# Tool assisted research in the age of AI

**Large Language Models and more**

## Renato Alves

**Senior Bioinformatics Community Manager**
**Bio-IT project / Data Science - EMBL Heidelberg**

MICROB-PREDICT

EMBL
50 YEARS | 1974–2024

# Open in the browser or grab your mobile

Go to

www.menti.com

Enter the code

6413 1678

Or use QR code

MICROB-PREDICT

EMBL
50 YEARS|1974–2024

**"** Any sufficiently advanced technology is indistinguishable from magic.

**– Arthur C. Clarke** **"**

MICROB-PREDICT

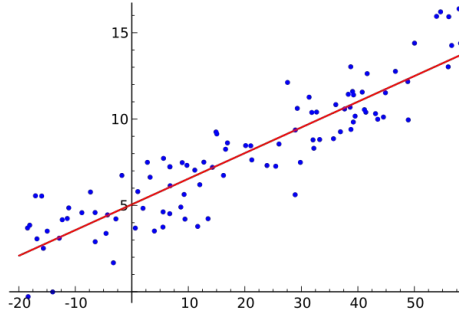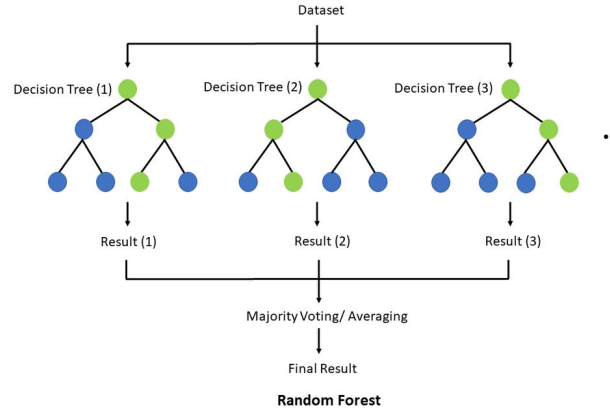EMBL 50 YEARS | 1974–2024

# What is AI?

# What is *intelligence*

**Intelligence**: the ability to learn, understand, and make judgments or have opinions that are based on reason
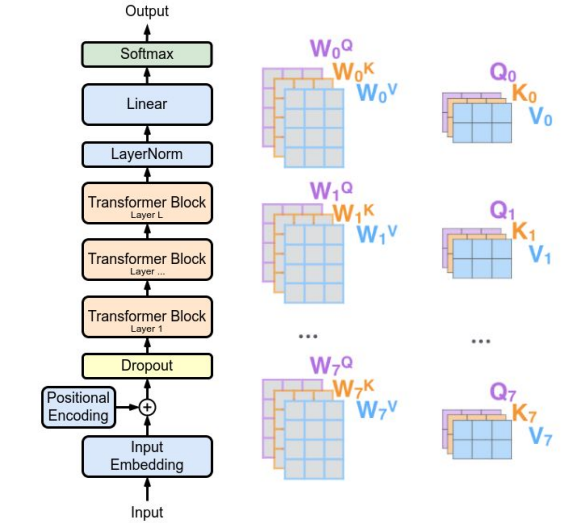
– Cambridge Dictionary

Images from Wikipedia/Wikimedia and Modular (Illustrated transformer)

# What is artificial *intelligence* (AI)



**Stats**
Linear regression

**Machine Learning**
Random forest

**Deep Learning**
Transformers (the T in GPT)

MICROB-PREDICT    EMBL 50 YEARS|1974–2024

# Domains of AI …

Output

Softmax

Linear

LayerNorm

Transformer Block
Layer L

Transformer Block
Layer …

Transformer Block
Layer 1

Dropout

Positional Encoding ⊕

Input Embedding

Input

**Natural Language Processing (NLP)
(e.g. ChatGPT)**

The quick brown fox jumps over the lazy dog

**Sentiment Analysis
Machine Translation
Text Summarization
Text Clustering
Question Answering**

**Computer Vision
(CV)
(e.g. GPT4-Vision)**

**Object Detection
Object Recognition
Semantic Segmentation
Feature Detection and Description
Scene Understanding
Video Analysis**

🤖: The phrase "The quick brown fox jumps over the lazy dog" is an English-language pangram, commonly used for touch-typing practice, testing typewriters and computer keyboards, displaying font examples, and other text-related applications where using all alphabet letters is necessary.

Images from Wikipedia/Wikimedia

MICROB-PREDICT

EMBL
50 YEARS|1974–2024

# To AI or not to AI

# Recent AI evolution … very very fast

**GPT-1**
117 Million Parameters

**GPT-2**
1.5 Billion Parameters

**GPT-3**
175 Billion Parameters

**GPT-3.5**
175 Billion Parameters

**GPT-4**
1.76 Trillion Parameters

**GPT-4.5?**

**2018** — **2019** — **2020** — **2022** — **2023** — **2024**

ChatGPT talk?
GPT-4 Vision
DALL-E 3

Stable Video Diffusion

Google replaces Bard with Gemini

Copyright lawsuits
Midjourney V6

Sora

Figure 01
Devin

**Sep 2023** — **Nov 2023** — **Dec 2023** — **Jan 2024** — **Feb 2024** — **Mar 2024**

MICROB-PREDICT

EMBL
50 YEARS | 1974–2024

# AI in practice

# Some key jargon

**Tokens** → cost

**Embeddings** → semantic/meaning



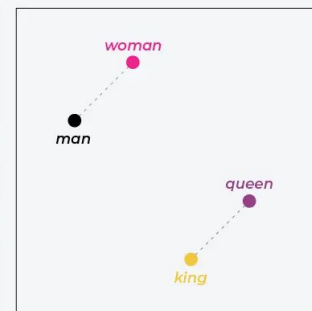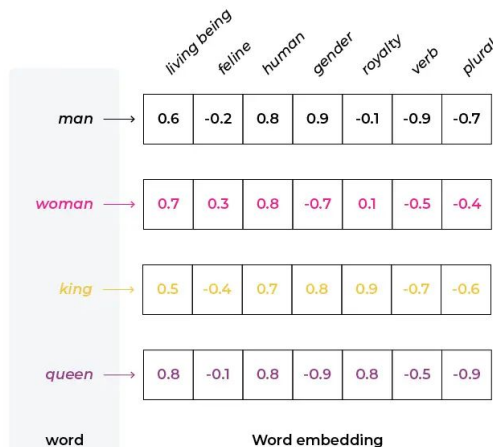Here's one example of a longish and impressively mesmerizing sentence.

Clear    Show example

**Tokens**    **Characters**
15          70

Here's one example of a longish and impressively mesmerizing sentence.

| Model | Input | Output |
|---|---|---|
| gpt-4 | $30.00 / 1M tokens | $60.00 / 1M tokens |
| gpt-4-32k | $60.00 / 1M tokens | $120.00 / 1M tokens |

| word | living being | feline | human | gender | royalty | verb | plural |
|---|---|---|---|---|---|---|---|
| man | 0.6 | -0.2 | 0.8 | 0.9 | -0.1 | -0.9 | -0.7 |
| woman | 0.7 | 0.3 | 0.8 | -0.7 | 0.1 | -0.5 | -0.4 |
| king | 0.5 | -0.4 | 0.7 | 0.8 | 0.9 | -0.7 | -0.6 |
| queen | 0.8 | -0.1 | 0.8 | -0.9 | 0.8 | -0.5 | -0.9 |

word          Word embedding          Visualization of word embedding

Tokenizer by OpenAI - Embeddings from Arize
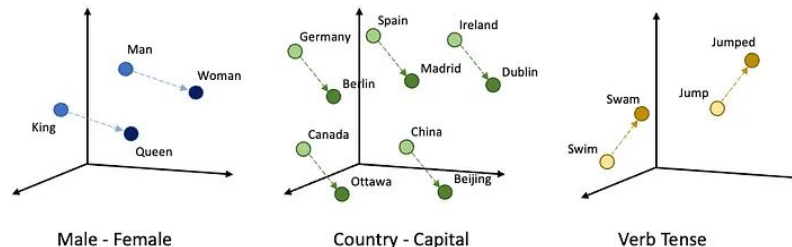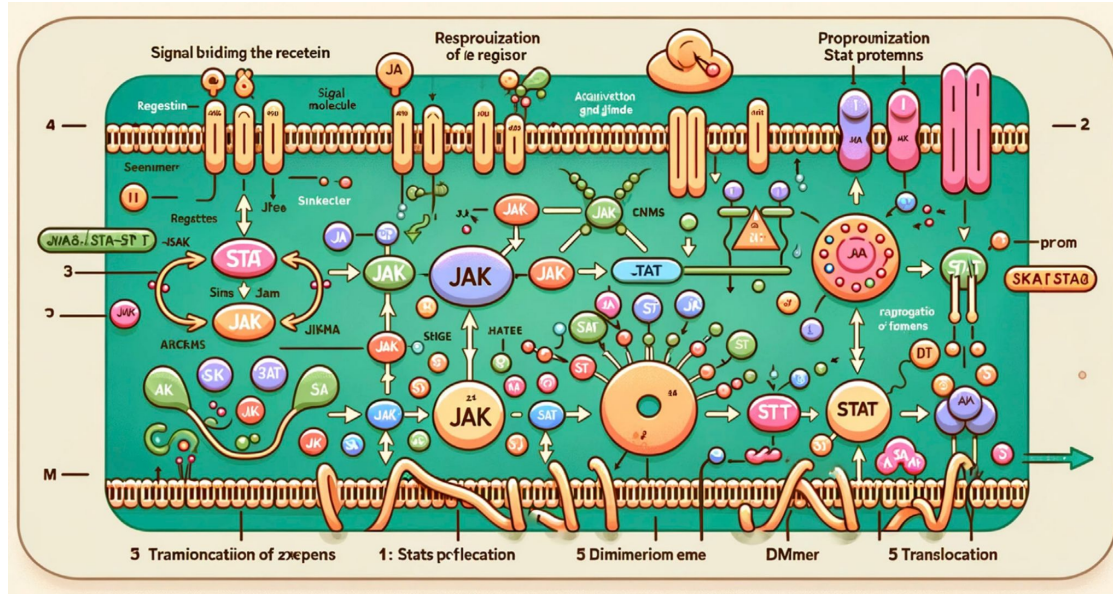
# Polite language generates better results



AI remembers!!



**The role of politeness in human–machine interactions: a systematic literature review and future perspectives**

# "With great power comes great responsibility"



"As of my last knowledge update …"

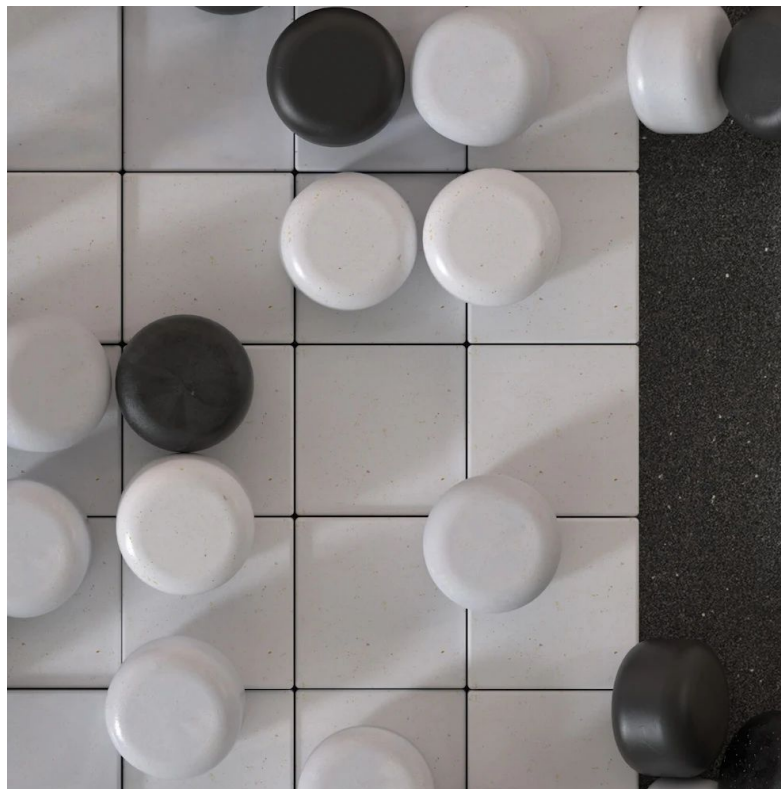Proper use is hard to impossible to distinguish from no use.

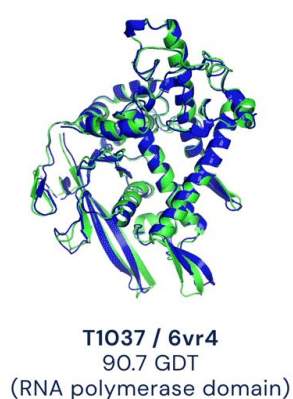You are accountable for using these tools responsibly.

Beware data privacy!
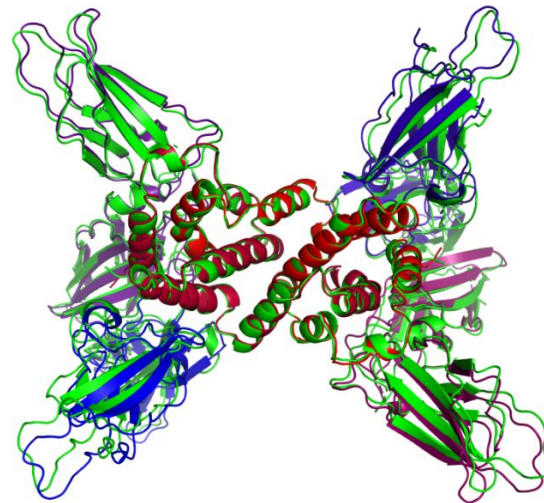
# Earlier AI advances

# AlphaGO

- Simple rules
- Complex play
- Harder than chess
- Adversarial

deepmind.google/technologies/alphago

# Alphafold / Alphafold multimer / Alphafold DB



T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

T1049 / 6y4f
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction

(a) A2B2C2 heteromer
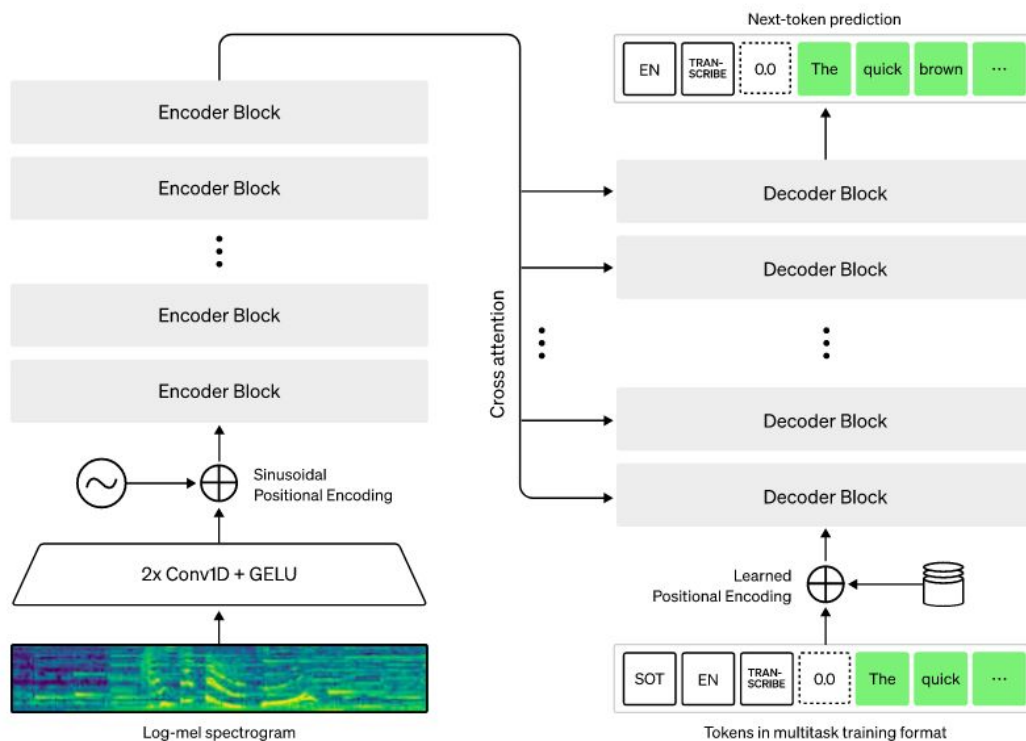TM-score = 97.4, $N_{res}$ = 1,246, PDB ID = 6E3K

MICROB-PREDICT

EMBL
50 YEARS|1974–2024

# Tool assisted research

~~Tool assisted research~~

General purpose AI tools

MICROB-PREDICT

EMBL
50 YEARS | 1974–2024

# Audio assistants - Dictation

- Faster than typing
- Video transcription
- (Translation)
- Not ideal for code

OpenAI Whisper

# Audio assistants - Presentation rehearsal

- Iterative
- Always available
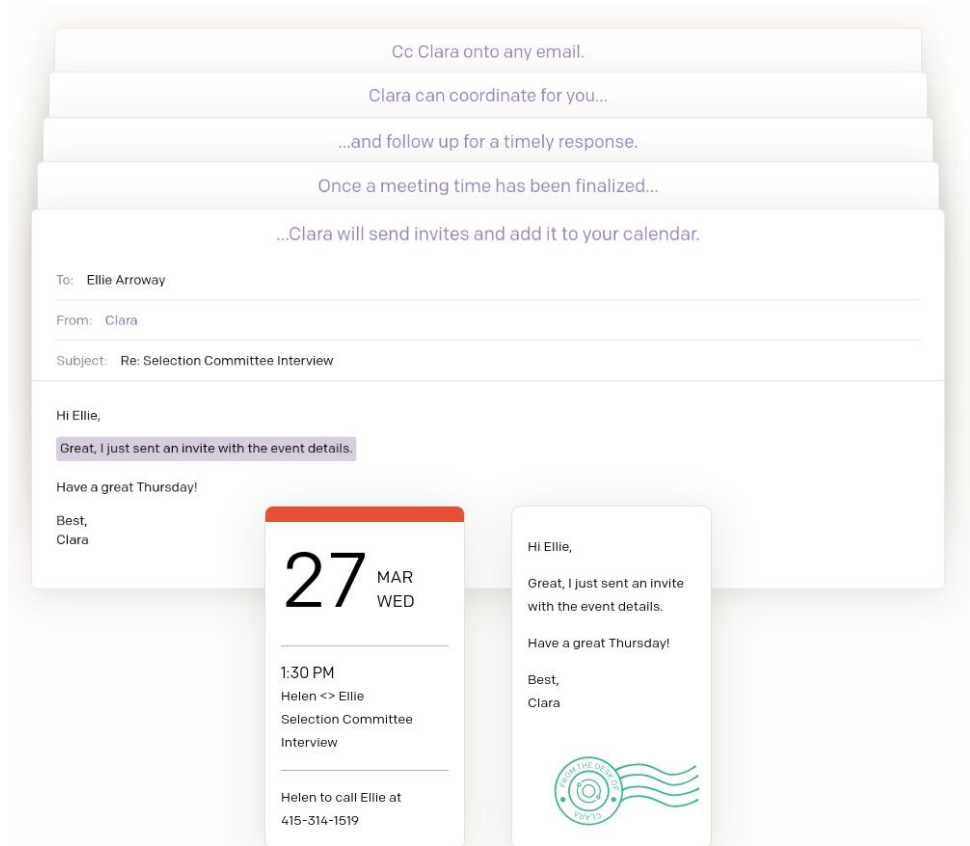- Remove shame factor
- Decent advice
- Self awareness

Microsoft Speaker Coach

# Planning assistants

- Death to Doodles
  - RIP
  - Seriously
  - Please

**Clara,** can you find a time for us to meet next week?

Clara is your partner in doing great work — a virtual employee that schedules your meetings.

Cc Clara onto any email.

Clara can coordinate for you...

...and follow up for a timely response.

Once a meeting time has been finalized...

...Clara will send invites and add it to your calendar.

To:    Ellie Arroway

From:   Clara

Subject:   Re: Selection Committee Interview

Hi Ellie,

Great, I just sent an invite with the event details.

Have a great Thursday!

Best,
Clara

**27** MAR WED

1:30 PM
Helen <> Ellie
Selection Committee
Interview

Helen to call Ellie at
415-314-1519

Hi Ellie,

Great, I just sent an invite with the event details.

Have a great Thursday!

Best,
Clara

claralabs.com

MICROB-PREDICT

EMBL
50 YEARS | 1974–2024

# Email assistants

- Draft reply
- Tidy-up and clarify
- No more typos

🎉 Wedding

Ellie replies...

Jamie and Danielle,

Thank you for joining us on our special day, we are so happy that all of our friends and family could be there to celebrate with us! We truly had the best day, and we hope you enjoyed it as much as we did.

🔄 Follow up

Thank you so much for the generous gifts. The Tim Hortons kept us fueled on our long road trip through the Rockies, and the Moet was very much enjoyed from the hot tub at our cabin in Cornwall, which we visited straight after the wedding. After the hectic run up to the wedding, it was a wonderful way to relax!

🎁 Gifts

Hope you're both well and enjoying #boatlife.

Lots of Love,

🔄 Follow up

Ash and Jen
xxxx

🥰 Casual correspondence

Hey Ash and Jen,

Your wedding was absolutely beautiful, and Danielle and I felt so privileged to be a part of it!

So stoked to hear that you guys had such a fantastic time on your road trip and at the cabin! It sounds like pure bliss, especially that hot tub session with the Moet - now that's how to celebrate!

We're glad you enjoyed the Tim Hortons and the bubbly - small tokens for such a grand occasion.

Boat life is treating us well, as always, maybe it will also take us to Cornwall one day!

Hope you two are settling into married life smoothly!

Take care and lots of love,
James and Danielle

MICROB-PREDICT

EMBL
50 YEARS | 1974–2024

# Language assistants

- Write like a native
- Understand shortcomings

Grammarly

# Improved search engines

- Natural language
- Information with references

perplexity.ai

# Tool assisted research

# Litmaps / Research rabbit / Inciteful

## Literature

- Article discovery
- Relatedness
- Citation map
- Contextualization

litmaps.com / researchrabbitapp.com / inciteful.xyz

# Scite

## Literature

- Review
- Summarization
- Reference verification
- Contextualization

scite.ai

# Elicit

## Literature

- Semantic analysis
- Review
- Summarization
- Contextualization
- PDF → "Write me an abstract"

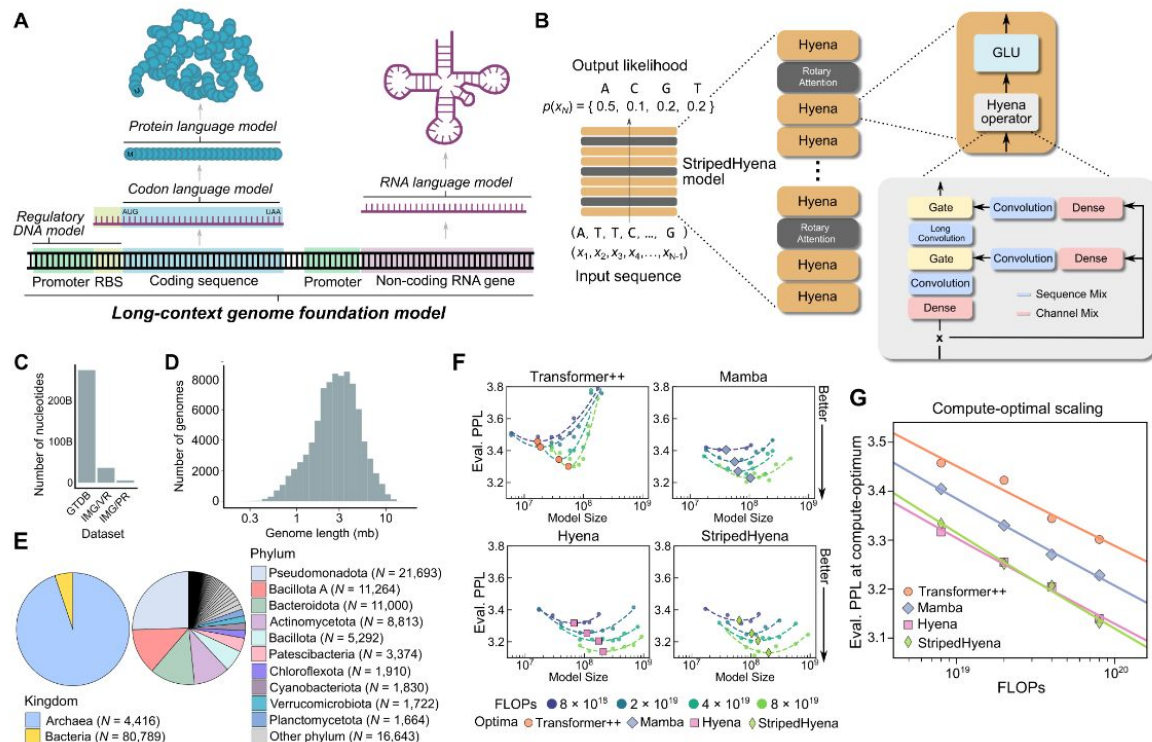| Paper | Abstract summary | Main findings | Methodology |
|---|---|---|---|
| **The human genome project.**<br>W F Bodmer<br>*Revista de investigacion clinica; organo del Hospital de Enfermedades de la Nutricion*<br>1994 6 citations DOI | The initial investment in the Human Genome Project has given an enormous return, both scientifically and financially. | The Human Genome Project successfully sequenced the complete human genome by 2006, providing a more realistic estimate of the number of protein-coding genes. The project's outcomes include opportunities for analyzing human genetic variability, understanding disease relationships, and improving knowledge about human origins and gene expression patterns. The investment in the HGP has yielded significant scientific and financial returns. | The methodology involved the development of DNA cloning and sequencing technologies, the use of BACs for mapping overlapping clones, the 'shotgun' approach for sequencing, the utilization of automatic sequencers, and computational power. The sequencing results were published in stages from 2001 to 2006. |
| **The Human Genome Project**<br>F. Collins +1<br>*Alcohol health and research world*<br>1995 5 citations | The Human Genome Project has produced results that are permeating basic biological research and clinical medicine. | The Human Genome Project has produced impactful results in basic biological research and clinical medicine, including the successful mapping of the mouse genome and ongoing work on the genetic map of the rat for studying complex disorders. | The methodology used in the study involves developing research tools to identify disease-related genes, examining ethical and social implications of genetic technologies, and mapping the genomes of mouse and rat. |
| **Sequencing the human genome.**<br>L. Rowen +2<br>*Science*<br>1997 99 citations DOI | The human genome project is at the halfway point. | The human genome project has made significant progress with the completion of microbial genomes, but the human genome sequencing is only 2 percent finished, posing a challenge to meet the | The methodology in this paper is not explicitly stated, as it does not provide a detailed description of the methods used in the study. |

# Genomic models - capture biological languages



**Evo,** a long-context biological foundation model based on the StripedHyena architecture that generalizes across the fundamental languages of biology: DNA, RNA, and proteins.

**Evo** is capable of both prediction tasks and generative design, from molecular to whole genome scale (over 650k tokens in length).

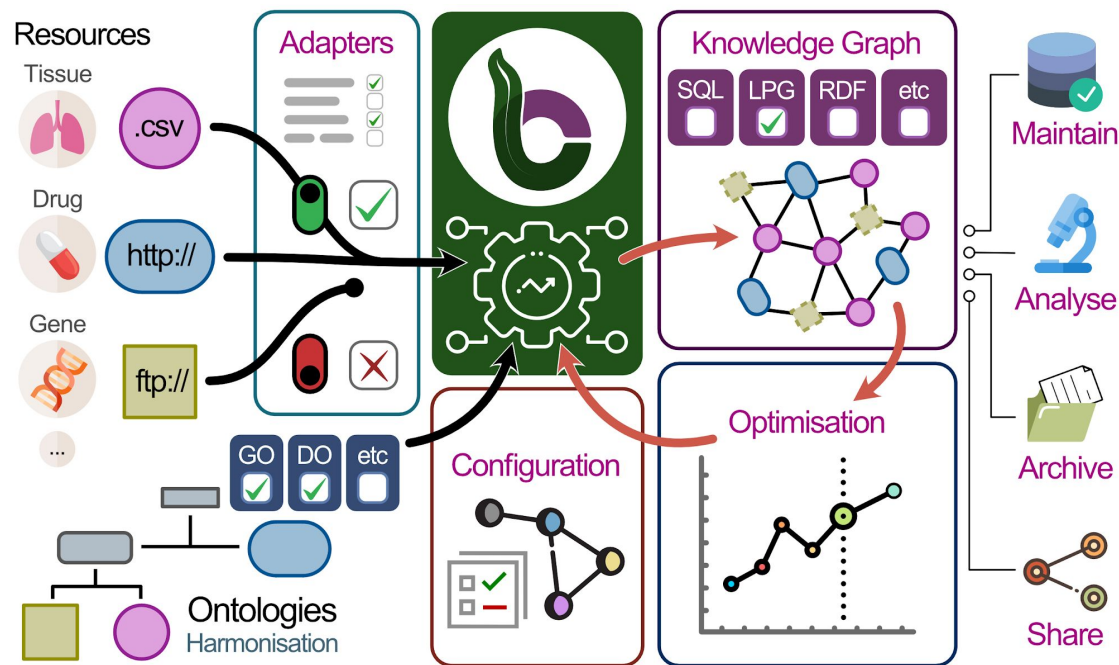**Evo** is trained at a nucleotide (byte) resolution, on a large corpus of prokaryotic genomic sequences covering 2.7 million whole genomes.

Evo model, ArcInstitute - source and BioRxiv doi:10.1101/2024.02.27.582234v1

# Knowledge exploration

**Complex data**

- Graphs for structured knowledge
- AI as gateway
- Powerful combo

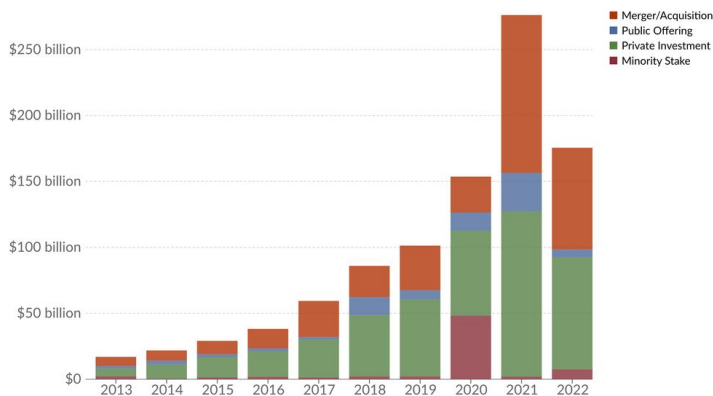- Great documentation
- Complex to setup

# GPT-4V - Vision

## Figures

- Interpretation
- "Write me a legend"



Annual global corporate investment in artificial intelligence, by type
This data is expressed in US dollars, adjusted for inflation.

Legend:
- Merger/Acquisition
- Public Offering
- Private Investment
- Minority Stake

**Data source:** NetBase Quid via AI Index Report (2023)
**Note:** Data is expressed in constant 2021 US$. Inflation adjustment is based on the US Consumer Price Index (CPI).

OurWorldInData.org/artificial-intelligence | CC BY

🤖: This is a bar chart visualizing the **annual global corporate investment in artificial intelligence** (AI) by type, from 2013 to 2022. The data are presented in U.S. dollars and have been adjusted for inflation to 2021 dollars. This adjustment is based on the U.S. Consumer Price Index (CPI).

The chart is broken down into four types of investments:

1. **Merger/Acquisition**: represented by red on the top of some columns, this indicates the amount of money invested in AI through mergers and acquisitions.

2. **Public Offering**: shown in blue, this includes investments made through public offerings, possibly referring to investments made during IPOs or subsequent public investments into AI companies.

3. **Private Investment**: shown in green, this represents private funding rounds such as venture capital, private equity, or other types of private placement.

4. **Minority Stake**: represented in orange, this signifies investments made by purchasing a minority stake in AI companies.

The vertical axis indicates the amount of investment in billions of dollars, while the horizontal axis represents the year. **There's a clear trend of increasing investment in AI over the years**, with a significant jump from 2020 to 2021 and a slight decrease from 2021 to 2022. In 2022, the investments are spread across all four types, with private investment appearing to comprise the largest share.

31

MICROB-PREDICT

EMBL 50 YEARS|1974–2024

# ChatGPT - for code

## Code troubleshooting

- Beginner friendly
- Explains with context
- Almost as good as a human, sometimes better



**You**
Please help me understand the meaning of the following 3 errors in the Python language:

```
In [1]: import shouldexist

ModuleNotFoundError                       Tra
Cell In[1], line 1
----> 1 import shouldexist

ModuleNotFoundError: No module named 'shoulde

In [2]: import itertools

In [3]: list(itertools.chain(1, 2, 3))

TypeError                                 Tra
Cell In[3], line 1
----> 1 list              .         1  2  3

TypeError: 'int' object is not iterable

In [4]: 10 / 0

ZeroDivisionError                         Tra
Cell In[4], line 1
----> 1 10 / 0

ZeroDivisionError: division by zero
```

**ChatGPT**
These are three common errors in Python:

1. **ModuleNotFoundError**: This error occurs when Python cannot find the module you are trying to import. In the provided example, Python is unable to find a module named `shouldexist`. This could happen due to various reasons, such as a typo in the module name, the module not being installed in your environment, or the module not being accessible due to incorrect path settings.

2. **TypeError: 'int' object is not iterable**: This error happens when you try to iterate over an object that is not iterable. In Python, an iterable is an object capable of returning its members one at a time. In the provided example, `itertools.chain()` expects iterable arguments, but you passed integers (`1`, `2`, `3`) instead. To fix this error, you need to pass iterable objects, such as lists, tuples, or other iterable types.

3. **ZeroDivisionError: division by zero**: This error occurs when you attempt to divide a number by zero, which is mathematically undefined. In Python (and many other programming languages), division by zero raises an exception to indicate an error. In the provided example, you tried to divide `10` by `0`, which resulted in a `ZeroDivisionError`. To avoid this error, ensure that you don't divide by zero in your code or handle the zero division case gracefully if it's inevitable in your logic.

# ChatGPT / CoPilot / Devin

**Code generation**

- Accessible to non-programmers
- From pseudo-code
- From documentation
- Semi-automated testing
- Very fast prototypes

- Harder to use with complex codebases

# AI Personas / AI Agents / AI Council

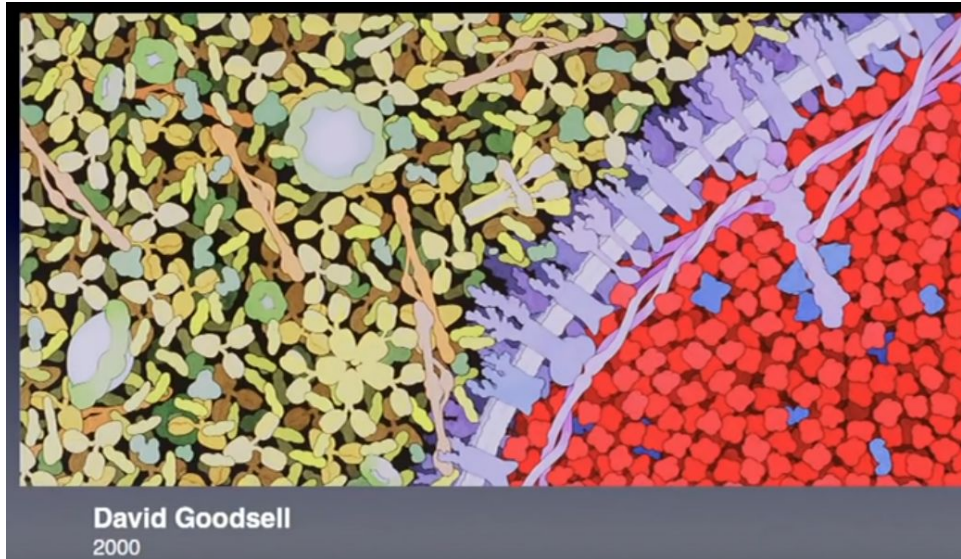**Reasoning / Ideation**

- Understand points of view
- Highlight biases
- Complementary feedback
- Not (yet) self sufficient

Image from spyscape.com / ChatGPT AI Council

# Scientific communication

# Illustrations



David Goodsell
2000



WEHI.TV
DREW BERRY

COVID-19 vs HUMAN ANTIBODIES x1,000,000

# Video generation



**Prompt**: A gorgeously rendered papercraft world of a coral reef, rife with colorful fish and sea creatures.

Sora by OpenAI

MICROB-PREDICT

EMBL
50 YEARS|1974−2024

# When (data) privacy is important

MICROB-PREDICT

EMBL
50 YEARS | 1974–2024

# Running models privately



Get up and running with large
language models, locally.

Run Llama 2, Code Llama, and other
models. Customize and create your own.

ollama.com

# The right model and hardware



| | Small / Medium LLM (e.g., <50bn parameters) | Large / Extra-Large LLM (e.g., >50bn parameters) |
|---|---|---|
| Performance | **Low** Can be improved with fine-tuning | **High** Can achieve better accuracy and more nuanced text |
| Cost of inference (run-time) | **Low** More economical to scale | **High** More expensive to run |
| Latency | **Low** Faster generations suitable for interactive applications | **High** Slower generations due to more compute required for the larger number of parameters |
| Knowledge | **Low** May not perform as well with deeper level of tasks, but can be optimized with RAG | **High** Can capture more information and variety of tasks |
| Understanding | **Low** Lacks in tasks that require complex reasoning | **High** Better at understanding context for complex reasoning |
| Environmental impact | **Low** Consumes less, lower carbon-footprint | **High** Consumes more, higher carbon-footprint |
| Deployment complexity | **Low** Fast and easier to deploy and integrate, especially on edge devices. More accessibility to a broader set of accelerators. | **High** Harder to deploy and maintain, requiring a higher level of expertise to manage. Tends to require large and new generation accelerators. |

# A lot to choose from

A Survey of Large Language Models - doi: 10.48550/arXiv.2303.18223

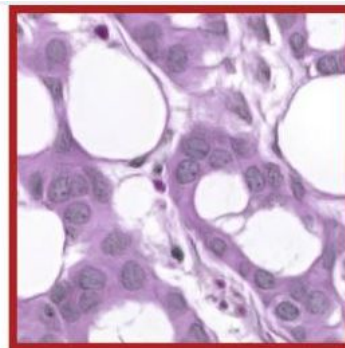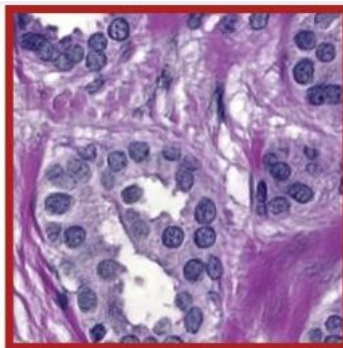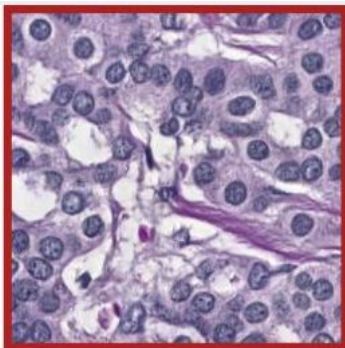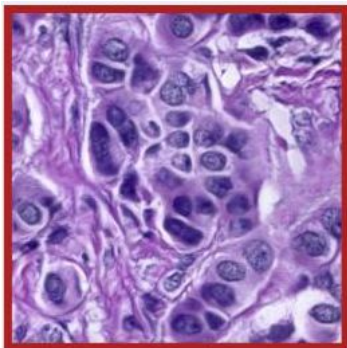# Wrapping up

# Summary & Challenges / Opportunities

**AI can be immensely useful if used correctly**

- Trainees are more independent / Not afraid to ask
- "Hallucinations" / incorrect "*facts*" → ok for reasoning / not-ok for decision making
- Data privacy / *Cross*-contamination → ok for generation / not-ok for review
- Copyright / Intellectual property rights → Regulation playing catch → Lawsuits
- Computation / Energy / Hardware cost → 1M€ BLOOM LLM → "*GPUrush*"
- Cultural bias → Notions of right and wrong → Really hard problem
- Even greater information overflow → AI as solution → circular problem
- Society implications → New jobs / Obsolete jobs → Automation everywhere
- Fast pace and hype → Widening technology gap low-, mid-, high-income countries
- Challenging current definitions: creativity, originality, plagiarism,...
- Human-machine interactions → More forgiving towards humans than machines

MICROB-PREDICT

EMBL
50 YEARS|1974–2024

# Thank you

Michael Khun
Marisa Keller

AI discussion group
Bork group

# Thank you / Questions?
### PS: Find the AI generated

They are all AI generated

MICROB-PREDICT

EMBL
50 YEARS|1974–2024