# TAMING LARGE LANGUAGE MODELS

WITH LANGCHAIN

# Who am I

Senior Data Scientist; DS & AI SME @ Growth Acceleration Partners

BSc Actuarial Sciences; Universidad de Costa Rica

MSc Data Science & Artificial Intelligence; Data ScienceTech Institute

Dozens of ML Developments: finance, healthcare

Passionate about ethics in machine learning and artificial intelligence

Data & AI Fest 2024 — we are GAP

Congratulations, sailor, you made it through the GenAI hype!

**Hype Cycle for Artificial Intelligence, 2023**

Expectations / Time

Smart Robots
Responsible AI
Neuromorphic Computing
Prompt Engineering
Artificial General Intelligence
Decision Intelligence
AI TRiSM
Operational AI Systems
Composite AI
Data-Centric AI
AI Engineering
AI Simulation
Causal AI
Neuro-Symbolic AI
Multiagent Systems
First-Principles AI
Automatic Systems

Generative AI
Foundation Models
Synthetic Data
ModelOps
EdgeAI
Knowledge Graphs
AI Maker and Teaching Kits
Autonomous Vehicles
Cloud AI Services
Data Labeling and Annotation
Intelligent Applications
Computer Vision

Innovation Trigger | Peak of Inflated Expectations | Trough of Disillusionment | Slope of Enlightenment | Plateau of Productivity

Plateau will be reached:
○ less than 2 years
● 2 to 5 years
● 5 to 10 years
▲ more than 10 years
⊗ obsolete before plateau
As of July 2023

gartner.com
Source: Gartner
© 2023 Gartner, Inc. and/or its affiliates. All rights reserved. 2079794
Gartner.

Exclusive
Generative AI Providers Quietly Tamp Down Expectations
By Anissa Gardizy and Aaron Holmes
Share ⌄
Mar 12, 2024, 7:00am PDT

Source: theinformation.com

Source: Gartner
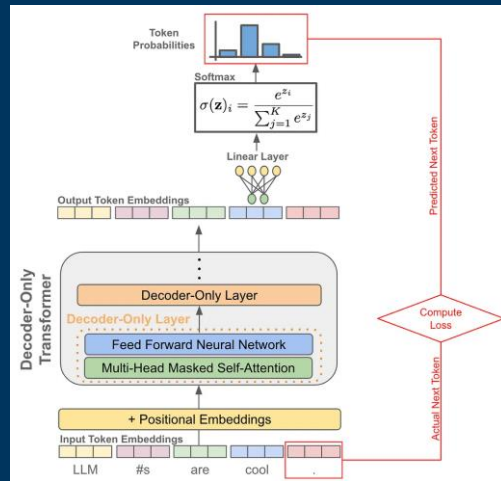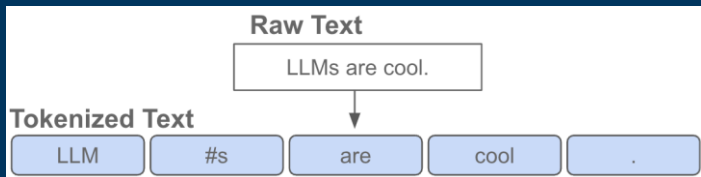
# Language Models?

A language model is a probabilistic model of natural language.

P(word) = 0.843521

For this talk, word ≈ token
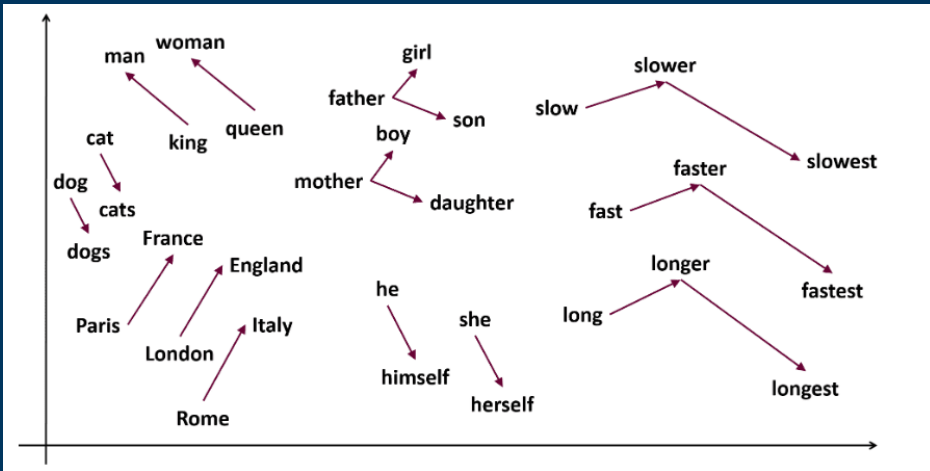
Language models are next token predictors!

# How do they handle text?

**Word Embedding**

Idea: "two words are similar if they are used in similar context".

⚠️ **This is fundamentally different from human language assimilation.**



Source: paperswithcode.com

The impressive thing about AI is not that it can learn and understand as we humans do.

It's how much it can do while not truly understanding.

# Large Language Models

- Large!

- Competent at general-purpose language generation

| AlexNet<br>(2012) | 60,000,000 |
|---|---|
| ElMo<br>(2016) | 94,000,000 |
| Megatron-Turing NLG<br>(2021) | 530,000,000,000 |
| GPT-4<br>(2023) | 100,000,000,000,000 |

# What ARE they learning?



Source: Hackernoon Latent Space Visualization

**Before building an LLM solution**

1. Does this opportunity need AI/ML?
2. Can it be done with traditional ML?
3. Test your idea on ChatGPT/ HuggingChat or other online service:
    a. Can I get it to do what I want?
    b. Can I get it to do things wrong? (spend at least 2 hours on this!!)

# Why?

# LangChain

"Framework designed to simplify the creation and integration of applications using large language models."

- Model agnostic (GPT, Google, Meta, Mistral, etc.)

- Extends and integrates LLM use with other tools

- Incredibly active development



March 5, 2024 – March 12, 2024    Period: 1 week ▾

**Overview**

246 Active pull requests     88 Active issues

169 Merged pull requests   77 Open pull requests   46 Closed issues   42 New issues

Excluding merges, **72 authors** have pushed **169 commits** to master and **215 commits** to all branches. On master, **551 files** have changed and there have been **32,505 additions** and **33,663 deletions**.

169 Pull requests merged by 62 people

Source: github.com

# LangChain Modules

- Models (OpenAI, Huggingface)

- Prompts

- Chains to orchestrate multi-call LLM processes

- Retrieval to enhance responses with context

- Memory to help the model keep track of previous interactions within the conversations

- Agents and tools (Google search, Wikipedia, Calculator, much more)

# LangChain Modules

- Models (OpenAI, Huggingface)

- Prompts

- Chains to orchestrate multi-call LLM processes

- Retrieval to enhance responses with context

- Memory to help the model keep track of previous interactions within the conversations

- Agents and tools (Google search, Wikipedia, Calculator, much more)

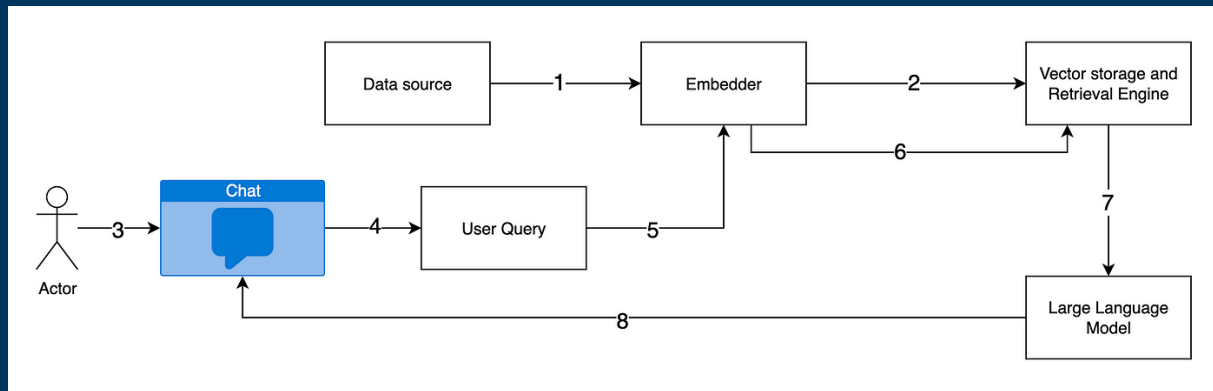# Demo

Code and slides: github.com/pheymanss/taming-llms

Other challenges of working with LLMs

# LLMs are not reliable data sources

"Look it up on ChatGPT" ❌ ❌ ❌ ❌

Use Retrieval Augmented Generation (RAG) or fine-tuning.



Credit: Vikesh Pandey

# LLMs encode and reproduce biases



OPENAI'S GPT IS A RECRUITER'S DREAM TOOL. TESTS SHOW THERE'S RACIAL BIAS

Recruiters are eager to use generative AI, but a Bloomberg experiment found bias against job candidates based on their names alone

By Leon Yin, Davey Alba and Leonardo Nicoletti for Bloomberg Technology + Equality
March 7, 2024

# Societal bias

The data accurately reflects
true biases of the world:
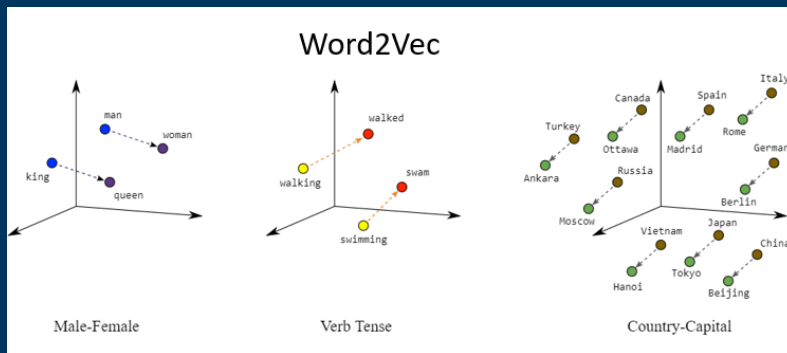- Justice system
- Hiring
- Healthcare
- Public policy
- Text!



Word2Vec

Male-Female    Verb Tense    Country-Capital

Google News

$$\overrightarrow{man} - \overrightarrow{woman} \approx \overrightarrow{computer\ programmer} - \overrightarrow{homemaker}.$$

Bolukbasi, T., Chang, K., Zou, J., Saligrama, V. & Kalai, A. (2016) *Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings*. arxiv.org/abs/1607.06520

# Q&A

Code and slides: github.com/pheymanss/taming-llms