

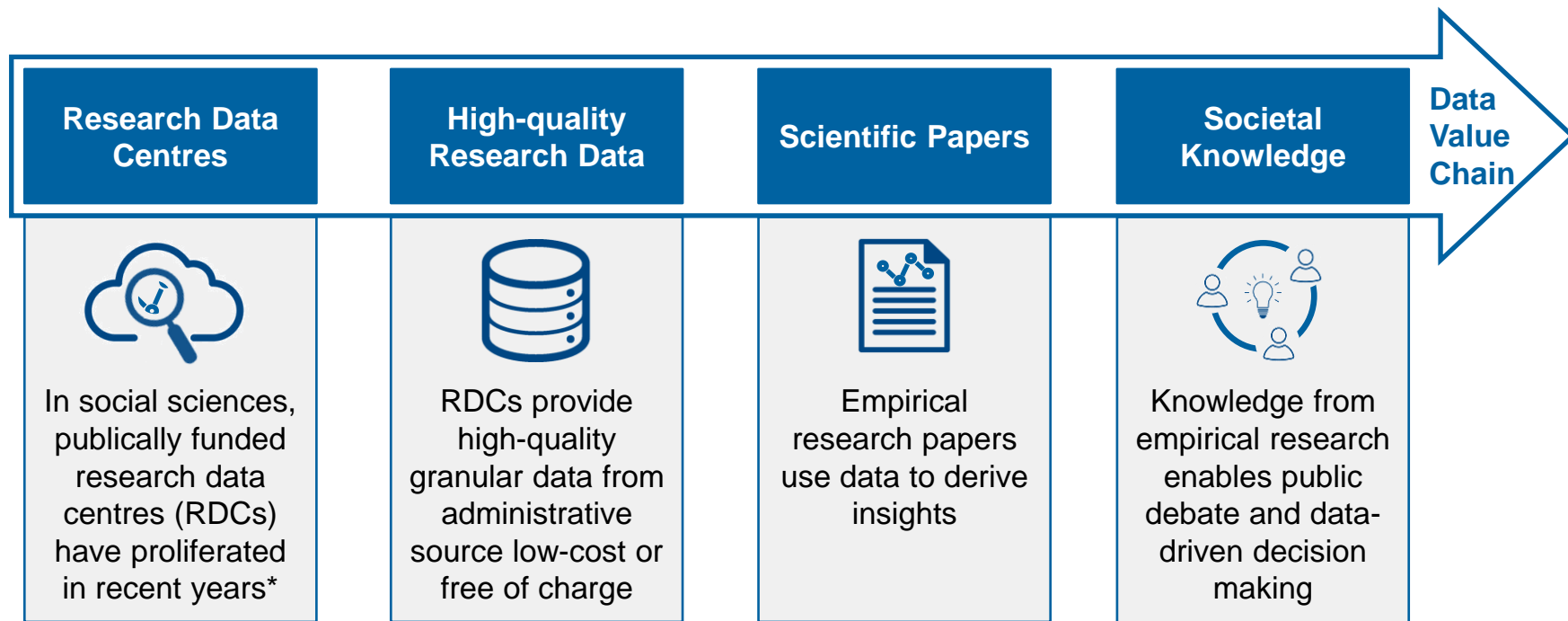
Leveraging Large Language Models to Extract Data Citations

3rd IFC Workshop on Data Science in Central Banking

Sebastian Seltmann, Emily Kormanyos, Hendrik Christian Doll, Shir Frank, and Kilian Graef – Deutsche Bundesbank, Data Service Centre

The authors would like to thank colleagues in Bundesbank's Data Service Centre for valuable comments and feedback.
All views expressed in this report are personal views of the authors and do not necessarily reflect the views of Deutsche Bundesbank or the Eurosystem.

The Data Value Chain



Since public funds are used, the public has an interest in knowing the value created by the investment

Issue: It is hard to measure data impact, as data citations in scientific papers are not standardised

Extracting Data Citations from Academic Papers



Theory

Blaschke & Hirsch (2023)¹ attempt to **measure the value of an RDC**, leveraging on manually collected info on projects and publications



Related Work

Polak & Morgan (2023)² propose **ChatExtract to extract materials data in research papers**, leading to accurate results when identifying known properties



Gap in Literature

Most of the literature to date seems to **focus on the natural sciences or ML/NLP papers** specifically



Our Contribution

Recent advances in natural language processing (NLP) enable us to flexibly **detect and connect data source descriptions from academic papers using GPT-3.5**

¹ Blaschke & Hirsch (2023)

² Polak & Morgan (2023)

The Assistant-Style Language Model Approach



User Messages

I want you to identify and list all datasets or data sources from the following text of a scientific paper.

Instruction

Okay!

*Non-technical summary
Research Question
The financial crisis showed that a sound capital base is a necessary, but not sufficient condition for banks to be resilient to major shocks: sound liquidity buffers to withstand short-term liquidity shocks and a sound stable funding base to withstand*

Text data

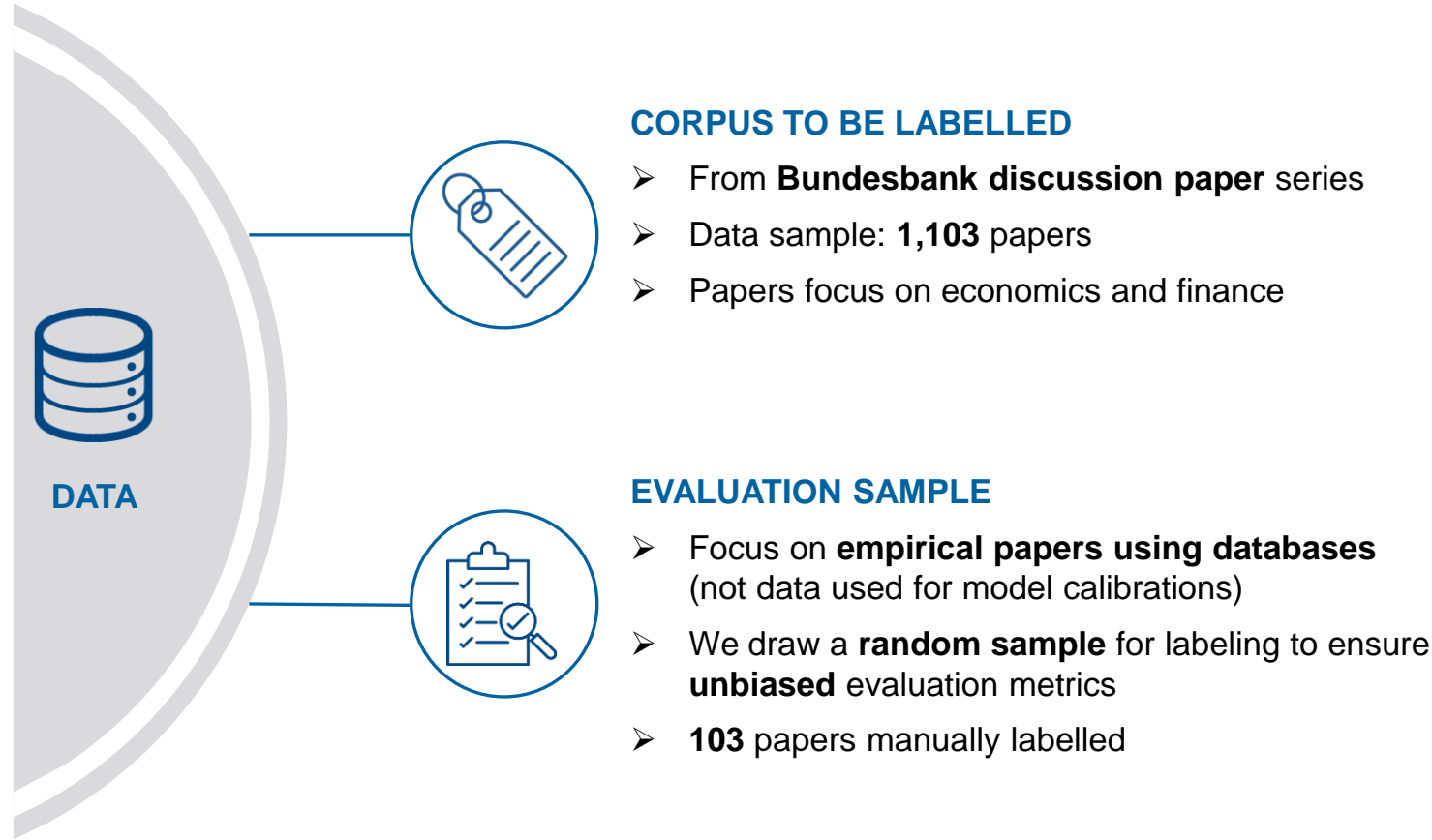
The provided paper contains the following datasets/datasources:
- Eikon
- Bloomberg
- Deutsche Bundesbank's balance of payments statistics
- BAKIS

Answer

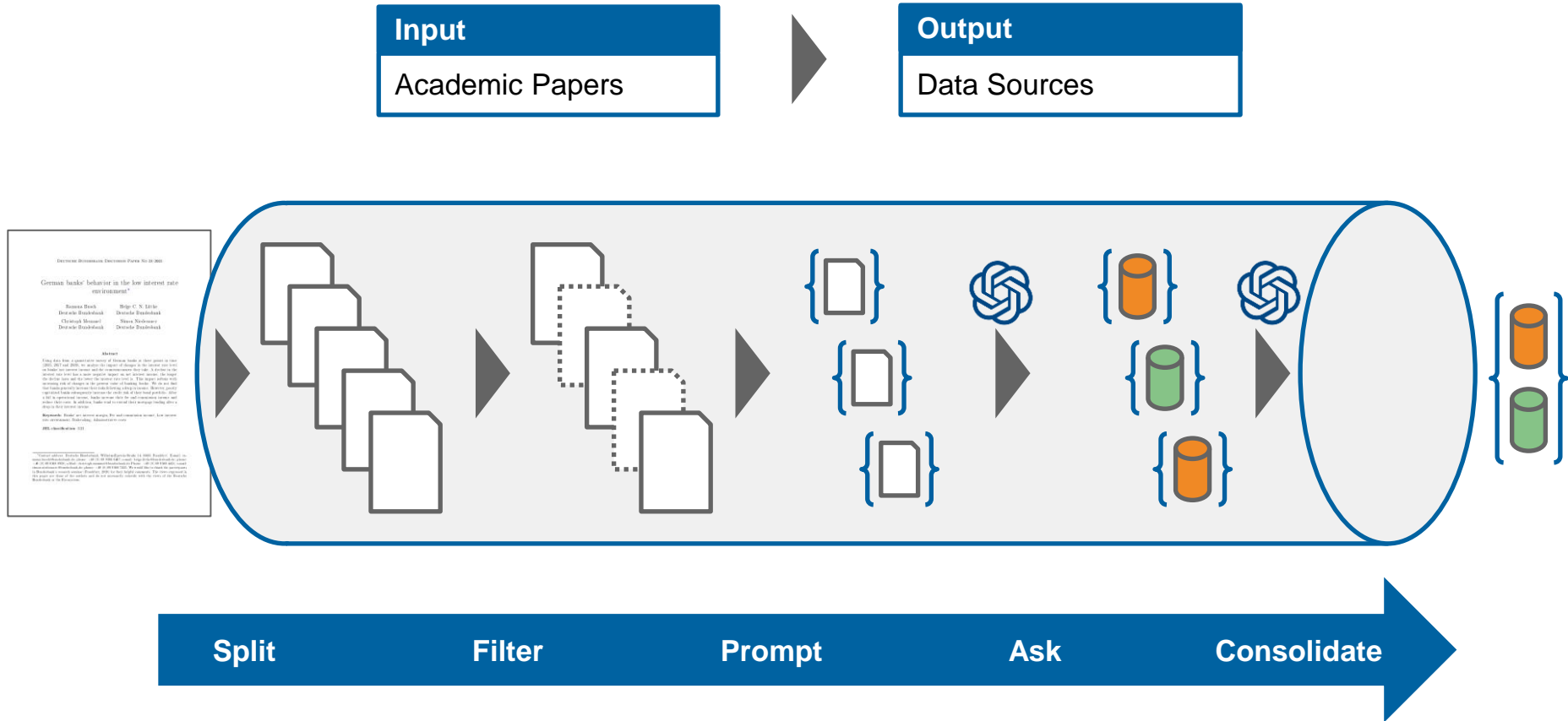


Assistant Model Messages

Paper Corpus and Evaluation Sample



Dataset Citation Extraction Pipeline

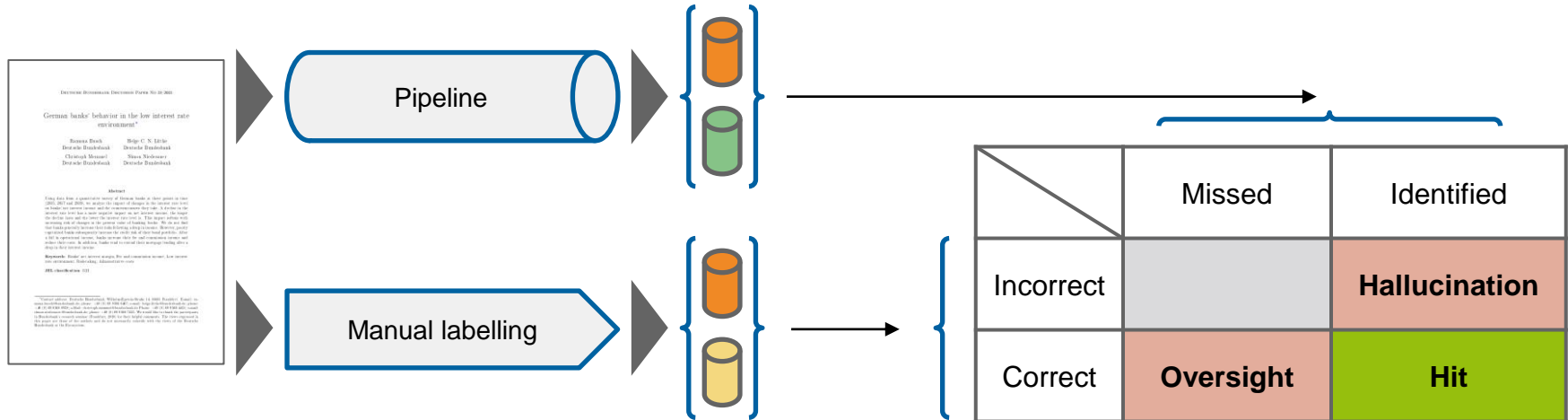
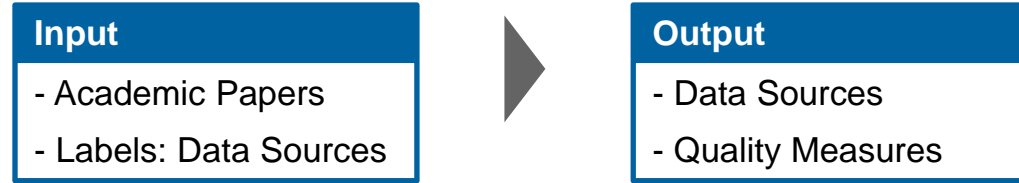


Leveraging Large Language Models to Extract Data Citations

October 2023

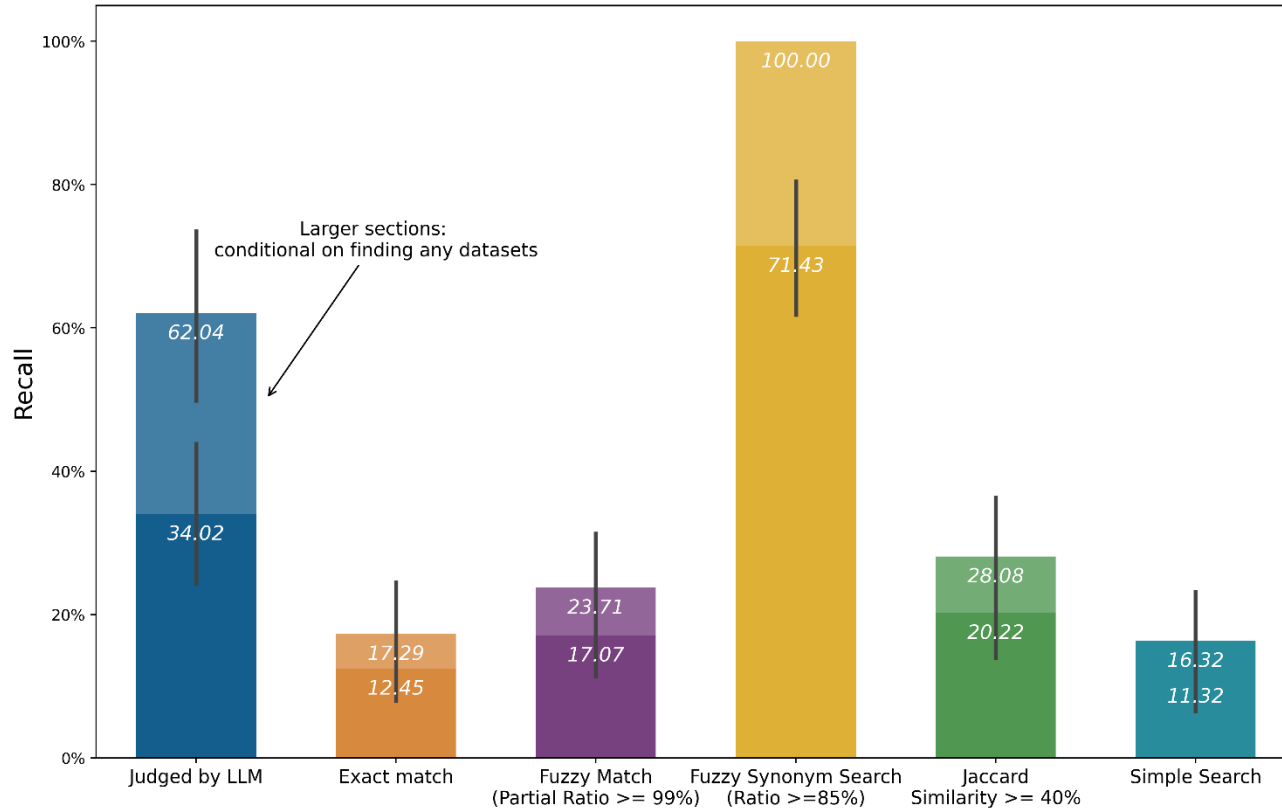
Seite 6

Assessing the Quality of the Assistant-Model's Output



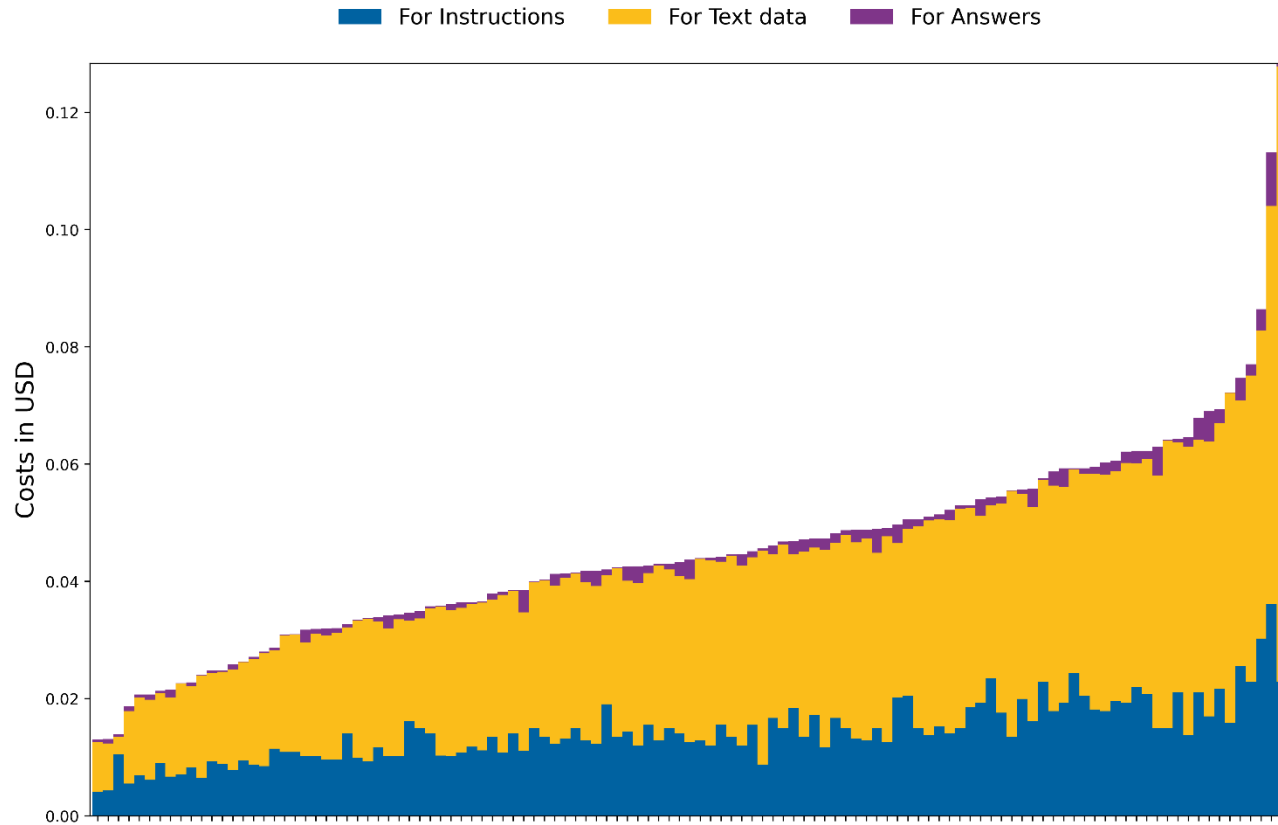
Results

Different measures of apparent performance



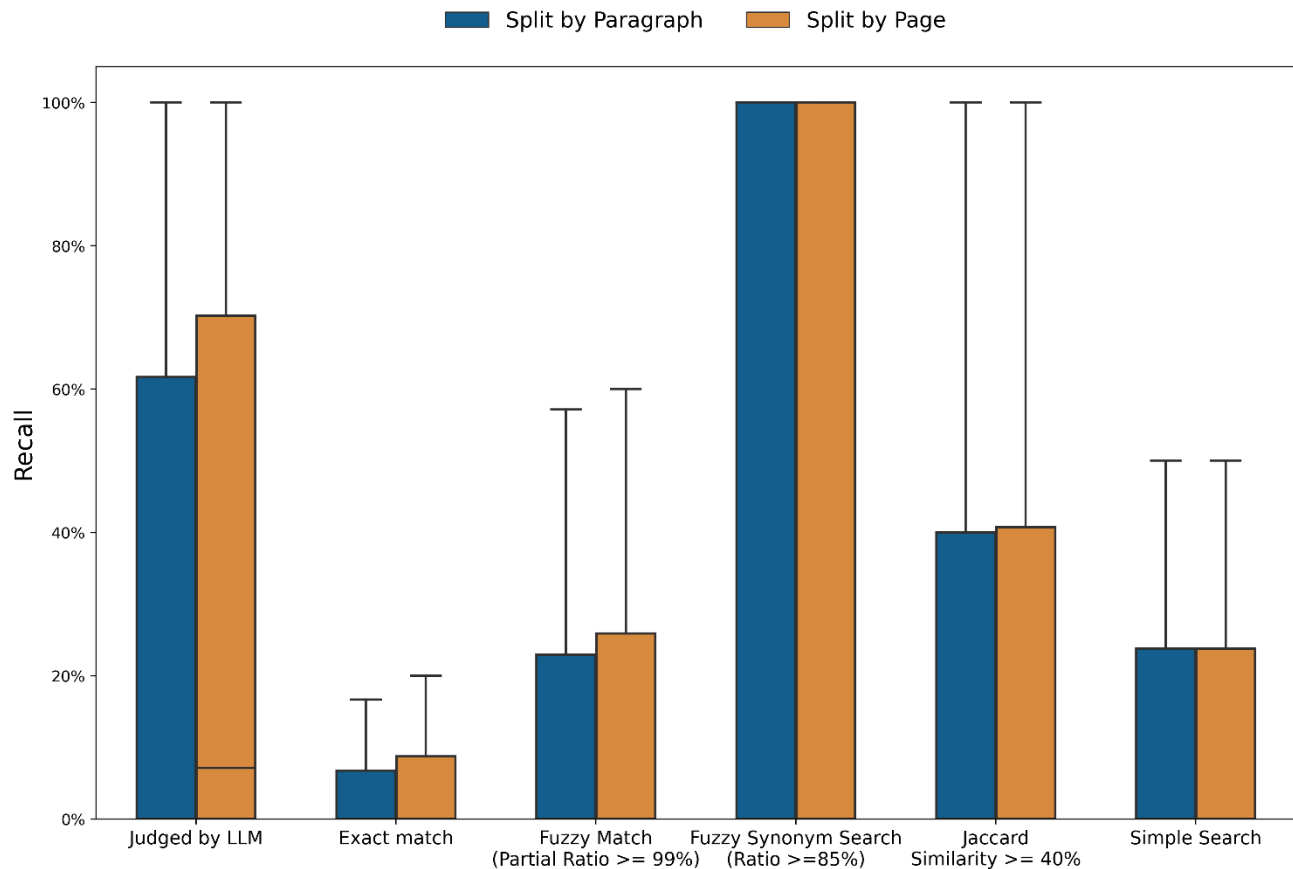
Results

API Costs per paper processed



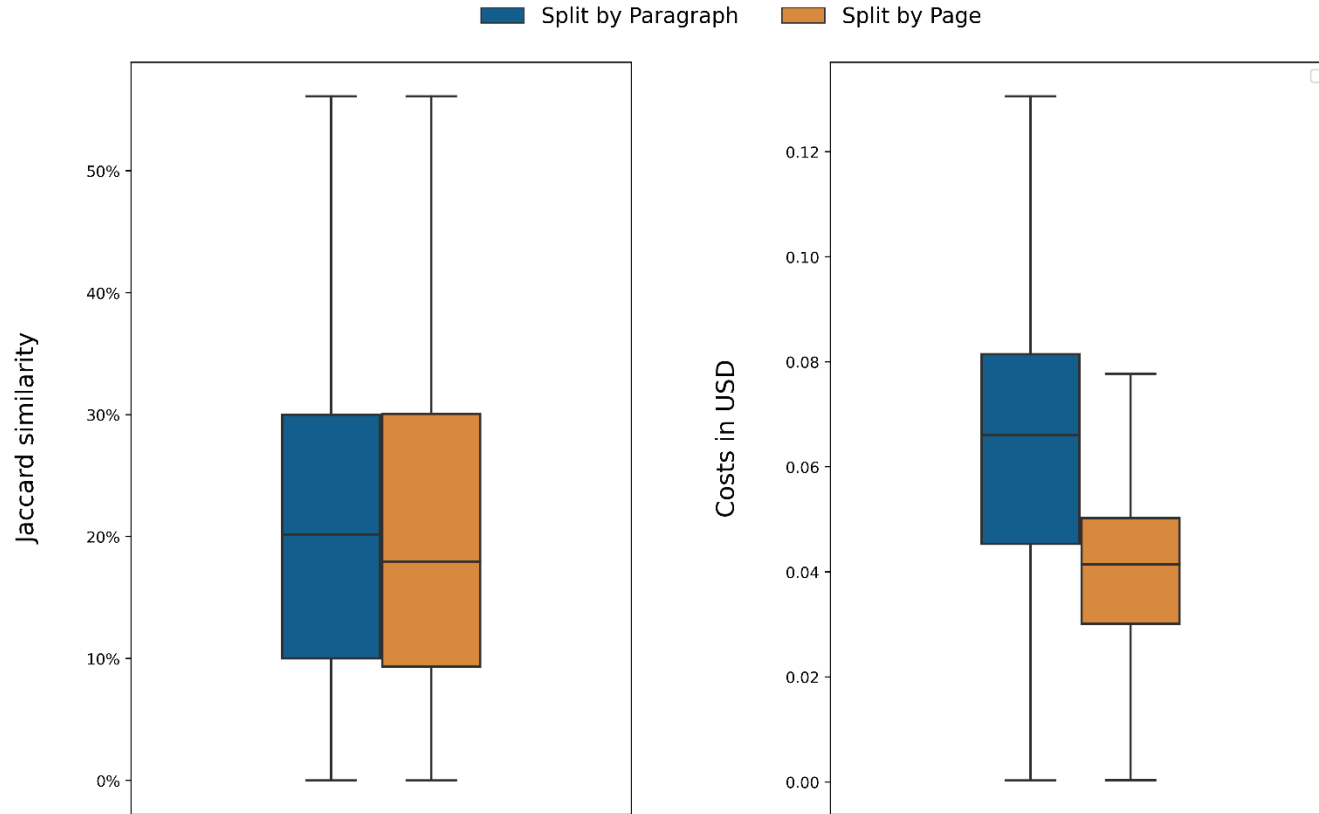
Results

Larger Contexts provide minor performance improvement



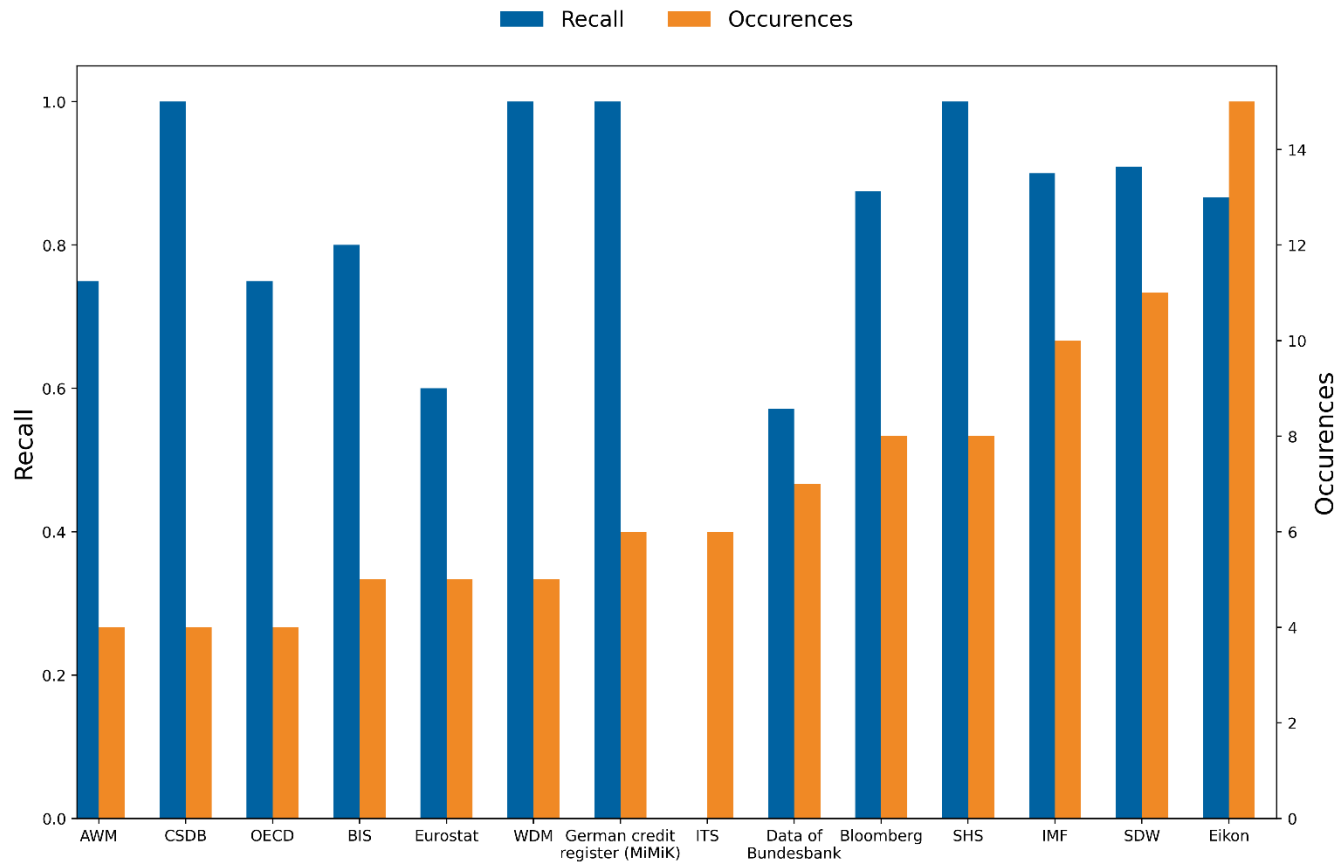
Results

Larger Contexts significantly reduce overall costs



Results

Recall for individual datasets



Challenges & Mitigations



DEPENDENCY

Reliance on the OpenAI API is a point of failure
Unforeseen outages are possible



- OpenAI transparently reports system status
- Pipeline can plug-in alternative models

UNSTRUCTUREDNESS

Assistant-Style language models produce plain text, strict output-syntax is merely a suggestion



- Zero “temperature” improves reliability

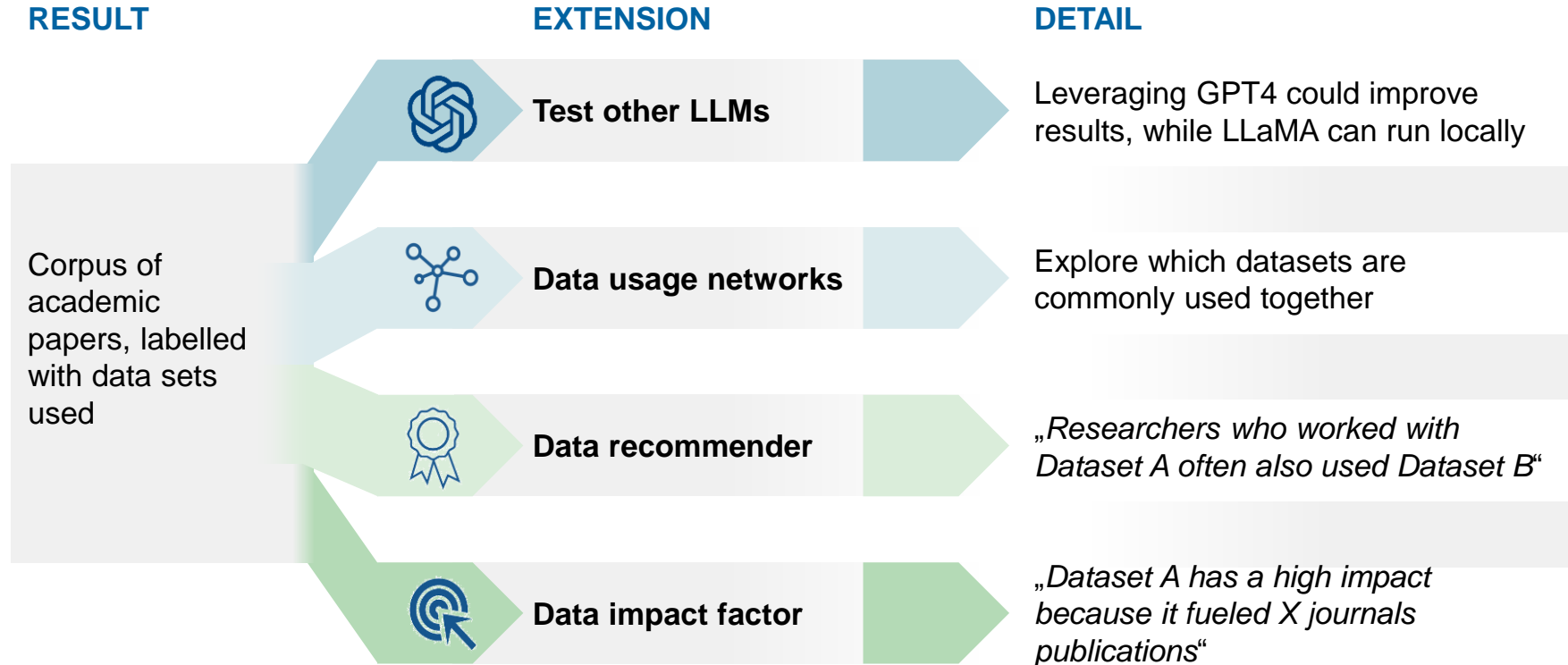
TESTING

The definition of success itself is up to debate and interpretation



- Pipeline tracks as many potentially relevant dimensions as possible

Future Extensions and Applications



Key Take-Aways



CAPABILITY

Assistant-style language models can effectively extract specific types of information from large bodies of unstructured text



EVALUATION

Think carefully about how to measure success and track experiments systematically



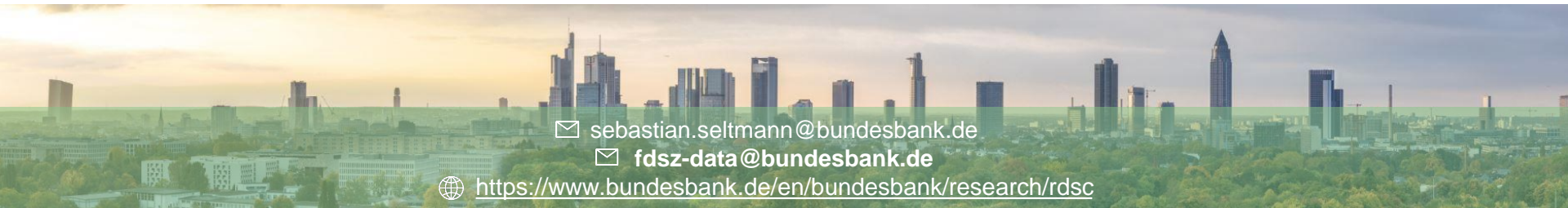
GUIDANCE

Give the model clear and simple tasks, yet as much context as possible



References

- Blaschke, J. & C. Hirsch (2023). On the value of data sharing: Empirical evidence from the Research Data and Service Centre, Technical Report 2023-08 – Version 1.0. Deutsche Bundesbank, Research Data and Service Centre. Retrieved 2023/08/21 from <https://www.bundesbank.de/resource/blob/863758/7ebe74476186cd3364a11b3869ada80a/mL/2023-08-value-data.pdf>
- Polak, M. P. & D. Morgan (2023). Extracting accurate materials data from research papers with conversational language models and prompt engineering—example of chatgpt. arXiv preprint arXiv:2303.05352

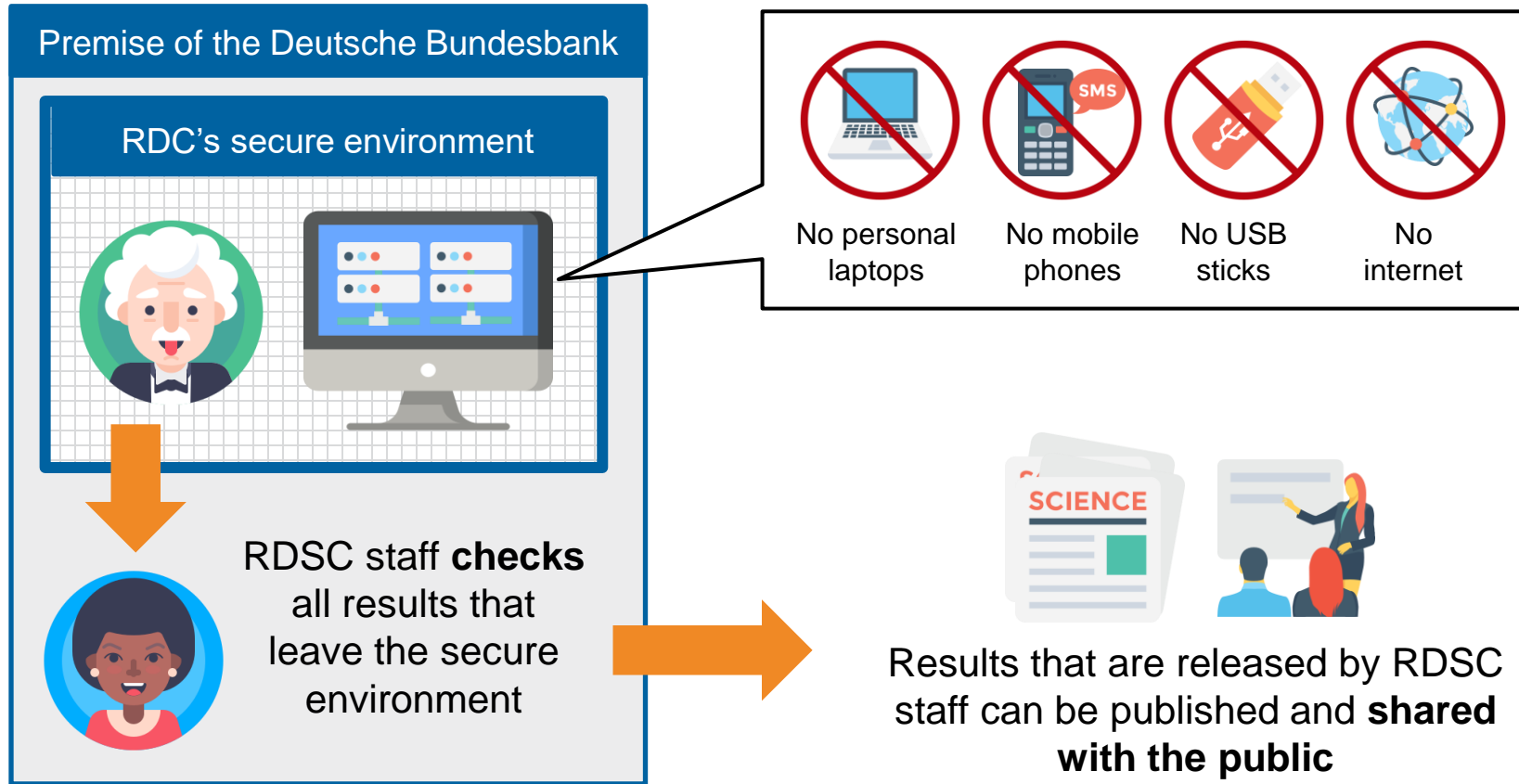


✉ sebastian.seltmann@bundesbank.de

✉ fdsz-data@bundesbank.de

🌐 <https://www.bundesbank.de/en/bundesbank/research/rdsc>

Appendix: RDCs allow to securely share confidential granular data from administrative sources with external researchers



Appendix: RDCs proliferate in recent years, both nationally and internationally, enabling high-quality research



RDCs in other institutions*



* *selected examples*



Appendix: The Prompts

Identification Prompt

A dataset is a collection of structured or unstructured data that is organized and grouped together for a specific purpose. It typically consists of multiple data points or observations related to a particular topic or subject. A dataset can include various types of information such as numerical values, text, images, audio, video, or any other form of data. It is often used in the context of data analysis, machine learning, and statistical research, where the data is utilized to extract insights, train models, or draw conclusions. Datasets can be generated through various means, including surveys, experiments, observations, or by gathering existing data from different sources.

The text after the empty lines is a scientific paper excerpt

According to the definition on the first line, search for any mentions of datasets or data-sources used in the paper's research.

If you find any, please compile a list of all mentioned datasets in this excerpt.

If there aren't any, please reply with 'None'.

Appendix: The Prompts

Consolidation Prompt

You will be provided one or more lists of datasets.
Each new list starts with '=>'.
You need to combine all the lists into a single one by removing redundant entries. Delimit the final list with simple '-' bullet points.

Recall Prompt

You will be provided with text delimited by triple quotes that is supposed to be a list of datasets.
Check if the following true datasets are directly contained in the answer:

...

For each of these true datasets perform the following steps:

- 1 - Restate the true dataset
- 2 - Write 'yes' if the true dataset is mentioned in the answer, otherwise write 'no'

Finally, provide a count of how many 'yes' findings there are. Provide this count as {"count":<insert count here>}.