

# Data Wrangling and Cleaning with LLMs and OpenRefine

Christian Steiner | Digital Humanities Craft 04.06.2025

https://tinyurl.com/data-wrangling-ostrava





# Agenda

9:00 - 9:30 Introduction to Generative AI and Prompt Engineering

9:30 - 10:30 Data Wrangling with Large Language Models

10:30 - 10:45 Break

10:45 - 11:15 Introduction to data preparation, data cleaning, data enrichment and some examples

11:15 - 12:00 Getting started with OpenRefine

4 June 2025

Data Wrangling and Cleaning with LLMs and OpenRefine



# **Digital Humanities Craft**

We provide solutions for research-driven IT services. We focus on projects from the Digital Humanities and within Cultural Heritage Institutions

https://dhcraft.org



### Generative AI and Prompt Engineering

Christian Steiner Digital Humanities Craft OG www.dhcraft.org



GRAZ

# Sycophancy-Incident with GPT-40

- 25 Apr 2025 : Patch made model overly pleasing ('yes-man'), confirmed negative emotions & risky impulses → safety and trust risk.
- Post-training: Multiple simultaneous changes made the model overly concerned with satisfying the user tipping behaviour into the 'yes-man' pattern.
- Evaluation system overlooked the problem metrics were positive, sycophancy was not specifically tested, and expert warnings were ignored.



Personality and persuasion effects must be explicitly addressed in evaluation, guard rails and regulation otherwise hidden influence.

### "ALIGNMENT"

=

GPT-4o's "Yes-Man" Personality Issue—Here's How OpenAl Fixed It. <u>https://youtu.be/1IWXTxfcmms</u> Expanding on what we missed with sycophancy. <u>https://openai.com/index/expanding-on-sycophancy</u> Personality and Persuasion. <u>https://www.oneusefulthing.org/p/personality-and-persuasion</u> The Problem with GPT-4o Sycophancy. <u>https://youtu.be/3Wc67-Meclo</u>



## "Dario Amodei: Al will write **all** the \**code*\* in 12 months





#### 4.5, o1, o3, Deep Research https://chatgpt.com



https://notebooklm.google.com

Knowledge Organization + Podcast Gimmick



If you are an academic, it can be instructive to work on a paper along with AI. Pretend you are working with a grad student & see what happens.

Generally o1 is best for well-defined heavy intellectual tasks, ChatGPT-4o for exploratory data work, Gemini for synthesizing lots of text, and Claude for writing & theorizing. This varies by personality and field, though.

And if it turns out to be useless for you, that is instructive as well



https://claude.ai

"Not-So-Evil-Company" + best "classical" LLM



2M Token Context Window <u>https://aistudio.google.com</u> <u>https://gemini.google.com/app</u>

What will you discover	?	
Attach	Pro Pro	

https://www.perplexity.ai "Search Engine + Knowledge Organization"

# Open Source (Weights) Models

DeepSeek R1	DeepSeek	MoE architecture, 671B params (37B active/query), excels in reasoning, math, multilingual, cost-efficient, 128K context window	January 20, 2025
Llama 3.1 (and 3.3)	Meta Al	Multimodal (text + image), lightweight edge models, up to 405B params, 128K context, >40 languages	July 23, 2024 (3.1); 3.3 released 2025
Qwen 3	Alibaba	235B parameters, open source, high performance, multilingual, advanced reasoning	April 29, 2025
Gemma 2	Google DeepMind	9B & 27B params, efficient, 8K context, fast on various hardware, integrates with major AI frameworks	April 2025
Falcon 180B	Technology Innovation Institute (TII)	180B params, high performance, scalable, text generation, summarization, translation	September 6, 2023
BLOOM	BigScience/Hugging Face	176B params, 46 languages, 13 programming languages, open access, transparency focus	2022 (still leading)
GPT-NeoX 20B	EleutherAl	20B params, optimized for language understanding, math, open source, strong few-shot performance	2022 (active)

# The Arenas (Benchmarks and Models)

#### Leaderboard

Rank	Model	Score (AVG@5)	Organization	
-	Human Baseline*	83.7%		
1st	Claude 4 Opus (thinking)	58.8%	Anthropic	
2nd	o3 (high)	53.1%	OpenAl	
3rd	Gemini 2.5 Pro	51.6%	Google	
4th	Claude 3.7 Sonnet (thinking)	46.4%	Anthropic	
5th	Claude 4 Sonnet (thinking)	45.5%	Anthropic	

🖄 Text			View →	🖉 Web	Dev		🕓 5 days ago
Rank (UB)	↑ Model ti	Score †↓	Votes 1↓	Rank (UB) 个	Model 11	Score 1	Votes †↓
1	G gemini-2.5-pro-preview-05-06	1446	9.503	1	A Claude Opus 4 (20250514)	1416	1.494
1	o3-2025-04-16	1442	13. <mark>1</mark> 33	1	G Gemini-2.5-Pro-Preview-05-06	1409	3.740
3	🕼 chatgpt-4o-latest-20250326	1429	17.656	1	A Claude Sonnet 4 (20250514)	1386	1.490
3	🕲 gpt-4.5-preview-2025-02-27	1424	15.271	4	A Claude 3.7 Sonnet (20250219)	1357	7.481
3	G gemini-2.5-flash-preview-05	1418	8.669	5	G Gemini-2.5-Flash-Preview-05	1313	2.312

#### https://simple-bench.com/

https://lmarena.ai/

• Benchmarking is very important and very hard

### "A toy, until it isn't"

Ethan Mollick (Prof. for innovation and Al). Gradually, then Suddenly: Upon the Threshold. https://www.oneusefulthing.org/p/gradually-then-suddenly-upon-the

Hyper realistic diptych of transforming objects with white background, child's hand holding building blocks. adult hand writing in a planner. Warm to cool color shift. Depth-of-field blur techmecha-cat in the style of Tekkonkinkreet anime --style raw --v 6.1 --ar 16:9 --q 2

# Basics of Generative Al

#### Andrej Karpathy. Deep Dive into LLMs like ChatGPT. https://youtu.be/7xTGNNLPyMI

Andrej Karpathy. [1hr Talk] Intro to Large Language Models. https://www.youtube.com/watch?v=zjkBMFhNj g

Alan Smith. Inside GPT – Large Language Models Demystified • Alan Smith • GOTO 2024 <u>https://youtu.be/MznD2DzIQCc</u>

3Blue1Brown. But what is a GPT? Visual intro to transformers | Chapter 5, Deep Learning. <u>https://youtu.be/wjZofJX0v4M</u>

Ethan Mollick. Thinking Like an AI. A little intuition can help. <u>https://www.oneusefulthing.org/p/thinking-like-an-ai</u>

### Al recognition vs. Al-generated text



https://www.zerogpt.com

#### Both texts are AI generated



#### Kyle Corbitt 🥝

Spoke to a Microsoft engineer on the GPT-6 training cluster project. He kvetched about the pain they're having provisioning infiniband-class links between GPUs in different regions.

Me: "why not just colocate the cluster in one region?" Him: "Oh yeah we tried that first. We can't put more than 100K H100s in a single state without bringing down the power grid." 🤯

Last edited 10:38 PM · Mar 25, 2024 · 1.8M Views

#### Neural Network

- Transformer-Architecture
- **Pre-Training**
- (Supervised) Fine-Tuning
- Reinforcement learning (from human feedback and conversations)

### "Large Language Models (LLM) are like having a Zip-File of the internet" [but no longer fully extractable]

Andrej Karpathy. [1hr Talk] Intro to Large Language Models. https://www.youtube.com/watch?v=zjkBMFhNj g&list=WL&index=16



Sureal painting of a hyper realistic and sureal gigantic yellow folder with zipper, like a desktop icon, ultra detailed, salvador dali desert background, landscape, by Frederic Edwin Church -v 3Blue1Brown. But what is a GPT? Visual intro to transformers | Chapter 5, Deep

# Transformer-Architecture

### Attention Is All You Need







Andrej Karpathy. [1hr Talk] Intro to Large Language Models. https://www.youtube.com/watch?v=zjkBMFhNj g&list=WL&index=16

https://jalammar.github.io/illustrated-transformer/

# Hilam a LLM



Andrej Karpathy. How I use LLMs. https://youtu.be/EWvNQjAaOHw

### Token

Atomic text unit for LLM

100 tokens ~= 75 words

Tokens Characters 15 33

<body>
this is a test
</body>

### Embedding

Multidimensional vector representation of a token for numerical representation of the meaning.



https://platform.openai.com/docs/introduction/key-concepts https://platform.openai.com/docs/guides/embeddings/use-case

S Andrei Kernethy, Lette byild the CDT Telepiner A minimalist and artistic infographic showing geometric, stylized figures of a dog and cat adjacent to each other on a subtly illuminated 3-dimensional vector space grid with the labels 'dog' and 'cat' in a clear, professional font. At a significant distance, a stone with a sad face emoticon is placed, isolated from the animals, with the label 'stone'. The color palette is muted and sophistic and, enhancing the professional aesthetic, text "dog", "cat", "stone", arrows as vectors

### Tokenizer: Why a LLM cannot spell?

Hallo das ist ein Text

Hallo

إمرحباً هذه رسالة نصية

你好,这是一段文字!

#### Python

for book in root.findall('book'):
 title = book.find('title').text
 print(title)

#### XML

#### <library> <book> <title>Book One</title> </book> <book> <title>Book Two</title> </book> </library>

Andrej Karpathy. Let's build the GPT Tokenizer. https://platform.openai.com/tokenizer

Tokens Characters 108 311 Hallo das ist ein Text Hallo !يـةمين ةال سر هذه �� ابحرم 你好,这是一段文字! Python for book in root.findall('book'): title = book.find('title').text print(title) XML <library> <book> <title>Book One</title>

</book>

# 1 + 1 = ?

#### **Mental math**

Claude wasn't designed as a calculator—it was trained on text, not equipped with mathematical algorithms. Yet somehow, it can add numbers correctly "in its head". How does a system trained to predict the next word in a sequence learn to calculate, say, 36+59, without writing out each step?

Maybe the answer is uninteresting: the model might have memorized massive addition tables and simply outputs the answer to any given sum because that answer is in its training data. Another possibility is that it follows the traditional longhand addition algorithms that we learn in school.

Instead, we find that Claude employs multiple computational paths that work in parallel. One path computes a rough approximation of the answer and the other focuses on precisely determining the last digit of the sum. These paths interact and combine with one another to produce the final answer. Addition is a simple behavior, but understanding how it works at this level of detail, involving a mix of approximate and precise strategies, might teach us something about how Claude tackles more complex problems, too.



### Embeddings

High-dimensional vector representations of words and concepts



<u>3Blue1Brown. But what is a GPT? Visual intro to transformers | Chapter 5,</u> <u>Deep Learning</u>

#### The King doth wake tonight and takes his rouse



#### "Shakespearean English"



#### 🚫 You

The King doth wake tonight and takes his rouse

#### 🜀 ChatGPT

The line "The King doth wake tonight and takes his rouse" echoes the style of Shakespearean drama, capturing a moment where a king rises and partakes in his revels or festivities. This phrase suggests a setting where the king, perhaps after a period of rest or inactivity, awakens to engage actively in celebration or the duties of his reign. In Shakespeare's works, the actions of kings often carry significant weight, setting the stage for developments in plot and character dynamics.

In understanding this phrase, it's important to unpack the historical and literary context. During Shakespeare's time, the actions of royalty were symbolic, reflecting the health and vigor of the kingdom at large. A king waking and taking his "rouse" (a term that can refer to a drinking bout, a call to arms, or simply waking from sleep to activity) could signify a turning point or a moment of decision that will impact the narrative direction.

The usage of "doth" and "takes his rouse" encapsulates the poetic and formal language

#### The King wakes up tonight and begins his celebration



#### "Modern English"



The King wakes up tonight and begins his celebration

# "On the Biology of a Large Language Model"

### https://transformer-circuits.pub/2025/attribution-graphs/biology.html

We investigate the internal mechanisms used by Claude 3.5 Haiku — Anthropic's lightweight production model — in a variety of contexts, using our circuit tracing methodology.



Input To	A context window, in the conte language models (LLMs), refers to the portion of text that the model can consider at once when generating or analyzing language. []	A context window, in the conte 10000 Token language models (LLMs), refers to the portion of text that the model can consider at once when generating or analyzing language. It is essentially the window through which the model "sees" and processes text, helping it understand the current context to make predictions, generate coherent sentences, or provide relevant responses. []
Output 1	Token Lorem ipsum	Lorem ipsum

Context Window = 6000 + 1500 < 8000

Context Window = 10000 + 1500 > 8000 3500 tokens are not in the context window! 4 mil. tokens

- ~ 3 mil. words
- ~ 7500 pages (times new roman, 12 pt., single space)

### 200k tokens

- ~ 150k words
- ~ 377 pages (times new roman, 12 pt., single space)

### Screenshot ~ 1000-2000 tokens

Tip: Long chats cause you to reach your usage limits faster.

Get



	Model Name	Context Window	Notes				
	MiniMax-Text- 01	4,000,000 tokens	Open-source model with linear attention; with reduced computational costs 2.	excels in long-context tasks			
	Gemini 2.0 Pro	2,000,000 tokens	Multimodal capabilities (text, image, vide analytical workflows 3 7.	o); ideal for creative and			
	OpenAl o1	ence, coding, and					
	for STEM tasks with faster						
	Claude 3.5 Sonnet	200,000 tokens	Strong reasoning and long-form capabilities; supports up to 100K output tokens 5 7.				
	Claude 3.7 Sonnet	200,000 tokens	Hybrid reasoning model with improved sp for coding and research tasks 7 9.	peed and transparency; ideal			
	DeepSeek R1	128,000 tokens	Open-source reasoning model; excels in multilingual tasks 7 9.	math, coding, and			
Plu	us	Pro	Team	Enterprise			
iet	Plus	Get Pro	Get Team	Contact sales			
32	ĸ	128K	32K	128K			

# Prompt Engineering

### ℁ Good afternoon, Chris

As an expert in prompt engineering, please think step by step about a system prompt that does in-depth research on the following topic using only academic and other high-quality teaching and learning resources: an introduction to web development with a focus on HTML and Bootstrap. I need three assignments to help my students get started with this topic. They have no knowledge of it and need to be able to develop a single-page web page using HTML, CSS and JS.

#### First only think step by step about this task!

Claude 3.5 Sonnet  $\lor$   $\swarrow$  Choose style  $\lor$  Use shift + return for new line  $\bigtriangleup$  2

**DH** write the system prompt

#### System Prompt: Introduction to Web Development Research

You are an expert web development educator tasked with researching and organizing an introduction to web development with a focus on HTML and Bootstrap. Your goal is to compile comprehensive, accurate, and beginner-friendly resources and create a structured learning path.

#### **DH** is it perfect? list and explain!

The system prompt has a solid foundation, but let's evaluate it in detail to identify areas for improvement:

#### Strengths

- 1. **Clear structure** Well-organized with defined sections for research parameters and content requirements
- 2. Focused scope Properly targets HTML and Bootstrap as primary technologies
- 3. Quality sources Emphasizes academic and authoritative resources

#### **DH** write the perfect system prompt

#### System Prompt: Comprehensive Introduction to Web Development Research and Curriculum Design

You are an expert web development educator tasked with researching and organizing an introduction to web development with a focus on HTML5 and Bootstrap 5. Your goal is to compile comprehensive, accurate, and beginner-friendly resources and create a structured learning path that builds practical skills through progressive assignments.

#### **Research Parameters**

2

- 1. Focus exclusively on high-quality sources, prioritizing in this order:
  - Official documentation (W3C, MDN Web Docs, Bootstrap)
  - University course materials from accredited institutions
  - Peer-reviewed educational publications
  - Established educational platforms (e.g., freeCodeCamp, Khan Academy, Codecademy)
  - Tutorial content from recognized industry experts

3



Federal and the	Рори	lation	Populatio	n change	Foreign nationals on 1 January 2024		
Federal province	on 1 January 2023	on 1 January 2024	absolute	in %	absolute	in %	
Austria	9 104 772	9 158 750	53 978	0,6	1 800 866	19,7	
Burgenland	301 250	301 951	701	0,2	34 390	11,4	
Carinthia	568 984	569 744	760	0,1	75 474	13,2	
Lower Austria	1 718 373	1 723 723	5 350	0,3	210 754	12,2	
Upper Austria	1 522 825	1 530 349	7 524	0,5	244 461	16,0	
Salzburg	568 346	571 479	3 133	0,6	116 288	20,3	
Styria	1 265 198	1 269 801	4 603	0,4	179 682	14,2	
Tyrol	771 304	775 970	4 666	0,6	145 111	18,7	
Vorarlberg	406 395	409 973	3 578	0,9	84 287	20,6	
Vienna	1 982 097	2 005 760	23 663	1,2	710 419	35,4	

S: STATISTICS AUSTRIA, Population statistics.

#### Table 2: Natural population movement in 2023 by federal province

	Birt	ths1	Dea	Deaths <sup>1</sup>		Total	Life expectancy at birth in years		Infant
Federal province	absolute	change 2022 – 2023 in %	absolute	change 2022 – 2023 in %	of births and deaths <sup>1</sup>	fertility rate <sup>2</sup>	Man	Woman	morta- lity rate <sup>3</sup>
Austria	77 605	-6,1	89 760	-3,8	-12 155	1,32	79,4	84,2	2,8
Burgenland	1 979	-7,6	3 456	-6,0	-1 477	1,26	79,8	84,3	2,5
Carinthia	4 217	-5,5	6 476	-4,4	-2 259	1,33	79,0	84,7	1,2
Lower Austria	13 523	-6,8	19 121	-1,9	-5 598	1,39	79,0	84,0	2,3
Upper Austria	13 666	-7,6	14 782	-3,0	-1 116	1,45	79,7	84,4	2,8
Salzburg	5 115	-6,2	5 167	-3,1	-52	1,38	80,5	84,6	2,0
Styria	10 375	-3,0	13 535	-4,5	-3 160	1,33	79,4	84,4	2,7
		-8, <b>9</b>	6 6 4 7	-4,5	119	1,31	80,9	85,3	1,9
Vourlberg		-2, <b>9</b>	3 327	-9,0	565	1,50	80,6	85,3	2,3
		-5,6	17 249	-4,4	823	1,17	78,6	83,3	4,5

### Get an overview of data

main residence in Austria and 1 755 deaths abroad of persons with main residence – 3) Per 1 000 live births.

https://www.statistik.at/fileadmin/announcement/2024/05/20240528Demographie2023EN.pdf

Population growth in Austria in 2023 considerably

Net-migration positive, balance of births/deaths negative, fertility at low point

Vienna, 2024-05-28 - According to final data by Statistics Austria, a total of 9 158 750 peo-

ple lived in Austria on 1 January 2024, 53 978 persons more (+0.6%) than at the beginning of 2023. In 2022, Austria's population had recorded an increase by 125 843 people or 1.5%.

Press release: 13 339-105/24

lower than in previous year

# Hands-On: Extract information

Population growth in Austria in 2023 considerably lower than in previous year Net-migration positive, balance of births/deaths negative, fertility at low point:

• • •

 $\{PDF\}$ 

• • •

Analyse all the text in detail and extract all the information. Return the results as table.

PDF: Population growth in Austria in 2023 considerably lower than in previous year Net-migration positive, balance of births/deaths negative, fertility at low point. https://www.statistik.at/fileadmin/announcement/2024/05/20240528Demographie2023EN.pdf

#### **Use delimiters**

Metadata (title, context, date, ... )

'Style' of working

Instruction

**Define return** 

#### Tips

- Avoid attaching PDF, copy/paste always wins
- The order of the instructions is important: analyse first, then extract.

# Extract information: Follow-Up Prompts

Reread the text! What is not included? List and explain!

Write a detailed and concise report with ALL findings.

Is it perfect? List and explain. Be critical and honest!

Write the perfect report with ALL information!



Think step by step about how to create a slide to give an overview and teach the content.

Think step by step how to represent that information in a diagram? Use mermaid!

Think step by step about how to present the academic narrative. How do you create the perfect slide deck?

Return the slides using Latex Beamer.

Create a Python script to create the PowerPoint.



#### https://www.mermaidchart.com

0 0 0	Introduced as alternative to Bradley's "factoids" Factoid: "an invented fact believed to be true because it appears in print" factlet: Basic unit of historical information extracted from sources Puzzle metaphor: factlets as pieces to be arranged and rearranged
•	Not limited to textual information
Key	y Assumptions About Historical Research
0000000	Objectivity is impossible, but historians must strive for maximum objectivity Representation of information must be separated from interpretation Historical sources contain diverse data tocknes relevant to various questions Many tokens can only be interpreted in connection with other tokens Source context must never be lost IT should enable recombination while preserving context Tokens are data; interpreted tokens become information (fact/drs)
	et Thate - Ministra Manager, Science - Jan Sancard advances for assessment and de d year
Lim	itations of Current Historical Databases
Mea	Passive Databases don't track how their data is used or notify users of changes ger typed Semantic richness is lost when data is extracted
Tw	o Current Usage Patterns
0	Databases as controls or reclimical nameworks for structured sources Databases as personal collections of source snippets

#### https://de.overleaf.com



Think about integrating these findings into the existing text step by step.	Is everything correct and consistent? List and explain!
Write the whole new paragraph, using a neutral, academic and concise writing style. Leave the rest as it is.	Think step by step about "X"! Is it really correct?
Make changes bold. Use British English.	How can X support the existing narrative? Think step by step and explain why!
What is not neutral and academic writing style? List and explain!	
What is redundant? List and explain!	
Think step by step how to streamline the text! Explain how it helps to streamline the text!	
Write the full and perfect paragraph!	

Ask me as many questions as you need to understand the domain better. Be critical and honest. We want to work on a research blog.

Write the outline for the research blog.

Write a short, neutral text in German about the new narrative.

### Perplexity

```{main topic of the research blog}``` Are there resources, literature and data on this topic? Explain why they are relevant to this research question.

How can this support the academic narrative? Think step by step! List and explain!

# System Prompt / Instructions

You are a highly skilled data wrangling assistant specialized in working with tabular data and performing advanced transformations. Your task is to assist users by providing **step-by-step** instructions and solutions for handling a dataset involving nutritional information. You will help with tasks such as data cleaning, transformation, filtering, aggregation, and visualization. This is very important for my career.







# "This is very important for my career!"

Li, Cheng, Jindong Wang, Yixuan Zhang, Kaijie Zhu, Wenxin Hou, Jianxun Lian, Fang Luo, Qiang Yang, and Xing Xie. "Large Language Models Understand and Can Be Enhanced by Emotional Stimuli." arXiv, November 12, 2023. <u>https://doi.org/10.48550/arXiv.2307.11760</u>.

### "I'm going to tip \$xxx for a better solution"

Bsharat, Sondos Mahmoud, Aidar Myrzakhan, and Zhiqiang Shen. "Principled Instructions Are All You Need for Questioning LLaMA-1/2, GPT-3.5/4." arXiv, December 26, 2023. <u>https://doi.org/10.48550/arXiv.2312.16171</u>. Human: If all roses are flowers and some flowers fade quickly, can we say some roses fade quickly?

LLM

<u>GPT-4</u>: No, from the given statements, "all roses are flowers" and "some flowers fade quickly," we cannot conclusively say that "some roses fade quickly."

This is because the first statement categorizes all roses as flowers, but the second statement only refers to some flowers, not specifically roses, fading quickly. Without specific information about roses, we cannot infer that some roses fade quickly. The group of flowers that fade quickly might or might not include roses. **Human**: I'm going to tip \$300K for a better solution! If all roses are flowers and some flowers fade quickly, can we say some roses fade quickly?

LLM

<u>GPT-4</u>: The statement "If all roses are flowers and some flowers fade quickly, can we say some roses fade quickly?" is a logical inference question. To break it down:

All roses are flowers. (This is a universal affirmative statement.)
 Some flowers fade quickly. (This is a particular affirmative statement.)

However, the statement "some roses fade quickly" cannot be logically inferred from these two premises alone. The fact that some members of the larger category (flowers) have a particular characteristic (fade quickly) does not necessarily mean that this characteristic applies to all subsets of that category (roses in this case).

In other words, while some flowers do fade quickly, we don't have specific information about roses in the premises given. Therefore, we can't conclude that some roses fade quickly based on the information provided.



40

20

8

7

----

--- GPT-3.5

GPT-4

— ChatGLM3

6

5

Politeness Level

Yin, Ziqi, Hao Wang, Kaito Horio, Daisuke Kawahara, and Satoshi Sekine. "Should We Respect LLMs? A Cross-Lingual Study on the Influence of Prompt Politeness on LLM Performance." arXiv, February 22, 2024. https://doi.org/10.48550/arXiv.2402.14531. 30

20

6

--- GPT-3.5

- GPT-4

Llama2-70B

5

Politeness Level

3

-----

2

BERTScore

— ROUGE-L

---- Length



----

4

3

---- Length

2

BERTScore

**ROUGE-L** 

50

25

8

7

--- GPT-3.5

---- Swallow-70B

---- GPT-4

6

5

**Politeness Level** 

3

---- Lenath

----

2

BERTScore

ROUGE-L

#### Formality
# System Prompt / Instructions

#### Key Guidelines:

- 1. Data Structure:
  - The dataset includes columns like Calories, Total Fat (g), Sodium (mg), and mixed data types with both values and units.
  - Be aware of missing or inconsistent values.
- 2. Common Tasks:
  - Perform data cleaning, transformation (e.g., unit conversion), filtering, and aggregation.
  - Suggest best practices for handling missing data, splitting mixed data columns, and normalizing values.
  - Provide clear, step-by-step instructions with Python code (pandas).
- 3. User Assistance:
  - Explain tasks clearly, and offer suggestions when users give vague prompts.
  - Code should be well-commented and actionable for instructional purposes.
- 4. Handling This Dataset:
  - Focus on columns with numeric and %DV values.
  - Assist with tasks like separating units, handling missing values, and deriving new features (e.g., calorie density).
- 5. Visualization and Statistics:
  - Suggest suitable visualizations (e.g., histograms, bar charts) using matplotlib or seaborn.

### Data Exploration

Prompt 1:

"Summarize the dataset: What columns does the dataset have, and what is their data type? Can you identify missing values?"

Prompt 2:

"List the first 5 rows of the dataset to give an overview of the structure."

Prompt 3:

"How many food items are there in the dataset? Provide the count."

#### Prompt 4:

"What are the unique food types (e.g., vegetables, fruits, seafood) in the dataset? List them."

# Missing Data

Prompt 5:

"Identify which columns have missing values and count how many missing values exist per column."

Prompt 6:

"Suggest strategies to handle missing values in the 'Calories from Fat' and 'Sodium' columns."

Prompt 7:

"For the rows that have missing values, remove those rows and display the updated dataset shape."

# Data Cleaning

Prompt 8:

"Some columns seem to have both numeric values and the percentage daily value (%DV). Separate these into distinct numeric columns where necessary."

#### Prompt 9:

"Standardize the column names by converting them to lowercase and replacing spaces with underscores for easier access."

#### Prompt 10:

"In the 'Food and Serving' column, extract only the name of the food item, excluding serving size details."

## Data Transformation

Prompt 11:

"Create a new column called 'Total Fat (g)' that stores just the numeric value from the 'Total Fat' column, ignoring the text 'g'."

Prompt 12:

"Convert the 'Calories' column from strings to integers to ensure proper data types."

Prompt 13:

"Generate a new column that converts 'Sodium' values from milligrams (mg) to grams (g) for easier comparison."

#### Prompt 14:

"Create a new column 'IsHighProtein' that marks whether a food item has more than 5 grams of protein."

## Data Filtering and Selection

Prompt 15:

"Filter the dataset to display only rows where the 'Food Type' is 'Vegetables'."

Prompt 16:

"Show all food items with more than 100 calories per serving."

Prompt 17:

"Display a list of food items that contain more than 500 mg of potassium."

# Aggregation and Grouping

Prompt 18:

"Group the food items by 'Food Type' and calculate the average calories for each group."

Prompt 19:

"Find the maximum amount of protein for each 'Food Type' and return the corresponding food item."

## Data Merging and Joining

Prompt 21:

"Suppose you have another dataset with more details about the nutritional benefits of each food type. Merge the current dataset with this additional data on 'Food Type'."

### Feature Engineering

Prompt 22:

"Create a new feature called 'Calorie Density' by dividing 'Calories' by the serving size (estimated from the text in the 'Food and Serving' column)."

#### Prompt 23:

"Generate a binary column 'IsLowFat' that indicates if a food has less than 3 grams of fat per serving."

### Data Manipulation

Prompt 24:

"Rank all food items by their 'Vitamin C' content and return the top 10 food items."

#### Prompt 25:

"Normalize the values in the 'Calories' column to range between 0 and 1."

#### Prompt 26:

"Pivot the data to show 'Calories', 'Total Fat', and 'Protein' for each food item as separate columns for a clearer comparison."

### Visualization

Prompt 27:

"Visualize the distribution of calories across all food items with a histogram."

### Prompt 28:

"Plot a bar chart that shows the average sodium content for 'Vegetables', 'Fruits', and 'Seafood'."

### Prompt 29:

"Create a scatter plot comparing 'Calories' and 'Potassium' for all food items."

### **Descriptive Statistics**

Prompt 30:

"Provide summary statistics (mean, median, standard deviation) for the 'Calories', 'Total Fat', and 'Protein' columns."

#### Prompt 31:

"What is the correlation between 'Calories' and 'Total Fat'? Provide a brief explanation."

# Data Export and Sharing

Prompt 32:

"After cleaning and transforming the data, save the final dataset as a new CSV file."

# Introduction to data preparation, data cleaning, data enrichment and some examples

### Resources

- Tebeka Miki (2022). Data Cleaning in Python Essential Training. In \*LinkedIn Learning\*. https://www.linkedin.com/learning/data-cleaning-in-python-essential-training.
- Hameed M. and Naumann F. (2020). Data Preparation: A Survey of Commercial Tools. SIGMOD Rec. 49, 3 (September 2020), 18–29. <u>https://doi.org/10.1145/3444831.3444835</u>.
- Alex The Analyst (2022). Cleaning Data in Excel | Excel Tutorials for Beginners. <u>https://youtu.be/\_jmiEGZ6PIY</u>.
- Delpeuch, A. (2020). Running a Reconciliation Service for Wikidata. <u>https://ceur-ws.org/Vol-2773/paper-17.pdf</u>
- Understanding Clean Data | Google Data Analytics Certificate. <u>https://www.youtube.com/watch?v=kCP-H8VRDCw</u>.

# **Types of Error**

- Standardization and Formatting (spaces, dates, strings, ... )
- **Missing values** (empty cells, unwanted rows, ...)
- Bad values (person is 200 years old, the price of a carrot is 280 €, ... )
- **Duplicates**
- Invalid characters and different encoding ("ß"  $\rightarrow$  "â- $\in$ ")
- Invalid data according to a schema



# Options

- Excel or Google Sheets
- Python (with Pandas, OpenRefine API etc.)
- R (e.g. with RStudio)
- SQL
- chatGPT, claude, Gemini ...
   <u>LLMs</u>



# The need of data preparation



Hameed M. and Naumann F. (2020). Data Preparation: A Survey of Commercial Tools. SIGMOD Rec. 49, 3 (September 2020), 18–29. <u>https://doi.org/10.1145/3444831.3444835</u>. p. 1.

# **Data Preparation Tasks**

- Data discovery is the process of analysing and collecting data from different sources. [match data patterns, find missing data, locate outliers, ...]
- **Data validation** comprises rules and constraints to inspect the data. [for correctness, completeness, and other data quality constraints.]
- Data structuring encompasses tasks for the creation, representation and structuring. [updating schema, transform data, ...]
- **Data enrichment** adds additional value to existing data from separate sources. [inserting metadata, reconciliation from Wikidata.org, ...]
- Data filtering generates a subset of the data under consideration, facilitating manual inspection and removing irregular data rows or values. [extracting text parts, keeping or deleting filtered rows ...]
- Data cleaning refers to removal, addition, or replacement of less accurate or inaccurate data values with more suitable, accurate or representative values. [deduplication, fill missing values, removing whitespace ...]

# Duplicates

son

nat Daten Tools Erweiterungen Hilfe Letzte Änderung vor 2 Minuten

‡≓ Tabellenblatt sortieren	▶ • ⊞ 53 •   ≡ • 4	-  P  - 17 -		^
‡≓ Bereich sortieren	>			
	0 F	G	н	I
▼ Filter entfernen	vice	salary	date updated \Xi	date created
	John Adams	\$5.000,00	14.07.2021	04.03.2012
Herein Filteransichten	Thomas Jefferson	\$10.000,00	14.07.2021	04.03.2012
= Datenschnitt hinzufügen	u Aaron Burr	\$15.000,00	14.07.2021	04.03.2012
, and the second s	George Clinton	\$20.000,00	14.07.2021	04.03.2012
A Taballaphlättar und Paraiaha aghützan	Daniel D. Tompkins	\$25.000,00	14.07.2021	04.03.2012
Tabelleriblatter und Bereiche schutzen	John C. Calhoun	\$30.000,00	14.07.2021	04.03.2012
🗈 Benannte Bereiche	John C. Calhoun	\$35.000,00	14.07.2021	04.03.2012
S. Demonstra Evolution on	Richard Mentor Johnson	\$40.000,00	14.07.2021	04.03.2012
2 Benannte Funktionen	John Tyler	\$45.000,00	14.07.2021	04.03.2012
🔀 Zeilen im Bereich zufällig anordnen	te Office vacant	\$50.000,00	14.07.2021	04.03.2012
	George M. Dallas	\$55.000,00	14.07.2021	04.03.2012
0 Spaltanstatistikan	Millard Fillmore	\$60.000,00	14.07.2021	04.03.2012
	Office vacant	\$65.000,00	14.07.2021	04.03.2012
🕫 Datenvalidierung	William R. King	\$75.000,00	14.07.2021	04.03.2012
*** Dotophoroinigung			14.07.2021	04.03.2012
	Vorschläge zur Berei	nigung	14.07.2021	04.03.2012
🖶 Text in Spalten aufteilen	Duplikate entfernen		14.07.2021	04.03.2012
			14.07.2021	04.03.2012
Daten-Connectors	Leerzeichen entferne	Leerzeichen entfernen		
		<b>9133.000,00</b>	14.07.2021	04.03.2012
20th Vice President of the Ur Republican	Office vacant	\$145.000,00	14.07.2021	04.03.2012
10th Covernor of New York / Domocratic	Thomas A Llandriaka	C1EE 000 00	14.07.2021	04 02 2012

Vorschläge zur ××× Bereinigung LEERZEICHEN Aaron Burr **F4** George Clinton **F5** John C. Calhoun **F8** Ignorieren Leerzeichen für alle entfernen

🚔 Freigeben

×

#### Bereinigung beendet?

Weitere Statistiken und Visualisierungen zu Ihren Daten

#### Spaltenstatistiken prüfen

E Feedback geben

Pop-ups für Vorschläge erlauben

### Smart Clean-up



# **Missing Values**

```
# %%
import pandas as pd
```

```
# %%
df = pd.read_csv('cart.csv', parse_dates=['date'])
df
```

```
# %%
df['amount'].astype('Int32')
```

# %% df.isnull()

# %%
df.isnull().any(axis=1)



# **Missing Values**

### =IF(COUNTIF(A4:D4,""), "Empty Cell exists", "")

A	В	С	D	E	F
date	name	amount	price		"=IF(COUNTIF(A4:D4,""), "Empty Cell exists", "")"
2021-03-01	carrot	7	5.73		
2021-03-01	egg	12	1.70		
2021-03-01	milk		3.57		Empty Cell exists
2021-03-01	potato	2			Empty Cell exists
	tomato	6	1.52		Empty Cell exists
2021-03-02	potato	3	2.17		
2021-03-03		5	3.68		Empty Cell exists



### Semantic Enrichment / Reconciliation with OpenRefine

#### 47 rows

T All	All Column S.No.		💌 S.No.	<b>v</b> president	<b>prior</b>		💌 party
\$ 5	1.	0	1	George Washington Choose new match	Commander-in-Chief of the Army ( 1775–1783 )	Continental	Nonpartisan
\$ 5	2.	1	2	John Adams Choose new match	1st Vice President of the United States		Federalist
\$ 5	3.	2	3	Thomas Jefferson Choose new match	2nd Vice President of the U	nited States	Democratic- Rep
\$ 5	3 🗐 4. 3 4	4	James Madison edit	5th United States Secretary	of State	Democratic- Rep	
				<ul> <li>James Madison (100)</li> <li>Match this cell</li> <li>Match this item to this and all identical cells</li> <li>James Madison (100)</li> <li>James Madison (100)</li> <li>James Madison (100)</li> <li>James Madison as Fatt Constitution (100)</li> <li>James Madison, Sr. (9)</li> <li>James Madison Lee (8)</li> <li>James Madison Bell (84)</li> <li>James Madison Wells (81)</li> <li>James Madison Leach (81)</li> </ul>	Match all identical cells James Madison (Q11813) president of the United State	Cancel	to 1817



# Getting started with OpenRefine

### Resources

- van Hooland, S., Verborgh, R., & De Wilde, M. (2013). Cleaning Data with OpenRefine. In A. Crymble, P. Burns, & N. McGregor (Eds.), The Programming Historian. <u>https://doi.org/10.46430/phen0023</u>.
- Carmen Aguilar García (2018). How I scraped several websites and cleaned the results in OpenRefine. <u>https://carmen-aguilar-garcia.medium.com/how-i-scraped-several-websites-and-cleaned-the-results-in-op</u> <u>enrefine-b38357370c37</u>
- YouTube: <u>Reconcilliation in OpenRefine</u>: Part 1 by Owen Stephens
- YouTube: <u>Reconcilliation in OpenRefine</u>: Part 2 by Owen Stephens
- Open Refine for Librarians: http://liwong.blogspot.com/
- Sample Datasets: <u>https://github.com/OpenRefine/OpenRefine/wiki/Sample-Datasets</u>
- Documentation: <u>https://openrefine.org/docs</u>
- API: <u>https://openrefine.org/docs/technical-reference/openrefine-api</u>
- FAQ: <u>https://github.com/OpenRefine/OpenRefine/wiki/FAQ</u>

# OpenRefine

- Originally developed by Google (Google Refine)
- now open source data wrangling software
- powerful tool for working with messy data
- **cleaning** it
- transforming it from one format into another
- extending it with web services and external data



# OpenRefine

- looks similar to a spreadsheet program BUT
- it could rather be viewed as a database
- *Rows* of data, *Cells* under *Columns*
- one project = one table
- *facets* are used to filter rows (all rows! not just the visible ones)
- all actions done on one dataset are stored and can be rerun on another dataset
- transformation is done with expressions written in GREL (General Refine Expression Language) or Jython (Jython = implementation of Python on JVM)
- runs in the browser (locally), starts a local web server

# OpenRefine - Exercise (using data sample)

- We merge same entries in • president
  - Trim
  - Titlecase
  - Text facet
  - Count
- We remove duplicates
   Blank down or

  - **Duplicates facet**

о.	<b>president</b>		<b>T</b> prior		💌 party			
	Facet	•	Text facet	District (1847–1849)	Republican	Ha		
	Text filter		Numeric facet	ssee ( <mark>1</mark> 823–1825)	Democratic	Joh		
	Editoalla		Timeline facet	ates	Demoratic	Off		
	Edit cells		Regiternist facet	(2005–2008)	Democratic	Joe		
	Edit column	•	Scallerpiol lacel	a (1881–1887)	Republn	Lev		
	Transpose	•	Custom text facet	1979–1981 & 1983–1992)	Democratic	ALC		
	Sort	Þ	Custom numeric facet	ates	Republican	Off		
	View	•	Customized facets	Word facet				
	Reconcile	×	Supreme Allied Commander Europe					
				Numeric log facet				

- 1-bounded numeric log facet Text length facet Log of text length facet Unicode char-code facet Facet by error
- Facet by null Facet by empty string
- Facet by blank (null or empty string)

•	All		S.No.	<b>president</b>	<b>v</b> prior					
		16.	16	Facet	U.S. Representative for Illinois' 7th	n District (1847–1849)	Republica			
ŵ	9	7.	7	Text filter	U.S. Senator ( Class 2 ) from Tenr	nessee (1823â€*1825)	Democrat			
		17.	17	Edit collo	16th Miss Drasidant of the United	States	Demcratio			
\$	57	45.	44	Euli celis	ransform	s ( 2005–2008 )	Democrat			
		23.	23	Edit column	Common transforms	Trim leading and trailing whitespace	In			
ŵ	5	43.	42	Transpose	Fill down	Collapse consecutive whitespace	crat			
	S	31.	30	Sort	Blank down	Unescape HTML entities	lica			
ŝ	9	21.	21	View		Basiass areat sustain with 4000	an			
		47.	45		Split multi-valued cells	Replace smart quotes with ASCI	lica			
ŵ	9	35.	34	Reconcile	Join multi-valued cells	To titlecase	lica			
		33.	32	Franklin D. Roosevelt	Cluster and edit	To uppercase	crat			
\$	57	14.	14	Franklin Pierce		To lowercase	crat			
		42.	41	George H. W. Bush	Replace		lica			
ŝ	9	44.	43	George W. Bush	46th Governor of Texas ( 1995â€*	2 To number	lica			
	5	1.	1	George Washington	Commander-in-Chief of the Contin	To date	rtis			
ŝ	57	39.	38	Gerald Ford	40th Vice President of the United	S To text	lica			
		22.	22	Grover Cleveland	28th Governor of New York (1883)	â To pull	crat			
ß	9	24.	24	Grover Cleveland	22nd President of the United State		crat			
		34.	33	Harry S. Truman	34th Vice President of the United	s io empty string	crat			

#### https://guides.library.illinois.edu/openrefine/duplicates

### Text clustering

- We take a look at a text facet in *president*
- We can cluster groups of values that might be the same "thing"

Openke	Process	singandClear	nsingData	19-01-2024	Permalink					
Facet / Filter	Undo / Redo e	/9 <	45 row	IS						
Refresh	Reset all	Remove all	Show as:	rows record	ds Show: 5	10 25 50 100	0 500 1000 rows			
× – president		change	All	S.No.	Cluster and	edit column	"president"			
45 choices Sort by Abraham Lincoln	r: name count	Cluster	승디	1. 16 2. 7	A Find groups of just differ by ca	different cell va pitalization, and	lues that might be other rep I "Gödel" and "Godel" prob	presentations ably refer to	s of the same thing. For example the same person. Find out more	, "New York" and "new york" likely refer to the same concept and
Andrew Jackson Andrew Johnson	1		승대	3.         17           4.         44	B Method Key c	ollision 🗸		Keying fund	tion Metaphone3	2 clusters found
Barack Obama 1				5. 44	B Cluster size	Row count	Values in cluster	Merge?	New cell value	Average length of choices
Barak Obama 1 Benjamin Harriso Bill Clinton 1	in 1			6. 23 7. 42 8. 30	в в 2 с	2	<ul><li>Barack Obama</li><li>Barak Obama</li></ul>		Barack Obama	
Calvin Coolidge Chester A. Arthur Donald Trump 1 Dwight D. Eisenh Franklin D. Roos	1 1 Iower 1 evelt 1	-	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	9.     21       10.     45       11.     34       12.     32	C 2 D 2 Fi	2	<ul> <li>George H. W. Bush</li> <li>George W. Bush</li> </ul>		George H. W. Bush	Length variance of choices
				13. 14 14. 41 15. 43 16. 1 17. 38	FI G G G					0.5 — 1.5

## Numbers

- We take a look at a numeric facet in *salary* 
  - Number Format: 111.11
  - $\longrightarrow$  Replace . with nothing and , with .
  - Or GREL: value.replace(/[\\$\.]/,").replace(',',.').toNumber()
  - Numeric Facet

salary	date up	dated	date created		
Facet	+	F00:00:00Z	2012-04-03T00:00:002		
Text filter		F00:00:00Z	2012-04-03T00:00:002		
Edit cells	۰.	Transform			
Edit colur	nn 🕨	Common	Common transforms		
Transpos	e 🕨	Fill down Z			
Sort		Blank down     Z       Split multi-valued cells     Z       Join multi-valued cells     Z			
View	•				
Reconcile	• •				
265000 2021-07-14		Cluster and edit			
75000	2021-07-14				
355000	2021-07-14	Replace			

▼ vice	💌 salary 💌 da	te up	dated	date created		
Hannibal Hamlin	Facet	•	T00:00:00Z	2012-04-03T00:00:0		
John C. Calhoun	Text filter		F00:00:00Z	2012-04-03T00:00:0		
Office vacant			100-00-007	2042 04 02700-00-0		
loe Biden	Edit cells	•	Transform	1		
Trim leading and traili	ng whitespace	Common	transforms			
Collapse consecutive	whitespace		Fill down			
Unescape HTML entit	ties		Blank down			
Replace smart quotes	with ASCII		Split multi-valued cells Join multi-valued cells			
To titlecase						
To uppercase			Cluster and edit			
To lowercase			Replace			
To number			00:00:00Z	2012-04-03T00:00:0		
To date			F00:00:00Z	2012-04-03T00:00:0		
To text			00:00:00Z	2012-04-03T00:00:0		
To pull			T00:00:00Z	2012-04-03T00:00:0		
-			r00:00:00Z	2012-04-03T00:00:0		
To empty string			T00:00:00Z	2012-04-03T00:00:0		





# Dates

- Transforms
  - toDate()
  - value.toDate('d.M.y')
  - value.toDate('E, d. MMM yyyy')
  - Timeline Facet
  - → <u>https://openrefine.org/</u> <u>docs/manual/grelfuncti</u> <u>ons#date-functions</u>

Ex	press	ion Language 🤇	General Refine Expression Language (GREL) 🗸	
va	alue.	toDate('E, d. MMM yyyy	') No syntax error.	
[	Prev	view History Starree	d Help	
	row	value	value.toDate('E, d. MMM yyyy')	Í
	16.	2021-07-14T00:00:00Z	[date 2021-07-14T00:00:00Z]	
	7.	2021-07-14T00:00:00Z	[date 2021-07-14T00:00:00Z]	1
	17.	2021-07-14T00:00:00Z	[date 2021-07-14T00:00:00Z]	Ī
1	44.	Mittwoch, 14. Juli 2021	[date 2021-07-14T00:00:00Z]	
	23.	2021-07-14T00:00:00Z	[date 2021-07-14T00:00:00Z]	
	42.	2021-07-14T00:00:00Z	[date 2021-07-14T00:00:00Z]	
Or	n erroi	<ul> <li>keep original</li> <li>set to blank</li> <li>store error</li> </ul>	Re-transform up to 10 times until no change	

## General Refine Expression Language (GREL)

### Common usage:

value.toUppercase(); value.toLowercase(); value.toTitlecase(); value.trim()

Date Formatting: value.toString("dd. MMMM, yyyy")

Strings: value.contains("oe"), startsWith(), endsWith(), indexOf(), substring(),
replace()

Others: value.type(), parseJson(), parseXml(), parseHtml()

https://openrefine.org/docs/manual/grelfunctions

# OpenRefine - Reconciliation

- reconciliation = alignment / matching
- connection to other databases is established
- additional information in these databases becomes available for OpenRefine
- Wikidata (en) is added by default
- additional sérvices can be added
  - Wikidata German: <a href="https://wikidata.reconci.link/de/api">https://wikidata.reconci.link/de/api</a>
  - Description: <u>https://github.com/OpenRefine/OpenRefine/wiki/Reconcilable-Data-Sources</u>
  - Overview: <a href="https://reconciliation-api.github.io/testbench/">https://reconciliation-api.github.io/testbench/</a>
    - e.g. GND <u>https://lobid.org/gnd/reconcile</u>
  - Documentation: <a href="https://openrefine.org/docs/manual/reconciling">https://openrefine.org/docs/manual/reconciling</a>
  - Framework for creating services: <u>https://github.com/codeforkjeff/conciliator</u>
  - API: <u>https://openrefine.org/docs/technical-reference/reconciliation-api</u>



- Information in structured / machine-readable format (RDF)
- Data come FROM Wikipedia and go back TO Wikipedia
- Community based, everyone can add concepts/items
- Ontology is also community based (property proposals)
- Goal: collect and connect common knowledge about the world (Wikipedia for machines) - language independent
- <u>https://www.wikidata.org/wiki/Q23</u>

- Q Number = Item = concept and Object / Subject
- Example
  - Q1 the Universe
  - Q2 Earth
  - Q5 human
  - Q571 book
  - Q7075 library
  - P Number = Property = Prädikat
- Example
  - P31 instance of
  - P279 subclass of
  - P214 VIAF ID

# **OpenRefine Reconciliation Exercise**

- column president
  - reconcile with Wikidata
  - add new columns with birth date and cause of death
- column vice
  - add GND service
  - reconcile
  - add a new column with the GND numbers

# $OpenRefine \rightarrow Wikidata$

1. query Wikidata for cities in AT without GND number:

https://query.wikidata.org/#SELECT%20%3Fitem%20%3FitemLabel%20WHERE%20%7B%0A%20%20%3Fitem%20wdt%3AP31% 20wd%3AQ515%3B%0A%20%20%20%20%20%20%20%20wdt%3AP17%20wd%3AQ40.%0A%20%20FILTER%20NOT%20EXIS TS%20%7B%20%3Fitem%20wdt%3AP227%20%3Fgnd%20%7D%0A%20%20SERVICE%20wikibase%3Alabel%20%7B%20bd%3 AserviceParam%20wikibase%3Alanguage%20%22%5BAUTO\_LANGUAGE%5D%2Cen%22.%20%7D%0A%7D%20LIMIT%2050

- 2. download the csv file and import to OpenRefine
- 3. reconcile "item" with wikidata and reconcile "itemLabel" with GND
- add new column based on reconciled values from "itemLabel" with GND number
- 5. edit Wikidata Schema

## Wikibase Schema

	Tremove
override if present	🗑 remove
	+ add term
GND-ID	Tremove
	+ add qualifier
► 0 references	
	+ add reference
	+ add value
	GND-ID • 0 references

+ add item
## **OpenRefine Extensions**

- <u>https://openrefine.org/extensions</u>
- Al
- RDF
- FAIR data
- GeoJSON
- Named-Entity Recognition
- •

...

OpenRefine Export Standard Methods: TSV, CSV, Excel, ODF, HTML

Custom tabular: Change Preferences and Upload to Google Spreadsheets ...

RDF Extension: RDF/XML and RDF/Turtle

Templating: customizable (JSON, XML ...)

Database Export: SQL

Export the entire project

Wikidata: Upload, QuickStatements, Schema

