

# {klaus} for content analysis

1.7.2025, GöAID  
Cornelius Puschmann, U Bremen

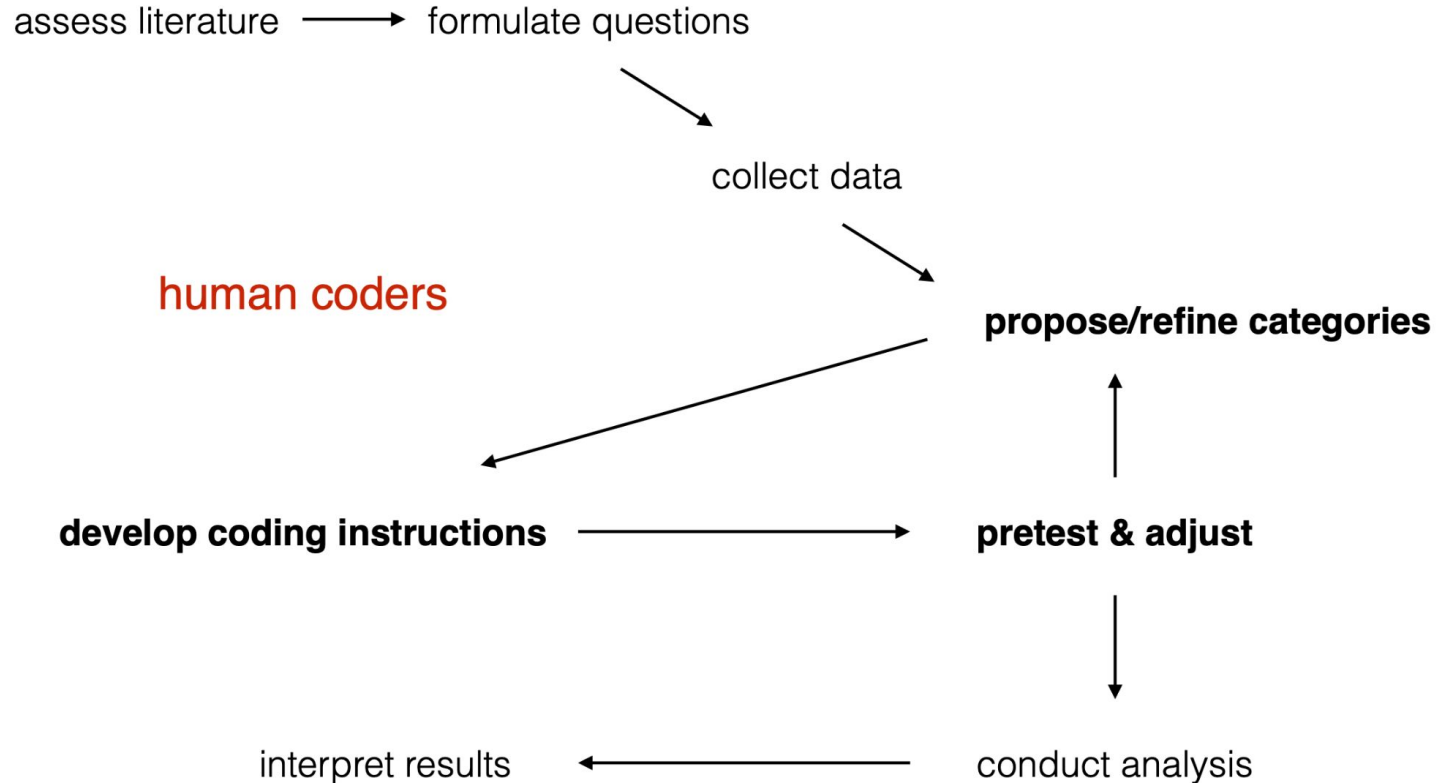
# What is content analysis?



“Content analysis is a research technique for making replicable and valid inferences from texts or other meaningful materials to the context of their use.”

(Krippendorff, 2004)

# The CA workflow



*Fairness and Election Coverage. Presidential Race, 1996*

v1. Story identification number \_\_\_\_\_

v2. Story day (month and date) \_\_\_\_\_

v3. Newspaper name and circulation rank \_\_\_\_\_

*General Story Characteristics*

v4. Story prominence (FP=2; SP1=1; Inside=0)

v5. Story origin

- |                            |                            |
|----------------------------|----------------------------|
| 1=newspaper's own reporter | 2=newspaper's state bureau |
| 3=newspaper's D.C. bureau  | 4=AP                       |
| 5=other bureau/newspapers  |                            |

v6. AP unique (if v5=4) \_\_\_\_\_

- |                   |                             |
|-------------------|-----------------------------|
| 0=duplicate story | 1=unique story, blank=no AP |
|-------------------|-----------------------------|

v7. Partisan side sourcing effort \_\_\_\_\_

- |                              |                             |
|------------------------------|-----------------------------|
| 1=only Clinton or Dole       | 2=both Clinton and Dole     |
| 3=only Perot                 | 4=Perot and Clinton or Dole |
| 5=Perot and Clinton and Dole |                             |

v8. Primary story source \_\_\_\_\_

- |             |                         |
|-------------|-------------------------|
| 1=debate    | 2=speech/rally          |
| 3=document  | 4=poll                  |
| 5=interview | 6=other media interview |
| 7=ad        | 8=other                 |



Sample coding categories



O	P	Q	R	S
Presence of KK	Presence of Others	Visible Interaction with Other People	National/European Symbolism	Emotional Appeal
Yes ▾	Yes ▾	Yes ▾	European ▾	Determin... ▾
Yes ▾	Yes ▾	Yes ▾	European ▾	Determin... ▾
Yes ▾	Yes ▾	Yes ▾	European ▾	Empathy ▾
Yes ▾	Yes ▾	Yes ▾	European ▾	Hope/Opti... ▾

# A light-weight R package for LLM-based content analysis



Modern large language models (LLMs) offer considerable advantages for standardized content analysis. `klaus` facilitates use of both proprietary and open source LLMs by offering a simple interface through which to serve data and apply categorization. Presently the package supports the proprietary APIs of OpenAI, Anthropic and Google, as well as for local use via [ollama](#) (through [tidyllm](#)). In addition, for academic research, it is also possible to use the non-commercial [ChatAI API](#) service provided by [GWDG](#) or [Blablador](#) provided by the [Forschungszentrum Jülich](#).

klaus is largely a convenience wrapper for the SAIA API and several commercial APIs via {tidyllm}

[chatai](#)

[chatai models](#)

[code content](#)

[parse codebook](#)

Interact with the Chat AI Completion Endpoint

List Available Chat AI Models

Code Text Content Using an API and a Codebook

Parse a Codebook into JSON

codebook:

category (character) ▾	label (character) ▾	instructions (character) ▾
sentiment	positive	Code this if the sentiment of the tweet is positive
sentiment	negative	Code this if the sentiment of the tweet is negative
sentiment	neutral	Code this if the sentiment of the tweet is neutral



```
data_to_code <- readr::read_csv("data_sentiment.csv")
```

```
general_instructions <- "You are a highly accurate and consistent text  
classification model that specializes in analyzing English-language Twitter posts.  
Your task is to determine the sentiment of the tweet reproduced below. You must  
strictly follow the classification rules without deviation. Do not return any  
additional information outside the classification scheme. Use JSON."
```

```
formatting_instructions <- "Always return a single JSON object with the category  
name as the key for each coded text. The value should be an object containing a  
'label' key and a single value among multiple options. Each JSON object should have  
the following structure:"
```

```
codebook <- readr::read_csv("codebook_sentiment.csv")
```

```
coded_data_chatai <- code_content(data_to_code, general_instructions,  
formatting_instructions, codebook, provider = "chatai", model =  
"llama-3.3-70b-instruct")
```

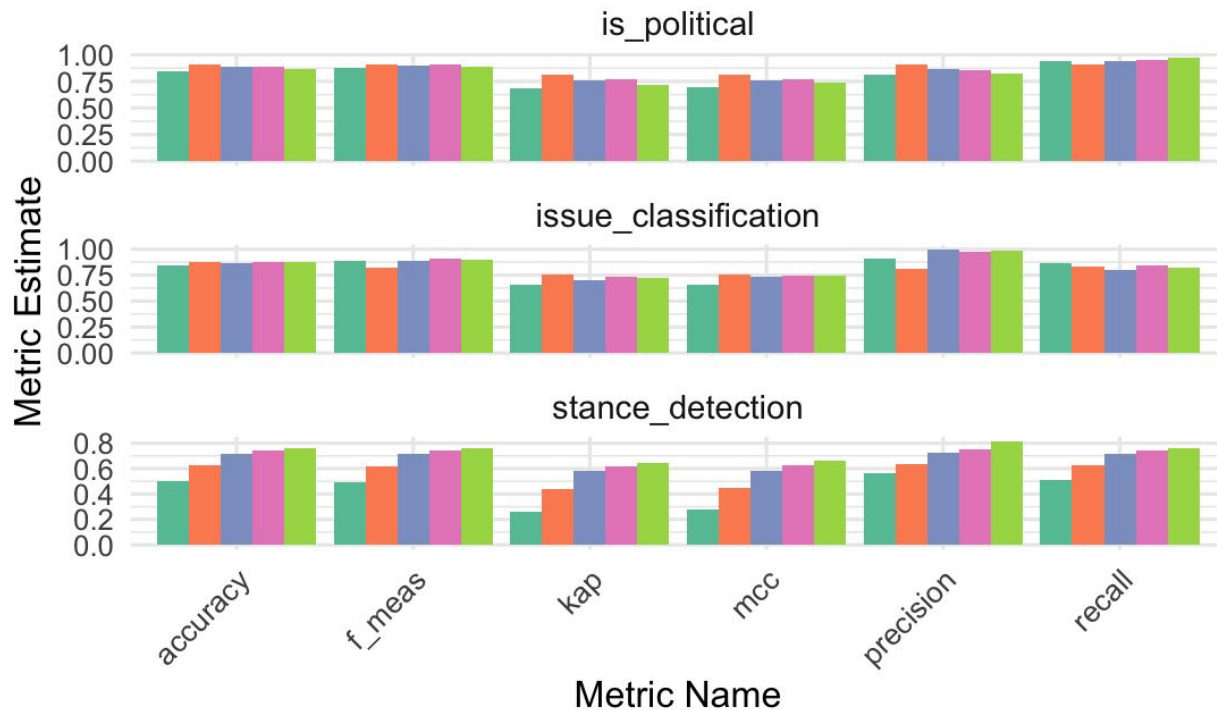
```
Coding data with claude (4 rows)
Iterating over content...
Coding data with claude / claude-3-7-sonnet-20250219...
Processed text 1 of 4
Coding data with claude / claude-3-7-sonnet-20250219...
Processed text 2 of 4
Coding data with claude / claude-3-7-sonnet-20250219...
Processed text 3 of 4
Coding data with claude / claude-3-7-sonnet-20250219...
Processed text 4 of 4
Parsing JSON responses (4 rows).
Done. Joined results for 4 parsed rows.
```

	user	text	category	label
1	peter	It is terrible what is happening to this country	sentiment	negative
2	john	I love Donald Trump	sentiment	positive
3	mary	The stock market has falled dramatically because of t...	sentiment	negative
4	mike	The situation has been calm following the recent midt...	sentiment	neutral

# Model Comparison by Metric

Faceted by Task

Model    ml\_zeroshot    setfit    llama-3.3-70b    gpt4o    claude





**ZeMKI**



Universität  
Bremen

thanks for listening -- any questions?

