

# Machine learning for physicists

<https://github.com/wangleiphy/ml4p>

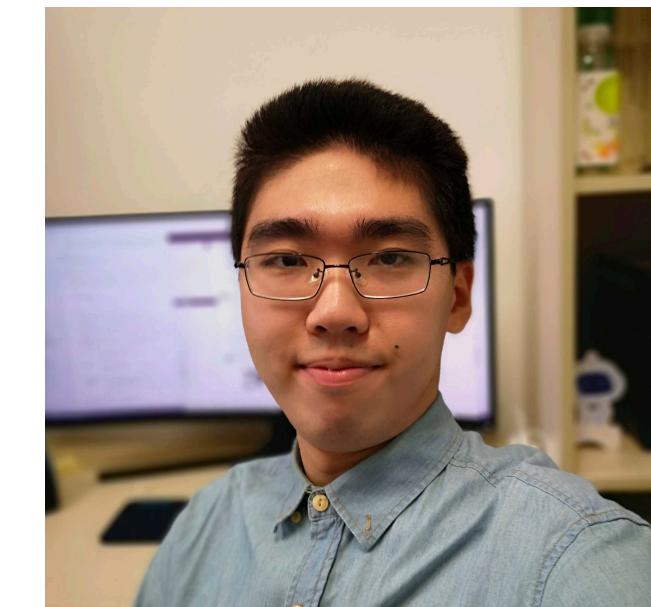
每周四上午10点

课程微信群

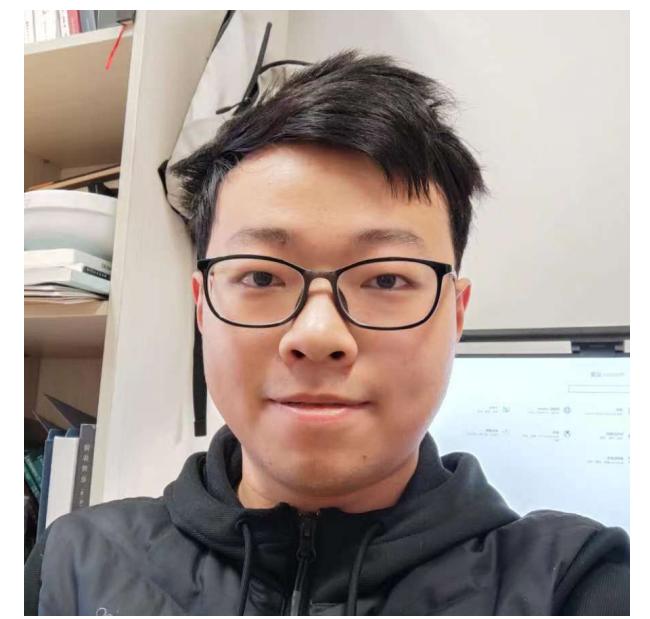
2.23	Overview
3.2	Machine learning practices
3.9	A hitchhiker's guide to deep learning
3.16	Research projects hands-on
3.23	Symmetries in machine learning
3.30	Differentiable programming
4.6	Generative models-I
4.13	Generative models-II
4.20	Research projects presentation
4.27	AI for science: why now ?



助教



李子航



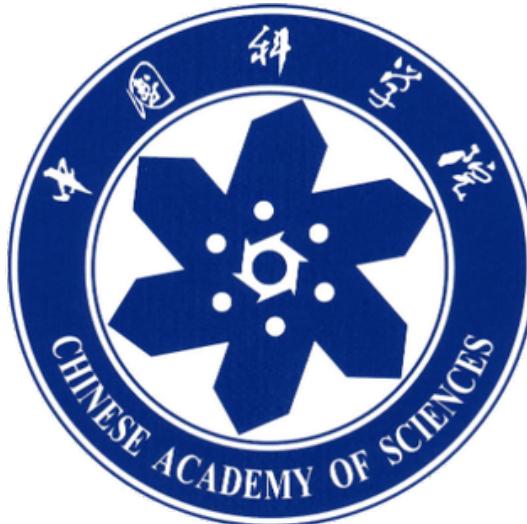
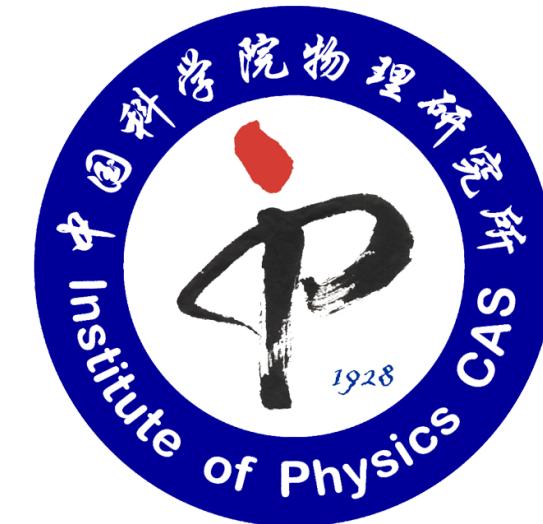
李扬帆

考核方式: project + presentation (1学分)

# A hitchhiker's guide to deep learning

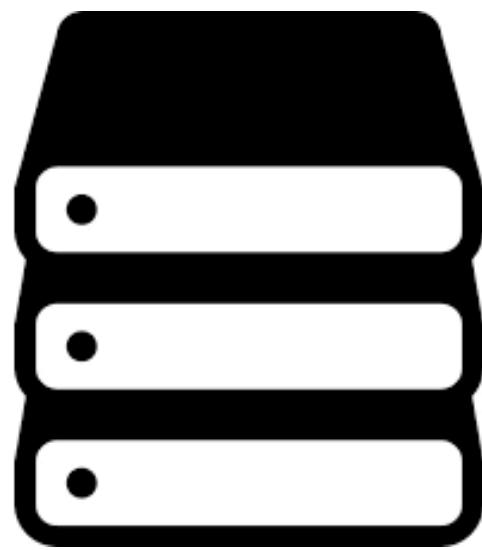
Lei Wang (王磊)

Institute of Physics, CAS  
<https://wangleiphy.github.io>

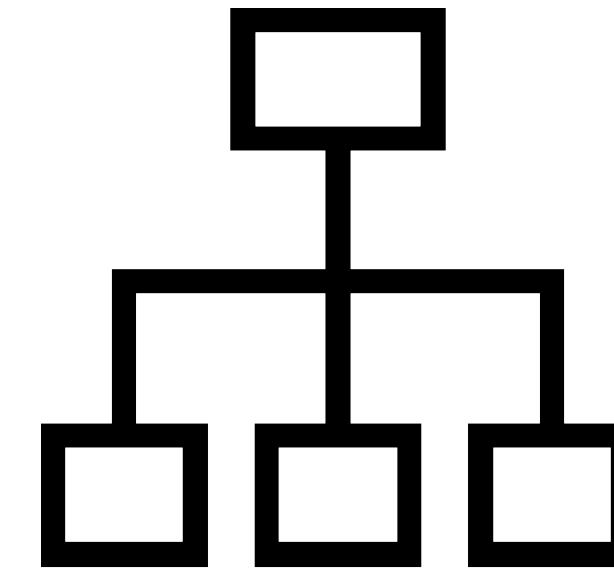


# Four pillars

Data



Model



Cost function



Optimization

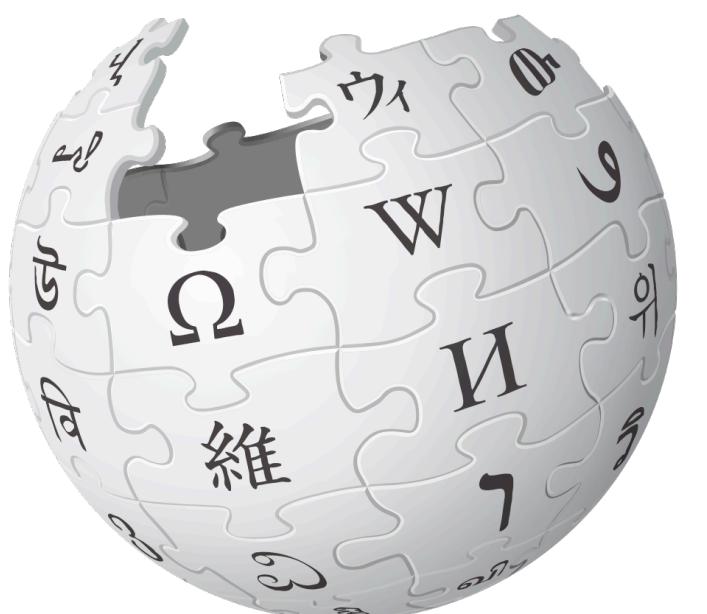
$\hat{\partial}$

# Data modalities

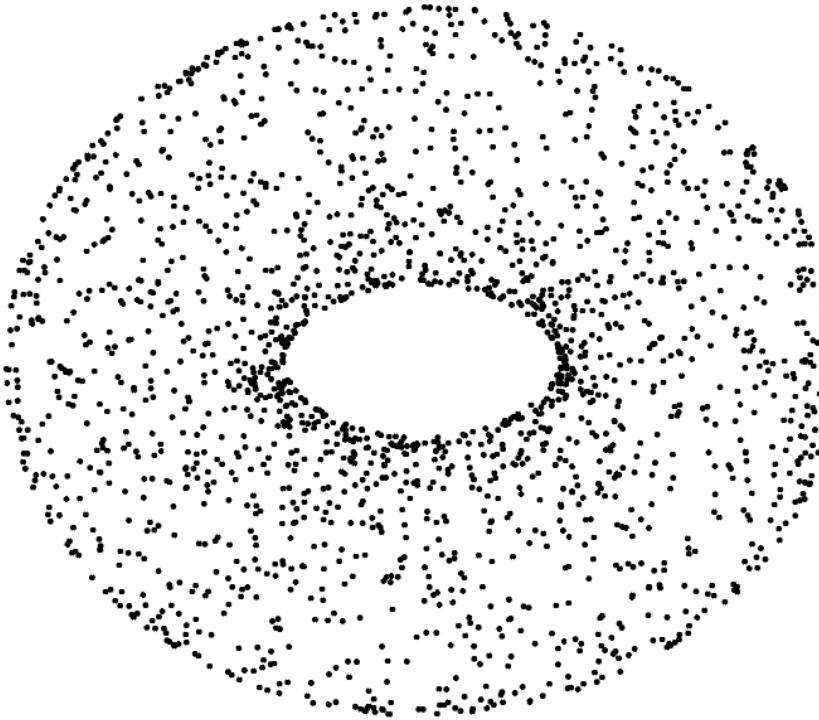
Image

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

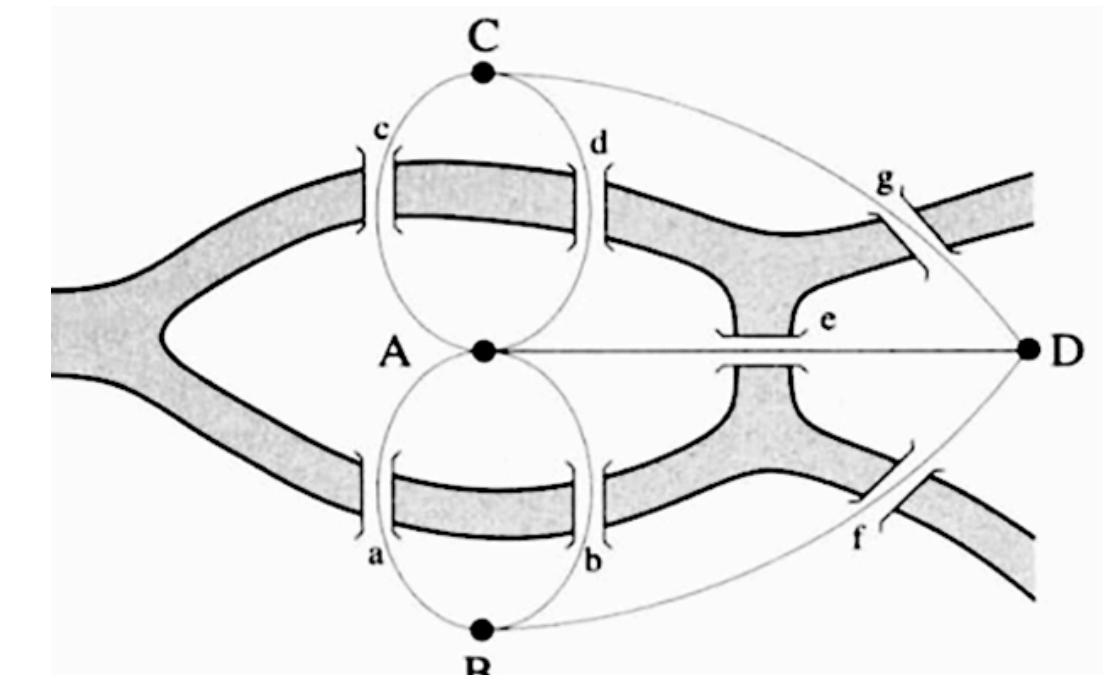
Language



Point cloud

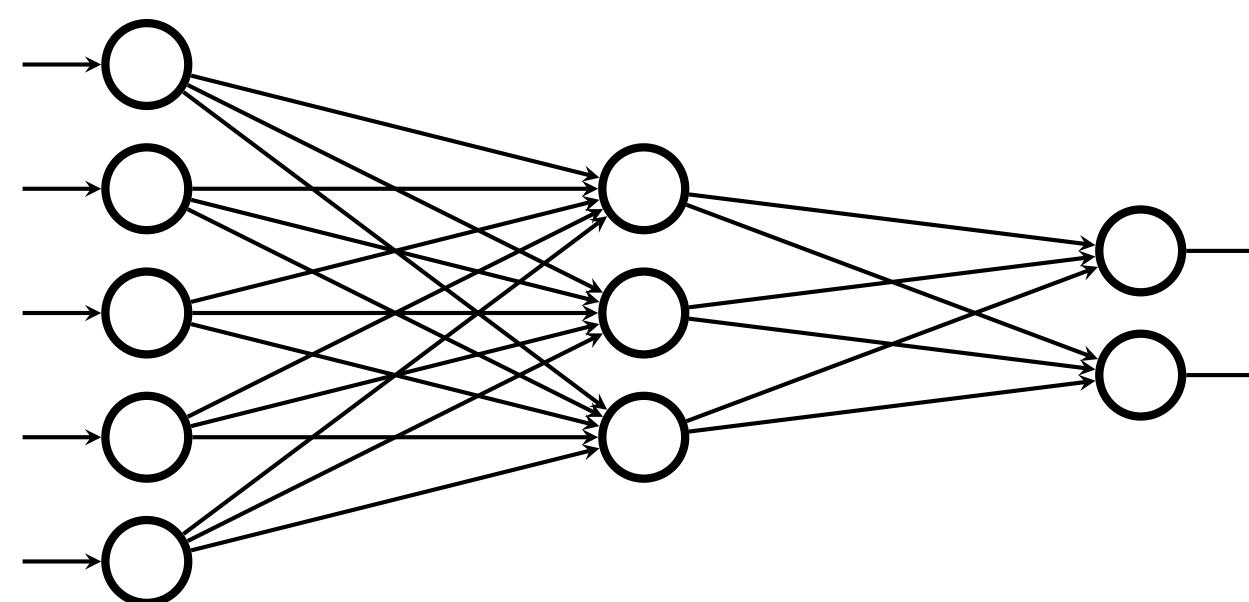


Graph

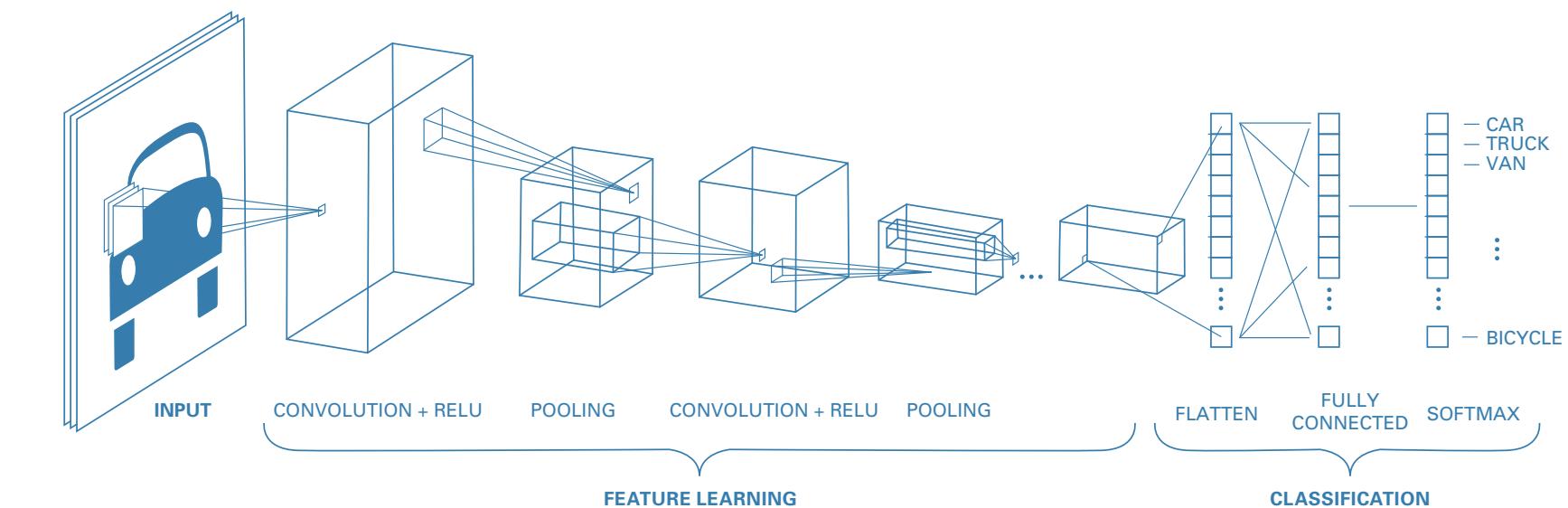


# Models

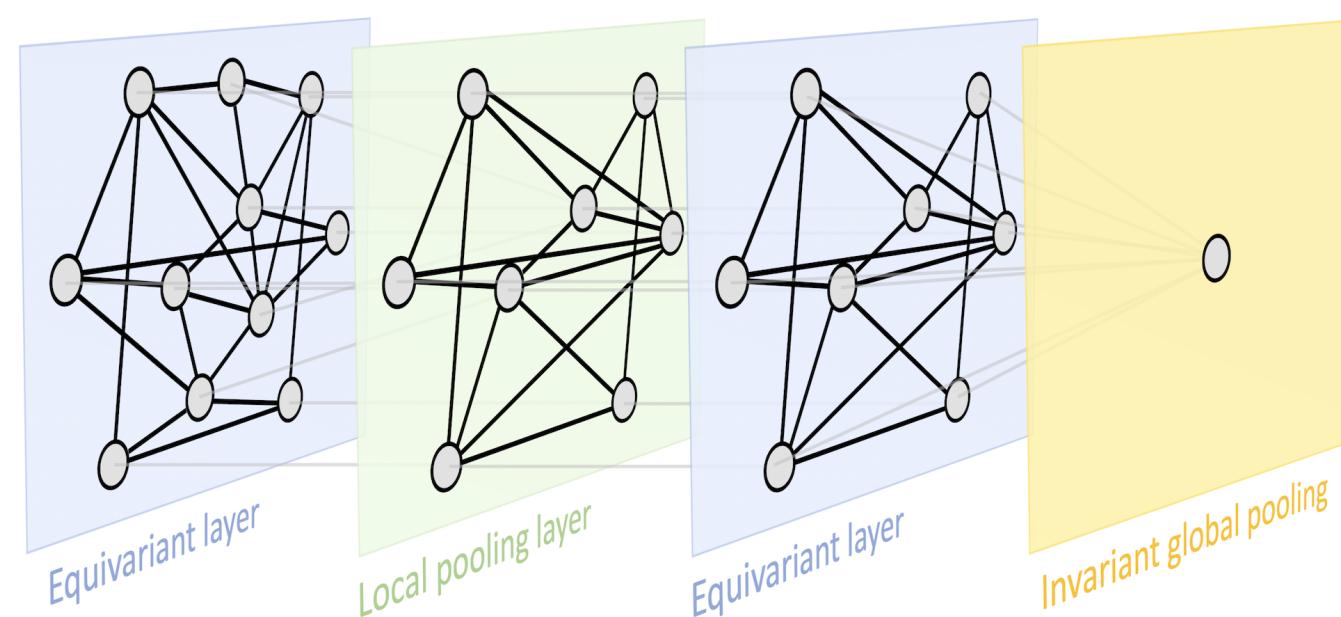
Multi-layer perceptron



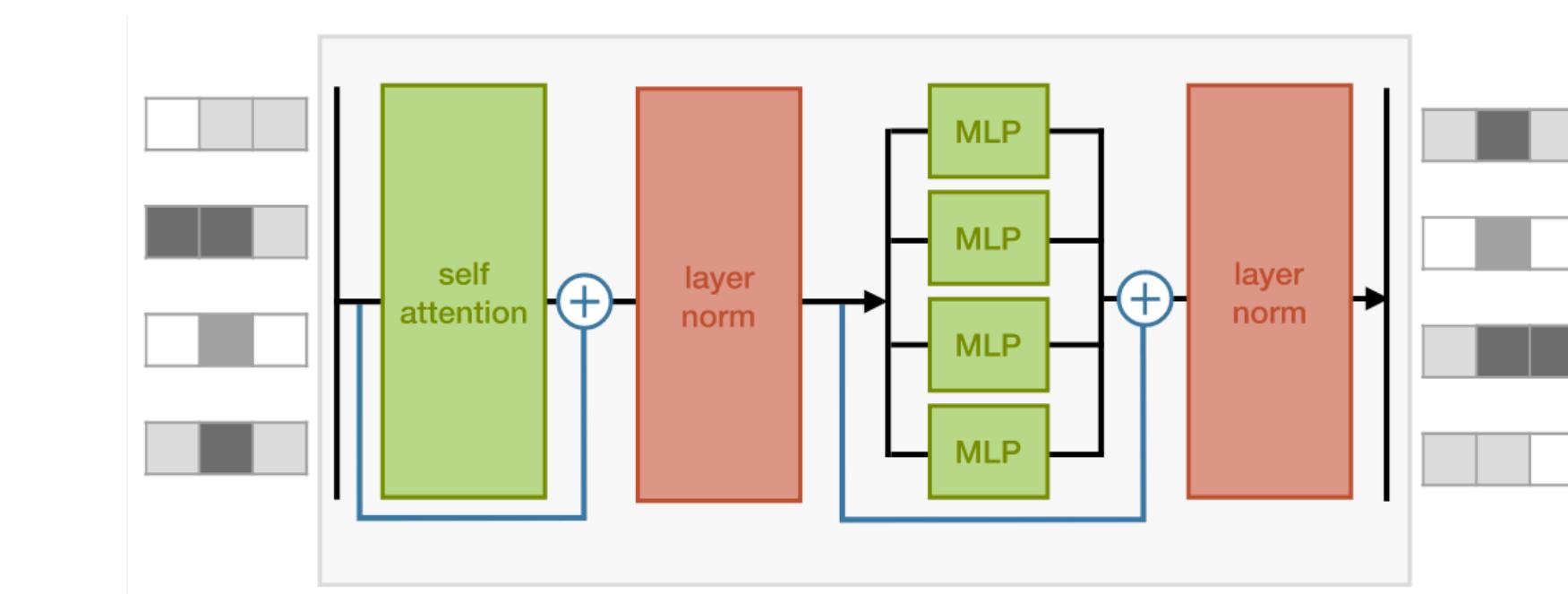
Convolutional neural network



Graph neural network



Transformer



# Multiple-layer perceptrons

$$y = \sigma(Wx + b)$$

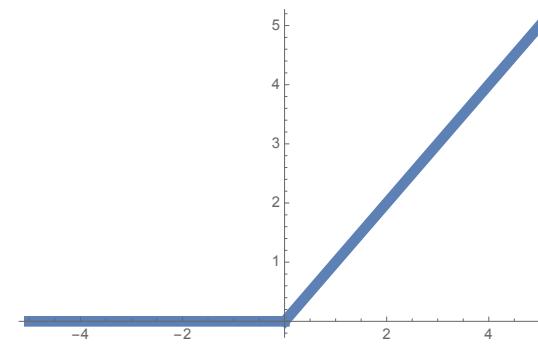
$$x \in \mathbb{R}^{F_x}$$

$$W \in \mathbb{R}^{F_y \times F_x}$$

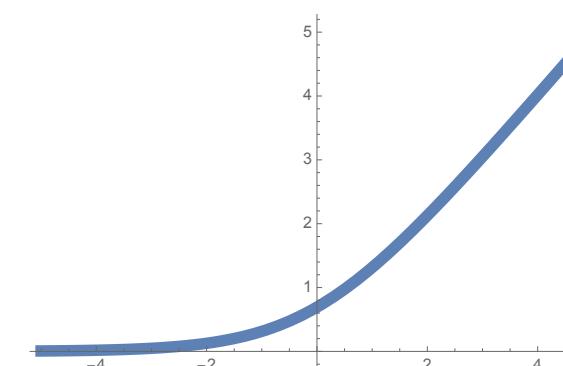
$$b \in \mathbb{R}^{F_y}$$

$\sigma$ : elementwise activation

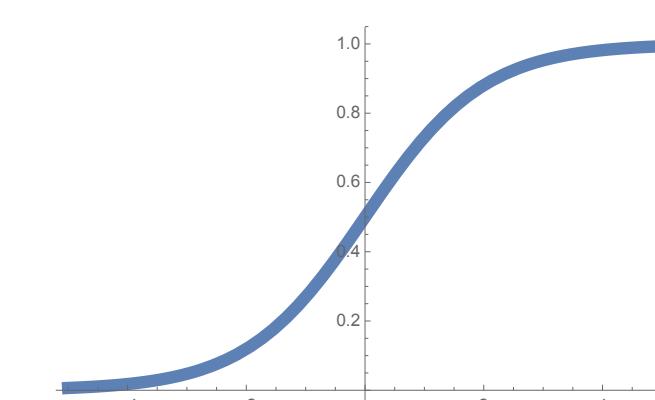
ReLU =  $\max(0, x)$



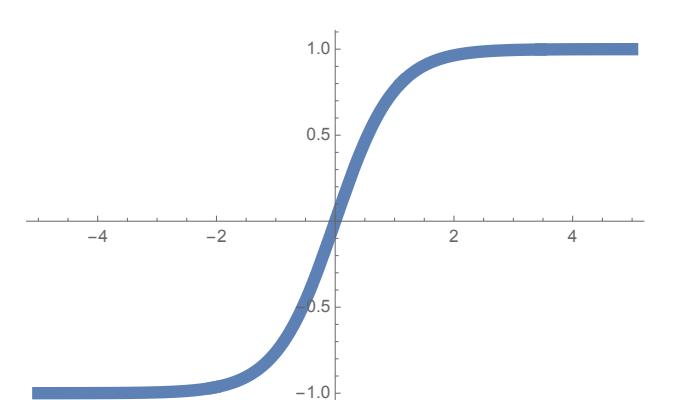
Softplus =  $\ln(1 + e^x)$



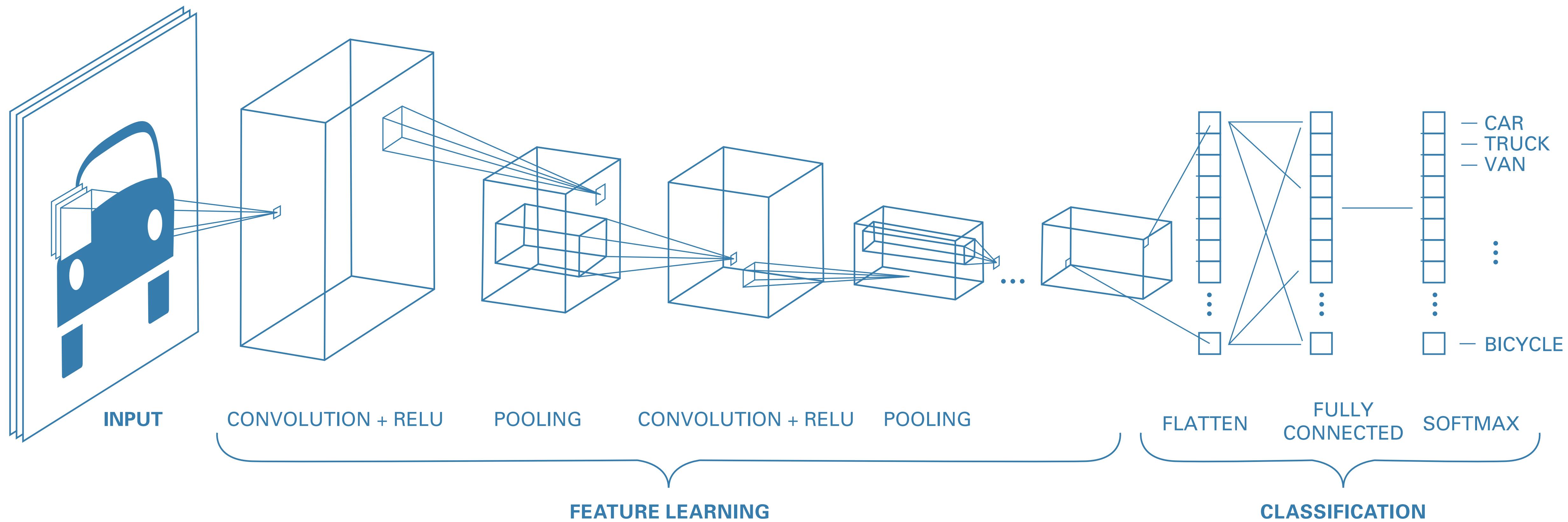
Sigmoid =  $1/(1 + e^{-x})$



tanh



# Convolutional neural networks



# 2d convolution (cross-correlation)

$$y_{ij} = (W \star x)_{ij} = \sum_{m=0}^{K-1} \sum_{n=0}^{K-1} W_{mn} x_{i+m, j+n}$$

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

$$\begin{aligned}
 & \begin{pmatrix} w_1 & w_2 \\ w_3 & w_4 \end{pmatrix} \star \begin{pmatrix} x_1 & x_2 & x_3 \\ x_4 & x_5 & x_6 \\ x_7 & x_8 & x_9 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \end{pmatrix} \\
 & \equiv \begin{pmatrix} w_1 & w_2 & 0 & w_3 & w_4 & 0 & 0 & 0 & 0 \\ 0 & w_1 & w_2 & 0 & w_3 & w_4 & 0 & 0 & 0 \\ 0 & 0 & 0 & w_1 & w_2 & 0 & w_3 & w_4 & 0 \\ 0 & 0 & 0 & 0 & w_1 & w_2 & 0 & w_3 & w_4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \end{pmatrix}
 \end{aligned}$$

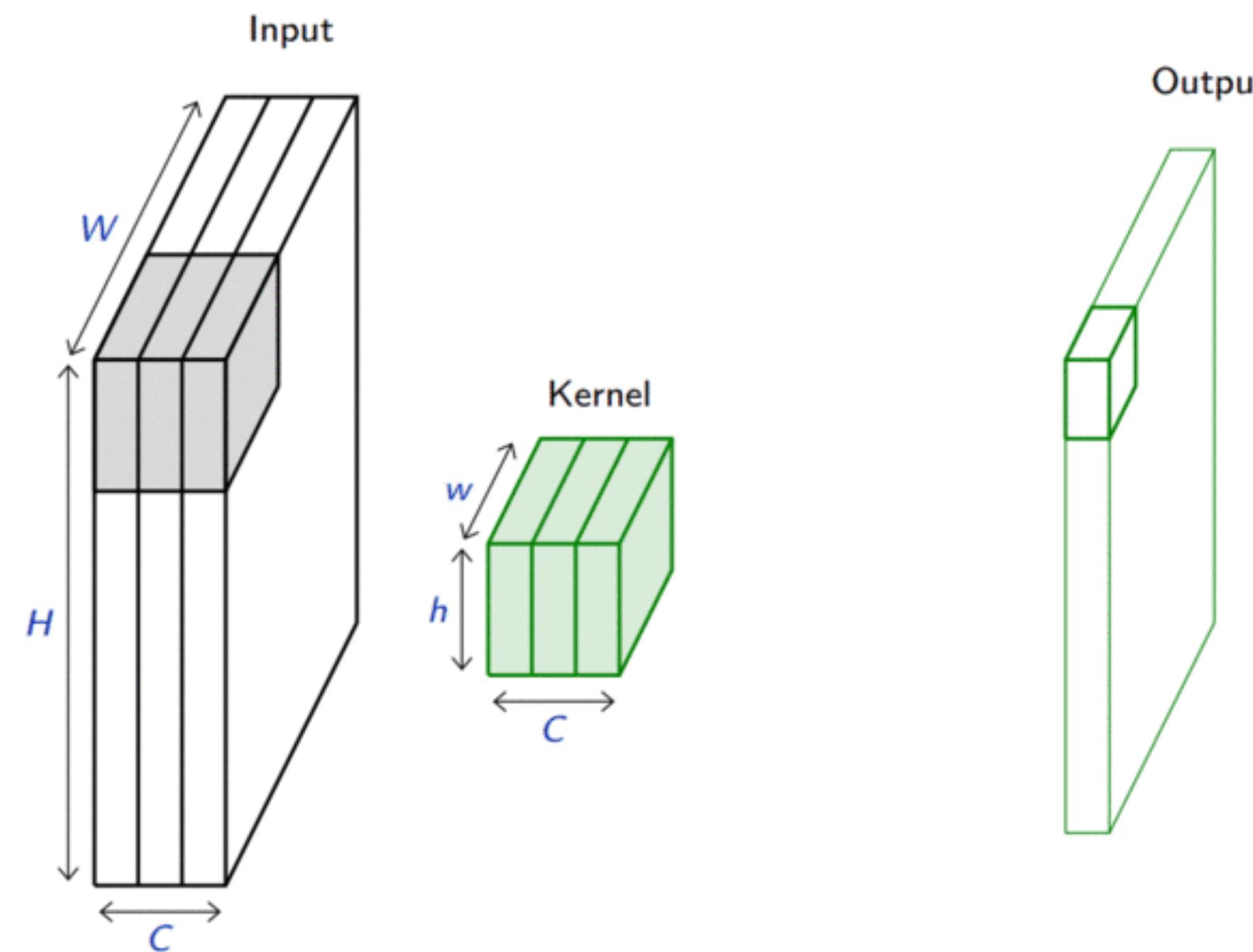
Equivalent to multiplying a sparse (doubly block circulant) weight matrix to flattened feature vector

Check [1603.07285](#) for more variations in padding, stride, dilation, ...

# Convolutional layer

$$y_{i,j,\mu} = \sum_{m=0}^{K-1} \sum_{n=0}^{K-1} \sum_{\nu=0}^{F_x-1} W_{m,n,\mu,\nu} x_{i+m,j+n,\nu} + b_\mu$$

$$x \in \mathbb{R}^{H \times W \times F_x} \quad W \in \mathbb{R}^{K \times K \times F_y \times F_x} \quad b \in \mathbb{R}^{F_y}$$



# Why is convolution useful?



Cat



Cat

$$\text{shift} \circ \text{conv} = \text{conv} \circ \text{shift}$$

Translational equivariance

(Ignore edge effect: infinite size or on a torus)

# Pooling layer

0	1	2
3	4	5
6	7	8

2 x 2 Max  
Pooling

4	5
7	8

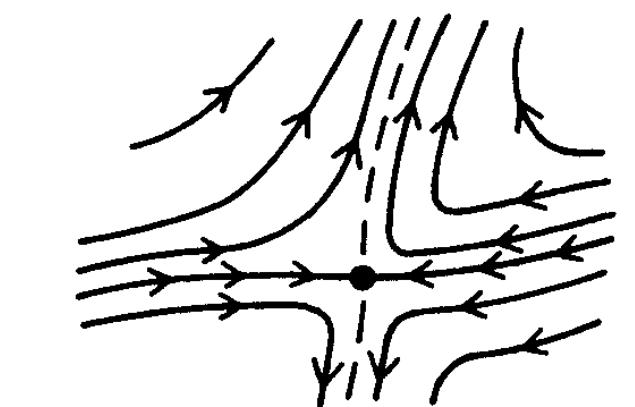
$$y_{i,j,\mu} = \max_{0 < m < K} \max_{0 < n < K} x_{i+m,j+n,\mu}$$

Detects a signal irrespective of its exact location

Kadanoff's block spin  
transformation

↑	↑	↓	↓	↓	↑	↑	↑
↑	↓	↓	↑	↓	↓	↑	↑
↓	↓	↓	↓	↓	↑	↑	↑
↑	↓	↓	↓	↓	↓	↑	↑
↑	↑	↑	↓	↑	↓	↑	↑
↑	↑	↑	↑	↑	↓	↑	↑
↑	↑	↓	↑	↑	↑	↑	↑
↑	↓	↑	↑	↑	↑	↑	↑

Is the “connection” to  
renormalization group useful ?



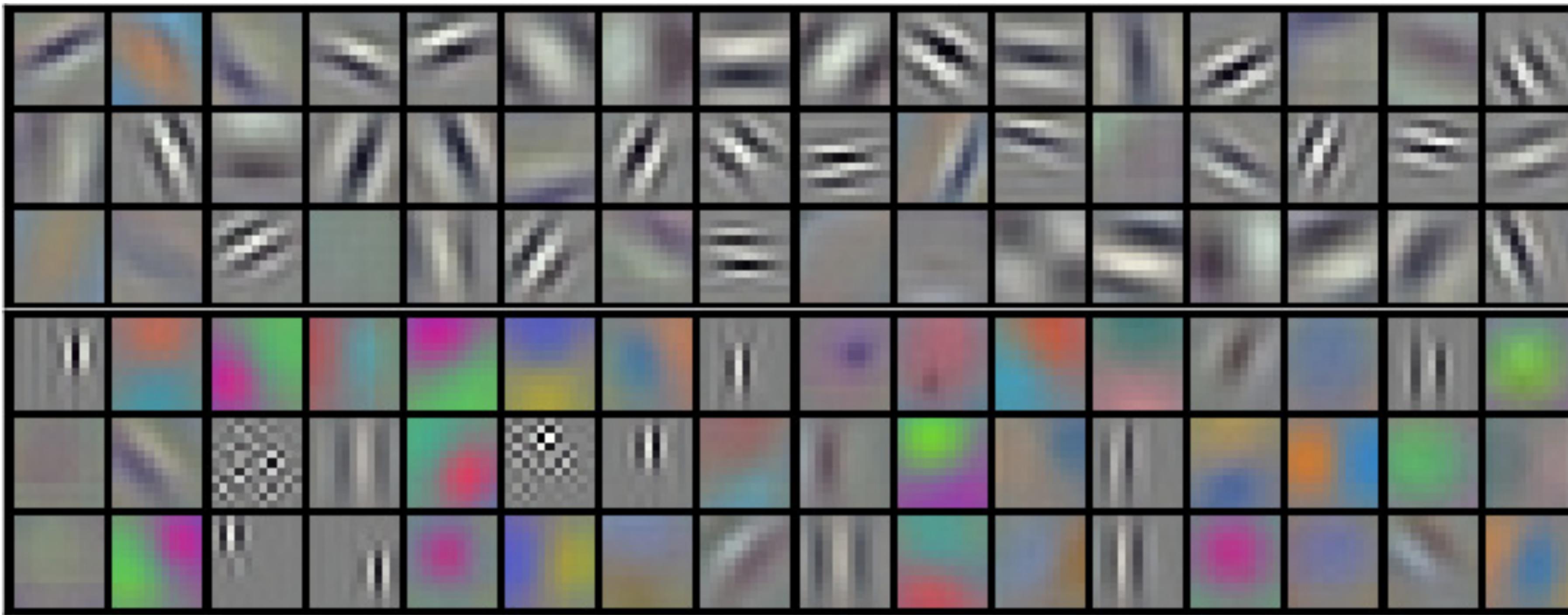
Many discussions, not conclusive yet.

Bény, 1301.3124, Mehta, Schwab, 1410.3831

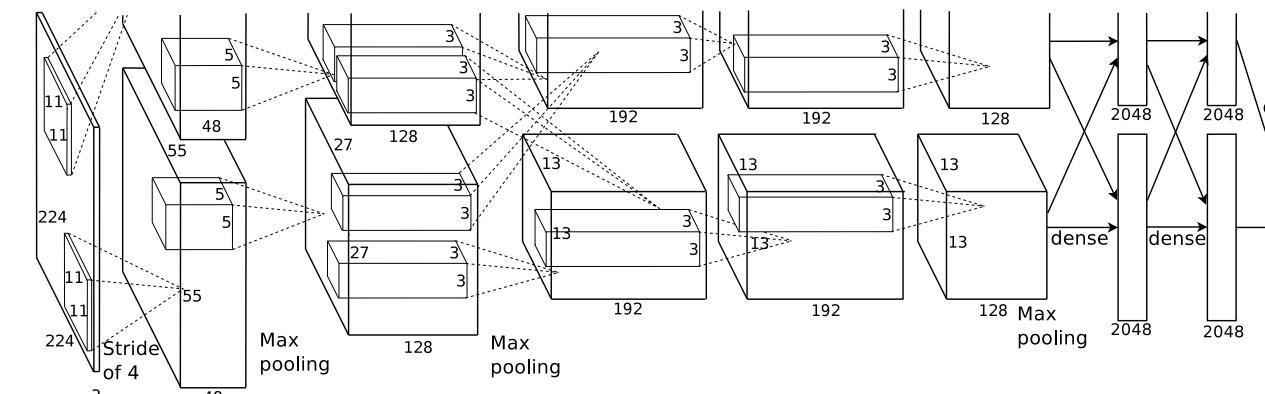
Lin, Tegmark, 1608.08225, Kenway, 1803.06111, ...

# Learned feature detectors

$$W \in \mathbb{R}^{11 \times 11 \times 96 \times 3}$$



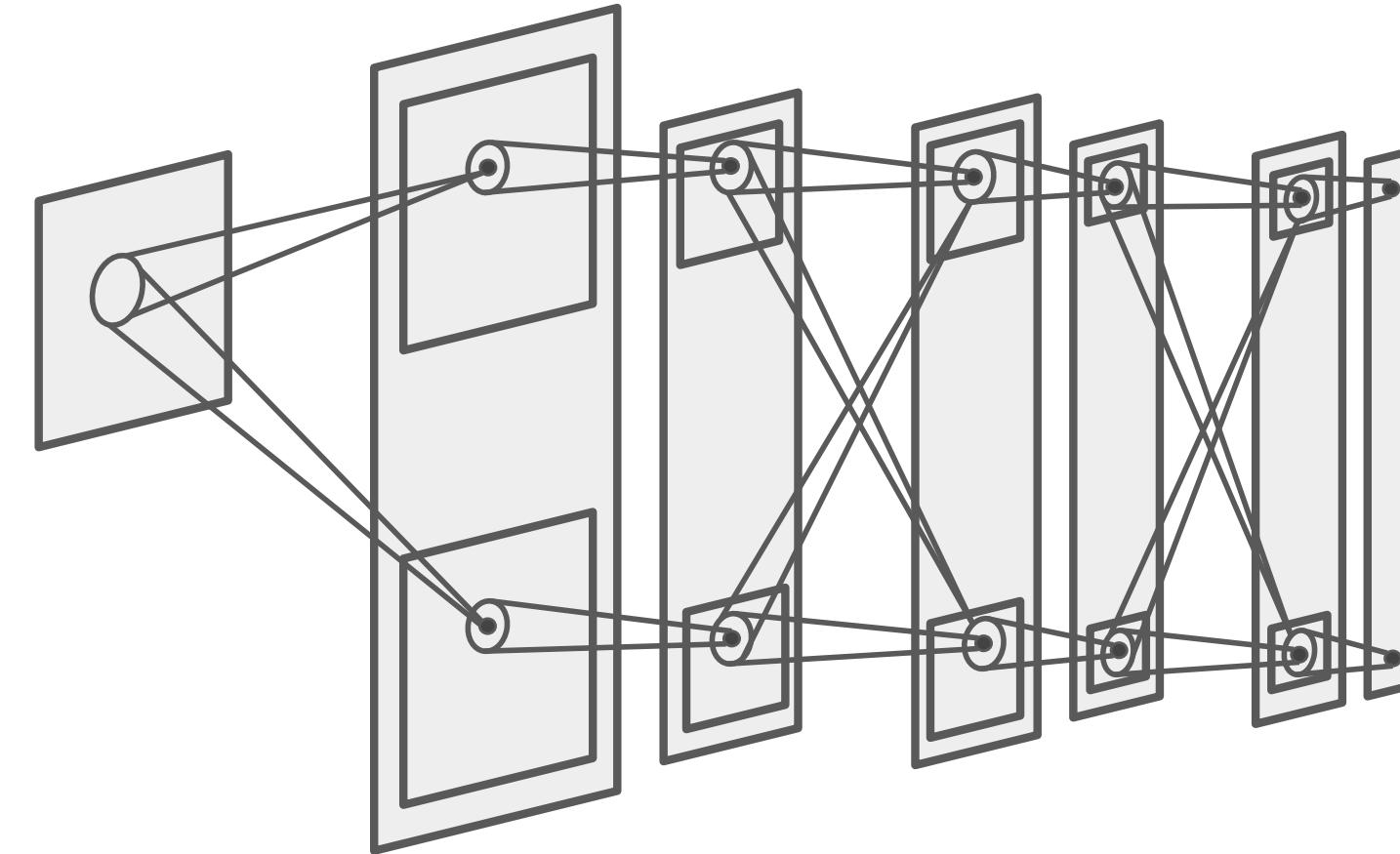
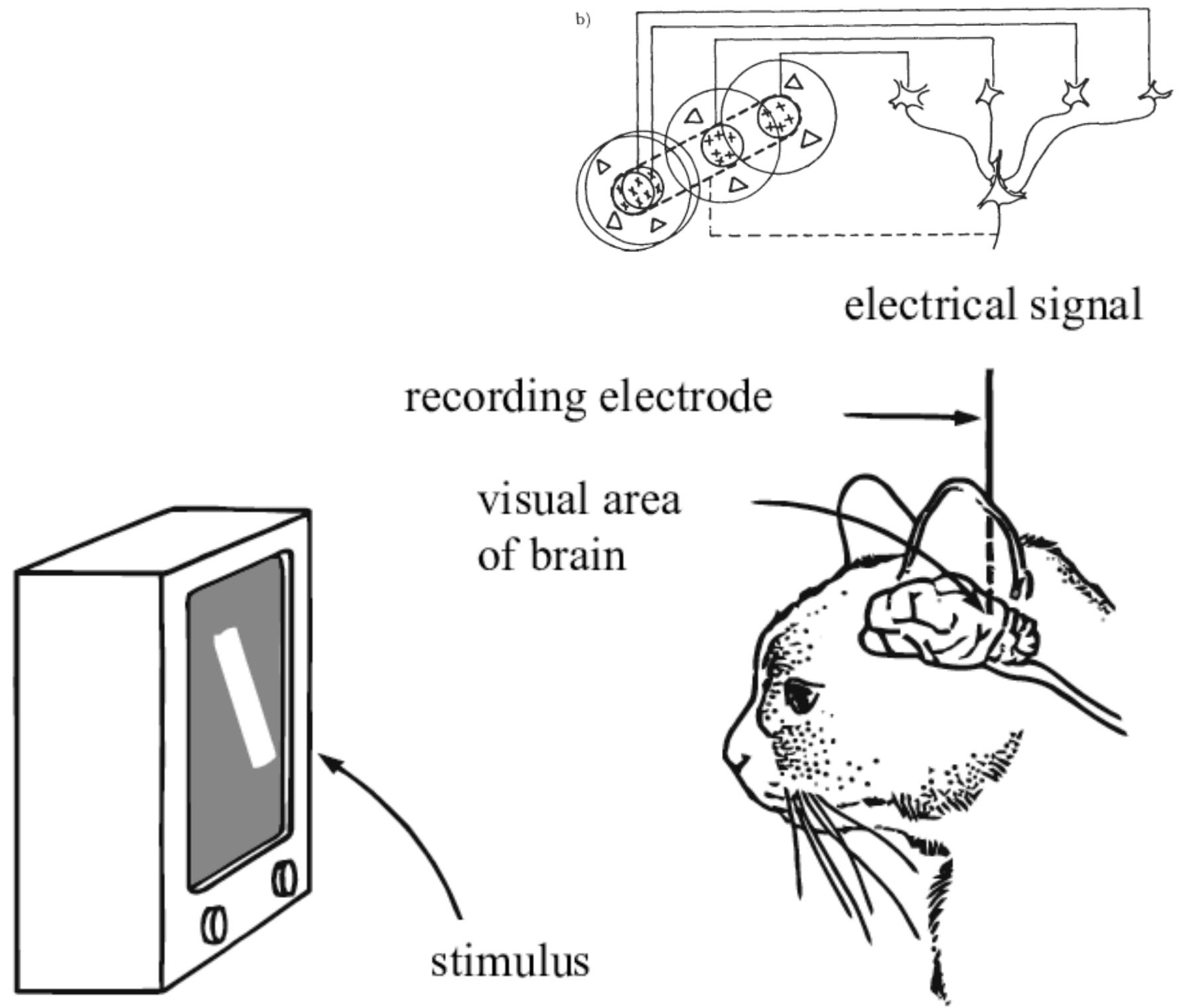
The kernels of the first convolutional layer of AlexNet



Trained on 2 GTX 580 GPU  
with 3GB memory each

Krizhevsky, Sutskever, Hinton, 2012

# CNN and visual system



Hubel & Wiesel, 1959  
Nobel prize 1981

“Neocognitron” Fukushima 1979





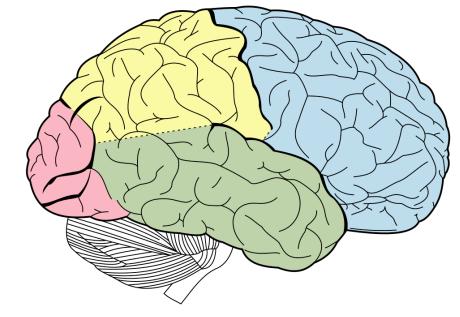
This event has sometimes been held up as an example of the importance of “accident” in science. We have never felt that it was an accident. If there is something there to discover one has to take the time to find it, and one has to be relaxed enough about the way one works so as not to foreclose the unexpected. Two other groups failed to discover orientation selectivity because they were too scientific, in a simplistic sense of that word: one group built a device to generate horizontal bright bars, the other group, vertical, in both cases so that they could explore the retina more efficiently than with a roving spot. In a certain early phase of science a degree of sloppiness can be a huge advantage.

David H. Huber

<https://www.sfn.org/-/media/SfN/Documents/TheHistoryofNeuroscience/Volume-1/c9.pdf>

# Frequently asked questions

## Biologically-inspired learning?

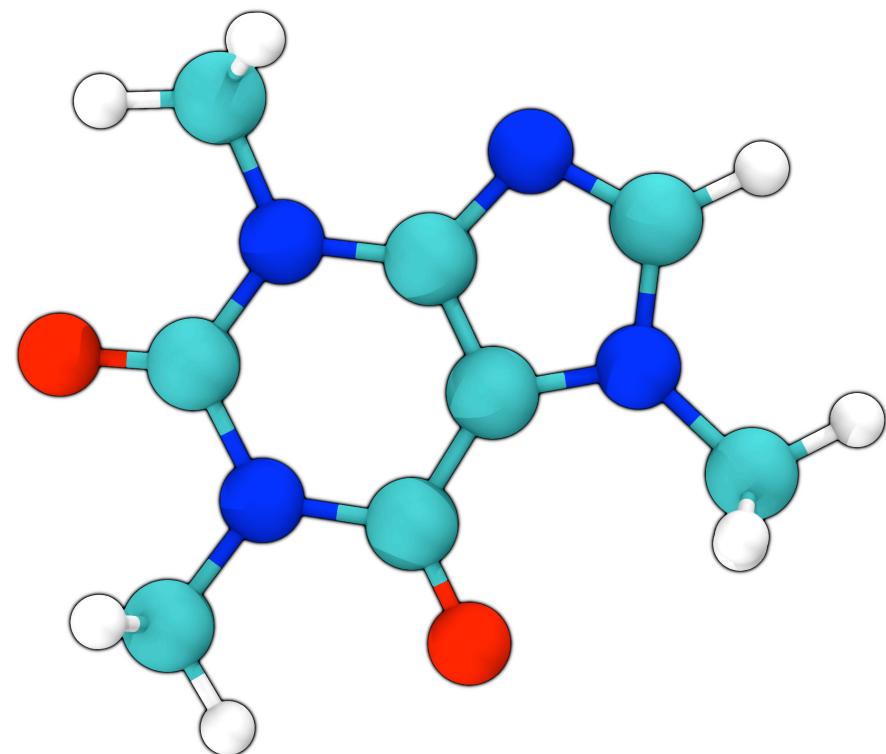
		
Power	10 W	300 W
Frequency	< kHz	GHz
Learning	Spikes	Gradient

Human intelligence is an existence proof of **artificial** intelligence

*"Can a plane fly, can a submarine swim ?"*

# Learning on graphs

Molecule graph

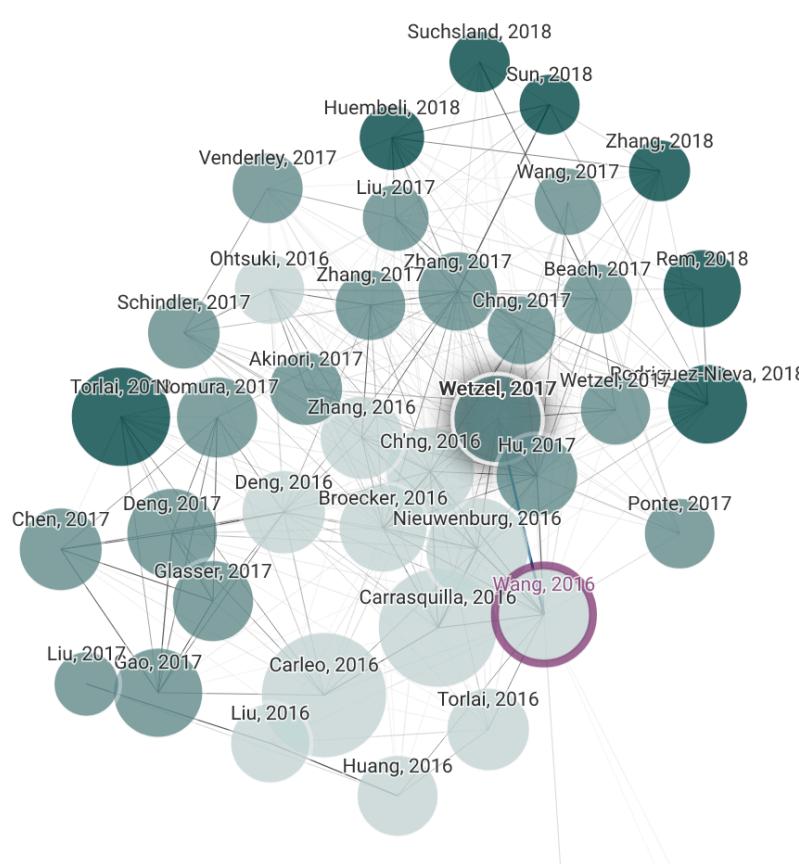


tasks

Graph level task

How does the molecule  
smell ?

Citation graph



Node level task

Force on each atom ?

Social network

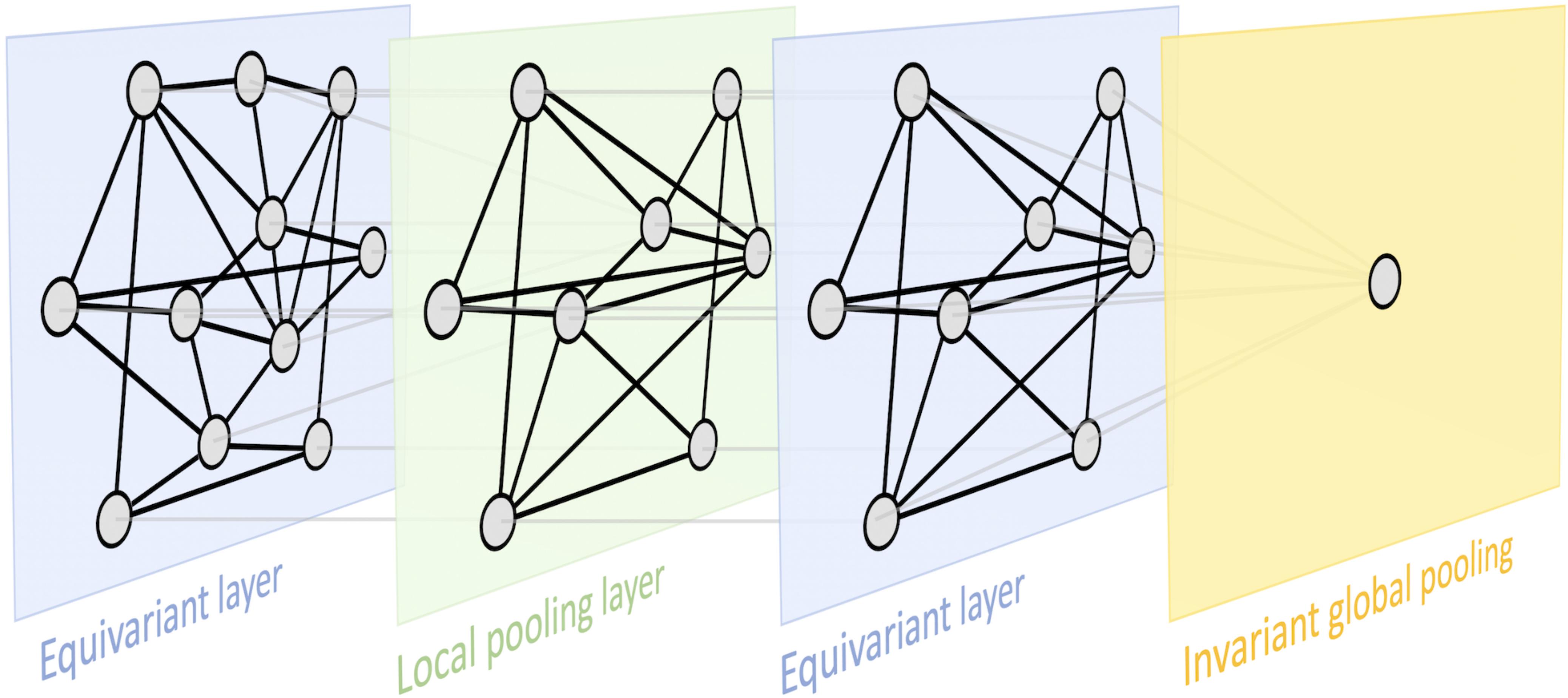


Edge level task

Hamiltonian of the  
molecule?

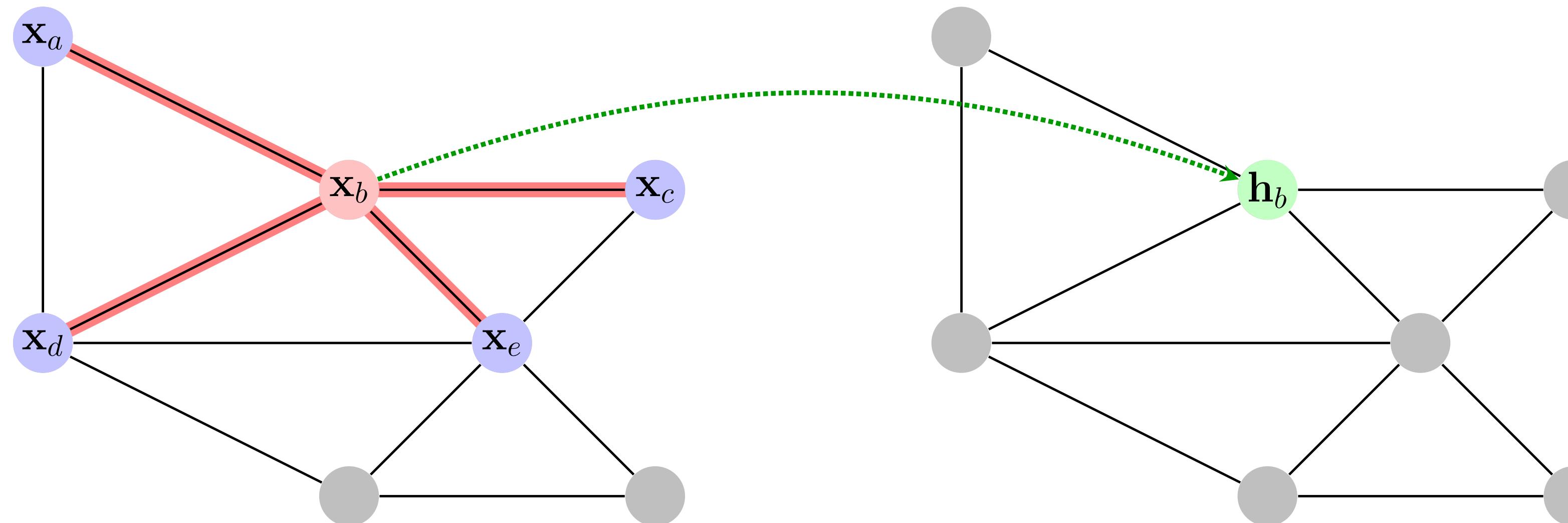


# Graph neural network



# Neural message passing

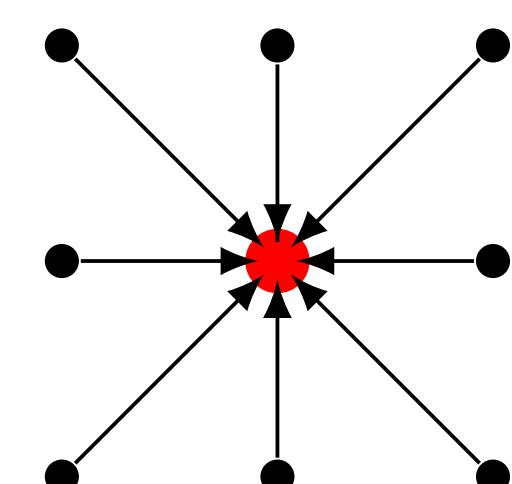
Gilmer et al, 1704.01212



$$h_b = \phi \left( x_b, \bigoplus_{v \in \mathcal{N}(b)} \psi(x_b, x_v) \right)$$

aggregation over neighbors

CNN passes message  
on a square grid



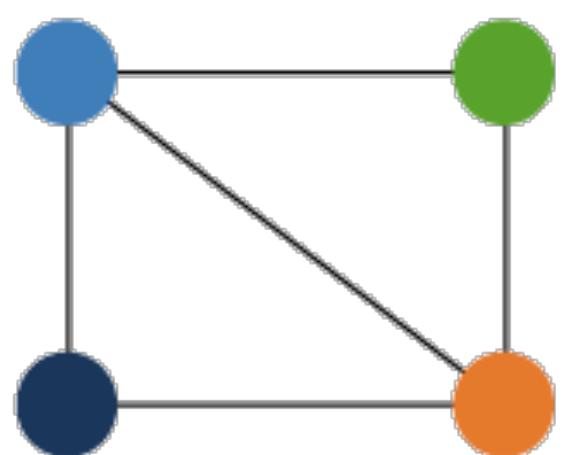
# Graph convolution

Kipf, Welling, 1609.02907

$$h = \phi(AxW)$$

Adjacency matrix + I

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



Weight matrix

$$A \in \mathbb{R}^{N \times N}$$

$$W \in \mathbb{R}^{F_x \times F_h}$$

Input feature

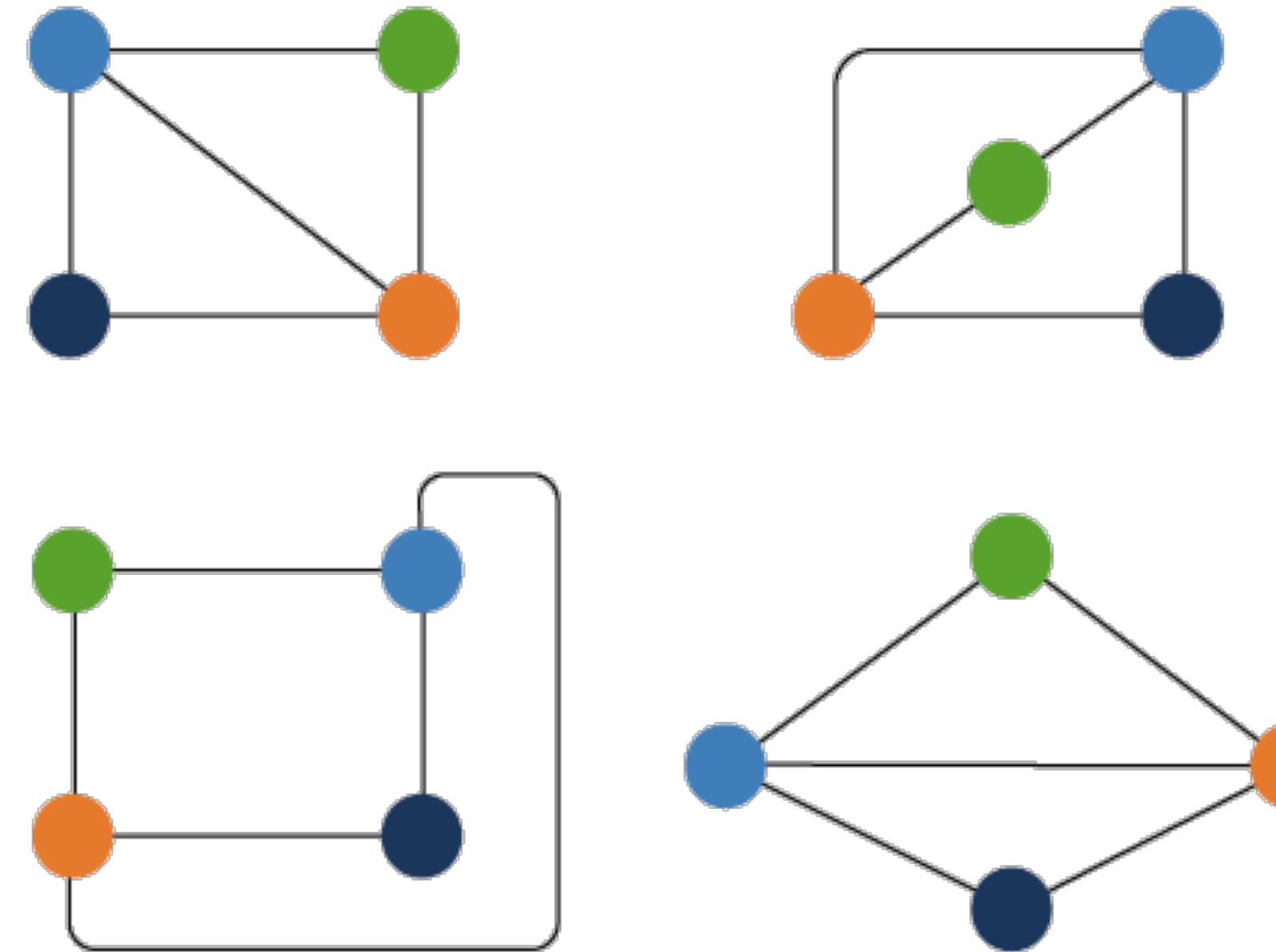
$$x \in \mathbb{R}^{N \times F_x}$$

Output feature

$$h \in \mathbb{R}^{N \times F_h}$$

$\phi$  Node wise activation

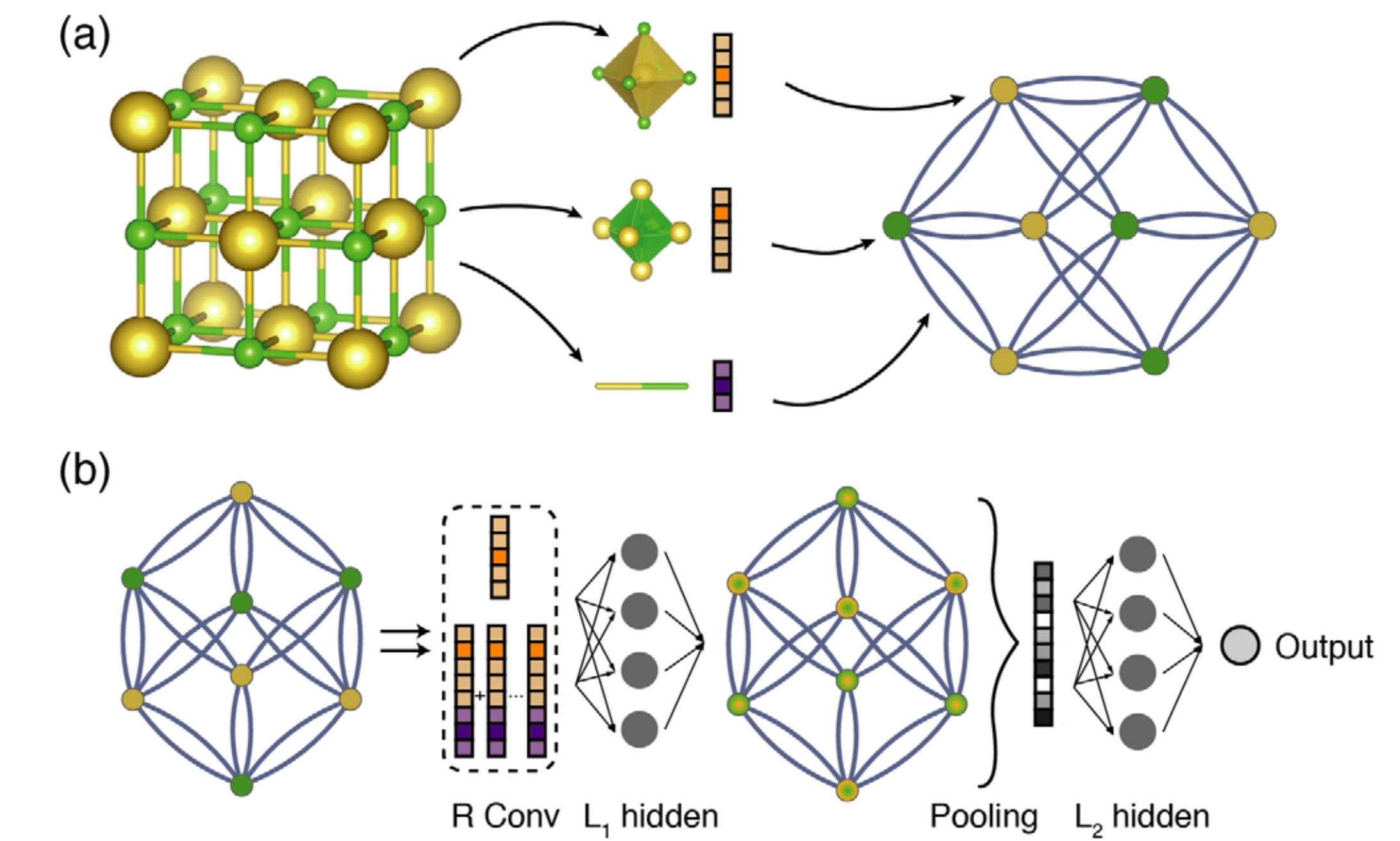
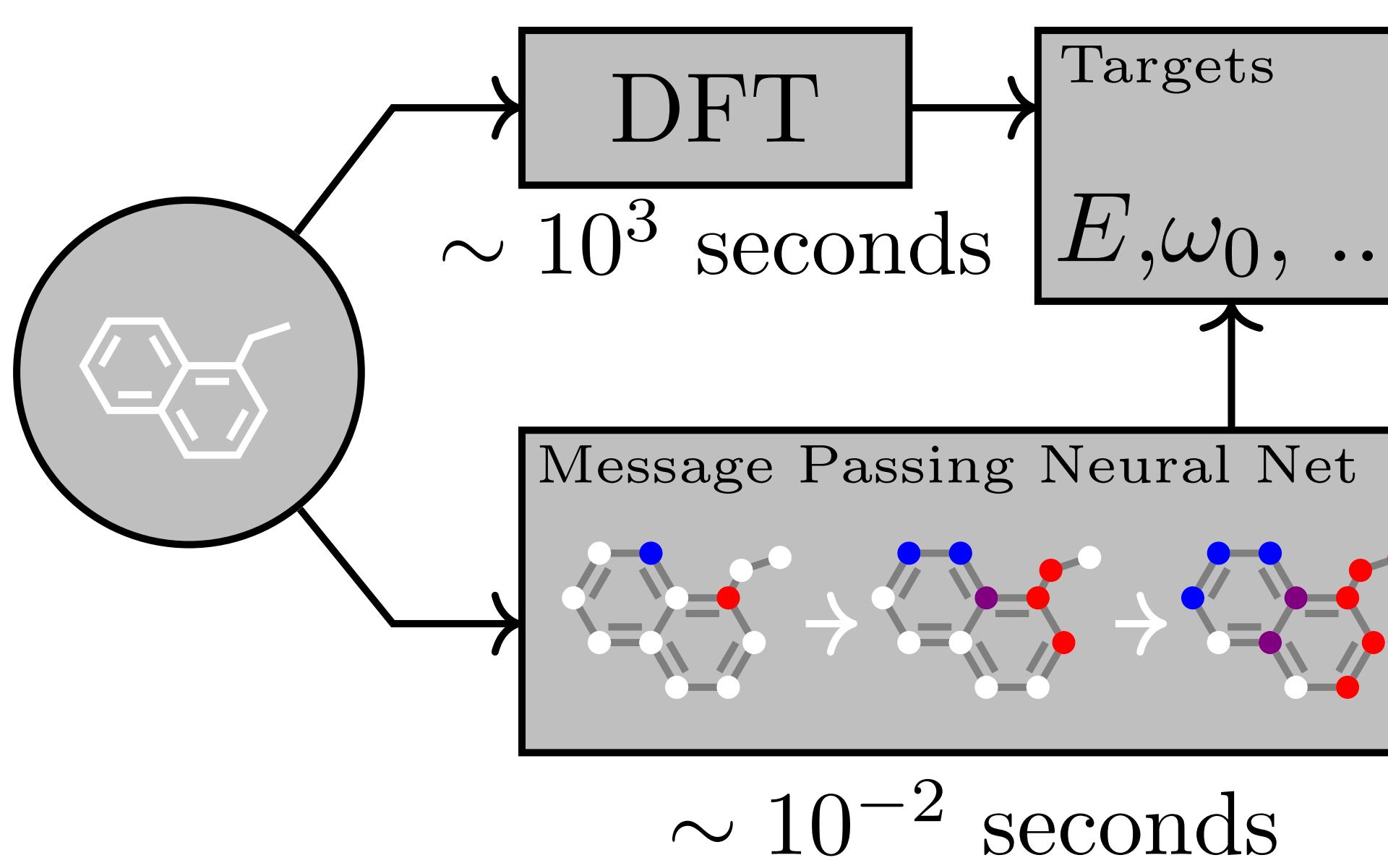
# Graph Isomorphism



Graph convolution is equivariant to node permutation

$$Ph = P\phi(AxW) = \phi(PAP^T PxW)$$

# GNN for chemistry and material



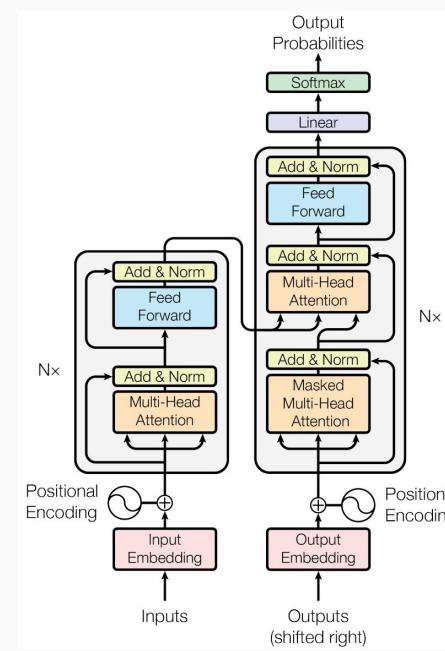
Gilmer et al, 1704.01212

Xie et al, PRL '18

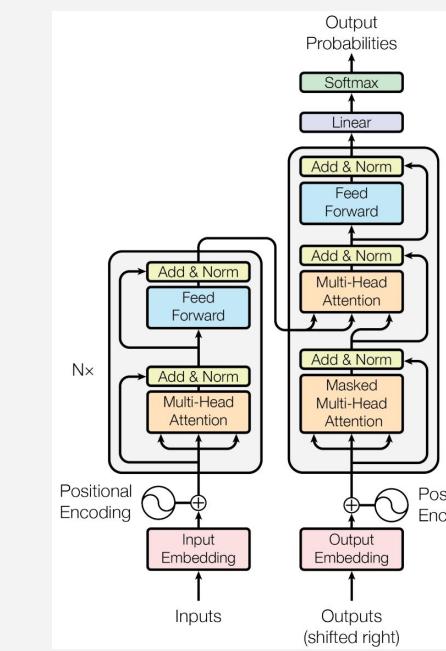
# Attention is all you need

Vaswani et al, 1706.03762

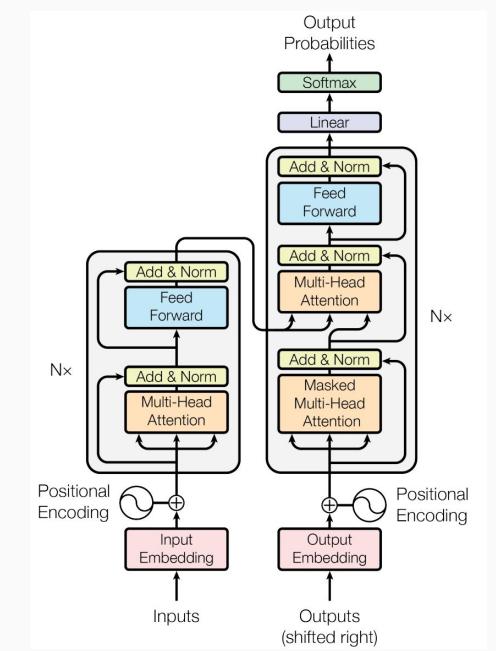
## Computer Vision



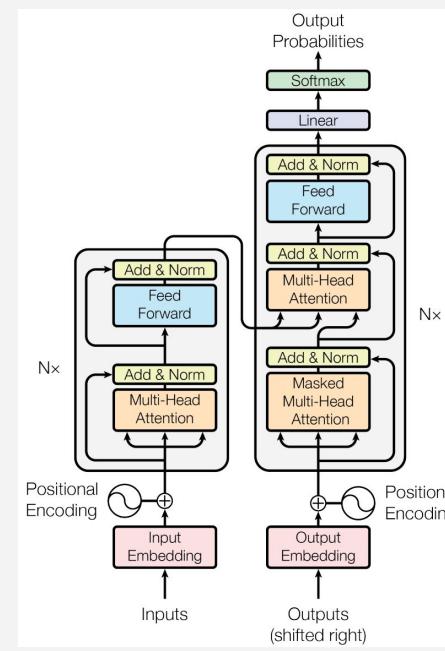
## Natural Lang. Proc.



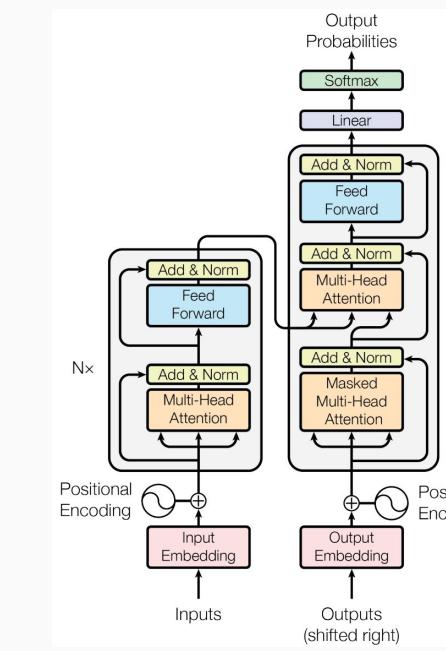
## Reinf. Learning



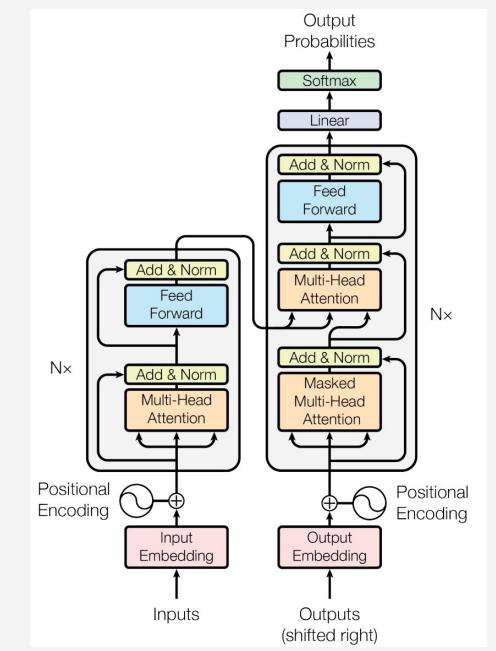
## Speech



## Translation



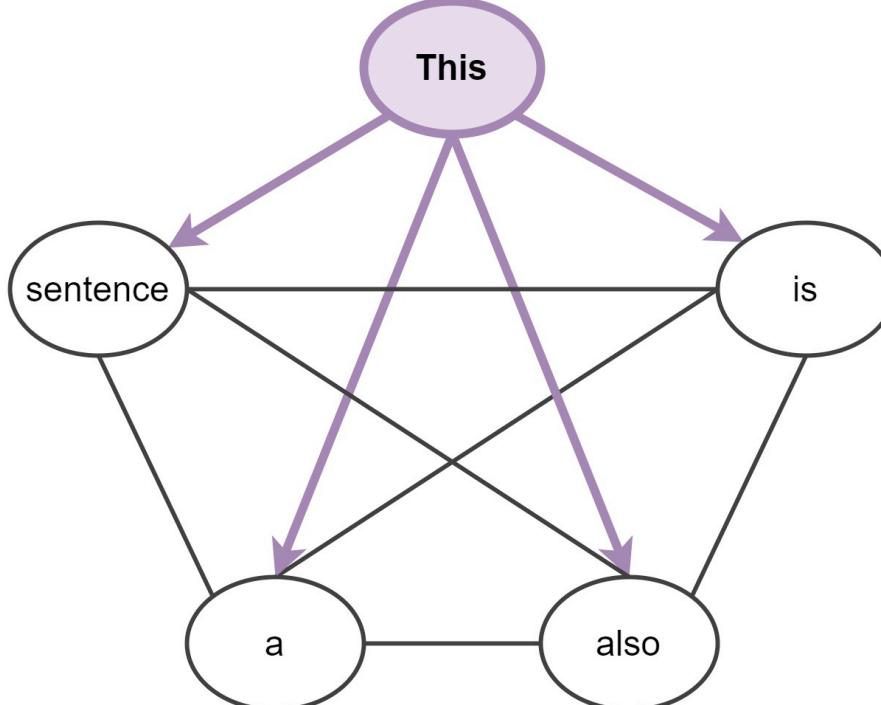
## Graphs/Science



# Attention mechanism

Up to now, we've seen models (MLP, CNN, GNN) process information depending on its **relative locations**.

Attention mechanism process information depending on its **actual content**.



Attention is “soft” dictionary

$$\sum_j \alpha(q_i, k_j) v_j \quad \text{attention weight}$$

key	value
a	1
b	2
c	3

$d = \{'a' : 1, 'b' : 2, 'c' : 3\}$

$d['b'] = 3$

↑ key  
↑ value

query ↪

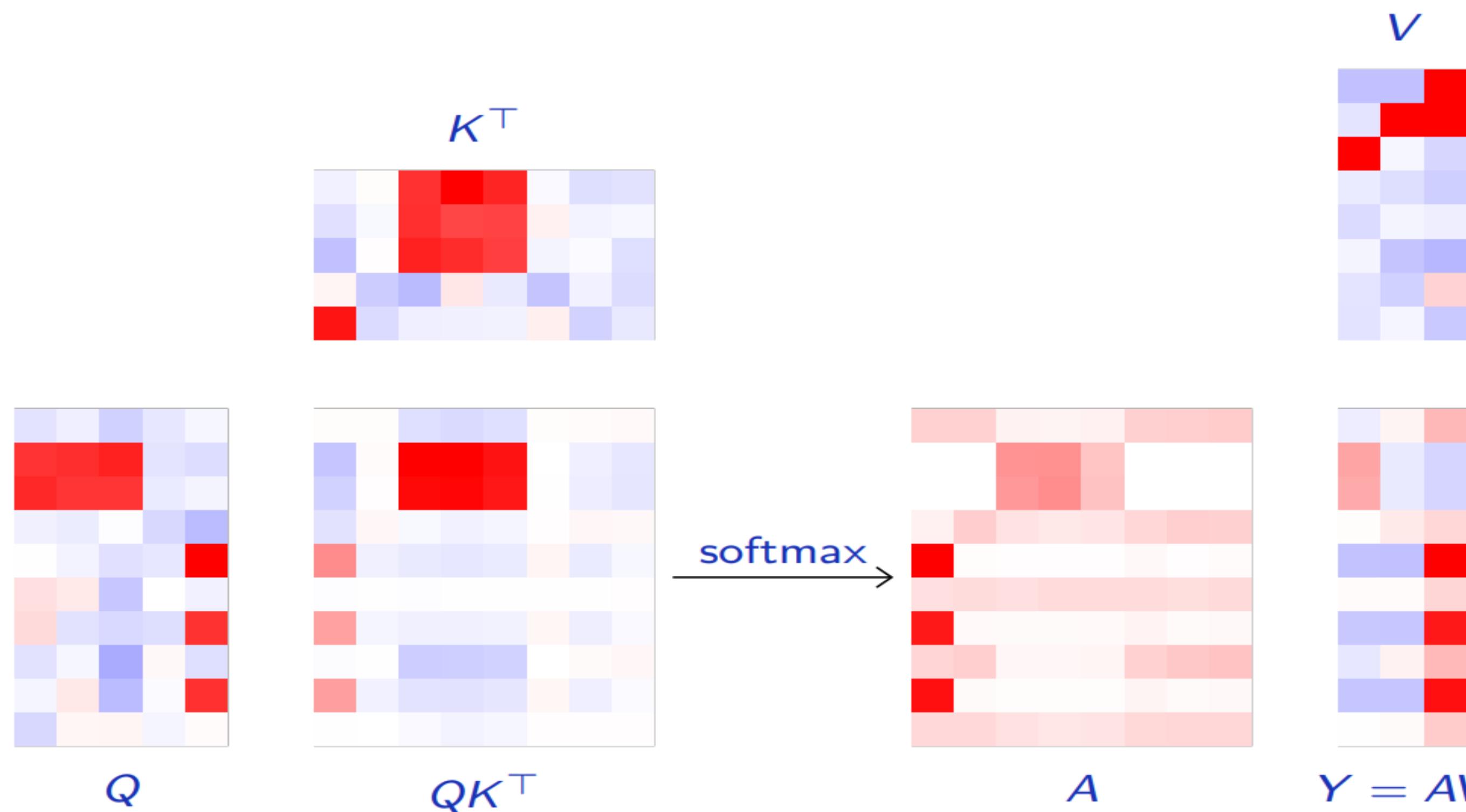
# Scaled dot product attention

$$Q \in \mathbb{R}^{N \times d}$$

$$K \in \mathbb{R}^{N \times d}$$

$$V \in \mathbb{R}^{N \times d_v}$$

$$\text{attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$

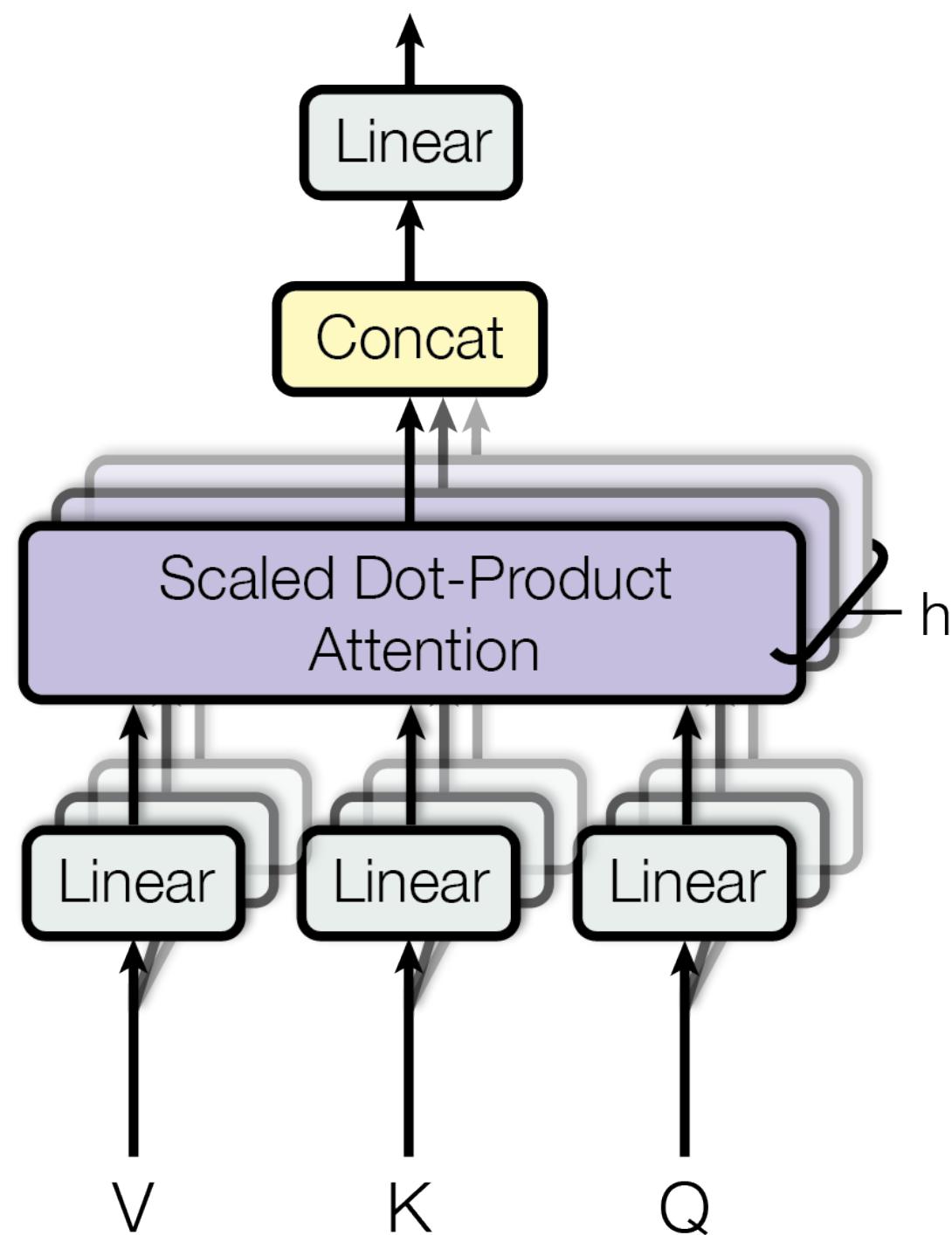


# Multihead attention

Project Q, K, V to h heads,

Run multiple attention layers in parallel

Concatenate the output together



$$\text{mha}(Q, K, V) = [\text{head}_1, \dots, \text{head}_h]W^O$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

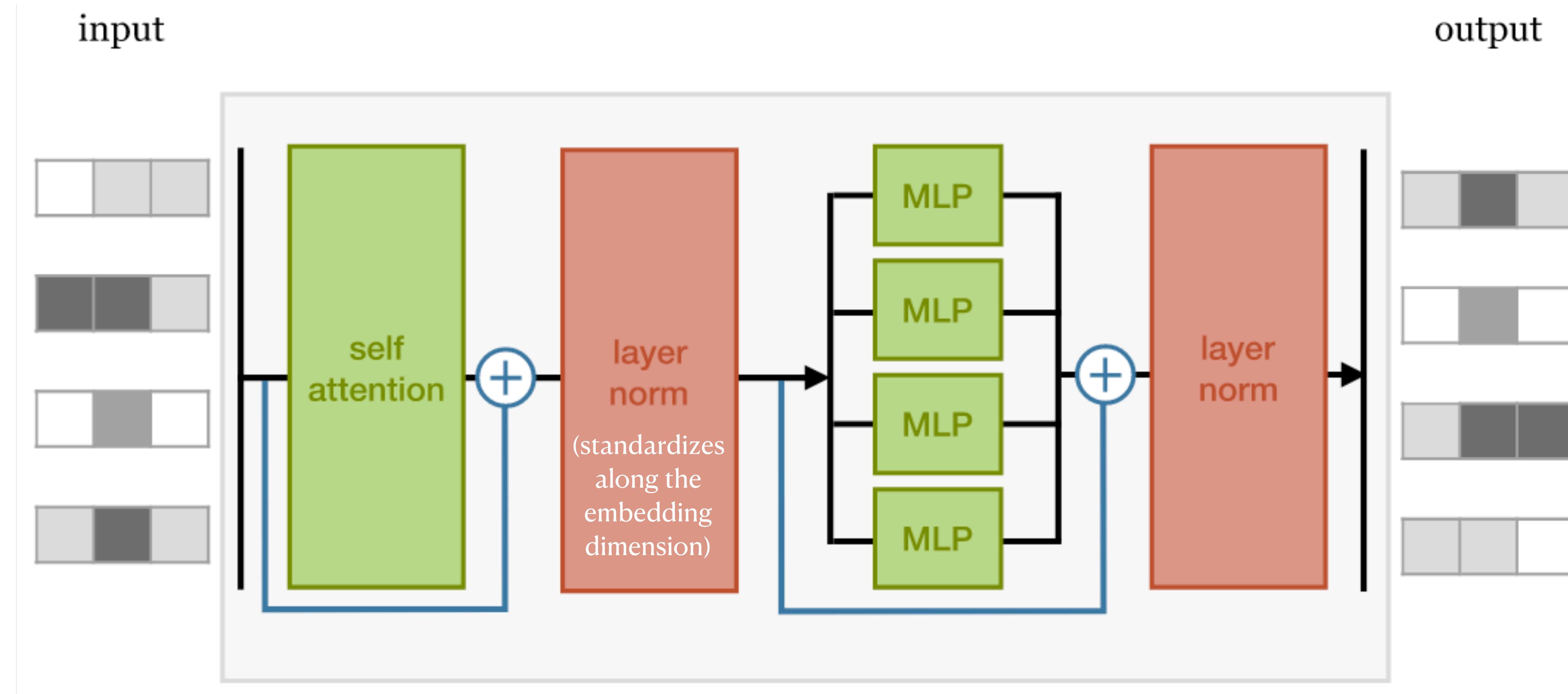
$$W_{1\dots h}^Q \in \mathbb{R}^{F \times d}$$

$$W_{1\dots h}^K \in \mathbb{R}^{F \times d}$$

$$W_{1\dots h}^V \in \mathbb{R}^{F \times d_v}$$

$$W^O \in \mathbb{R}^{h \cdot d_v \times d_{out}}$$

# The transformer block



permutation equivariant  $P \circ T = T \circ P$

add position embedding if position is needed

# Why transformer?

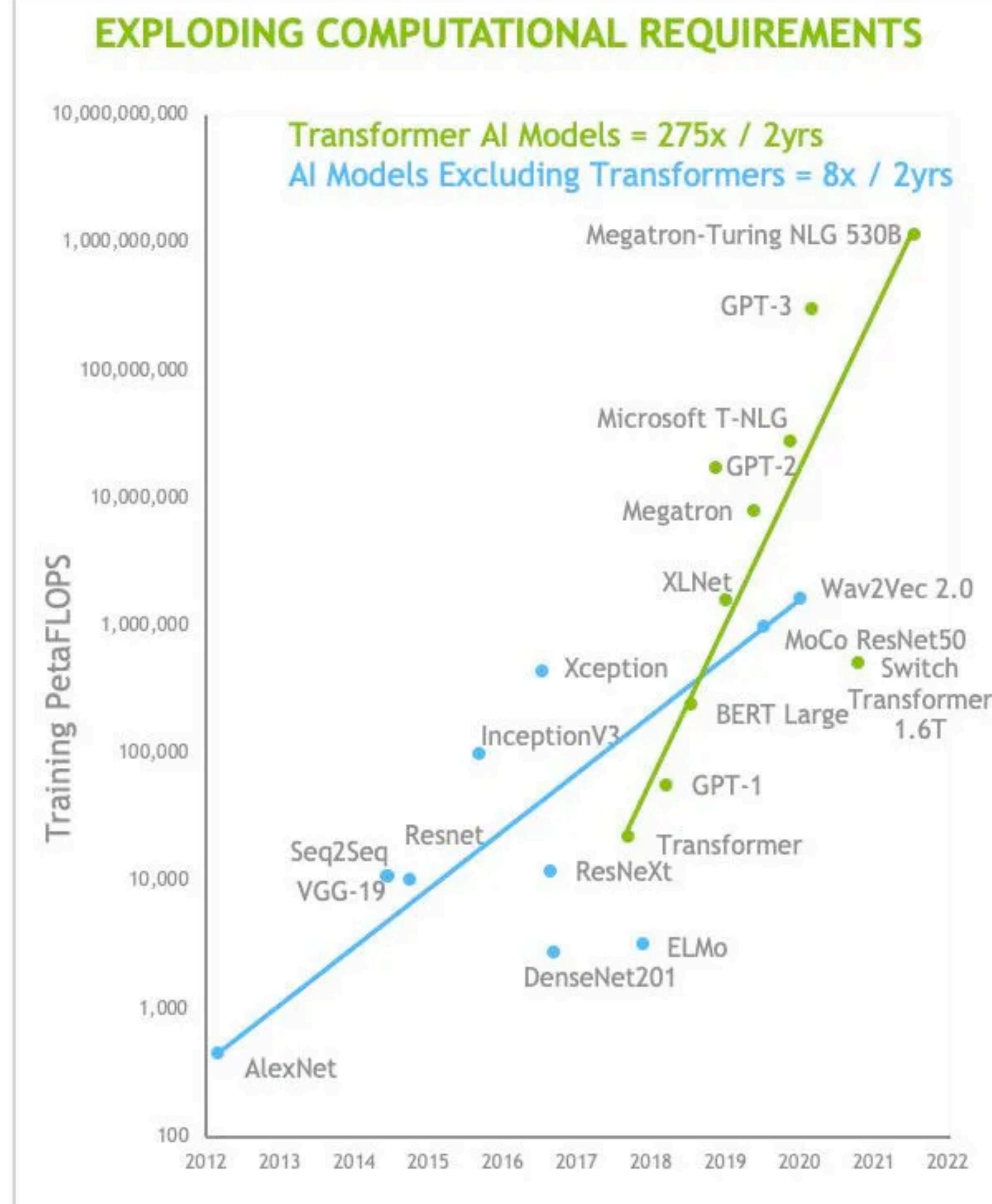
Layer type	Complexity	Sequential ops.	Max. path length
Self-attention	$O(n^2d)$	$O(1)$	$O(1)$
Recurrent	$O(nd^2)$	$O(n)$	$O(n)$
Convolutional	$O(knd^2)$	$O(1)$	$O(\log_k n)$

$n$ : sequence length

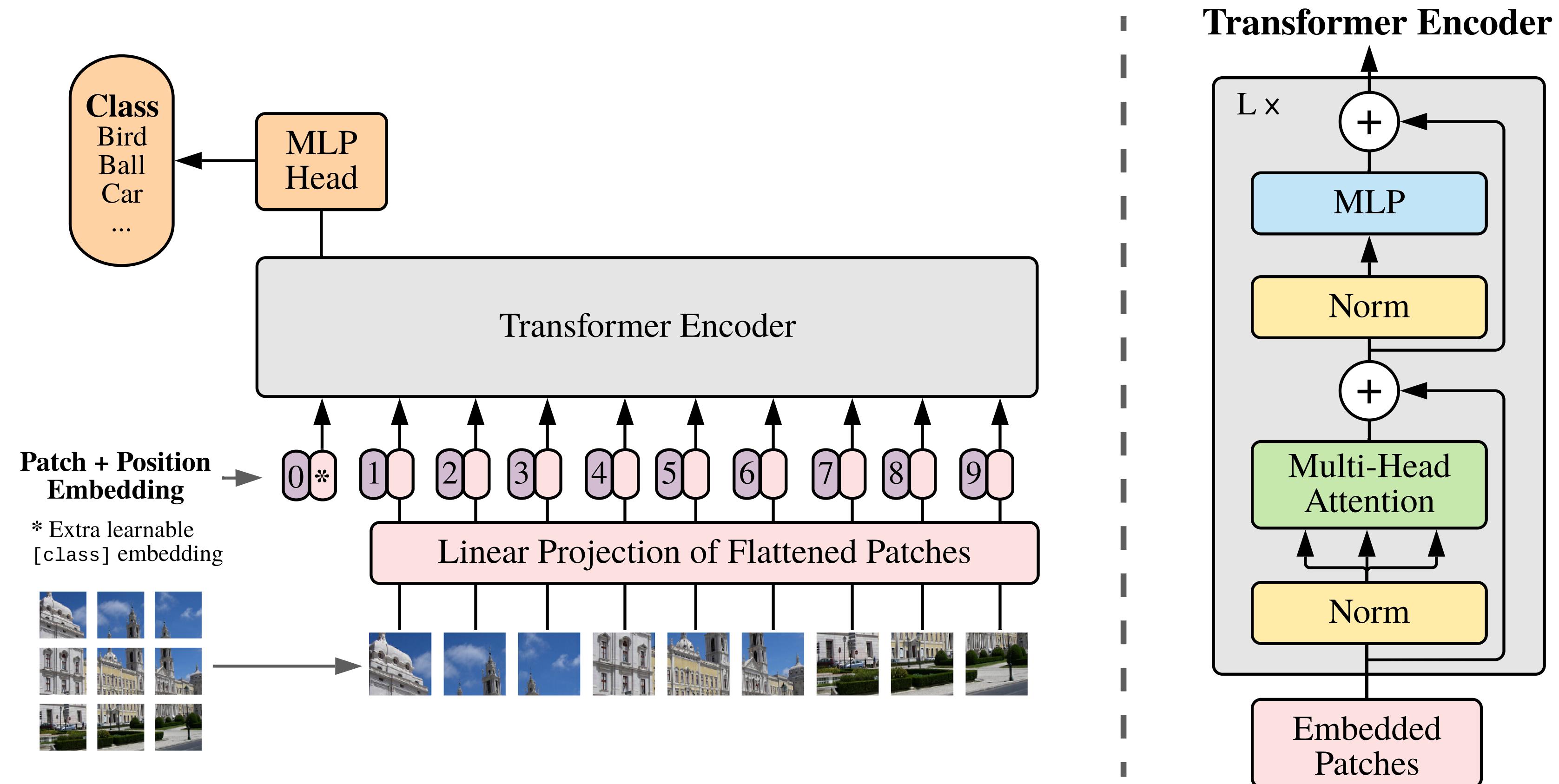
$d$ : feature dimension

$k$ : kernel size

- Directly captures *dynamical* long-range dependence (+transformer circuit math <https://transformer-circuits.pub/2021/framework/index.html>)
- Friendly to backpropagation
- Great usage of hardware accelerators

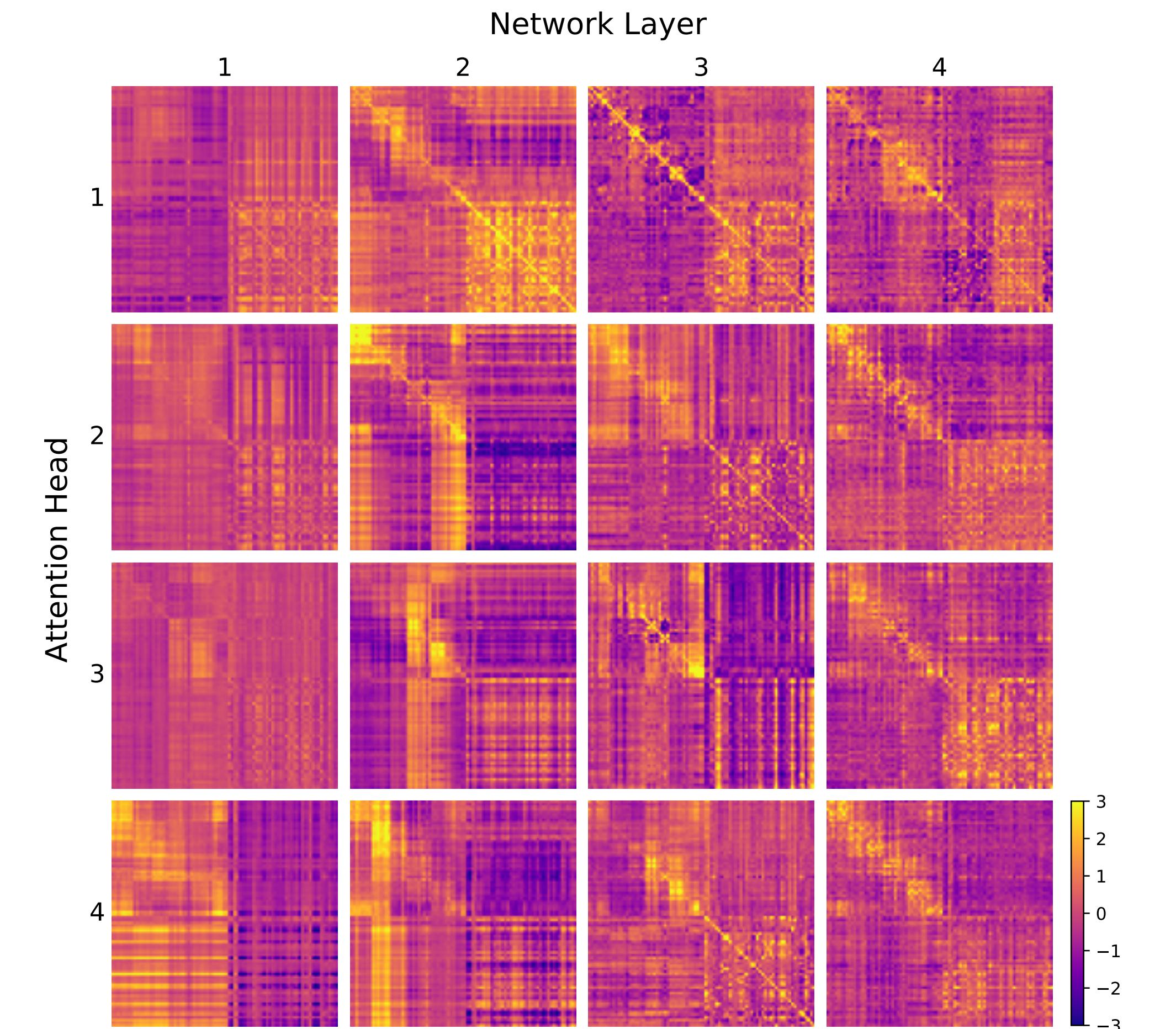
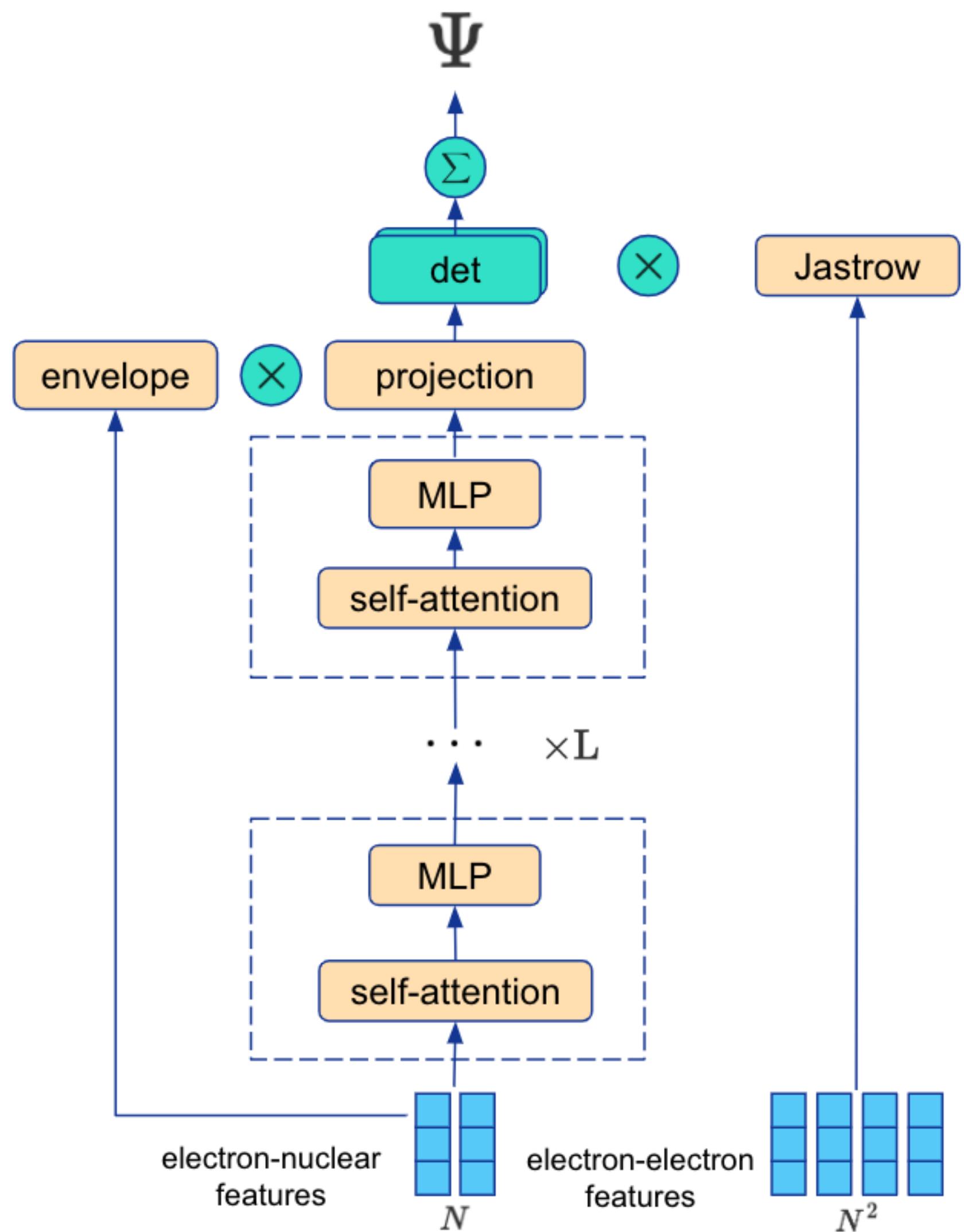


# Vision transformer



“Vision Transformer attains excellent results compared to state-of-the-art convolutional networks ”

# Psiformer: a self-attention quantum state



Attention map of electron of a benzene dimer  
“Electrons attend more to other electrons  
around the same molecule”

# Cost function

## Maximum likelihood estimation

In general, we can think of the neural network as representing a function  $f(\mathbf{x}; \boldsymbol{\theta})$ . The outputs of this function are not direct predictions of the value  $\mathbf{y}$ . Instead,  $f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\omega}$  provides the parameters for a distribution over  $y$ . Our loss function can then be interpreted as  $-\log p(\mathbf{y}; \boldsymbol{\omega}(\mathbf{x}))$ .

Goodfellow, Bengio, Courville, Deep Learning

Discriminative learning

$$\mathcal{L} = - \mathbb{E}_{(\mathbf{x}, y)} [\ln p(y | \mathbf{x})]$$

Generative learning

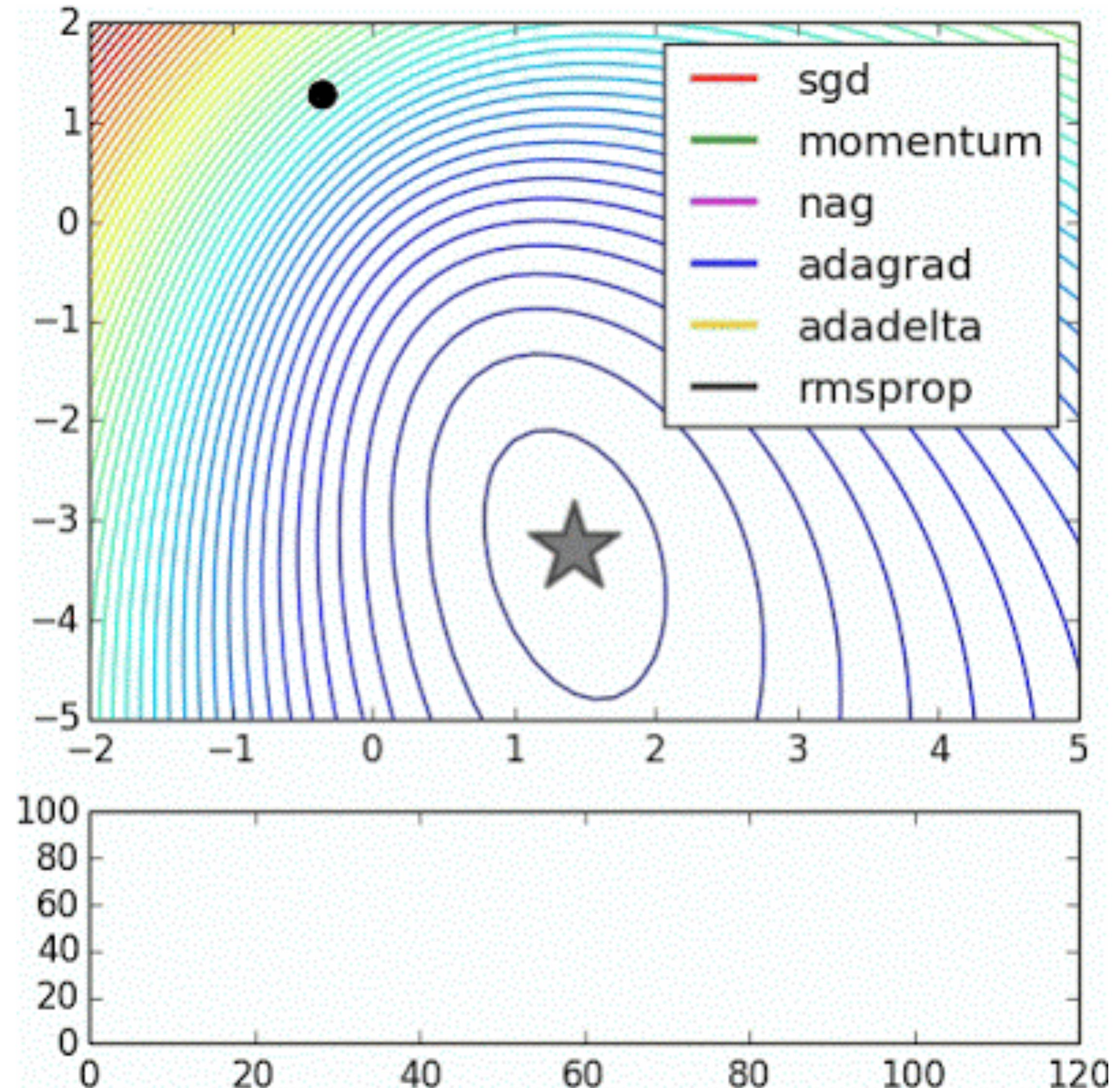
$$\mathcal{L} = - \mathbb{E}_{\mathbf{x}} [\ln p(\mathbf{x})]$$

# Optimization

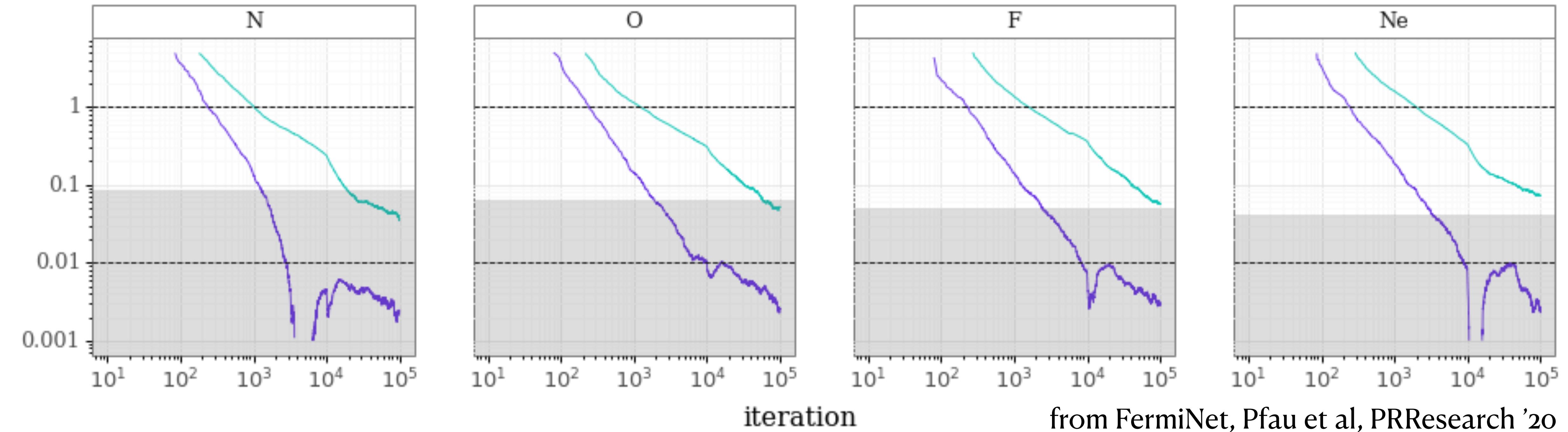
Stochastic gradient descent

$$\theta \leftarrow \theta - \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

$$\mathcal{L} = \underset{x}{\mathbb{E}} [\ell(x)]$$



# Higher-order optimizations



— ADAM Kingma & Ba, 2014  
— KFAC Martens & Grosse 2015

Natural gradient, Amari 1998  
Stochastic reconfiguration, Sorella 2005  
Hessian free optimization, Martens 2015

$$\theta \leftarrow \theta - \eta \mathbf{G}^{-1} \frac{\partial \mathcal{L}}{\partial \theta}$$

[https://www.cs.toronto.edu/~rgrosse/courses/csc2541\\_2022/](https://www.cs.toronto.edu/~rgrosse/courses/csc2541_2022/)

