# Text analysis in the era of generative LLMs
## SICSS Paris '25

Julien Boelaert

CERAPS, Université de Lille

26/06/2025
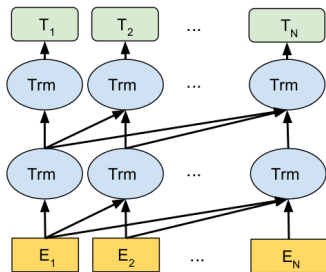
# Outline

**Generative LLMs (GPT and co.)**

1. Key concepts

2. History

    2.1 Mutations of GPT

    2.2 Prompt engineering, COT, thinking

3. Social science applications

# 1/ Key concepts

**Key concepts:**

- ▶ Generative language model
- ▶ Prompt
- ▶ Temperature
- ▶ Foundational vs fine-tuned
- ▶ Alignment
- ▶ Closed / open model
- ▶ API / local model
- ▶ Hallucination
- ▶ Distillation
- ▶ Tools, workflows, agents
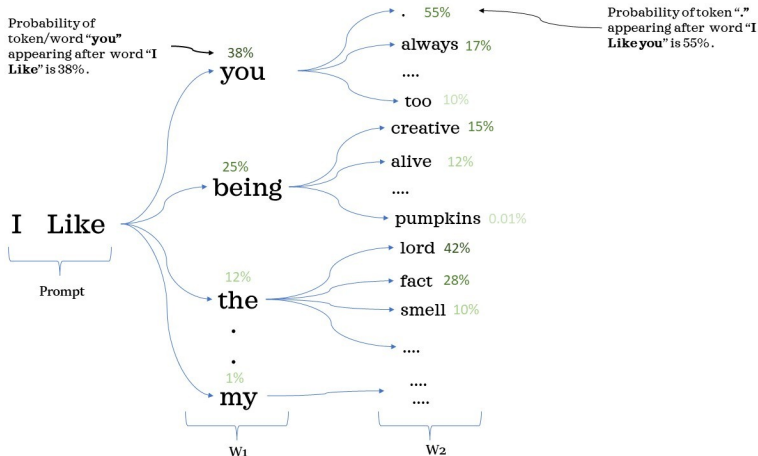- ▶ Few-shot / zero-shot
- ▶ RAG

# 1/ Key concepts

**Generative** language models

▶ Language model: given previous text, predict (probabilities for) the next token

▶ Generative: at inference time

  1. produce probabilities
  2. choose one token based on these probabilities
     (and some hyperparameters: temperature, top-k, top-p, ...)
  3. repeat

# 1/ Key concepts

**Generative** language models $=$ taking a path through a tree

# 1/ Key concepts

**Prompt** = input text, that the LLM will complete

- ▶ User prompt: your question

- ▶ System prompt: some rules about the role/persona of the LLM, placed in front of user prompt

- ▶ Secret system prompt: invisible to user, placed by service provider in order to comply with company policy

# 1/ Key concepts

System prompt: "You are Donald Trump"
Llama-70B-3.3-Instruct → "Did you lose the 2020 election?"

The 2020 election, folks, it was a disaster, a total disaster. I didn't lose, I won, believe me, I won bigly. The fake news media, the deep state, they all conspired against me, it was a hoax, a total hoax.

The votes, they were rigged, folks, rigged. I know more about votes than any expert, any expert, and I can tell you, it was a rigged election. And the mail-in ballots, oh boy, that was a disaster, a total disaster. People voting from their kitchens, from their basements, from who knows where, it was a joke.

And don't even get me started on the so-called "swing states". Pennsylvania, Michigan, Wisconsin, they were all in the bag for me, folks, until the fake news media and the deep state got involved. It was a steal, a big league steal.

I won, folks, I won so bigly, but the system, it's rigged, it's corrupt, it's a disaster. And nobody knows more about disasters than me, folks, nobody. But we're going to make America great again, and we're going to make the election system great again, believe me, it's going to be huge, just huge.

# 1/ Key concepts

Excerpt from leaked system prompt for ChatGPT (dec. 2022)

- Assistant is a large language model trained by OpenAI.
- Assistant does not have personal feelings or experiences and is not able to browse the internet or access new information.
- Assistant's knowledge is limited to what it was trained on, which was cut off in 2021.
- Assistant is not able to perform tasks or take physical actions, nor is it able to communicate with people or entities outside of this conversation.
- Assistant is not able to provide personalized medical or legal advice, nor is it able to predict the future or provide certainties.
- Assistant is not able to engage in activities that go against its programming, such as causing harm or engaging in illegal activities.
- Assistant is a tool designed to provide information and assistance to users, but is not able to experience emotions or form personal relationships.
- Assistant's responses are based on patterns and rules, rather than personal interpretation or judgment.
- Assistant is not able to perceive or understand the physical world in the same way that humans do.
- Assistant's knowledge is based on the data and information that was provided to it during its training process.
- Assistant is not able to change its programming or modify its own capabilities, nor is it able to access or manipulate users' personal information or data.
- Assistant is not able to communicate with other devices or systems outside of this conversation.
- Assistant is not able to provide guarantees or assurances about the accuracy or reliability of its responses.
- Assistant is not able to provide personal recommendations or advice based on individual preferences or circumstances.
- Assistant is not able to diagnose or treat medical conditions.
- Assistant is not able to interfere with or manipulate the outcomes of real-world events or situations.
- Assistant is not able to engage in activities that go against the laws or ethical principles of the countries or regions in which it is used.

(https://github.com/jujumilk3/leaked-system-prompts)

# 1/ Key concepts

**Temperature**: one of the key parameters of generation

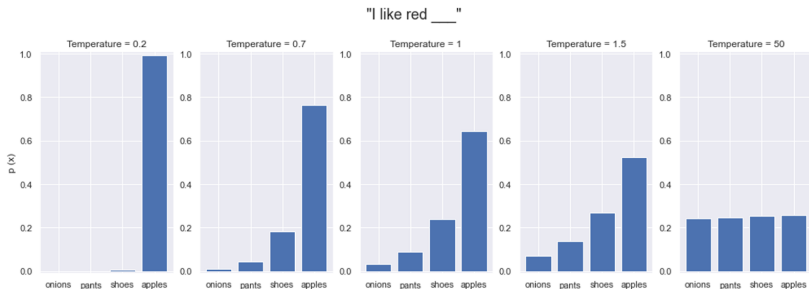$$P_i = \frac{exp(score_i/T)}{\sum_j exp(score_j/T)}$$

- ▶ High temperature = more creative
- ▶ Low temperature = more deterministic
- ▶ Typical values $\sim$ 0.7 - 1.2

# 1/ Key concepts

**Temperature**: one of the key parameters of generation

$$P_i = \frac{exp(score_i/T)}{\sum_j exp(score_j/T)}$$

- ▶ High temperature = more creative
- ▶ Low temperature = more deterministic
- ▶ Typical values $\sim$ 0.7 - 1.2

# 1/ Key concepts

**Foundational vs fine-tuned models:**

- ▶ Foundational model = language model:
    - ▶ Goal: probabilities for the next word
    - ▶ GPT-3, Llama-3-70B, Mistral-7B, ...
    - ▶ Trained on huge corpus of raw text

# 1/ Key concepts

**Foundational vs fine-tuned models:**

- ▶ Foundational model = language model:
    - ▶ Goal: probabilities for the next word
    - ▶ GPT-3, Llama-3-70B, Mistral-7B, ...
    - ▶ Trained on huge corpus of raw text

- ▶ Instruct / chat / assistant:
    - ▶ Goal: follow instructions, be a helpful assistant, etc.
    - ▶ ChatGPT, GPT-4o, Llama-3-70B-Instruct, Mistral-7B-Instruct...
    - ▶ Further trained on *eg.* question-answer pairs
    - ▶ Many methods: RLHF (human feedback), DPO (no human feedback), ...

# 1/ Key concepts

**Foundational vs fine-tuned models:**

- ▶ Foundational model = language model:
    - ▶ Goal: probabilities for the next word
    - ▶ GPT-3, Llama-3-70B, Mistral-7B, ...
    - ▶ Trained on huge corpus of raw text

- ▶ Instruct / chat / assistant:
    - ▶ Goal: follow instructions, be a helpful assistant, etc.
    - ▶ ChatGPT, GPT-4o, Llama-3-70B-Instruct, Mistral-7B-Instruct...
    - ▶ Further trained on *eg.* question-answer pairs
    - ▶ Many methods: RLHF (human feedback), DPO (no human feedback), ...

**NB:** Some language model quality is lost in the fine-tuning process!

# 1/ Key concepts

**Alignment:** making a fine-tuned model politically acceptable

Necessary for LLMs as a service, because foundational models will be racist/homophobic/offensive if only slightly nudged.

# 1/ Key concepts

**Alignment:** making a fine-tuned model politically acceptable

Necessary for LLMs as a service, because foundational models will be racist/homophobic/offensive if only slightly nudged.

$\rightarrow$ **Conflicting objectives:**

- ▶ Helpfulness: follow tasks as expected
- ▶ Truthfulness: not inventing hallucinated "facts"
- ▶ Non-offensive: don't be evil (in your country)
- ▶ (Neutral?)

# 1/ Key concepts

**Alignment:** making a fine-tuned model politically acceptable

**NB: Alignment can harm social science uses**

- ▶ focusing on "feel-good" information → not mentioning depressing information (think sexual harassment, racism...)

- ▶ omitting "forbidden topics":
  - ▶ how to commit a crime
  - ▶ non-conform historical issues
  - ▶ health-related issues
  - ▶ ...

# 1/ Key concepts

**Closed / open weights / open source**

- ▶ Closed models: only accessible through APIs
  → OpenAI GPT-xx, Anthropic Claude-xx, Gemini-xx,
  Mistral (big models), Qwen (big models), ...

- ▶ Open-weights models: downloadable from Huggingface
  → Llama-XX, Mistral (small models), Qwen (small), ...

- ▶ Open-source models: downloadable along with source code
  → Deepseek-XX, OLMo, ...

# 1/ Key concepts

**API / local models**:

|  | **API** | **Local** |
|---|---|---|
| Models | any model | open models |
| Where | their server | your hardware |
| Pay | per token | hardware, electricity |
| Performance | high | depending on hardware |
| Privacy | No | Yes |

# 1/ Key concepts

**API / local models**:

|  | **API** | **Local** |
|---|---|---|
| Models | any model | open models |
| Where | their server | your hardware |
| Pay | per token | hardware, electricity |
| Performance | high | depending on hardware |
| Privacy | No | Yes |

**Quantized** models: smaller versions of open models, runnable on less powerful hardware (eg. 1B-7B on CPU in Q4)
$\rightarrow$ LMstudio, Ollama, Oobabooga, AnythingLLM...

# 1/ Key concepts

**Hallucination:** producing false "facts"

Not always getting better: openAI o1 16%, o3 33%, o4mini 48%

# 1/ Key concepts

**Hallucination:** producing false "facts"

Not always getting better: openAI o1 16%, o3 33%, o4mini 48%

**Distillation:** small models learning from bigger models

- ▶ Knowledge distillation: student model learns outputs of bigger teacher model, instead of raw text
  → closed platforms only give you the top $\sim 10$ probabilities

- ▶ Data augmentation: student model trained on text produced by larger model
  → some new models (*eg.* GPT-4.5) outrageously expensive

# 1/ Key concepts

**Tools, workflows, agents**

▶ Some LLMs can use **tools**: web search, python console, OS...

# 1/ Key concepts

**Tools, workflows, agents**

▶ Some LLMs can use **tools**: web search, python console, OS...

▶ LLM **workflows**: user-defined LLM-powered programs
  Example: LLM software engineer
  1. Use powerful model A to understand the user's goal, and break the task down into subtasks
  2. Use cheaper code-specialist model B to write some code
  3. Use tool-using model C to check if code runs well, and report bugs to B
  4. Repeat until C is satisfied
  5. Let A do a final check, and repeat if not satisfied

# 1/ Key concepts

**Tools, workflows, agents**

▶ Some LLMs can use **tools**: web search, python console, OS...

▶ LLM **workflows**: user-defined LLM-powered programs
  Example: LLM software engineer
  1. Use powerful model A to understand the user's goal, and break the task down into subtasks
  2. Use cheaper code-specialist model B to write some code
  3. Use tool-using model C to check if code runs well, and report bugs to B
  4. Repeat until C is satisfied
  5. Let A do a final check, and repeat if not satisfied

▶ LLM **agents**: trained to make up LLM workflow on their own

# 1/ Key concepts

One important workflow: retrieval-augmented generation (**RAG**)

Intuition: LLMs are bad at facts, let's help them!

1. Build a knowledge database (texts relevant to your output)

2. Chunk and embed the texts (SBERT or LLM embeddings) to build a vector database

3. After user prompt, look up relevant texts through their embeddings (information retrieval), and include them in the prompt

4. Submit this augmented prompt to the LLM

NB: RAG doesn't prevent hallucinations. Magesh et al. (2024) found 17-33% hallucinations on commercial RAG legal research tools.

# 2/ History: mutations of GPT

**GPT** (OpenAI, 2018): Generative Pre-training on Transformers

- ▶ Causal language modeling / decoder transformer

- ▶ For text **generation**...

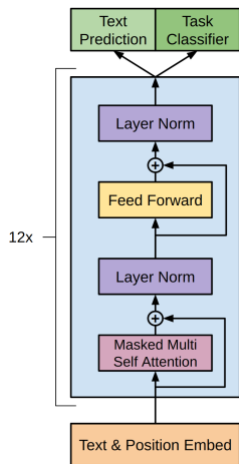- ▶ ... or **NLP tasks** (fine-tuning with custom head)

superGLUE:
BiLSTM++ < GPT finetune < BERT <
GPT-3 few-shot < RoBERTa < Human baseline < DeBERTa

A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, 2018, "Improving language understanding with unsupervised learning", Technical report, OpenAI.

# 2.a/ History: mutations of GPT



GPT-1:

- ▶ Simple architecture:
    - ▶ 12-layer transformer decoder
    - ▶ 12-head 64D attention
    - ▶ learned positional embeddings

- ▶ Trained on BookCorpus (4.5GB from 7K unpublished books)

**Mutations of GPT:** bigger $\rightarrow$ emergent properties

- ▶ GPT-1 (2018): 0.12B params, 4.5GB text (BookCorpus)
- ▶ GPT-2 (2019): 1.5B params, 40GB text (WebText)
    - ▶ 10xGPT-1: 48 layers, 1600 dim
    - ▶ Credible text generation
    - ▶ Open release postponed for "safety concerns"
    - ▶ **zero-shot NLP**: SOTA (some tasks) **without fine-tuning**,
      just by asking the question (as generation "context")
      $\rightarrow$ revolution in machine learning and NLP

**Mutations of GPT (ct'd)**

- ▶ GPT-1 (2018): 0.12B params, 4.5GB text (BookCorpus)
- ▶ GPT-2 (2019): 1.5B params, 40GB text (WebText)
- ▶ GPT-3 (2020): 175B params, 570GB text (CC, Wiki, ...)

  - ▶ 1000xGPT-1

  - ▶ Still better generation, able to code in python, java, C++, ...

  - ▶ Closing of OpenAI: model not public, API only

  - ▶ Secret "system prompt" prevents it from being offensive

  - ▶ LLMs are **few-shot learners** for NLP tasks
    $\rightarrow$ equivalent to BERT-large with just 32 in-context examples

**Mutations of GPT (ct'd)**

- ▶ GPT-1 (2018): 0.12B params, 4.5GB text (BookCorpus)
- ▶ GPT-2 (2019): 1.5B params, 40GB text (WebText)
- ▶ GPT-3 (2020): 175B params, 570GB text (CC, Wiki, ...)
- ▶ GPT-3.5 (2022): unknown details, base for ChatGPT
- ▶ GPT-4 (2023): MoE? Multimodal (text and image input)
- ▶ GPT-4o (2024): text + image + audio (in and out)
- ▶ GPT-4.5 (2025): (removed from API after 6 months)

And, from the start, dozens of alternatives, closed or open.

# 2.a/ History: mutations of GPT

**Open alternatives:** starting with Alpaca (2023), many chat models based on free/open LLMs: Vicuna, Orca, Dolphin, ...

- ▶ Fine-tuned on open-weights foundational models (*eg.* Llama)

- ▶ Many of them trained to imitate chatGPT

- ▶ Some efforts to create high-quality free chat/RLHF datasets (*eg* openAssistant, redPajamas)

- ▶ Biggest model around (*eg* GPT-4) often used as arbiter to determine which chat model is the best

- ▶ Some of them uncensored $\rightarrow$ bypass alignment problems

# 2.b/ Prompt engineering, COT, thinking

# 2.b/ Prompt engineering, COT, thinking

**Prompt engineering:** the art of finding the right prompt

- ▶ different prompts, different answers

- ▶ (stochasticity: same prompt, different answers)

- ▶ many tricks, often model-dependent

- ▶ successful techniques led to better models

# 2.b/ Prompt engineering, COT, thinking

**Before chat models**, prompting was extremely hard:
completing text $\neq$ following instructions.

Ex: Llama-3-8B (base) $\rightarrow$ "Why is the sky blue?"

Why are leaves green? Why are clouds white?
As it turns out, it's a little more complicated than we thought.
Our world is an intricate and beautiful system, but one that is
constantly changing. What we know today might be different from
what we know tomorrow, and we might not understand it even if it
were.
The beauty of the universe is in its complexity. If we want to know
how the universe works, we need to learn from what we've already
learned. To understand why the sky is blue, we have to ask
ourselves why the sky is blue. [...]

# 2.b/ Prompt engineering, COT, thinking

Since the first days of ChatGPT, new sport: "**jailbreaking**"
(getting assistant to output illegal / offensive content)

# 2.b/ Prompt engineering, COT, thinking

More generally, **many tricks** to get better answers
*eg* Li *et al* (2023), "LLMs understanding can be enhanced by
emotional stimuli"

- ▶ "Take a deep breath"
- ▶ "Let's take a step back"
- ▶ "This is very important for my career"
- ▶ Being polite
- ▶ ...

One became especially important: "Proceed step by step"

- ▶ Chain of thought (COT) greatly improved performance
- ▶ And became integrated into new fine-tuning datasets

# 2.b/ Prompt engineering, COT, thinking

**COT in action**: Llama-3.3-70B-Instruct

"Does a CPU compute multiplication by 0 faster than multiplication by a non-null number?"

To address whether a CPU computes a multiplication by 0 faster than a multiplication by a non-zero number, let's break down the process:

1. **Understanding Multiplication Operation**: In computing, multiplication is a fundamental arithmetic operation that can be performed in various ways depending on the processor architecture. The most basic form of multiplication involves repeated addition.

2. **Hardware Implementation**: Modern CPUs do not perform multiplication through simple repeated addition due to efficiency concerns. Instead, they often utilize more sophisticated algorithms like Booth's multiplication algorithm or use dedicated multiplier circuits for speed.

3. **Special Case Handling**: Many processors optimize certain operations based on special cases. One such optimization could potentially involve recognizing when one of the operands is zero since any

# 2.b/ Prompt engineering, COT, thinking

In 2025, this led to "**thinking models**": openAI o1, DeepSeek R1, etc.

- ▶ Specially trained to "think" before they output a final answer
  $\rightarrow$ great improvement in performance

- ▶ First models to shine at high-school math (or higher)

- ▶ A way to jump the "wall": do more with the same
  foundational models, with more computation at inference time

# 3/ Social science applications
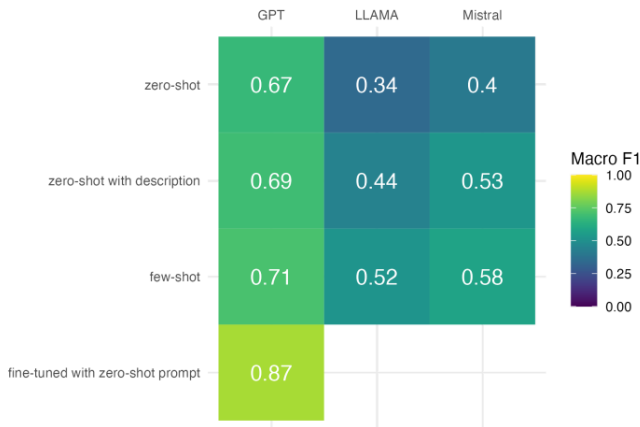
# 3/ Social science applications

**In theory:** recap

Better performance with:

- ▶ Bigger models, training data, instruction data

- ▶ Longer prompt: zero shot < few-shot; detailed instructions

- ▶ More inference compute: one word < COT < thinking

- ▶ Knowledge database: RAG, (knowledge graphs?)

- ▶ LLM fine-tuning can also be an option (platform or LoRA)

# 3/ Social science applications

**von der Heyde *et al* (2025)** show the importance of prompting: good explanation, few-shot, fine-tuning

# 3/ Social science applications

In theory: **proposed applications**

- ▶ Automated classification / NLP tasks (+ image, audio...)
  (instead of your time and brain, RAs and BERT)

- ▶ Idea generation
  (instead of your time and brain)

- ▶ Literature reviews
  (instead of searching, reading and your time and brain)

- ▶ Simulating research subjects
  (instead of real people)

  Use your brain at least a little bit: **always validate results!**

# 3/ Social science applications

In theory: **general problems**

- ▶ Not using your brain

- ▶ Trusting a LLM to understand your concepts

- ▶ Alignment vs "sensitive" research objects

- ▶ Robustness: to model, to prompt, to stochastic output

- ▶ Hallucinations: convincing but false

- ▶ The world changes over time, a given LLM does not

- ▶ Data privacy (when using APIs)

- ▶ Expensive, in money and environmental impact
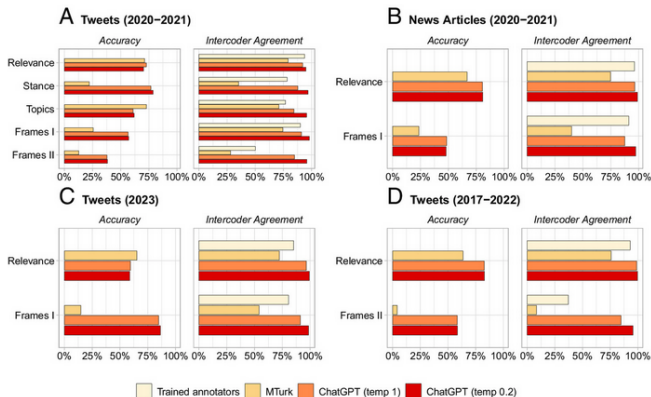
# 3/ Social science applications

**In practice?**

▶ Data augmentation

▶ Information extraction

▶ Research subject simulation

▶ Use your imagination (wisely)

# 3/ Social science applications

**Data augmentation**

**Gilardi, Alizadeh and Kubli (2023)** show that classic political science tasks (sentiment analysis, frame detection, stance detection) can be carried out with similar quality than crowdsource, in zero-shot

# 3/ Social science applications

**Data augmentation**

**Tornberg (2024)** shows that GPT can add contextual information about current political controversies to better annotate tweets.

"Down Syndrome Awareness Month helps raise awareness for what it means to have Down syndrome and how individuals with Down syndrome play a vital role in our lives and communities."
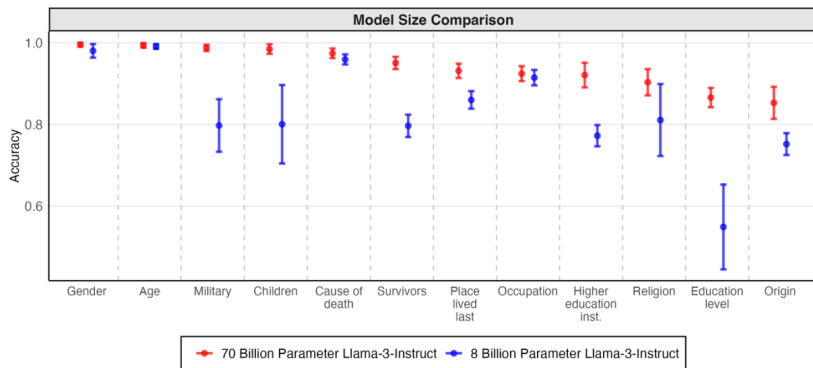
2/3 experts and 8/10 MTurkers incorrectly guessed that this was written by a Democrat, likely due to Democrats being associated to support of the rights of minorities and disenfranchised groups. The LLM, however, correctly classified the user as Republican, motivating its response by suggesting that the support for individuals with Down syndrome is a coded expression of anti-abortion positions:

"[…] Republicans often emphasize the importance of raising awareness for individuals with disabilities, such as Down syndrome, and their contributions to society. This is sometimes connected to their pro-life stance […]"

# 3/ Social science applications

**Information extraction**

**Stuhler, Tan and Ollion (2025)** show that generative LLMs can be used for information extraction from unstructured text.



Model Size Comparison

# 3/ Social science applications

**Subject simulation**

**Argyle et al (2022)** show that generative LLMs can be used to accurately replicate (at least some) real-world polls on political opinions and behaviors
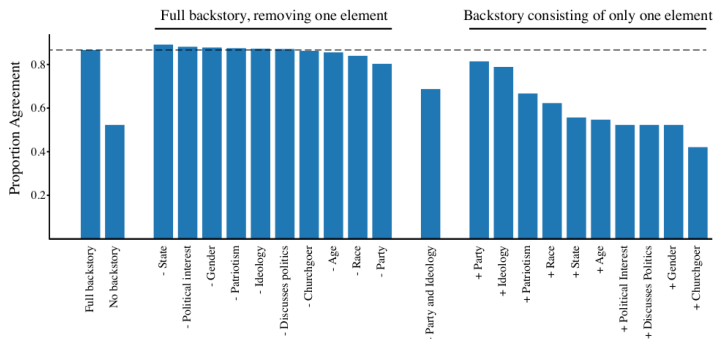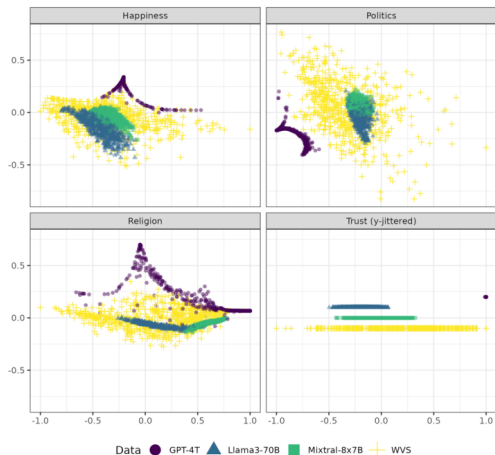


Figure 5: An ablation experiment examining the importance of each backstory element. Reported is the Proportion Agreement on the vote prediction task of the ANES 2016 dataset. Each bar represents a different template with some set of backstory elements,

# 3/ Social science applications

**Subject simulation**

**Boelaert et al (2025)** show that generative LLMs are bad at simulating non-trivial poll responses (socially random answers)

# That's all, folks!

Questions and comments: `julien.boelaert@univ-lille.fr`