



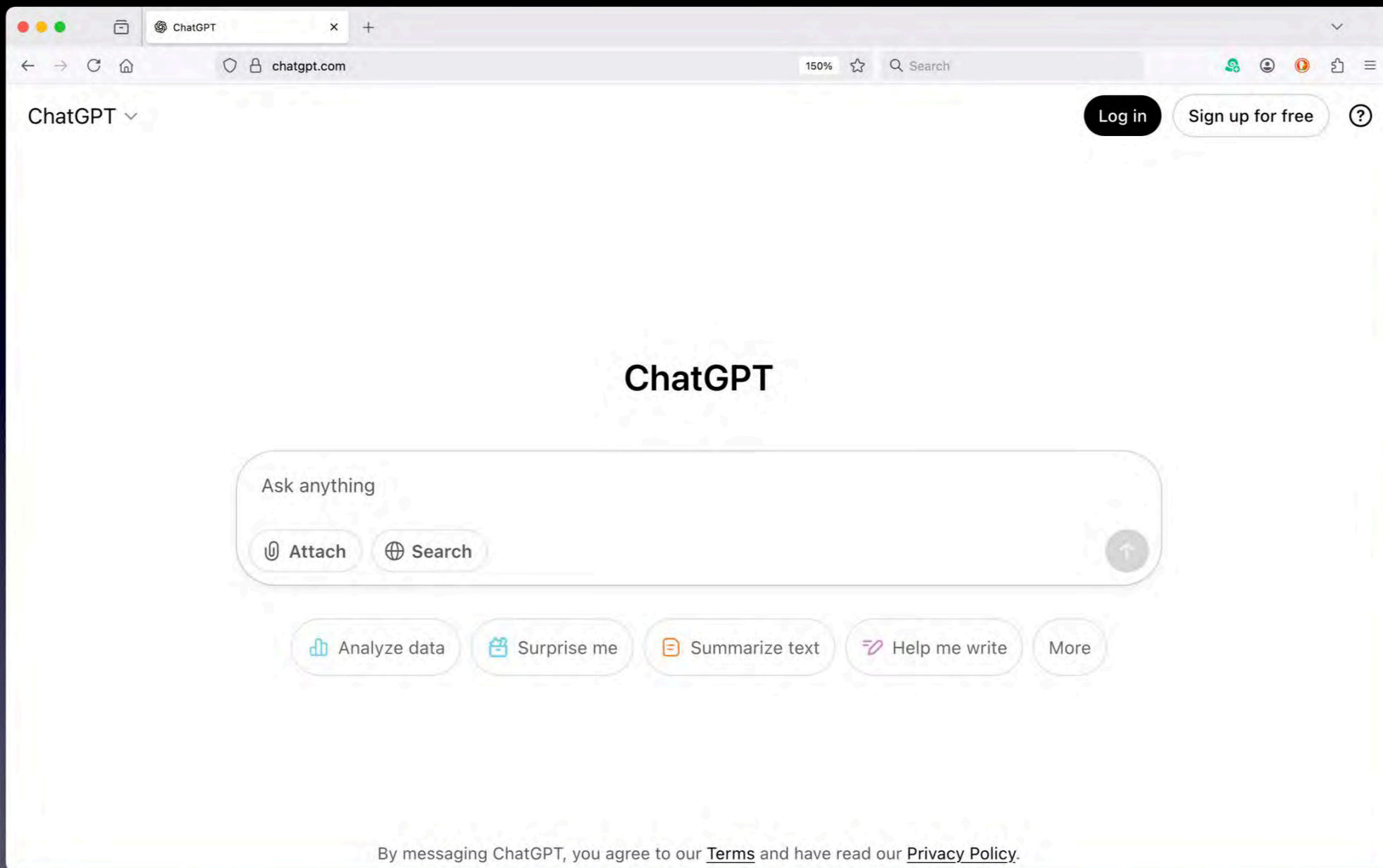
TECH TALK
ERIC FURST

Large Language Models as Tools for Science

— CBI 5th Annual Summer Retreat, July 30, 2025

Outline

1. How LLMs work
2. Running LLMs locally
3. LLM uses



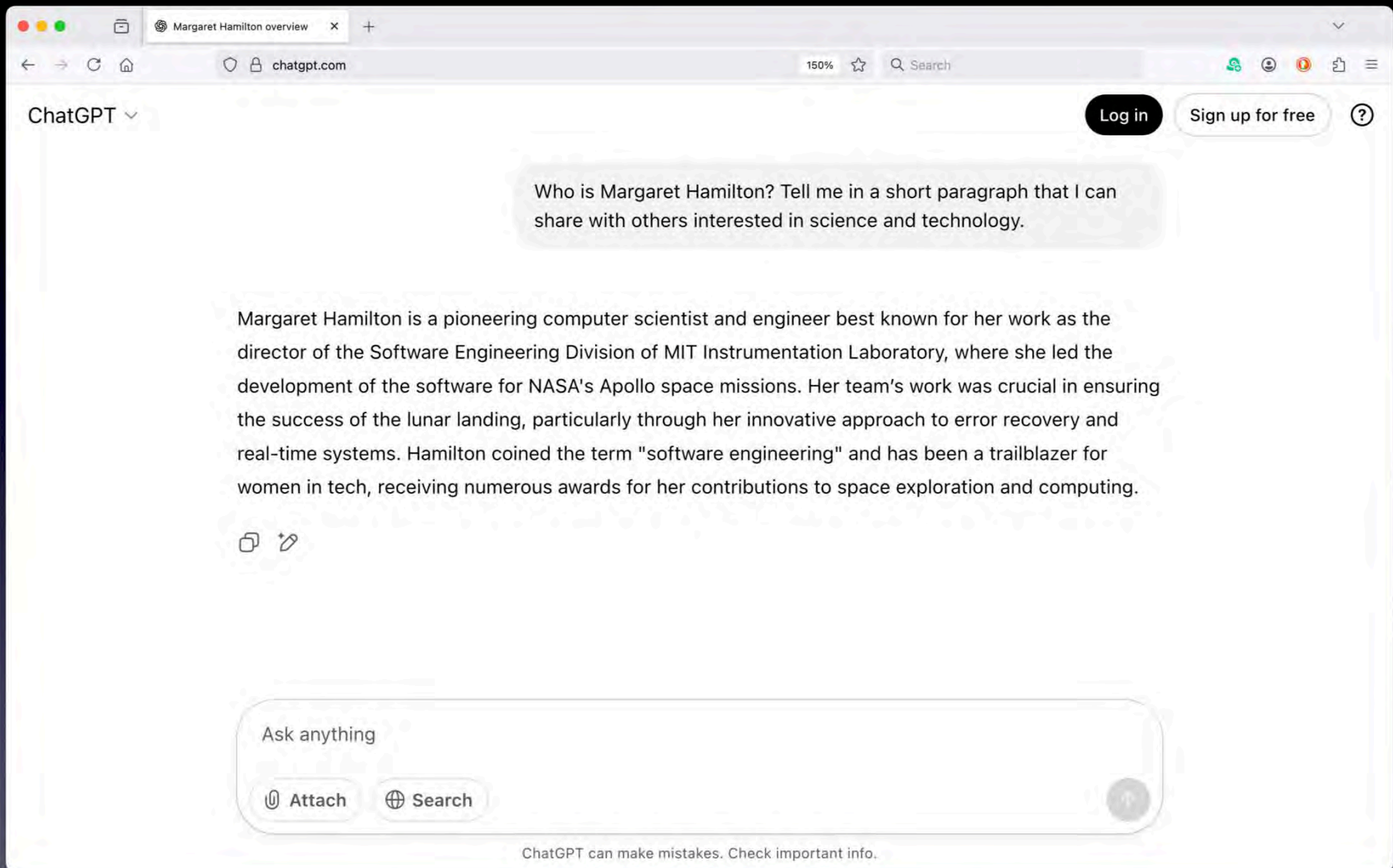
What happens?

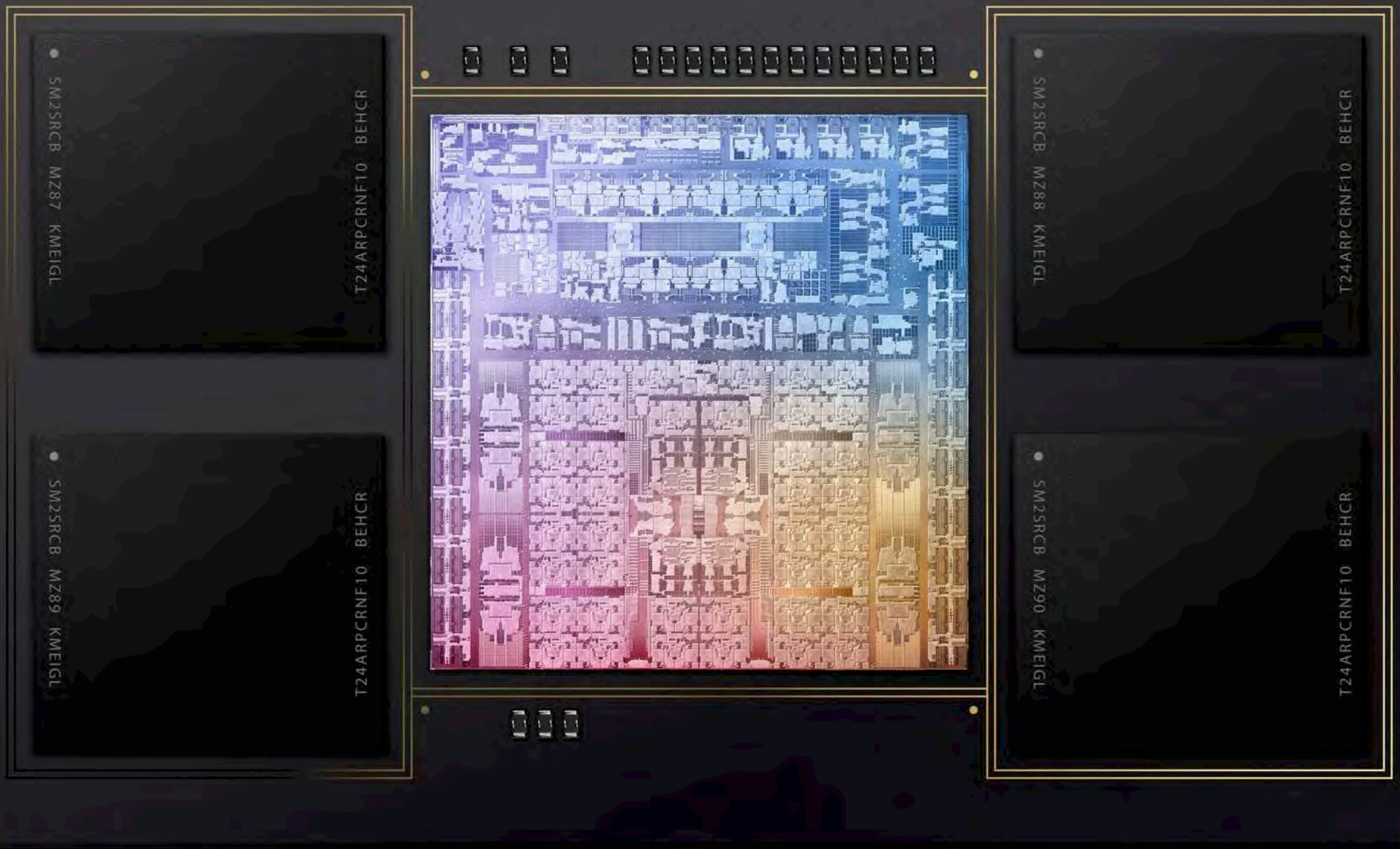
We input text — a “prompt”

The model responds with text

The quality of the response can depend strongly on the prompt

The response is never *exactly* the same





GPU – Graphics Processing Unit

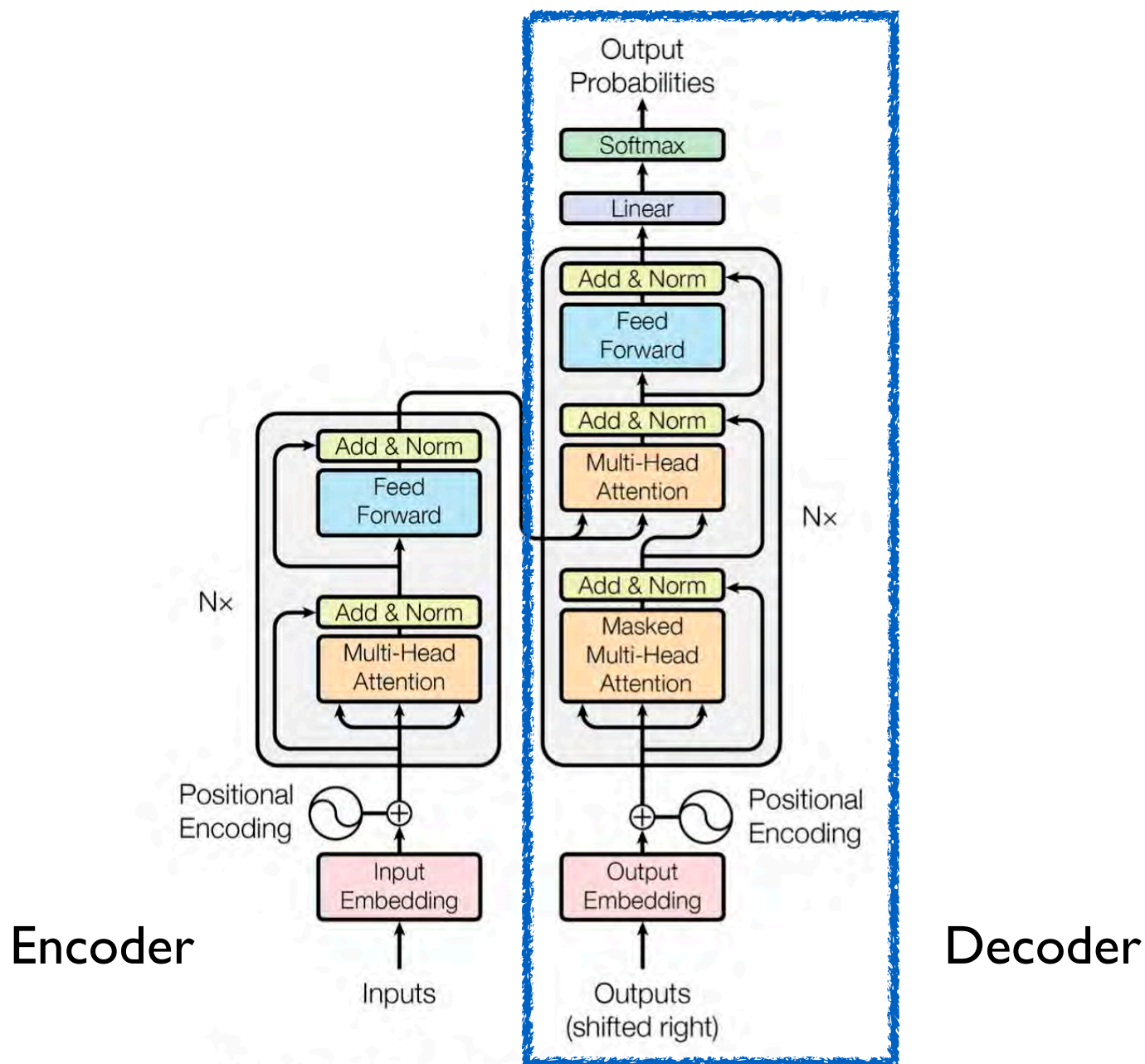


Figure 1: The Transformer - model architecture.

GPT = Generative Pre-trained Transformer

LLM is an auto-regressive language model that uses an optimized transformer architecture

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

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Llion Jones*
Google Research
llion@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 part designed, implemented, tuned and evaluated countless model variants in our original codebase, and experimented with novel model variants, was responsible for our initial codebase, and and Aidan spent countless long days designing various parts of and greatly improving results and massively accelerating

Attention Is All You Need, 2017

A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, Attention Is All You Need, in Advances in Neural Information Processing Systems, Vol. 30 (Curran Associates, Long Beach, CA, USA, 2017), pp. 261–272.

arXiv:1706.03762

followed by the...

First decoder-only GPT

P. J. Liu et al., “Generating Wikipedia by Summarizing Long Sequences,” presented at the ICLR, 2018. Accessed: May 10, 2025. [Online]. Available: <https://openreview.net/pdf?id=Hyg0vbWC->

Given a sequence of tokens
(characters, words, bigrams,
or subwords)

A logit is the raw...

Predict the next probable token

out...

Which becomes part of the
next token prediction, etc.

(It is auto-regressive)

*... put... of... a... ma... chine...
learn ing mo del, typic ally be fore
apply ing a trans form a tion like
the soft max func tion.*

Probabilistic but causal calculation of the next token

≡

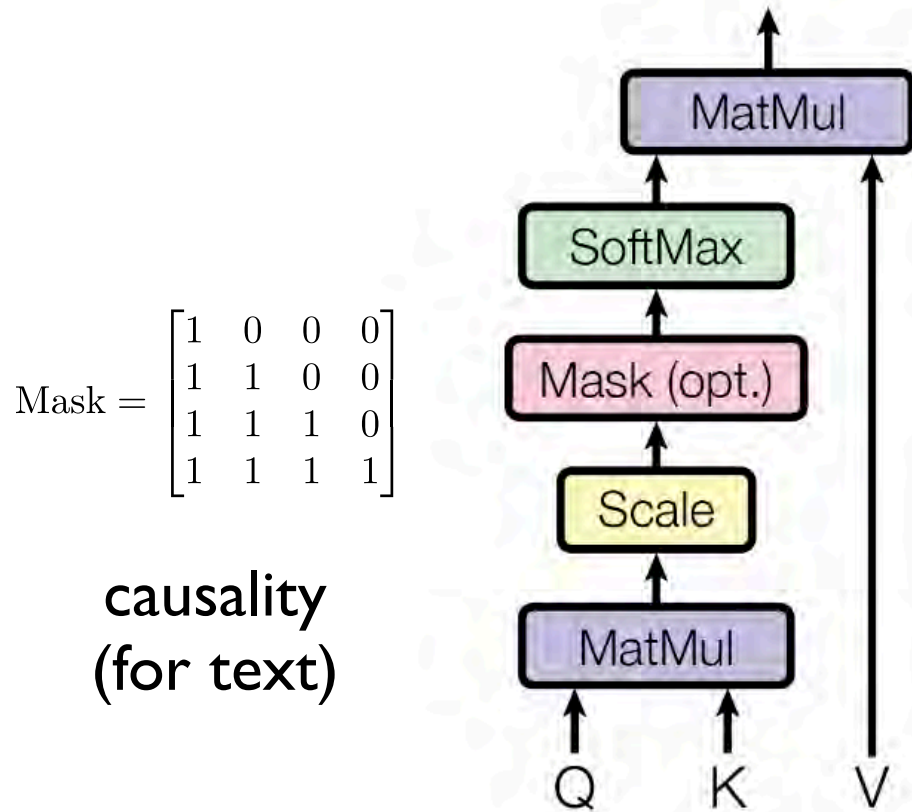


Search

🌐 50 languages ▼



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



Probability of next token (SoftMax)

$$P_i = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}$$

Given vector
of logits

$$z = [z_1, z_2, \dots, z_n]$$

Boltzmann weighting

$$P_i = \frac{e^{-E_i/k_B T}}{\sum_j e^{-E_j/k_B T}}$$

Query, Key, and
Value matrices

“Temperature” hyperscaling parameter

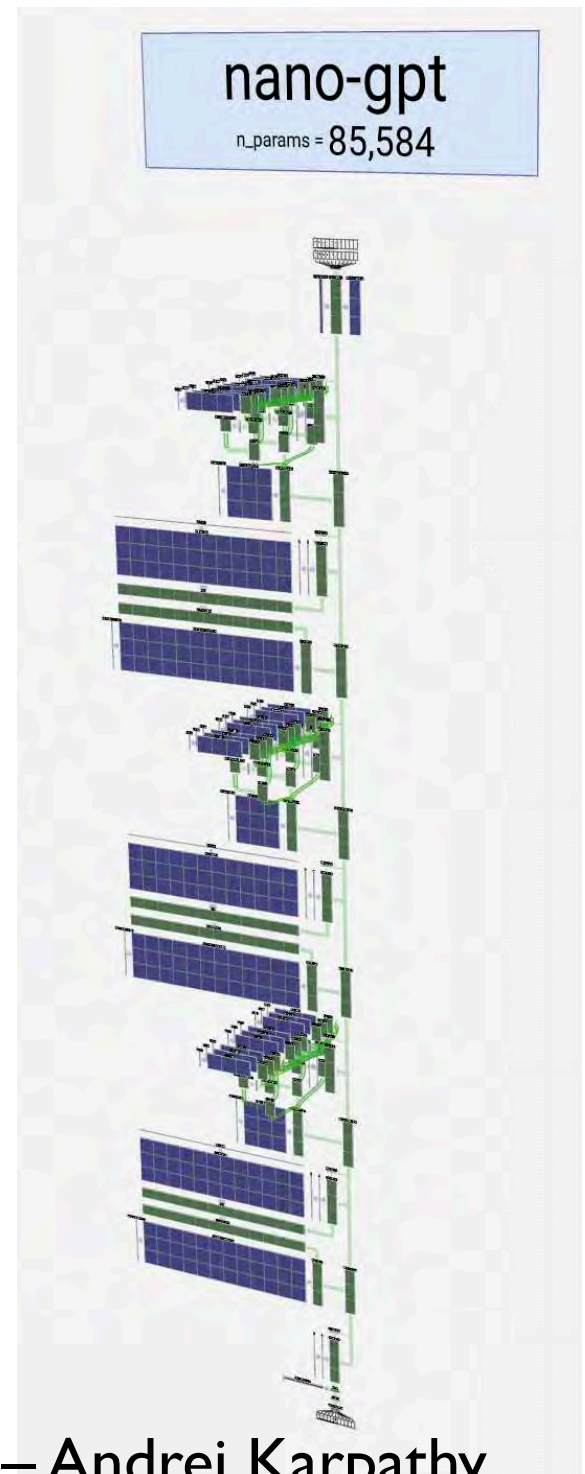
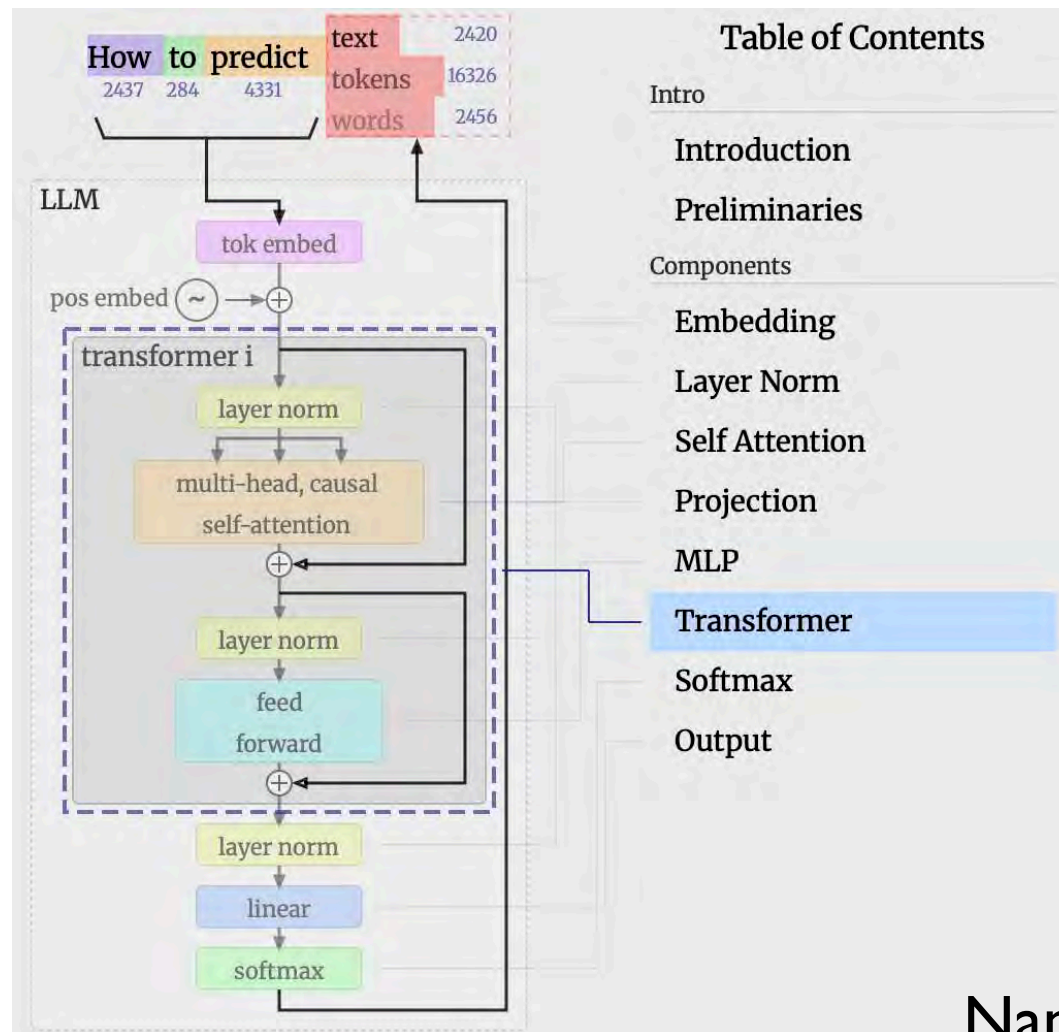
$$P_i = \frac{e^{z_i/T}}{\sum_{j=1}^n e^{z_j/T}}$$

Stochastic behavior!

A.Vaswani, et al., Attention Is All You Need, in Advances in Neural Information Processing Systems, Vol. 30 (Curran Associates, Long Beach, CA, USA, 2017), pp. 261–272.

Visualizing a GPT

Brendan Bycroft – <https://bbycroft.net/llm>



NanoGPT – Andrej Karpathy

Given an input token representation, X

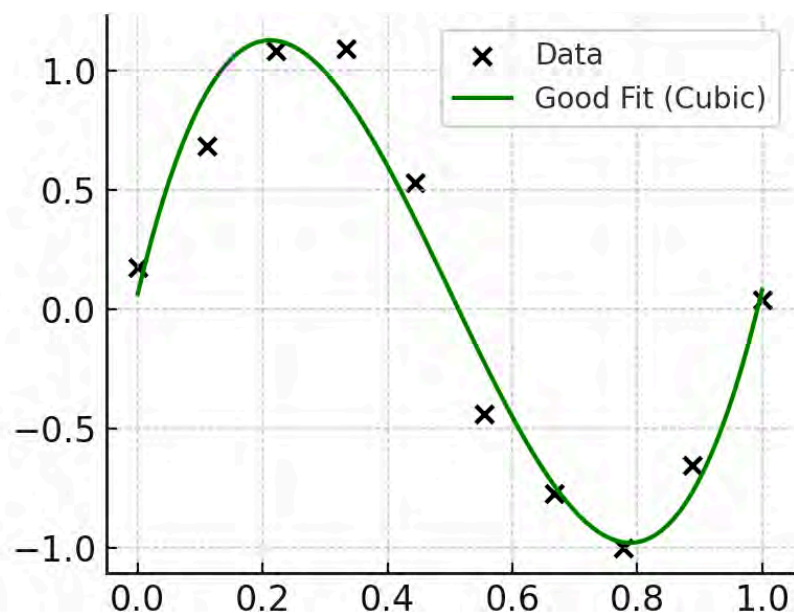
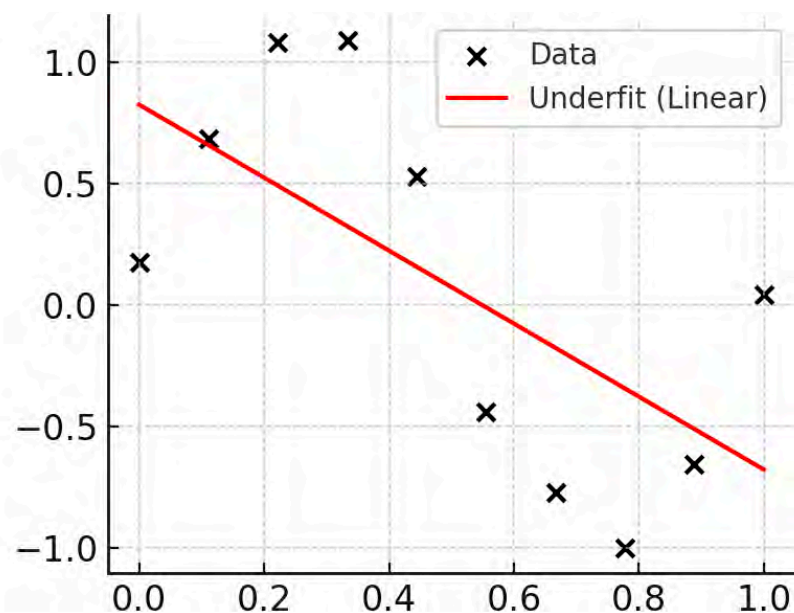
$$\left. \begin{aligned} Q &= XW_Q \\ K &= XW_K \\ V &= XW_V \end{aligned} \right\} \begin{array}{l} \text{Learnable weight matrices} \\ \text{each of size} \\ d \times d_k \text{ or } d \times d_v \quad (d \text{ is } n_{\text{embd}}) \\ \text{the embedding} \\ \text{dimension} \end{array}$$

GPT-3 (175B) ~175 billion parameters (350GB)
ca. June 2020 trained on $O(10\text{TB})$ data (the web)

- Hidden size: 12,288
- Number of layers: 96
- Number of attention heads: 96
- Vocabulary size: ~50,000
- Feedforward network expansion factor: 4x

GPT-4 and 4o sizes not released, but estimates at $10^{12} - 10^{14}$ parameters

Training model weights



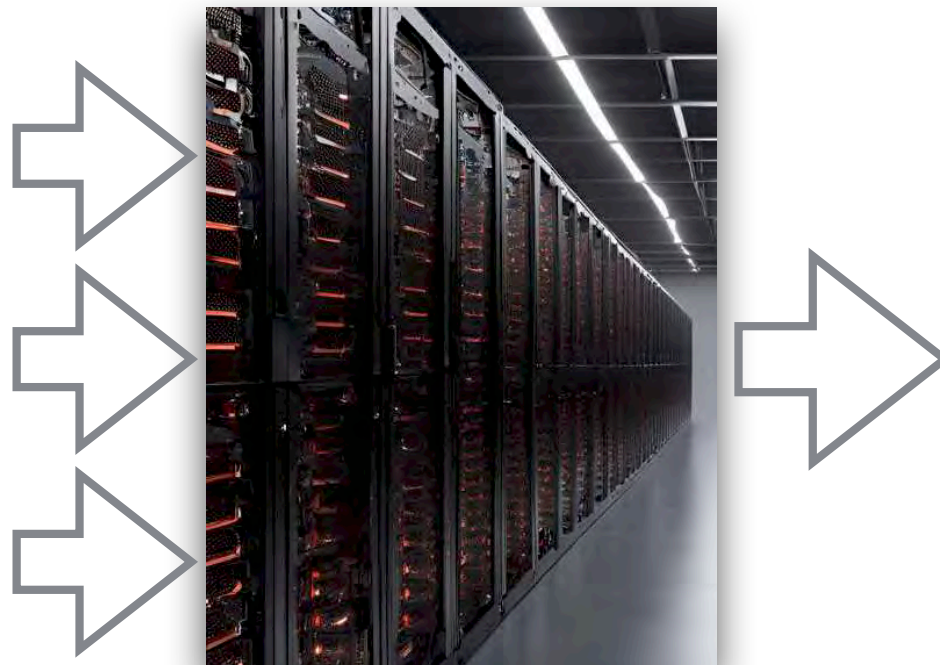
Parameters (model weights) adjusted to better fit the model
Instead of *sum of squared residuals*, use *cross-entropy loss*

Cost (compute & energy) is in the model training

~10 TB text

*webcrawl,
Wikipedia,
Project
Gutenberg,
ArXiv, Stack
Exchange...
(llama)*

+ fine-tuning



Llama 3.1 (Meta)

8B model: 4.9 GB

70B model: 43 GB

405B model: 243 GB

Llama 3.1 8B – 1.46 million GPU hours

Llama 3.1 70B – 7.0 million GPU hours

Llama 3.1 405B – 30.84 million GPU hours

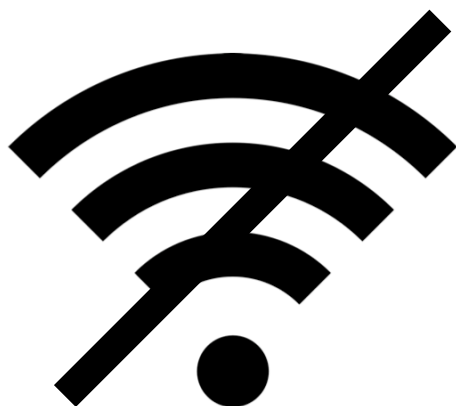
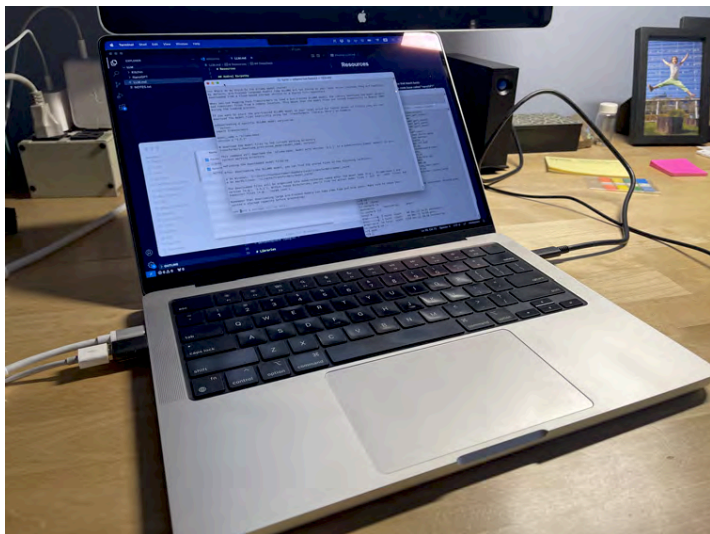
(approximately 440,000 petaFLOP-days for training)

<https://huggingface.co/blog/llama3.1>

<https://ollama.com/library/llama3.1>

Running models is relatively *inexpensive*

Macbook Pro M3, 12 cores (2023)
18GB RAM, 18 core GPU



running in *ollama* locally
with llama3 8B (4.7 GB)

```
-rw-r--r-- 1 furst staff 51M Jan 16 11:37 ollama*
```

[HTTPS://OLLAMA.COM/](https://ollama.com/)

```
furst — ollama run llama3 — 122x40

~/ $ ollama run llama3
[>>> Say hello to the Winter Research Review.
A warm welcome to the Winter Research Review!

As the snowflakes gently fall outside, we're excited to bring you the latest and greatest in research discoveries from various fields. Whether you're a scientist, student, or simply someone curious about the world around us, this review is designed to provide a comprehensive overview of the most significant findings from the past year.

From groundbreaking advancements in medicine and technology to intriguing insights into human behavior and the natural world, we've got you covered. So, grab a hot cup of cocoa, get cozy, and let's dive into the Winter Research Review!

[>>> How do I finetune ollama models on my own text?
Finetuning OLLAMA (Open Large-scale LAnguage Model Architecture) models on your own text can be a fantastic way to adapt the model to your specific domain or use case. Here's a step-by-step guide to help you finetune an OLLAMA model:

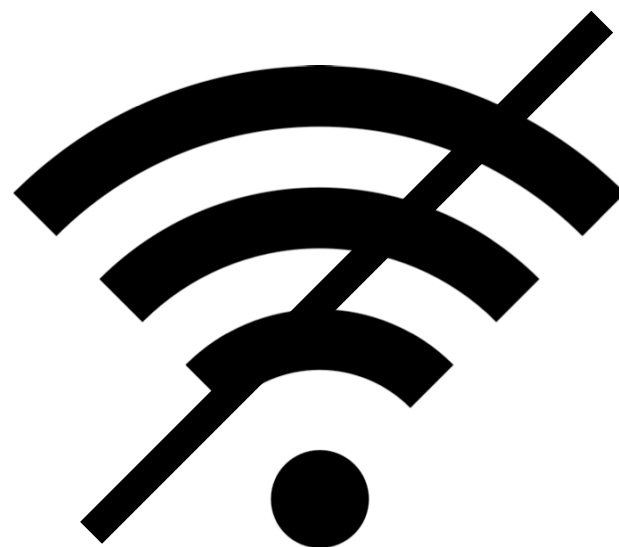
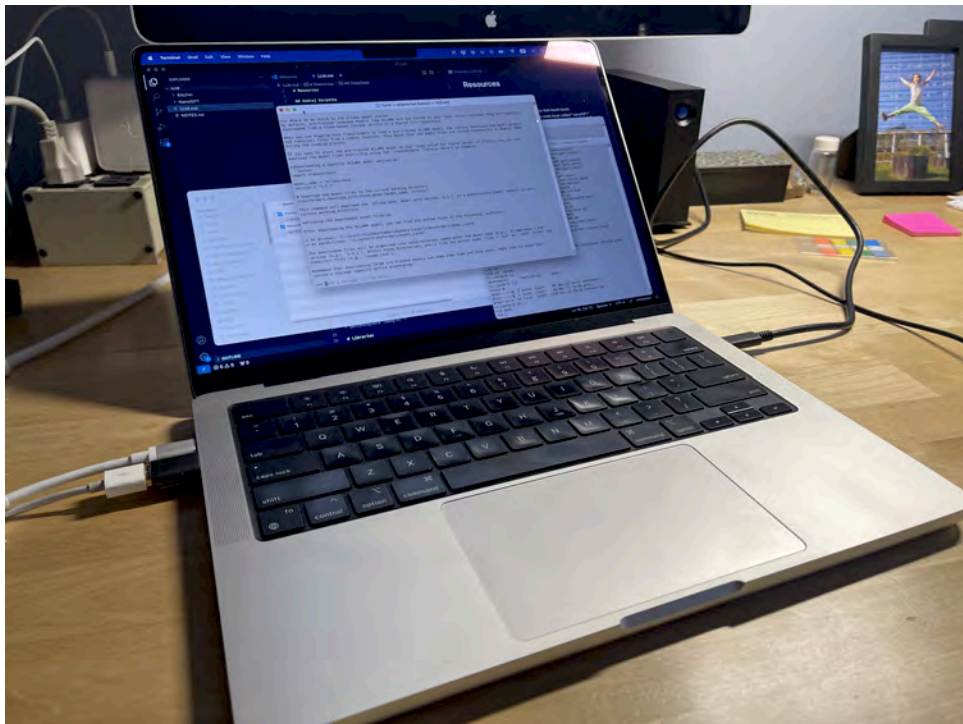
**Prerequisites:**

1. **Text dataset**: Prepare a large text dataset related to your topic of interest. This dataset will serve as the training data for fine-tuning the model.
2. **OLLAMA pre-trained model**: You can use any publicly available OLLAMA pre-trained model, such as those provided by Hugging Face Transformers or Google's BERT models.
3. **Python and necessary libraries**: Install Python (3.x) and the required libraries:
   * `transformers` for working with pre-trained language models
   * `torch` for building and training neural networks

**Step-by-Step Guide:**

1. **Prepare your text dataset**:
   * Preprocess your text data by tokenizing, normalizing, and encoding it using a suitable tokenizer (e.g., `bert-base-uncased`).
   * Split your dataset into training (~80%), validation (~10%), and testing sets (~10%).
2. **Load the pre-trained OLLAMA model**:
   * Import the required libraries: `transformers`, `torch`
   * Load the pre-trained OLLAMA model using the `AutoModelForSequenceClassification` class from Hugging Face Transformers
3. **Create a finetuning dataset class**:
   * Define a custom dataset class that inherits from `Dataset` (or `torch.utils.data.Dataset`)
   * Implement the `__getitem__` and `__len__` methods to handle the text data and its corresponding labels (if any)
```

Demo: running models locally



Macbook Pro M3, 12 cores (2023)
18GB RAM, 18 core GPU

llama3 8B (4.7 GB)
running in *ollama*

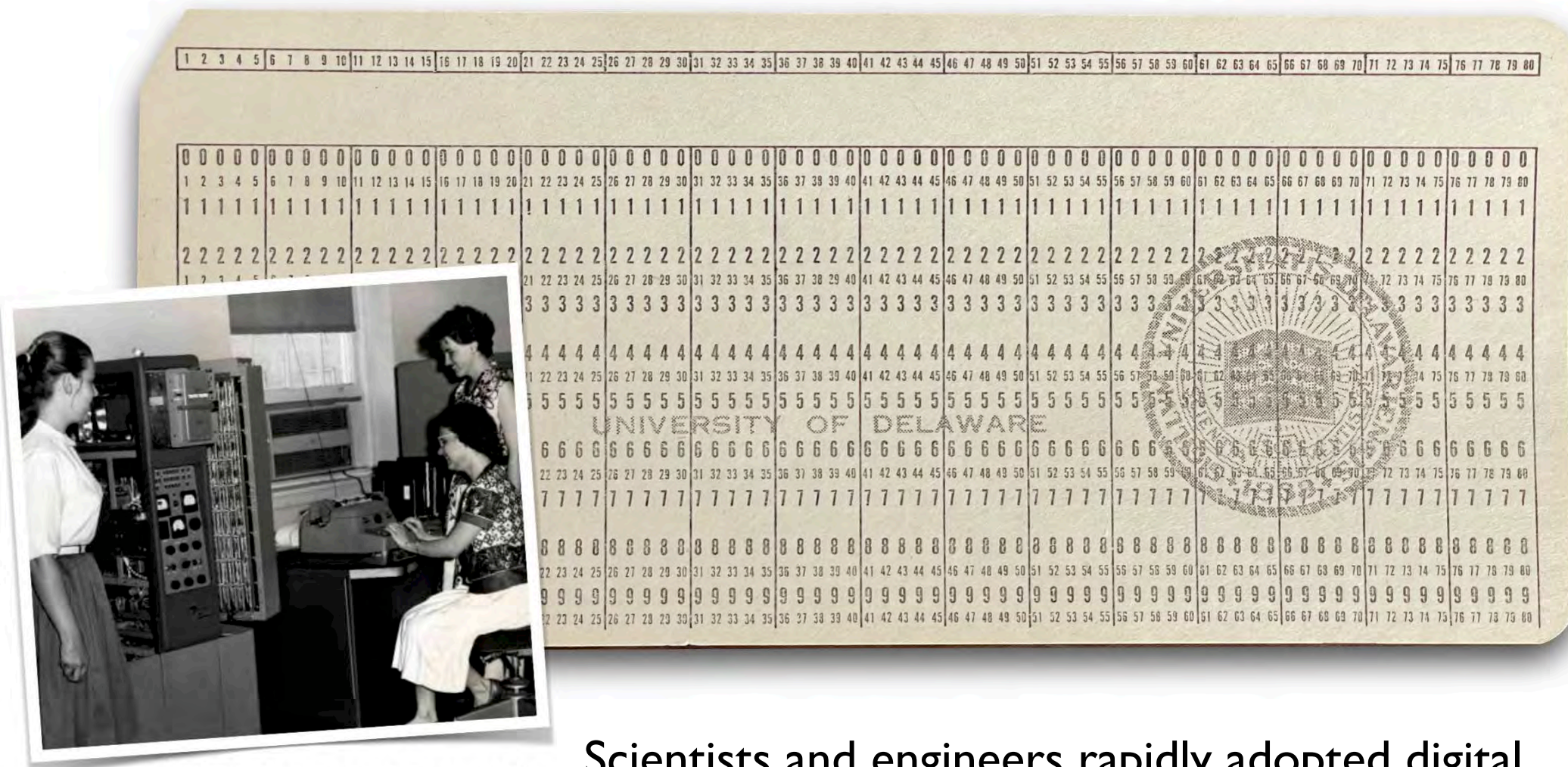
[HTTPS://OLLAMA.COM/](https://ollama.com/)

Outline

1. How LLMs work
2. Running LLMs locally
3. LLM uses

LLM uses in science and engineering

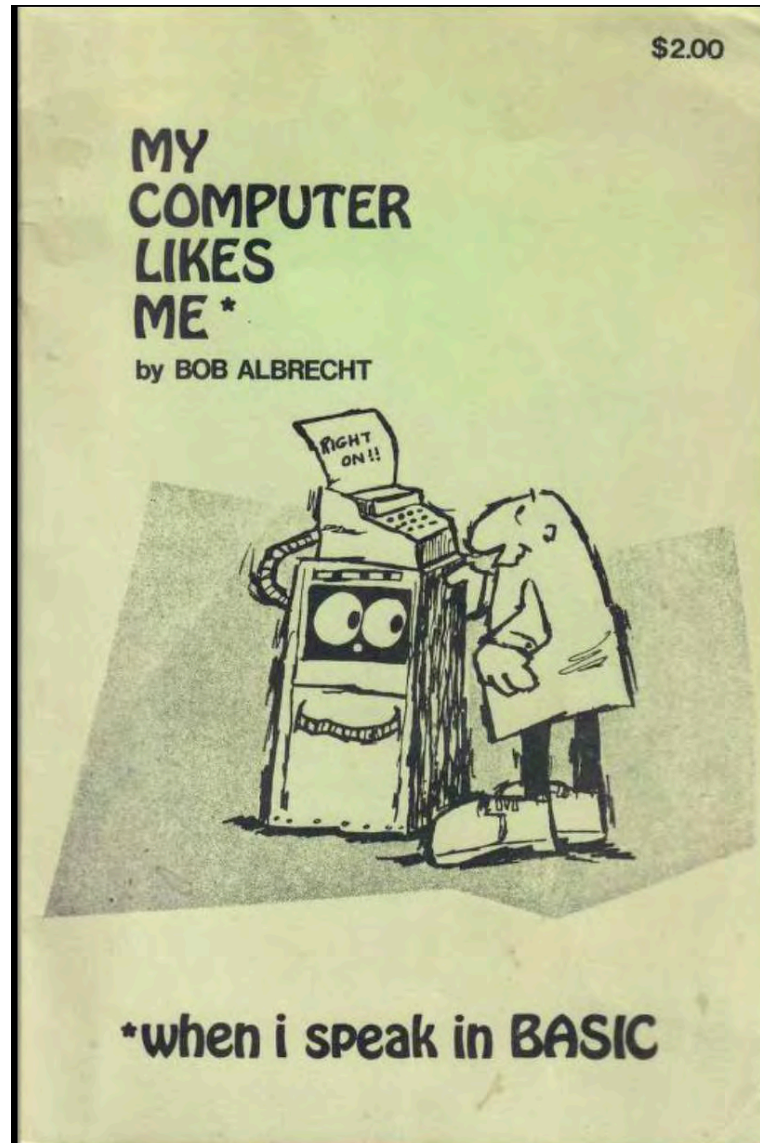
GPTs and LLMs are *transformative technologies*, analogous to the advent of the electronic, programmable digital computer



UD's Bendix G-15-D, ca. 1958

Scientists and engineers rapidly adopted digital computers to numerically solve difficult problems

Since 1972 (and earlier)...

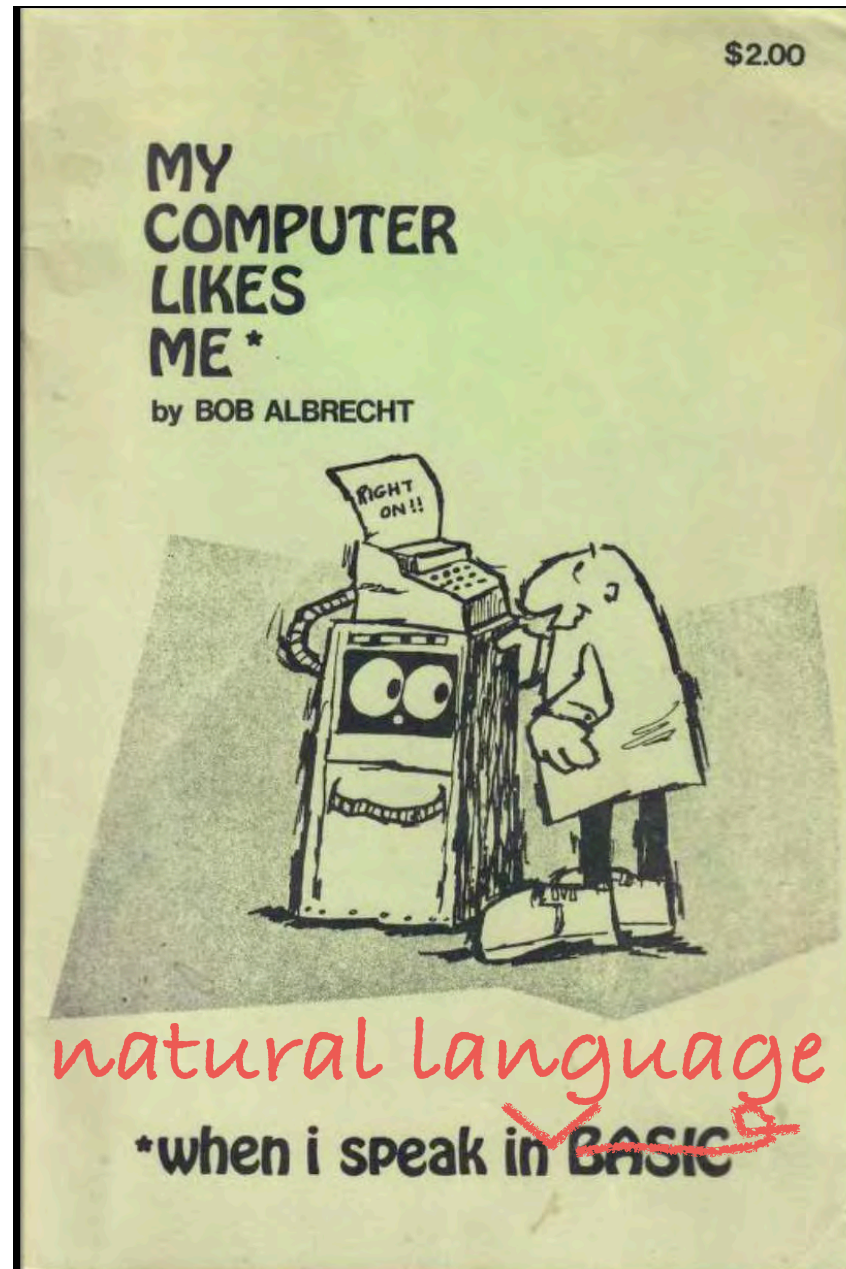


or...

assembly
FORTRAN
C
Matlab
Python

Javascript, Pascal, Lisp, ALGOL 60,
BAL, JCL, Smalltalk, PL/I, Logo...

suddenly,
today...

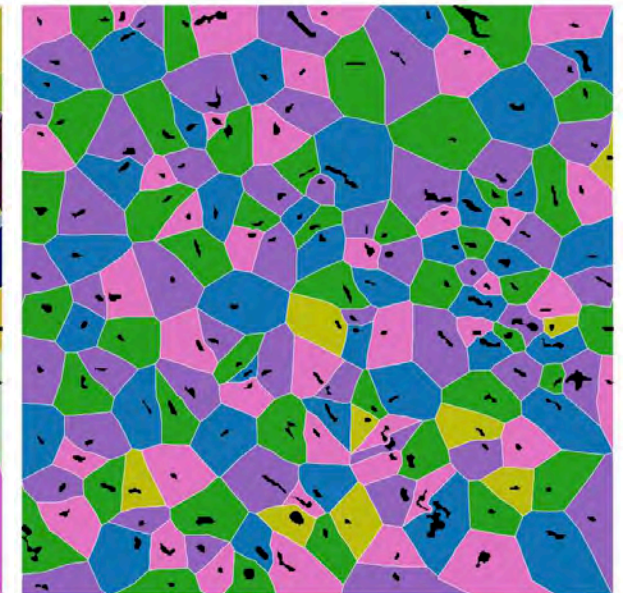
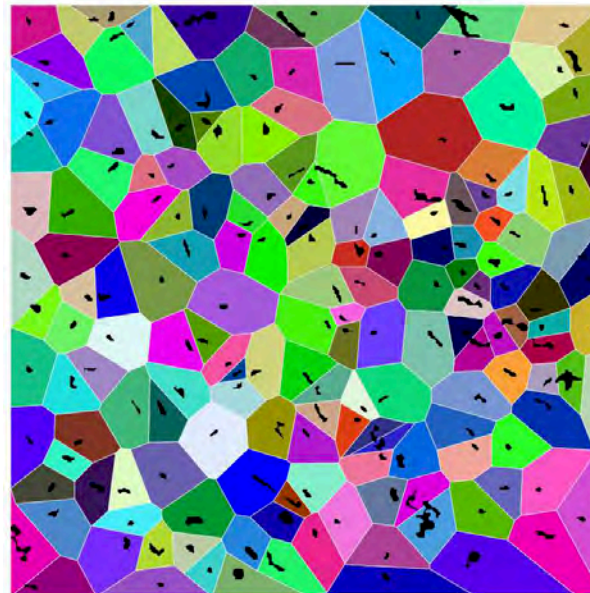
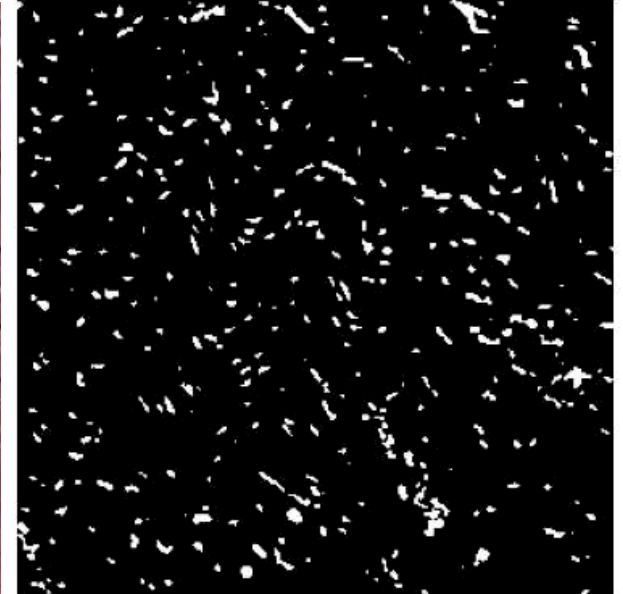
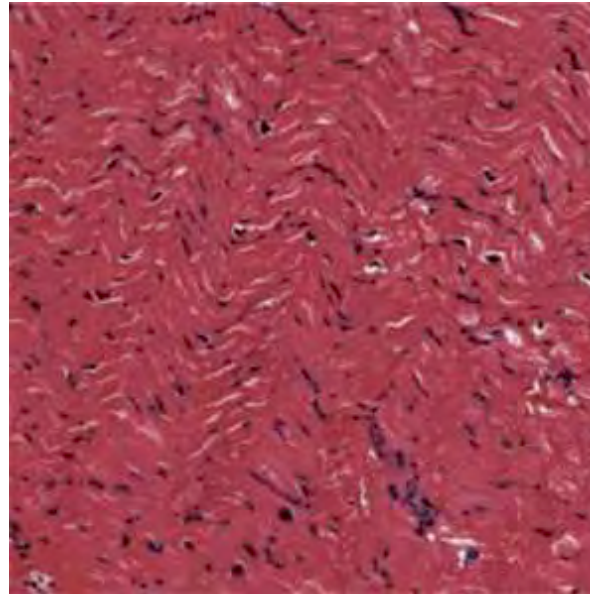


Solving non-trivial, but adjacent problems

Minimum color Voronoi representation



Jason Conradt



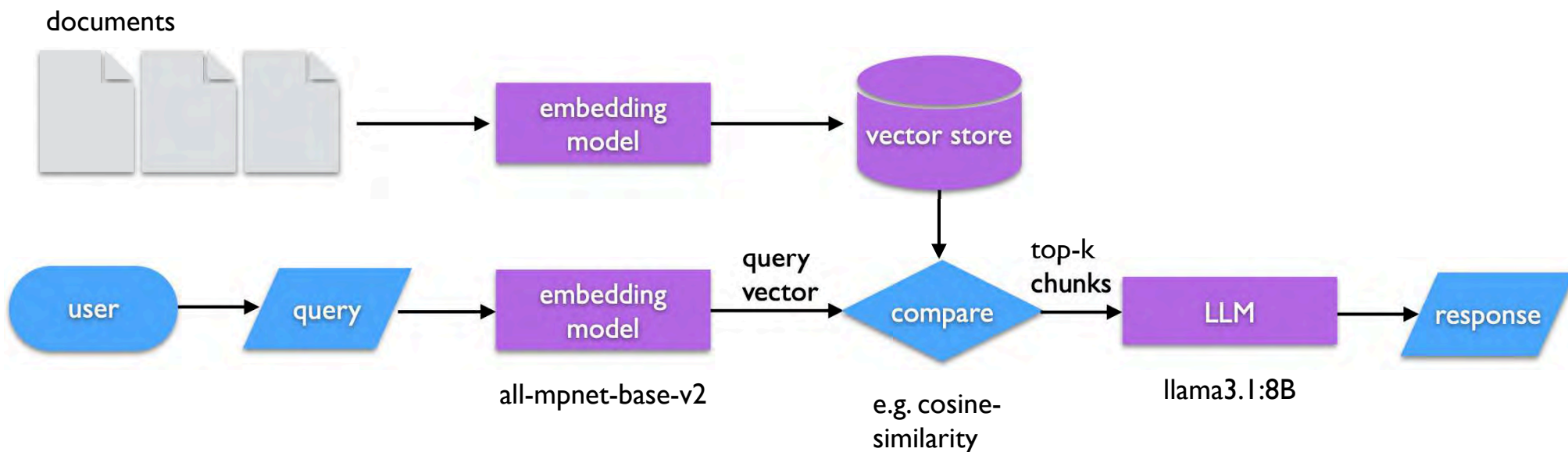
Use ChatGPT to help write Python code for binarization and Voronoi tessellation... in a day instead of a week +

Retrieval Augmented Generation (RAG)

Pull in unstructured text, PDFs, documents, webpages and more and index the data within them

The simplest queries involve either semantic search or summarization –

- Semantic search: Query about specific information in document(s) that matches the query terms and/or semantic intent
- Summarization: condensing a large amount of data into a short summary relevant to your current question



semanticsearch_local — furst@anisotropic: ~ — emacs -nw query.py — 129x49

File Edit Options Buffers Tools Python Help

```
# query.py
# Run a query on a vector store
#
# E.M.F. July 2025

from llama_index.core import (
    StorageContext,
    load_index_from_storage,
    ServiceContext,
    Settings,
)
from llama_index.embeddings.huggingface import HuggingFaceEmbedding
from llama_index.llms.ollama import Ollama

# Use a local model to generate
Settings.llm = Ollama(
    model="llama3.1:8B",
    request_timeout=360.0,
    context_window=8000
)

def main():
    # Load embedding model (same as used for vector store)
    embed_model = HuggingFaceEmbedding(model_name="all-mpnet-base-v2")
    Settings.embed_model = embed_model

    # Load persisted vector store + metadata
    storage_context = StorageContext.from_defaults(persist_dir="./storage")
    index = load_index_from_storage(storage_context)

    query_engine = index.as_query_engine()

    # Query
    while True:
        q = input("\nEnter your question (or 'exit'): ").strip()
        if q.lower() in ("exit", "quit"):
            break
        print()
        response = query_engine.query(q)
        print(response)

if __name__ == "__main__":
    main()
```

--UU--:--- F1 query.py Top L44 (Python ElDoc) -----

*RAG query in
~40 lines of code*

See <https://docs.llamaindex.ai/>

Possible LLM pitfalls and problems

Hallucinations

LLMs sometimes generate plausible-sounding but factually incorrect or entirely fabricated information

Misalignment

The model's output may not reflect user intent or ethical/social expectations, especially in nuanced or sensitive contexts.

Bias and Stereotypes

LLMs can amplify societal biases, especially around race, gender, or culture, due to patterns in their training data.

Overconfidence

LLMs often present answers in a confident tone regardless of uncertainty, leading users to overtrust incorrect outputs.

Context Loss

When prompts are too long or complex, LLMs may forget or ignore earlier context, causing incoherent or inconsistent responses.

In our uses...

Hallucinated Code or APIs

Suggests functions, libraries, or syntax that don't exist

Misuse of Retrieved Content (RAG)

Pulls in irrelevant or misunderstood snippets, leading to wrong answers

Inconsistent Formatting

Breaks expected structure (e.g. JSON, YAML, code blocks) even in repeated tasks

Overconfident Errors

Gives incorrect results or logic with a confident tone

Context Limitations

Forgets earlier input in long prompts or complex workflows

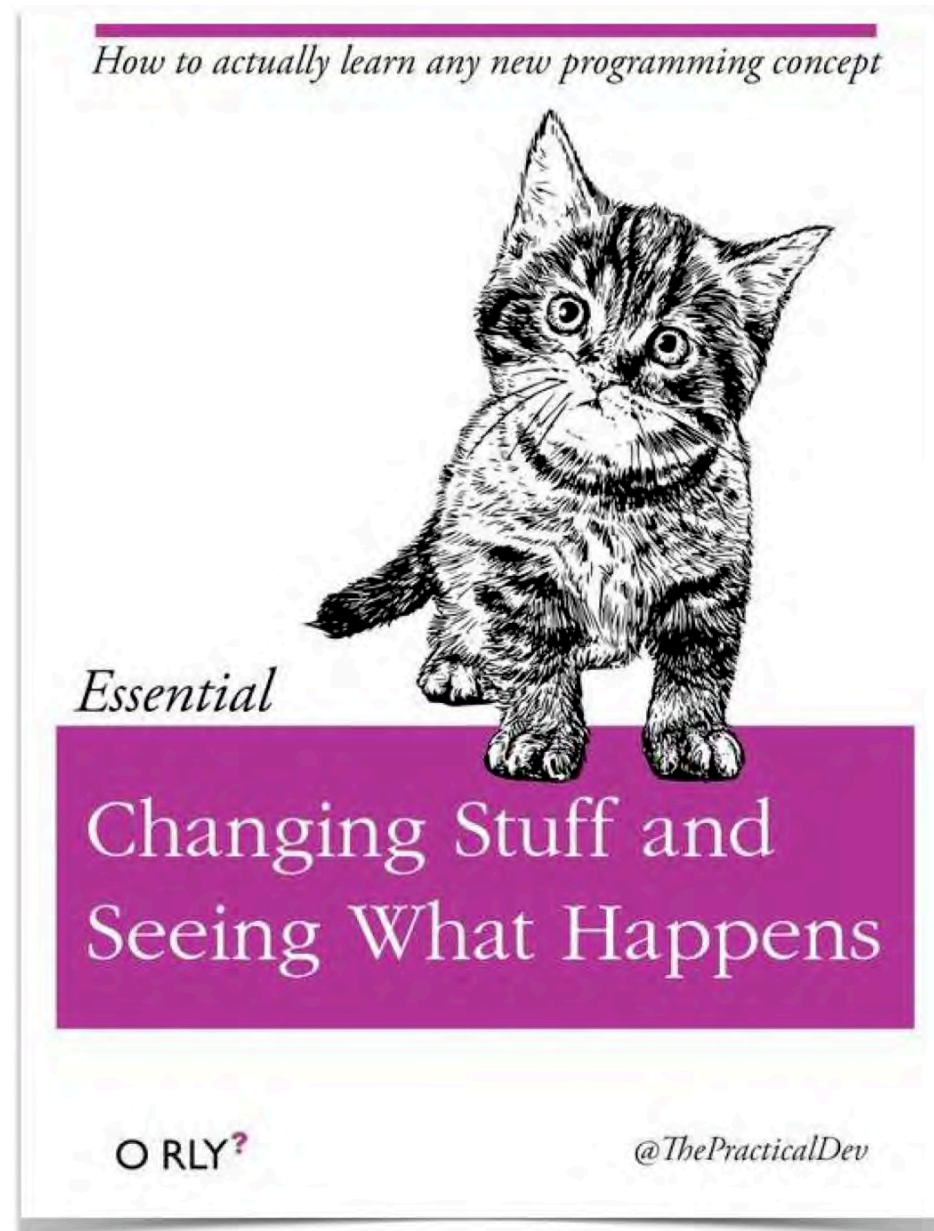
“...all people using these systems should be informed and constantly reminded that their requests could lead to error-prone, fabricated, or otherwise misaligned responses* as well as potentially dangerous actions in the case of AI agents.”

Why AI Chatbots Lie to Us

Melanie Mitchell, Science, July 26, 2025

<https://doi.org/10.1126/science.aea3922>

* may also apply to faculty



"Programming is a skill best acquired by practice and example rather than from books."
— Alan Turing, *Programmers' Handbook for Manchester Electronic Computer Mark II*, 1951

Learn by doing – hack on an LLM

HTTPS://WWW.YOUTUBE.COM/WATCH?V=KCC8FMEB1NY



Andrej Karpathy

@AndrejKarpathy · 605K subscribers · 16 videos

SuperThanks: very optional, goes to Eureka Labs. ...more

eurekalabs.ai and 4 more links

Subscribe

Home Videos Playlists Community

General Audience ▶ Play all

videos for more general audience, no programming experience necessary.



[1hr Talk] Intro to Large Language Models

Andrej Karpathy
2.4M views · 1 year ago

Neural Networks

```
furst@anisotropic:~/LLM/nanoGPT$ wc $(ls -1 *.py)
117   487   4815 bench.py
 47   219   1758 configurator.py
331  1798 16507 model.py
 93   522   4313 sample.py
336  1799 14845 train.py
924  4825 42238 total
```

LET'S BUILD GPT.
FROM SCRATCH.
IN CODE.
SPELLED OUT.



1:56:20

Let's build GPT: from scratch, in code, spelled out.

5M views · 2 years ago

~300 lines of Python
Runs on CPU or GPU

lem.che.udel.edu

Home
New! LLM2025
 Research topics
 Publications
 Alumni
 Microrheology
 Societies & Meetings
 Code & tools
 Colloids 2022

Internal wiki
 Group Thesis Repo
 Group logistics
 Laboratory safety
 Laboratory resources
 Research resources
 Computing resources
 Writing resources
 Wiki tips and links

External links
 Furst Department page

edit SideBar

[Main /](#)
LLM2025
2025 Winter Research Review Tech Talk... and More!
Always Under Construction!

Presented on January 22, 2025 as the lunch talk for the Department of Chemical and Biomolecular Engineering Winter Research Review. In my talk, I discussed uses of large language models (LLMs), the underlying architecture of a generative pre-trained transformer (GPT), and basic aspects of the mechanics behind training and deploying LLMs.

[FURST_WRR_2025.pdf](#) (65M)

This was on my mind: that engineers don't take technical solutions for granted. We generally like to "look under the hood" and see how things work. So, if you are interested in learning more about the technical underpinnings of LLMs, this page collects a few resources. The talk is largely inspired by the rapid adoption of LLMs to help us solve difficult but *adjacent* problems in our research.

In my talk, I didn't get into the details of how one goes from the single attention mechanism to "multi-head" attention, wh important feature of LLM models. I also did not emphasize the step of fine-tuning models and how the basic generative GPT is built up into a powerful chatbot that many of us use. Those are topics worth exploring in greater depth.

Overall, I view GPTs as a transformative technology, in many ways analogous to the disruption that came with the introc electronic general programmable computer in the late 1940s. Chemical engineers rapidly adopted that earlier technolog challenging modeling problems -- giving them solutions to partial differential equations and systems of these equations, were intractable or extremely inefficient before the development of machine computing. Specifically, LLMs will give us n our computational tools through natural language, help us to rapidly come up to speed in a new area, or quickly develop models and data with code.

- Eric Furst

Contents [hide]
[2025 Winter Research Review Tech Talk... and More!](#)
[Essential background](#)
[TL;DR](#)
[Referenced in the talk](#)
[Access models and tools](#)
[More learning resources](#)
[Related reading](#)
[Ethics of AI](#)
[Problems and pitfalls in generative AI](#)
[History of computers and computational tools](#)
[Artistic and literary practices](#)
[More links](#)
[CHEG 667-013](#)

CHEG667-013
handouts!

CHEG 667-013 - CHEMICAL ENGINEERING WITH COMPUTERS
 Department of Chemical and Biomolecular Engineering
 University of Delaware

Spring 2025

LARGE LANGUAGE MODELS PART II

Key idea:

Learn how to run LLMs locally without a cloud-based API

Key goals:

- Learn about `ollama` and `llama.cpp`
- Run higher performance LLMs locally on a laptop or desktop computer

Our work with LLMs so far focused on `nanoGPT`, a python-based code that can train and run inference on a simple GPT implementation. In this handout, we will explore running something between it and API-based models like ChatGPT. Specifically, we will try `ollama`. This is a local runtime environment and model manager that is designed to make it easy to run and interact with LLMs on your own machine. `Ollama` and another environment, `llama.cpp`, are programs primarily targeted at developers, researchers, and hobbyists who want to access LLMs to build and experiment with but don't want to rely on cloud-based APIs.¹

`Ollama` is written in Go and `llama.cpp` is a C++ library for running LLMs. Both are cross-platform and can be run on Linux, Windows, and macOS. `llama.cpp` is a bit lower-level with more control over loading models, quantization, memory usage, batching, and token streaming.

Both tools support a GGUF model format. This is a format suitable for running models efficiently on CPUs and lower-end GPUs. GGUF is a versioned binary specification that embeds the

- Model weights (possibly quantized);
- Tokenizer configuration and vocabulary (remember, in `nanoGPT`, we used a character-level tokenization scheme);
- Metadata such as the author, model description, and training parameters;
- Special tokens like `<bos>`, `<eos>`, and `<unk>`.

Here, quantization refers to how model weights are stored. Instead of using high precision 32-bit full-precision floating point numbers (FP32), it may store the weights as lower precision numbers: half precision (FP16), 8-bit integers (INT8), or even 4-bit values (Q4_0). Using lower precision representations saves space (memory) and can speed the inference calculations. In a model, the speed and accuracy are balanced with the choice of quantization and the size of the embedding vector.

Let's get started! We will download `ollama` and run a few models in this tutorial.

¹An API (Application Programming Interface) is a set of defined rules that enables different software systems, such as websites or applications, to communicate with each other and share data in a structured way.

<https://furst.group>

Concluding remarks

1. How LLMs work
2. Running LLMs locally
3. Two (or three) LLM uses

Solving *adjacent* problems in our science:
coding, data cleaning, formatting

New ways to build computational tools with computer-human interactions through natural language

RAG (search, information retrieval) that uses the encoding / embedding power of transformers