

LLMs for Preparing Data in R

Guest lecture for ETC5512: Wild Caught Data

Cynthia A. Huang

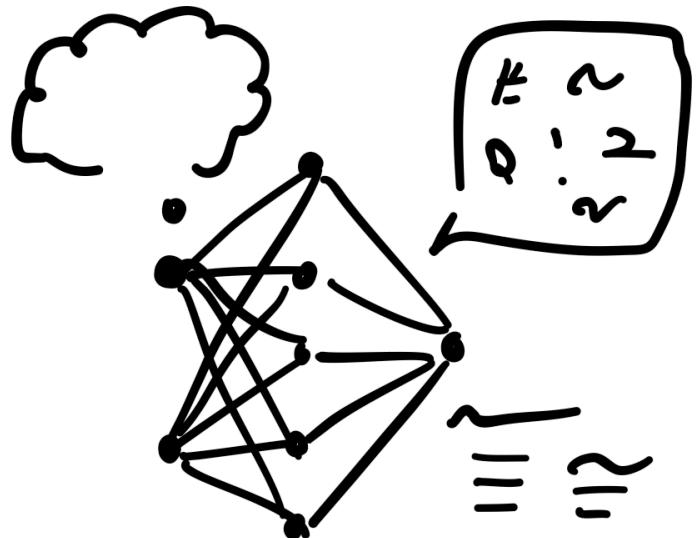
19 May 2025

Introduction

Learning goals

- develop your understanding of:
 - what LLMs are,
 - different types of “wild caught data” tasks that LLMs can help with,
 - how to use and check LLMs for specific data preparation tasks,
 - how to interact with LLMs in R using {ellmer}

Lecture outline



- About Me!
- About LLMs
- Using LLMs for “Wild Caught Data” tasks
- Using LLMs in R with `{ellmer}`



ANTHROP\IC

About Me!

Who am I?

What do I work on?

About Large Language Models

Generative AI and LLMs

Generative AI refers to:

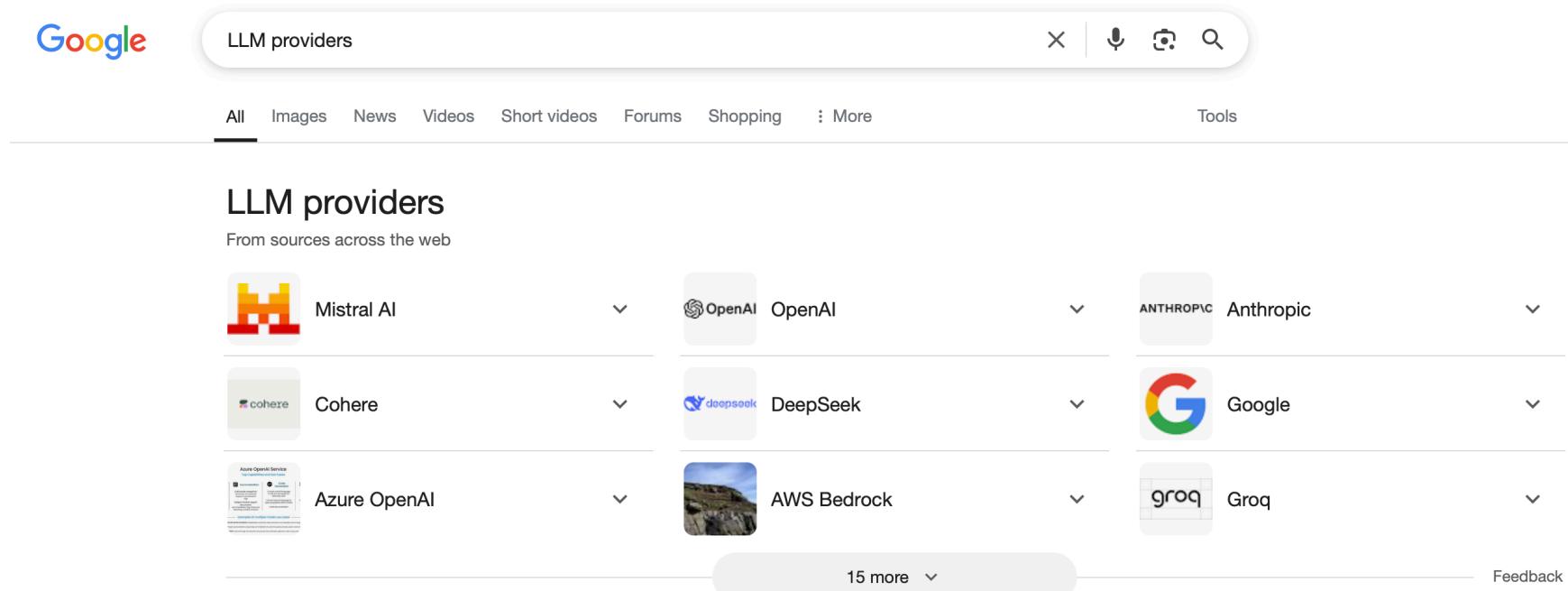
- computer algorithms and systems
- that can generate content such as text, images and sound
- based on patterns learnt from existing data

What are Large Language Models?

LLMs are...

Who makes LLMs?

LLM providers develop and offer access to large language models and systems



A screenshot of a Google search results page for the query "LLM providers". The search bar at the top shows the query. Below it, the standard Google search navigation bar includes "All", "Images", "News", "Videos", "Short videos", "Forums", "Shopping", "More", and "Tools". The main content area is titled "LLM providers" and "From sources across the web". It lists several LLM providers with their logos and names: Mistral AI, OpenAI, Anthropic, Cohere, DeepSeek, Google, Azure OpenAI, AWS Bedrock, and Groq. Each provider entry has a dropdown arrow to its right. At the bottom of the list is a button labeled "15 more" with a dropdown arrow. On the far right, there is a "Feedback" link.

Provider	Logo
Mistral AI	
OpenAI	
Anthropic	
Cohere	
DeepSeek	
Google	
Azure OpenAI	
AWS Bedrock	
Groq	

Why are there so many different models?

LLM providers offer paid and free access to multiple models:

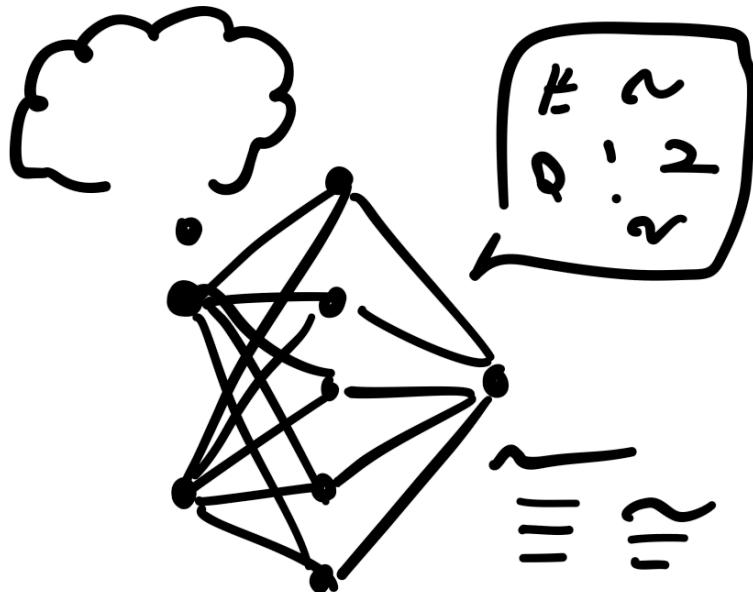
- **OpenAI**: GPT-3 and 4, o-series,
- **Anthropic**: Claude Haiku, Sonnet and Opus
- **Google**: Gemini Flash and Pro
- **Meta**: Llama 3, Llama 4
- **Alibaba**: Gwen 2.5, 3, Max, Plus and Turbo

Model differentiation

Learn more about picking the right tool:

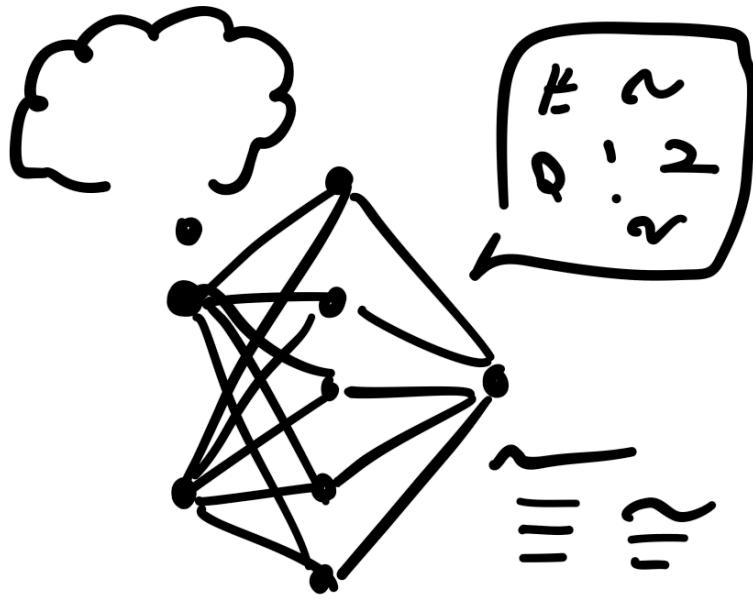
-  [Beyond ChatGPT: THE RAPIDLY EVOLVING LANDSCAPE OF AI](#)
-  [Suggestions on provider/model choice in ‘ellmer’ docs](#)
-  [OpenAI Model Selection Guide](#)
-  [Anthropic Claude Model Comparison Table](#)

How can we interact with LLMs?



ANTHROPIC

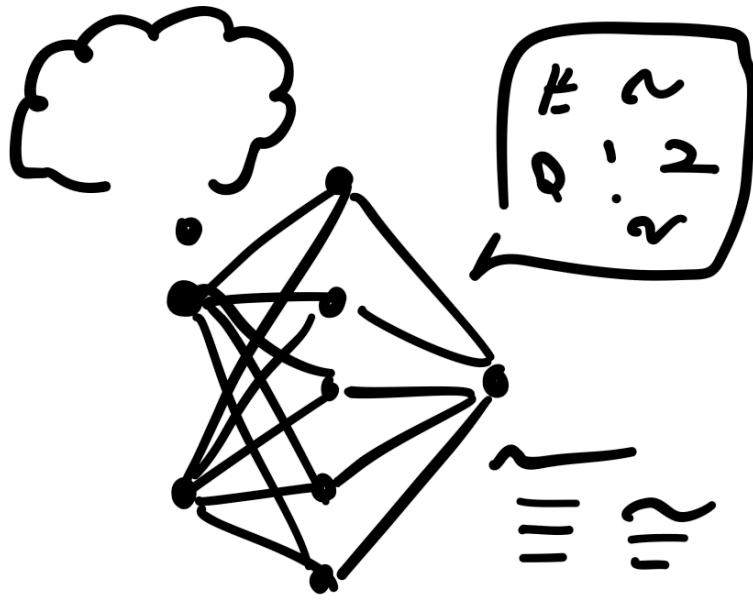
How can we interact with LLMs?



ANTHROPIC

- 1. Web-based Chat Interface**
- 2. Programmatic Access**
- 3. Other interfaces**
 - mobile (chat) applications
 - voice assistants
 - embedded LLMs (e.g. suggestions in Gmail)

How can we interact with LLMs?



ANTHROPIC

- 1. Web-based Chat Interface**
- 2. Programmatic Access**
- 3. Other interfaces**

Today, we will use web-based chat via:

- [chatgpt.com](#)
- [claude.ai](#),

as well as programmatic access via:

- [{ellmer}](#)

How do we verify what LLMs are doing?

For a dishwasher, we consider:

Using LLMs for “Wild Caught Data” Tasks

What are common WCD tasks?



Discuss in groups [5 mins]

- What tasks are involved in preparing Wild Caught Data for analysis? *List at least 3.*
- Which of these might be more or less suitable for addressing with LLMs? Why?

Using LLMs for Data Preparation Tasks

Using LLMs for Data Preparation Tasks

- 1. Generating data wrangling code**
- 2. Transformating and generating new data**

Modifying and augmenting data

Prompts for Data Preparation Tasks

Starting Prompts

- **'look up' facts:** “What nationality is *the author* <author name>?”
- **suggesting values:** “What is the likely volume unit *of a beverage* <can> *of* with a volume of <300>?”
- **harmonise different abbreviations:** “Convert the following list of Australian states to all use three-letter state codes: <list>”
- **comparison:** “How similar are these two occupations: <occupation A>, <occupation B>?”
- **summary:** “Summarise the following survey response: <text>”
- **extraction:** “What movie is the follow review about”: <review text>

Requesting different output formats

Verifying success

Verification is the most important skill when using LLMs, but it requires:

- Clearly defined tasks and expected outcomes
- Ways of checking the outcomes have been achieved



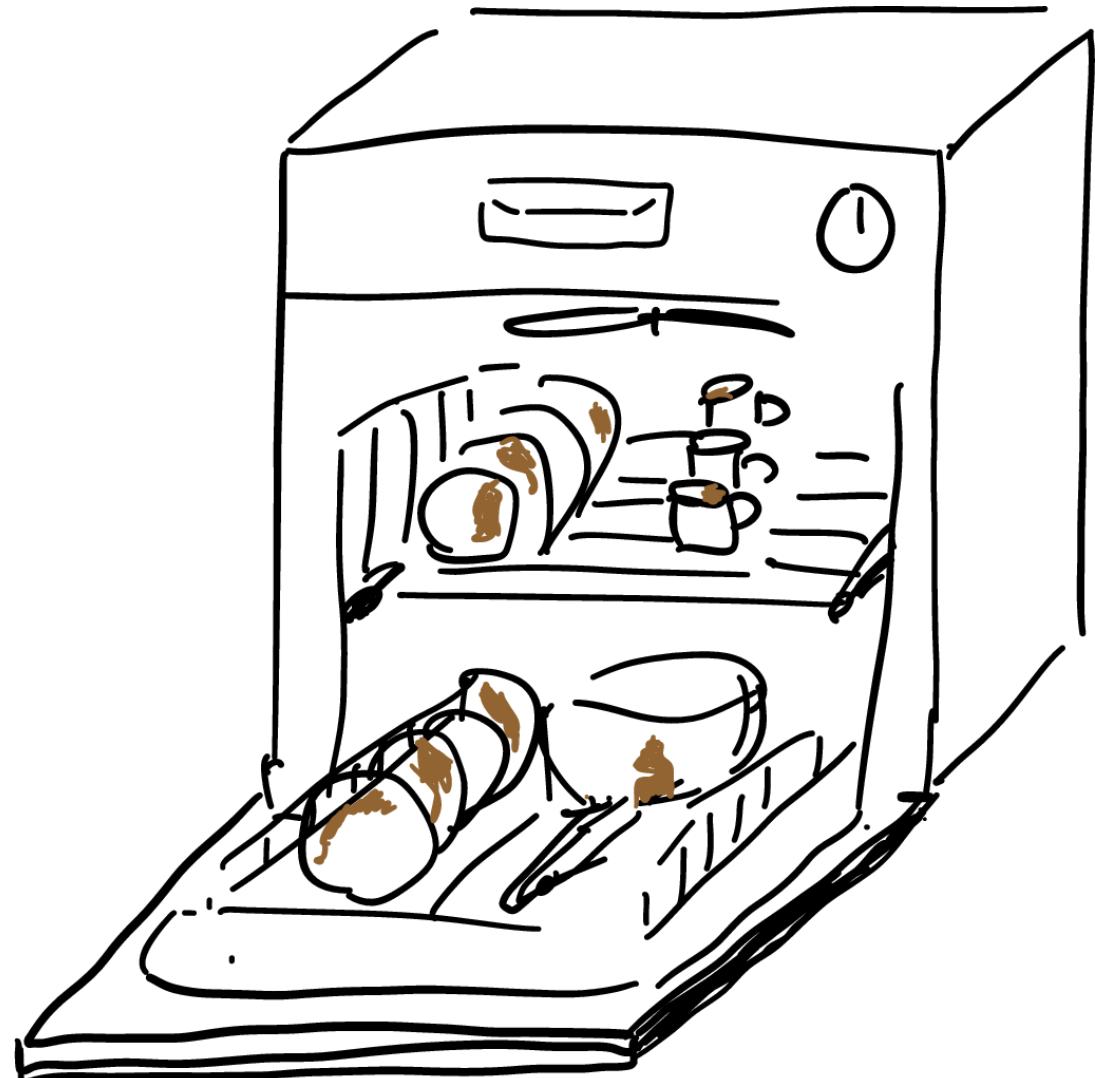
Approaches to verification

There are multiple ways to verify outcomes match expectations.

Checking on the dishes

There are multiple ways to verify outcomes match expectations.

- **Positive verification:**
- **Negative verification:**
- **Trust-based verification:**



Using LLMs in R with {ellmer}

Beyond web-based interfaces

Different interfaces mean different data preparation workflows:

Using the `{ellmer}` R package we can:

Connecting {ellmer} to an LLM

The basic steps:

LIVE DEMO: Chatting via {ellmer}

```
1 library(ellmer)
2 ## A session is like a chat conversation
3 session <- chat_openai()
4
5 question <- "
6   How can I pick a random letter from A-Z.
7 "
8
9 ## send a question to the 'chat'
10 session$chat(question)
11
12 ## clarify your request
13 session$chat("Return R code only")
14
15 ## inspect all turns in the session so far
16 session
```



LIVE DEMO: System Prompts

```
1 library(ellmer)
2
3 session_tidy_expert <- chat_openai(system_prompt =
4   You are an expert R programmer
5   who prefers the tidyverse.
6   Only return code without explanation.
7   ")
8
9 session_tidy_expert$chat(question)
10
11 session_tidy_expert
```

Example adapted from [ellmer docs](#)



Sessions and system prompts

Revisiting Author Nationalities

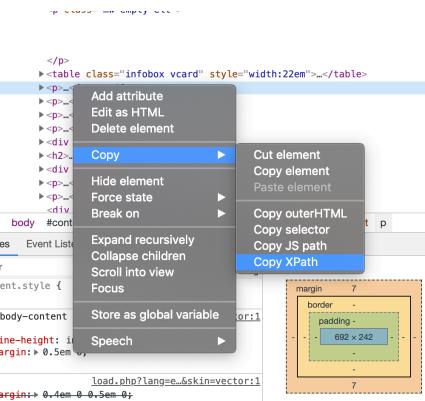
Jane Austen

From Wikipedia, the free encyclopedia

 692 × 242 · confused with Jane G. Austin.

Jane Austen (16 December 1775 – 18 July 1817) was an English novelist known primarily for her six major novels, which interpret, critique and comment upon the British landed gentry at the end of the 18th century. Austen's plots often explore the dependence of women on marriage in the pursuit of favourable social standing and economic security. Her works critique the novels of sensibility of the second half of the 18th century and are part of the transition to 19th-century literary realism.^{[2][3]} Her use of biting irony, along with her realism, humour, and social commentary, have long earned her acclaim among critics, scholars, and popular audiences alike.^[4]

With the publications of *Sense and Sensibility* (1811), *Pride and Prejudice* (1813), *Mansfield Park* (1814) and *Emma* (1816), she achieved success as a published writer. She wrote two additional novels, *Northanger Abbey* and *Persuasion*, both published



WCD 2025, S1, Week 10 - Slide 29

```
# A tibble: 11 x 3
  author_name    author_links
  <chr>          <chr>
  <chr>
  1 Rick Warren  https://en.wikipedia.org/wiki/R... <NA>
  2 William Griffith Wilson
  https://en.wikipedia.org/wiki/B... <NA>
  3 F. Scott Fitzgerald
  https://en.wikipedia.org/wiki/F... <NA>
  4 John Green
  https://en.wikipedia.org/wiki/J... <NA>
  5 Sam McBratney
```

nationality <chr>

LIVE DEMO: Extract Nationalities

```
1 text <- "Jane Austen (/ˈæstɪn, 'ɔ:stɪn/ OST-in, Aw-stin; 16 December 1775 – 18 July 1817) wa
2
3 session_read <- chat_openai(system_prompt = "You are a data entry assistant.")
4
5 session_read <- chat_openai()
6
7 nationality_prompt <- "Nationality of person"
8
9 session_read$extract_data(text, type = type_string(description = nationality_prompt))
10
11 std_prompt <- "Extract structured data of the nationality of person. Return only ISO 3-digit
12
13 session_read$extract_data(text, type = type_string(description = std_prompt))
```

LIVE DEMO: Ask for Nationalities

```
1 library(dplyr)
2
3 author_df <- readr::read_csv('example_data/week10-author_df.csv')
4
5 short_prompt <- "Nationality of person only"
6 session_lib <- chat_openai(system_prompt = "You are a librarian with expert knowledge of pop
7
8 ## let's ask about multiple authors
9 author_df |>
10   tail(6) |>
11   rowwise() |>
12   mutate(nationality_llm =
13     session_lib$clone()$extract_data(author_name,
14                                         type = type_string(short_prompt))
15   )
```

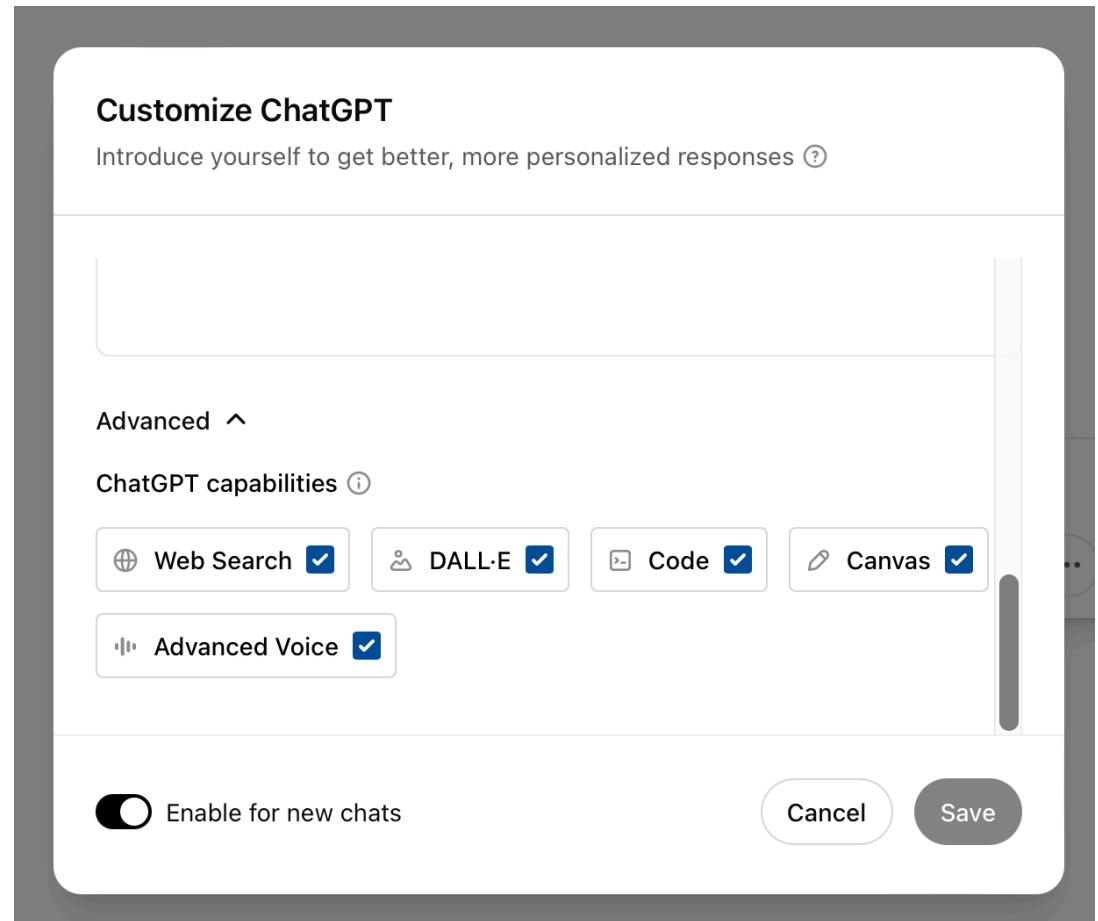


Evaluation through agreement

Another way to verify data quality is via **consensus**.

How much do LLMs ‘know’?

LLM agents, web search and more



Final Comments

Ethics and AI safety

! Generative AI acknowledgement

I used AI in the following ways:

- i. generate definitions and suggested explanations for key concepts covered in this lecture. I used Claude AI to suggest definitions for terms like ‘Generative AI’, and ‘System Prompt’, and to generate lists of “top LLM providers in 2025” and “ways of interacting with LLMs”.

Key takeaways

What we've learnt...

