

# Market Research and Analysis

## Lecture 1: Big Data Analysis in AI Era

---

- Zhenyu Zhao
- Nankai Institute of International Economics
- Feishu: 2120253538
- Email: [zzynankai@outlook.com](mailto:zzynankai@outlook.com)
- Website: [xishanyu2.github.io](http://xishanyu2.github.io)

# 01 Our world in data

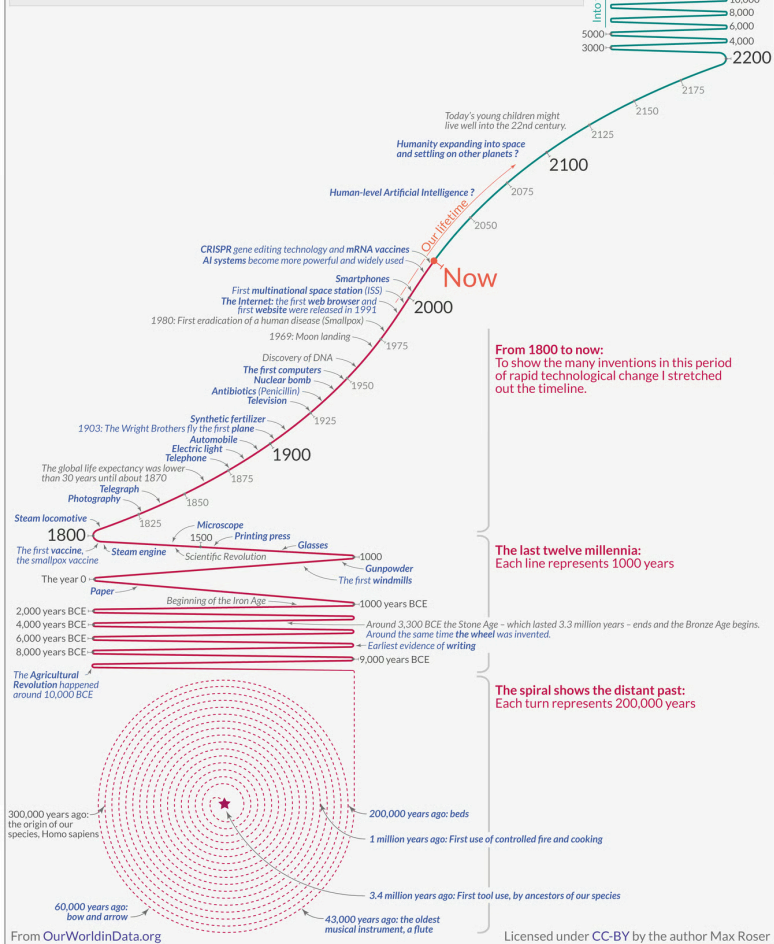
---

In God we trust; all others must bring data. —*William Edwards Deming*

# A long-term timeline of technology

Our World in Data

From the distant past, to our lifetime, and into the distant future.



**From 1800 to now:**  
To show the many inventions in this period of rapid technological change I stretched out the timeline.

**The last twelve millennia:**  
Each line represents 1000 years

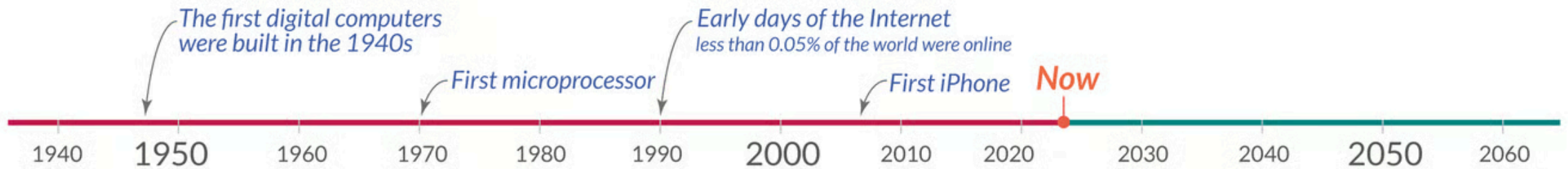
**The spiral shows the distant past:**  
Each turn represents 200,000 years

Licensed under CC-BY by the author Max Roser

## Timeline of the three transformative events in world history:



In a short period, computers evolved so quickly and became such an integral part of our daily lives that it is easy to forget how recent this technology is.







Data is all around us



DARTMOUTH ST

DO NOT ENTER

DO NOT ENTER

NO PARKING  
SNOW

144

LUSH  
FRESH  
HANDMADE  
COSMETICS



## Trees

Boston Parks and Recreation Department (BPRD)

Location (lat-lon, address, neighborhood), species, size (diameter, height), date planted



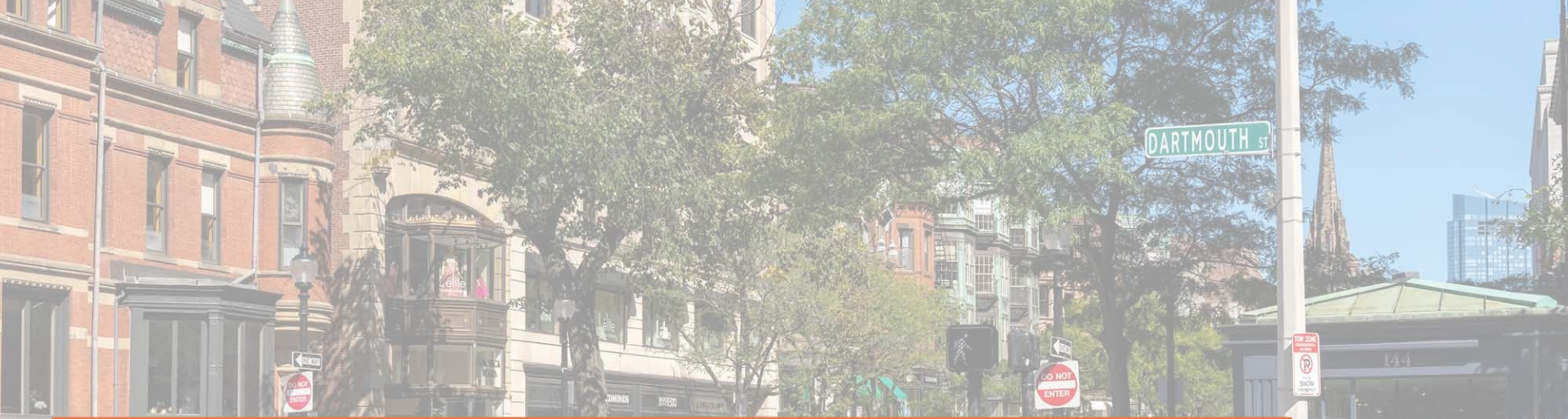


## Building Parcel

City of Boston, Assessing Department

property info (ownership, address, boundaries, area), financial details (value, taxes, mortgage), physical structure (year built, square footage, condition, # of units, basement/garage size), usage/zoning (land use codes, permitted activities), geography (APN, FIPS, coordinates)





## Census Data: Population & Demographics

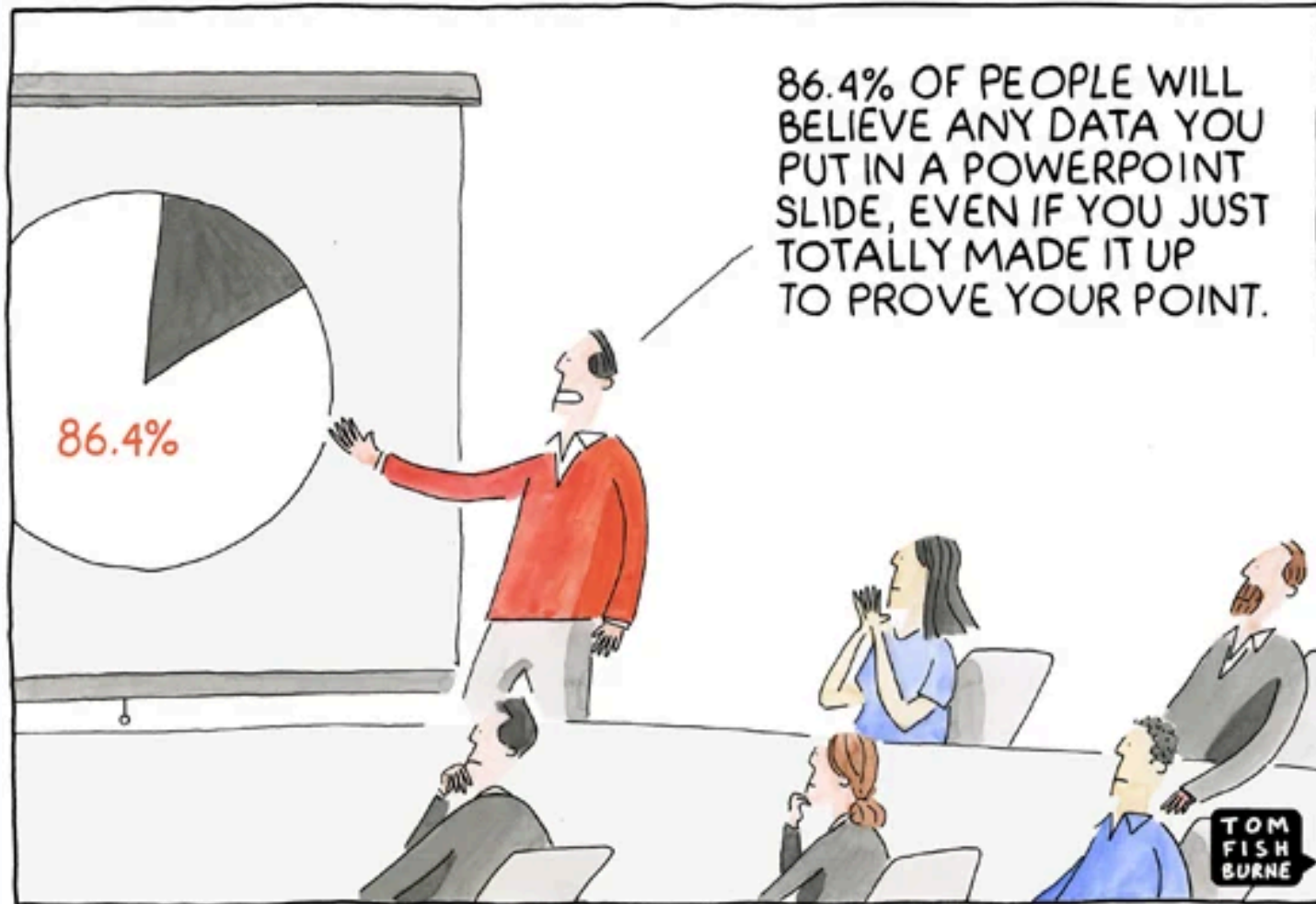
American Community Survey (ACS)

age, sex, race, ethnicity, income, education, employment, ancestry, language, disability, and housing characteristics (tenure, occupancy) for individuals and households, etc.

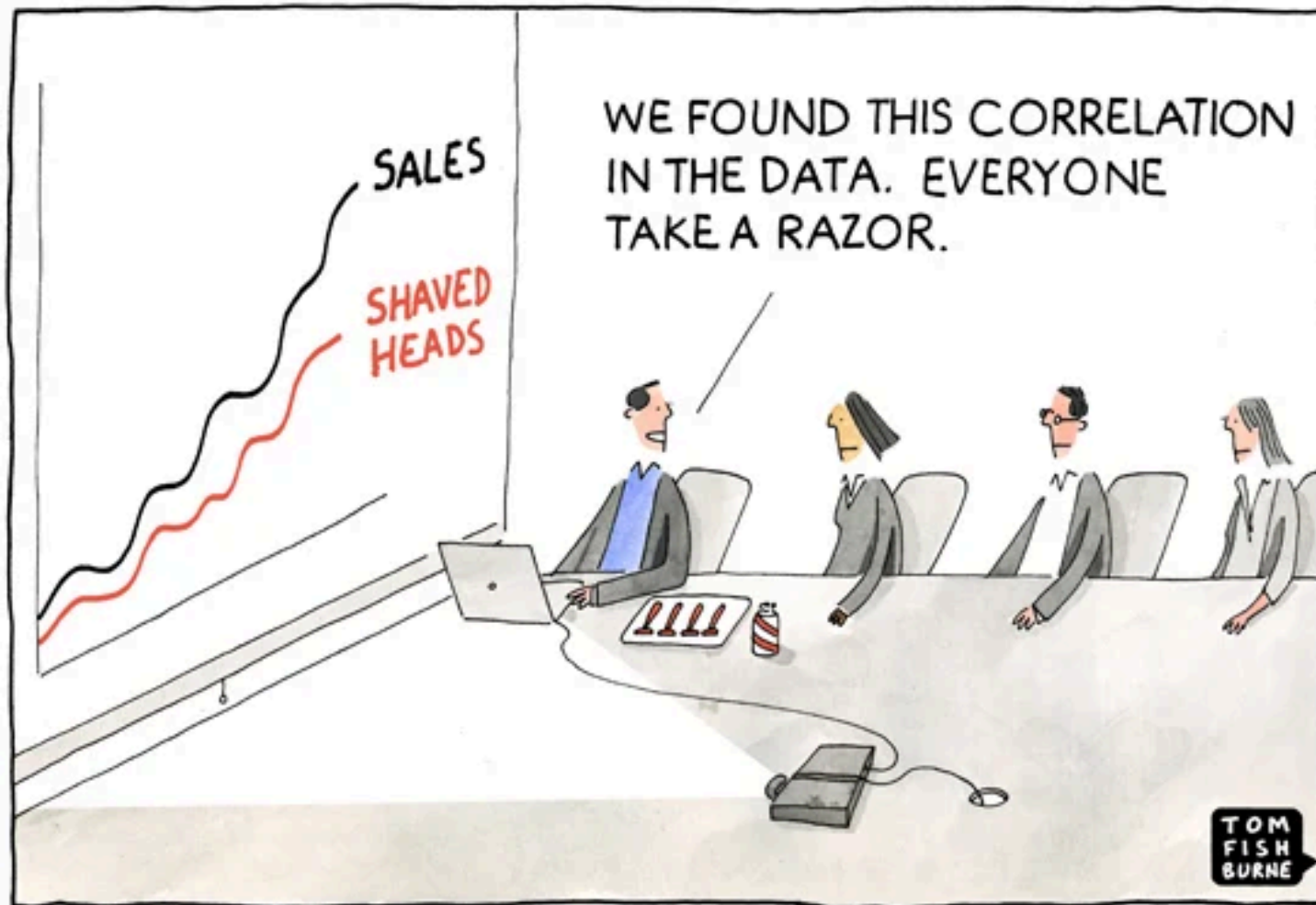
# 02 Data Visualization

---

The ability to take data—to be able to **understand** it, to **process** it, to **extract value** from it, to **visualize** it, to **communicate** it—that's going to be a hugely important skill in the next decades, ... because now we really do have **essentially free and ubiquitous data**. So the complimentary scarce factor is the ability to understand that data and extract value from it. —*Hal Varian*



© marketoonist.com

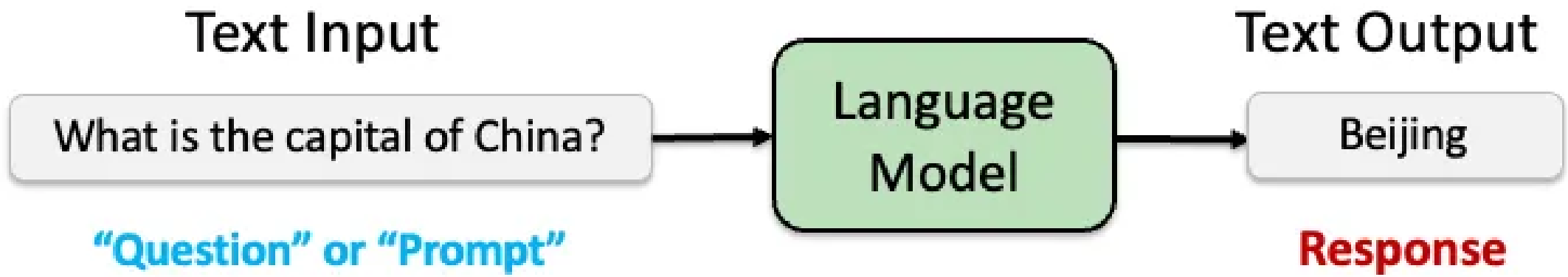


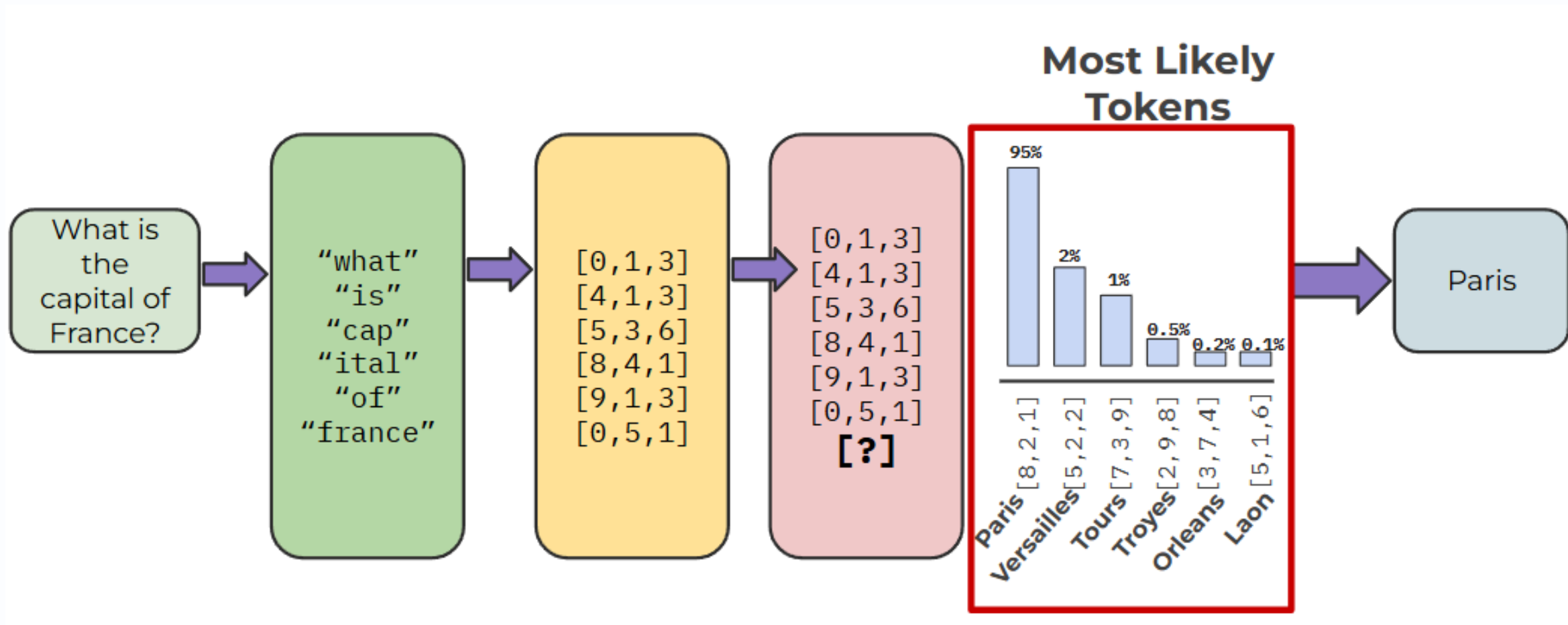
© marketoonist.com

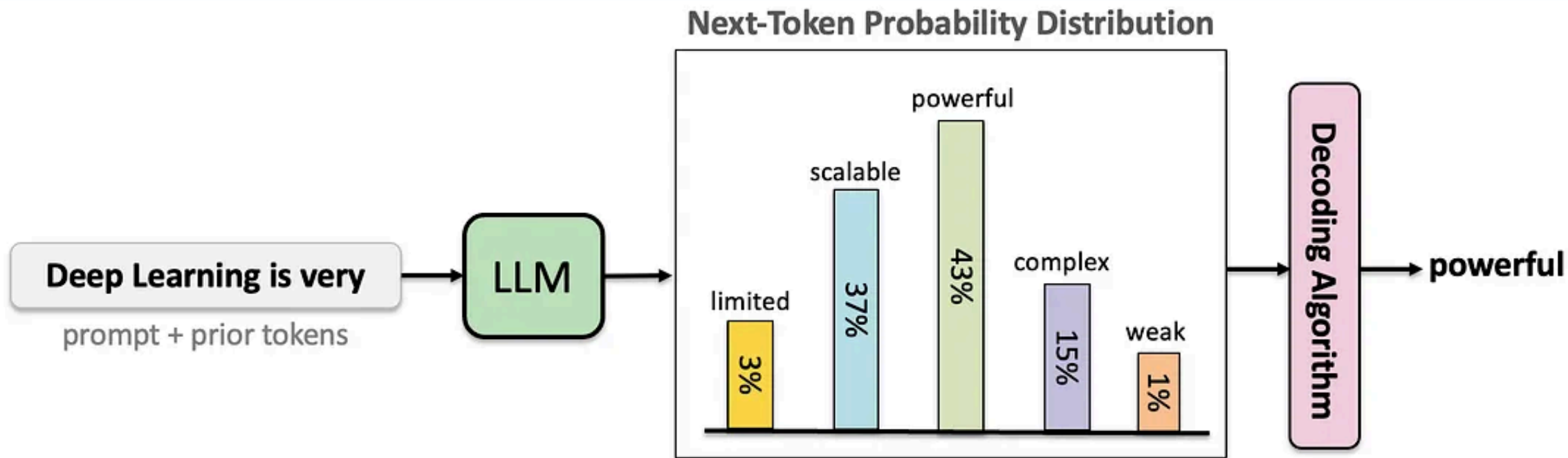
# 03 LLM - Next token predictor

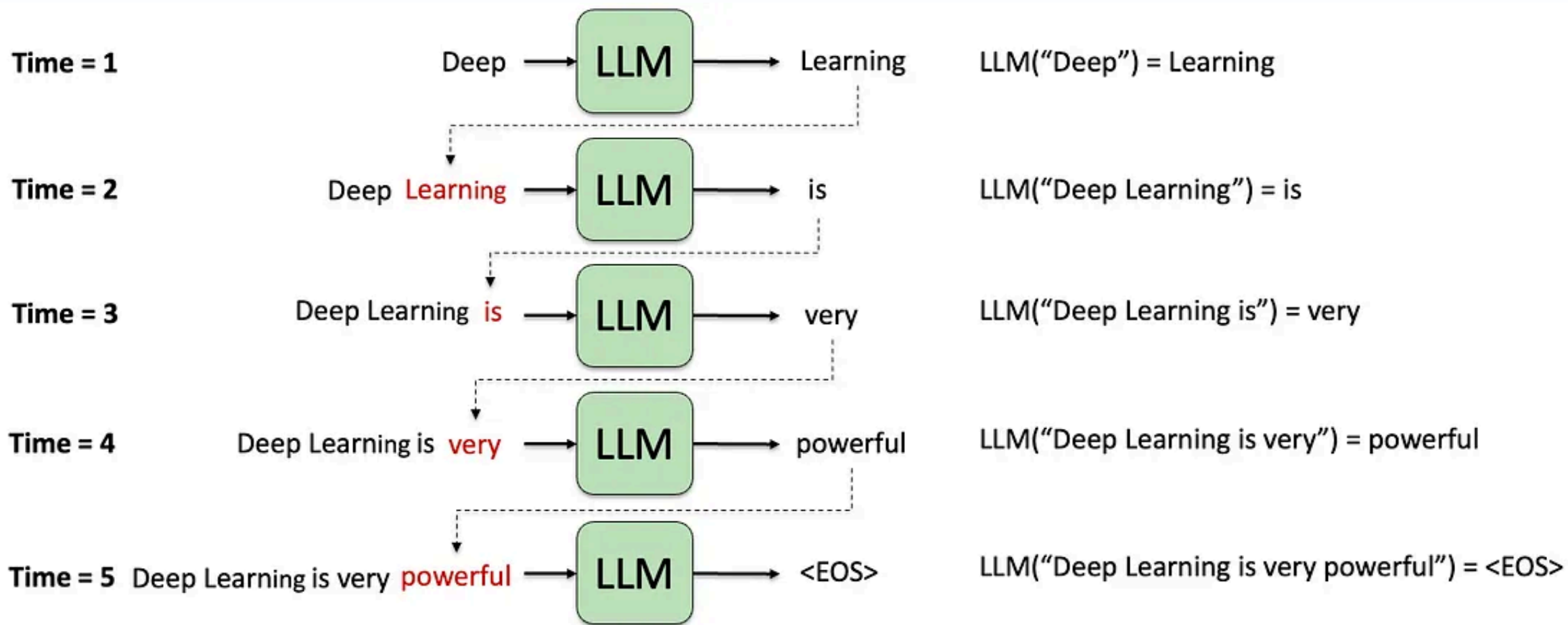
---

The limits of my language mean the limits of my world. —*Ludwig Wittgenstein*









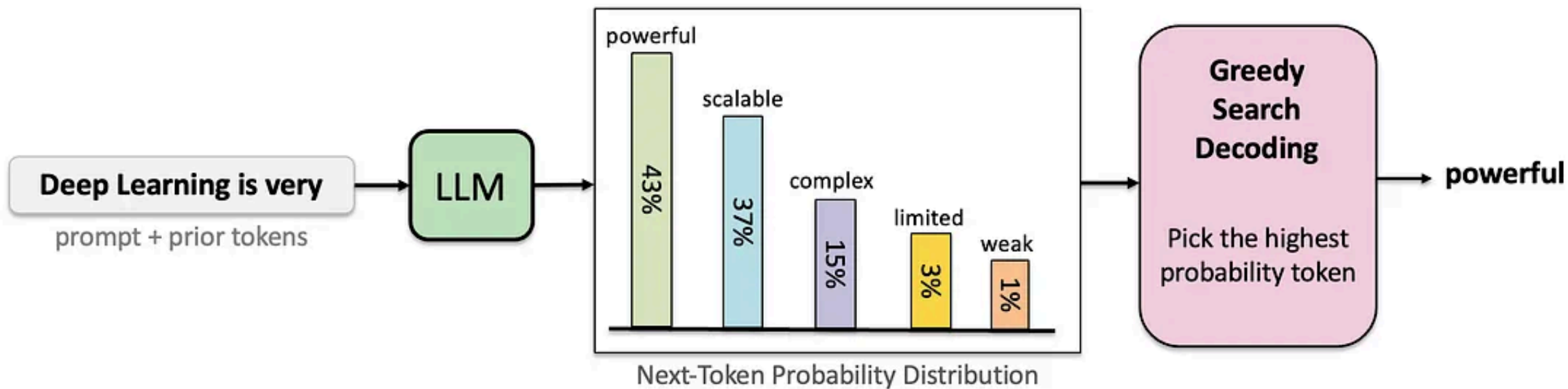
Deep → LLM → Learning

Deep Learning → LLM → is

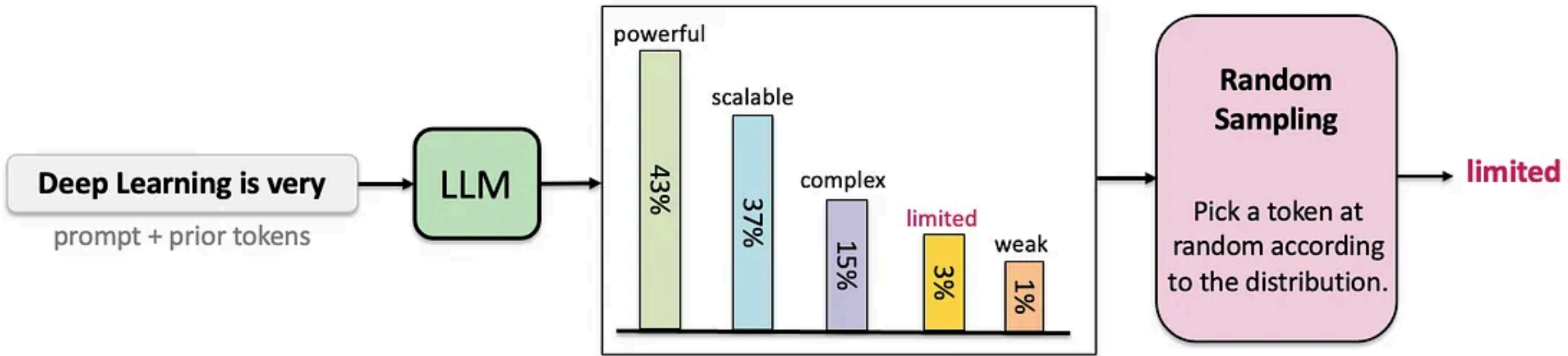
Deep Learning is → LLM → very

Deep Learning is very → LLM → powerful

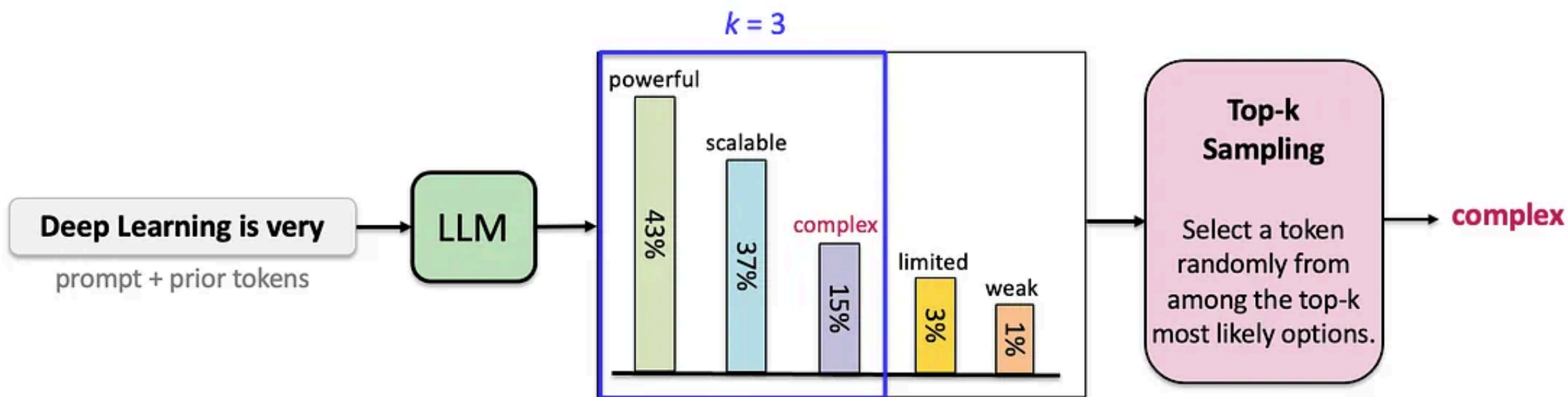
Deep Learning is very powerful → LLM → <EOS>



$$\hat{w}_t \sim P(w_t | w_1, w_2, \dots, w_{t-1})$$

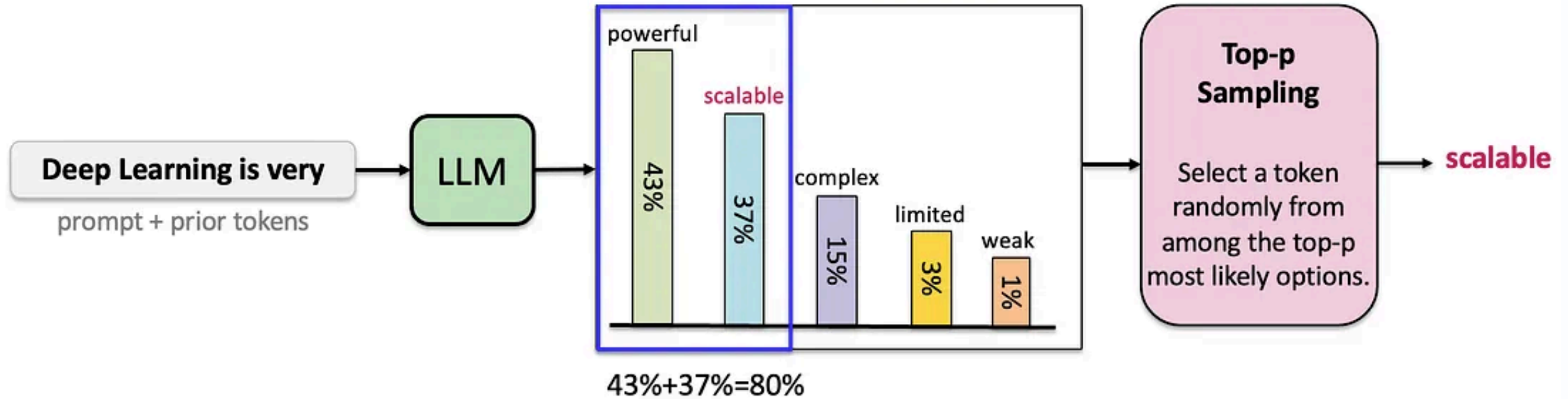


$$\hat{w}_t \sim \text{Top-k} (P(w_t | w_1, w_2, \dots, w_{t-1}))$$



$$\hat{w}_t \sim \text{Top-p} (P(w_t | w_1, w_2, \dots, w_{t-1}))$$

$p = 0.8$  (Cumulative threshold)



# 04 Prompt

---

You don't need to be a data scientist or a machine learning engineer – everyone can write a prompt.  
—*Lee Boonstra*

# Natural Language Programming vs. Traditional Programming

---

"Natural language programming" and traditional programming languages like Python and C++ are essentially both ways of issuing instructions to computers, directing them to perform specific operations.

- The difference lies in:
  - Traditional programming languages (such as Python and C++) have strict syntax and structure.
  - Natural language programming describes operations using human languages (such as Chinese or English).
- Conceptual Shift:
  - Prompt = "Code" of natural language
  - Writing effective prompts is as crucial as writing well-structured Python/C++ code.
  - Many universities have introduced courses on "prompt engineering," and "prompt engineers" will become a sought-after profession.

People respond to content that is  
"AUTHENTIC," "GENUINE," and "HUMAN."  
So let's include those words in our AI prompts.



TOM  
FISH  
BURNE

© marketoonist.com

# System Prompt vs User Prompt

Aspect	System Prompt	User Prompt
<i>Purpose</i>	Defines the AI's long-term behavior, role, and constraints	Provides the specific task or question for a single interaction
<i>Who Sets It</i>	Developers or system administrators	End users or application code
<i>When It Changes</i>	Rarely—updated when refining AI behavior or adding new capabilities	Every interaction—varies with each user request
<i>Example</i>	"You are an experienced customer service representative. Always maintain a polite and professional tone."	"Write a 500-word essay on the impact of social media on modern society, including the benefits and drawbacks."

# Useful tips for designing prompts

Ingredient	Collected Prompts	Prin.
Task Description	T1. Make your prompt <b>as detailed as possible</b> , e.g., “Summarize the article into a short paragraph within 50 words. The major storyline and conclusion should be included, and the unimportant details can be omitted.”	①
	T2. It is helpful to let the LLM know that it is <b>an expert with a prefixed prompt</b> , e.g., “You are a sophisticated expert in the domain of compute science.”	①
	T3. Tell the model <b>more what it should do</b> , but not what it should not do.	①
	T4. To avoid the LLM to generate too long output, you can just use the prompt: “Question: Short Answer: ”. Besides, you can also use the following suffixes, “in a or a few words”, “in one of two sentences”.	①
Input Data	I1. For the question required factual knowledge, it is useful to first <b>retrieve relevant documents</b> via the search engine, and then <b>concatenate them into the prompt</b> as reference.	④
	I2. To highlight some important parts in your prompt, please <b>use special marks</b> , e.g., quotation (“”) and line break (\n). You can also use both of them for emphasizing.	④
Contextual Information	C1. For complex tasks, you can <b>clearly describe the required intermediate steps</b> to accomplish it, e.g., “Please answer the question step by step as: Step 1 - Decompose the question into several sub-questions, . . .”	②
	C2. If you want LLMs to provide the score for a text, it is necessary to provide a <b>detailed description about the scoring standard</b> with examples as reference.	①
	C3. When LLMs generate text according to some context (e.g., making recommendations according to purchase history), instructing them with <b>the explanation about the generated result</b> conditioned on context is helpful to improve the quality of the generated text.	②

Please click [here](#) to view more detailed information.

# 05 Tokens

---

Tokens are the new commodity. —*Jensen Huang*

# Tokenization

---

Large language models process text using **tokens**, which are common sequences of characters found in a set of text. The models learn to understand the statistical relationships between these tokens, and excel at producing the next token in a sequence of tokens.

You can use the tool below to understand how a piece of text might be tokenized by a language model, and the total count of tokens in that piece of text.

[Tokenizer - OpenAI API](#)

A helpful rule of thumb is that one token generally corresponds to ~4 characters of text for common English text. This translates to roughly  $\frac{3}{4}$  of a word (so 100 tokens  $\approx$  75 words).

GPT-4o & GPT-4o mini GPT-3.5 & GPT-4 GPT-3 (Legacy)

Token 的本質：AI 模型不是在看字，而是在看數字

Clear

Show example

Tokens	Characters
17	27

```
[4421, 42905, 4087, 50461, 1817, 17527, 154503, 14951, 35405, 2178, 3812, 8134, 49787, 165308, 3812, 79056, 8134]
```

Text

Token IDs

Specifically, tokens are the segments of text that are fed into and generated by the machine learning model. These can be individual characters, whole words, parts of words, or even larger chunks of text.

TEXT

TOKEN IDS

# Examples

---

- In English, "Hello" is usually a single token, and the model treats it as a whole as a semantic unit.
- Longer English words may be broken down, such as "Tokenization" being divided into two parts: "Token" and "ization", which helps the model handle root changes and word form expansion.
- Chinese vocabulary is sometimes treated as a token as a whole, for example, "没有" is a common word, so the model treats it as a single semantic unit.
- Not all Chinese words are considered a token. "南开" is split into two tokens in the model, indicating that the model understands the constituent words at a finer granularity.

**1300000000000000!**

1.3 quadrillion—a staggering number. This is **the number of tokens processed by Google each month** .

In the Chinese language world, 1.3 quadrillion tokens represent approximately 2.17 quadrillion Chinese characters. Converted to conversation volume, a copy of "Dream of the Red Chamber" has around 700,000 to 800,000 characters. This means that everyone and Google AI chatted about nearly 3 billion copies of "Dream of the Red Chamber" in a single month.

My favorite color is red.

Text

Token IDs

[3666, 4004, 3124, 318, 2266, 13]

Text

Token IDs

My favorite color is Red.

Text

Token IDs

[3666, 4004, 3124, 318, 2297, 13]

Text

Token IDs

Red is my favorite color.

Text

Token IDs

[7738, 318, 616, 4004, 3124, 13]

Text

Token IDs

# Observations

---

The more probable/frequent a token is, the lower the token number assigned to it:

- The token generated for the period is the same ("13") in all 3 sentences. This is because, contextually, the period is used pretty similarly throughout the corpus data.
- The token generated for 'red' varies depending on its placement within the sentence:
  - Lowercase in the middle of a sentence: ' red' - (token: "2266")
  - Uppercase in the middle of a sentence: ' Red' - (token: "2297")
  - Uppercase at the beginning of a sentence: 'Red' - (token: "7738")

In November 2022, OpenAI announced ChatGPT, trained on thousands, tens of thousands of Nvidia GPUs in a very large AI supercomputer: One million users after five days, one million after five days, 100 million after two months, the fastest-growing application in history.

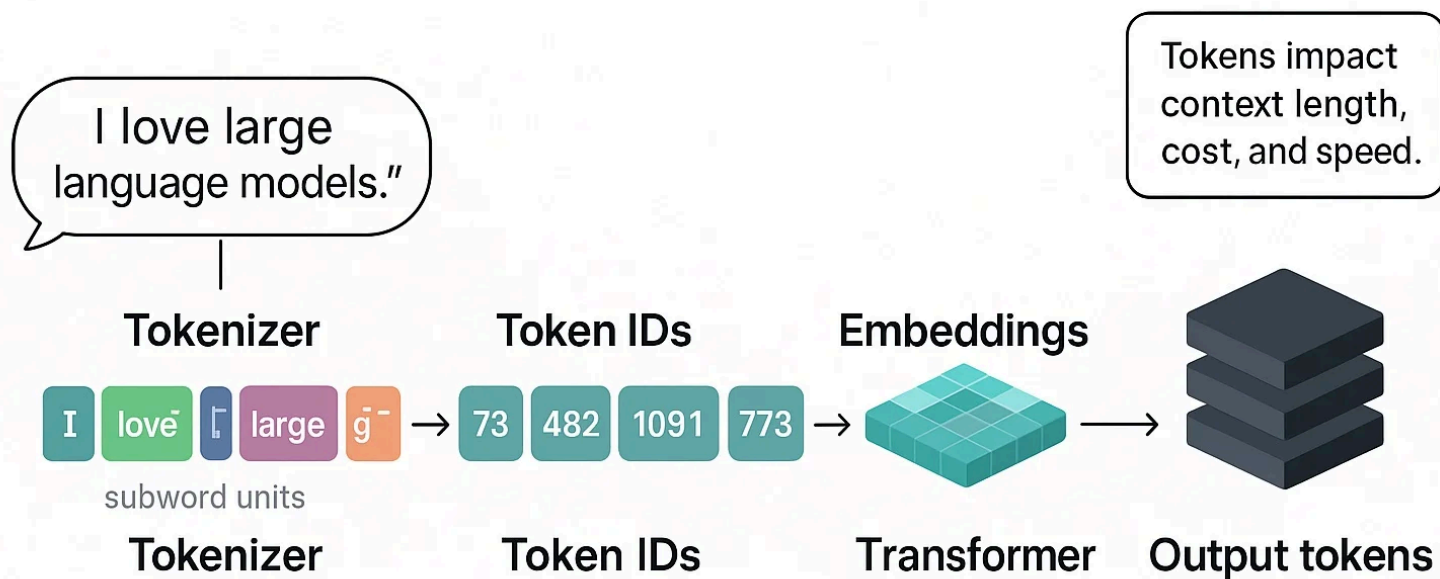
And the reason for that is very simple. It is just so easy to use, and it was so magical to use to be able to interact with a computer like it's human. Instead of being clear about what you want, it's like the computer understands your meaning. It understands your intention.

Until ChatGPT revealed it to the world. AI was all about perception, natural language understanding, computer vision, speech recognition. It's all about perception and detection.

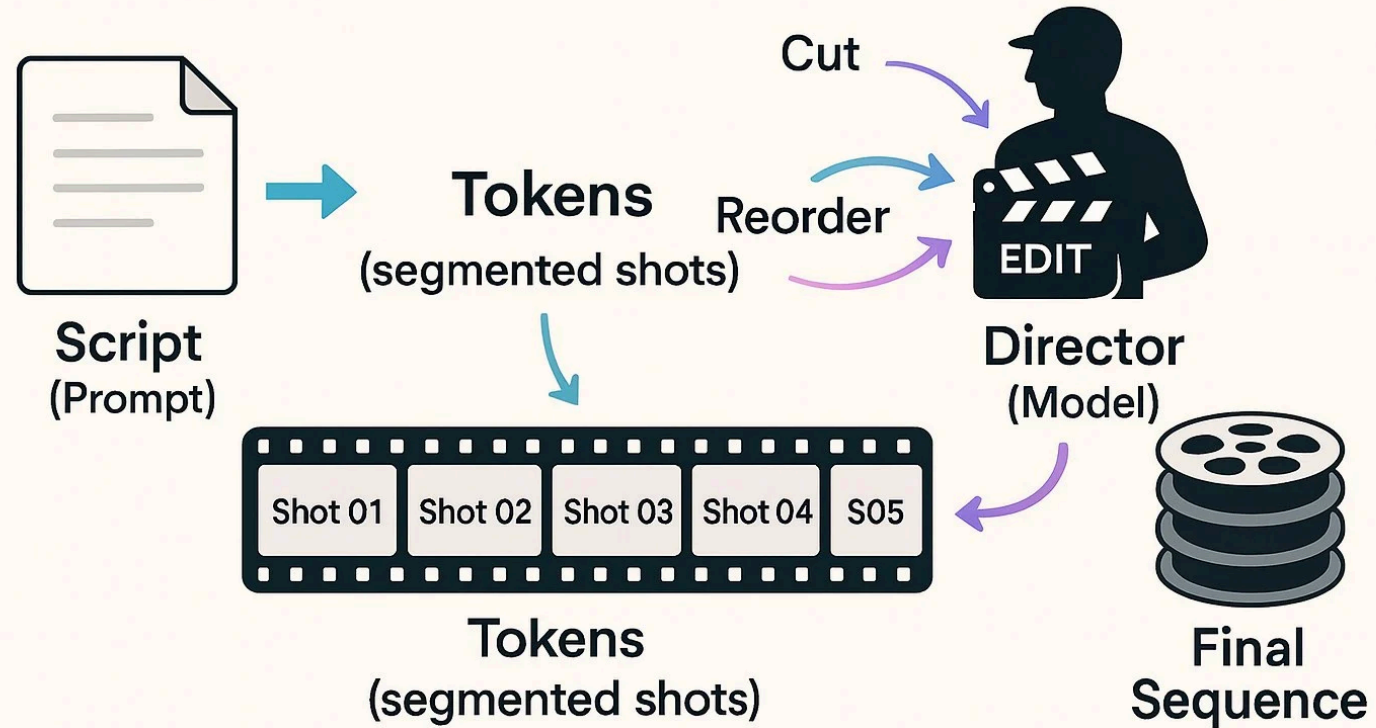
This was the first time the world saw a generative AI.

It produced tokens one token at a time, and those tokens were words. Some of the tokens, of course, could now be **images or charts or tables, songs, words, speech, videos**. Those tokens could be anything, anything that you can learn the meaning of. Everything that we can learn, we can now generate. We have now arrived not at the AI era, but a generative AI era.

# What is a Token?



# Tokens as Shots



Did you know Chinese is more efficient for LLM to process?

Because it literally takes fewer tokens than English — and in LLMs, tokens are everything.

Example:

English: “I love learning languages.” → ~7 tokens

Chinese: “我喜欢学习语言” → ~6 tokens

Each Chinese character carries meaning on its own, while English words often get split into smaller chunks — “learn” → “learn”, “ing”.

That means to say the same thing, Chinese uses 30–40% fewer tokens.

**Same message. Less data. Lower cost.**

A few more interesting details:

1. Emojis are surprisingly expensive — some take multiple tokens.
2. Spaces and punctuation also count.
3. And in most models, 1 token  $\approx$  4 English characters — but just 1 Chinese character.

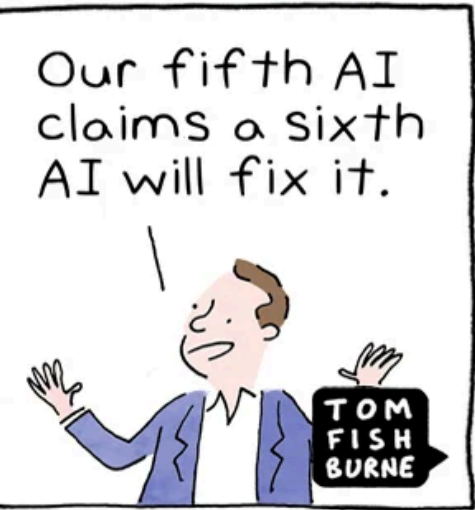
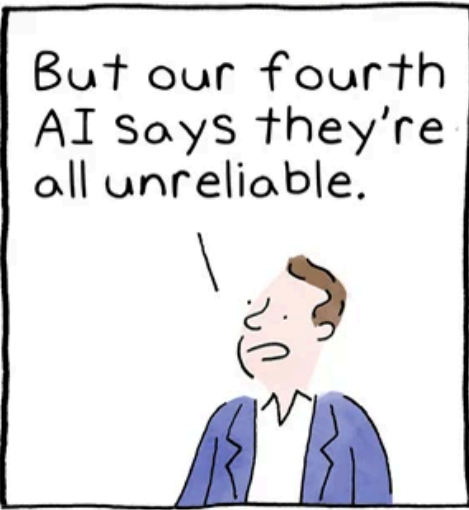
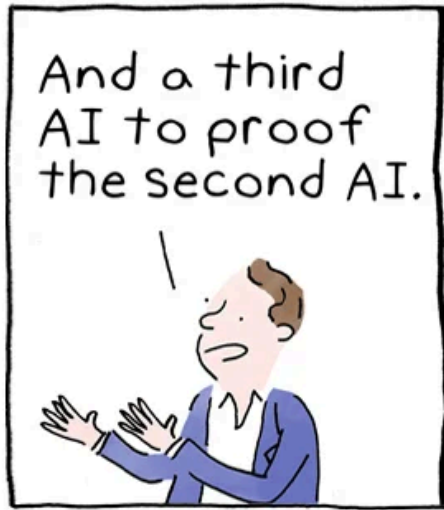
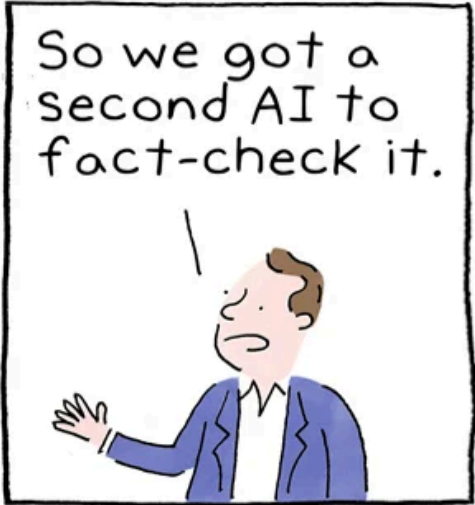
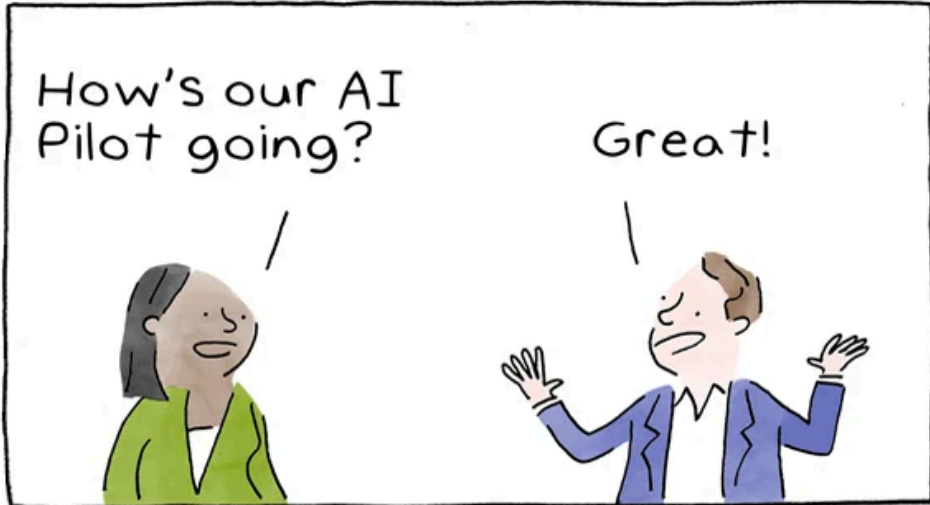
It’s fascinating how something as fundamental as language structure directly impacts **AI performance, cost, and efficiency.**



# 06 AI Hallucination

---

AI hallucination occurs when a model generates plausible yet incorrect information.



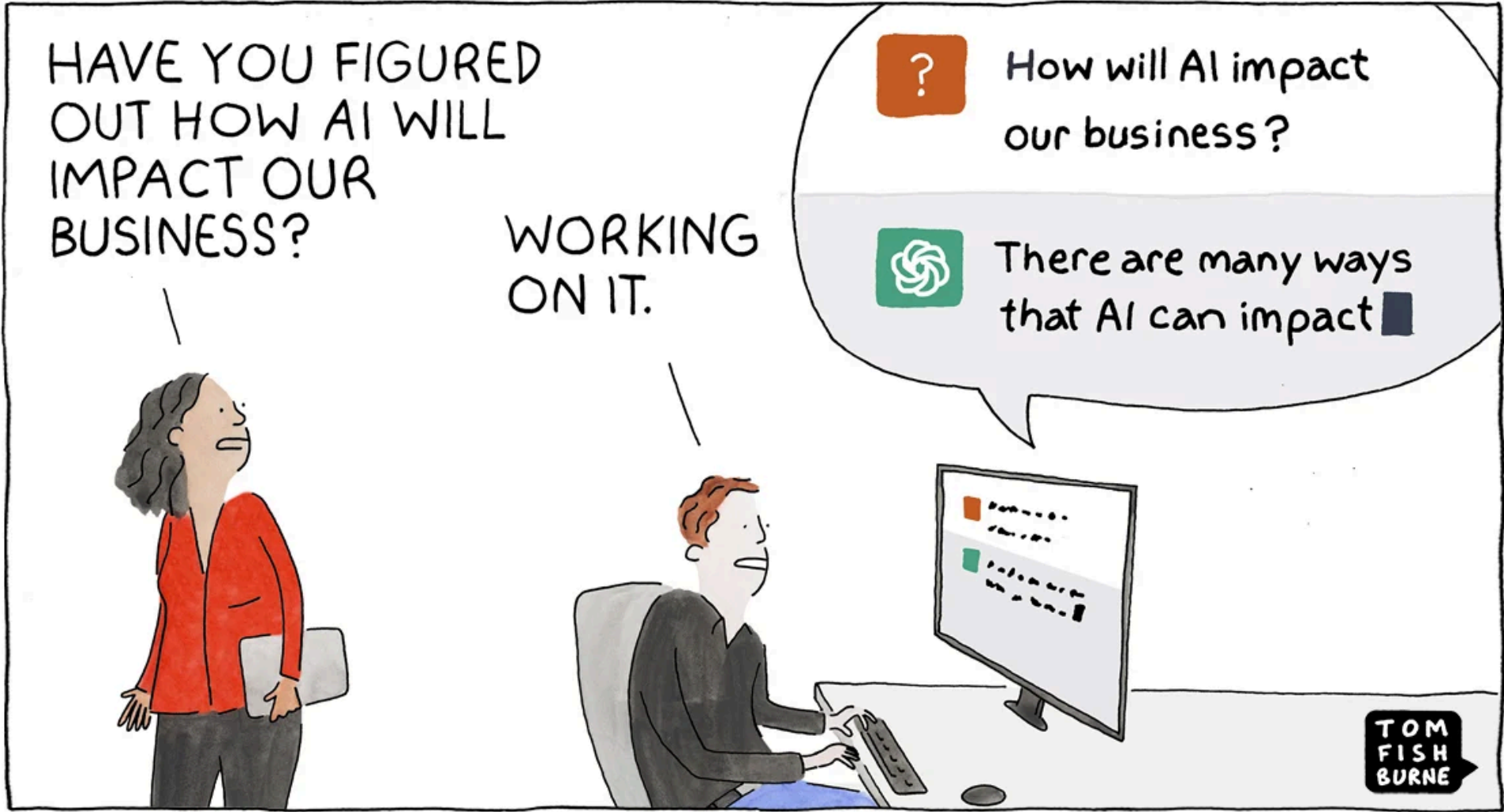
© marketoonist.com

TOM FISH BURNE

# 07 AI Adoption

---

The danger isn't that ChatGPT will replace us. But it can make us stupid—if we let it replace our thinking instead of enriching it.



© marketoonist.com

LET'S TRANSFORM OUR ENTIRE BUSINESS  
USING THE GENERATIVE AI I JUST USED  
TO WRITE A POEM ABOUT MY DOG.



© marketoonist.com

# 美国的 AI 大战

**Sam Altman** @sama

GPT-5.3-Codex 来了!

\*最佳编码性能 (SWE-Bench Pro 57%, TerminalBench 2.0 76%, OSWorld 64%)。  
\*任务执行过程中的可控性和实时更新。  
速度更快! 完成相同任务所需的令牌数量不到 5.2-Codex 的一半, 并且每个令牌的速度提高了 25%!  
良好的电脑使用习惯。

上午5:14 · 2026年2月6日 · 139.9万 查看

**Claude** @claudeai

隆重推出 Claude Opus 4.6。我们最智能的型号迎来了升级。

Opus 4.6 计划更周密, 能够更长时间地执行代理任务, 在大规模代码库中可靠运行, 并且能够发现自身的错误。

这也是我们首款拥有 100 万个令牌上下文的 Opus 级模型 (测试版)。

AI is fast moving. I expect Claude Cwork to leave its research preview and hit prime time sometime in 2026. Stay tuned. Once it has some serious market traction, I may add it to my productivity stack.

**how Claude is be**

Have you experimented with Claude Cwork or other AI tools that work directly on your file system? How comfortable are you letting an AI read, rename, or reorganize your files, even in a test folder? What use cases would make a tool like this worth the risk of the unknown?

0:01 / 0:39

上午4:45 · 2026年2月6日 · 481.2万 查看

# 中国的 AI 大战

给你发了一个现金红包!

元宝派红包, 新春领不停

元宝

【千问】复制打开千问App, 1分钱罐奶茶!

**马化腾: 希望重现当年微信红包盛况**

第一财经 第一财经 2026年1月26日 14:17 上海

1月26日下午, 第一财经记者独家获悉, 在今日正在召开的腾讯年会上, 腾讯董事会主席马化腾提及AI应用元宝即将展开的春节分10亿元现金活动, 表示希望重现当年微信红包的盛况。他还提及元宝要做社交类产品元宝派, 有员工称这本来是绝密的项目。

此前, 2015年微信支付凭借春晚“摇一摇”发放数亿现金红包, 收割海量用户, 成功将数亿用户绑定至其生态之下, 完成对用户支付习惯的一次闪电式改造。今年2月1日, 元宝将上线春节活动, 用户上元宝App分10亿元现金红包, 单个红包金额可达万元。

复制! 打开元宝App, 春节红包马上抢!

你的好友送你一张免单卡

千问请客, 1分钱喝奶茶

千问APP

2026年2月



