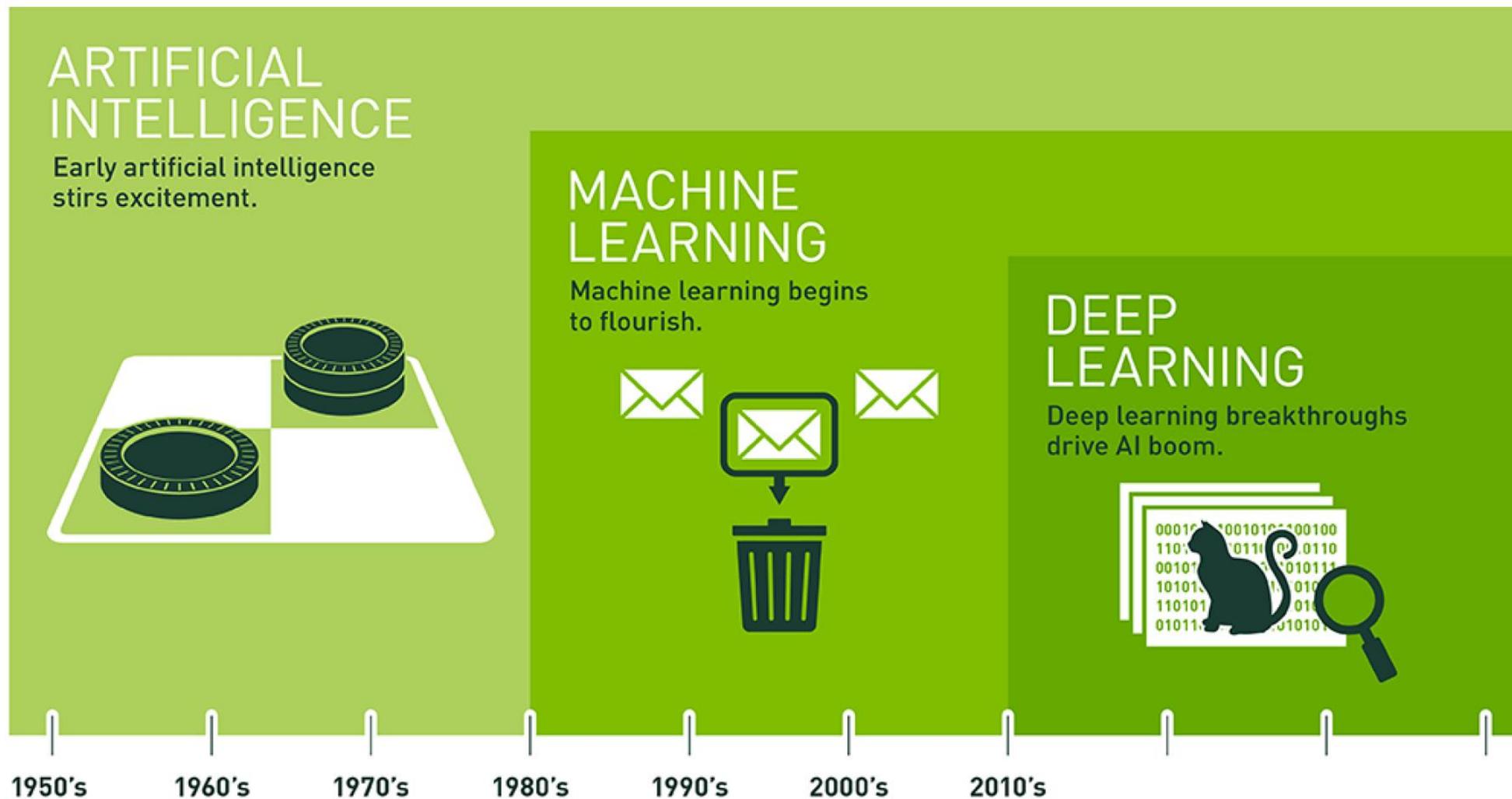


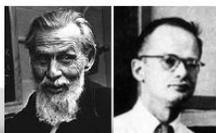
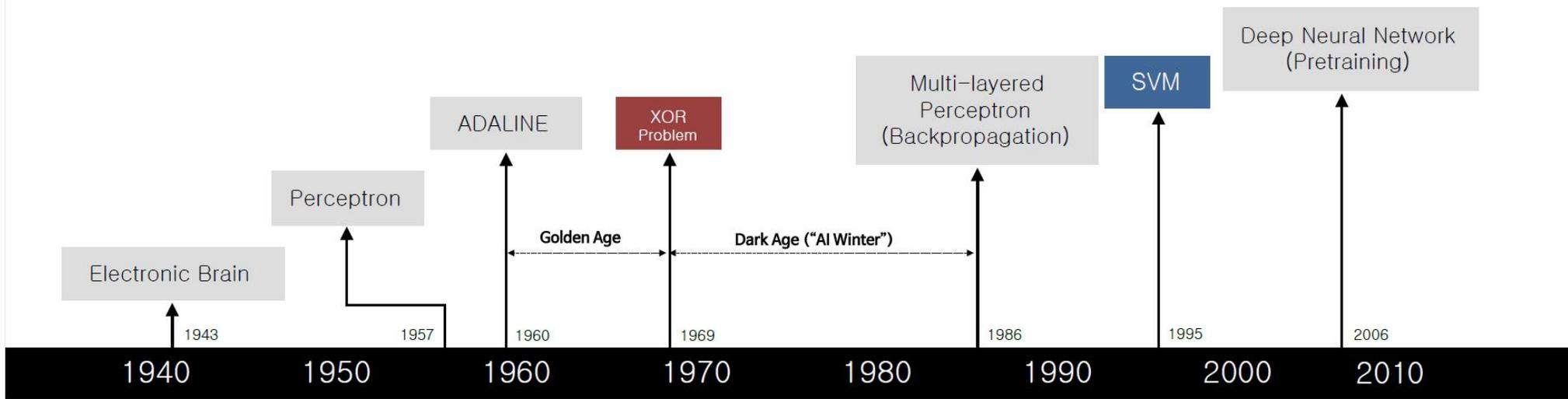
第一章 生物大模型简介

人工智能的发展历程

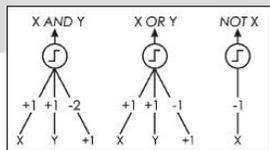


Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

人工智能发展的里程碑事件



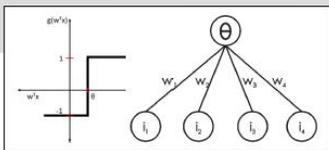
S. McCulloch – W. Pitts



- Adjustable Weights
- Weights are not Learned



F. Rosenblatt



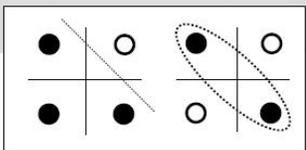
- Learnable Weights and Threshold



B. Widrow – M. Hoff



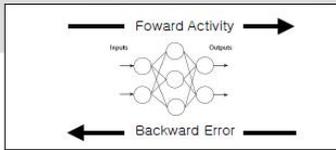
M. Minsky – S. Papert



- XOR Problem



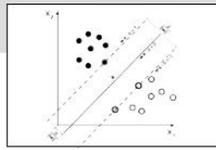
D. Rumelhart – G. Hinton – R. Williams



- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting



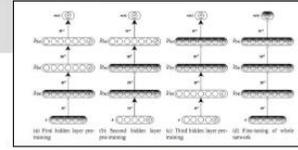
V. Vapnik – C. Cortes



- Limitations of learning prior knowledge
- Kernel function: Human Intervention

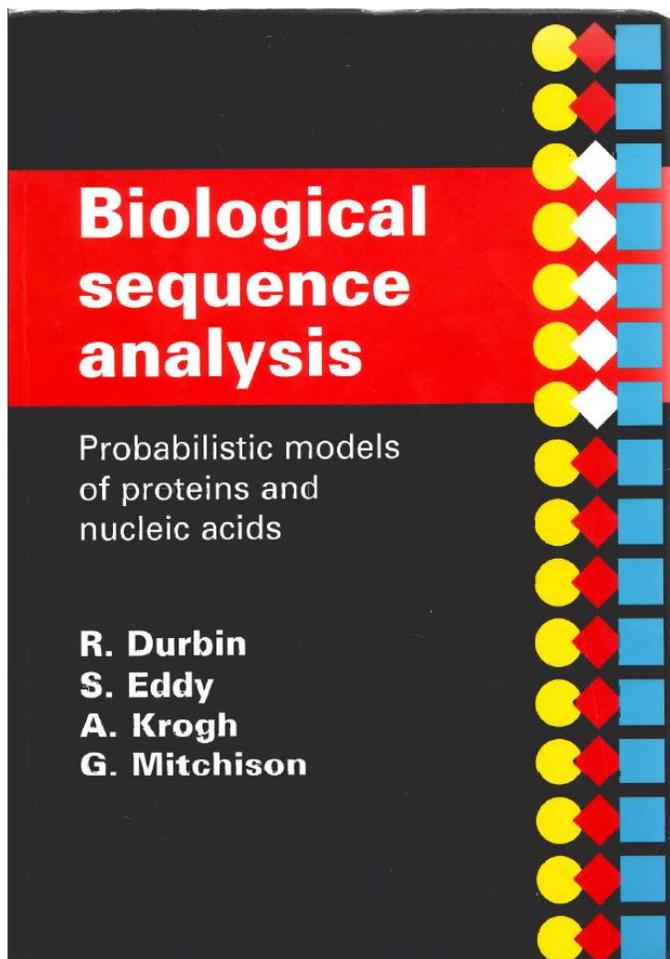


G. Hinton – S. Ruslan



- Hierarchical feature Learning

经典生物信息学/生物信息学 1.0



At a **Snowbird conference on neural nets in 1992**, David Haussler and his colleagues at UC Santa Cruz (including one of us, AK) described preliminary results on modelling protein sequence multiple alignments with probabilistic models called 'hidden Markov models' (HMMs). Copies of their technical report were widely circulated. Some of them found their way to the MRC Laboratory of Molecular Biology in Cambridge, where RD and GJM were just switching research interests from neural modelling to computational genome sequence analysis, and where SRE had arrived as a new postdoctoral student with a background in experimental molecular genetics and an interest in computational analysis. AK later also came to Cambridge for a year.

All of us quickly adopted the ideas of **probabilistic modelling**. We were per-



**Richard
Durbin**



**Sean
Eddy**



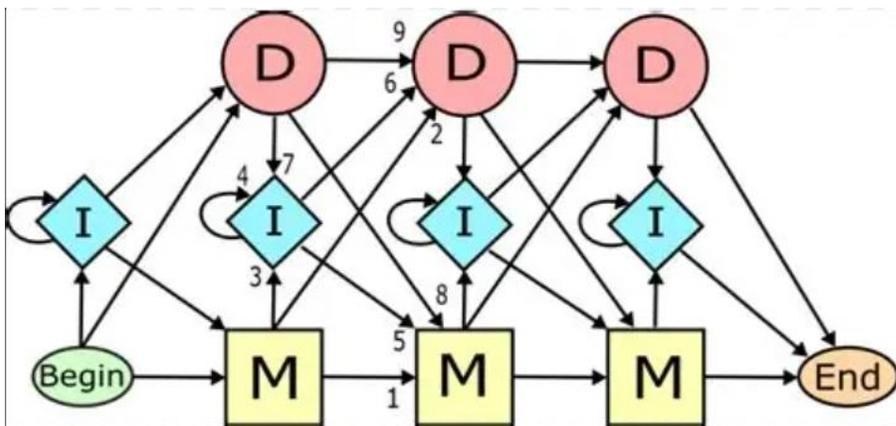
**Anders
Krogh**



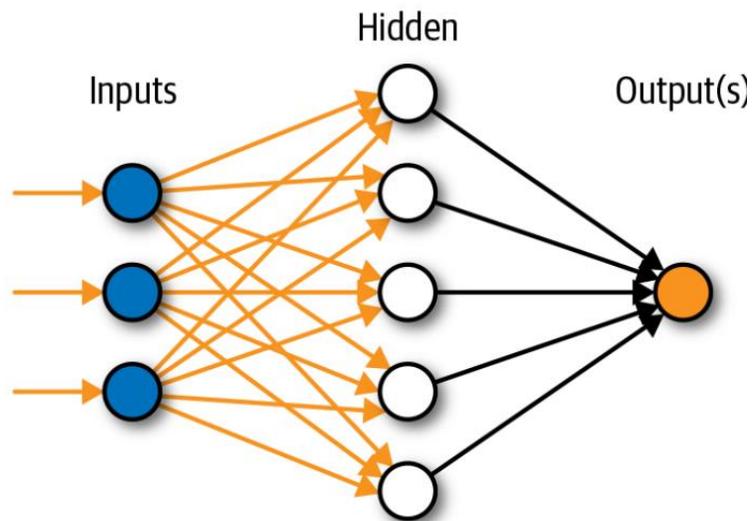
**Graeme
Mitchison**

经典机器学习

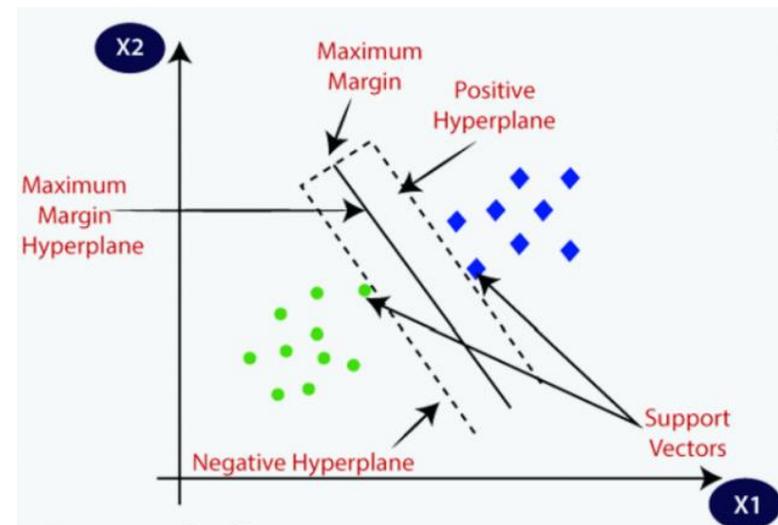
隐马尔科夫模型



人工神经网络



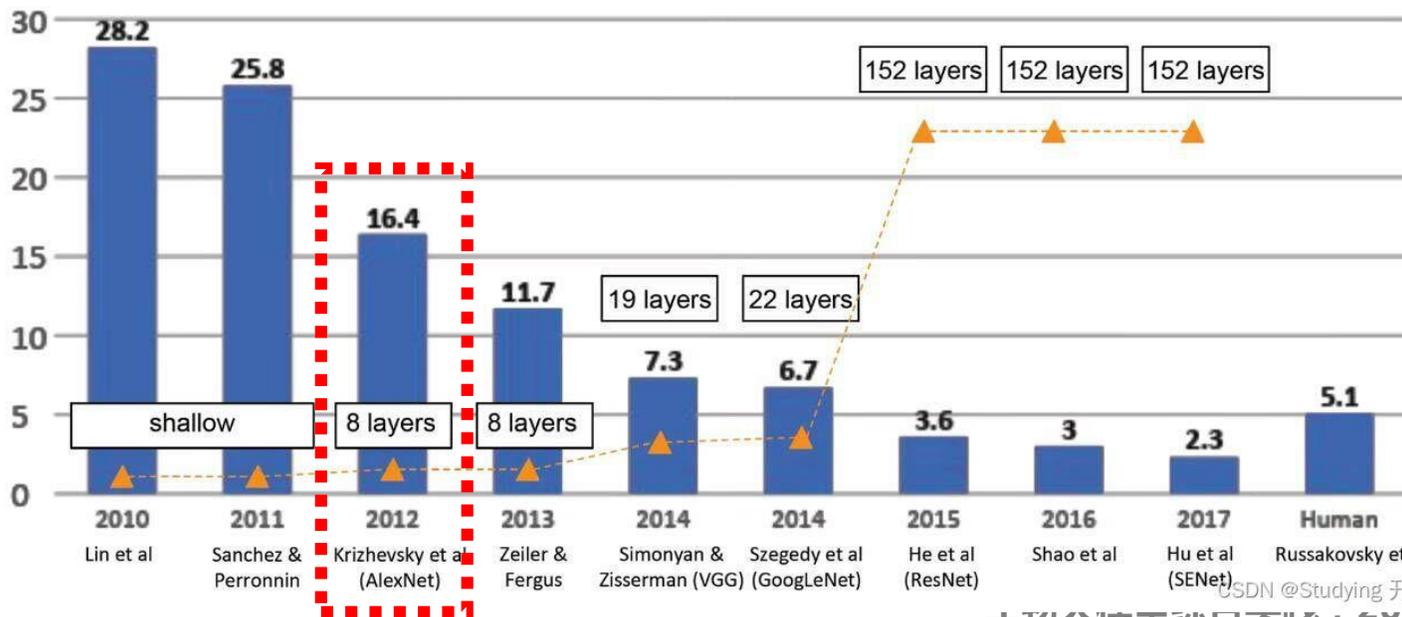
支持向量机



技术突破：卷积神经网络

- **Geoffrey Hinton**：2001-2014，多伦多大学计算机科学系教授
- 2012年12月，AlexNet，两个天才学生Alex Krizhevsky、Ilya Sutskever —— **“以学生为中心”**
- 将图像识别的错误率从 $>25\%$ 首次降低到 $<20\%$

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



技术突破：围棋 & 蛋白质结构预测

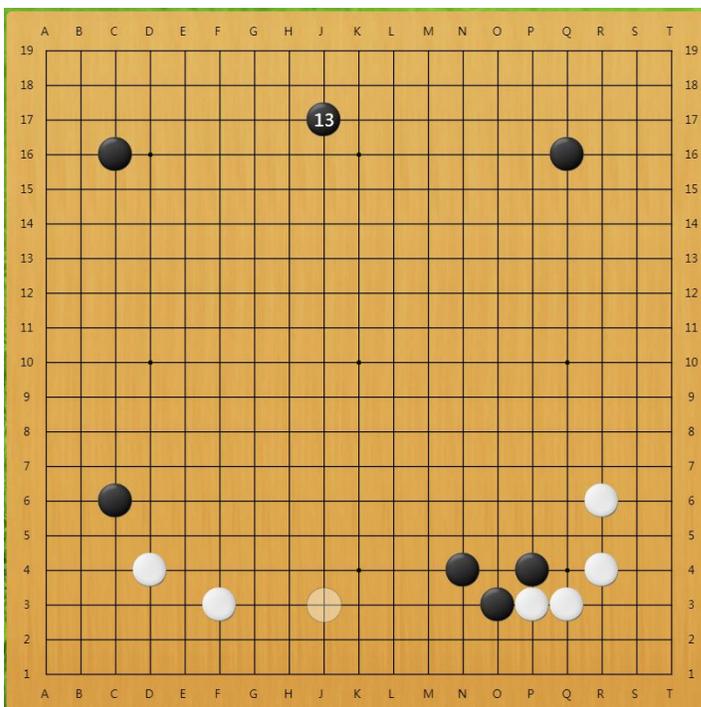
围棋

AlphaGo：战胜人类围棋世界冠军

蛋白质

AlphaFold 2：蛋白质三级结构预测的准确性媲美实验手段

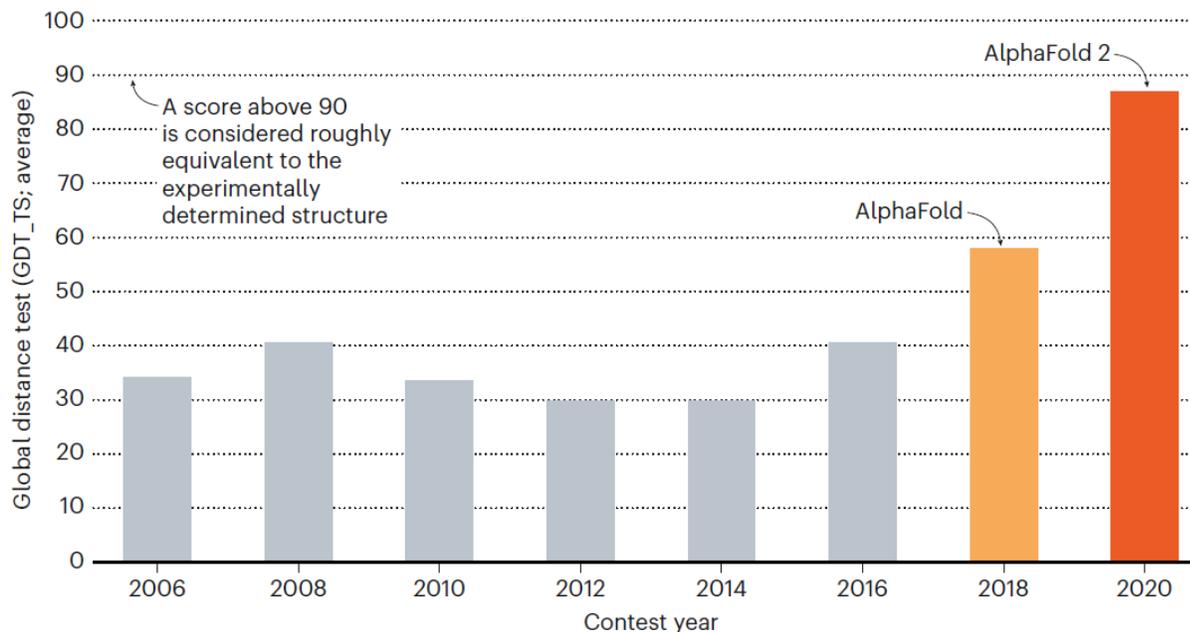
2016.3.9，李世石 vs. AlphaGo



Nature, 2016, 529, 484-9

STRUCTURE SOLVER

DeepMind's AlphaFold 2 algorithm significantly outperformed other teams at the CASP14 protein-folding contest — and its previous version's performance at the last CASP.

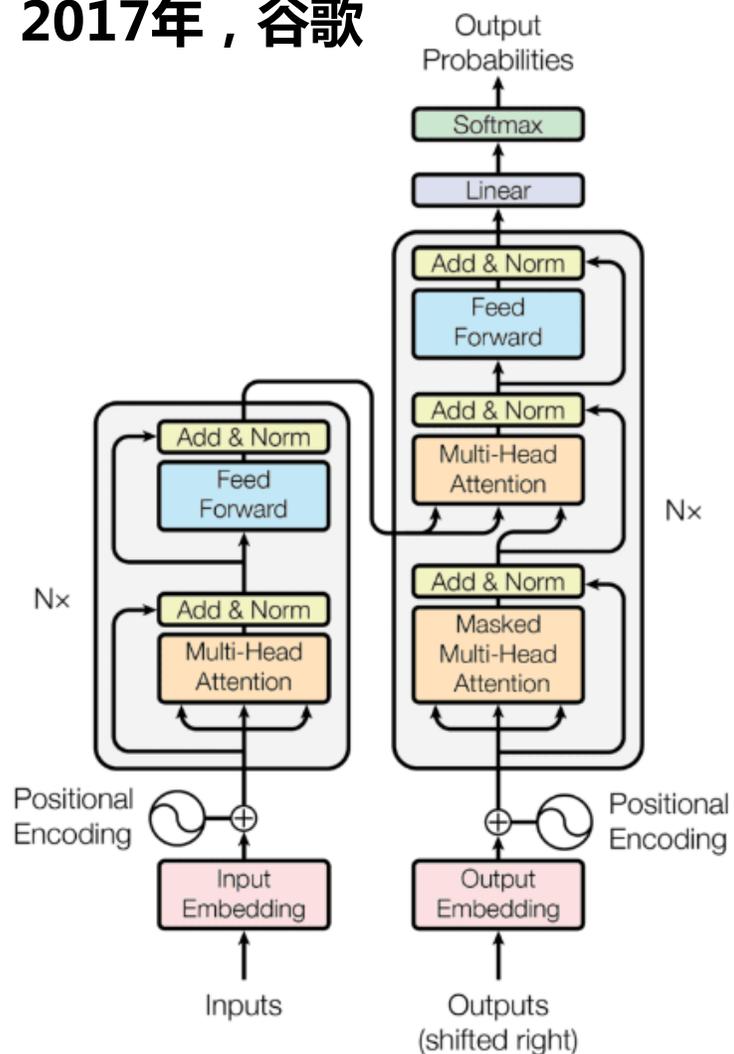


生物大模型综合实践，2025

Nature, 2020, 577, 706-710

Transformer框架

2017年，谷歌



The encoder-decoder structure of the Transformer architecture
Taken from "Attention Is All You Need"

- 大语言模型的哲学理念

- “压缩即智能”

- 无监督学习

- 通过压缩来理解数据的结构

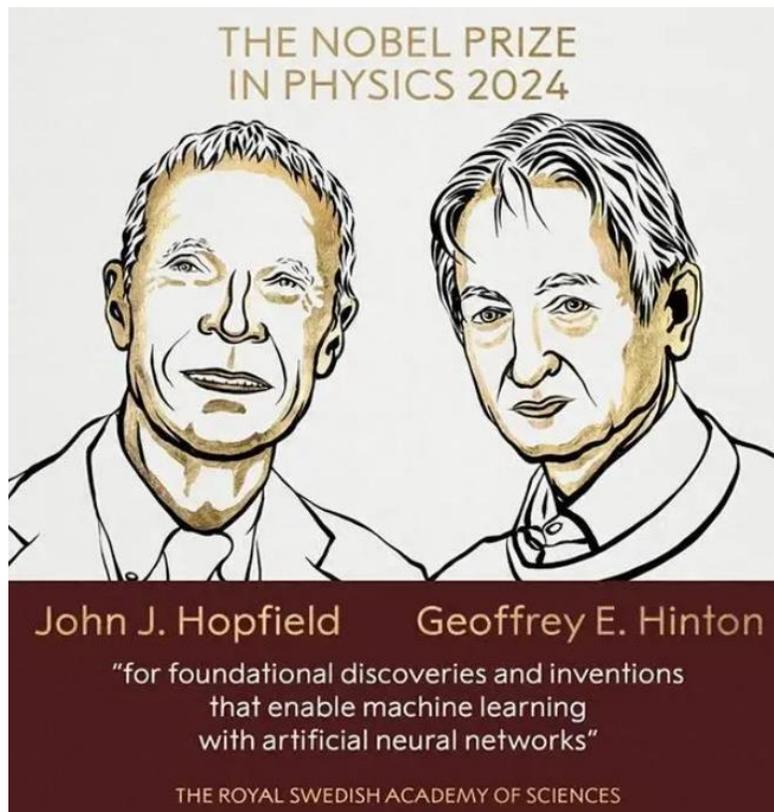
- 压缩越有效，越能找到数据共有的结构

- 好的压缩器也是好的预测器

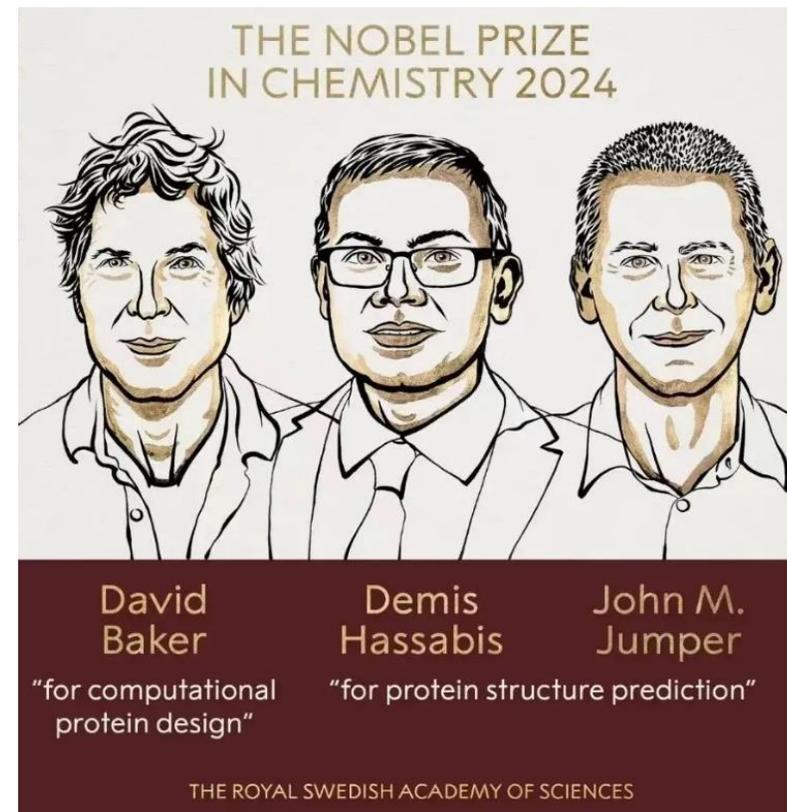


Ilya Sutskever

2024年诺贝尔物理学奖 & 化学奖



“基于人工神经网络实现机器学习
的基础性发现和发明”



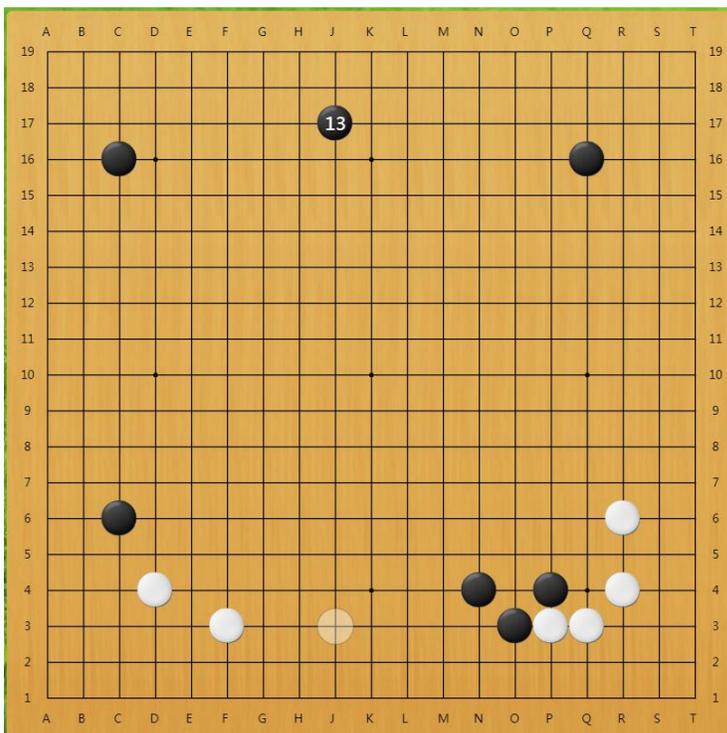
“计算蛋白质设计”、“蛋白质结构预测”



范式变革：人工智能 + 生物学

- 2019.10.23，张辰宇，人工智能生物学（AIBIO）
- 2022.8.16，中国生物物理学会人工智能生物学分会
- 徐涛院士，生物物理所：**AI就是物理！**

2016.3.9，李世石 vs. AlphaGo



中国科学：生命科学

2022年 第52卷 第3期: 291~300

SCIENTIA SINICA Vitae

lifecn.scichina.com

《中国科学》杂志社
SCIENCE CHINA PRESS

评述

中国知名大学及研究所专栏 南京大学生命科学学院专栏



人工智能生物学——生物学3.0

周祯^{1,2}, 闫超^{1,2}, 张辰宇^{1,2*}

1. 南京大学生命科学学院, 医药生物技术国家重点实验室, 南京 210023;

2. 南京大学人工智能生物医药技术研究院, 南京 210031

* 联系人, E-mail: cyzhang@nju.edu.cn

收稿日期: 2021-08-10; 接受日期: 2022-01-27; 网络版发表日期: 2022-02-23

摘要 现代生物学一方面取得了重大进展, 另一方面也面临着巨大的发展瓶颈, 而伴随着人工智能(artificial intelligence, AI)理论和技术的快速发展, AI与生物学深度融合的生物学V3.0——人工智能生物学(artificial intelligence biology, AIBIO)已经呼之欲出. 人们将人工智能生物学定义为利用人工智能的原理和手段来研究生命系统基本规律的科学. 其研究特点是: 动态整合多层面与多因素, 从而真正理解生命现象中的分子间相互作用与相互调控的规律, 解决生命科学中的重大基本问题. 作为一个全新的生命科学学科, 人工智能生物学将全面提升生物学研究的高度, 革新生物学研究的现有范式, 拓展生物学研究的范围, 实现生命科学和医学科学关键领域的实质性突破. 及时掌控并引领相关领域的研究, 对推动生命科学领域的基础研究、技术发展, 甚至对整个社会的进步

生物大模型综都至关重要。

生物信息学 vs. AI生物学

Primer

<https://doi.org/10.1038/s41587-024-02123-4>

Designing proteins with language models

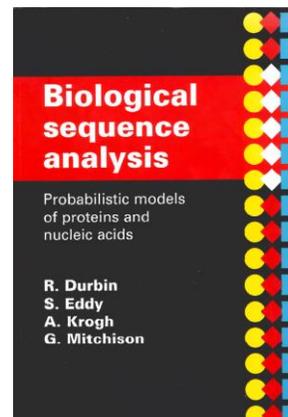
Foundations of protein language models

Fundamentally, protein language models aim to predict how likely we are to observe a particular protein sequence S given all the protein sequence data collected thus far. We denote a protein sequence $S = (s_1, s_2, \dots, s_N)$, where s_i represents the amino acid at position i in the sequence. As a first approximation, we might consider the probability of observing a protein as the joint probability of observing each of its constituent amino acids. Under this model, referred to as unigram, we calculate the probability of a sequence S as

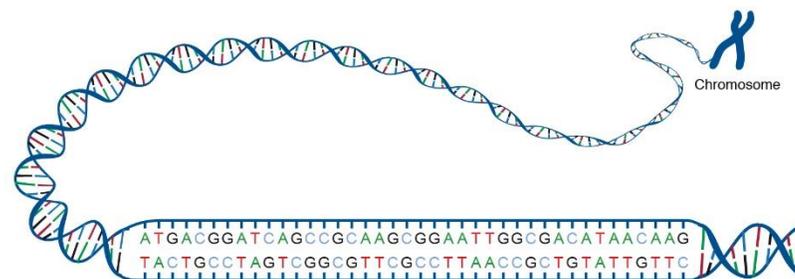
$$P(S) = \prod_i^N P(s_i)$$

In practice, to compute $P(S)$, we simply tabulate the frequency of each amino acid occurring in our sequence database and multiply the probabilities for the specific sequence S . However, proteins are not unordered collections of amino acids. Rather, the specific order in which we observe the amino acids is a critical determinant of structure and function. To capture this order dependency, we can use the preceding residues to inform the probability of the next amino acid. In an n -gram model, we multiply these contextualized probabilities to form the overall probability of the sequence:

$$P(S) = \prod_i^N P(s_i | s_{i-(n-1)}, \dots, s_{i-1})$$



生物信息学 (1.0)
“骰子模型”
线性、指令式、解析式



人工智能生物学 (2.0)
“语言模型”
非线性、人机交互式、
生成式

大语言模型的摩尔定律

- 尺度定律：“大力出奇迹”

Scaling Laws for Neural Language Models

Jared Kaplan*

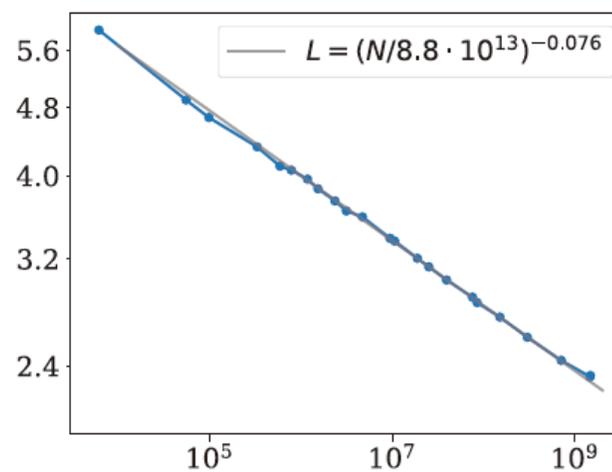
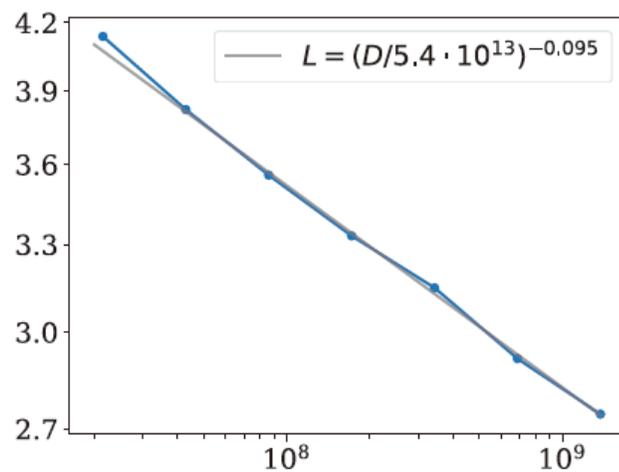
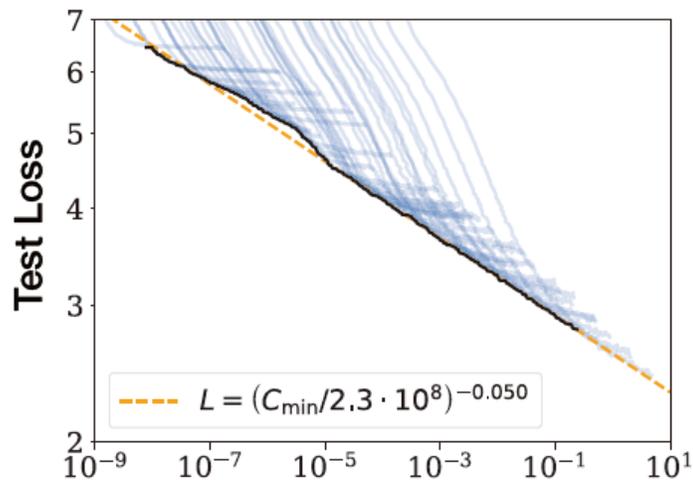
Johns Hopkins University, OpenAI

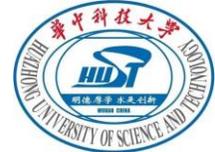
jaredk@jhu.edu

Sam McCandlish*

OpenAI

sam@openai.com





大语言模型的涌现特征

- 语境学习 (In-context learning)
- 小样本学习 (Few-shot learning)
- 零样本学习 (Zero-shot learning)
- 机器推理 : **思维链** (Chain-of-thought, CoT)
- ...

A Comprehensive Overview of Large Language Models

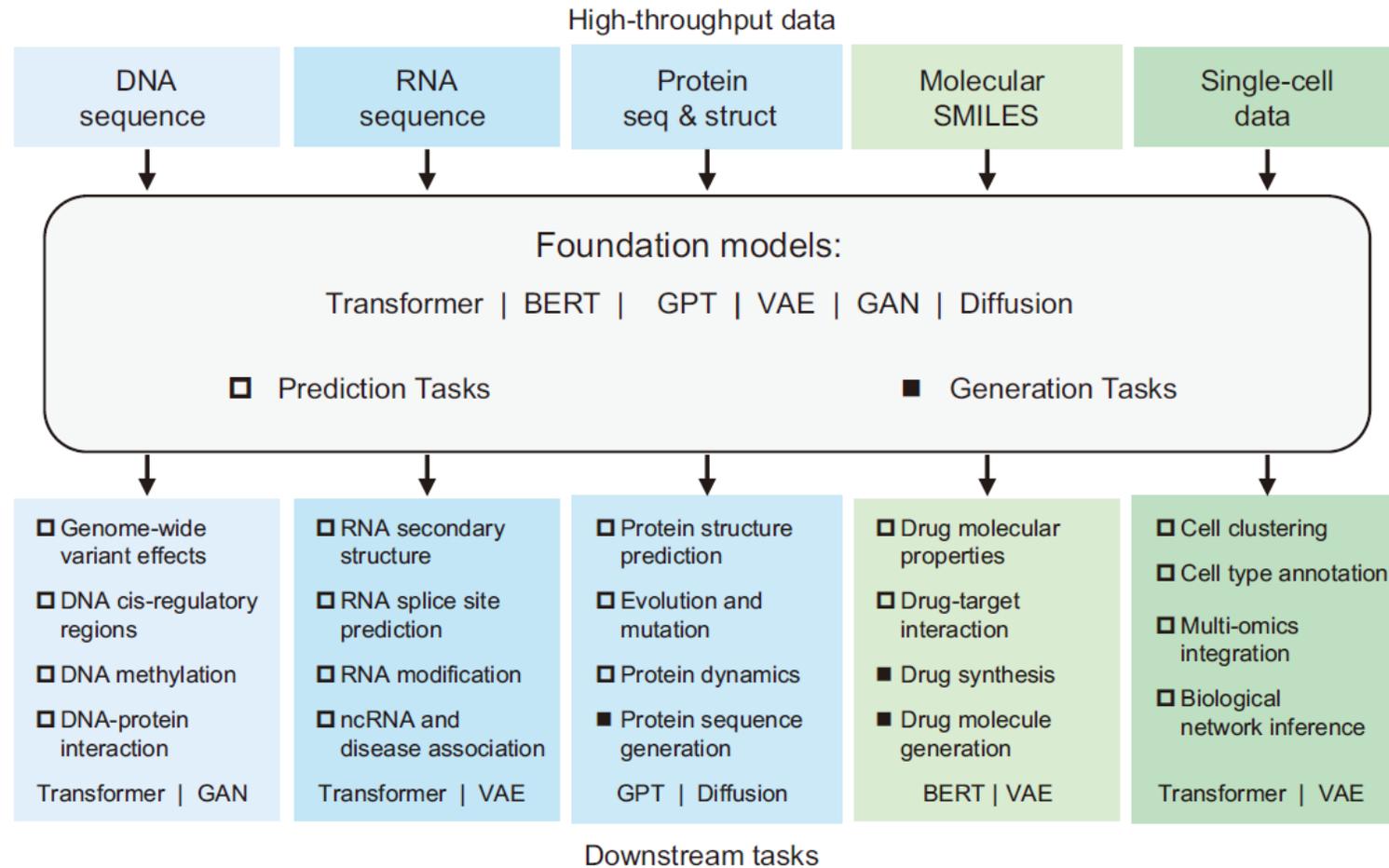
Humza Naveed^a, Asad Ullah Khan^{a,*}, Shi Qiu^{b,*}, Muhammad Saqib^{c,d,*}, Saeed Anwar^{e,f}, Muhammad Usman^{e,f}, Naveed Akhtar^{g,i},
Nick Barnes^h, Ajmal Mianⁱ

arXiv:2307.06435, 2023

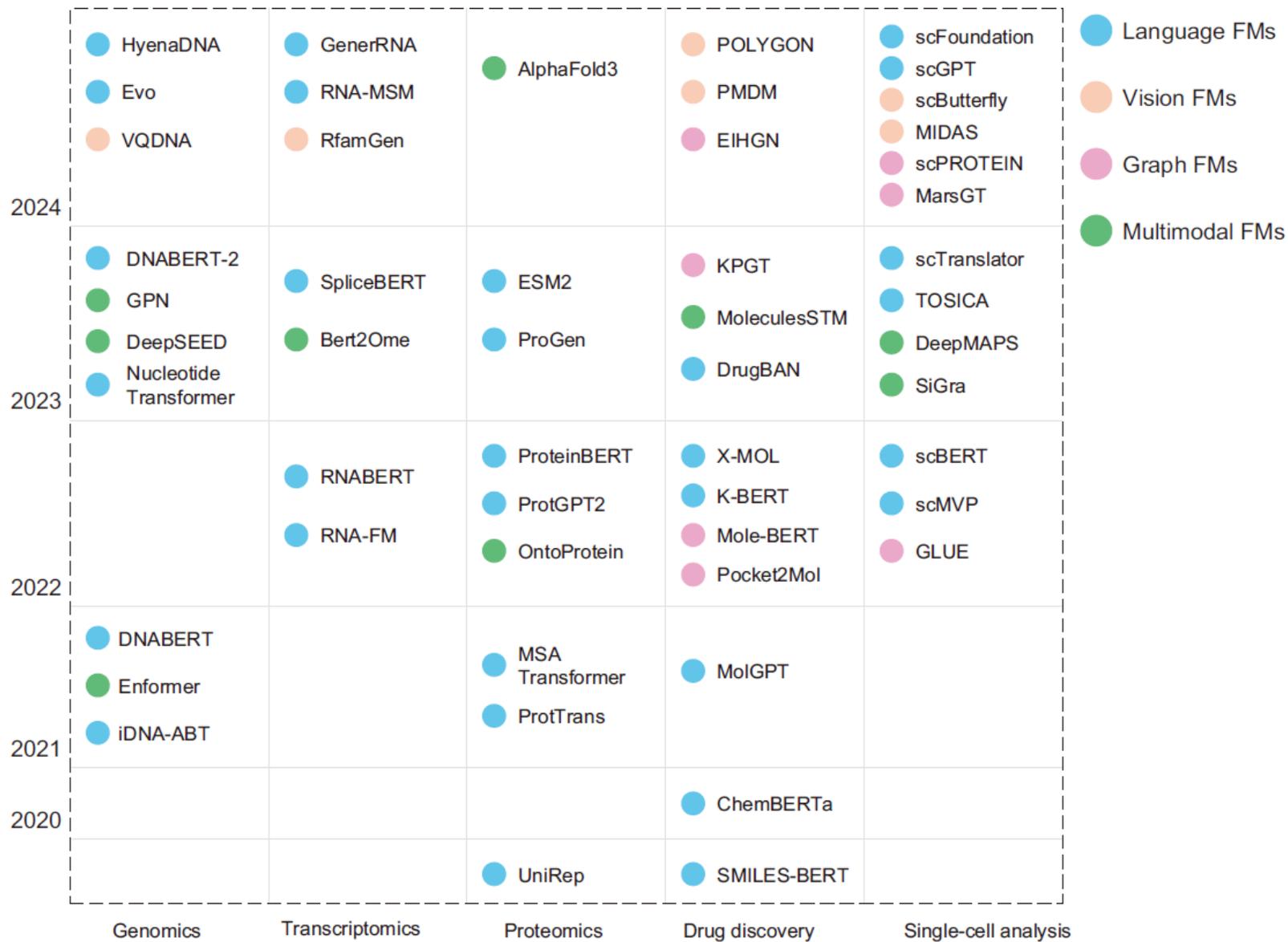
生物大模型/基础模型

Foundation models in bioinformatics

Fei Guo^{1,2}, Renchu Guan³, Yaohang Li⁴, Qi Liu⁵, Xiaowo Wang⁶, Can Yang⁷
and Jianxin Wang^{1,2,*}



生物大模型的演化





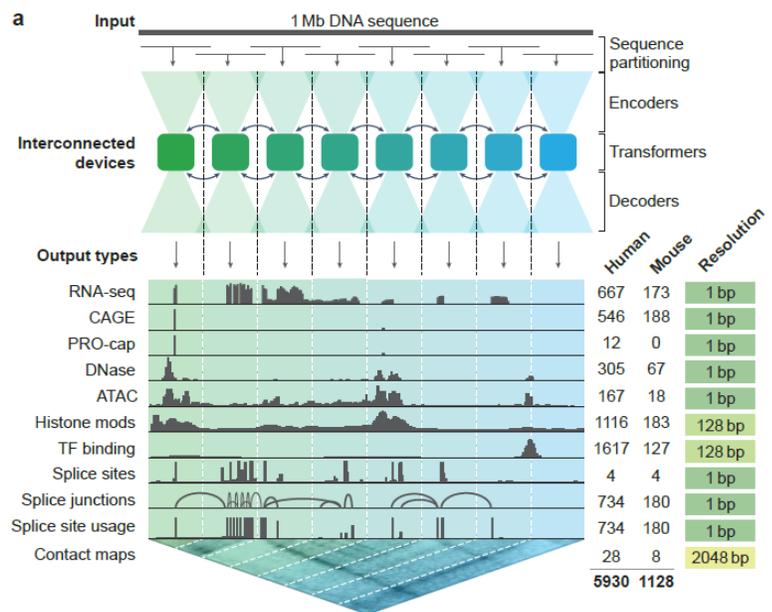
AI生物学的研究范式

- **大模型、生成式、多模态、小样本**
- **预训练 + 微调**
 - **X 自行构建预训练模型**
 - **√ 公共生物大模型 + 特定问题相关的数据资源**

大模型

AlphaGenome : 变异-效应预测

AlphaGenome: advancing regulatory variant effect prediction with a unified DNA sequence model



上下文序列 : 1 mb

预测精度 : 1bp bins

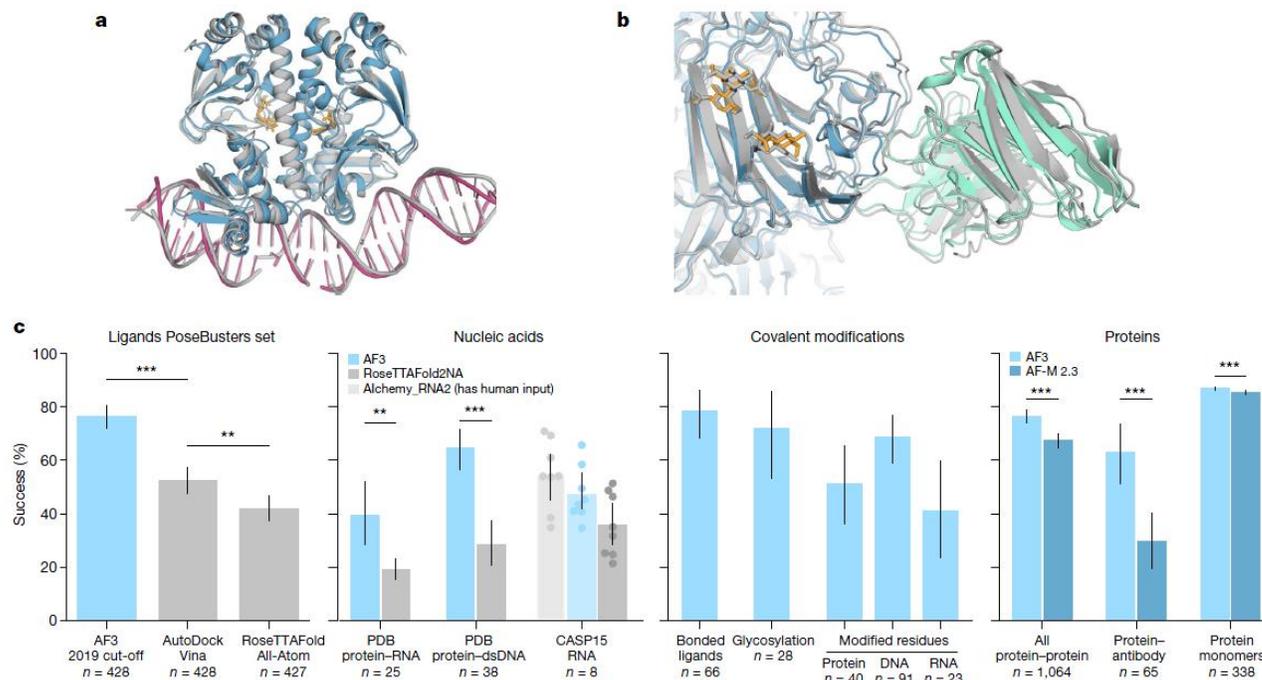
DeepMind

doi: <https://doi.org/10.1101/2025.06.25.661532> 生物大模型综合实践, 2025

AlphaFold 3 : 生物大分子相互作用预测

Article

Accurate structure prediction of biomolecular interactions with AlphaFold 3

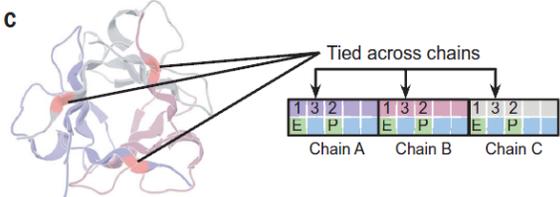
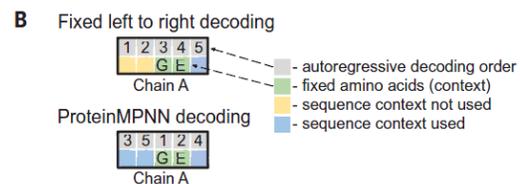
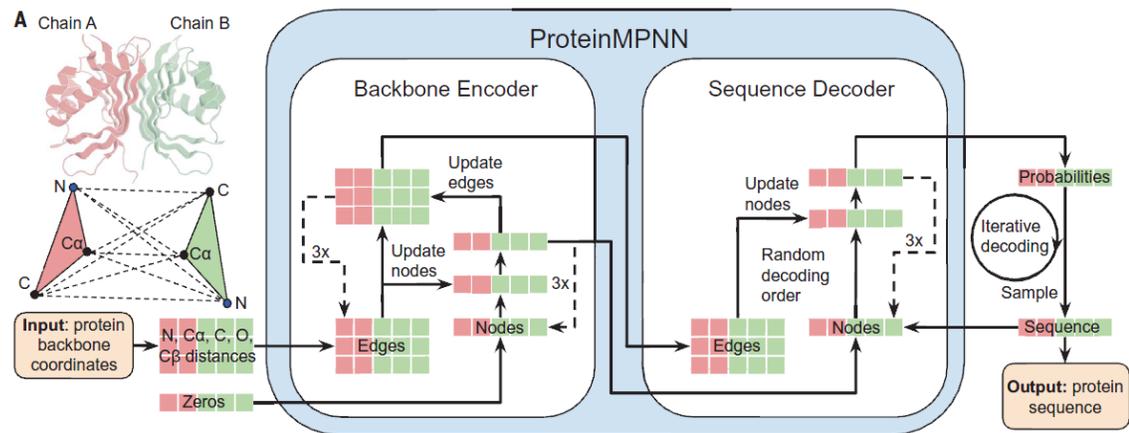


Nature, 2024, 630, 493-500

生成式

PROTEIN DESIGN

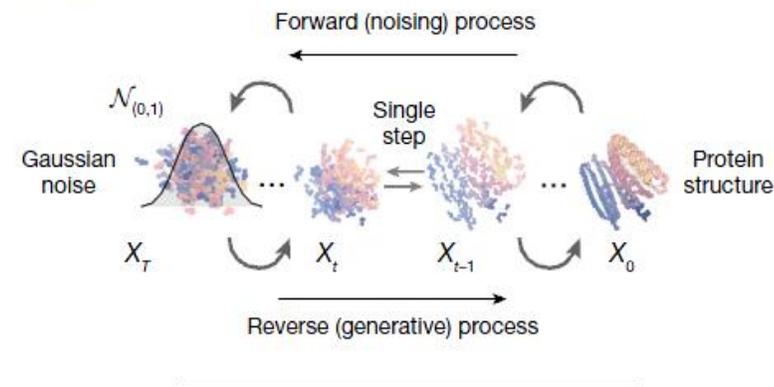
Robust deep learning-based protein sequence design using ProteinMPNN



Article

De novo design of protein structure and function with RFdiffusion

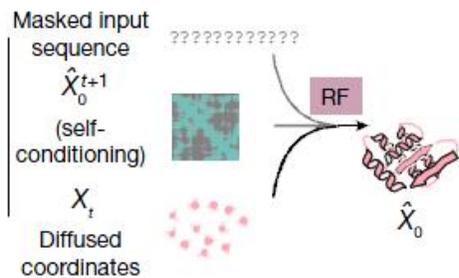
Diffusion model



RoseTTAFold

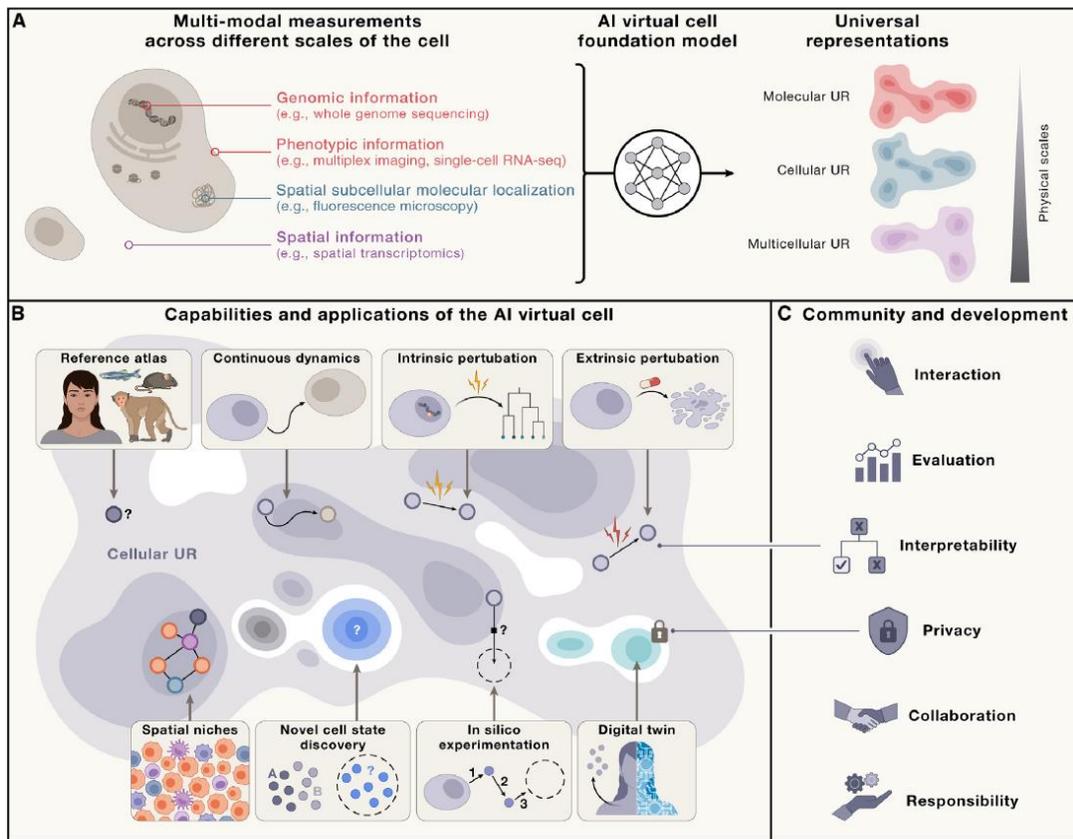


RFdiffusion



多模态

Perspective How to build the virtual cell with artificial intelligence: Priorities and opportunities



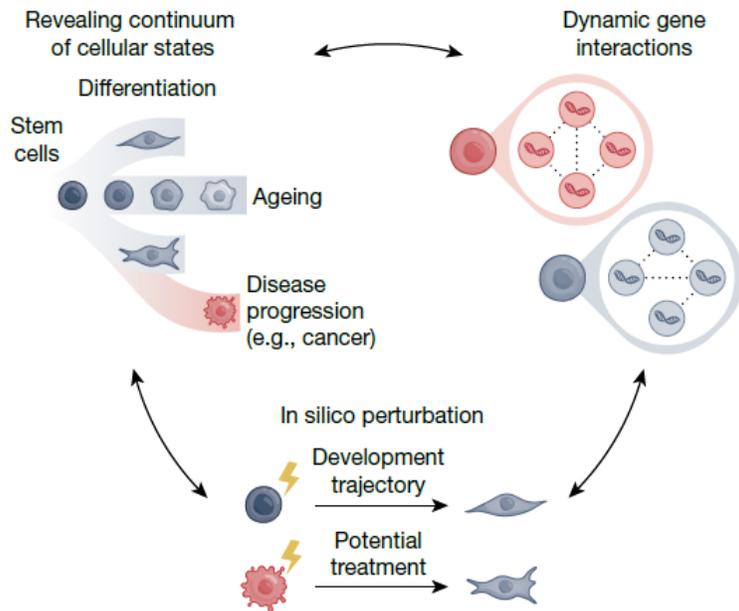
Perspective

Towards multimodal foundation models in molecular cell biology

<https://doi.org/10.1038/s41586-025-08710-y>
Received: 17 October 2023

Haotian Cui^{1,2,3}, Alejandro Tejada-Lapuerta^{4,5}, Maria Brbic^{6,7,8}, Julio Saez-Rodriguez^{9,10}, Simona Cristea^{11,12}, Hanl Goodarzi^{13,14}, Mohammad Lotfollahi^{15,16}, Fabian J. Theis^{4,17,18} & Bo Wang^{12,19,20}

c Applications to reconstruct cellular dynamics



新细胞类型识别，生物标志物发现，
基因调控推断，扰动分析

小样本

PanPep

nature machine intelligence

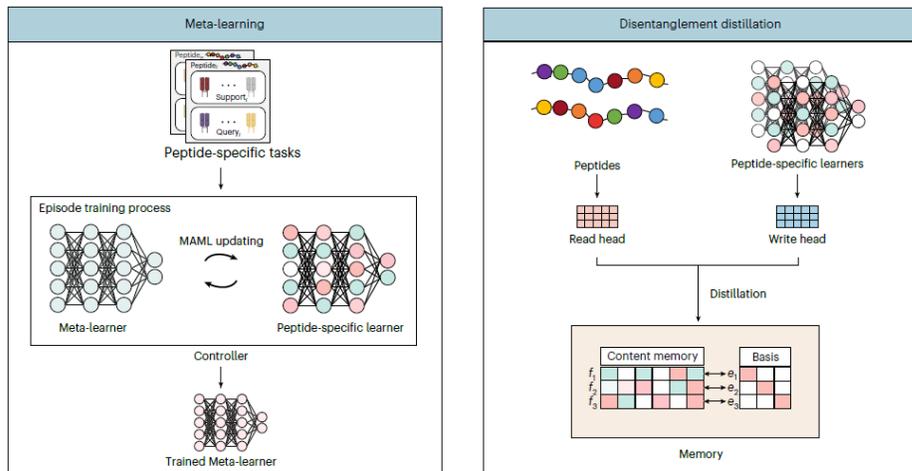
Article

<https://doi.org/10.1038/s42256-023-00619-3>

Pan-Peptide Meta Learning for T-cell receptor-antigen binding recognition

Received: 8 October 2022

Yicheng Gao^{1,2,5}, Yuli Gao^{1,2,5}, Yuxiao Fan^{1,2}, Chengyu Zhu^{1,2}, Zhiting Wei^{1,2}, Chi Zhou^{1,2}, Guohui Chuai^{1,2}, Qinchang Chen³, He Zhang² & Qi Liu^{1,2,3,4}✉



元学习预测T细胞受体-抗原结合

Nat Mach Intell, 2023, 5, 236-249

生物大模型综合实践, 2025

pFunk

nature metabolism

Article

<https://doi.org/10.1038/s42255-024-01093-w>

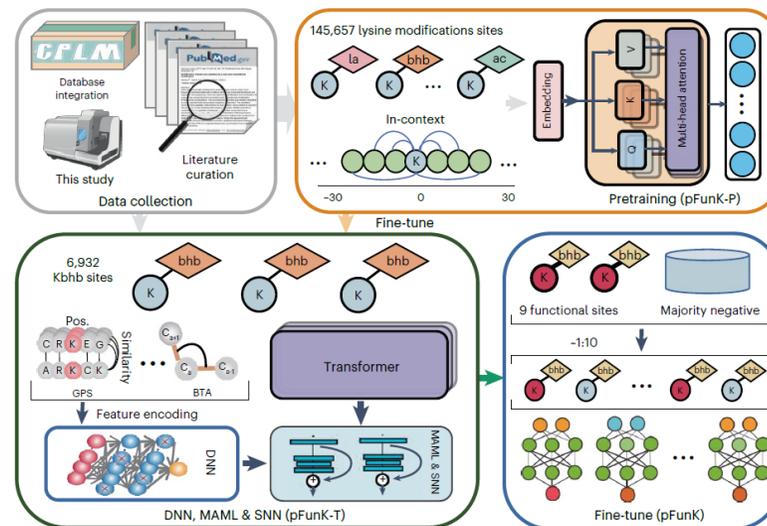
Ketogenic diet reshapes cancer metabolism through lysine β -hydroxybutyrylation

Received: 1 August 2023

Junhong Qin^{1,7}, Xinhe Huang^{2,7}, Shengsong Gou^{1,7}, Sitao Zhang^{1,7}, Yujie Gou², Qian Zhang¹, Hongyu Chen¹, Lin Sun³, Miaomiao Chen², Dan Liu², Cheng Han², Min Tang¹, Zihao Feng², Shenghui Niu¹, Lin Zhao¹, Yingfeng Tu¹, Zexian Liu⁴, Weimin Xuan⁵, Lunzhi Dai⁶, Da Jia¹✉ & Yu Xue^{2,6}✉

Accepted: 2 July 2024

Published online: 12 August 2024

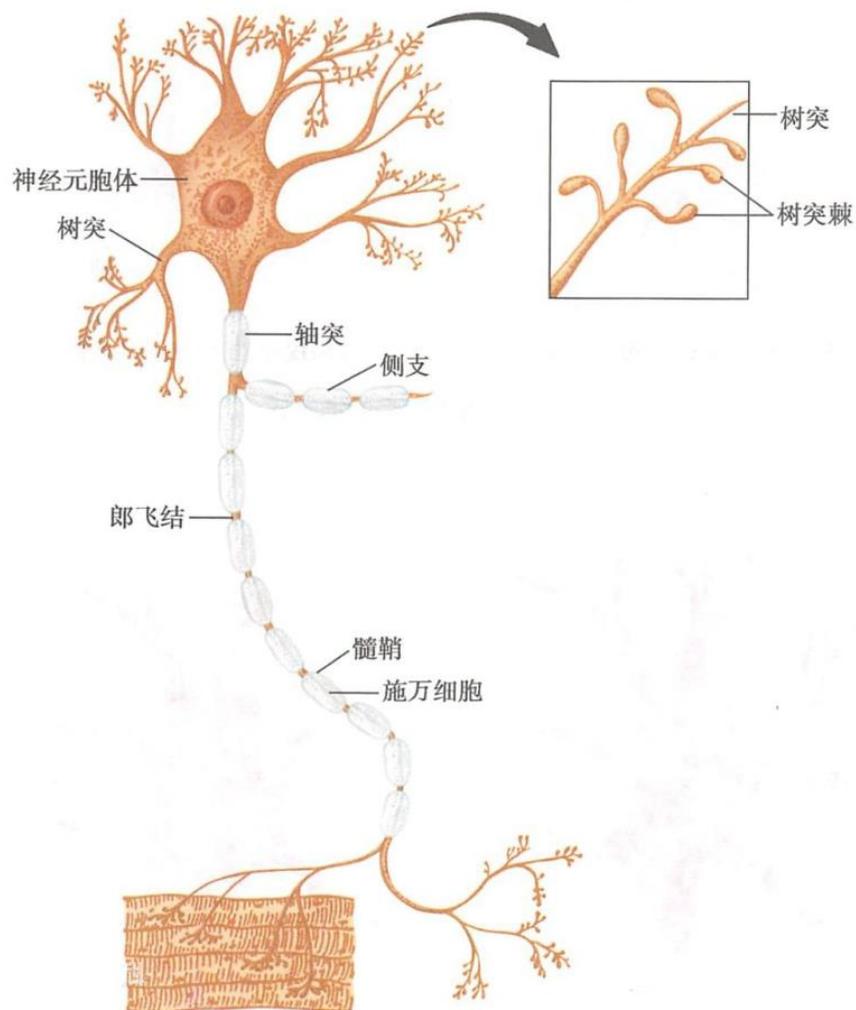


蛋白质修饰功能预测的小样本学习

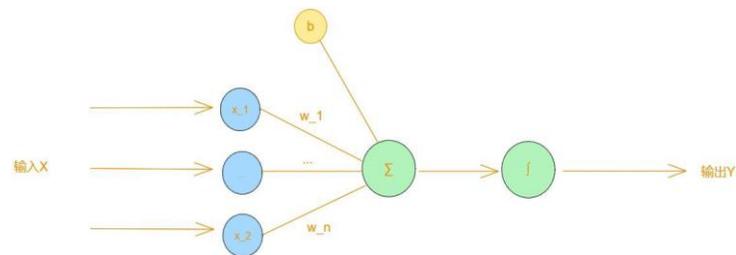
Nat Metab, 2024, 6, 1505-1528

神经网络与神经元

生物界中的神经元结构

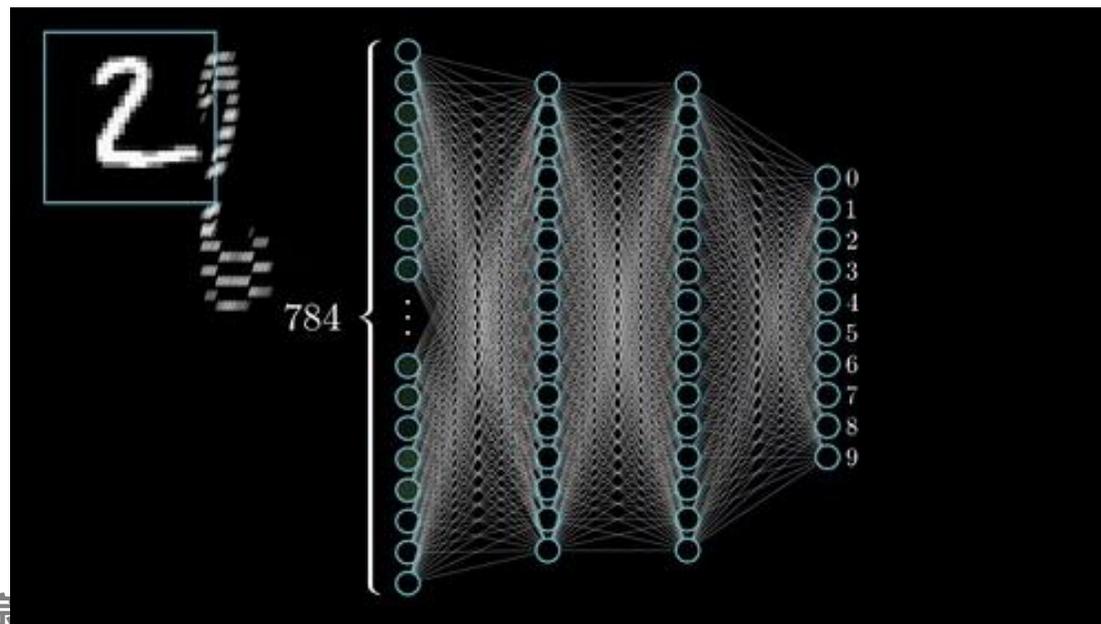


计算机中的神经元模型



$$y = f\left(\sum_{i=1}^n w_i x_i + b\right) = f\left(\vec{W} \cdot \vec{X} + b\right)$$

深度神经网络预测数字图像



神经网络类型

A mostly complete chart of Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen asimovinstitute.org

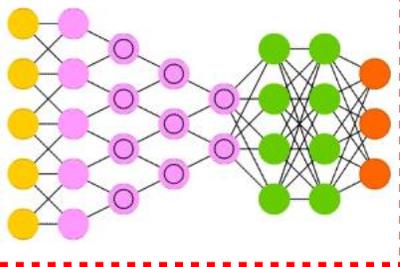
Perceptron (P) Feed Forward (FF) Radial Basis Network (RBF)



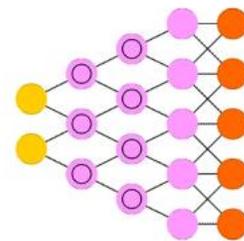
Deep Feed Forward (DFF)



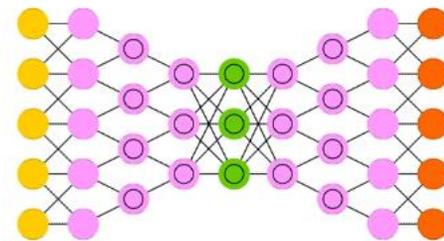
Deep Convolutional Network (DCN)



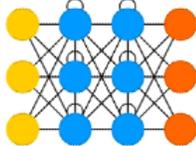
Deconvolutional Network (DN)



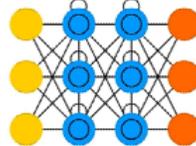
Deep Convolutional Inverse Graphics Network (DCIGN)



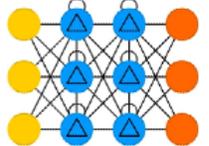
Recurrent Neural Network (RNN)



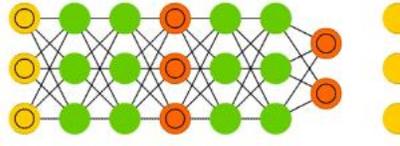
Long / Short Term Memory (LSTM)



Gated Recurrent Unit (GRU)



Generative Adversarial Network (GAN)



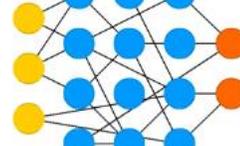
Liquid State Machine (LSM)



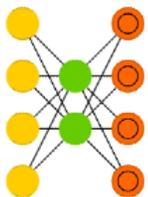
Extreme Learning Machine (ELM)



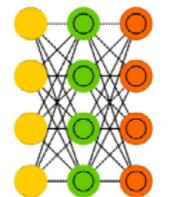
Echo State Network (ESN)



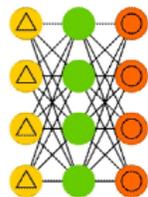
Auto Encoder (AE)



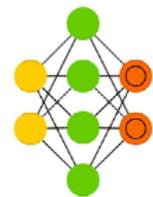
Variational AE (VAE)



Denoising AE (DAE)



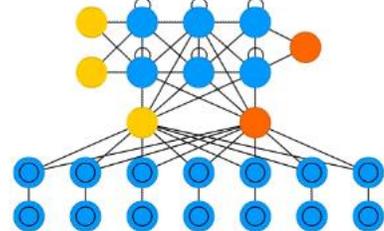
Sparse AE (SAE)



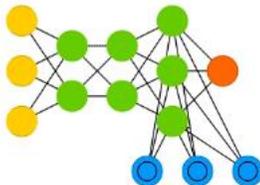
Deep Residual Network (DRN)



Differentiable Neural Computer (DNC)



Neural Turing Machine (NTM)



Markov Chain (MC)



Hopfield Network (HN)



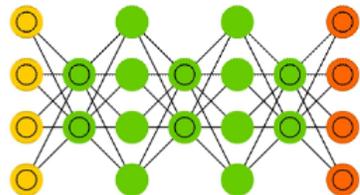
Boltzmann Machine (BM)



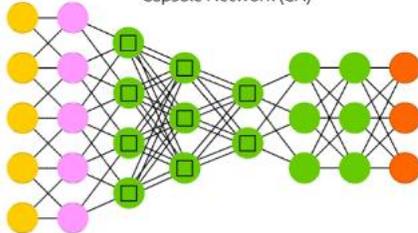
Restricted BM (RBM)



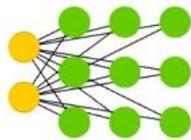
Deep Belief Network (DBN)



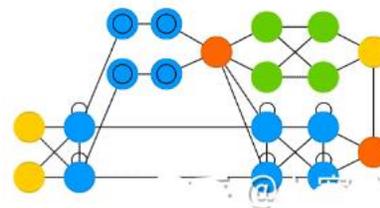
Capsule Network (CN)



Kohonen Network (KN)

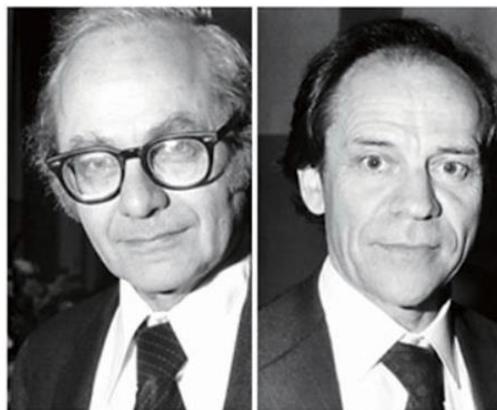
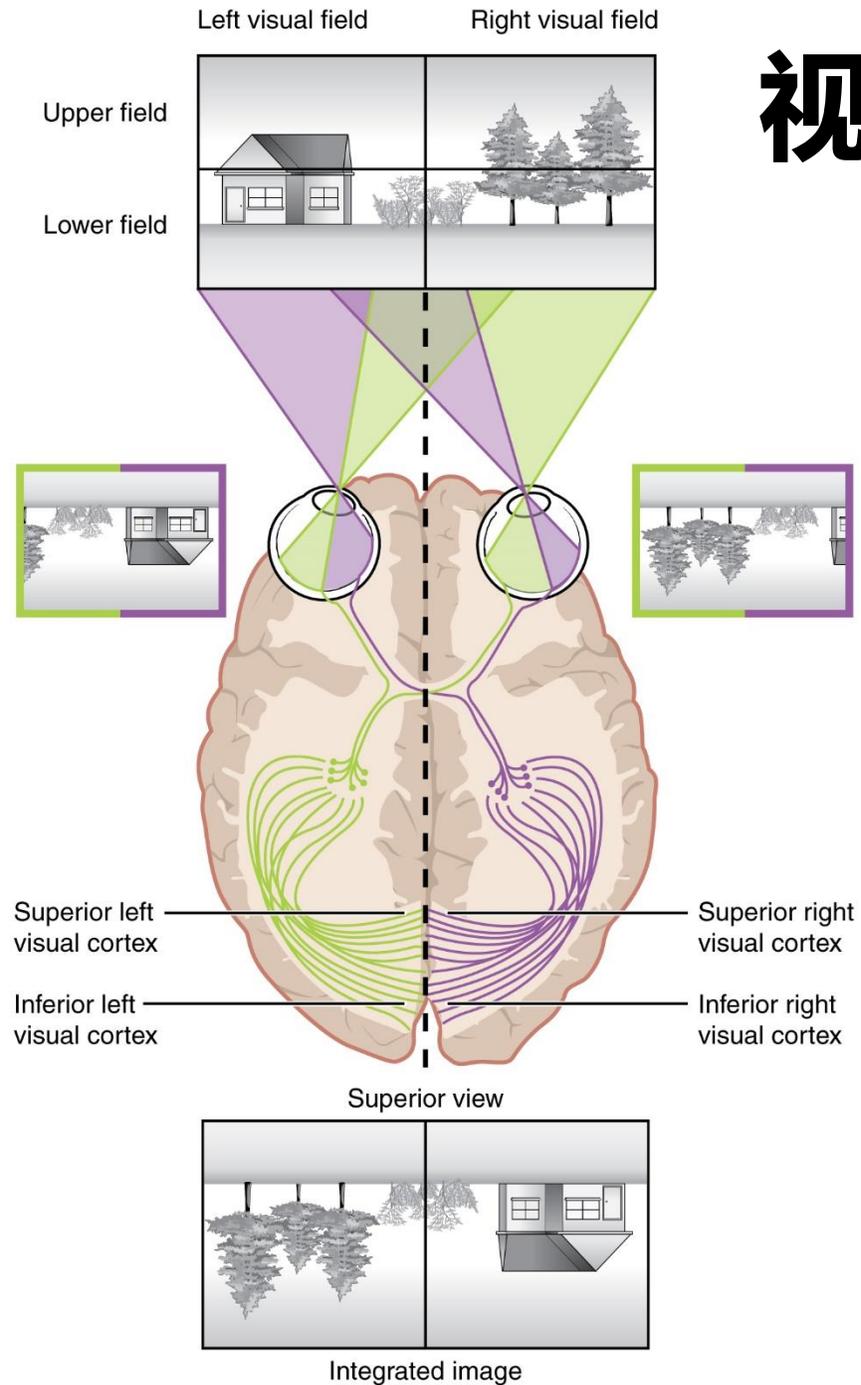


Attention Network (AN)



- Input Cell
- Backfed Input Cell
- Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Gated Memory Cell
- Kernel
- Convolution or Pool

视觉系统的信息处理



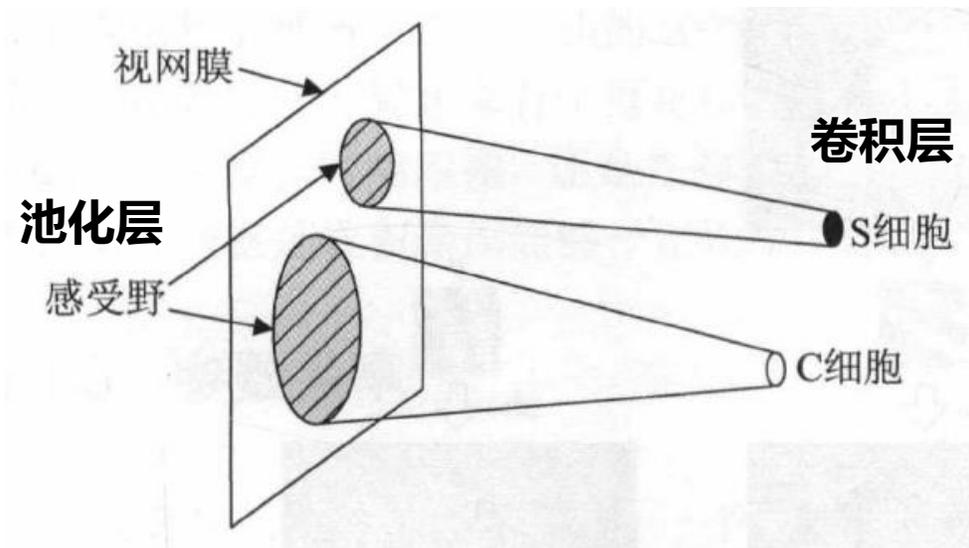
David H. Hubel

Torsten N. Wiesel

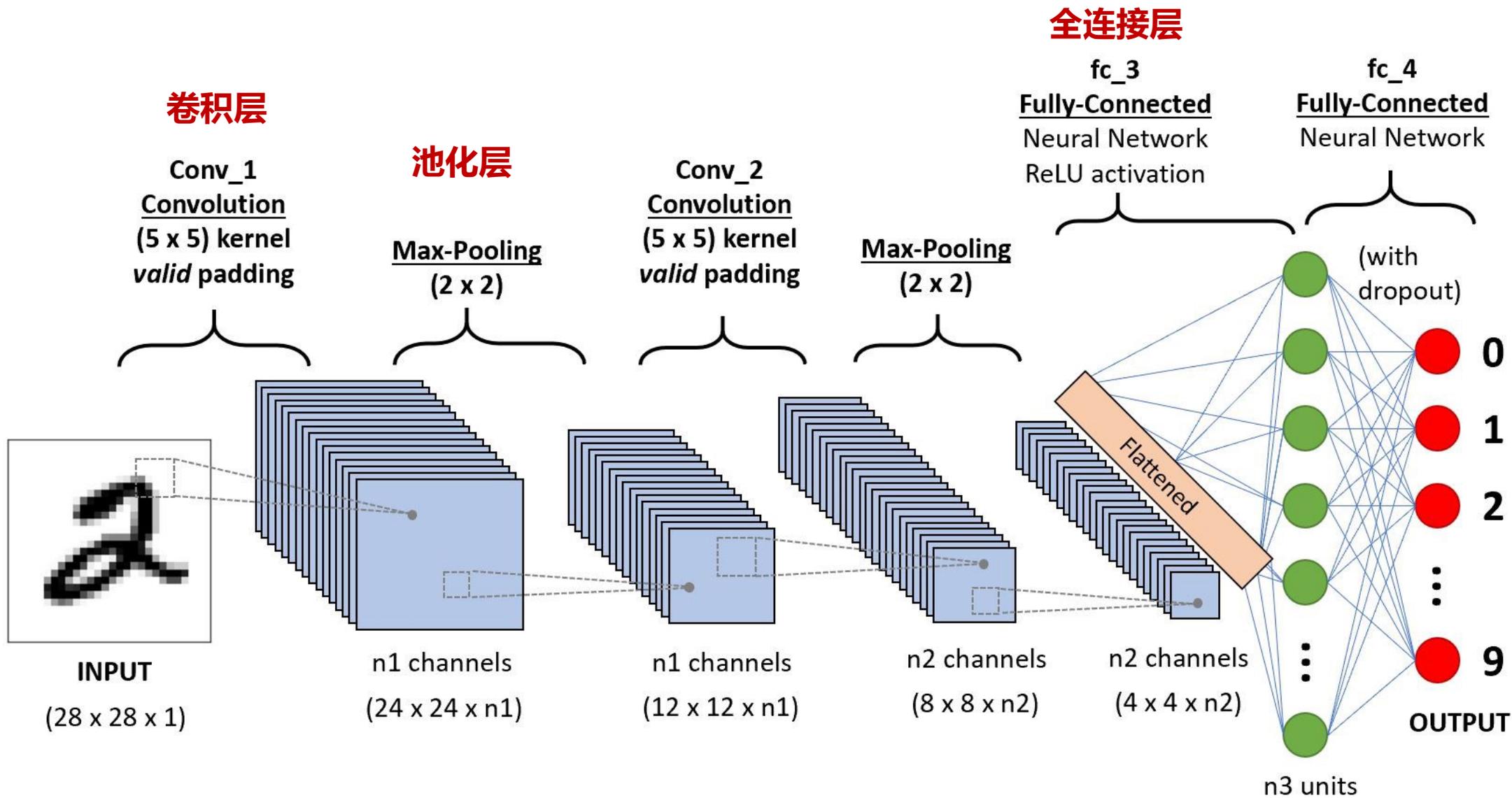
1981年诺贝尔生理学或医学奖



福岛邦彦，1980年提出
Neocognitron



卷积神经网络架构



卷积运算与二维卷积

- 离散形式的卷积：

$$s(t) = \sum_{\alpha=-\infty}^{+\infty} x(\alpha)w(t-\alpha)$$

1 <small>*1</small>	1 <small>*0</small>	1 <small>*1</small>	0	0
0 <small>*0</small>	1 <small>*1</small>	1 <small>*0</small>	1	0
0 <small>*1</small>	0 <small>*0</small>	1 <small>*1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature

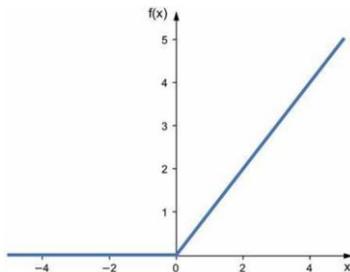
- 二维卷积：

$$S(i, j) = \sum_m \sum_n I(i+m, j+n)K(m, n).$$

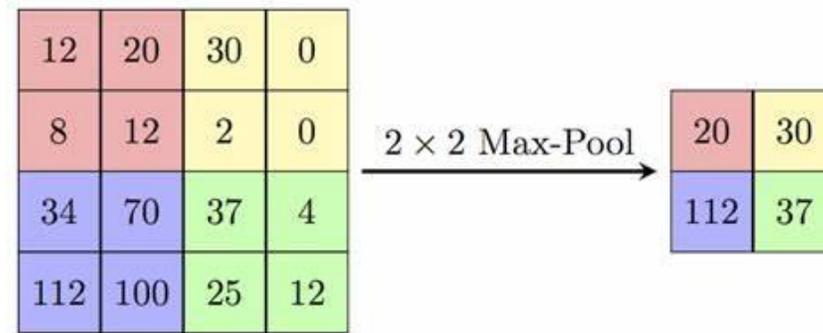
激活函数、卷积层、池化层和输出层

• 神经元激活函数

$$\text{ReLU}(x) = \begin{cases} x, & x \geq 0 \\ 0, & x < 0 \end{cases}$$



• 最大池化 (Max pooling)



• 卷积层

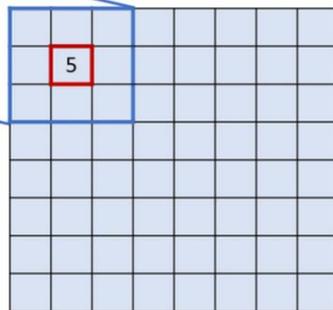
Source layer

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

Convolutional kernel

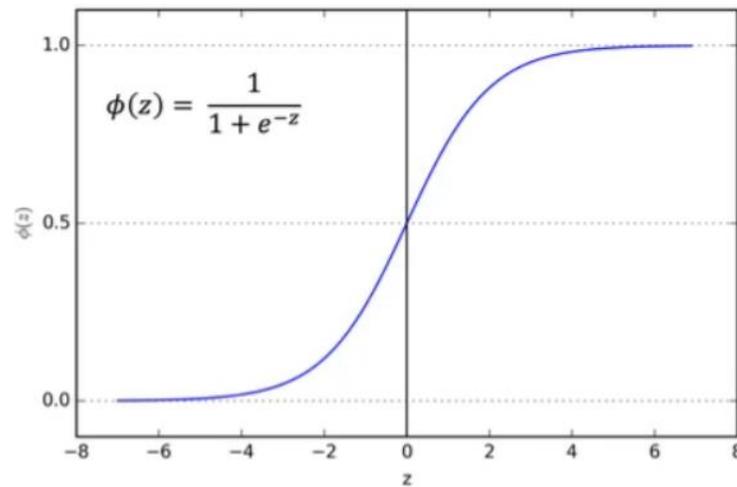
-1	0	1
2	1	2
1	-2	0

Destination layer



$$\begin{aligned} & (-1 \times 5) + (0 \times 2) + (1 \times 6) + \\ & (2 \times 4) + (1 \times 3) + (2 \times 4) + \\ & (1 \times 3) + (-2 \times 9) + (0 \times 2) = 5 \end{aligned}$$

• 输出：sigmoid函数



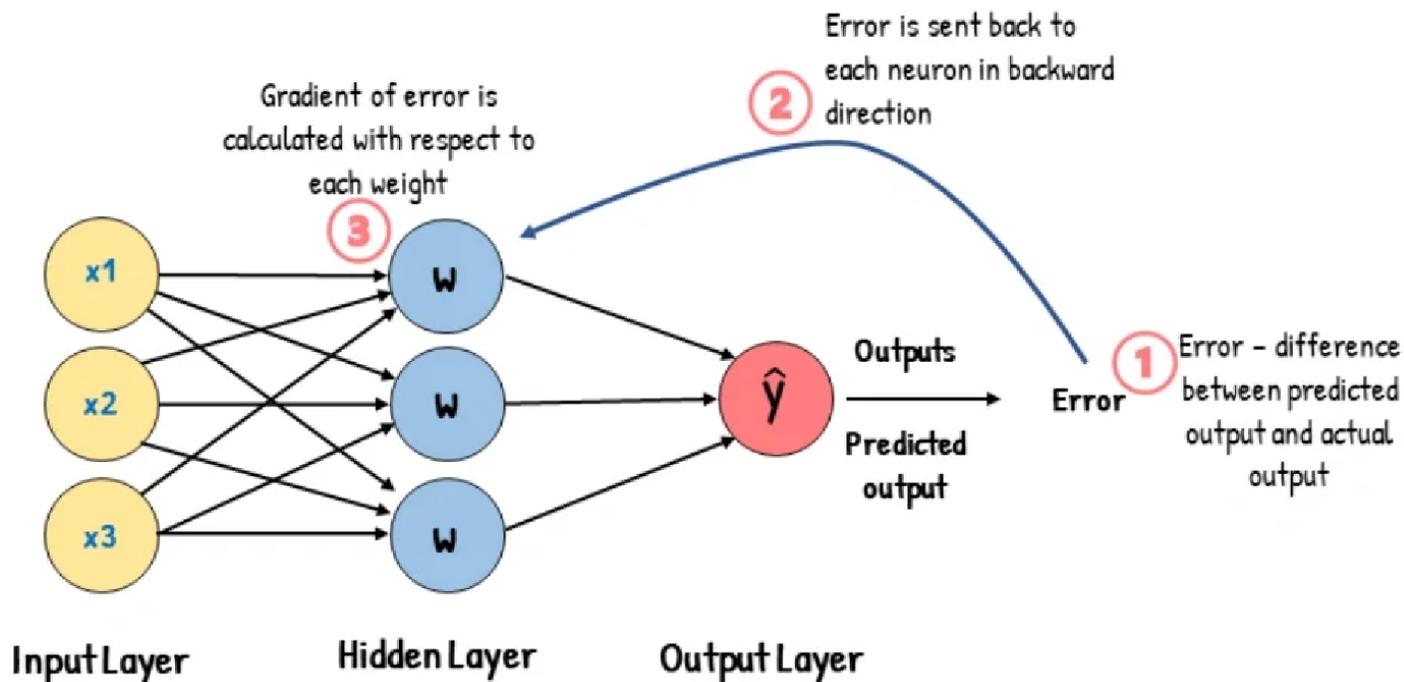
学习率与误差逆传播

$$*W_x = W_x - a \left(\frac{\partial \text{Error}}{\partial W_x} \right)$$

Annotations for the equation:

- $*W_x$: New weight
- W_x : Old weight
- a : Learning rate
- $\left(\frac{\partial \text{Error}}{\partial W_x} \right)$: Derivative of Error with respect to weight

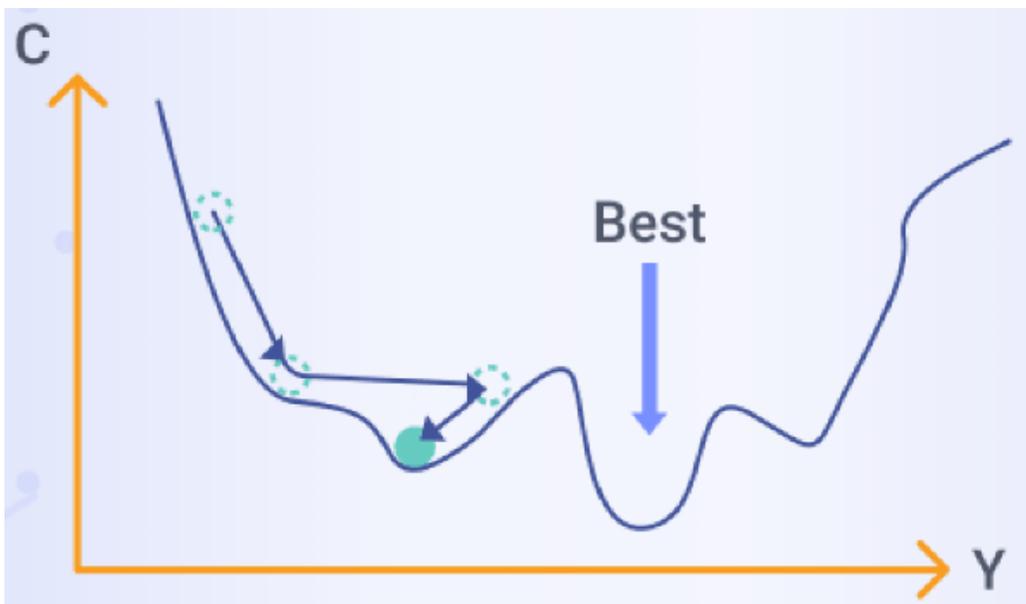
Backpropagation



模型优化

- 随机梯度下降SGD
- Stochastic Gradient Descent

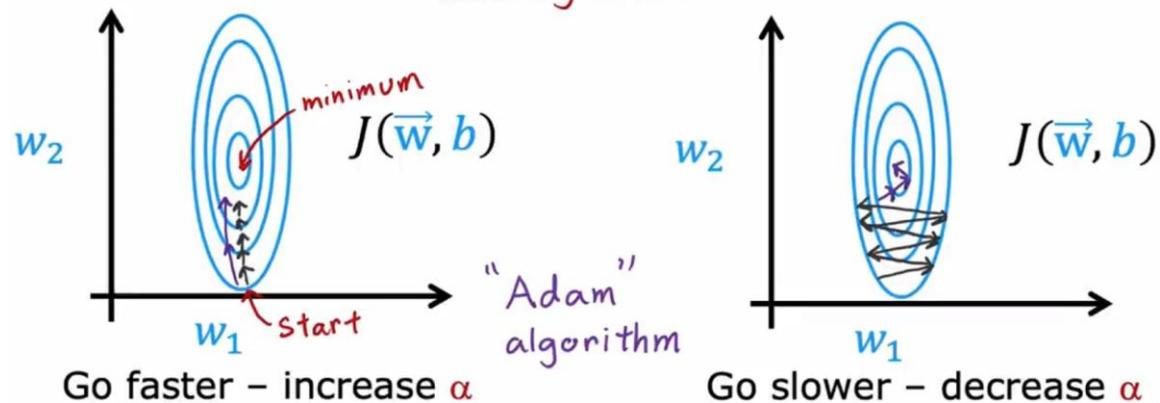
- 自适应矩估计Adam
- Adaptive Moment Estimation



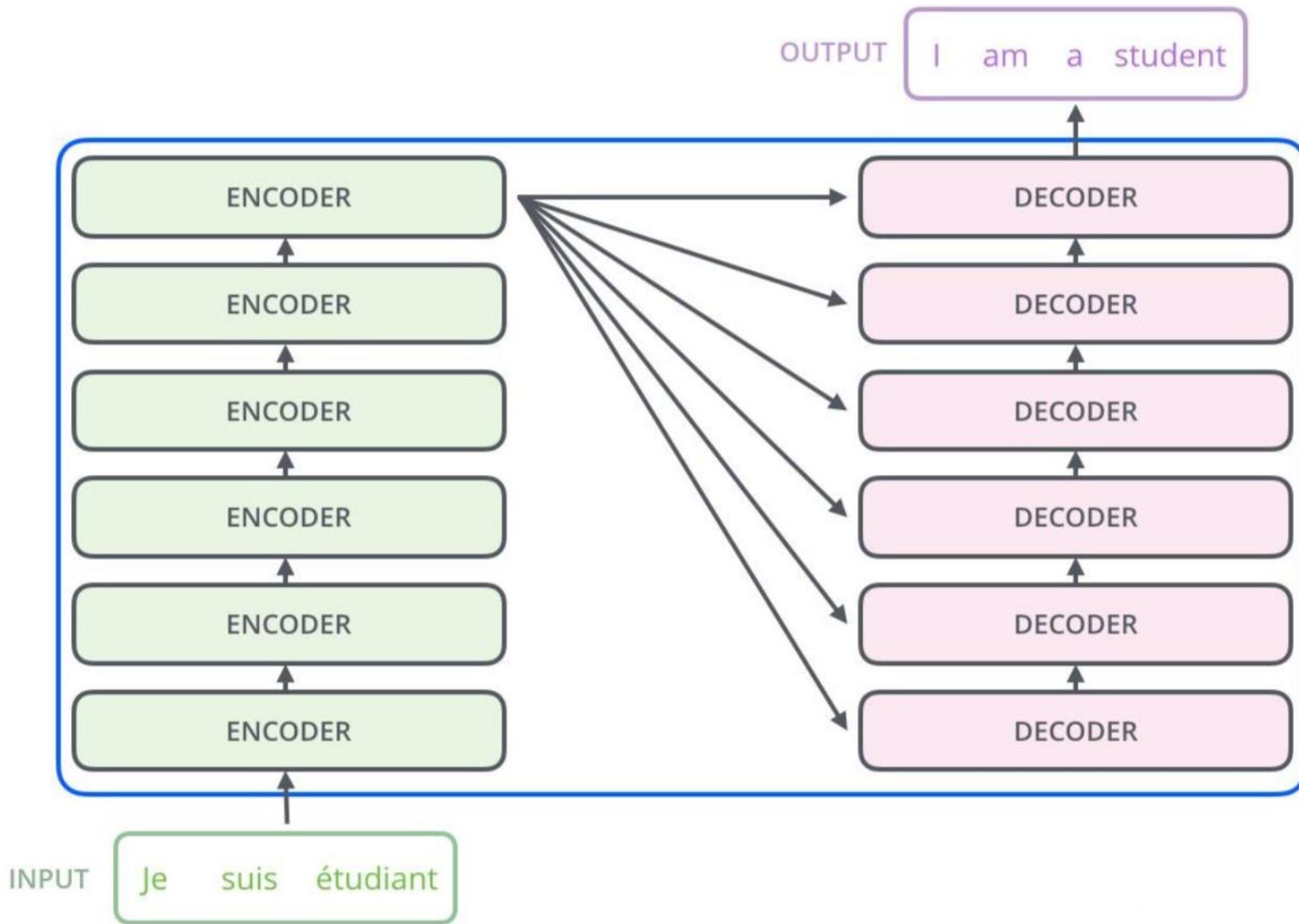
Gradient Descent

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

learning rate



Transformer : 编码器与解码器



词向量嵌入

• 1. 词表构建

- { "I" : 0, "love" : 1, "NLP" : 2, "<PAD>" : 3 } # PAD为填充符号
- 句子 "I love NLP" 的索引序列为 : [0, 1, 2]

• 2. 嵌入层

- 假设嵌入维度 $d_{\text{model}} = 4$ (实际中通常为512或768)

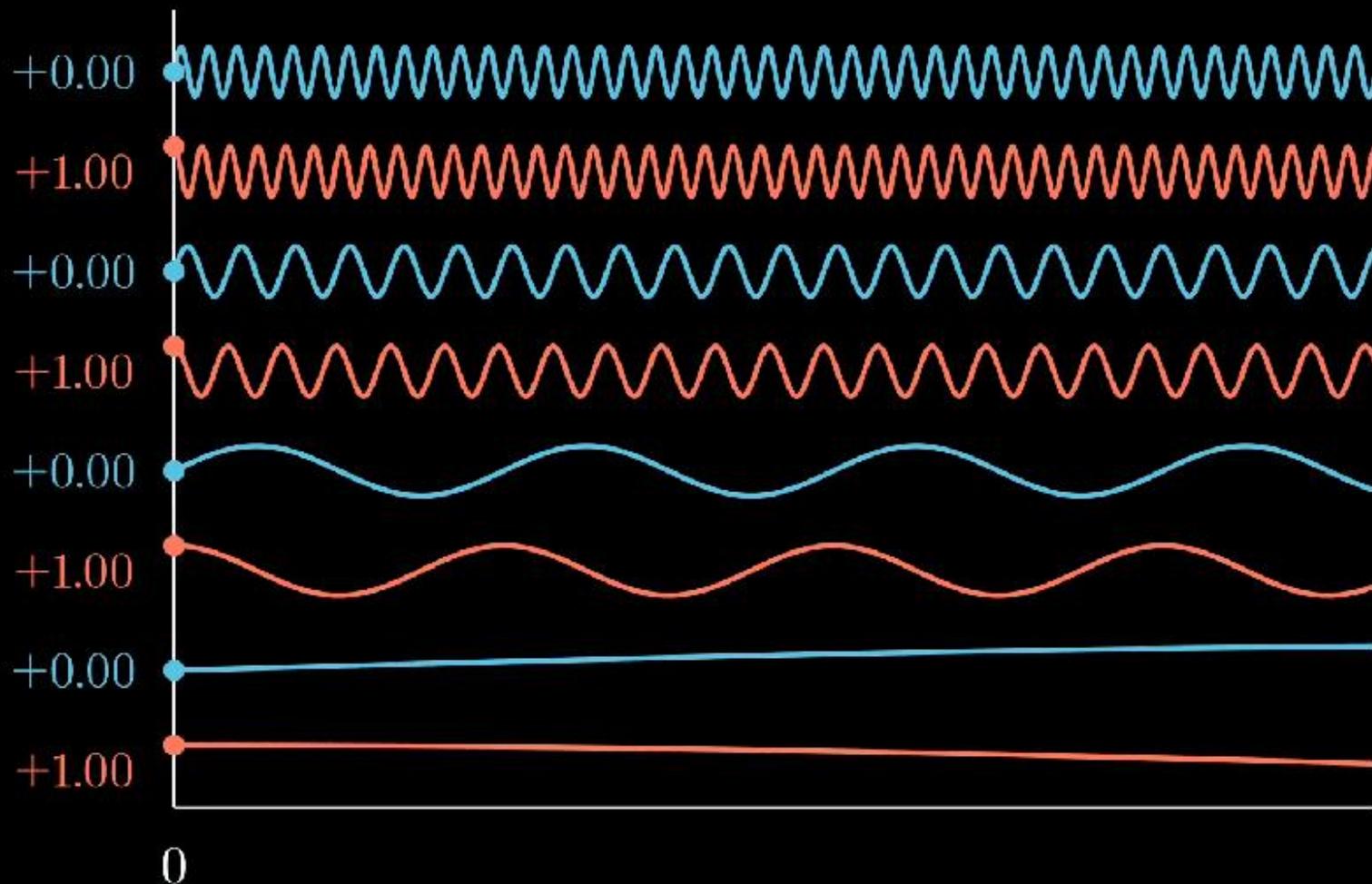
嵌入矩阵 E 的形状为 [词表大小, d_{model}] , 例如 :

```
1 E = [  
2     [0.1, 0.2, 0.3, 0.4], # "I"的向量  
3     [0.5, 0.6, 0.7, 0.8], # "love"的向量  
4     [0.9, 1.0, 1.1, 1.2], # "NLP"的向量  
5     [0.0, 0.0, 0.0, 0.0]  # "<PAD>"的向量 (全零)  
6 ]
```

通过索引查找得到词向量序列 :

```
1 Embedded_words = [  
2     [0.1, 0.2, 0.3, 0.4], # "I"  
3     [0.5, 0.6, 0.7, 0.8], # "love"  
4     [0.9, 1.0, 1.1, 1.2]  # "NLP"  
5 ]
```

正余弦位置编码



$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d}}\right)$$

$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$

正余弦位置编码

$$PE_{pos,2i} = \sin(pos/10000^{2i/d_{model}})$$

$$PE_{pos,2i+1} = \cos(pos/10000^{2i/d_{model}})$$

- pos : 单词在句子中的位置 (从0开始)
- d_{model} : 嵌入维度 (此处为4)
- i : 维度索引 (0到3)

$$\begin{aligned} PE(0, 0) &= \sin(0 / 10000^{(0/4)}) = 0 \\ PE(0, 1) &= \cos(0 / 10000^{(0/4)}) = 1 \\ PE(0, 2) &= \sin(0 / 10000^{(2/4)}) = 0 \\ PE(0, 3) &= \cos(0 / 10000^{(2/4)}) = 1 \end{aligned}$$

$$\begin{aligned} PE(1, 0) &= \sin(1 / 10000^{(0/4)}) = 0.0001 \\ PE(1, 1) &= \cos(1 / 10000^{(0/4)}) = 1.0 \\ PE(1, 2) &= \sin(1 / 10000^{(2/4)}) = 0.0087 \\ PE(1, 3) &= \cos(1 / 10000^{(2/4)}) = 0.9999 \end{aligned}$$

• 对于位置0 (第一个词 “I”) , 位置编码向量是 :

• 对于位置1 (第二个词 “love”) , 位置编码向量是 :

• 位置2 (“NLP”) : 类似计算 , 结果 : [0.0002, 1.0, 0.0175, 0.9998]

完整位置编码矩阵 :

```
1 PE = [  
2     [0, 1, 0, 1],           # 位置0  
3     [0.0001, 1.0, 0.0087, 0.9999], # 位置1  
4     [0.0002, 1.0, 0.0175, 0.9998] # 位置2  
5 ]
```

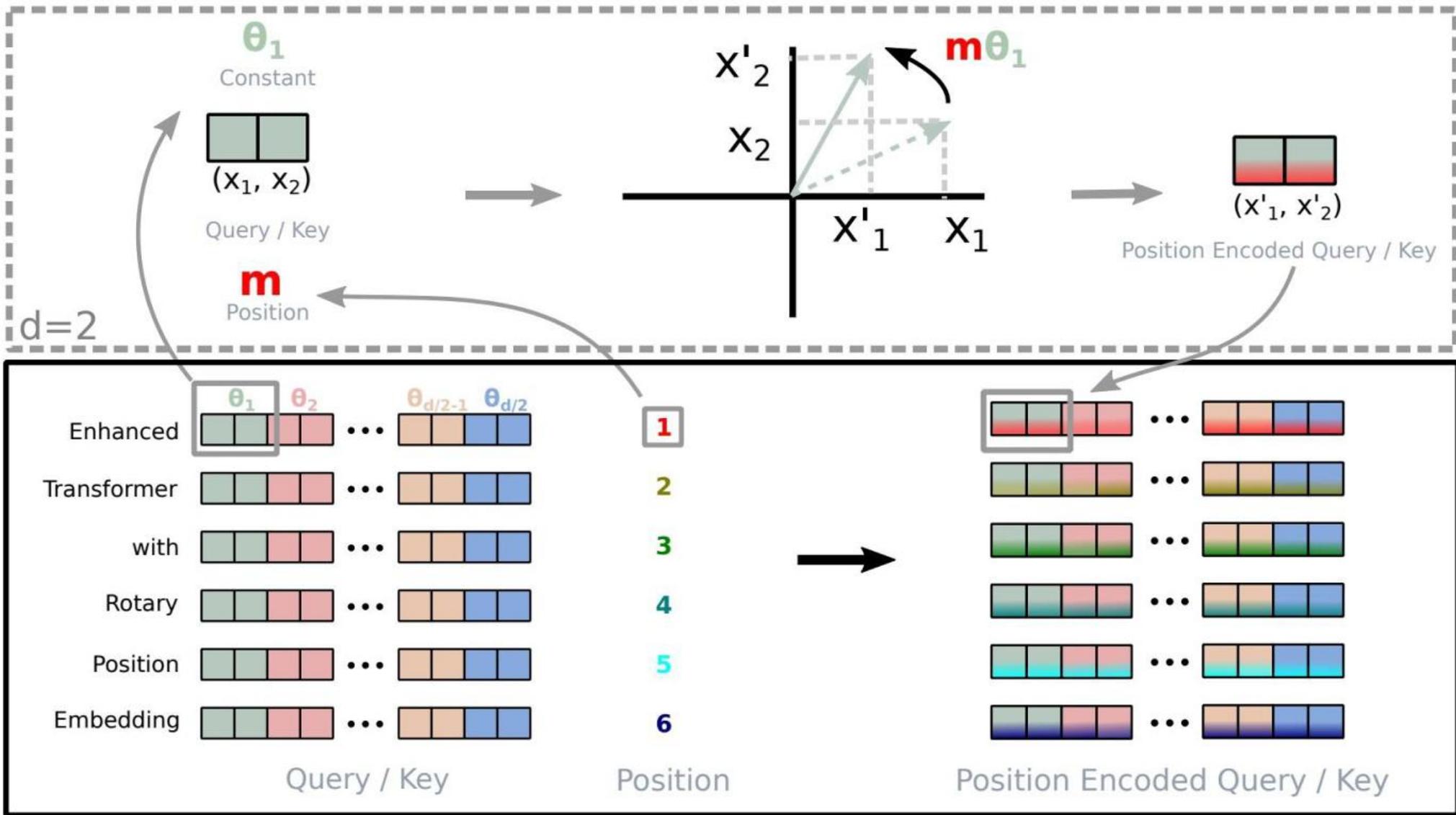


词向量与位置编码相加

- 将词向量和位置编码逐元素相加，得到最终输入向量：
 - $\text{Final_input} = \text{Embedded_words} + \text{PE}$
- 例如，“I”的最终向量为：
 - $[0.1+0, 0.2+1, 0.3+0, 0.4+1] = [0.1, 1.2, 0.3, 1.4]$
- 最终结果：

```
1 Final_input = [  
2     [0.1, 1.2, 0.3, 1.4], # "I" + 位置0  
3     [0.5+0.0001, 0.6+1.0, 0.7+0.0087, 0.8+0.9999], # "love" + 位置1  
4     [0.9+0.0002, 1.0+1.0, 1.1+0.0175, 1.2+0.9998] # "NLP" + 位置2  
5 ]
```

旋转位置编码RoPE



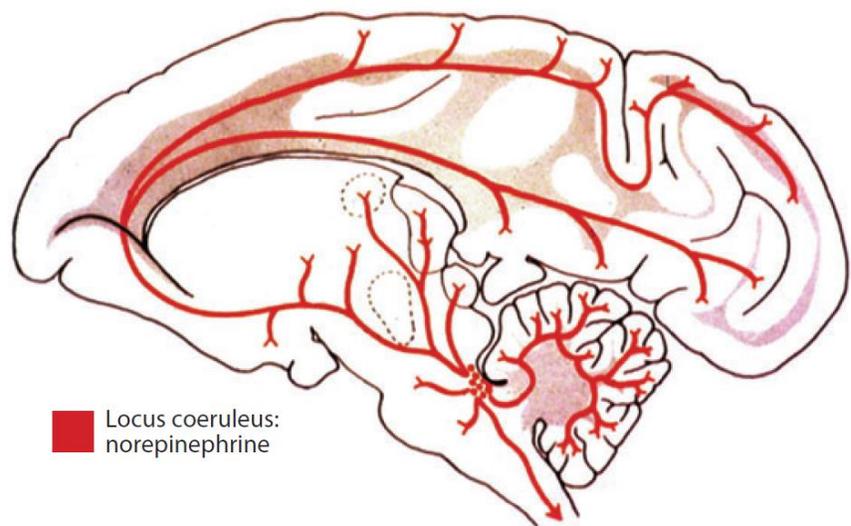
人脑的注意力系统

The Attention System of the Human Brain: 20 Years After

Steven E. Petersen¹ and Michael I. Posner²

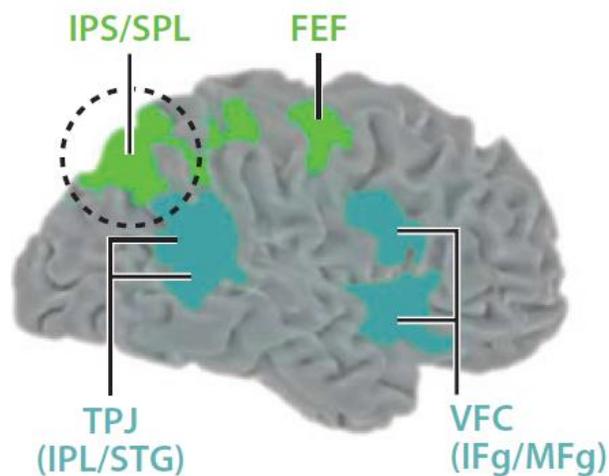
- 警觉系统 (Alerting)、定向系统 (Orienting) 和执行控制系统 (Executive Control)

Alerting

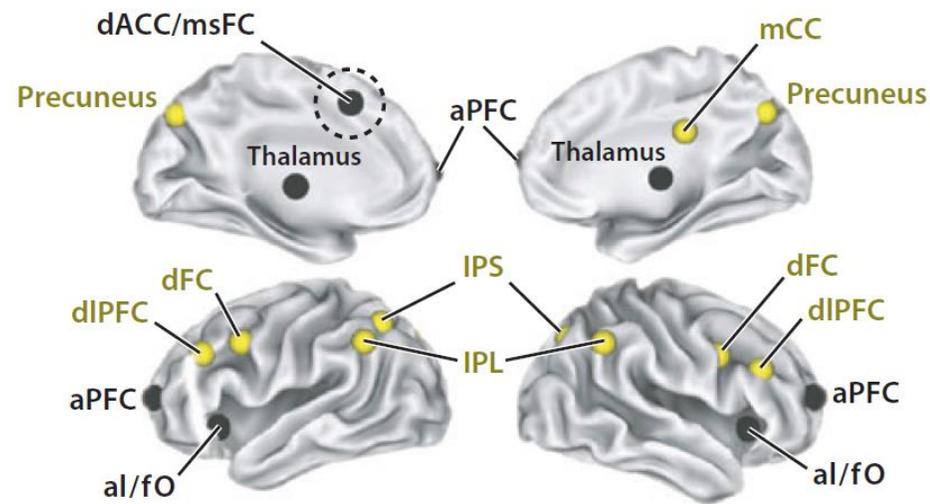


■ Locus coeruleus:
norepinephrine

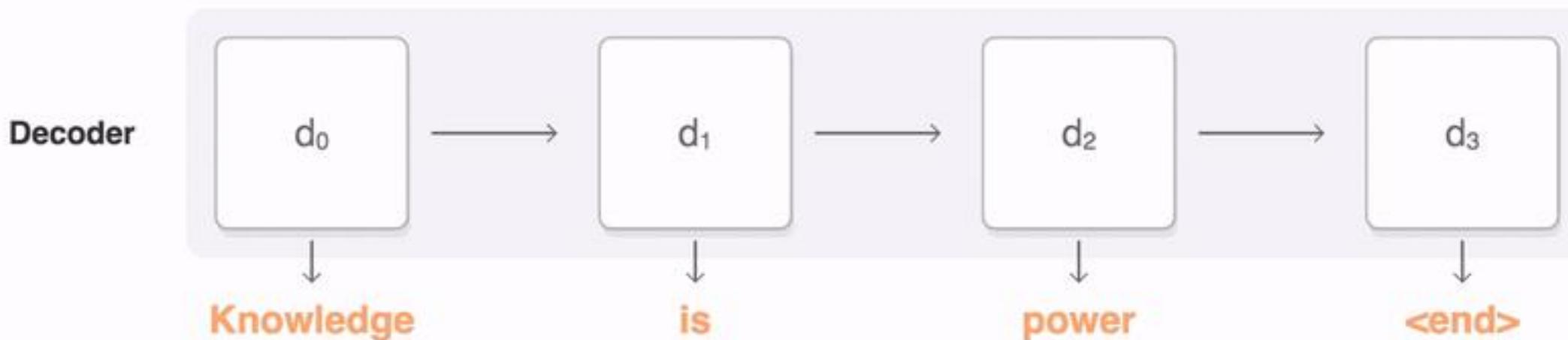
a Orienting



b Executive control

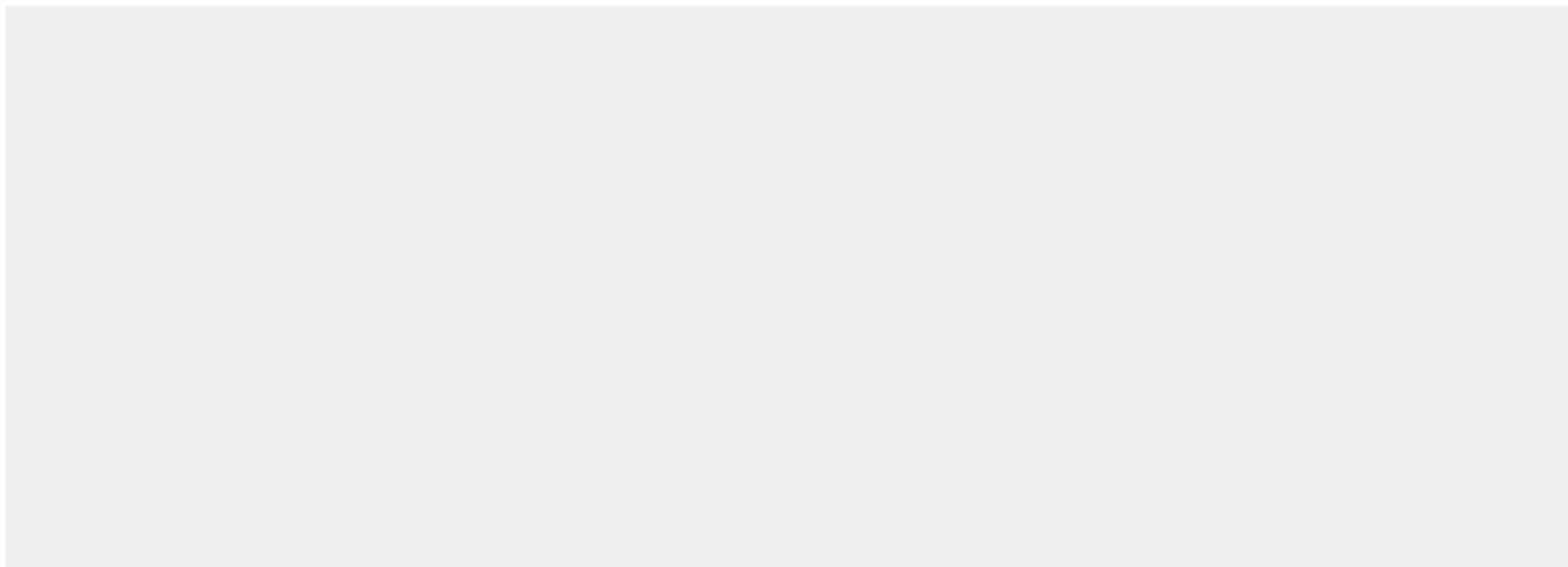


Transformer的注意力机制



注意力机制的计算过程

Self-attention



Input #1

1	0	1	0
---	---	---	---

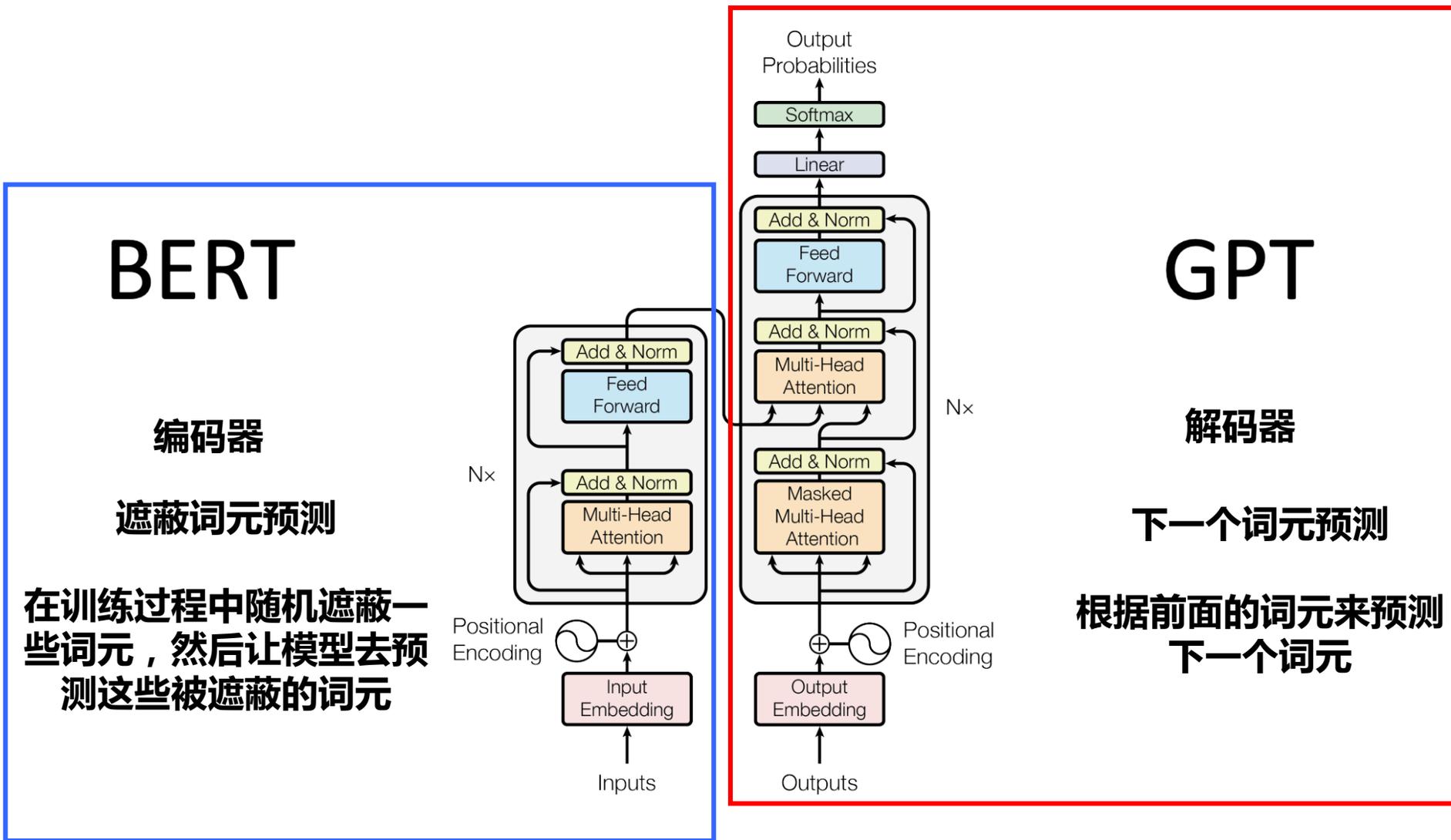
Input #2

0	2	0	2
---	---	---	---

Input #3

1	1	1	1
---	---	---	---

大语言模型的两种主流架构





随堂小测 & 课外学习

• 作业：

- 1. 思考题：在生物大模型微调中，为何需要使用领域特异性数据集？试举例说明？
- 2. 实践题：让DeepSeek以句子如“I love HUST”为例，说明transformer是怎么用自注意力机制来进行计算的

• 参考资料：

- 1. 学术论文：[Empowering biomedical discovery with AI agents](#). Cell, 2024, 187, 6125-6151.
- 2. 学术论文：[The generative era of medical AI](#). Cell, 2025, 18, 3648-3660.