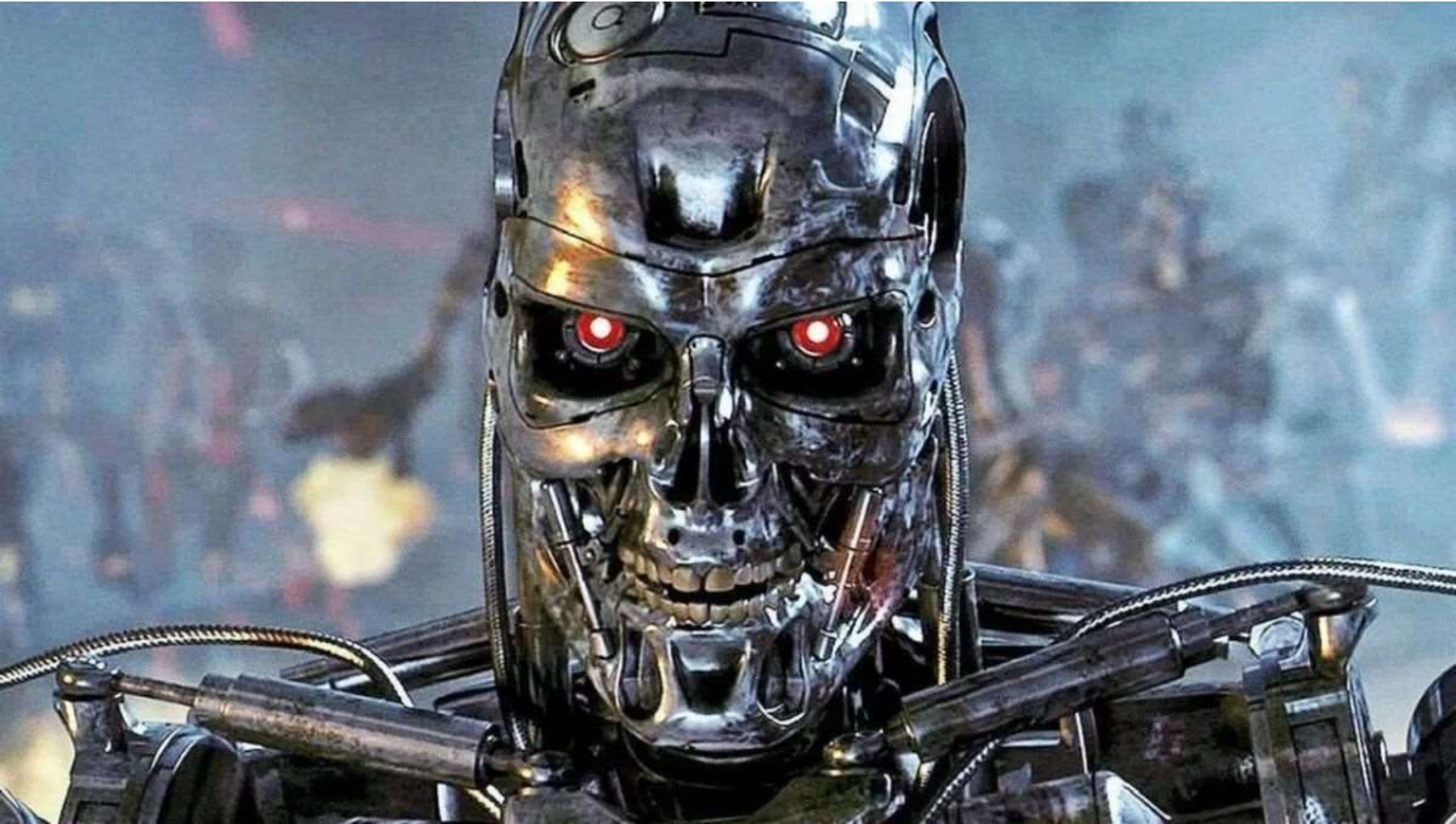


# Artificial Intelligence is here



October 1984

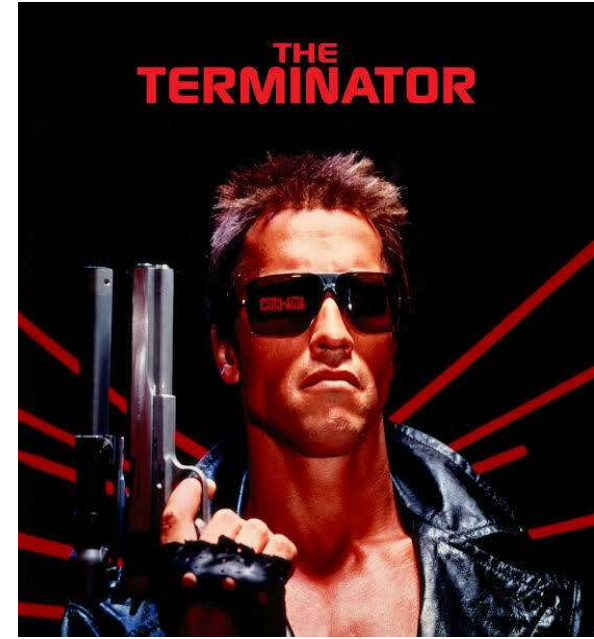
# Artificial Intelligence is here



October 1984



# Artificial Intelligence is here



## What Large Language Models (LLM) are?

- LLMs can be understood as advanced tools built on the principles of **statistical pattern recognition and prediction**.
- LLMs are designed to **predict the next most probable word ("token" )** in a sequence.
- A **"token" is the fundamental unit of text**. It can be a word, a character, or even a punctuation mark. On **average a single word translates to about 0.75 tokens**.

## How Do Large Language Models (LLMs) Work?

- The term "**sequence**" refers to the context or the "**window**" of text that the model considers when making its predictions. This could be a single sentence, a paragraph, or even a longer body of text like a book chapter.
- Models like ChatGPT-03, the maximum sequence length is 4096 tokens (**July 2022**), which is equivalent to a few pages of text.

# Introduction to Large Language

Model Version	Maximum Tokens (Context Window)	Notes	Year
GPT-3.5	4,096 – 16,385 tokens	Base model supports 4,096 tokens;	2022
GPT-4	8,192 – 32,768 tokens	Available in 8K and 32K token versions.	2023
GPT-4 Turbo	128,000 tokens	Enhanced version with a significantly larger context window.	May 2024
GPT-4o	128,000 tokens	Multimodal model supporting text, audio, and vision inputs.	Dec 2024
GPT-4.1	1,000,000 tokens	Extensive context window, suitable for processing large datasets.	Feb 2025
GTP-5.0	400,000 tokens	improved reasoning and reduced hallucinations	Aug 2025

## How a LLM problem is modelled?

Basically, a LLM is aimed to answer the question: **What is  $p(\text{text})$**

Given a sequence of tokens:  $(x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)})$  (1)

Then:

$P(x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(N)})$  is the probability of sequence (1)

For a sequence of **three tokens** the probability is equal to

$$P(x^{(1)}, x^{(2)}, x^{(3)}) = p(x^{(1)})p(x^{(2)} | x^{(1)})p(x^{(3)} | x^{(1)} x^{(2)})$$

# How a LLM problem is modelled?

The general case is defined as follows

$$P(x^{(1)}, \dots, x^{(N)}) = \prod_{i=1}^N p(x^{(i)} \mid x^{(1)}, \dots, x^{(i-1)})$$



al nassr jersey ronaldo



al nassr jersey ronaldo **kids**

al nassr jersey ronaldo **original**

al nassr jersey ronaldo **2025**

al nassr jersey ronaldo **price**

al nassr jersey ronaldo **youth**

al nassr jersey ronaldo

al nassr jersey ronaldo **nike**

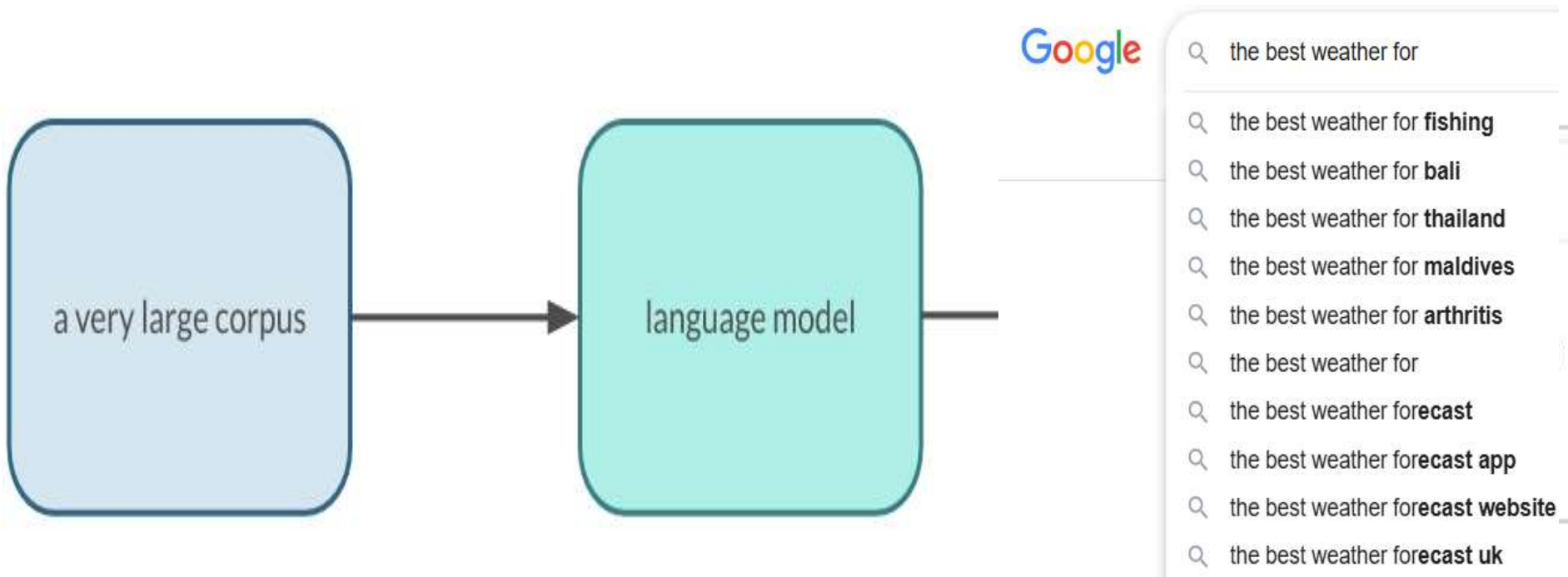
al nassr jersey ronaldo **blue**

al nassr jersey ronaldo **away**

al nassr jersey ronaldo **adidas**



# How a LLM problem is modelled?



"How are you this afternoon? Has your car been broken?"  $\rightarrow P(10^{-15})$

because the second sentence is **less typical or more unexpected** in a daily basis conversation.

"How are you this afternoon? It's good to see you"  $\rightarrow P(10^{-5})$

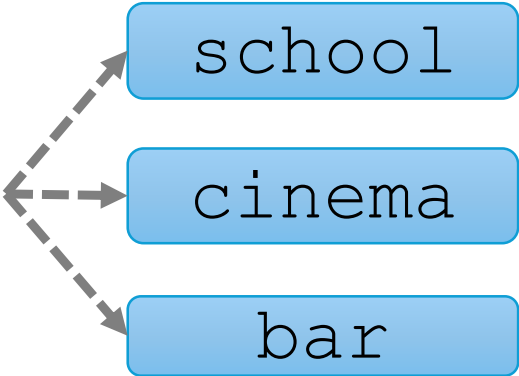
Because **it is more frequent of expected in daily bases conversations**, leading to a higher probability.

## Chain rule

[https://en.wikipedia.org/wiki/Chain\\_rule\\_\(probability\)](https://en.wikipedia.org/wiki/Chain_rule_(probability))

$$P(x^{(1)}, \dots, x^{(N)}) = \prod_{i=1}^N p(x^{(i)} \mid \underbrace{x^{(1)}, \dots, x^{(i-1)}}_{\text{Context}})$$

The students went  
to \_\_\_\_\_



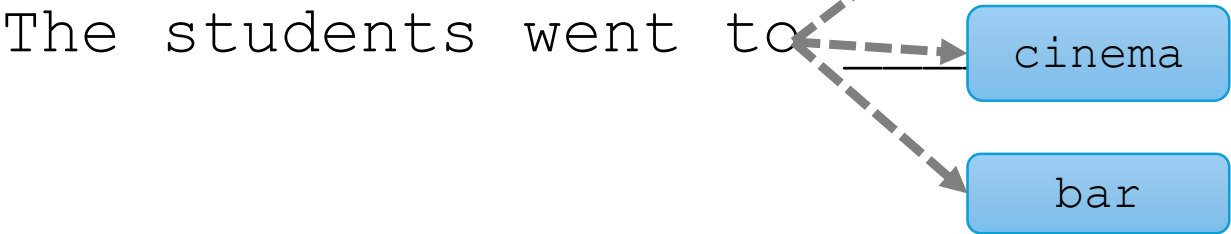
# Introduction to Large Language

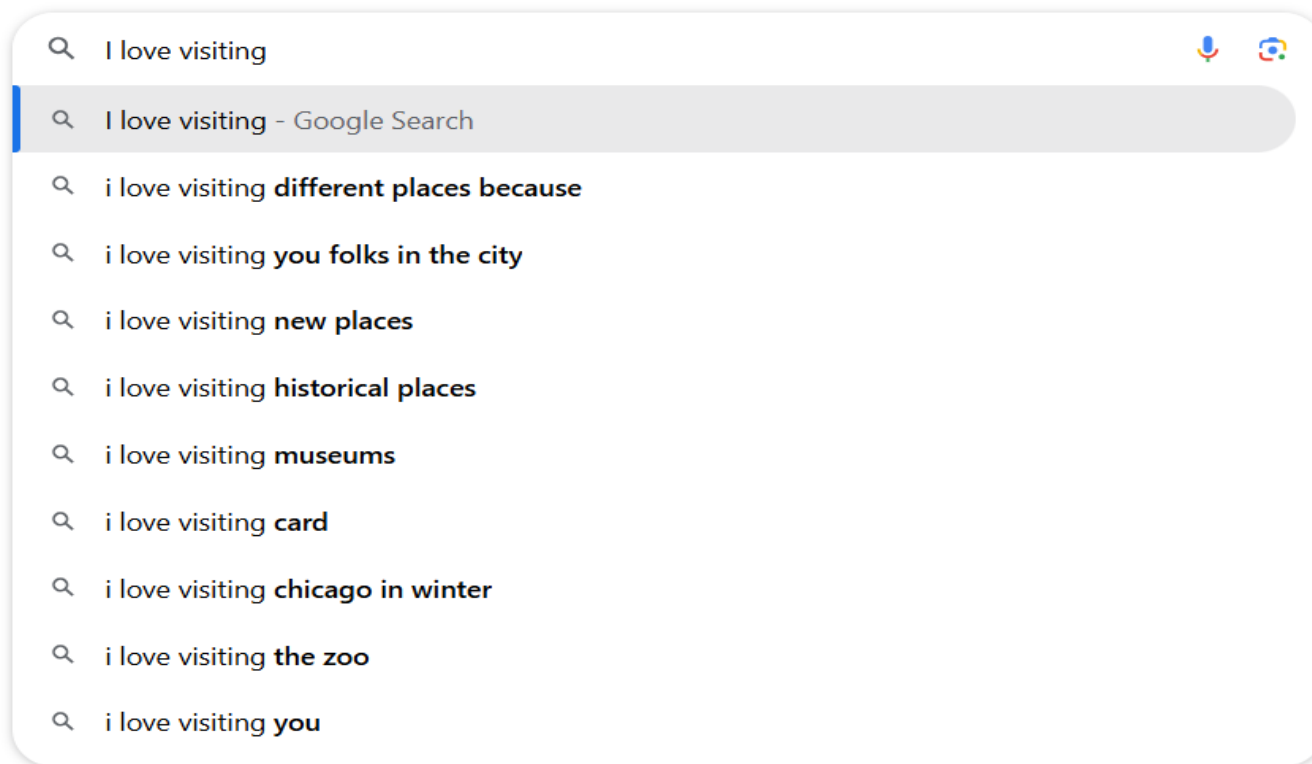
Every time a new token is added to the context, then a new probability distribution is calculated

The

The students \_\_\_\_\_

The students went \_\_\_\_\_

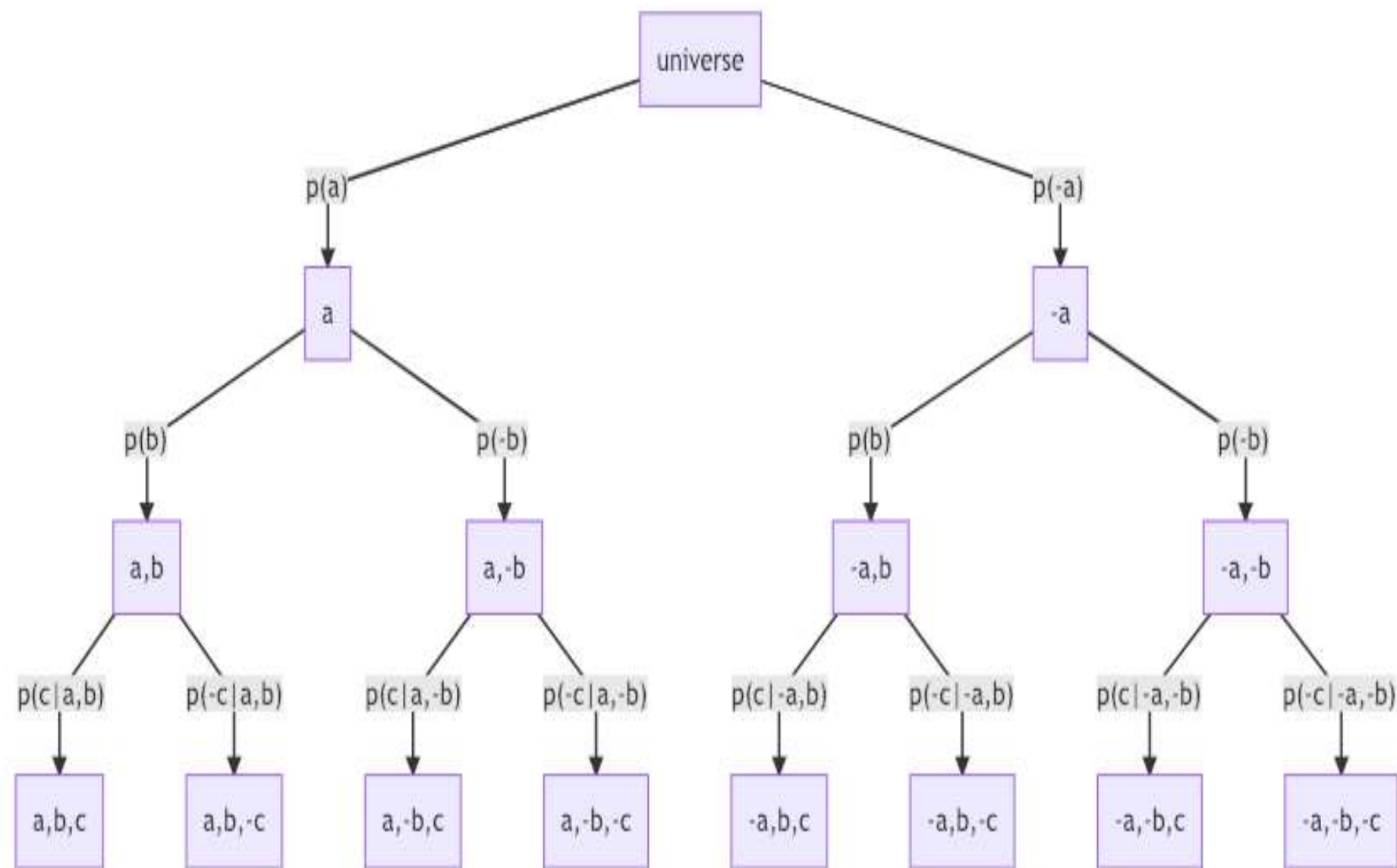




$$P(x^{(1)}, \dots, x^{(N)}) = \prod_{i=1}^N p(x^{(i)} \mid x^{(1)}, \dots, x^{(i-1)})$$



# Visual representation of the





How does

a

Language

Model

“learn”?



I have a duck whose name is Rocky. I have two cats. They like playing with Rocky.

- Corpus size: 17
- $P(\text{Rocky}) = 2/17$
- $P(\text{cats}) = 1/17$

I have a dog whose name is Rocky. I have two cats, they like playing with Rocky.

*Bigram probability (based on this corpus).*

$$P(A | B) = \frac{P(A,B)}{P(B)}$$

$$P(\text{have} | \text{I}) = \frac{P(\text{I have})}{P(\text{I})} = \frac{2}{2} = 1$$

$$P(\text{two} | \text{have}) = \frac{P(\text{have two})}{P(\text{have})} = \frac{1}{2} = 0.5$$

$$P(\text{eating} | \text{have}) = \frac{P(\text{have eating})}{P(\text{have})} = \frac{0}{2} = 0$$

$$P(w_2 | w_1) = \frac{C(w_1, w_2)}{\sum_w C(w_1, w)} = \frac{C(w_1, w_2)}{C(w_1)}$$



I have a dog whose name is Rocky. I have two cats, they like playing with Rocky.

*Trigram probability (based on this corpus).*

$$P(A | B) = \frac{P(A,B)}{P(B)}$$

$$P(a | \text{I have}) = \frac{C(\text{I have a})}{C(\text{I have})} = \frac{1}{2} = 0.5$$

$$P(\text{several} | \text{I have}) = \frac{C(\text{I have several})}{C(\text{I have})} = \frac{0}{2} = 0$$

$$P(w_3 | w_1 w_2) = \frac{C(w_1, w_2, w_3)}{\sum_w C(w_1, w_2, w)} = \frac{C(w_1, w_2, w_3)}{C(w_1, w_2)}$$

I have a dog whose name is Rocky. I have two cats, they like playing with Rocky.



*N-gram probability (based on this corpus).*

$$P(A \mid B) = \frac{P(A,B)}{P(B)}$$


$$P(w_i \mid w_1, w_2, \dots, w_{i-1}) = \frac{C(w_1, w_2, \dots, w_{i-1}, w_i)}{C(w_1, w_2, \dots, w_{i-1})}$$





- Prompt engineering means tailoring your questions and input so you can get the most out of an LLM
- The output of a LLM is determined by both what the **system has been trained** on and **what information you give it**
- Prompts can take many forms, from instructing the LLM to take on a role (e.g. a helpful teacher, a pirate) or guiding the way it should process its output (e.g. "chain of thought" or a particular method)





ChatGPT 5 ▾


 New chat


 Search chats


 Library


 Codex


 Sora


 GPTs


 AI PDF Drive: Chat, Create, Org...


 New project


 Applications

 QS Rankings

 Predicting citations

 Values trades off

 Labs\_Medical\_exams


 See more

Chats

Update model table

Comment multiple lines RStudio

OECD investment data



Igor Barahona

Unlike the API, ChatGPT still caps context based on these tiers. [Horror](#)

 Share

## Summary

- **API (all GPT-4.1 variants):** Up to **1,000,000 tokens** context window.
- **ChatGPT UI:** Tier-based limits—8K (Free), 32K (Plus/Team), 128K (Pro/Enterprise).

Let me know if you'd like specifics on usage per variant (e.g., mini or nano), or how this compares to other models like GPT-5!

      Sources

GPT-5.0 Maximum Tokens (Context Window)?

Here's the latest on **GPT-5's maximum context window**, as of early September 2025:

## GPT-5 Context Window Specs

### 1. API (Developers & Businesses)

- Up to **400,000 tokens total**, typically comprising:
  - ~272,000 tokens for input, and
  - ~128,000 tokens for output

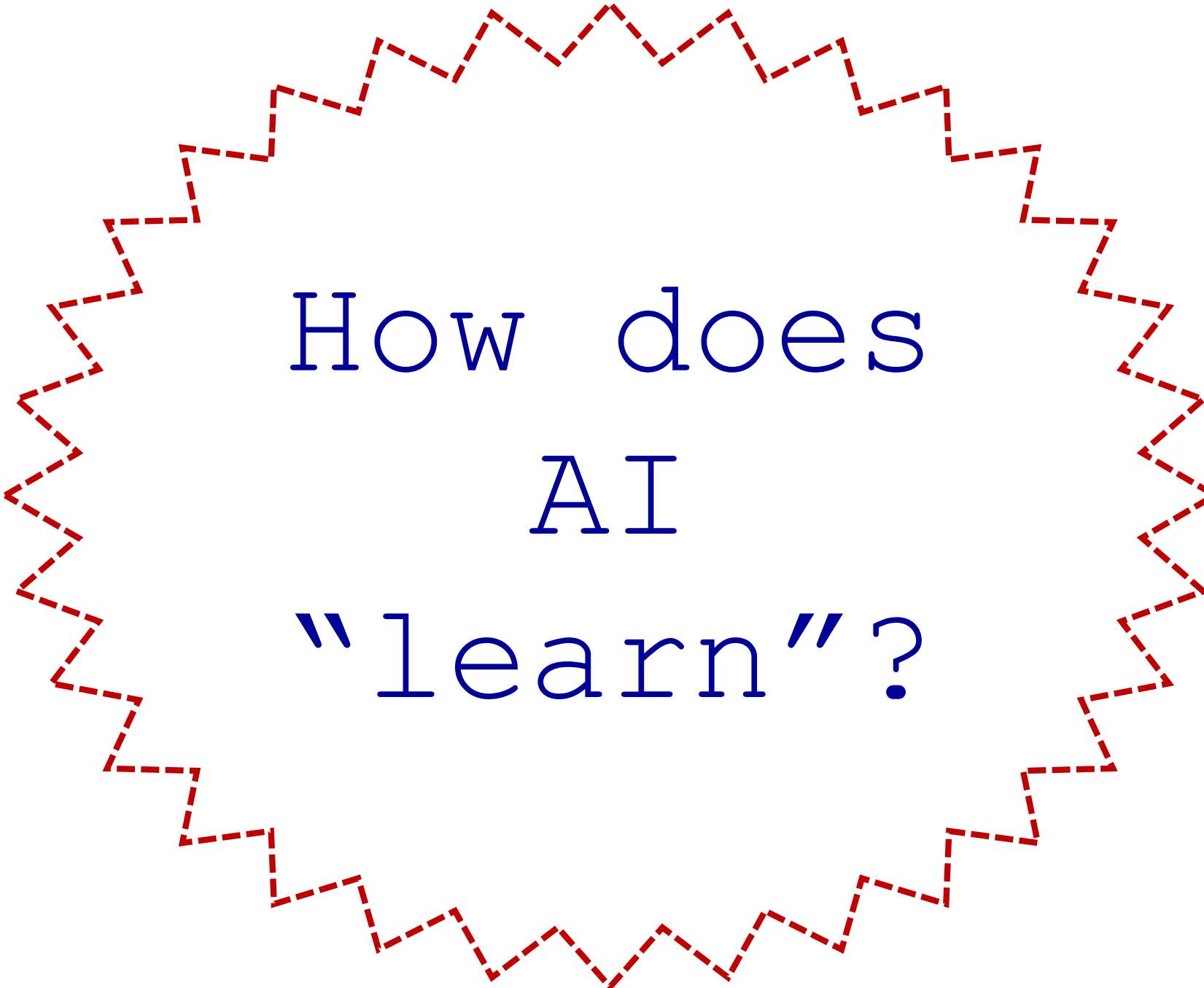
These figures are confirmed in multiple developer-focused sources [All Things Hw](#) [AI Agents for Customer Service](#)



### 2. ChatGPT (User Interface)

+ Ask anything





How does  
AI  
“learn”?





## Sustainability (vs) Cost

What do you guys  
might select?



## VALUE TRADE-OFFS IN SOURCING DECISIONS



OPERATIONS  
MANAGER

**How operations managers make value trade-offs decisions that involve sustainability and cost?**

**How elicitation methods impact common biases, such as sensitivity?**



# Motivations

**Are managers capable to make consistent value trade-offs that reflect rational prioritizations?**

**Are importance weights and trade-offs considered consistent when attribute ranges are explicitly specified?**





Are managers' prioritizations between sustainability and cost conditioned on whether:

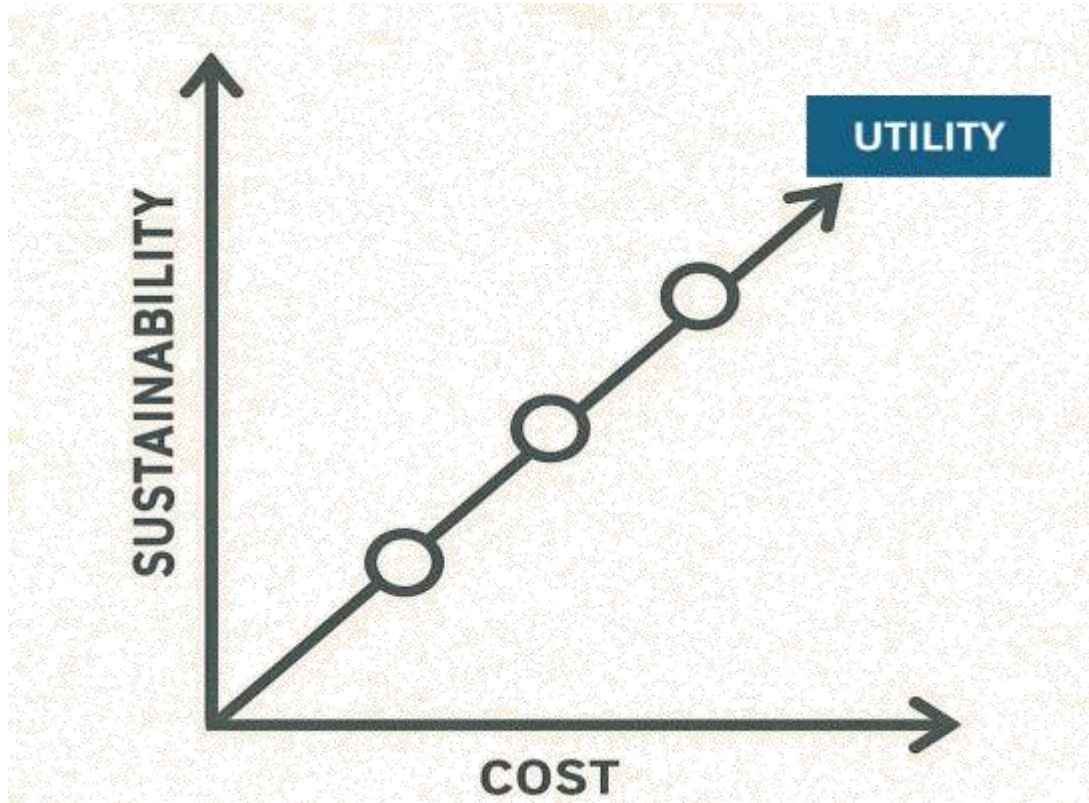
**a) Personal preferences or**

**b) Ranges of the underlying attributes.**



# Introduction

A naïve perspective assumes a linear relationship between **Sustainability** and **Cost**



Our experiments show that the relationship is far from linear

# Methodology



## COGNITIVE REFLECTION TEST

- CRT assesses a person's tendency to override an initial intuitive (but wrong) answer and engage in deeper, reflective reasoning.
- The classic CRT consists of 3 questions (Frederick, 2005), but has been expanded in later versions to address issues like low score variability.
- Performance can be influenced by numerical skills, leading to concerns about confusing effects in interpretation.

The organization has set up two strategic goals

- 1) Minimize emissions of the car fleet
- 2) Minimize acquisition costs of the fleet.



Minimizing  
emissions  
of the car fleet



Minimizing  
acquisition  
costs of the  
fleet



The range of costs per car is between \$30k and \$130k.

The range of emissions is between 0.0 g-CO<sub>2</sub>/km (for a fully electric vehicle) and 260.0 g CO<sub>2</sub>/km (for a large gasoline vehicle).





