# Harnessing the Power of LLMs to Construct Structured Datasets from Unstructured Sources

## David Zarruk

Disclaimer: All views expressed are my own.

# Empirical research stages

1. Dataset construction:

   - Access to raw files:
     - Available online
     - Administrative or bureaucratic process
     - Go to archives and find historical documents

   Unavoidable

   - Data processing:
     - Readily available excel files
     - Press releases, news articles, speeches
     - Historical archives (pdfs, physical paper, etc)
     - Manually-filled survey forms

   Easy processing

   Can take a significant amount of time (50%? 90%?)

   Difficult/slow to process

2. Data analysis
3. Paper writing

# Some examples of costly data processing

- Papers that study impact of speeches/press releases:
    - The Importance of Fed Chair Speeches as a Monetary Policy Tool – Swanson (2023)
    - Shocking language: Understanding the macroeconomic effects of central bank communication – Hansen and McMahon (2016)

- Datasets based on historical archives:
    - The Colonial Origins of Comparative Development: An Empirical Investigation – Acemoglu, Johnson and Robinson (2001)
    - History, Institutions, and Economic Performance: The Legacy of Colonial Land Tenure Systems in India – Banerjee and Lakshmi (2005)

# LLMs to the rescue

- LLMs excel at reading and understanding text and images

- Can easily transcribe images into text and extract relevant information

- Can analyze a text, find relevant parts, and analyze sentiments

# 2 examples

1. Construct a dataset with European Central Bank bulletins
   - Hard data (inflation, unemployment, etc)
   - Qualitative features:
     - Risks
     - Labor market conditions: improving, deterioriating, etc.
     - Financial conditions: tightening, easing, etc



| | report_date_yyyy_mm | hicp_current | hicp_core_current | unemployment_rate_latest | m3_growth_latest | covid_mentions_count |
|---|---|---|---|---|---|---|
| 0 | 2020-01 | 1.3 | 1.3 | 7.5 | 5.6 | 0 |
| 1 | 2020-03 | 1.2 | 1.2 | 7.4 | 5.2 | 19 |
| 2 | 2020-05 | 0.4 | 0.9 | 7.4 | 7.5 | 8 |
| 3 | 2020-06 | 0.1 | 0.9 | 7.3 | 8.3 | 48 |
| 4 | 2020-07 | 0.3 | 0.8 | 7.4 | 8.9 | 28 |
| 5 | 2020-09 | -0.2 | 0.4 | 7.9 | 10.2 | 35 |
| 6 | 2020-10 | -0.3 | 0.2 | 8.1 | 10.4 | 15 |
| 7 | 2020-12 | -0.3 | 0.2 | 8.4 | 10.5 | 29 |
| 8 | 2021-01 | -0.3 | 0.2 | 8.3 | 11.0 | 18 |
| 9 | 2021-03 | 0.9 | 1.1 | 8.1 | 12.5 | 12 |
| 10 | 2021-03 | 1.3 | 0.9 | 8.3 | 12.3 | 2 |
| 11 | 2021-06 | 2.0 | 0.9 | 8.0 | 9.2 | 8 |
| 12 | 2021-07 | 1.9 | 0.9 | 7.9 | 8.4 | 18 |

# 2 examples

2. Construct a dataset based on hand-written forms
   - Elections in Colombia

# Tips for good prompt engineering

1. Enclose sections in the prompt using <> keys

2. Define a set of rules that define boundaries for the LLM:
   - "Do not hallucinate"
   - "Return a JSON – nothing else"

3. Give examples to the LLM!

4. Ask the LLM to return citations or the thought process that support the output

5. ITERATE: writing a good prompt is an iterative process – usually the LLM never gets it right on the first attempt