)}$

$\Sigma$ ← $b_y$

$w_y$

$h_{(t)}$ For "hidden"

$a(= \sigma = g)$

For "gate"

$w_h$

$\Sigma$ ← $b_x$

$w_x$

$x_{(t)}$

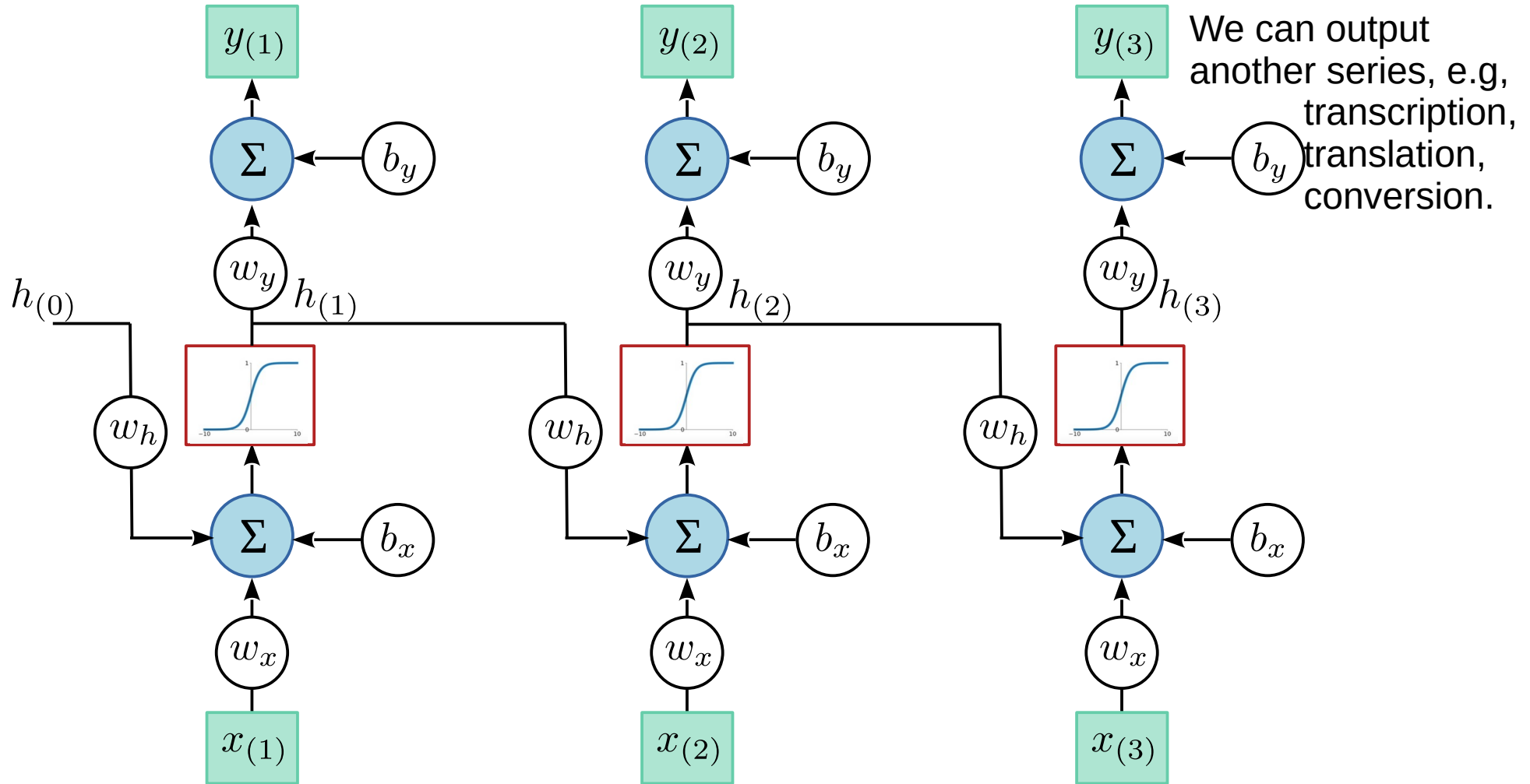NB: implicit "identity" activation function

$$y_{(t)} = w_y \times h_{(t)} + b_y$$

not the same timepoint

$$h_{(t)} = a(w_x \times x_{(t)} + b_x + w_h \times h_{(t-1)})$$
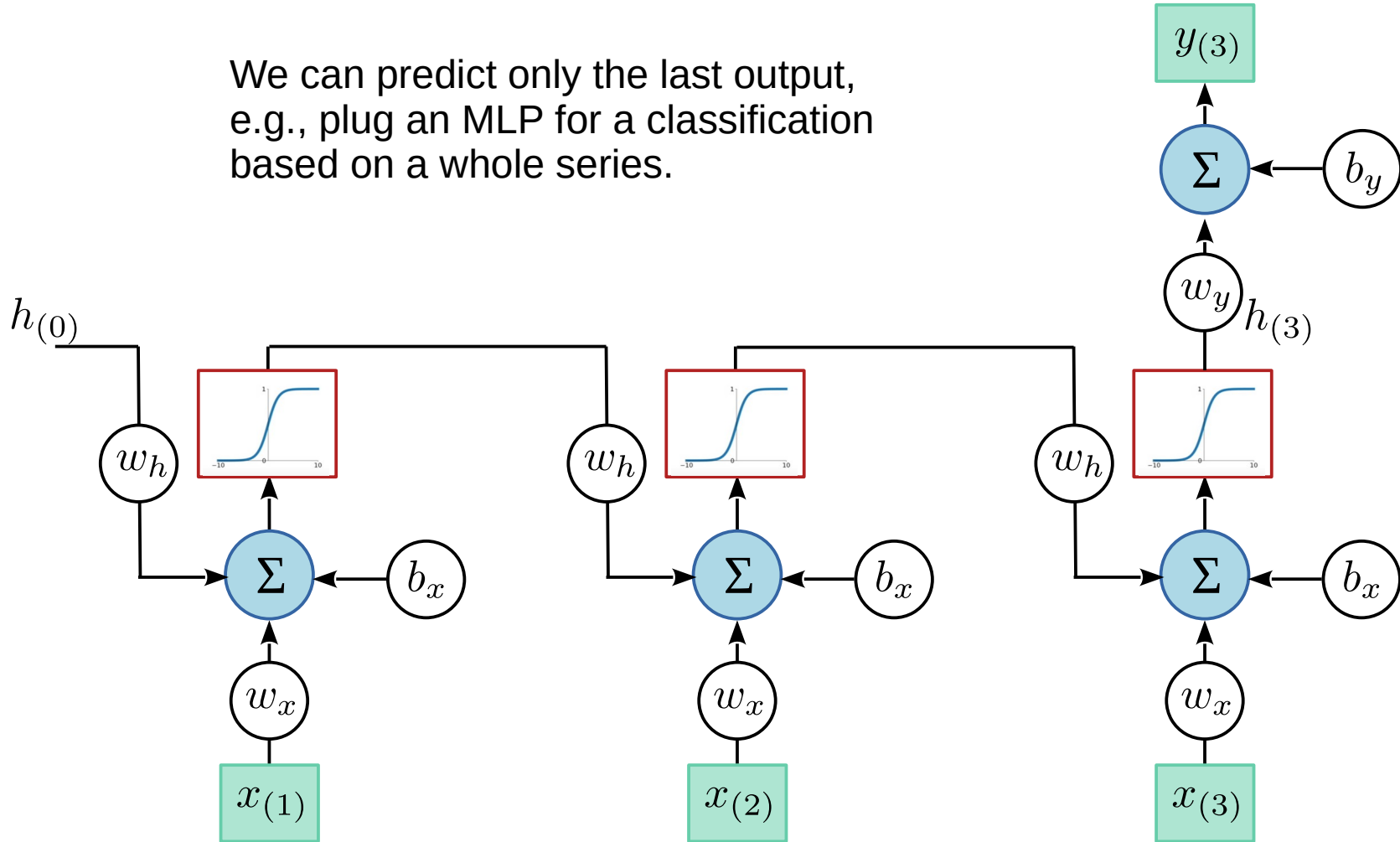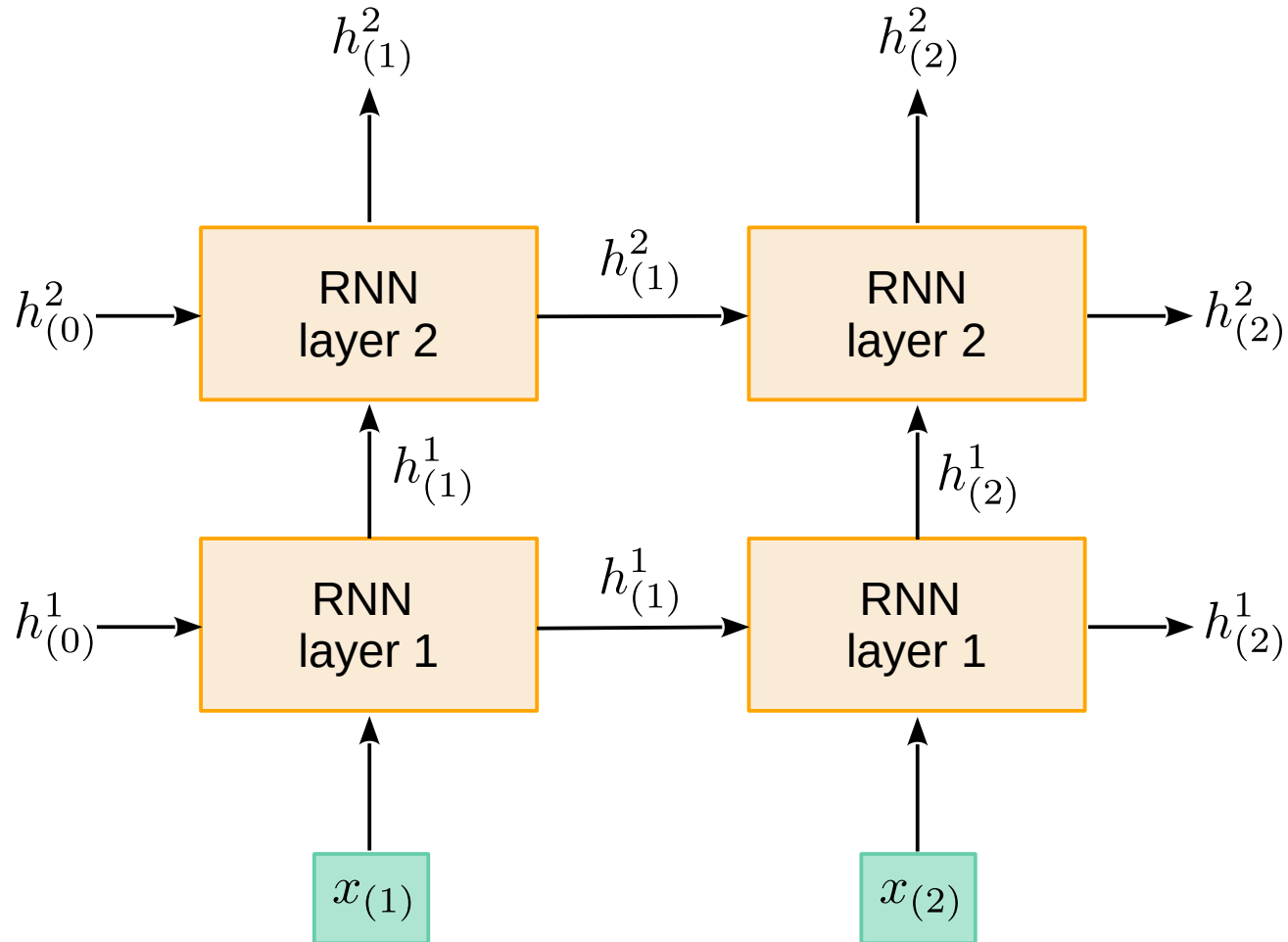
"memory"

# RNN: 1 cell - unfolding



We can output another series, e.g, transcription, translation, conversion.
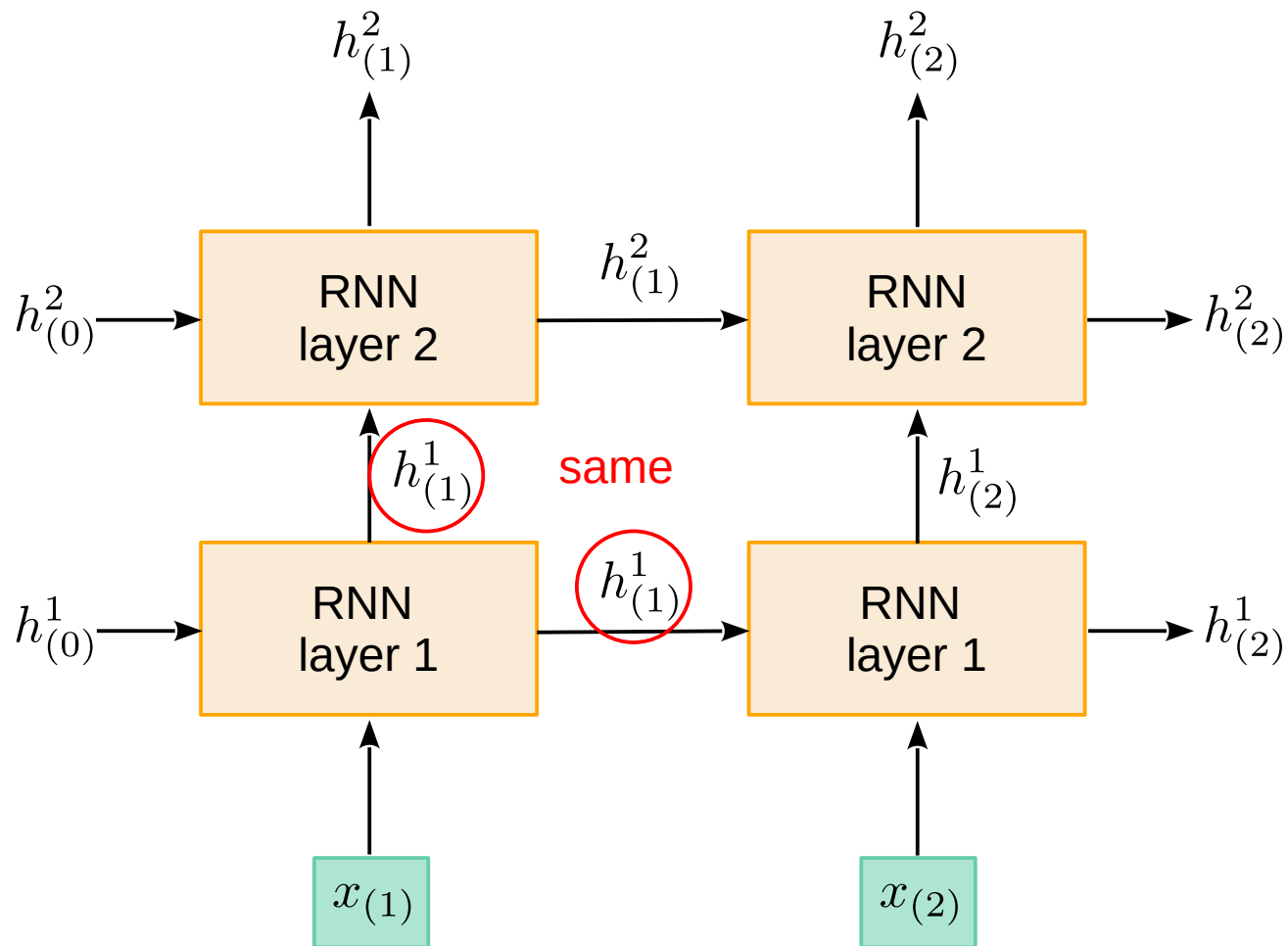
# RNN: 1 cell - unfolding



No feedback to the model

We can predict only the last output, e.g., plug an MLP for a classification based on a whole series.
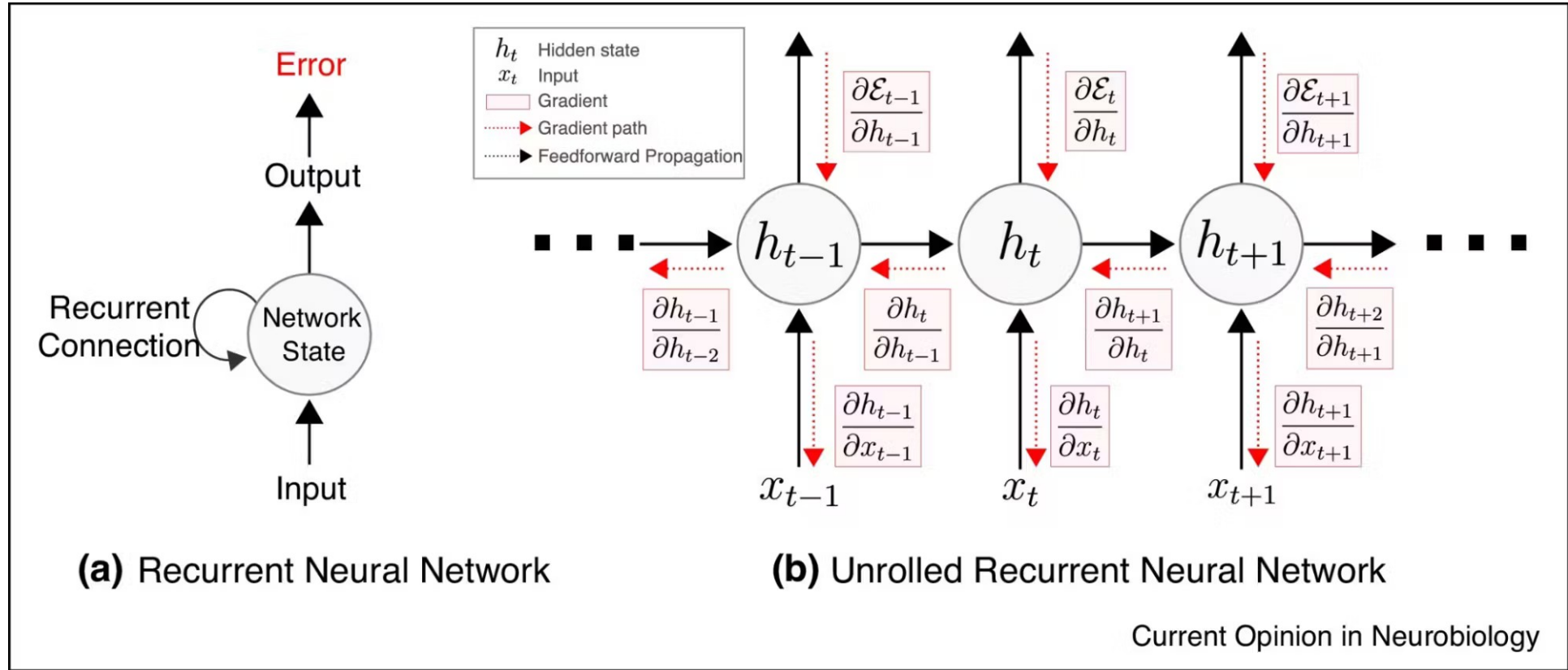
# Stacked RNNs

# Stacked RNN

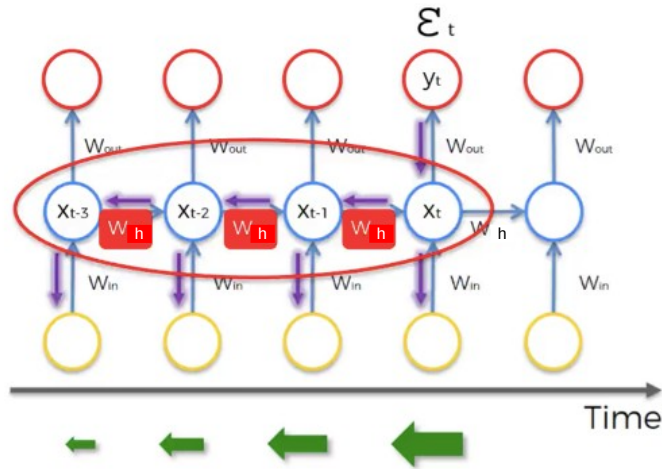# Several RNNs may learn different patterns in parallel

Same idea as different kernels in CNNs

No connexions between parallel cells of the same layer!

# Training RNNs by backpropagation in time



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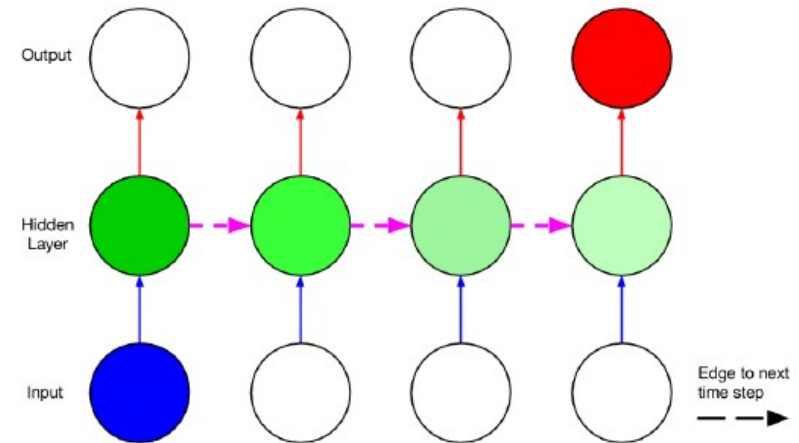
# Exploding and vanishing gradients

E.g.: activation function = ReLU

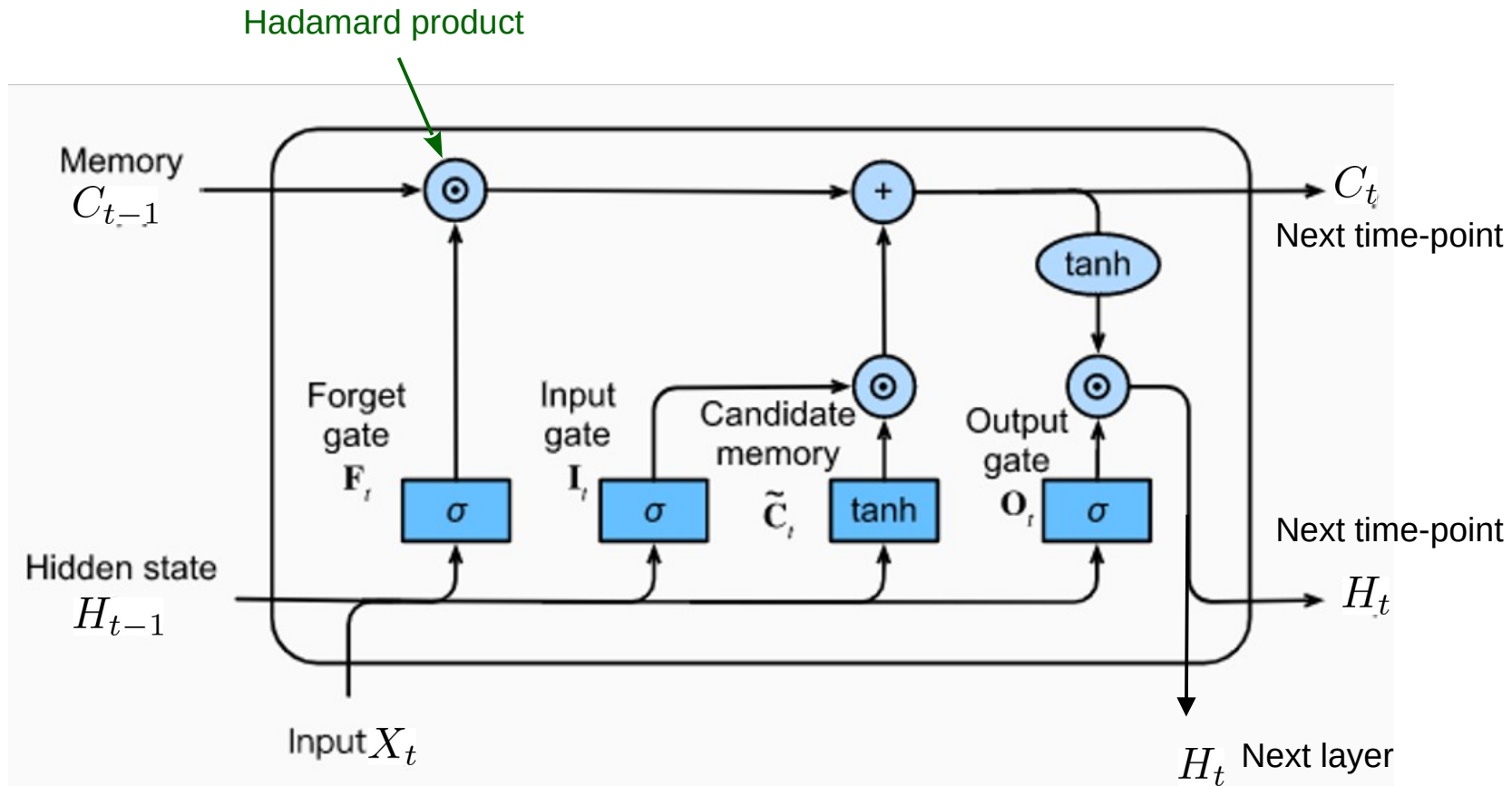$$\frac{\partial x_i}{\partial x_{i-1}} = w_h$$

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

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

$$\frac{\partial \mathcal{E}}{\partial \theta} = \sum_{1 \le t \le T} \frac{\partial \mathcal{E}_t}{\partial \theta}$$

$$\frac{\partial \mathcal{E}_t}{\partial \theta} = \sum_{1 \le k \le t} \left( \frac{\partial \mathcal{E}_t}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} \frac{\partial^+ \mathbf{x}_k}{\partial \theta} \right)$$

$$\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \prod_{t \ge i > k} \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_{i-1}} = \prod_{t \ge i > k} \mathbf{W}_{h \, rec}^T diag(\sigma'(\mathbf{x}_{i-1}))$$

W h ~ small ⟹ Vanishing

W h ~ large ⟹ Exploding

Green = sensitivity of output on input

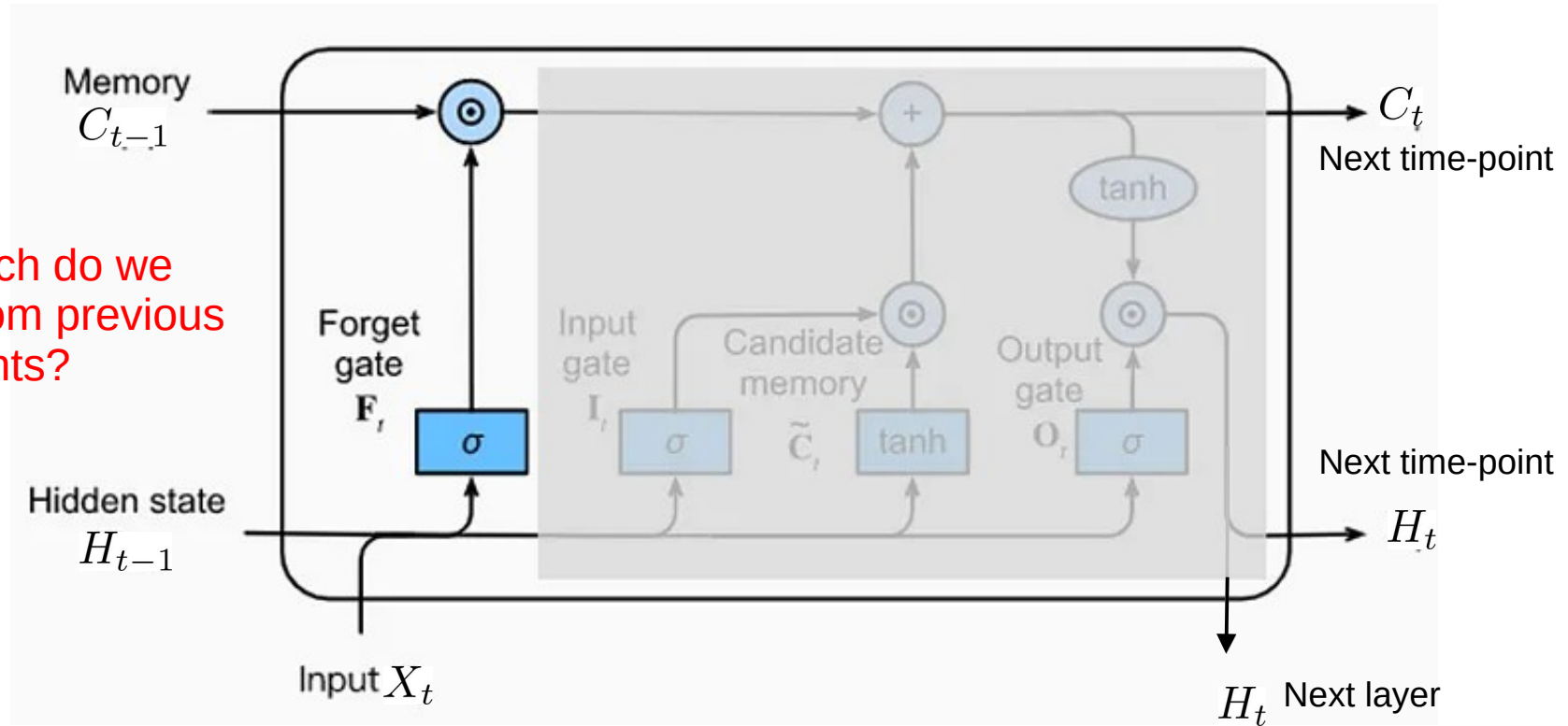# Solution: Long Short-Term Memory (LSTM)



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

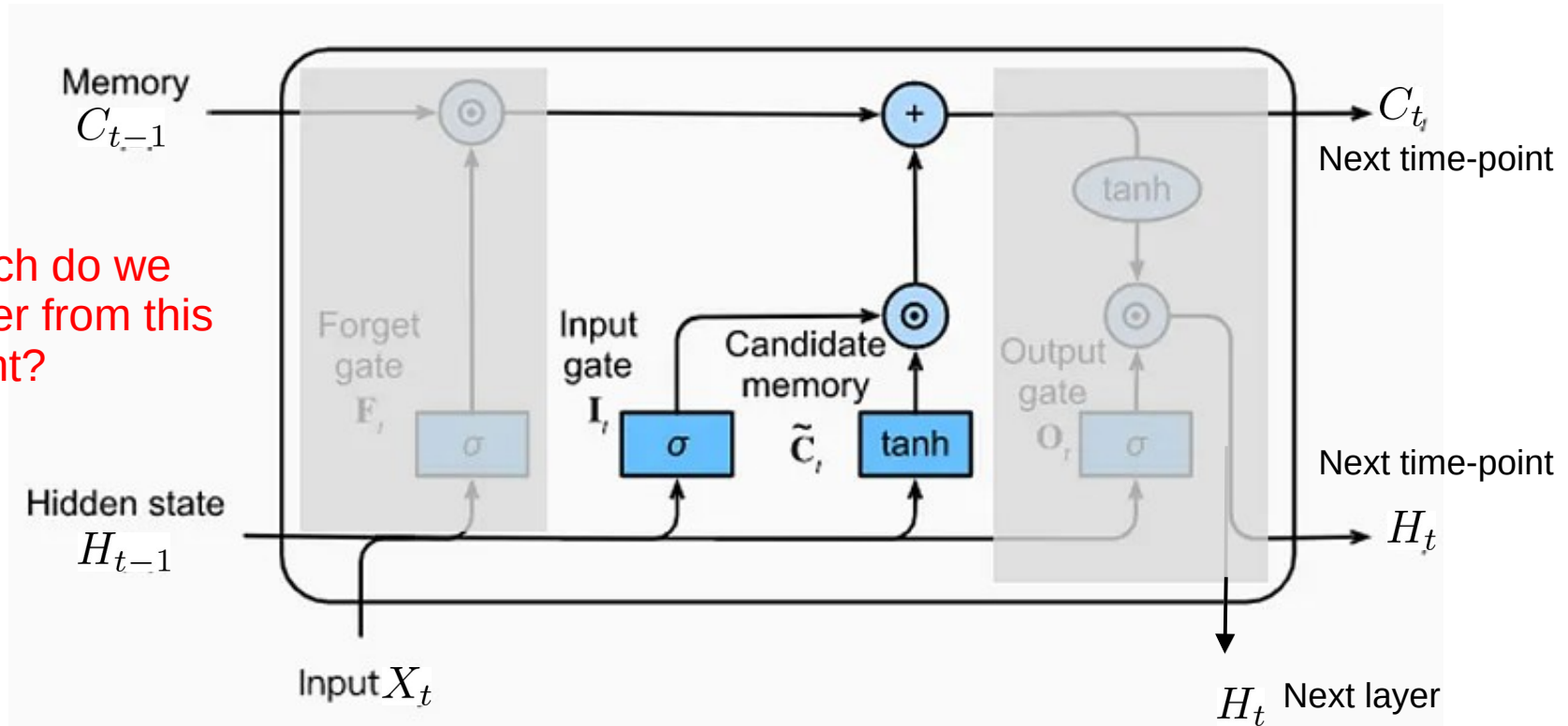Source: Ottavio Calzone (2002) An Intuitive Explanation of LSTM. https://medium.com/@ottaviocalzone

# Solution: Long Short-Term Memory (LSTM)

# Solution: Long Short-Term Memory (LSTM)



2-How much do we remember from this time-point?

3-How much do we transmit to the next LSTM cells?

Memory $C_{t-1}$

$C_{t}$
Next time-point

Long-term memory

tanh

3-How much do we transmit to the next LSTM cells?

Forget gate $F_t$

Input gate $I_t$

Candidate memory $\tilde{C}_t$

Output gate $O_t$

$\sigma$ $\sigma$ tanh $\sigma$

Short-term memory

Next time-point

Hidden state $H_{t-1}$

$H_t$

Input $X_t$

$H_t$ Next layer

# LSTMs for Encoder-Decoder



MSNovelist: de novo structure generation from mass spectra

Michael A. Stravs, Kai Dührkop, Sebastian Böcker & Nicola Zamboni ✉

*Nature Methods* **19**, 865–870 (2022) | Cite this article

# Examples in the biomedical domain

BMC Medical Informatics and
Decision Making
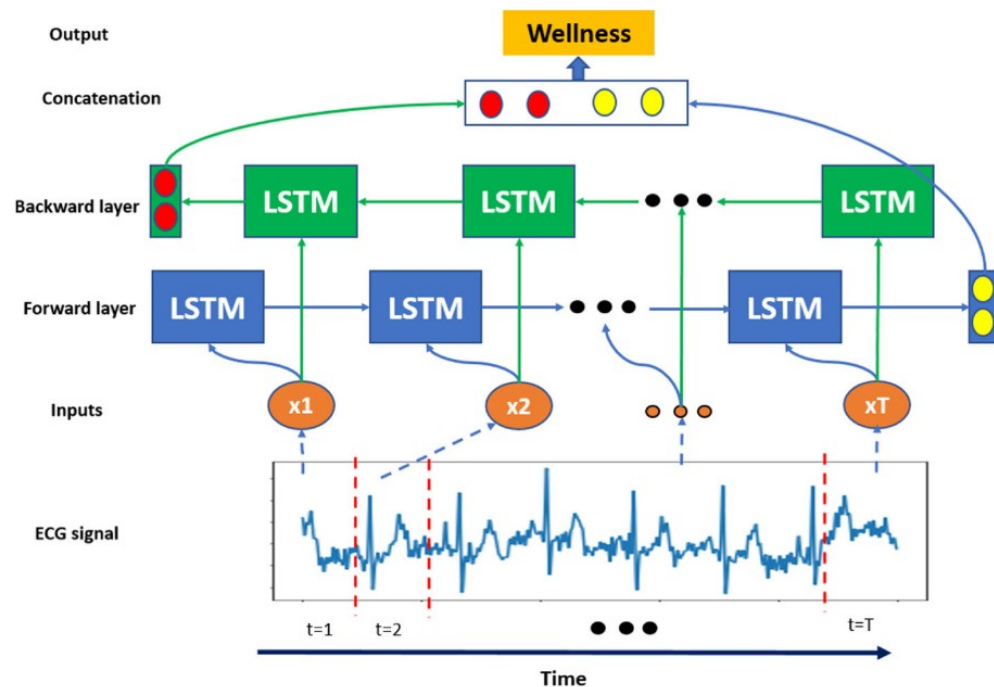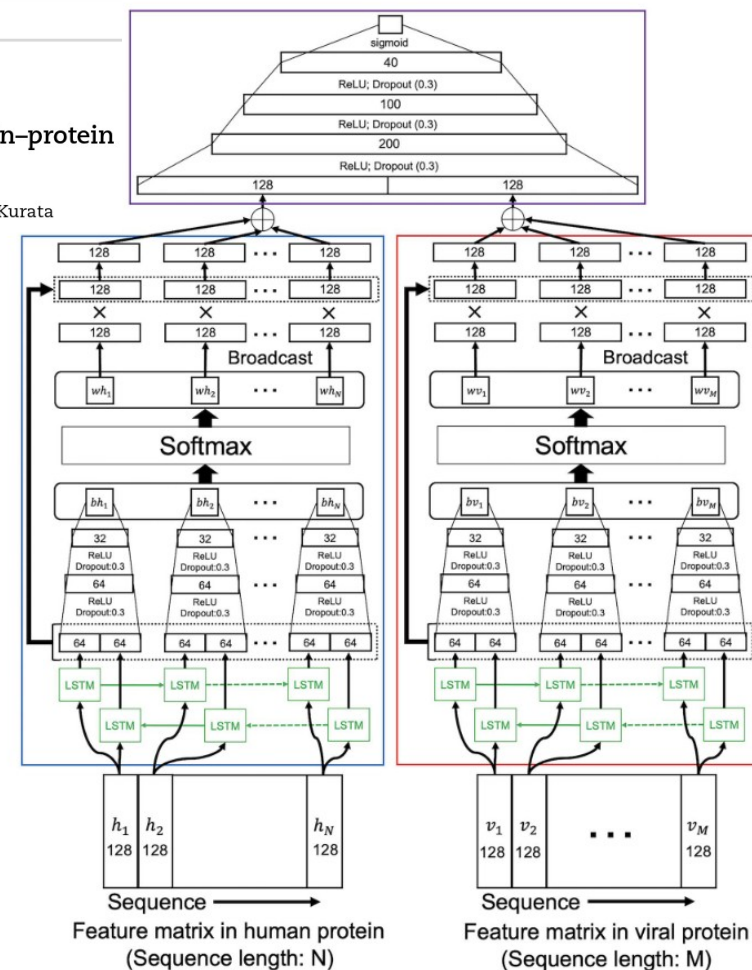
## Forecasting one-day-forward wellness conditions for community-dwelling elderly with single lead short electrocardiogram signals

Xiaomao Fan[1], Yang Zhao[2]* , Hailiang Wang[2] and Kwok Leung Tsui[1,2]

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

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

# Examples in the biomedical domain

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

# The paper that changed everything: the Transfomer

## 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 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.

## 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].

*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

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

All authors equal

*Equal contribution. Listing order is random.

Never published in a journal

# The paper that changed everything: the Transfomer

**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 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.

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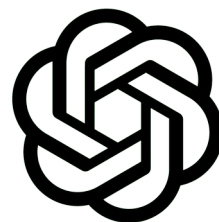
[*]Equal contribution. Listing order is random. Jakob proposed replacing RNNs with self-attention and started the effort to evaluate this idea. Ashish, with Illia, designed and implemented the first Transformer models and has been crucially involved in every aspect of this work. Noam proposed scaled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

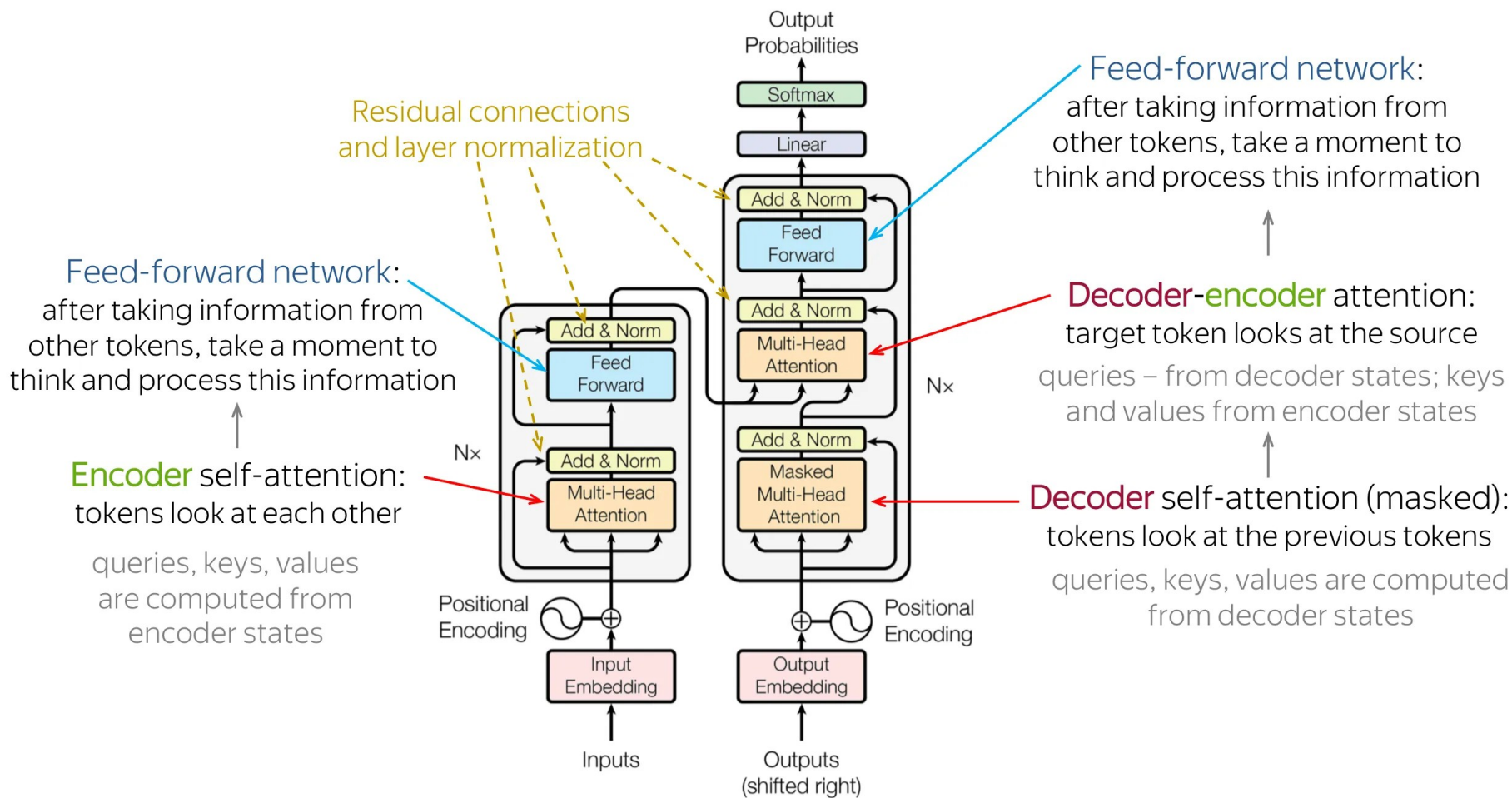[†]Work performed while at Google Brain.
[‡]Work performed while at Google Research.

# The Transformer: Memory + context = attention



**Feed-forward network:** after taking information from other tokens, take a moment to think and process this information

Residual connections and layer normalization

**Feed-forward network:** after taking information from other tokens, take a moment to think and process this information

**Decoder-encoder attention:** target token looks at the source

queries – from decoder states; keys and values from encoder states

**Encoder self-attention:** tokens look at each other

queries, keys, values are computed from encoder states

**Decoder self-attention (masked):** tokens look at the previous tokens

queries, keys, values are computed from decoder states

# The Transformer: Memory + context = attention

**I love AI**

Output Probabilities

**Residual connections and layer normalization**

Softmax

Linear

**Feed-forward network:** after taking information from other tokens, take a moment to think and process this information

**Feed-forward network:** after taking information from other tokens, take a moment to think and process this information

Add & Norm

Feed Forward

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

**Decoder-encoder attention:** target token looks at the source

queries – from decoder states; keys and values from encoder states

**Encoder self-attention:** tokens look at each other

queries, keys, values are computed from encoder states

Nx

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

**Decoder self-attention (masked):** tokens look at the previous tokens

queries, keys, values are computed from decoder states

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

**J'adore l'IA**

Inputs

Outputs (shifted right)

**<start> I love**

# Attention in the Transformer



Q = query, K = key, V = value

# Attention in the Transformer



Q = query, K = key, V = value

# Attention in the Transformer

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V = Z$$

## Scaled Dot-Product Attention

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q K V

## Multi-Head Attention

Linear

Concat

Scaled Dot-Product Attention — h

Linear Linear Linear

V K Q

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Masked Multi-Head Attention

Add & Norm

Multi-Head Attention

Nx

Nx

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

Q = query, K = key, V = value

# Attention in the Transformer

n = #tokens

d = embedding
dimension

Attention: aggregated attention

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V$$

Z

=

Attention matrix: Impact of each token
on all the others

n

n

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

Scaled Dot-Product Attention

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q   K   V

Multi-Head Attention

Linear

Concat

Scaled Dot-Product
Attention

Linear   Linear   Linear

V   K   Q

h

Add & Norm

Feed
Forward

Nx

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

Q = query, K = key, V = value

# Attention in the Transformer



Actual relevant attentions
1 per token

Attentions of all tokens on all tokens

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

The dot product between $Q$ and $K^T$
compute how aligned are the vectors
encoding two tokens (~cosine similarity)

# Self versus cross-attention

# 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



Pathways affected in cervical cancer

Random, GO BP, MSigDB hallmarks, Clustering

$X$ Gene expression, methylation, siRNA, etc.
e.g., TCGA cervical cancers

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

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)

# Cell2Sentence: gene names as token



Input single-cell data

Cell2Sentence and LLM fine-tuning

"CD4+ T-cell in Human PBMC CD3E CD4 IL2RA ... IFNG [MASK]"
"CD8+ T-cell in Breast Cancer CD8A GZMB PRF1 ... PD1 [MASK]"
"Monocyte in Lupus CD14 CCR2 TNF ... IL6 [MASK]"
"Hepatocyte in liver tissue ALB CYP3A4 CYP2E1 ... AFP [MASK]"
⋮
"Beta Cell in Type 1 Diabetes INS PDX1 GCK ... GLUT2 [MASK]"

Generate cell sentences via prompting

Model Input: "T-cell in multiple sclerosis"

Model Output: "CD3E CD4 IL17A IFNG TNF RORC TBX21 CCR6 CXCR3 STAT3 STAT4 FOXP3 GATA3 TIGIT PD1 CTLA4 ICOS CD28 CD25 CD127 IL2 IL7R IL12RB1 IL23R CD69 CD44 ..."

Generated single-cell data

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

Single-cells & multi-cells

800+ datasets, 57+ million cells

Biological annotations

CD8+ T cell
Lung tissue
Drug perturbation

Textual information

C2S

Cell generation

Single-cell | Multiple cells
IGKV4 MALAT1 ...

Label prediction
Most probable cell type : Monocyte

Perturbation prediction
Based on control cell, perturbed cell expression will be:
CD74 MALAT1 MT-CO1...

Generate insights
This single-cell data likely represents...

**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], Laust Hvas Mortensen[5,7,8], Ewan Birney[2], Tom Fitzgerald[2] & Moritz Gerstung[1,2,9,10,11,12]

Decision-making in healthcare relies on understanding patients' past and current health states to predict and, ultimately, change their future course[1-3]. 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 |
| 72.1: | J80 adult respiratory distress syndrome |
| 72.2: | No event |
| 72.7: | Death |

# Vision Transformer



CNN = local features
ViT = relations between distant features

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

# Patches are embedded by CNNs



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)

6 nearest neighbours in
the 4096 dimension space

input
image

# ViTs are replacing vanilla CNNs



brain cancer

lung cancer

non-small cell lung cancer

oral cancer

bladder cancer

skin cancer

# Most biological knowledge comes as graphs

# Graph Neural Networks (GNNs)

## CNNs

### The Graph Neural Network Model

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



Regular grid (same
number of neighbours)
Homogeneous kernels

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

# GNN can be heterogeneous

Building
Training
Explanation



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

# Example of GNN convolution



$$A \quad B \quad C \quad D$$

$$D = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{matrix} A \\ B \\ C \\ D \end{matrix}$$

Degree matrix

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

Adjacency matrix

$$L = D - A$$

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

Laplacian matrix

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

EPFL, Lausanne, Switzerland
{michael.defferrard,xavier.bresson,pierre.vandergheynst}@epfl.ch

# Example of GNN convolution



A   B   C   D

$$D = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{matrix} A \\ B \\ C \\ D \end{matrix}$$

Degree matrix

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

Adjacency matrix

$$L = D - A$$

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

Laplacian matrix

identity matrix
no influence from neighbours

influence from
immediate neighbours

influence from neighbours and
neighbours of neighbours

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

$$y = p_w(L)x$$

$$w = [w_0, w_1, w_2, \dots, w_k] \quad \text{kernel de convolution}$$

**Convolutional Neural Networks on Graphs
with Fast Localized Spectral Filtering**

Michaël Defferrard    Xavier Bresson    Pierre Vandergheynst

EPFL, Lausanne, Switzerland
{michael.defferrard,xavier.bresson,pierre.vandergheynst}@epfl.ch

# Example of GNN convolution

D only influences C

D influences A, B, C



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

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

$w = [1, 0.1, 0.01]$

layer 1

layer 2

# Example of GNN convolution

D only influences C

D influences A, B, C

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

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

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

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

layer 1

layer 2

$$\sum \text{layer1} = \sum \text{layer2}$$

# 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



Insights from GNN
→ cluster-specific molecular and ontology subgraphs

Training objective
→ disease clusters

# Highly accurate protein structure prediction with AlphaFold

John Jumper[1,4✉], Richard Evans[1,4], Alexander Pritzel[1,4], Tim Green[1,4], Michael Figurnov[1,4], Olaf Ronneberger[1,4], Kathryn Tunyasuvunakool[1,4], Russ Bates[1,4], Augustin Žídek[1,4], Anna Potapenko[1,4], Alex Bridgland[1,4], Clemens Meyer[1,4], Simon A. A. Kohl[1,4], Andrew J. Ballard[1,4], Andrew Cowie[1,4], Bernardino Romera-Paredes[1,4], Stanislav Nikolov[1,4], Rishub Jain[1,4], Jonas Adler[1], Trevor Back[1], Stig Petersen[1], David Reiman[1], Ellen Clancy[1], Michal Zielinski[1], Martin Steinegger[2,3], Michalina Pacholska[1], Tamas Berghammer[1], Sebastian Bodenstein[1], David Silver[1], Oriol Vinyals[1], Andrew W. Senior[1], Koray Kavukcuoglu[1], Pushmeet Kohli[1] & Demis Hassabis[1,4✉]

~93 million parameters (weights+biases)

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

# AlphaFold2: evoformer



Transformers+GNNs

between residues of one sequence

between sequences of the MSA

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



Layer 0 = local contacts

Head 2 of layer 2 recognises disulphide bonds

At layer 11, heads recognise distant contacts

Source: AlphaFold2 paper supplementary info

# RNNs are back. Rise of the Mamba



"Attention" = embedding size x input length
→ <u>linear growth with input length</u>
(not quadratic like transformers)

Prediction/classification

SSM = State Space Model

The model learns about variants' Interactions

→ Reusable foundation model

...AGGCTAGGATATCGATAAGCTGACTGAT...

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

# Mamba everywhere

## 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

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

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

Ali Hatamizadeh, Jan Kautz
NVIDIA
{ahatamizadeh, jkautz}@nvidia.com

## 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

## Graph Mamba: Towards Learning on Graphs with State Space Models

Ali Behrouz, Farnoosh Hashemi

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

Chloe Wang, Oleksii Tsepa, Jun Ma, Bo Wang
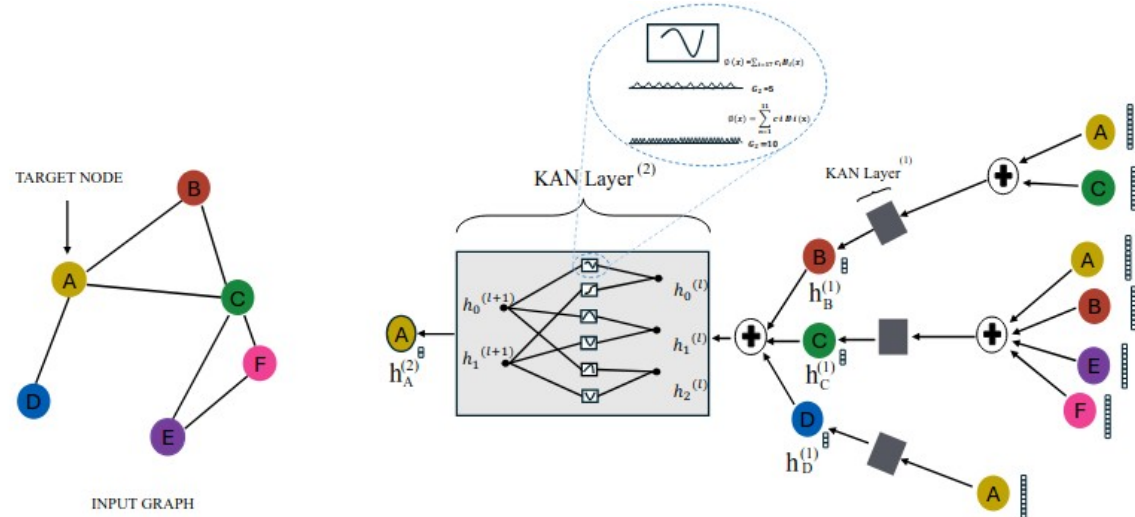
# Kolmogorov Arnold Networks



Liu *et al.*(2024) KAN: Kolmogorov-Arnold Networks. https://doi.org/10.48550/arXiv.2404.19756

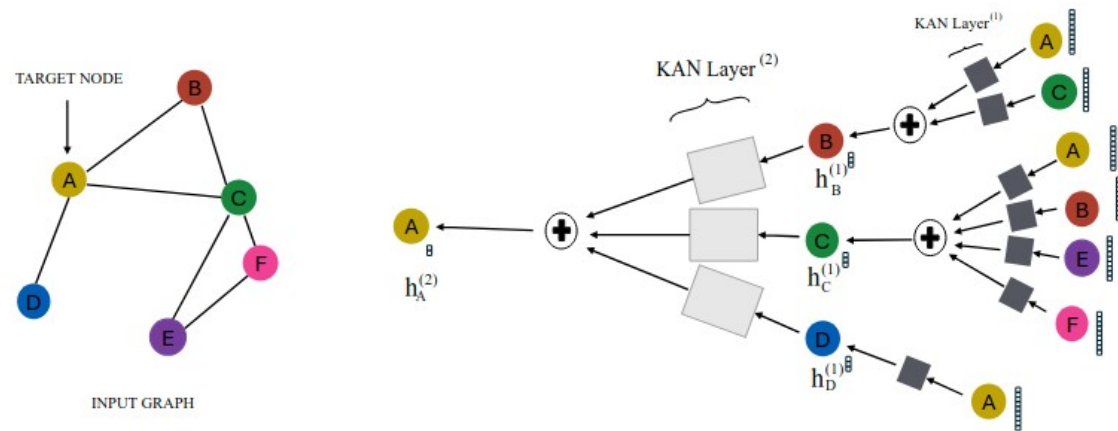# KANs to extract symbolic formulas

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

# How many architectures can you cram into one model?

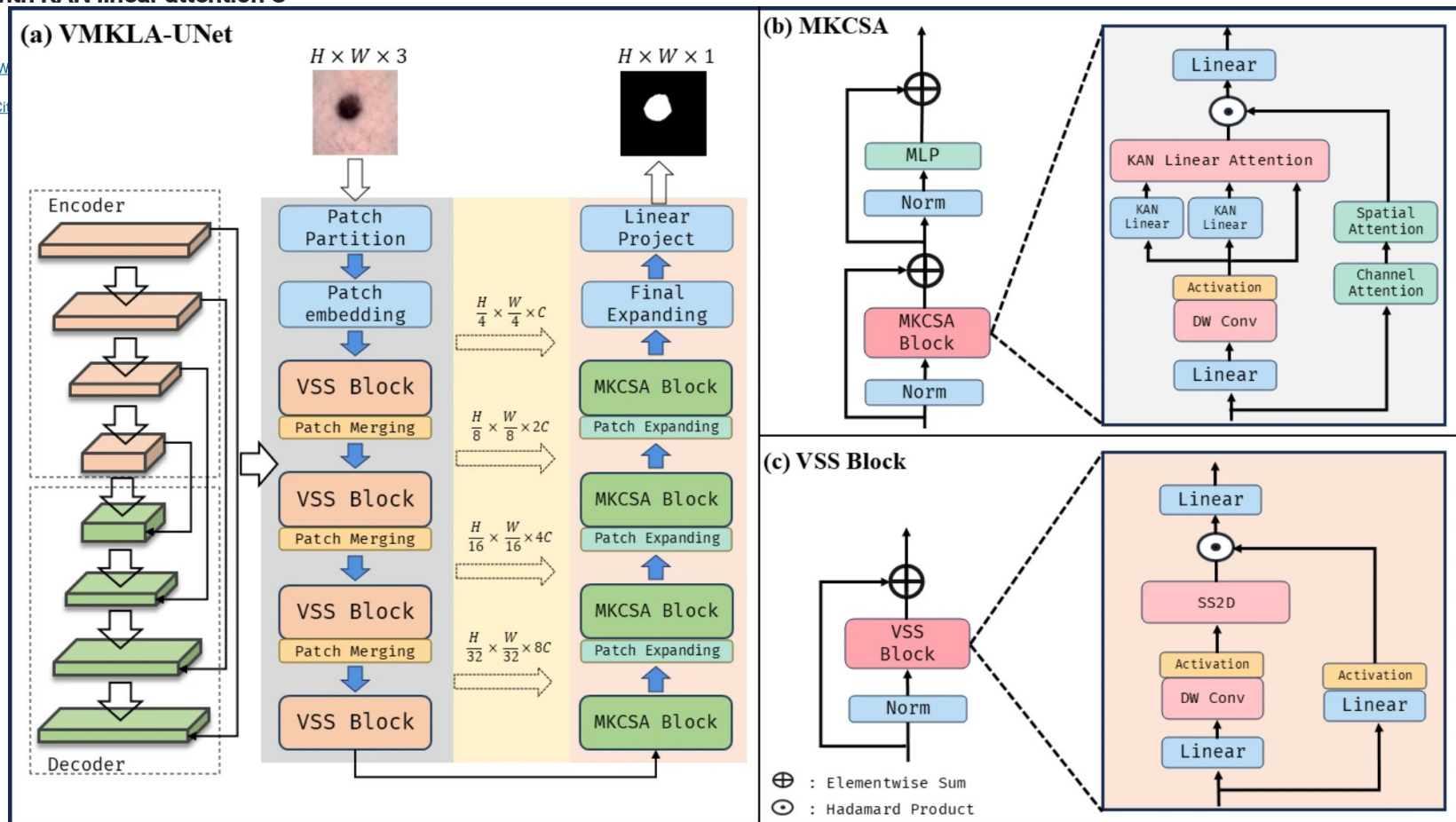## VMKLA-UNet: vision Mamba with KAN linear attention U-Net

Chenhong Su, Xuegang Luo, Shiqing Li, Li Chen & Juan W
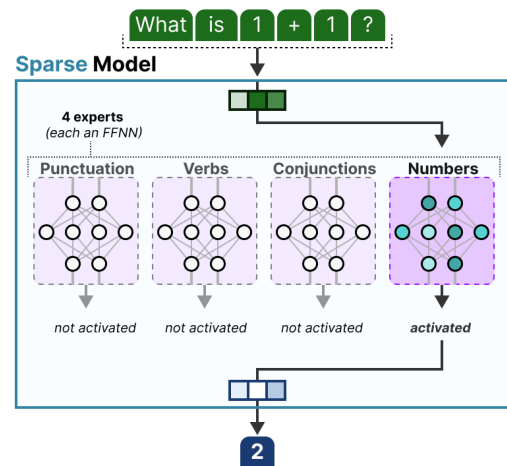
# Topics to explore



## 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/