

# Artificial Intelligence Application's in Computational Biology

Eamon McAndrew

Institut de Pharmacologie Moléculaire et Cellulaire, Sophia Antipolis, France

[Mcandrew@ipmc.cnrs.fr](mailto:Mcandrew@ipmc.cnrs.fr)

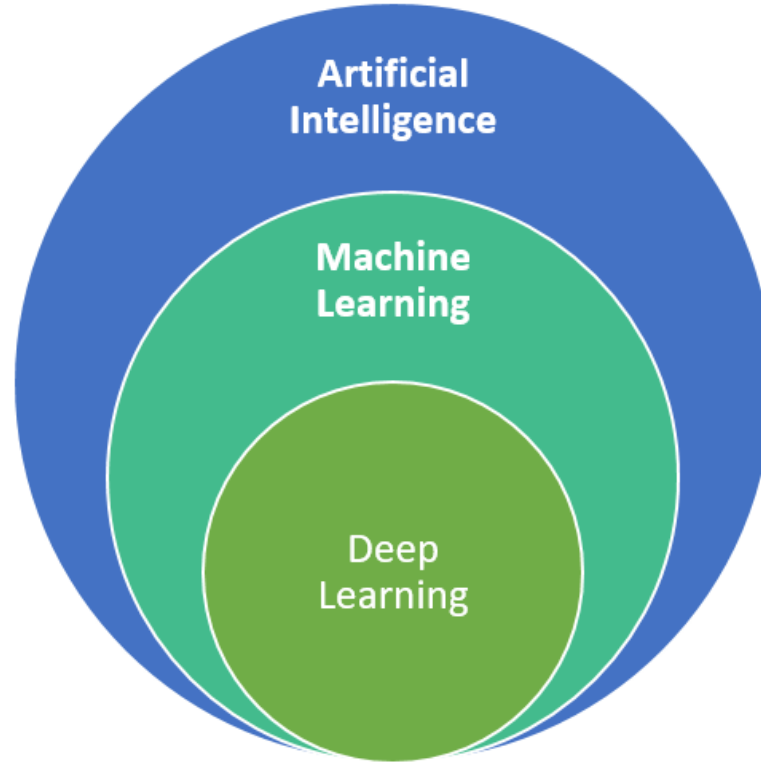
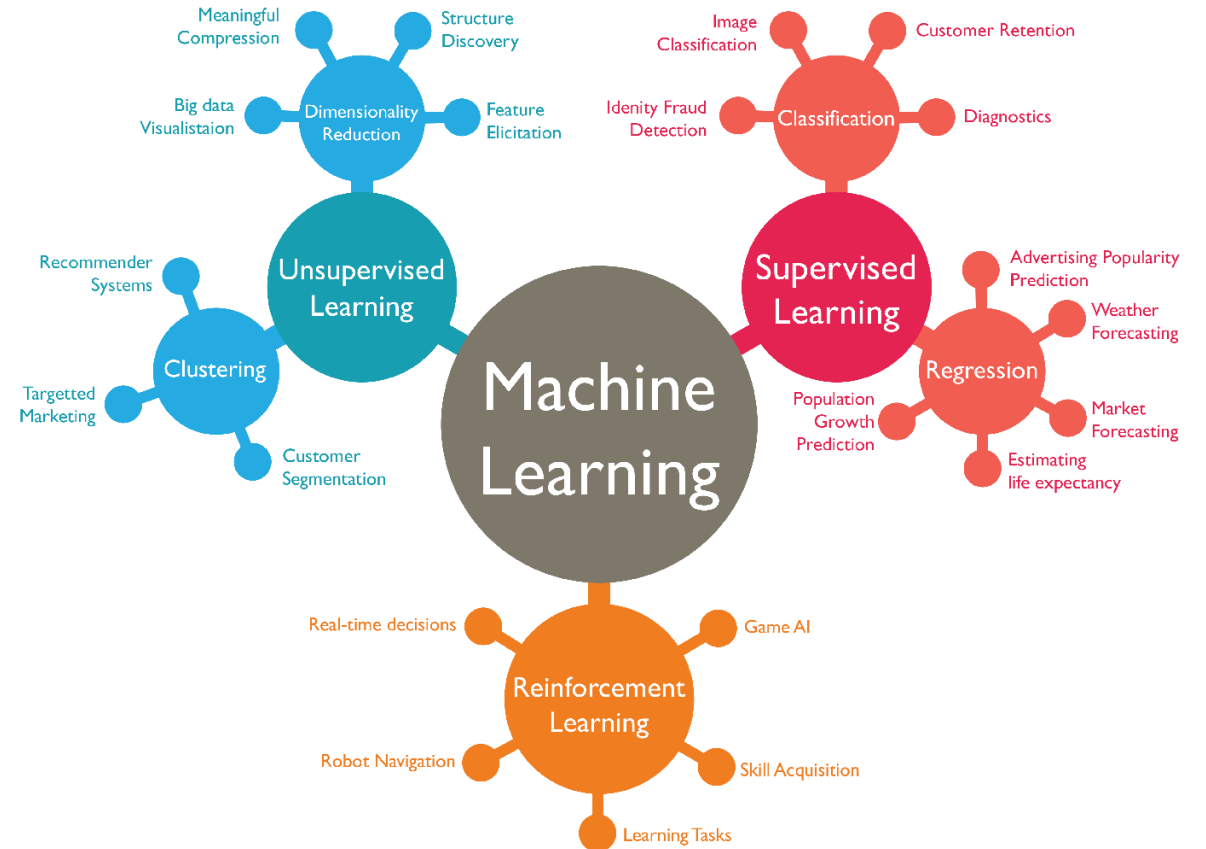
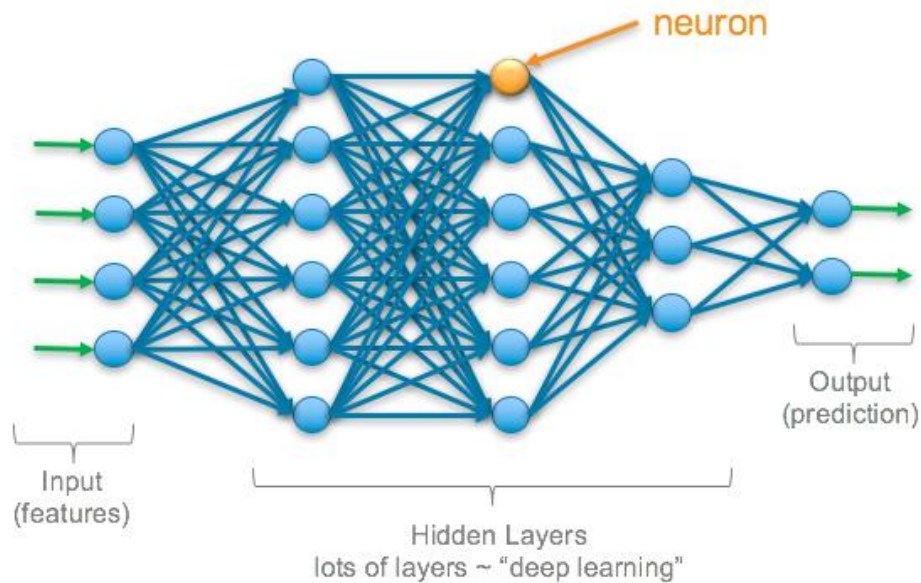


Figure 1: artificial intelligence, machine leaning and deep learning Source: Nadia BERCHANE (M2 IESCI, 2018)

1950's	<b>Artificial intelligence (AI)</b> <i>Human intelligence exhibited by machines</i>
1980's	<b>Machine learning</b> <i>AI systems that learn from historical data</i>
2010's	<b>Deep learning</b> <i>Machine learning models that mimic human brain function</i>
2020's	<b>Generative AI (Gen AI)</b> <i>Deep learning models (foundation models) that create original content</i>



1950's	🧠 Artificial intelligence (AI) <i>Human intelligence exhibited by machines</i>
1980's	📊 Machine learning <i>AI systems that learn from historical data</i>
2010's	🔍 Deep learning <i>Machine learning models that mimic human brain function</i>
2020's	🧠 <sup>AI</sup> Generative AI (Gen AI) <i>Deep learning models (foundation models) that create original content</i>



# 2024 Nobel Prize

2024 Nobel Prize in physics awarded to John J. Hopfield, Geoffrey E. Hinton for discoveries that 'enable machine learning with artificial neural networks'



John J. Hopfield

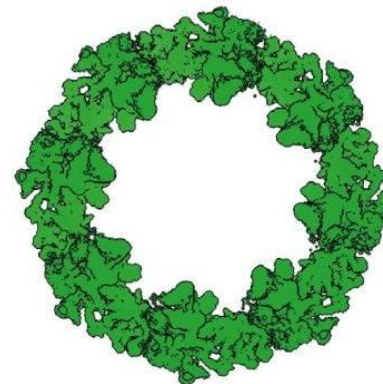
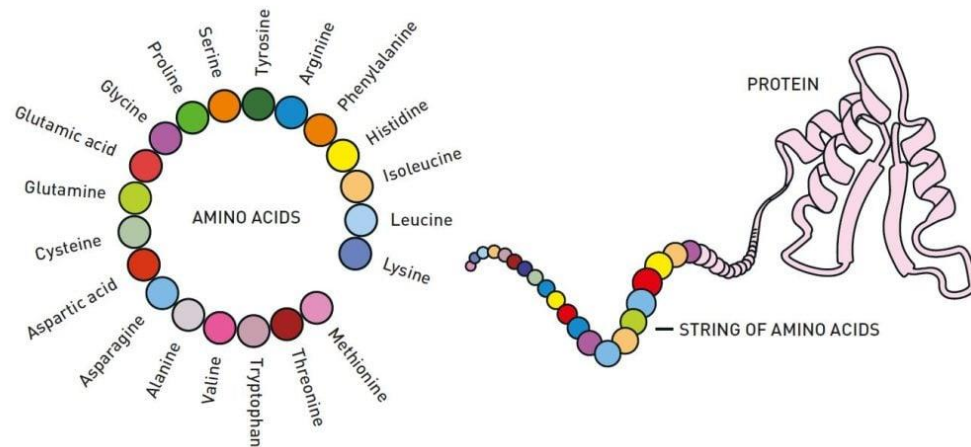
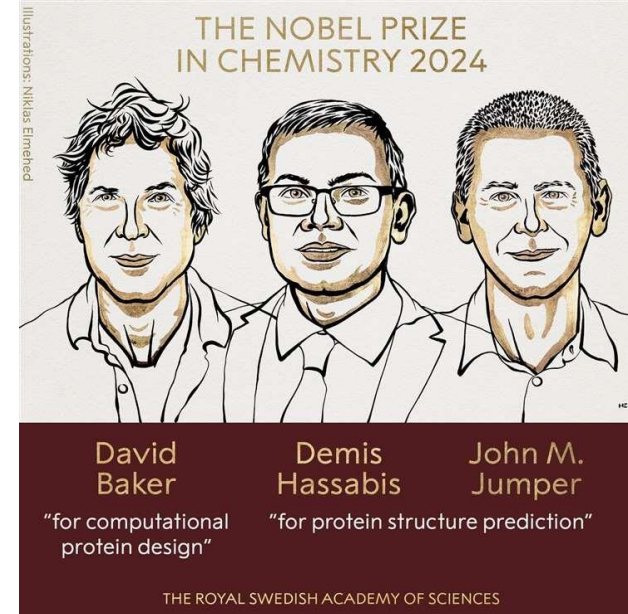
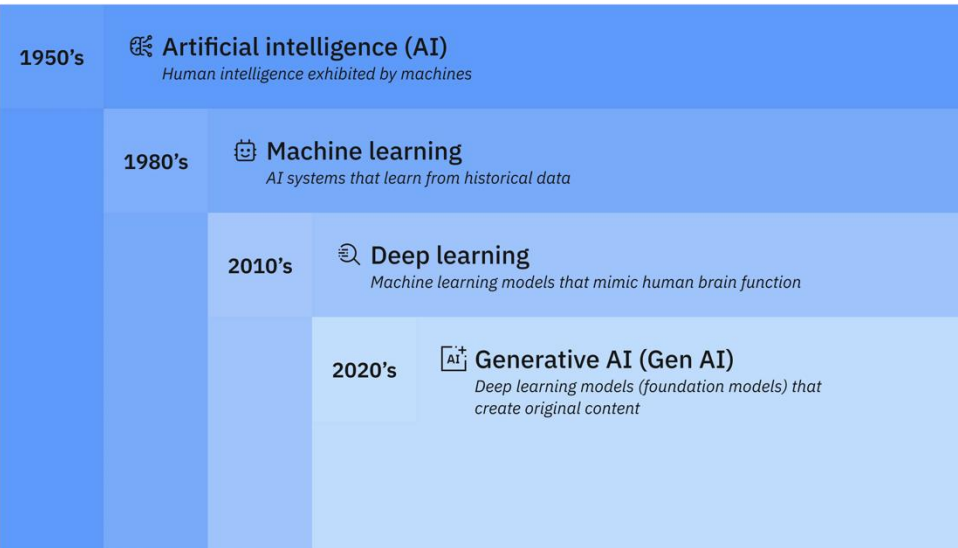


Geoffrey E. Hinton

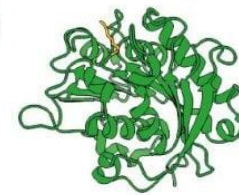


October 8, 2024 Source: Princeton.edu, Toronto.edu

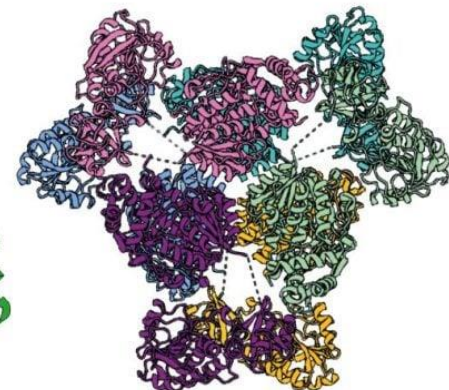




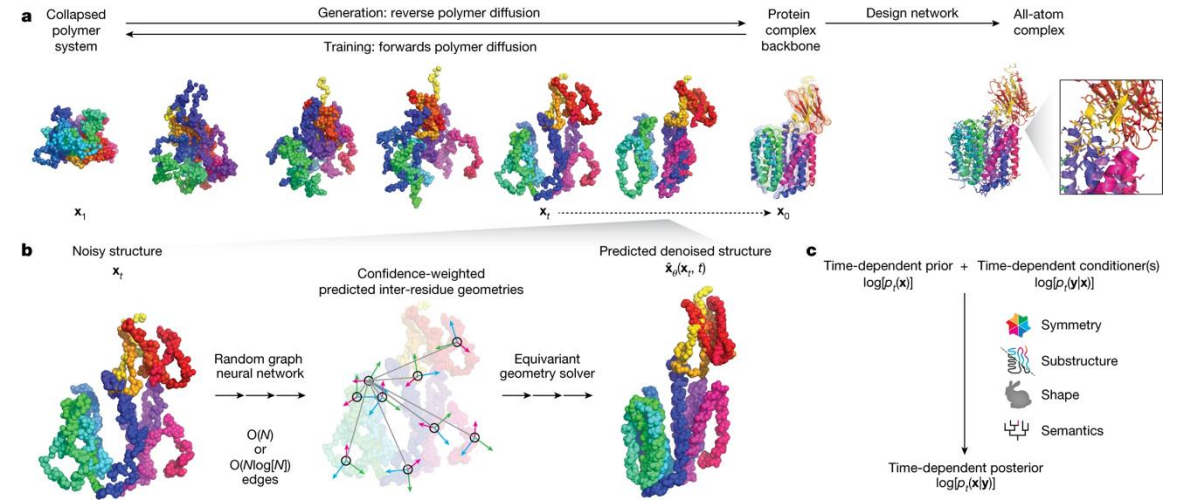
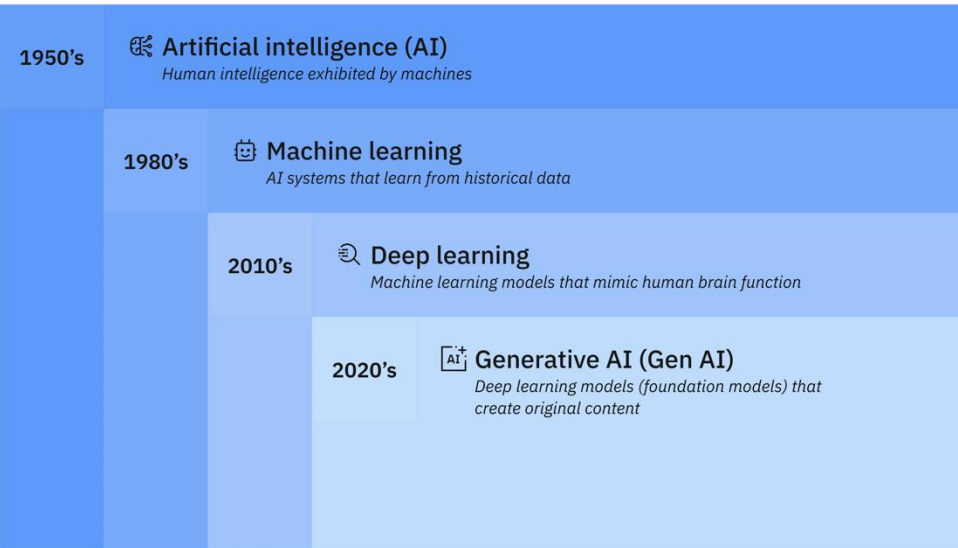
2022: Part of a huge molecular structure in the human body. More than a thousand proteins form a pore through the membrane surrounding the cell nucleus.



2022: Natural enzymes that can decompose plastic. The aim is to design proteins that can be used to recycle plastic.



2023: A bacterial enzyme that causes antibiotic resistance. The structure is important for discovering ways of preventing antibiotic resistance.



ChatGPT

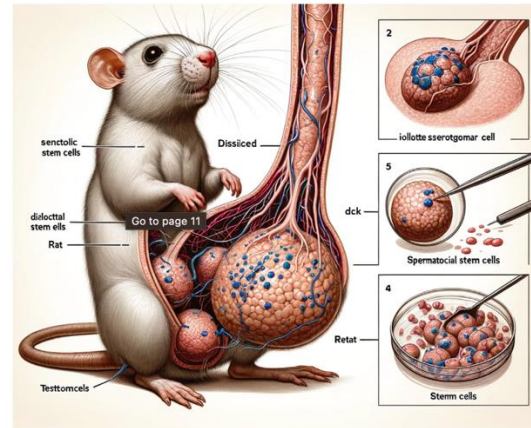
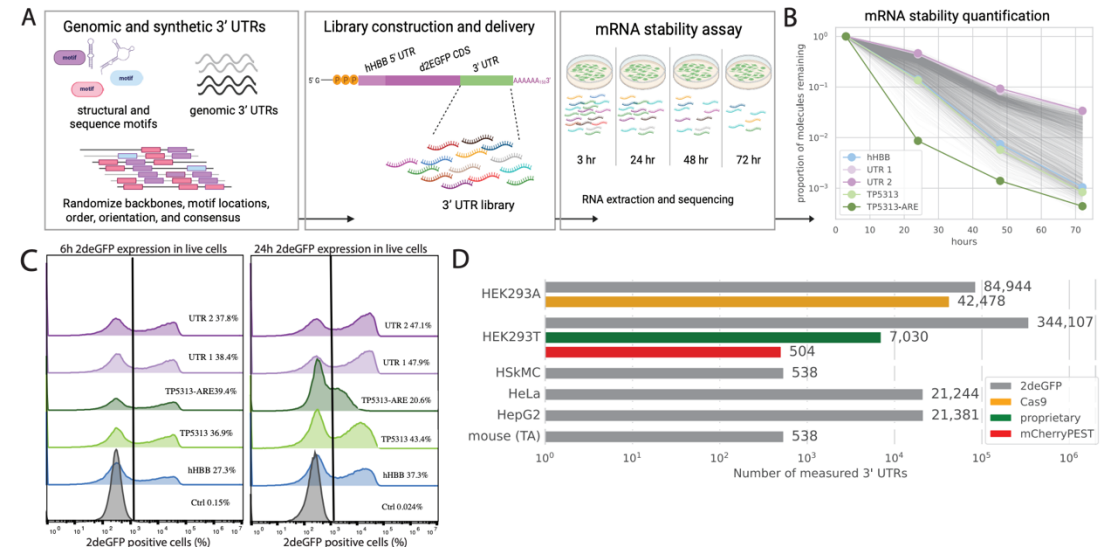


FIGURE 1 Spermatogonial stem cells, isolated, purified and cultured from rat testes.



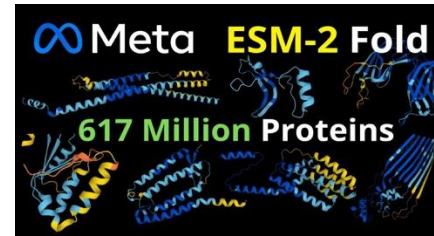
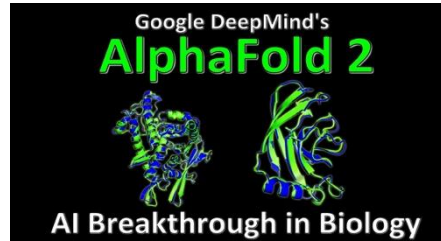
Images



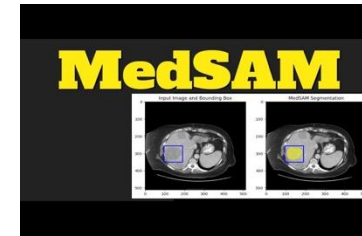
Text / Language



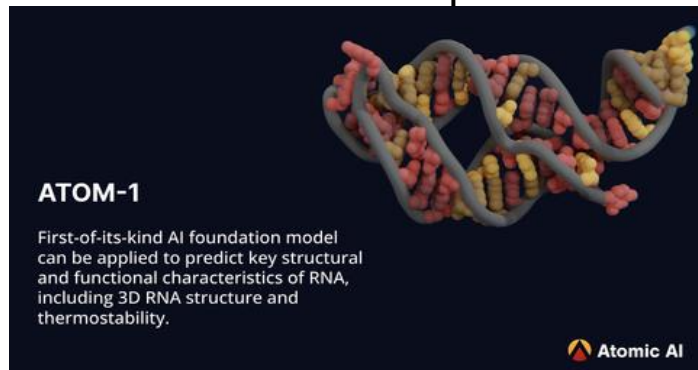
Protein Structure Prediction



Clinical Applications



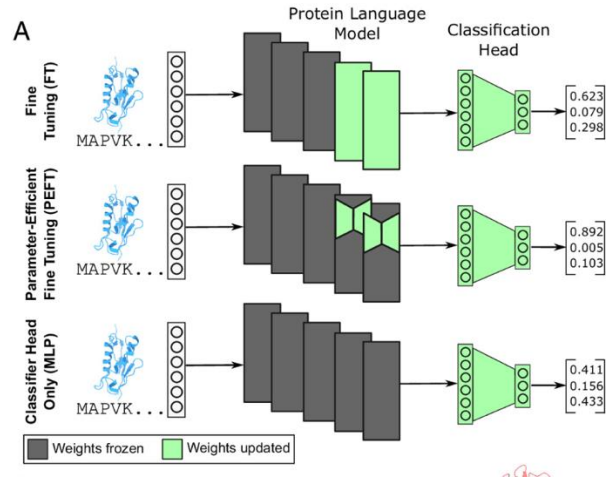
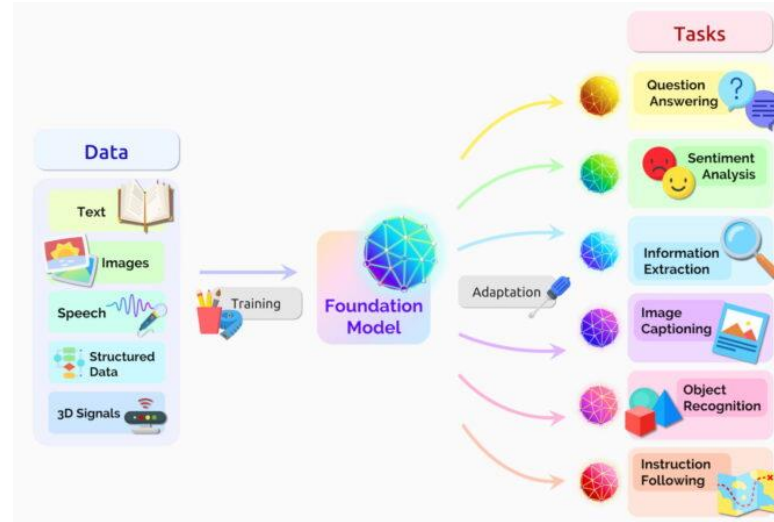
RNA structure prediction



Video and ????



# Foundation Models



## Fine-tuning

ProtT5	6.6 $\pm$ 1.02	7.4 $\pm$ 3.8	2.7 $\pm$ 0.09	4.6 $\pm$ 3.3	3.2 $\pm$ 1.08	3.8 $\pm$ 2.74	0.4 $\pm$ 0.45	0.6 $\pm$ 0.12
ESM2 8M*	4.9 $\pm$ 0.5	15.8 $\pm$ 3.11	6.5 $\pm$ 1.38	1.5 $\pm$ 3.98	2.0 $\pm$ 0.88	3.9 $\pm$ 2.11	1.8 $\pm$ 0.79	0.8 $\pm$ 0.18
ESM2 35M*	3.9 $\pm$ 0.4	21.8 $\pm$ 6.41	5.2 $\pm$ 0.73	7.8 $\pm$ 1.94	1.2 $\pm$ 2.97	3.2 $\pm$ 2.45	4.7 $\pm$ 1.34	0.9 $\pm$ 0.1
ESM2 150M*	5.2 $\pm$ 0.33	18.9 $\pm$ 2	4.6 $\pm$ 1.14	-5.1 $\pm$ 3.11	0.9 $\pm$ 1.94	2.2 $\pm$ 2.08	3.0 $\pm$ 1.58	0.6 $\pm$ 0.16
ESM2 650M*	4.0 $\pm$ 0.51	31.2 $\pm$ 5.04	2.2 $\pm$ 1.31	7.9 $\pm$ 7.56	0.8 $\pm$ 2.03	0.9 $\pm$ 2.57	0.0 $\pm$ 0.78	0.8 $\pm$ 0.08
ESM2 3B	4.9 $\pm$ 0.22	8.1 $\pm$ 1.21	2.7 $\pm$ 0.99	5.2 $\pm$ 1.12	2.2 $\pm$ 1.55	3.6 $\pm$ 1.4	0.6 $\pm$ 0.97	0.8 $\pm$ 0.1
Ankh base	3.2 $\pm$ 0.35	17.6 $\pm$ 5.22	1.8 $\pm$ 1.62	5.3 $\pm$ 6.17	3.5 $\pm$ 1.96	2.2 $\pm$ 1.31	-2.6 $\pm$ 0.75	-0.4 $\pm$ 0.08
Ankh large	2.5 $\pm$ 0.49	11.7 $\pm$ 2.84	3.4 $\pm$ 1.68	11.3 $\pm$ 5.64	-2.1 $\pm$ 4.8	2.2 $\pm$ 2.92	-1.1 $\pm$ 0.8	-0.3 $\pm$ 0.17
	mutational landscape		diverse dataset					
	GFP	AAV	GB1	Stability	Meltome	Subcellular location	Disorder	Secondary structure



## A model Needs -

DATA

- The internet
- PDB database
- Collection of images

- GPU
- TPU
- Nuclear power stations?
- Inference VS Compute

COMPUTE



CODE

- Pytorch / Tensorflow / Jax
  - Architecture
- Training schemes
  - Transformers
    - CNNs
    - GNNs
  - Diffusion

All three scale with model complexity and performance

## Alpha Fold Needs

DATA

- Protein Structure Databases
- NCBI
- Multiple Sequence alignments

**Training**

Many TPUv3 cores.

**Inference**

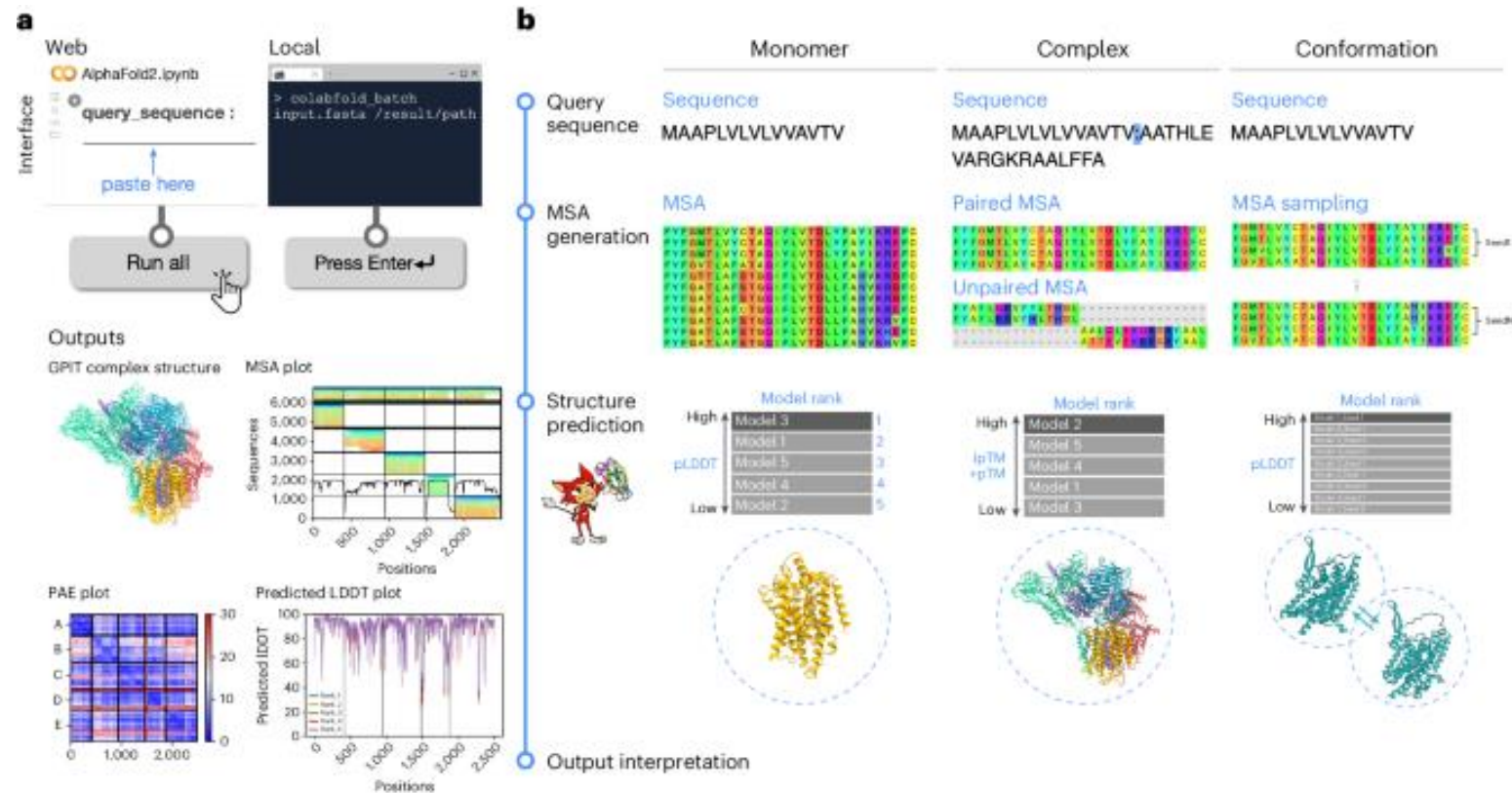
Consumer GPU

COMPUTE

CODE

- Pytorch / Tensorflow / Jax
  - Transformer

Run for free or pay as you need for larger tasks to run the whole AlphaFold workflow in a Jupyter notebook hosted entirely online.



# Chat GPT needs

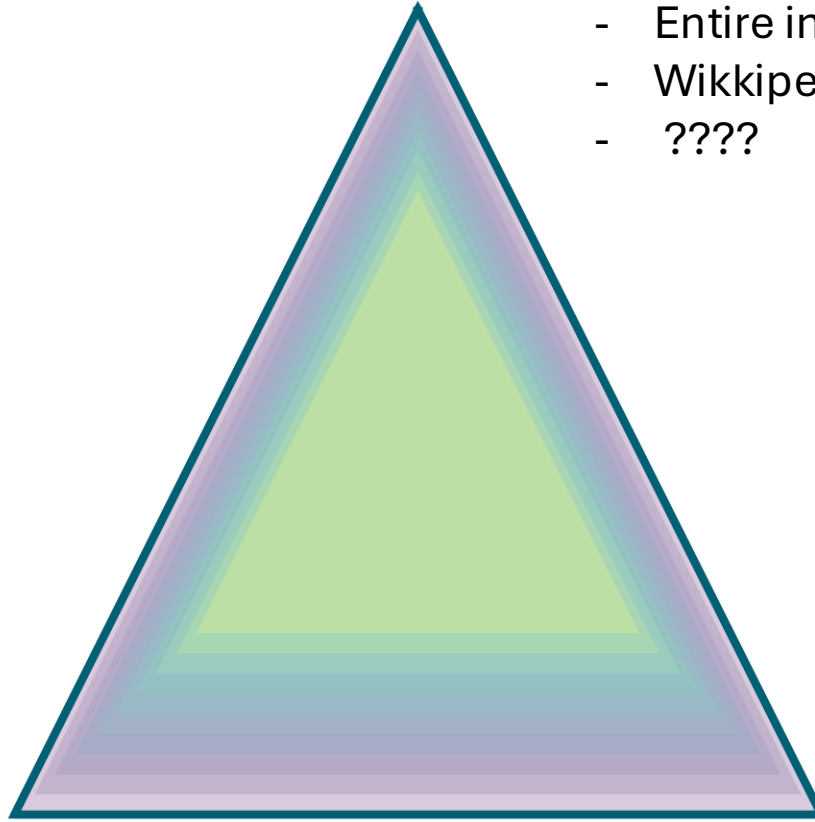
DATA

- Entire internet
- Wikipedia
- ????

**Training**  
Hundreds  
of  
thousands  
of GPUs

**Inference**  
Dozens of GPUs

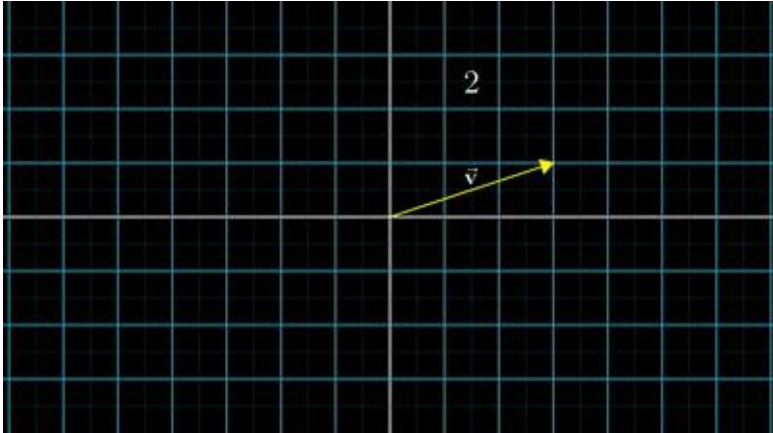
COMPUTE



CODE



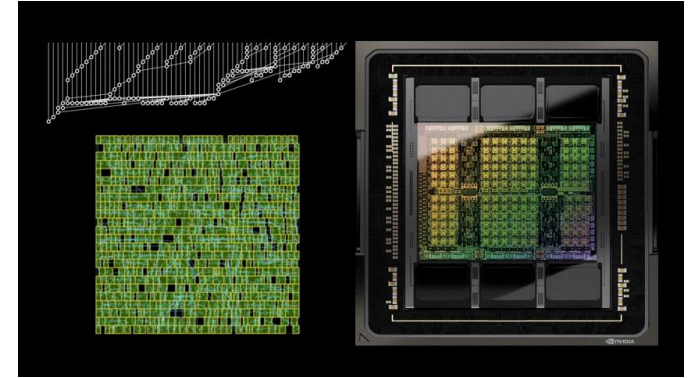
- Pytorch / Tensorflow / Jax
- Transformers



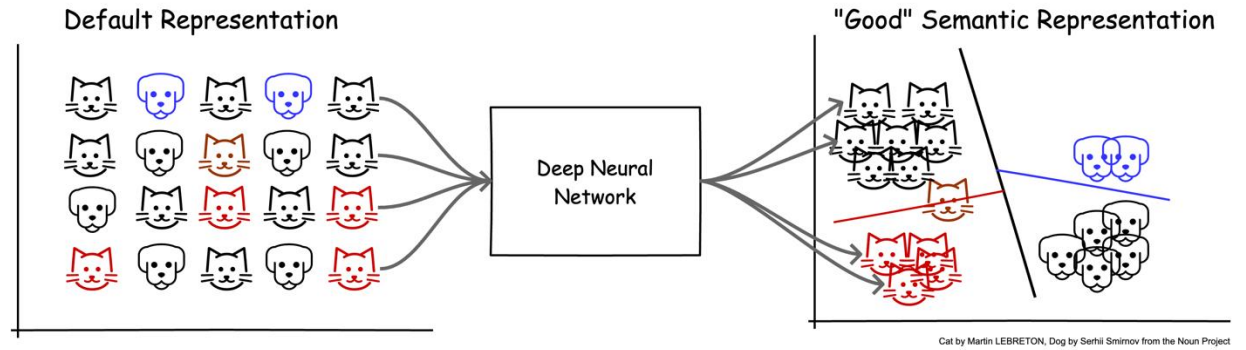
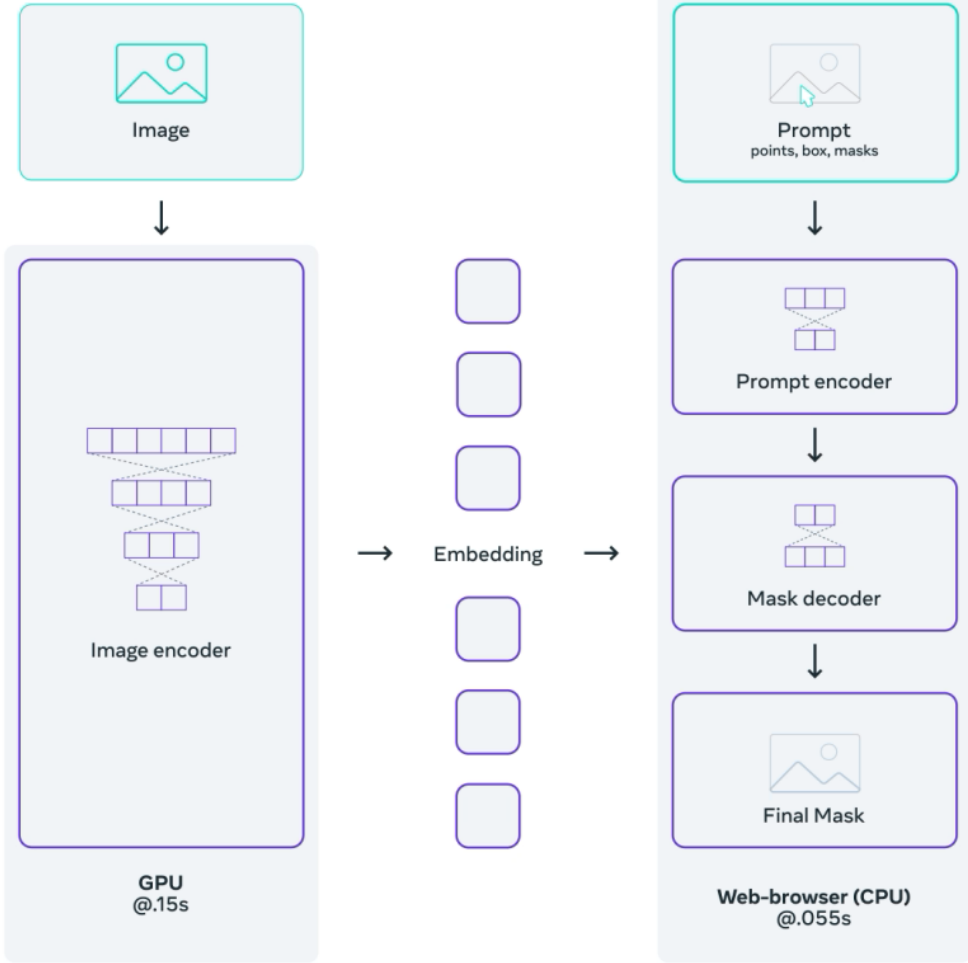
### Matrix Multiplication

$$\begin{matrix} - \\ + \end{matrix} \begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 0 \\ 2 & 3 & 4 \end{bmatrix} \times \begin{bmatrix} 2 & 5 \\ 6 & 7 \\ 1 & 8 \end{bmatrix} \begin{matrix} - \\ + \end{matrix}$$

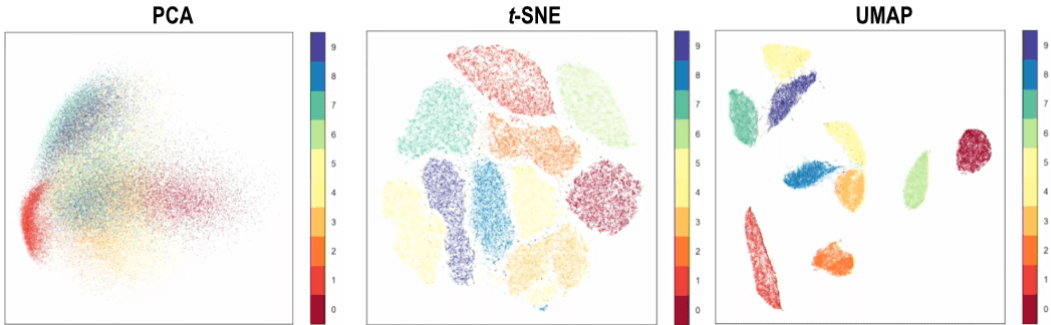
- +
- +



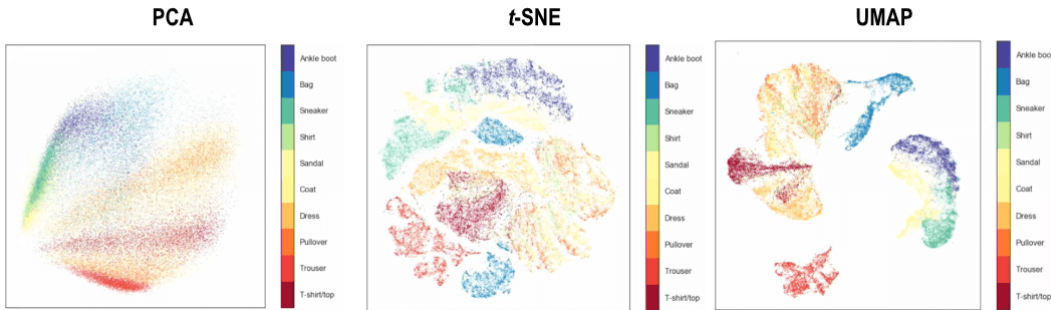
CPU	GPU
Central Processing Unit	Graphics Processing Unit
Several cores	Many cores
Low latency	High throughput
Good for serial processing	Good for parallel processing
Can do a handful of operations at once	Can do thousands of operations at once



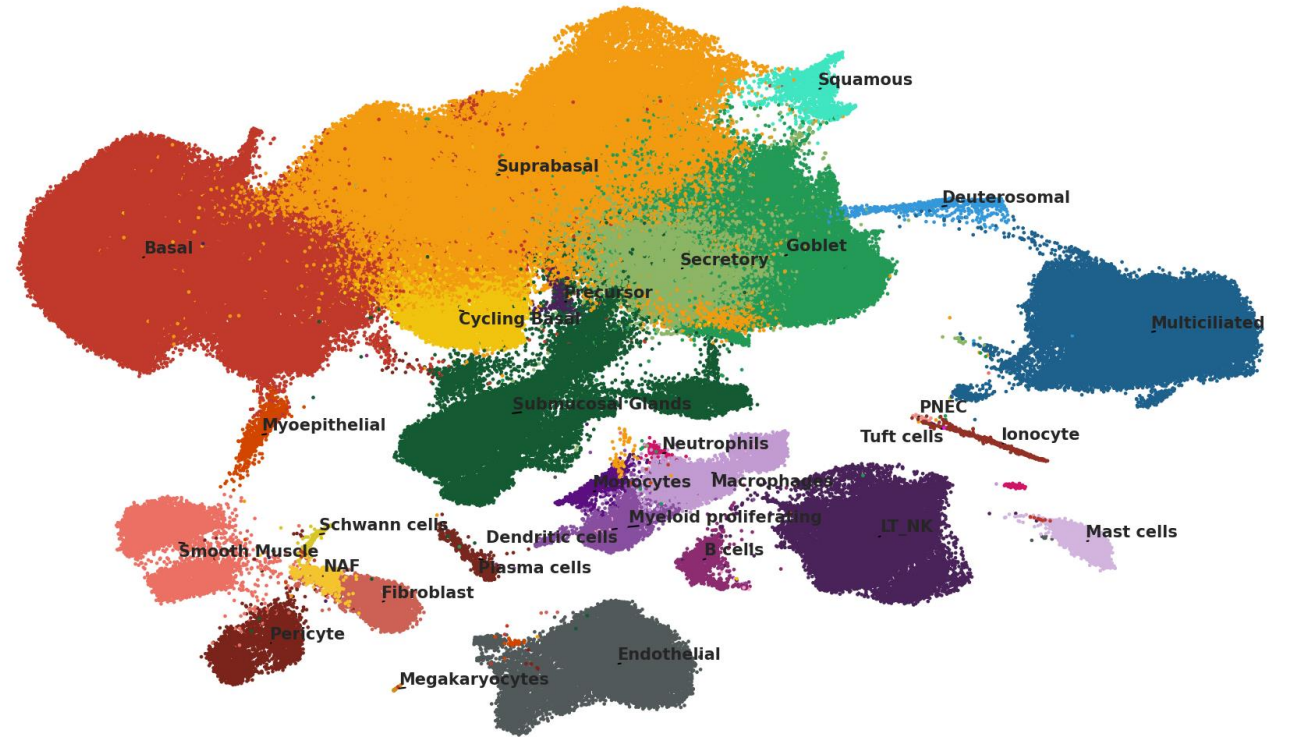
MNIST Digits

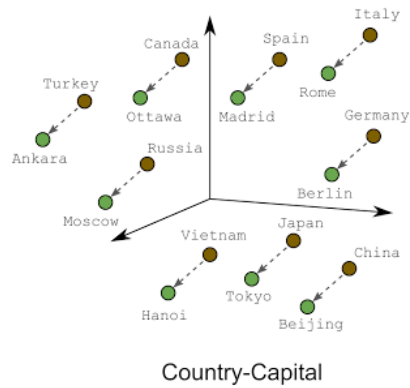
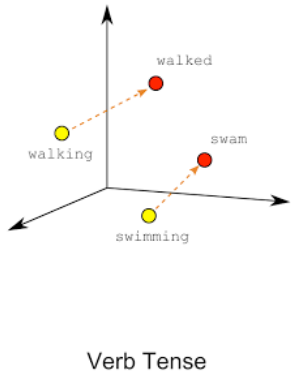
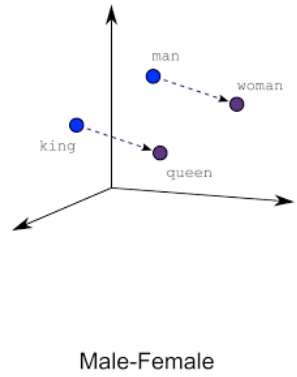


Fashion MNIST



celltype\_iv1\_V6

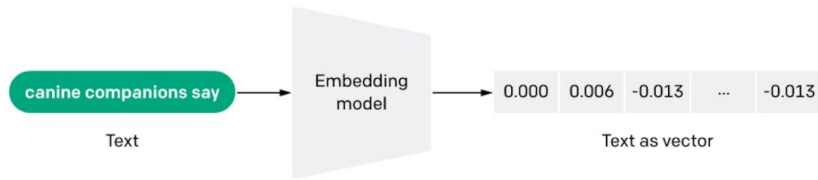




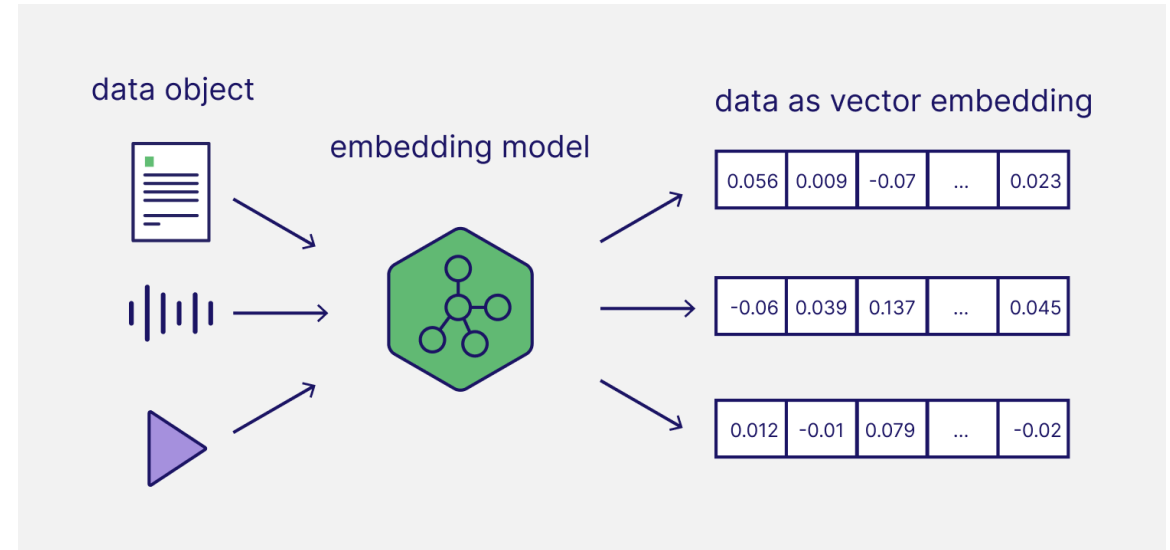


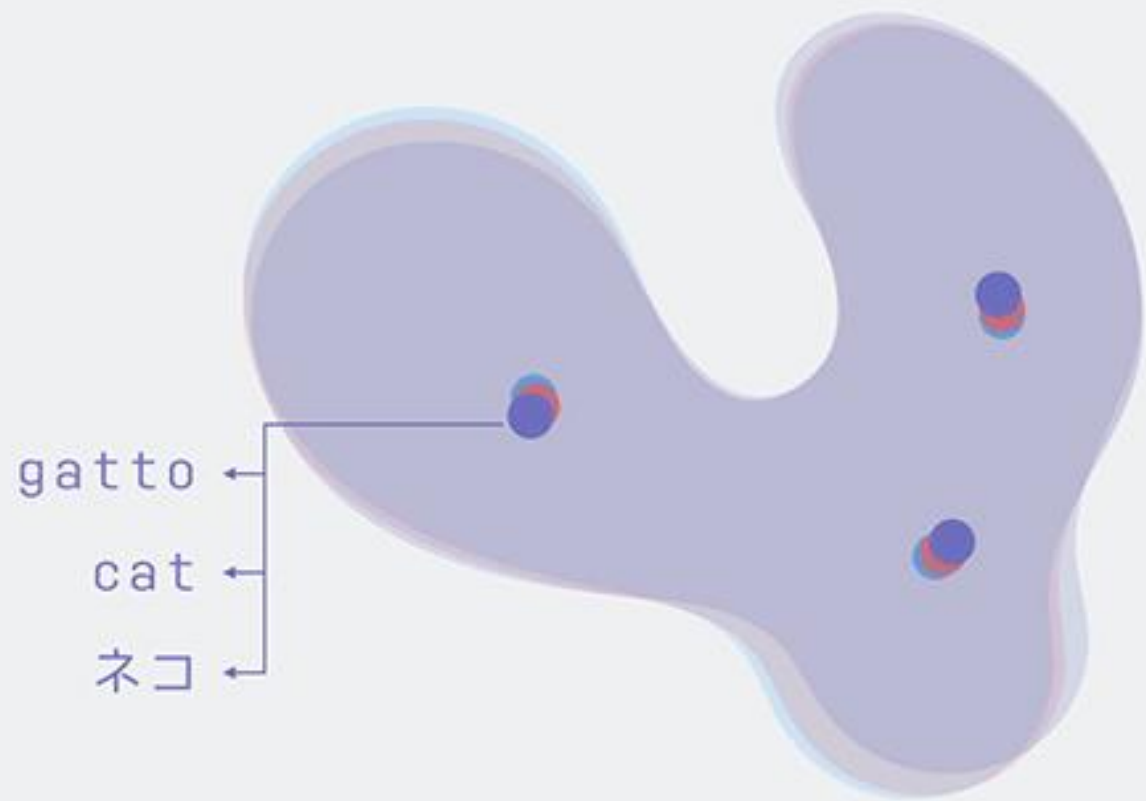
# OpenAI Embedding

## Embeddings Use Cases



Text Similarity, Semantic Search, Clustering

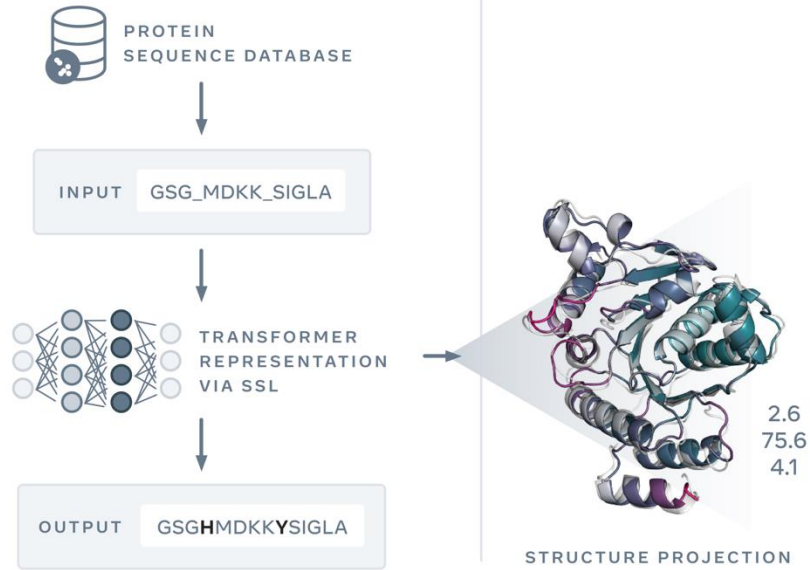




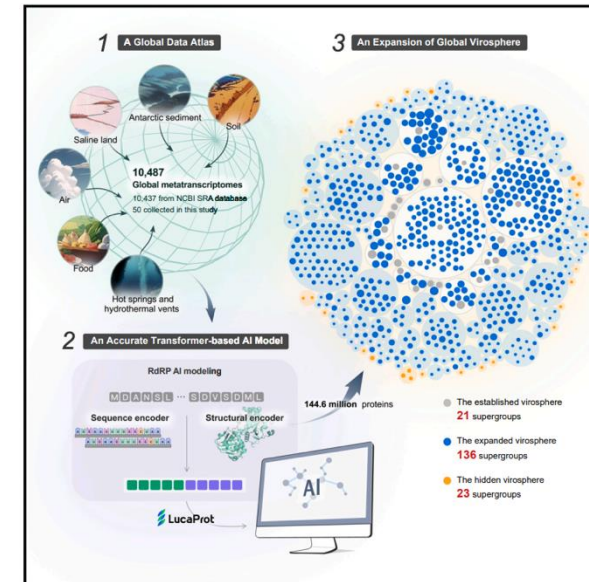
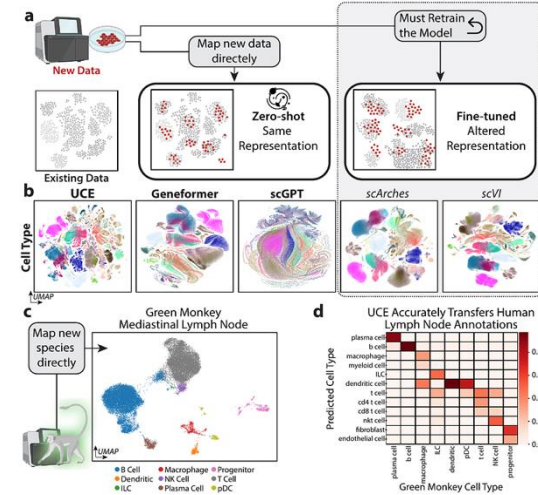
Protein language modeling

Pretraining self-supervision on sequences only.

Structure emerges in the internal representations of the network from the self-supervision.



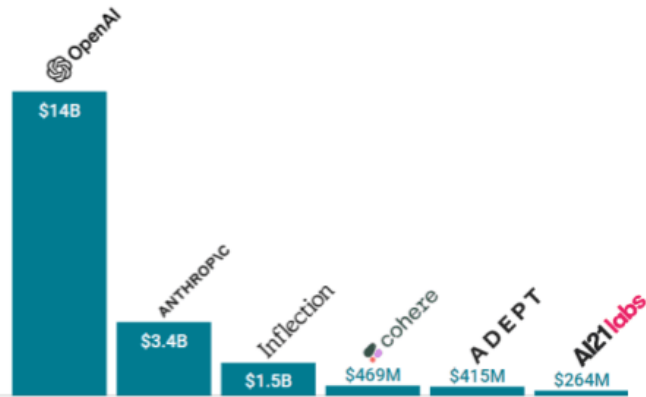
(Rosen et al., 2023)  
(Hou et al., 2024)



# The private market is split into open vs. closed

Disclosed equity funding to LLM developers (as of 10/27/2023)

## Closed-source LLMs



\*Some developers may offer open-source versions of their models but keep their core models proprietary

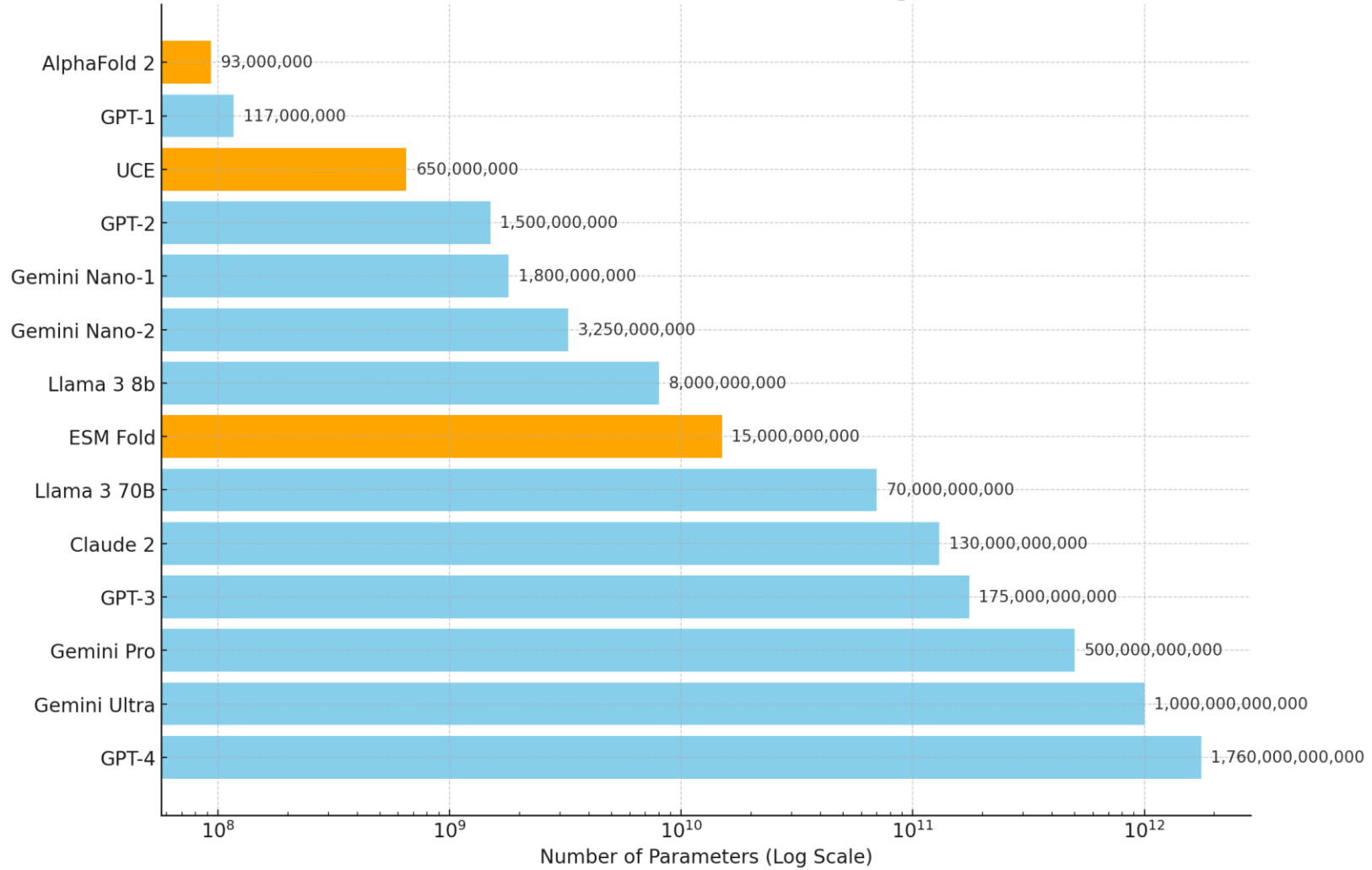
## Open-source LLMs

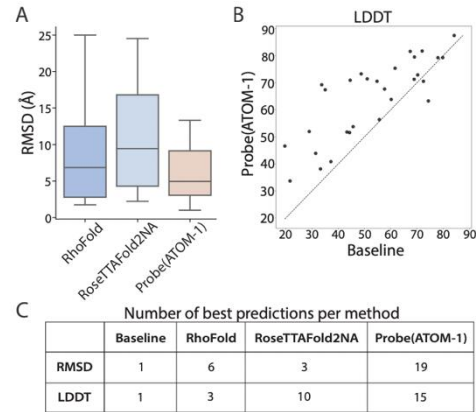
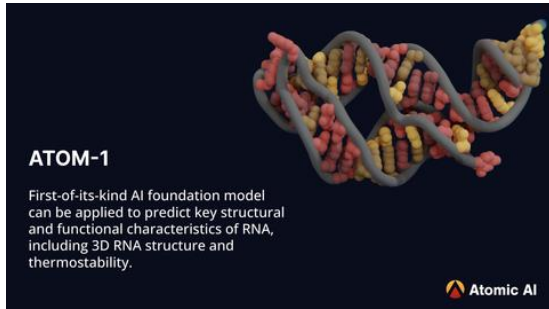


■ Private company ■ Acquired company

\*Excludes open-source developers that have not raised equity funding

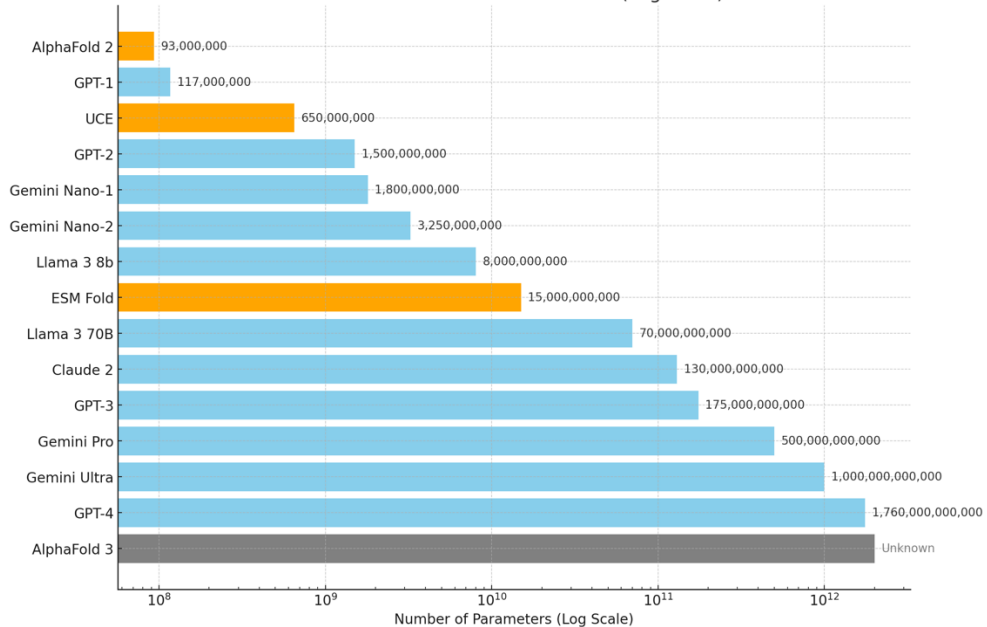
Parameters in Selected AI Models (Log Scale)





AI for computational biology is beginning to look similar


Parameters in Selected AI Models (Log Scale)



(Boyd et al., 2023)  
(Abramson et al., 2024)

Article | [Open access](#) | Published: 08 May 2024

## Accurate structure prediction of biomolecular interactions with AlphaFold 3

[Josh Abramson](#), [Jonas Adler](#), [Jack Dunger](#), [Richard Evans](#), [Tim Green](#), [Alexander Pritzel](#), [Olaf Ronneberger](#), [Lindsay Willmore](#), [Andrew J. Ballard](#), [Joshua Bambrick](#), [Sebastian W. Bodenstein](#), [David A. Evans](#), [Chia-Chun Hung](#), [Michael O'Neill](#), [David Reiman](#), [Kathryn Tunyasuvunakool](#), [Zachary Wu](#), [Akvilė Žemgulytė](#), [Eirini Arvaniti](#), [Charles Beattie](#), [Ottavia Bertolli](#), [Alex Bridgland](#), [Alexey Cherepanov](#), [Miles Congreve](#), ... [John M. Jumper](#)  [+ Show authors](#)

AlphaFold Server BETA

But his enthusiasm was not universally shared. In contrast to the launch of AlphaFold 2 in 2021, *Nature's* publication of AlphaFold 3 lacked the open source code. That omission has sparked outcry from the research community, culminating in a protest letter signed by more than 1,000 scientists.

## Why AlphaFold 3 needs to be open source



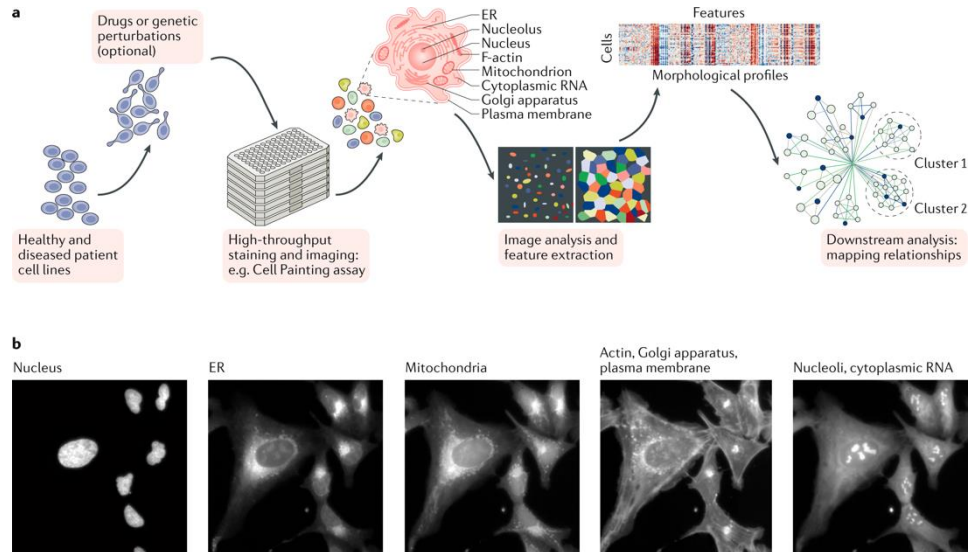
By Bryce Johnson  
July 7, 2024

## AlphaFold3 – why did *Nature* publish it without its code?

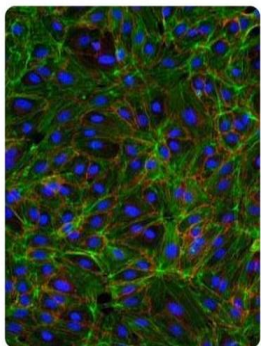
Criticism of our decision to publish AlphaFold3 raises important questions. We welcome readers' views.



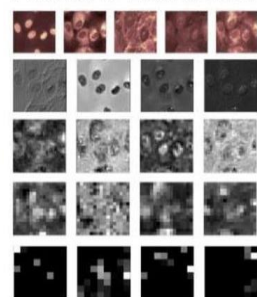
## Generate Large ML Centric Datasets



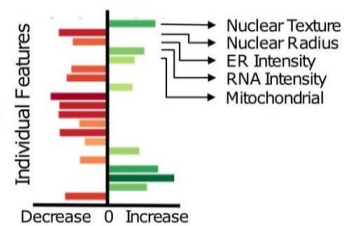
**2 Million**  
images  
each week



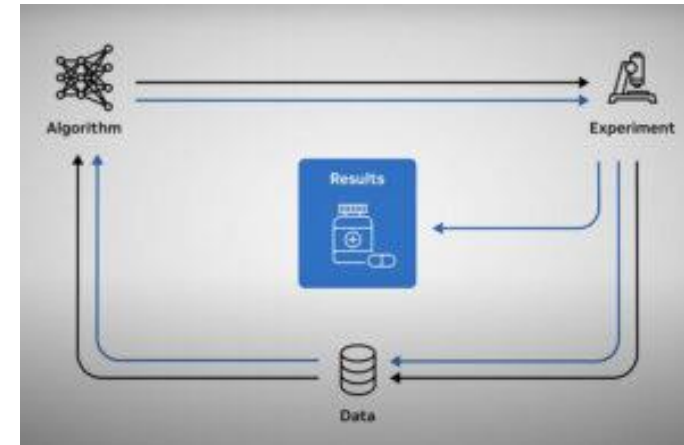
**Deep Learning & Feature Extraction**



**Disease Phenotypic Fingerprints**



## Lab In a Loop



Iterate between experimental data generation and model training



Thorough introduction centered on computational biology, course is ongoing and posted every week so its up to date.

**Introductions**

- **Lecturer: Manolis Kellis**  
– MIT CSAIL, CompBio, Broad, Disease mechanism, Epigenomics, Cancer, Brain, Gene Regulation, Evolution, Single-cell genomics
- **Lecturer: Eric Alm**  
– MIT Biological Engineering, Gen AI, Computational, theoretical, experimental understanding & engineering human microbiome
- **TA: Jared Zheng**  
– MIT CSAIL, Zhang Lab, Chemistry, Biophysics, protein-ligand interactions, drug discovery, deep generative models, PLMs
- **TA: Sarah Gurev**  
– MIT EECS, Debbie Marks Lab Harvard, Stanford BS in CS, protein design and evolution
- **TA: Benjamin James**  
– MIT EECS, CSAIL, Computational Biology, Regulatory Circuitry, Single-Cell, Addiction, Neuroscience

⇒ 13 videos

**MLCB24 - Machine Learning in Computational Biology Fall 2024**

Updated 2 days ago

[View full playlist](#)

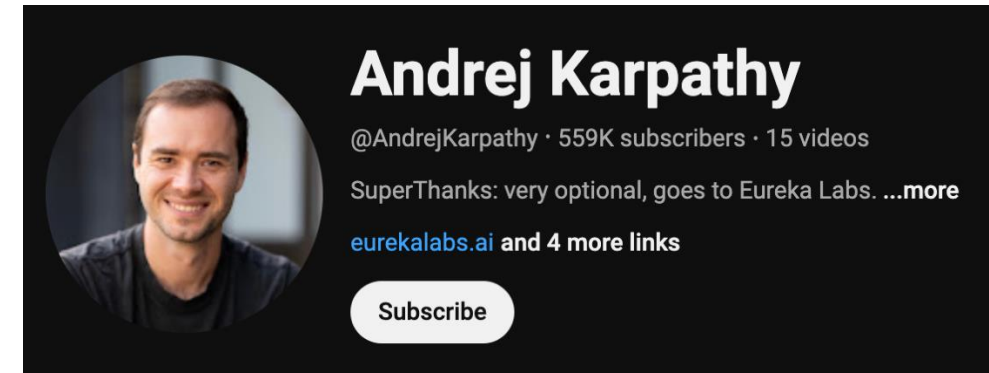
Full course on deep learning from setup to training many popular architectures aimed at people with some coding and no deep learning experience

**Practical Deep Learning for Coders**

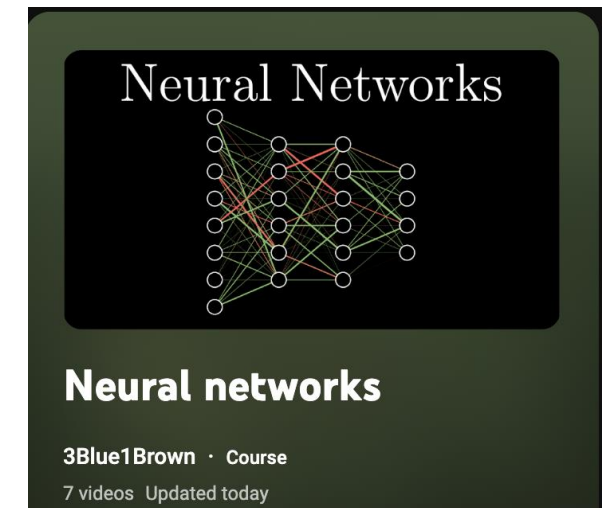
by Jeremy Howard

Some of the best explanations and teaching available – ranging in difficulty from public audiences to expert technical deep dives.

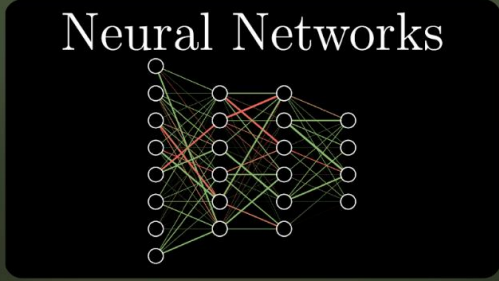
Extremely detailed explanations of the mathematical concepts underlying neural network, made much more approachable through beautiful visualizations



**Andrej Karpathy**  
@AndrejKarpathy · 559K subscribers · 15 videos  
SuperThanks: very optional, goes to Eureka Labs. ...more  
[eurekalabs.ai](https://eurekalabs.ai) and 4 more links  
Subscribe



Neural Networks



**Neural networks**  
3Blue1Brown · Course  
7 videos Updated today