# Text Analysis with Large Language Models

## An Introduction and Best-Practice Guide

Violeta Haas
violeta.haas@iast.fr

29th January 2025

Postdoctoral Research Fellow
Institute for Advanced Study in Toulouse
Toulouse School of Economics
University of Toulouse Capitole

# Introduction

About myself:

- Violeta Haas
- Research Fellow at IAST and TSE
- BA, MA & PhD at HU and DYNAMICS

Research Interests:

Comparative politics and political behavior: *how institutions and elites shape public perceptions and inter-group relations*, with a special focus on sexuality and gender. I also study electoral behavior, party competition, social norms and movements.

What we'll cover today:

- Short introduction to Large Language Models (LLMs).
- Step-by-Step Guide for text analysis with LLMs.
    - Model Selection
    - Prompt Engeneering and Parameter Tuning
    - Prompt Validation
    - Model Evaluation
- Applied example using GPT API.
- (How to avoid the next Replication Crisis)

Goal: Introduction and best-practice guide for text analysis with LLMs

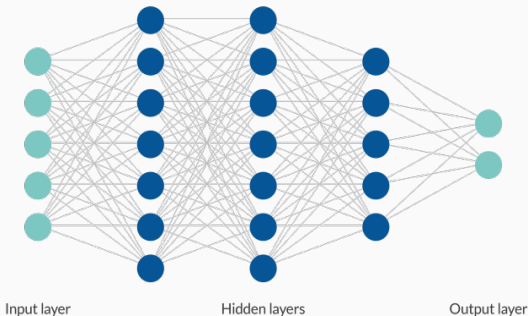Interactive: Feel free to ask questions throughout!

## Two Disclaimers

### 1. Introductory Session:

- This is an *introductory session*—we will only cover the basics.
- Advanced topics like fine-tuning and detailed parameter adjustments are beyond today's scope (but see "Appendix").

### 2. Rapidly Evolving Field:

- LLMs are a *rapidly evolving field.*
- Best-practice advice is often a work in progress.
- What is true today may change tomorrow as new techniques and models emerge.

Input layer    Hidden layers    Output layer

## What are LLMs?

- Deep learning models trained on massive text corpora to understand and perform human-like behavior.
- Built using transformer architectures consisting of neural networks with many layers (hence "deep" learning).
  *Examples:* GPT, BERT, LLaMa and their variants.

Applications in Social Sciences:

- Label topics and frames of documents (e.g., Gilardi et al., 2023)
- Analyzing open-ended responses (e.g., Mellon et al., 2024)
- Act as treatments in experiments (e.g., Argyle et al., 2023)
- Generate data by simulating respondents (e.g., Argyle et al., 2023)
  - → Not uncontested (Bisbee et al., 2024, Dominguez-Olmedo et al., 2023, von der Heyde et al., 2024)

LLMs have outperformed humans across a variety of classification and annotation tasks for a fraction of the cost (Bang et al., 2023, Qin et al., 2023, Goyal et al., 2022, Chiang and Lee, 2023, Grossmann et al., 2023, Ziems et al., 2024).

Ease-of-use, high accuracy and relatively low costs!

1. Model Selection
2. Prompt Engineering and Parameter Tuning
3. Prompt Validation
4. Model Evaluation

## Model Selection: Proprietary vs. Open-Source Models

Two Catgories of LLMs:

- *Propriatary:* LLMs developed and maintained by private entities with exclusive rights to the model's code, data, and usage.
- *Open-Source:* LLMs whose code, architecture, and sometimes training data are publicly accessible.

|  | Proprietary LLMs | Open-Source LLMs |
|---|---|---|
| *Ownership* | Retained by company | Community or developers |
| *Access* | Paid or restricted | Free or low-cost |
| *Transparency* | Limited | Full |
| *Customization* | Limited | High |
| *Examples* | GPT | LLaMA |

## Model Selection: Proprietary Models

Challenges with Proprietary Models:

- *Privacy concerns:* Commercial APIs require sending data to service providers, who may use it for training.
- *Lack of transparency:* Training data and methodologies are unknown, potentially leading to bias.
- *Unpredictable updates:* Models can change over time without notice or depreciate, altering results and making replication impossible.

## Model Selection: Open-Source Models

Challenges of Open-Source Models:

- *Performance Gaps:* May not match proprietary models' performance, especially for complex tasks.
- *Scalability Challenges:* Often not optimized for large-scale deployment.
- *Setup Costs:* May require access to powerful hardware and more expertise in machine learning and model optimization.

## Model Selection: Self-Hosting

Why host models on your own infrastructure?

- Control over model version and updates without relying on external services (improves reproducibility).
- Ensures sensitive or confidential information is handled securely without exposing it to third-party APIs.

## Not all downloadable models are open-source models!



Benchmarks for openness (Liesenfeld et al., 2023)
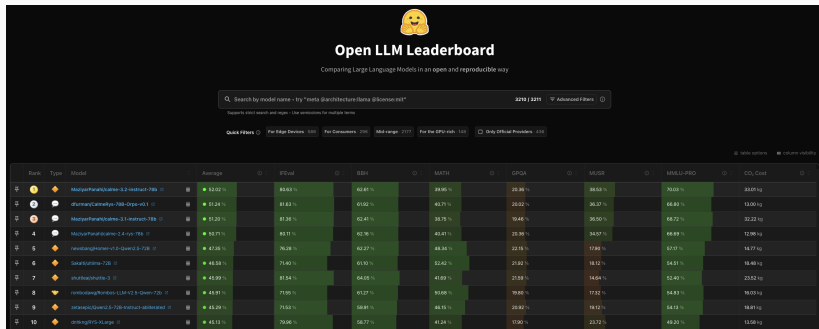Link to Github

Tracks carbon dioxide emission!



Link to Hugging Face

#### 6 Factors to consider (Part 1):

1. *Reproducibility:* Ensure results can be replicated by documenting and using a fixed LLM version.
2. *Ethics and legality:* Adhere to ethical and legal standards, respecting privacy and data regulations (e.g., not storing research data).
3. *Transparency:* Clearly document methodologies, assumptions, and limitations.

6 Factors to consider (Part 2):

4. *Culture and language:* Select models proficient in relevant languages and cultural contexts to avoid bias.
5. *Scalability:* Match the model to data size and resource constraints; consider semi-supervised options for large datasets.
6. *Complexity:* Use sophisticated models (e.g., GPT-4) for tasks requiring advanced reasoning or subtle meaning parsing.

$\rightarrow$ Simple decision tree in "Appendix"

Best Practice:

→ Use *self-hosted, open-source models* with *publicly known training data.*

→ But, *trade-offs* between factors 1-3 and 4-6 can complicate use of open-source models.

→ Always offer *motivations* for model choice.

The best model ultimately depends on the task at hand!

Never use proprietary models with sensitive, secret, or private data!

# Prompt Engineering and Parameter Tuning

**Step 1:** Draft the coding task, with instruction for coders and prompt.

**Step 2:** Take a small random sample and code with both human coder(s) and LLM.

**Step 3:** Compare coding and examine points of disagreement. Ask LLM to motivate its coding decision.

**Step 4:** Revise codebook and/or LLM prompt.

**Step 5:** Code a larger validation dataset and measure final LLM performance.

Repeat until desired performance achieved

Figure 1: Example of a systematic coding procedure.

Source: Törnberg (2024)

## Prompt Engineering

**What is Prompt Engineering?**

- Designing effective instructions (inputs) to guide LLM behavior.
- Requires multiple iterations of modification and testing.
- Poor theor. understanding of why some techniques work better.

**A good prompt contains these five elements (in this order):**

1. Context
2. Question/Instruction
3. Constraints
4. (Examples)
5. Input data

**Include background information:**

- Add context that helps the model "understand" its *role* and the *domain* of the input data.

### Example: Context

*You are an expert annotator for natural language processing tasks focusing on social media content analysis. Your role involves fact-checking Twitter messages related to the US 2020 election.*

Start with a simple Question/Instruction

- Define Concepts and Tasks
- Use direct and instructional language.

> **Example: Instruction**
>
> *Your expertise is critical in identifying misinformation that shares false or misleading content, which can distort facts, omit critical context, or misrepresent the truth. For each Twitter message, assess whether it contains misinformation regarding the US 2020 election by evaluating the factual accuracy of claims, sources, and context. Does this message contain misinformation?*

**Adapt specificity based on performance:**

- Provide additional guidelines

**Example: Additional Guidelines**

*Verify statements against reputable and publicly available sources like government websites, credible news outlets, and fact-checking organizations.*

- Provide specific instructions where the model fails

**Example: Failure Adaptation**

*Look for signs of altered or fabricated images, videos, or text like deepfakes, fake screenshots, or doctored evidence.*

# 3. Constraints

## Provide Categories and Constraints

- Give Classification Options
- Allow for Uncertainty

### Example: Options

*Options: Yes/No/Uncertain. If you are uncertain about the classification, choose 'Uncertain' and provide a rationale for this uncertainty.*

- Use JSON output format

### Example: Output format

*Provide your response in JSON format, as follows: { "contains_misinformation": "Yes/No/Uncertain", "justification": "Provide a brief justification for your choice." }*

*Zero-shot prompting might not always deliver good results!*

Improve prompt using in-context learning techniques:

- *Zero-shot prompting*: instruction without examples.
- *One-shot prompting:* instruction + one example.
- *Few-shot prompting:* instruction + two or more examples.
- *Negative prompting:* instruction + example of unwanted output

*For an Example Structure of each technique see "Appendix"

**Insert Input Data:**

- One-by-one or grouping multiple queries (i.e., Batching)

| Example: Input Data |
|---|
| *Twitter message: [MESSAGE]* |

\*For Batching see "Appendix"

## Parameter Tuning: Temperature

### What is Temperature?

- Controls randomness in output.
- Range: 0 (deterministic) to 1 (creative and diverse).
    - Low temperature (e.g., 0.0 or 0.2): Ensures consistency and accuracy.
    - High temperature (e.g., 0.8): Useful for subjective or exploratory tasks.

$\rightarrow$ Use low temperature for annotation tasks to ensure that similar inputs produce similar outputs (*consistency*) and that the model focuses on the most probable completion (*accuracy*).

*For other parameter settings, see section "Advanced"

**Me saying "thank you" to ChatGPT, so they spare my life when AI takes over the world**

$\rightarrow$ Impolite prompts reduce performance, while overly polite ones offer no guarantee of improvement.

$\rightarrow$ optimal politeness varies by language/culture (Yin et al., 2024).

**Additional Tips for Prompting:**

- Give an "I don't know" option.
- Break down tasks into several simpler intermediate steps.
- Use LLMs (like ChatGPT) to improve prompt language.
- Use affirmations (e.g., "do" and "don't" or capitalize "you MUST").
- Use delimiters for structuring (e.g., ###Context###, …).
- Request output in JSON format.
- Repeat important phrases multiple times.
- Avoid long prompts (i.e., balance brevity and specificity).

**Some Additional Resources:**

- Prompt Engineering Guide by DAIR.AI
- Prompt engineering by OpenAI
- Brex's Prompt Engineering Guide

## Two Common Problems

#### Hallucination:

- LLMs may generate incorrect or fabricated information, potentially compromising the validity of results.

#### "Guardrails":

- Models tuned to avoid controversy may refuse to annotate certain topics, posing challenges for tasks involving controversial content (e.g., radical political messages) or sensitive annotations (e.g., identifying an author's gender).

# Prompt Validation

Validate your prompt before undertaking any model validation steps, such as manually coding substantial amounts of data!

$\rightarrow$ Prompt validation safes ⏱ and 💵

## Prompt Validation

Asses Prompt Performance:

- Compare results against labeled data (i.e., small hand-coding sample) using classic performance metrics.
- Identify patterns of errors:
    - Does the output align with task expectations?
    - Are similar inputs yielding similar outputs?
    - Are the outputs unbiased and ethical?

LLMs can be very sensitive to trivial things like punctuation!

**Calculate Stability Scores:**

- Create several paraphrases of the prompt and run the analysis on a subset of the data.
- Estimate stability by comparing results using measures like Krippendorff's Alpha.
- → Barrie et al. (2024) have released a package for prompt stability scoring (Link to GitHub).

# Model Evaluation

### Must take place after prompt has been finalized!

#### How to evaluate LLM model outputs?

1. Label sufficient number of texts (min. 20–30 units per category)
   - Use larger validation datasets for imbalanced categories or high-stakes tasks
   - Power analysis recommended
2. Compare model results to manual labels using multiple performance metrics (e.g., Accuracy, F1 Score, Cohen's Kappa)
   - Validate across diverse subsets (e.g., languages, cultural contexts)
3. Examine and explain failures
   - Asses downstream risks for later analysis.

Applied Example 💻

# Replication with LLMs

"Research tools are not merely passive instruments, but active participants in research procedures (Latour & Woolgar, 2013)."

Exact replication is impossible with LLMs (same for hand-coding)!

Best Practice:

- Give motivation for model choice.
- Provide exact model version and parameter settings used.
- Save all versions and performance metrics of optimized prompts.
- Share random sample of input/output pairs as additional benchmarks.
- Store and share code/scripts for prompt automation.

Open Zotero Library: "LLMs in Social Sciences"

End of session 🎉
Thank you!

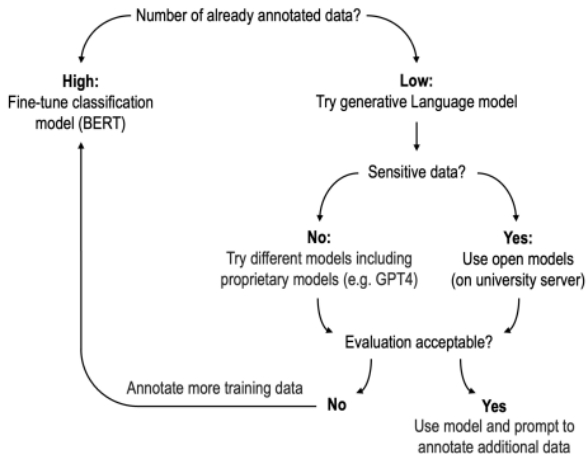*This presentation was designed with the help of ChatGPT

# Appendix

Figure 1: Decision tree for the use of generative LLM for text annotation

Source: Weber and Reichardt (2024)

## Advanced: Chain-of-Thought (CoT) Technique

### What is CoT?

- A technique for improving reasoning and instruction-following in LLMs.
- Breaks tasks into simpler intermediate steps to mimic human problem-solving.

### Why use CoT?

- Elicits step-by-step reasoning (Wei, Wang, et al., 2024).
- Enhances performance on complex tasks (Chung et al., 2024).
- Can be triggered with a prefix such as: *"Let's think step by step."*

## Other Parameters

### Top-P/K:

- Focuses on diversity by setting a probability- or top-K most probable tokens threshold for selecting tokens (requires temperature > 0)
  - Low Top-P/K (e.g., 0.2 to 0.4/ e.g., 5 or 10) ensures precision.
  - High Top-P/K promotes diverse outputs but may reduce accuracy.

### Max Tokens:

- Sets a limit on the length of the model's response.
- Helps prevent overly lengthy or irrelevant outputs.

### Best Practices:

- Start with low temperature for consistency.
- Adjust Top-P/K for controlled creativity.
- Use max tokens to prevent overly verbose responses.

*What exactly are Tokens? See Appendix

## Advanced: Fine-Tuning

### What is Fine-Tuning?

- Adapting a pre-trained model to specific tasks or domains using additional data.
- Requires labeled datasets and computational resources.

### When to Consider Fine-Tuning:

- Highly domain-specific tasks (e.g., medical or legal text).
- Need for custom outputs not achievable with prompt engineering.

## Advanced: Batching

### What is Batching?

- Grouping multiple input queries together for simultaneous processing.
- Reduces latency and improves throughput.
- Often used in high-volume workflows or API integrations.

### Advantages of Batching:

- **Efficiency:** Minimizes overhead for repeated queries.
- **Cost-Effectiveness:** Lowers cost per request when using APIs.
- **Consistency:** Ensures uniform settings across similar tasks.

# Example Structure for in-context learning techniques

| Prompting Strategie | Example Structure |
|---|---|
| Zero-shot | {"role": "system", "content": "Text of System Prompt"}, <br> {"role": "user", "content": "(Text to classify) + classification question"} |
| One-shot | {"role": "system", "content": "Text of System Prompt"}, <br> {"role": "user", "content": "(Example text) + classification question"}, <br> {"role": "assistant", "content": "Example classification"}, <br> {"role": "user", "content": "(Text to classify) + classification question"} |
| Few-shot | {"role": "system", "content": "Text of System Prompt"}, <br> {"role": "user", "content": "(Example text) + classification question"}, <br> {"role": "assistant", "content": "Example classification"}, <br> {"role": "user", "content": "(Example text) + classification question"}, <br> {"role": "assistant", "content": "Example classification"}, <br> … more examples <br> {"role": "user", "content": "(Text to classify) + classification question"} |
| Chain-of-Thought | {"role": "system", "content": "Text of System Prompt"}, <br> {"role": "user", "content": "(Text to classify) + reasoning question"}, <br> {"role": "assistant", "content": "Reasoning"}, <br> {"role": "user", "content": "Classification question"} |

Table 1: Example structures of different prompting strategies for annotation purposes

Source: Weber and Reichardt (2024)

*Chain-of-Thought*, see section "Advanced"

# Tokens

## What are Tokens?

The atomic unit for language models is a token. Tokens resemble syllables, with 750 words per 1,000 tokens, and include punctuation, sentence boundaries, and document ends.