

Extraction of tabular data from annual reports with LLMs

Using in context learning with open source models and RAG

submitted by

Simon Schäfer

Matr.-Nr.: 944 521

Department VI – Informatics and Media
Berliner Hochschule für Technik Berlin
presented Master Thesis
to acquire the academic degree

Master of Science (M.Sc.)

in the field of

Data Science

Date of submission September 1, 2025



Studiere Zukunft

Gutachter

Prof. Dr. Alexander Löser
Prof. Dr. Felix Gers

Berliner Hochschule für Technik
Berliner Hochschule für Technik

Abstract

Content of this thesis is a benchmark on information extraction from PDFs. The focus are annual reports of German companies. Special characteristic of the task is handling hierarchies in tables with financial data to prepare the data for import into a relational database.

The benchmark is composed of three sub tasks and the performance of different open source large language models is tested with different prompting approaches and compared to alternative methods.

This can be seen as a reimplementation study of “Extracting Financial Data from Unstructured Sources: Leveraging Large Language Models” - a paper published by Li et al. (2023). The key differences are the application on German documents using open source large language models.

Zusammenfassung

Gegenstand dieser Arbeit ist ein Benchmark zur Informationsextraktion aus PDF-Dateien. Dabei wird sich auf das Auslesen der Bilanzen und Gewinn- und Verlustrechnungen aus Jahresabschlüssen deutscher Unternehmen beschränkt. Ein besonderer Aspekt der Aufgabe ist die Berücksichtigung der Hierarchie innerhalb der Tabellen, um die Werte einem festen Schema zuzuordnen und so den Import in eine relationale Datenbank vorzubereiten.

Reading advices

The author recommends to read the thesis in its digital gitbook version instead of the PDF version. Furthermore, the author recommends to read the thesis (any version) on a screen that is larger than 21” and has at least full HD resolution¹. The more, the merrier.



Goals and Learnings

Achieved:

- thesis with bookdown
- docker image creation
- cluster orchestration
- llm usage
- guided decoding

¹Most of the time the thesis was inspected at a third of the authors 42” screen with 4k resolution. For inspecting the large overview graphics it is a very handy tool the author recommends every data scientist or software developer.

Missed:

- Administrating a k8s cluster
- Fine tuning a model
- using small language models
- training a lm
- using vllms

Contents

Contents	i
1 Introduction	1
1.1 Motivation	1
1.2 Objectives	1
1.3 Methodology (1 p)	2
1.4 Thesis Outline (0.5 p)	2
1.5 To place in chapters above	2
1.6 RHvB	3
1.7 Datenverfügbarkeit	3
1.8 Unstrukturierte Daten	3
1.8.1 Portable Document Format	3
2 Literature review (less than 10 p)	5
2.1 NLP history	5
2.2 Basic terms	5
2.3 Supervised Learning Approaches	5
2.3.1 Generalized Linear Models	5
2.3.2 Random Forest	5
2.3.3 Large Language Models	6
2.3.4 Information extraction	6
2.4 Data balancing	6
2.4.1 Under sampling, oversampling	6
2.5 Evaluation Metrics	6
2.5.1 For classification	6
2.5.2 For regression	6
2.6 Technological topic (related work)	6
2.7 Term frequency	7
2.7.1 Extraction of numeric values	7

2.8	optimal more topics like previous	7
2.9	Summary (0.5 p)	7
2.10	To place in chapters above	7
2.11	Table extraction tasks	7
2.11.1	Difficulties	7
2.12	Document Extrraction Process	7
2.12.1	Document Layout Analysis	7
2.12.2	8
2.13	Tools	8
2.13.1	TableFormer	8
3	Methods	9
3.1	Data	9
3.2	Page identification	9
3.2.1	Baselines	9
3.3	Table detection	10
3.3.1	LLM	10
3.3.2	Vision Model	10
3.3.3	Docling and Co	10
3.4	Information extraction	10
3.4.1	Baselines	10
3.4.2	Simple pipeline	10
3.4.3	Sophisticated approaches	10
4	Implementation (max 5p)	11
4.1	Speedup with vLLM and batching	11
4.2	Setup (Dockerfile and PV)	11
5	Results	13
5.1	Page identification	13
5.1.1	Baseline: Regex	15
5.1.2	Table of Contents understanding	17
5.1.3	Classification with LLMs	24
5.1.4	Term frequency based classifier	36
5.1.5	Comparison	38
5.2	Table extraction	38
5.2.1	Baseline: Regex	38
5.2.2	Extraction with LLMs	42
5.2.3	Comparison	54

6 Discussion	57
6.1 Limitations	57
6.1.1 table extraction	57
6.1.2 classification	57
6.2 Not covered	57
6.3 Outlook	57
6.3.1 Table extraction	58
7 Conclusion	59
References	61
List of Figures	63
List of Tables	65
Glossary	67
A Appendix	69
A.1 Local machine	69
System Details Report	69
A.2 Benchmarks	70
A.2.1 Text extraction	70
A.2.2 Table detection	70
A.2.3 Large language model process speed	74
A.3 Prompts	74
A.3.1 TOC understanding	74
A.4 Regular expressions	75
A.5 Figures	76
A.5.1 Page identification	76
A.5.2 Table extraction	76
A.6 Annual Comprehensive Financial Report Balance Sheet	92
A.7 Extraction framework flow chart	92
A.8 Table extraction with regular expressions	92

Chapter 1

Introduction

1.1 Motivation

- market: public administration, companies with data of special requirements for treating (secret and personal data (high risk data)) <- DSGVO, AI act
 - next market for hyper scalers might be public administration with local computing clusters
- whom is it helping
- why now: digital sovereignty, AI act; people want NLP AI products, frameworks get easier
- is the problem easier solvable then years ago? why?

missing law to access digital data and no law to choose the format of the data extensible Business Reporting Language as a standard changing from HGB to IFSR

Land Berlin							
Kredit- und Versicherungswirtschaft	Wohnungswirtschaft	Landesentwicklung und Grundstücksverwaltung	Verkehr und Dienstleistungen	Ver- und Entsorgungswirtschaft	Kultur und Freizeit	Wissenschaft und Ausbildung	Gesundheit und Soziales
IBB Unternehmensverwaltung Gewährträger: Berlin	degewo AG 100%	Berlinovo Immobilien Ges. mbH 100%	Amt für Statistik Berlin-Brandenburg, Gewährträger: Bln. u. Brandenbg.	BEN Berlin Energie und Netzholding GmbH 100%	BBB Infrastrukt. Verw. GmbH 100%	Dr. Film- u. Fernsehakadem. GmbH 100%	Berliner Werkst. f. Beh. GmbH 70%
	GESBAU AG 100%	BIM GmbH 100%	BEHALA GmbH 100%	Berl. Stadtreinigungsbetriebe, Gewährträger: Berlin	BBB Infrastrukt. GmbH & Co. KG 100 % Kommanditist: Berlin	Deutsches Zentrum f. Hochschul- u. Wiss.forschung GmbH 1,85%	Vivantes GmbH 100%
	Gewobag AG 96,69%	Berliner Städtegüter GmbH 100%	Berlin Tourismus & Kongress GmbH 15%	Berliner Wasserbetriebe, Gewährträger: Berlin	Berliner Bäder-Betriebe, Gewährträger: Berlin	Ferdinand-Braun-Institut gGmbH 100%	
	HOWOGE GmbH 100%	Campus Berlin-Buch GmbH 50,1%	Berliner Energieagentur GmbH 25%	Berlinwasser Holding GmbH 100%	Friedrichstadt-Palast GmbH 100%	FWU Institut für Film GmbH 6,25%	
	STADT U. LAND GmbH 100%	Grün Berlin GmbH 100%	Berliner Großmarkt GmbH 100%	MEAB GmbH 50%	Hebbel-Theater GmbH 100%	Helmholtz-Zentrum Bln. GmbH 10%	
	WBM GmbH 100%	Liegenschaftsfonds GmbH 100%	Berliner Verkehrsbetriebe, Gewährträger: Berlin	SBB Sonderabfall GmbH 25%	KuJ Wuhlheide gGmbH 100%	Wissenschaftszentrum gGmbH 25%	
		Liegenschaftsfonds KG 100 % Kommanditist: Berlin	BGZ GmbH 60%	Kulturprojekte Berlin GmbH 100%			
		Liegenschaftsfonds Projekt KG 100 % Kommanditist: Berlin	DEGES Dt. Einheit Fernrohren- planungs- u. -bau GmbH 5,91%	Kunsthalle BR Deutschland, GmbH 2,44%			
		Olympiastadion Berlin GmbH 100%	Deutsche Klassenlotterie, Gewährträger: Berlin	Musikboard Berlin GmbH 100%			
		Tegel Projekt GmbH 100%	Flughafen Berlin-Brandenburg GmbH 37%	Rundfunk-Orchester gGmbH 20%			
		Tempelhofer Projekt GmbH 100%	IT-Dienstleistungszentrum Berlin, Gewährträger: Berlin	Zoologischer Garten Berlin AG 0,03%			
		WISTA-Management GmbH 100%	Landesamt Schienenfahrzeuge Berlin, Gewährträger: Berlin				
			Messe Berlin GmbH 100%				
			Partner für Deutschland 1%				
			VBB GmbH 33,33%				

Figure 1.1: Overview of companies Berlin holds share at

1.2 Objectives

The sixth division at RHvB is auditing the companies Berlin is a stakeholder of. Basic information they have to process are the balance sheets and profit and loss accounting. Those information is provided via their

annual reports in form of PDF files. The provided annual reports often differ from the publicly available ones in matter of information granularity and design and are treated as non public information. Automate the extraction of those information would be a good starting point for AI assisted information retrieval from PDFs for the RHvB overall.

It is important to get numeric values totally accurate; numeric values are difficult to handle for language models

- special part of big problem? central question
- two sentences: why this problem? new problem or just a part in the big task? hard to solve of straight forward? research or application? what was not done and why?
- building a system? what task to solve? core functionality? typical use cases?

Research questions and hypotheses

Q1: Can a LLM (large language model) be used to efficiently extract financial information from German annual reports? Q2. Can LLMs be used to identify the page of interest automatically?

Q3: Can confidence scores be used to head up the human in the loop on which results to double check? (How can sources of the automatic extraction being communicated down stream in order to make double checking easy before making decisions?) Q4: Can contextual information from similar documents reduce errors made during table extraction? Q5: What are characteristics of financial tables that make it hard for LLMs to identify / extract them? (How does the length and complexity of financial documents (e.g., multi-column layouts, nested tables) affect table extraction performance?)

1.3 Methodology (1 p)

- how to solve the problem?
- what foundations to have in mind?
- proceeding?

Experimental / Comparative Research • Reimplementing framework(s) • Comparing / Benchmarking • Frameworks • Models • Methods • Use cases • Ablation test

1.4 Thesis Outline (0.5 p)

1.5 To place in chapters above

This master thesis is motivated by a use case from practical work at the Berlin court of audit (Rechnungshof von Berlin; RHvB). The auditors often are faced with the problem that they need information that is provided as natural language or in tables inside of unstructured documents, i.e. in PDF files. The goal of this thesis is benchmarking methods for automated information extraction from specific tables from PDF files.

Ideally, the data extraction pipeline is able to autonomously * identify the pages with the tables of interest. * identify the tables of interest on these pages. * extract the information as provided into a structured table (e.g. as JSON, a csv file or HTML code). * transform the data into a given schema, stripping all aggregated values.

It should extract the values without errors. It would be nice if the computation time and energy consumption is as low as possible.

A more realistic approach, that is also beneficial to satisfy the AI Act (keine Entscheidung ohne menschliche Beteiligung), is an assistant system, that helps extracting information. Key features to get the human into the loop already at the step of information extraction for such an assistant might be:

- showing the results together with the systems confidence.
- showing the results next to the values of the source.
- allowing in place adjustments to the extracted data.

A sound decision making is only possible if the information the decision is based on is valid.

1.6 RHvB

- what does the RHvB do
- why is this important
- what does it not do yet (because data source is missing)

1.7 Datenverfügbarkeit

- keine Regelung, in welcher Form der Rechungshof die Daten, die er benötigt, bereitgestellt zu bekommen hat

Das Gesetz zur Förderung der elektronischen Verwaltung (EGovG) wurde erlassen, "um die Verwaltung effektiver, bürgerfreundlicher und effizienter zu gestalten." (BMI, Referat O2, 2013)

§ 12 EGovG

- Vorhaben zur Datenkatalogisierung innerhalb der Verwaltung angestoßen, aber noch nicht richtig gestartet
- Vornehmlich für Bürger*innen Zugang

1.8 Unstrukturierte Daten

- Beispielbilder

1.8.1 Portable Document Format

- print optimized
- Table structure information gets lost
- Bild und Textextract

Chapter 2

Literature review (less than 10 p)

(5 to 10 lines)

- overview of subchapters
- relevance for reader (Gutachter)
- link to previous chapter
- relevant basic tasks
- parameter vs active parameter

2.1 NLP history

2.2 Basic terms

2.3 Supervised Learning Approaches

2.3.1 Generalized Linear Models

2.3.2 Random Forest

XGBoost not used finally, because calculation SHAP (SHapley Additive exPlanations) values for XGBoost model took to long for just a first glimpse on what might influence the extraction.

2.3.3 Large Language Models

2.3.3.1 Embeddings

2.3.3.2 Neural networks in NLP

2.3.3.3 Attention / Multi-Head

2.3.3.4 Transformers

2.3.3.5 Encoder

2.3.3.6 Decoder

2.3.3.7 BERT

2.3.3.8 Bi-Encoder

2.3.3.9 Mixture of Experts

2.3.3.10 Guided decoding

generation template strict (closed) vs open

2.3.3.11 Classification trained models (not used)

Soft max

2.3.3.12 Few-shot Learning

2.3.3.13 RAG

2.3.3.14 GPT (Generative Pretrained Transformers)

2.3.4 Information extraction

closed-domain vs open-domain

2.4 Data balancing

2.4.1 Under sampling, oversampling

2.5 Evaluation Metrics

2.5.1 For classification

2.5.2 For regression

2.6 Technological topic (related work)

- LLM generation

- structured output
- Fewshot
- context length can be harmful
- most important papers
- connection of papers (timeline)
- what used, what not?
- extending existing paper?

2.7 Term frequency

2.7.1 Extraction of numeric values

99.5 % or 96 % accuracy for extracting financial data from Annual Comprehensive Financial Reports (Li et al., 2023) In the untabulated test, GPT-4 achieved an average accuracy rate of 96.8%, and Claude 2 achieved 93.7%. Gemini had the lowest accuracy rate at 69%. (ebd.)

Too many hallucinated values when it was NA instead (Grandini et al., 2020)

2.8 optimal more topics like previous

2.9 Summary (0.5 p)

- lessons learned
- link to goal thesis
- link to next chapter

2.10 To place in chapters above

2.11 Table extraction tasks

2.11.1 Difficulties

- Beispielbilder

2.12 Document Extraction Process

2.12.1 Document Layout Analysis

An important step in the process of extracting information from documents is to recognize the layout of a document (Zhong et al., 2019).

Getting the order of texts correct align captions to tables and figure identify headings, tables and figures

One of the most popular datasets used for training and benchmarking is PubLayNet (see PubLayNet on paper-withcode.com). It contains over 360_000 document automatically annotated images from scientific articles publicly available on PubMed Central (Zhong et al., 2019, p. 1). This was possible, because the articles have been provided in PDF and XML format. For the annotations most text categories (e.g. text, caption, footnote) have been aggregated into one category. <- is this a problem for later approaches where a visual and textual model work hand in hand to identify e.g. table captions?

Manual annotated datasets often were limited to several hundred pages. Deep learning methods need a much larger training dataset. Previously optical character recognition (OCR) methods were used.

Identify potentially interesting pages with text / regex search. Check if there is a table present on this page.

Object detection

2.12.1.1 Vision Grid Transformer

2.12.2

2.13 Tools

2.13.1 TableFormer

SynthTabNet <- has it: - nested / hierarchical tables, where rows add up to another row? - identifying units and unit cols/rows

Chapter 3

Methods

norm gpu hours

3.1 Data

- companies Beteiligungsbericht
- number found Jahresberichte
- number used Jahresberichte first rows
- number used Jahresberichte Aktiva Tabellen

3.2 Page identification

Due to the imbalanced distribution of the classes the accuracy is not a good metric to compare the performance of the different methods. The number of pages of interest is much smaller than the number of irrelevant pages. Therefore, precision, recall and F1 score are presented as well.

3.2.1 Baselines

3.2.1.1 Regex based

results potentially dependend on package used for text extraction (Auer et al., 2024, p. 2 f.)

- PyMuPDF
- pypdf
- docing-parse
- pypdfium
- pdfminer.six

pdfminer informs that some pdfs should not be extracted based on their authors will (meta data field)

results dependend on regex pattern

start with pypdf backend and simple regex developed more sophisticated regex based on missed pages

took wrong identified pages as base for a table detection benchmark and n-shot base for llm classification (contrasts)

some tables can't be found without previous ocr; some pages hold image of table and machine readable text

LLM based**3.2.1.2 Term frequency based**

VLLM based was not implemented

3.3 Table detection

Can be used to narrow down set of possible pages

Can be used to focus only on the table content (measure if correct area was identified would be necessary)

Vision model as baseline

3.3.1 LLM

- table: yes/no
- akiva: yes/no
- multiclass

3.3.2 Vision Model

Yolo

3.3.3 Docling and Co

VLLM based was not implemented

3.4 Information extraction

3.4.1 Baselines

simple regex?

3.4.2 Simple pipeline

- extract text (if document can't be passed directly)
- query LLM directly

3.4.3 Sophisticated approaches

not implemented

- with pipelines
- Nougat
- maker
- Azure
- docling

Chapter 4

Implementation (max 5p)

4.1 Speedup with vLLM and batching

4.2 Setup (Dockerfile and PV)





Chapter 5

Results

This chapter presents the results for the two research questions of this thesis:

1. How can we use LLMs effectively to locate specific information in a financial report?
2. How can we use LLMs effectively to extract these information from the document?

Section 5.1 presents the results for the first research question. Section 5.2 presents the results for the second question.

Each section will start with an overview about the specific sub tasks as well about the models, methods and data used to investigate the research question. The subsections present the results of the sub tasks. At the end of each section all results get compared and summarized.

5.1 Page identification

The first research question asks, how LLMs can be used, to effectively locate specific information in a financial report. The task for this thesis is identifying the pages where the balance sheet (*Bilanz*) and the profit-and-loss-and-statement (*Gewinn- und Verlustrechnung, GuV*) are located. The balance sheet is composed of two tables showing the assets (*Aktiva*) and liabilities (*Passiva*) of a company. Often these two tables are on separate pages. Hereafter, the German terms **Aktiva**, **Passiva** and **GuV** (Gewinn- und Verlustrechnung) will be used.

Li et al. (2023) describes two ways to identify the relevant pages (see Figure A.19). For longer documents they propose to use the TOC (table of contents) to determine a page range that includes the information of interest. In addition, they develop target specific regular expressions and rules to filter out irrelevant pages¹. The result of this “Page Range Refinement” is then passed to the LLM to extract information from.

This section is presenting four approaches to identify the page² of interest.

- Subsection 5.1.1 presents the performance of a page range refinement using a list of key words with a regular expression.
- Subsection 5.1.2 presents the performance of a TOC understanding approach
- Subsection 5.1.3 presents the performance of a text classification using LLMs.
- Subsection 5.1.4 presents the performance of a term-frequency approach.

¹Personal opinion: Developing well performing regular expressions can be a very tedious and setting appropriate rules requires some domain knowledge. It can be worth the effort if there are a lot of documents with similar information to extract. For this thesis it took multiple months. At least, now there is kind of a pipeline one can reuse, exchanging the rules and key word lists. Thus the next similar task should be solved faster.

²In some cases the information of interest is spanning two pages. These rare cases are not covered from the approaches presented here, yet.

In subsection 5.1.5 the results get compared and summarized. Subsection @ref() proposes an efficient combination of approaches to solve the task of this thesis and discusses its limitations.

Figure 5.1 shows how the document base for most of the tasks in this section is composed³. Overall 74 annual reports from 7 companies are used. For this thesis the tables of interest are those that show **Aktiva**, **Passiva** and **GuV**. Among the 4981 pages 265 tables have to be identified on 252 pages. Figure 5.1 also gives an impression on how many pages the documents have. The documents of *IBB* tend to be longer. The documents of *Amt für Statistik Berlin-Brandenburg* tend to be shorter.

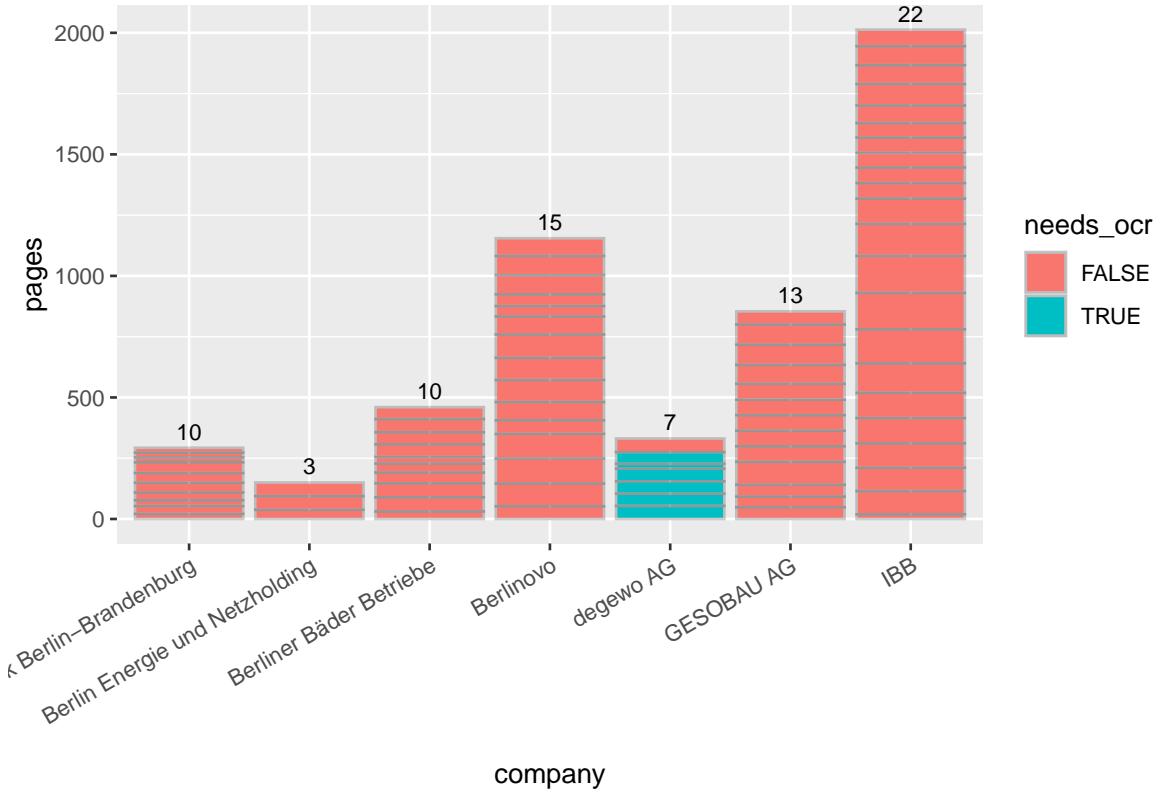


Figure 5.1: Showing the number of pages (bar height) and number of documents (number above the bar) per company for the data used for the page identification task. Some documents would require ocr before being processed and were not used.

Table @ref(tab:display_multiple_tables_per_type_and_document) shows how many documents have multiple target tables per type and how many target tables span two pages. In total 21 tables are distributed on two pages. In 8 documents there are multiple tables per type of interest. There are 13 pages with two target tables (**Aktiva** and **Passiva**) on it.

This task is broken down to a classification task all of the approaches presented in this section but the TOC understanding approach.

Thus, we prompt the LLM to classify if the text extract of a given page

for implementation: As described in A.2.1 open source libraries have been used to extract the text from the annual reports.

³I downloaded all publicly available annual reports for some of the companies shown in the first row of Figure 1.1. I assumed that this will give a representative sample of document structures for the other companies of the same type. Realizing that the degewo AG reports would require ocr preprocessing I additionally downloaded reports for GESOBAU AG. This approach could have been more systematic. For the second task I downloaded reports for all companies available and tried to use a balanced amount of reports per company.

Table 5.1: Showing the number of documents with multiple target tables per type and the number of target tables that span two pages.

type	multiple targets in document	target two pages long
Aktiva	7	1
GuV	8	20
Passiva	7	0

Table 5.2: Comparing page identification metrics for different regular expressions for each classification task by type of the target table.

method	type	precision	recall	F1
Aktiva				
simple regex	Aktiva	0.273 ± 0.005	0.788 ± 0.010	0.403 ± 0.005
exhaustive regex restricted	Aktiva	0.190	0.990	0.320
exhaustive regex	Aktiva	0.132 ± 0.004	0.997 ± 0.005	0.233 ± 0.008
Passiva				
simple regex	Passiva	0.400 ± 0.009	0.780 ± 0.009	0.530 ± 0.009
exhaustive regex restricted	Passiva	0.190	0.980	0.320
exhaustive regex	Passiva	0.130 ± 0.000	0.993 ± 0.010	0.230 ± 0.000
GuV				
simple regex	GuV	0.180 ± 0.006	0.938 ± 0.008	0.302 ± 0.010
exhaustive regex restricted	GuV	0.210	1.000	0.350
exhaustive regex	GuV	0.173 ± 0.008	1.000 ± 0.000	0.295 ± 0.012

5.1.1 Baseline: Regex

The first approach presented in this section is, to use a key word list and regex (regular expression) to filter out irrelevant pages. It is setting the performance baseline for the following approaches. Building a sound regular expression often is an iterative process. In a first approach a very *simple regex* was implemented. To increase the recall to 1.0 the regular expression was extended⁴. This second regex is called *exhaustive regex*. In a third attempt minor changes have been made to the *exhaustive regex* to increase the precision without decreasing the recall. This regular expression is called *exhaustive regex restricted*. The regular expressions can be found in the appendix (see section A.4).

Table 5.2 shows the mean performance for precision, recall and F1 for the three regular expressions for the three types of pages to identify⁵. It was possible to create a regular expression that has a high recall for all target types. The precision is low for all tested regular expressions and target types. Figure 5.2 gives insight into performance differences between the companies. There is only one document from *Berlin Energie und Netzholding* where the **GuV** is not identified except with the *exhaustive regex restricted*⁶.

The regular expressions have been tested on the texts extracted with multiple Python libraries. The reported standard deviations are very small. This means that there are no substantial differences in the extracted texts on a word level⁷. But table A.1 in section A.2.1 shows that there are differences in the extraction speed.

Code can be found at “benchmark_jobs/page_identification/page_identification_benchmark_regex.ipynb”

Todo: * look into details where they differ and if it is because of a line break or whitespace?

⁴The idea is that the regular expression approach is computationally cheap. If we can rely on the fact, that it keeps all relevant pages we can use additional, computationally more expensive approaches to further refine the page range.

⁵See Figure A.1 for a graphical representation.

⁶I don't understand why the restricted version is finding the page but the non-restricted regex is not.

⁷Since the results are not depending on the text extraction library, the *exhaustive regex restricted* ran only with the text extracted by the fastest extraction library: *pdflum*. This library is used for the most tasks in this thesis. Later faced issues with the text extracted by *pdflum* are discussed in @ref{}

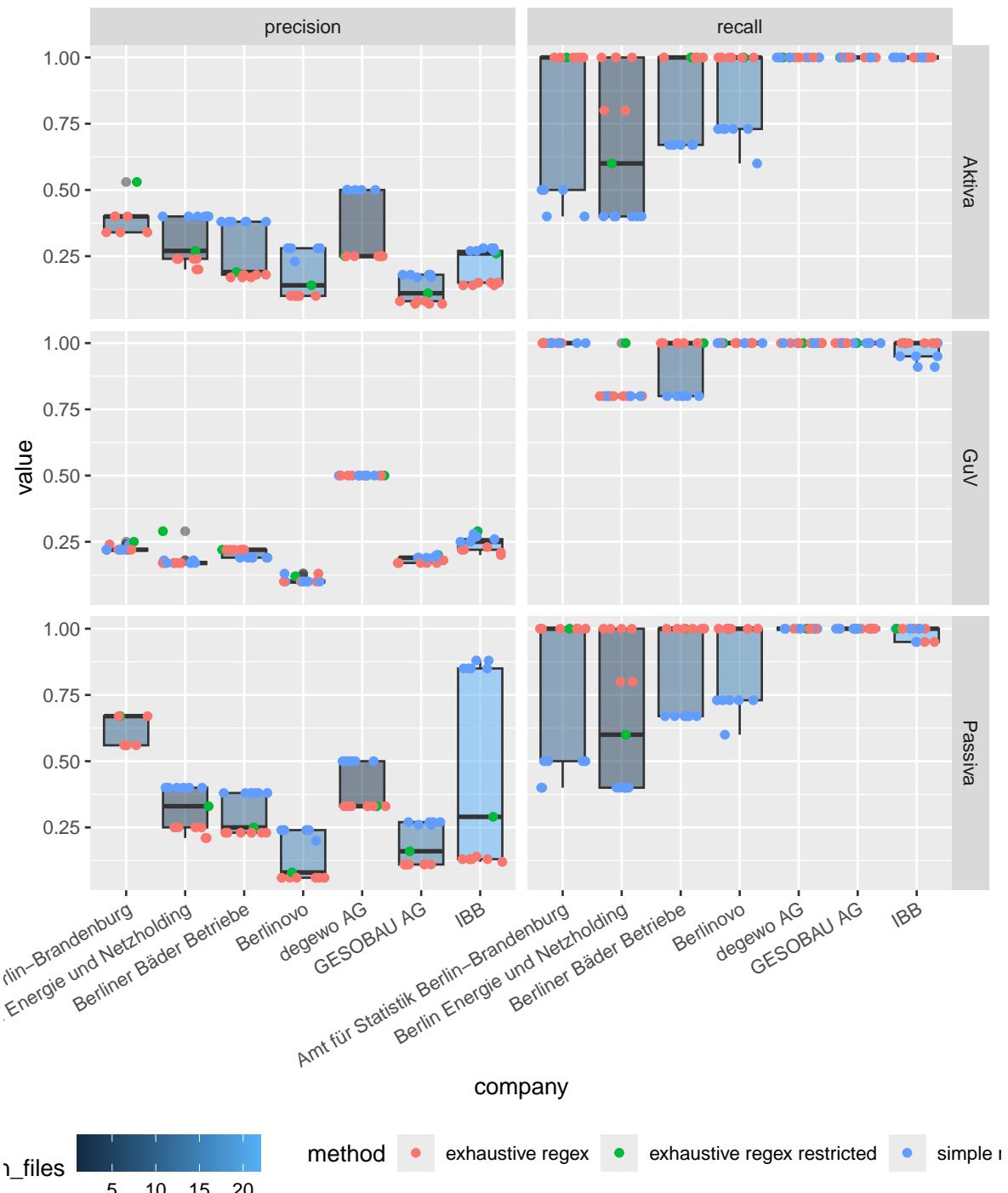


Figure 5.2: Comparing the performance among different companies.

5.1.2 Table of Contents understanding

The second approach presented in this section leverages the TOC understanding capabilities of LLMs. Li et al. (2023) use this approach with long documents as a first step to determine a page range of interest. If the predicted page range is correct and narrow, this approach is more efficient than processing the whole document with a LLM directly. The TOC in a PDF (Portable Document Format) document can be embedded in a standardized, machine readable format or be presented in varying, human readable forms of text on any page. Of course there are documents without any TOC.

Thus, the task is investigated based on two different input data formats In one case the LLM is provided with text extracted from the beginning of the document. In the other case the LLM is provided with the Markdown formatted version of the machine readable TOC embedded in the document. Subsection 5.1.2.1 shows the results for the text based approach. Subsection 5.1.2.1 shows the results for the approach, using the embedded TOC.

Additionally, each approach is performed three times with minor changes in the prompt. The prompts used for both approaches can be found at A.3.1. The prompt was adjusted two times to tackle shortcomings in the results. The first change adds the information, that assets and liabilities are part of the balance sheet. It is the balance sheet, that is listed in the TOC - not the assets or liabilities itself. The second change specifies the information, that assets and liabilities are often on separated pages, into, liabilities often are found on the page after the assets.

The code can be found in:

- “benchmark_jobs/page_identification/toc_extraction_mistral.ipynb”
- “benchmark_jobs/page_identification/toc_extraction_qwen.ipynb”

Discussion:

- Li et al. (2023) did not report any issues with this approach. They use few-shot learning and Chain-of-Thought techniques to help the LLM to understand the task. They ask just for one information at a time.
- ChatGPT 4 vs Mistral 2410 8B (huge parameter difference)
- For a lot of short annual reports one can find the tables of interest within the first eight pages as well.

5.1.2.1 Details for the approaches

Text based Li et al. (2023) used the TOC to identify the pages of interest. In their approach the table of contents is extracted from the text. Based on their observation, that the TOC in ACFR (Annual Comprehensive Financial Report)s is found within the initial 165 lines of the converted document (Li et al., 2023, p. 20), they use the first 200 lines of text.

My initial expectation was to find the TOC within the first five pages. Often there are way less than 200 lines of text on the five first pages (see Figure 5.3). In my approach the first step is to prompt the LLM to identify and extract the TOC in a given text extract^ [The prompt can be found in section A.3.1]. For the same documents Mistral 2410 8B finds^ [The strings extracted in this step have not been checked in detail.]

- 63 strings that should represent a table of contents among the first five pages.
- 68 strings that should represent a table of contents among the first 200 lines.

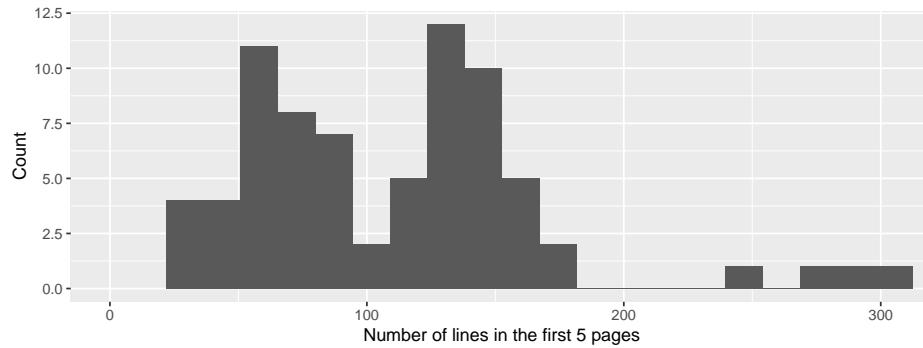


Figure 5.3: Histogram of the number of lines in the first 5 pages of the annual reports

Table 5.3: Comparing the number and percentage of correct identified page ranges among the approaches.

benchmark_type	type	n_correct	n	n_total	perc_correct	perc_correct_total
200 lines	Aktiva	9.0	63	82	14.3	11.0
200 lines	GuV	22.0	95	102	23.2	21.6
200 lines	Passiva	{6.0}	62	81	9.7	{7.4}
5 pages	Aktiva	7.0	58	82	12.1	8.5
5 pages	GuV	15.0	89	102	16.9	14.7
5 pages	Passiva	3.0	57	81	5.3	3.7
machine readable	Aktiva	{22.0}	35	82	{62.9}	{26.8}
machine readable	GuV	{28.0}	56	102	{50.0}	{27.5}
machine readable	Passiva	4.0	34	81	{11.8}	4.9

Machine readable TOC based I also tested to use the TOC representation embedded within the PDF files. First, this limits the text amount to process. Second, this hopefully increases the quality of the data passed to the LLM. 43 of the 80 annual reports have a machine readable embedded TOC. The embedded TOC is converted into markdown format before it gets passed to the LLM. Here is an example:

```
## | hierarchy_level | title | page_number | enumeration |
## |-----:|:-----|-----:|-----:|
## | 1 | Lagebericht | 5 | 1 |
## | 1 | Bilanz | 7 | 2 |
## | 1 | Gewinn- und Verlustrechnung | 10 | 3 |
## | 1 | Anhang | 13 | 4 |
## | 1 | Lagebericht | 17 | 5 |
## | 1 | Bilanz | 25 | 6 |
## | 1 | Anhang | 31 | 7 |
## | 1 | Anlagenspiegel | 39 | 8 |
## | 1 | Bestätigungsvermerk | 42 | 9 |
```

5.1.2.2 Results

Comparison of the different approaches: base prompt Table 5.3 shows that the machine readable TOC approach has the highest rate of correct page ranges for all types with the base prompt. It also predicts the most correct page ranges in absolute numbers for **Aktiva** and **GuV**. Thus, it also has the highest rate of correct page ranges based on the total number of page ranges to identify over all documents - no matter, if there was a TOC of any type in the document or not - for **Aktiva** and **GuV** of around 27 %.

Table 5.4: Comparing the number and percentage end pages prediction for Aktiva and Passiva that are equal.

benchmark_type	equal_end_page	n	perc_equal_end_page
200 lines	20	58	34.5
5 pages	26	53	49.1
machine readable	28	33	84.8

Figure 5.4 shows that the amount of correct predicted page ranges for **Passiva** is lowest for all approaches but can be improved by simply extending the predicted end page number by one the most. This improvement would be best for the machine readable TOC approach. This approach is the only one, where the number of correct page ranges **Aktiva** would not increase if we extend its range by one. Table 5.4 shows that this is the case, because the machine readable TOC approach predicts the same end page for **Passiva** as for **Aktiva** in 84.8 % of the cases, even though the prompt for all approaches included the information, that **Aktiva** and **Passiva** are on separate pages.

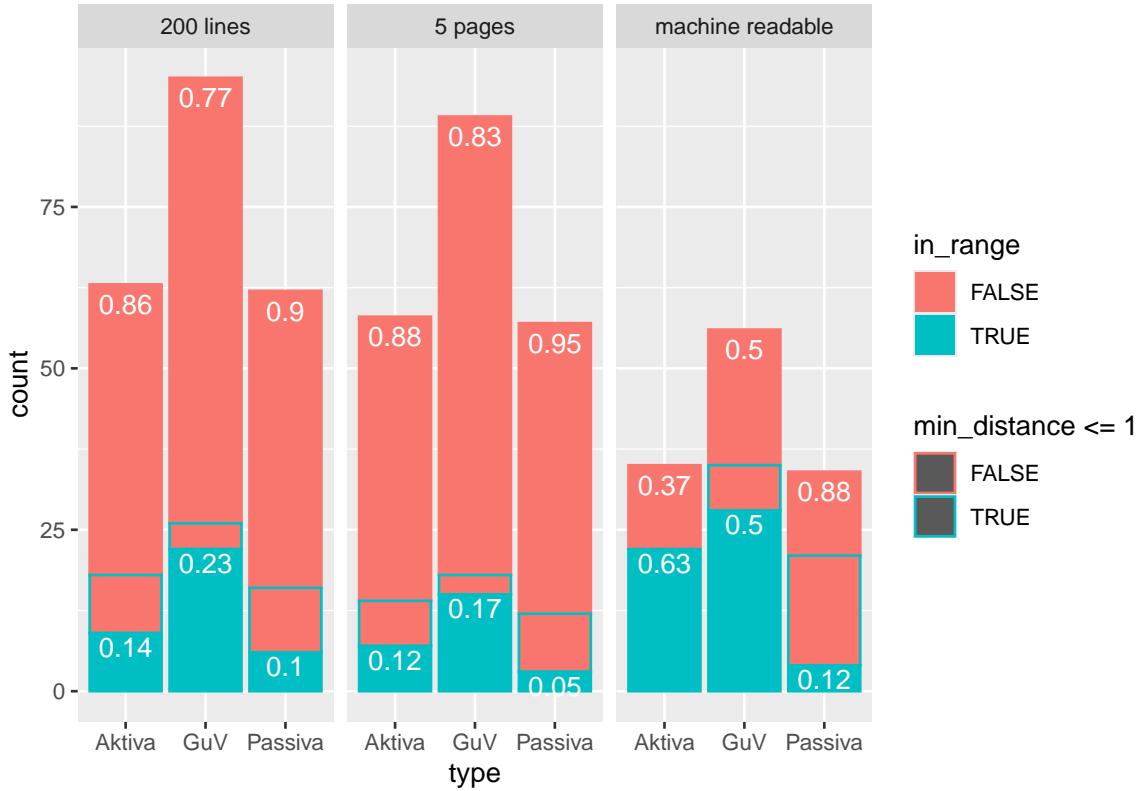


Figure 5.4: Comparing number of found TOC and amount of correct and incorrect predicted page ranges

Comparison of the different approaches: advanced prompts As a first attempt, to increase the correct page range rate for **Passiva** I tried to specify, that assets and liabilities are part of the balance sheet. This did work for the text based approaches, but not for the machine readable approach (see Figure A.3). Figure 5.5 shows that it is more successful, to explicit tell the LLM that the liabilities table is often on the page, after the assets table.

Table 5.5 shows the results from the final zero shot prompt. The machine readable TOC approach is now predicting best for all types. Nevertheless, a correct page range prediction rate below 60, 45, 50 % is still unsufficient to build downstream task on without human checkups. Table 5.6 shows, that the machine readable

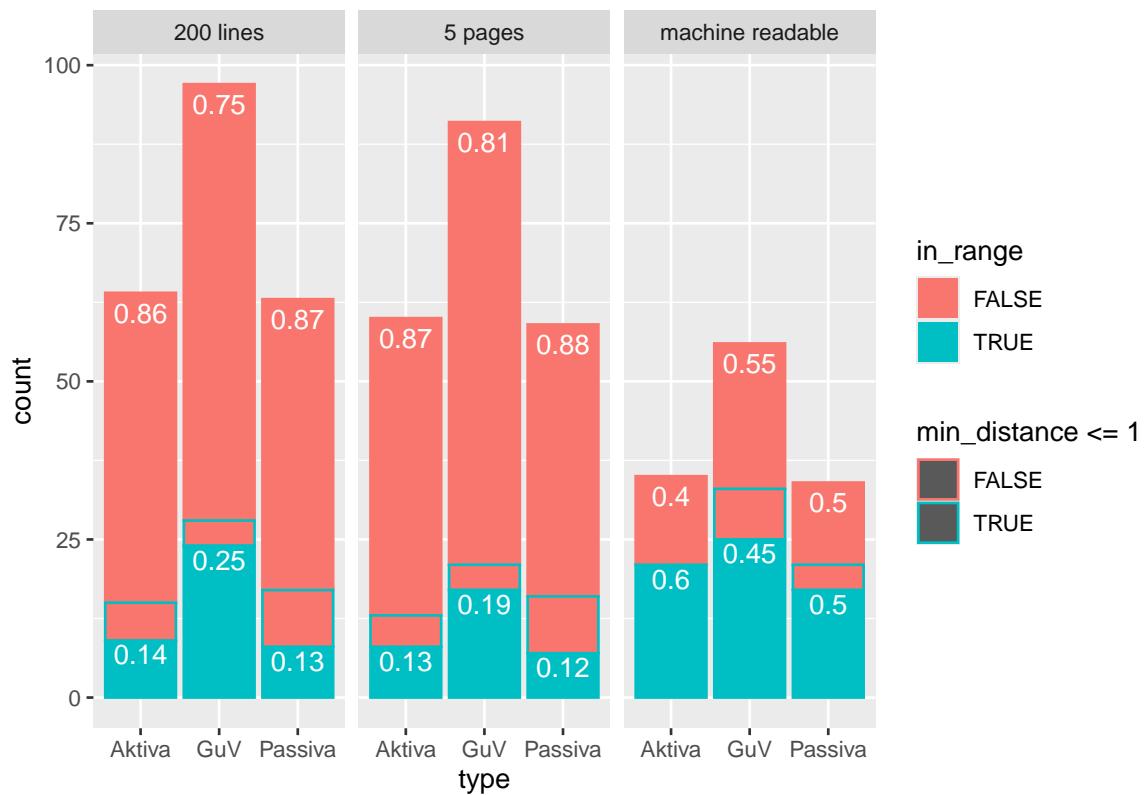


Figure 5.5: Comparing number of fount TOC and amount of correct and incorrect predicted page ranges

Table 5.5: Comparing the number and percentage of correct identified page ranges among the approaches.

benchmark_type	type	n_correct	n	n_total	perc_correct	perc_correct_total
200 lines	Aktiva	9.0	64	82	14.1	11.0
200 lines	GuV	24.0	97	102	24.7	23.5
200 lines	Passiva	8.0	63	81	12.7	9.9
5 pages	Aktiva	8.0	60	82	13.3	9.8
5 pages	GuV	17.0	91	102	18.7	16.7
5 pages	Passiva	7.0	59	81	11.9	8.6
machine readable	Aktiva	{21.0}	35	82	{60.0}	{25.6}
machine readable	GuV	{25.0}	56	102	{44.6}	{24.5}
machine readable	Passiva	{17.0}	34	81	{50.0}	{21.0}

Table 5.6: Comparing GPU time for page range prediction and table of contents extraction. Time in seconds per text processed.

Benchmark Type	Page range predicting	TOC extracting
200 lines	0.57	3.8
5 pages	{0.56}	{2.19}
machine readable	0.63	NA

TOC approach is the fastest as well.

Table 5.7 shows, that this advantage of the machine readable TOC approach is not coming from wide predicted page ranges. It has the smallest median range size among all approaches. Figure 5.6 shows, that especially the ranges for **GuV** are not normally distributed. Some far off lying range sizes are shifting the mean off from the median.

Figure 5.7 shows that the confidence of the LLMs responses is higher for the machine readable TOC approach as well. Besides a single group that was predicted far off, the page ranges are closer to the correct pages too. A linear regression of the correlation between minimal page distance and logistic probability shows that is has a similar slope for all approaches and target types.

Machine readable TOC approach specific results Figure 5.8 shows, that correct predictions for the page range are more probable when the embedded TOC has a medium number of entries. It is possible to drop documents with less than 9 without loosing a single correct prediction. This means that the LLM was not able to make a correct prediction for documents with TOC, that have less then 9 entries. This is not surprising since neither **Bilanz** nor **GuV** are mentioned there explicit.

Table 5.7: Comparing the mean and median page range sizes.

benchmark_type	type	mean_range	SD_range	median_range	MAD_range
200 lines	Aktiva	2.11	1.09	2	1.48
200 lines	GuV	4.25	3.29	4	2.97
200 lines	Passiva	1.7	0.59	2	0
5 pages	Aktiva	2.03	1.29	2	1.48
5 pages	GuV	3.15	2.17	2	1.48
5 pages	Passiva	1.64	0.89	2	0
machine readable	Aktiva	1.6	2.56	{1}	0
machine readable	GuV	3.89	5.75	{1}	0
machine readable	Passiva	1.24	0.74	{1}	0

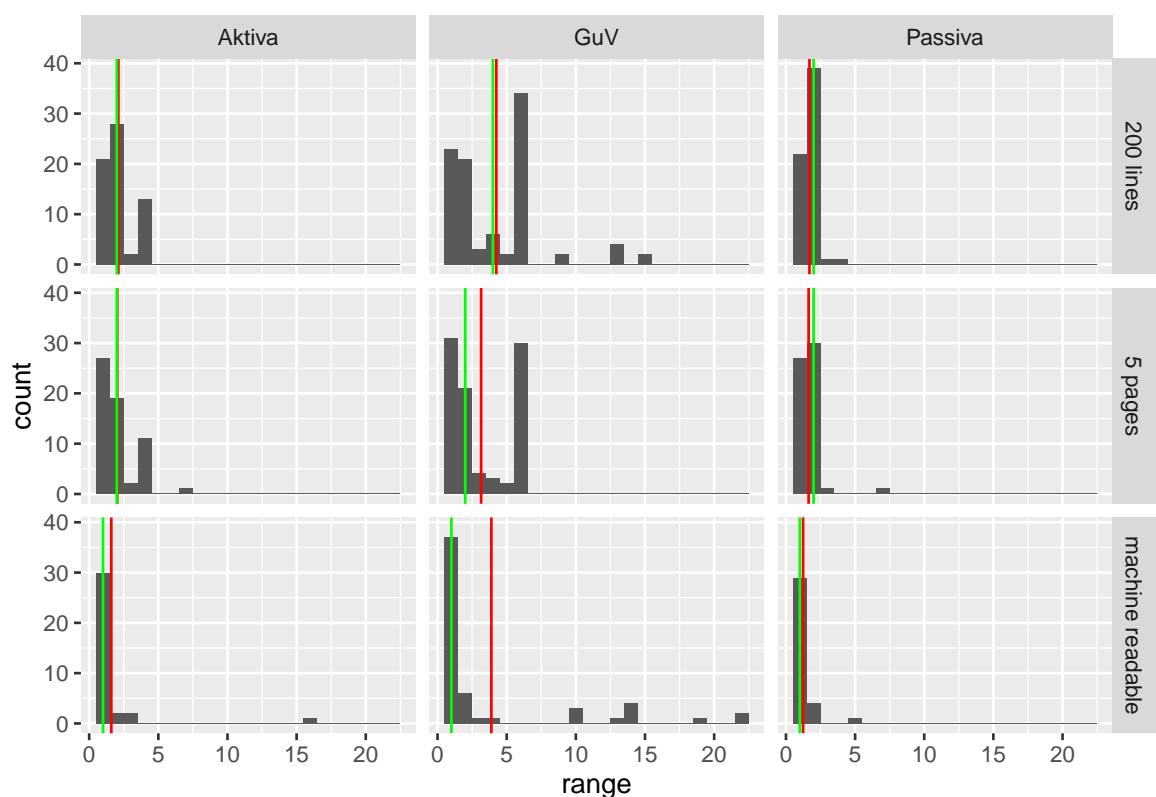


Figure 5.6: Comparing the predicted page range sizes. The red vertical line shows the mean and the green one shows the median of these sizes.

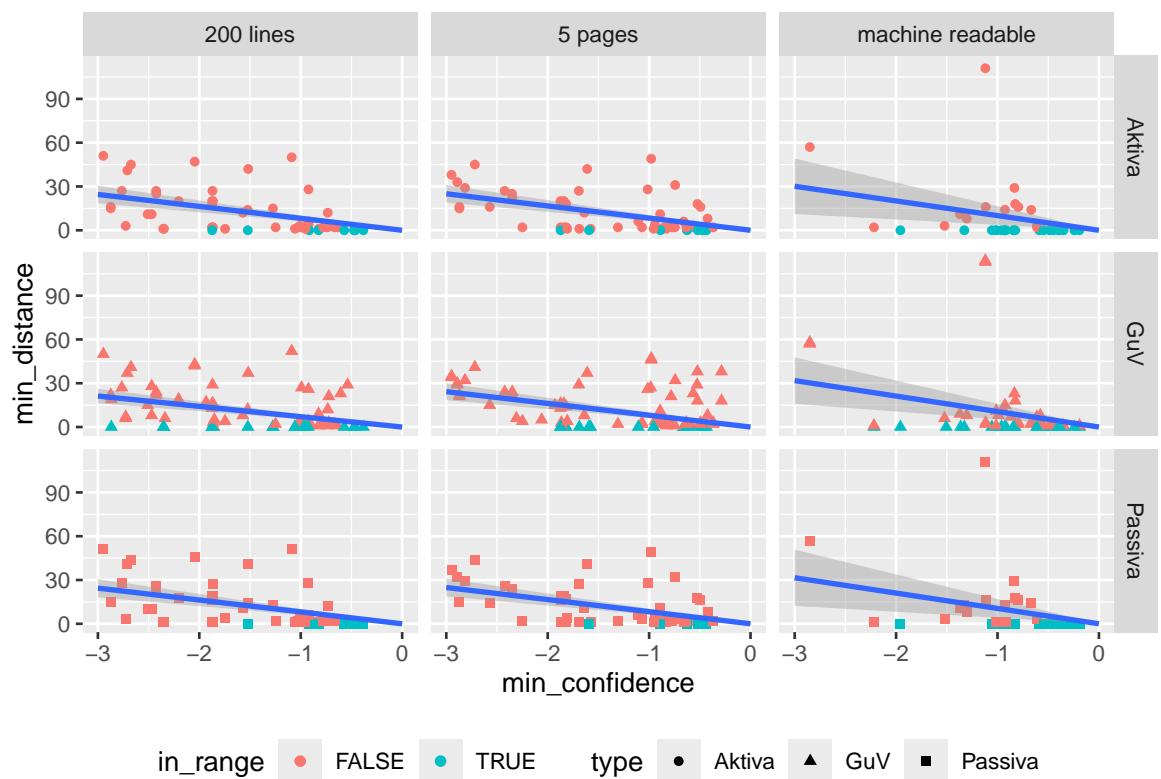


Figure 5.7: Showing the minimal distance of the predicted page range to the actual page number over the logprobs of the models response confidence.

It has no big influence on the predictions, if the TOC is passed formatted as markdown or json. With the json formatted TOC it found two more correct page ranges⁸. This was tested because the relation between heading and value for the column *page_number* might have been clearer⁹ in json for a one-dimensional working LLM.

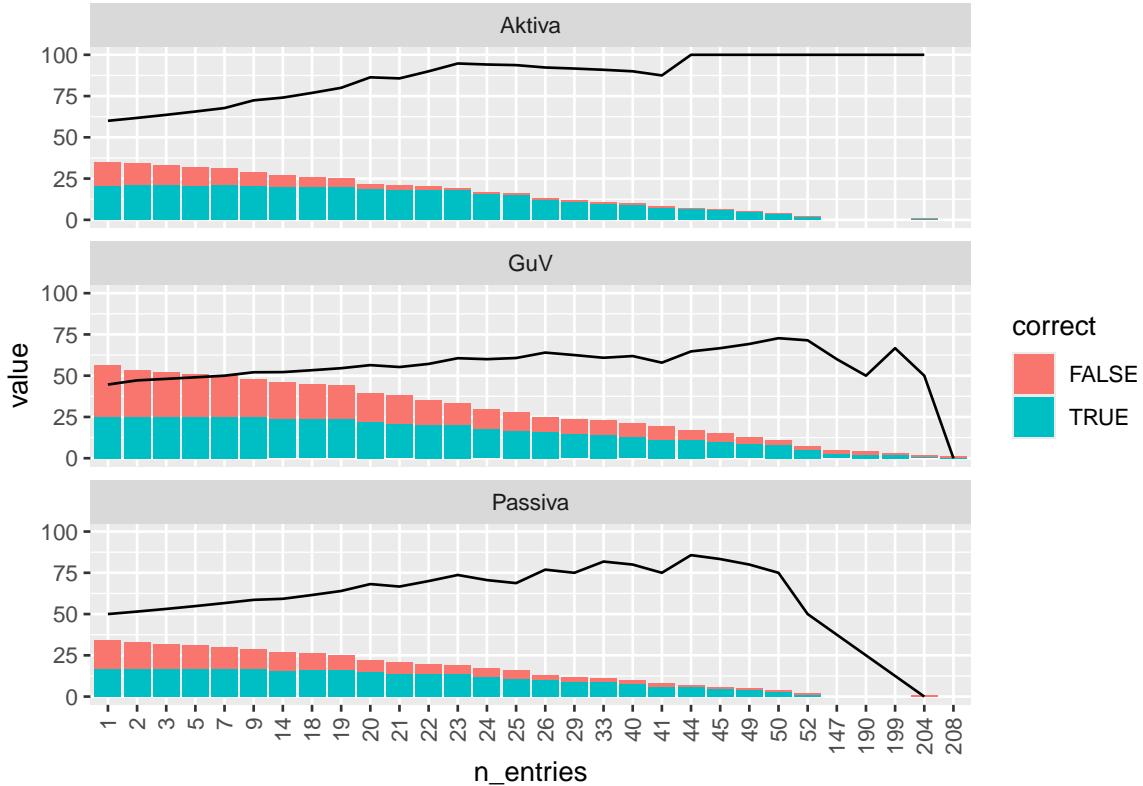


Figure 5.8: Showing the amount of correct and incorrect predicted page ranges (bars) and the percentage of correct predictions (black line).

Delete or place somewhere else?:

- Thus it is safer to go with the 200 lines approach. But it also takes longer. 5.6
- Values can be higher than 80, the total number of PDF files, since there can be multiple tables of interest for the same type in a single document or a table of interest can span two pages.

5.1.3 Classification with LLMs

structured outputs forcing to answer with a *yes* or *no* for binary task or with *Aktiva*, *Passiva*, *GuV* or *other* for multi classification task

top n accuracy

out of company vs in company rag

⁸This result is based on a single test run.

⁹With json the key *page_number* gets repeated every line, while it is just mentioned once in the beginning of the markdown formatted tables.

model_family	model	classification_type	method_family	n_examples	f1_score
mistralai	mistralai_Minstral8BInstruct2410	GuV	n_rag_examples	3	0.93
meta-llama	metallama_Llama4Scout17B16EInstruct	GuV	n_rag_examples	3	0.92
mistralai	mistralai_Minstral8BInstruct2410	Passiva	n_rag_examples	3	0.92
mistralai	mistralai_Minstral8BInstruct2410	Aktiva	n_rag_examples	3	0.92
Qwen	Qwen_Qwen2.532BInstruct	GuV	n_rag_examples	1	0.87
meta-llama	metallama_Llama4Scout17B16EInstruct	Passiva	n_rag_examples	3	0.85
Qwen	Qwen_Qwen2.532BInstruct	Aktiva	n_rag_examples	1	0.84
Qwen	Qwen_Qwen3235BA22BInstruct2507	Aktiva	n_rag_examples	3	0.84
meta-llama	metallama_Llama4Scout17B16EInstruct	Aktiva	n_rag_examples	3	0.83
Qwen	Qwen_Qwen2.532BInstruct	Passiva	n_rag_examples	1	0.79
microsoft	microsoft_phi4	Aktiva	law_context	1	0.68
microsoft	microsoft_phi4	Passiva	law_context	1	0.65
google	google_gemma327bit091	Passiva	n_rag_examples	1	0.54
google	google_gemma327bit091	Aktiva	n_rag_examples	1	0.51
tiuae	tiuae_Falcon310BInstruct	Passiva	n_random_examples	1	0.5
google	google_gemma327bit091	GuV	n_rag_examples	1	0.49
tiuae	tiuae_Falcon310BInstruct	Aktiva	n_rag_examples	1	0.44
tiuae	tiuae_Falcon310BInstruct	GuV	top_n_rag_examples	1	0.33

5.1.3.1 Binary classification

Could be more efficient to predict “is any of interest” and then which type, because dataset is highly imbalanced.

25 models from 6 have been benchmarked among 5 methods

Most models have been used up to 3 examples for the context.

The best combination of model and method for each method family is presented in the following table. It is clear that the Google Gemma models are performing worst.¹⁰ Surprisingly Mistral 2410 is the best performing model for all three prediction tasks even though it only has 8B parameters.

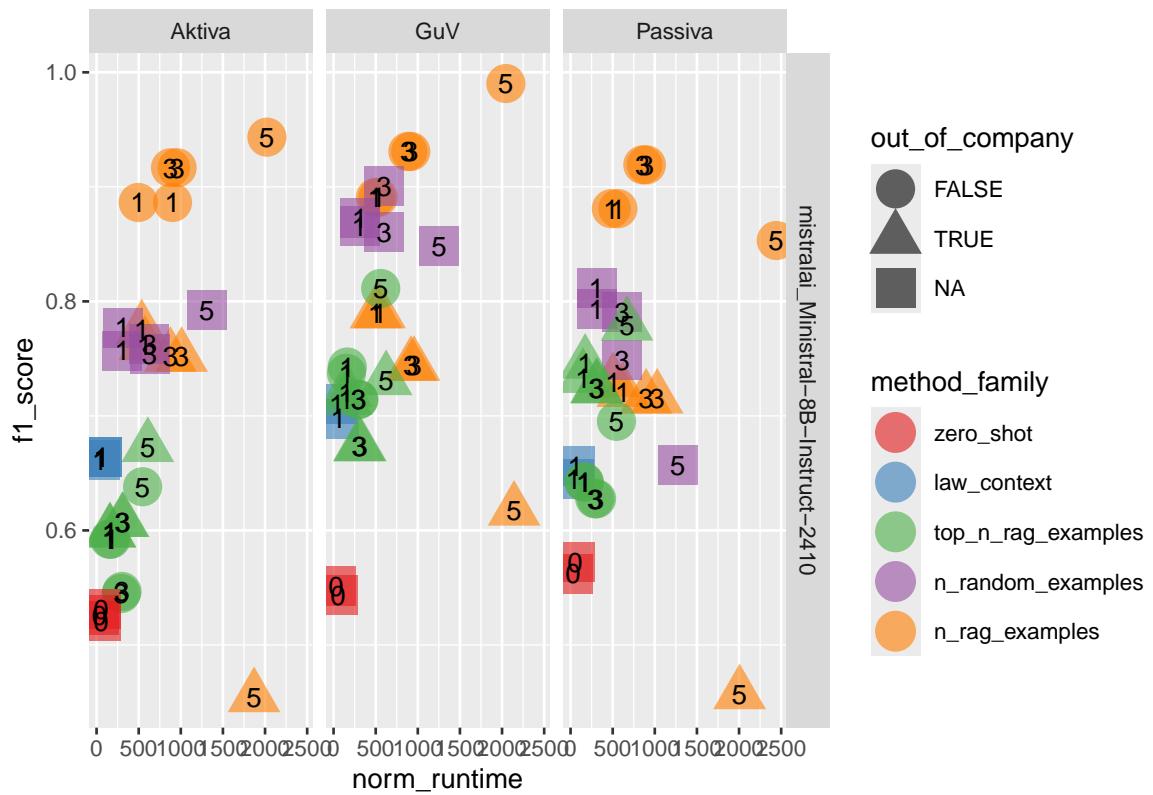
It is interesting that the predictions do not get better by providing more and more examples. Especially for the n-rag-example approach we find a significant drop in the F1 score if the examples pages come from different companies annual reports. This is caused by a severe recall drop. But also for the n-random-example approach we see this for the prediction of class Passiva.

Recall better with examples from same company. Precision better without.

We can also see that the prediction performance is stable.¹¹

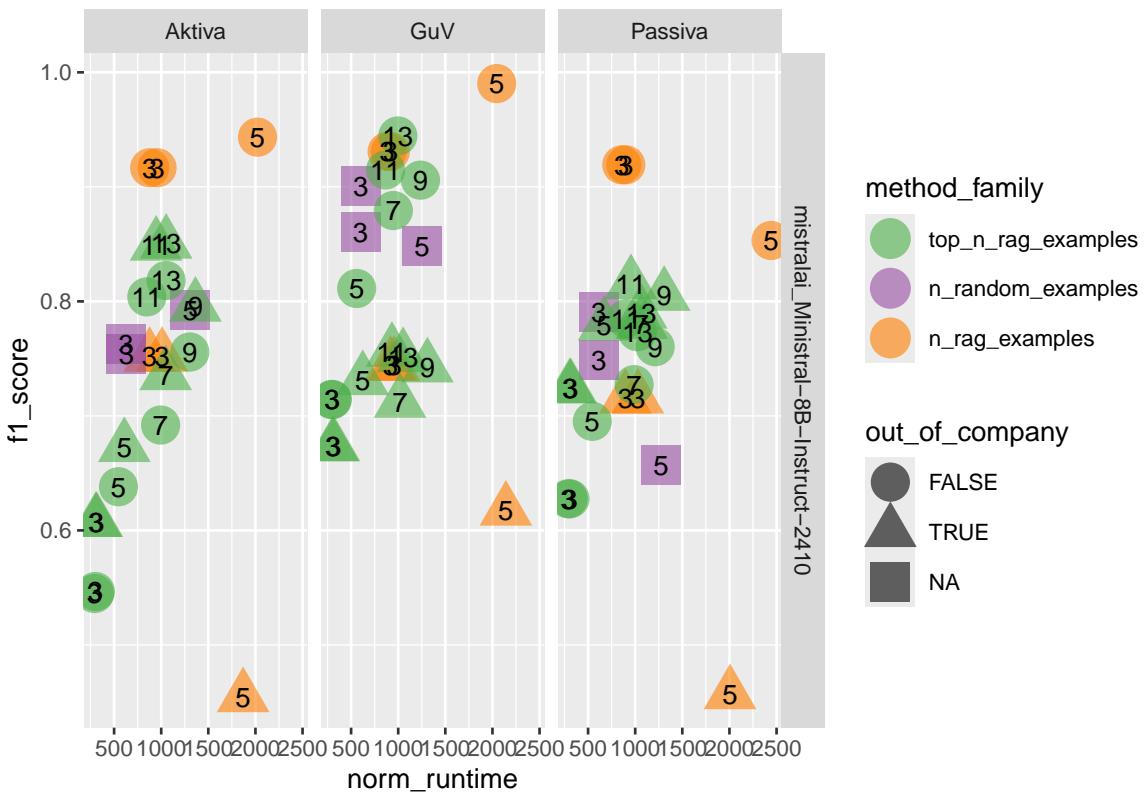
¹⁰This is not due to a temporary technical problems caused by a bug in the transformers version shipped with the vilm 0.9.2 image. Those problems have been overcome. The performance stays bad.

¹¹Earlier experiments on a subset of the pages have been run five times indicating stable results. Running the experiments up to three times in this very task indicate this as well.



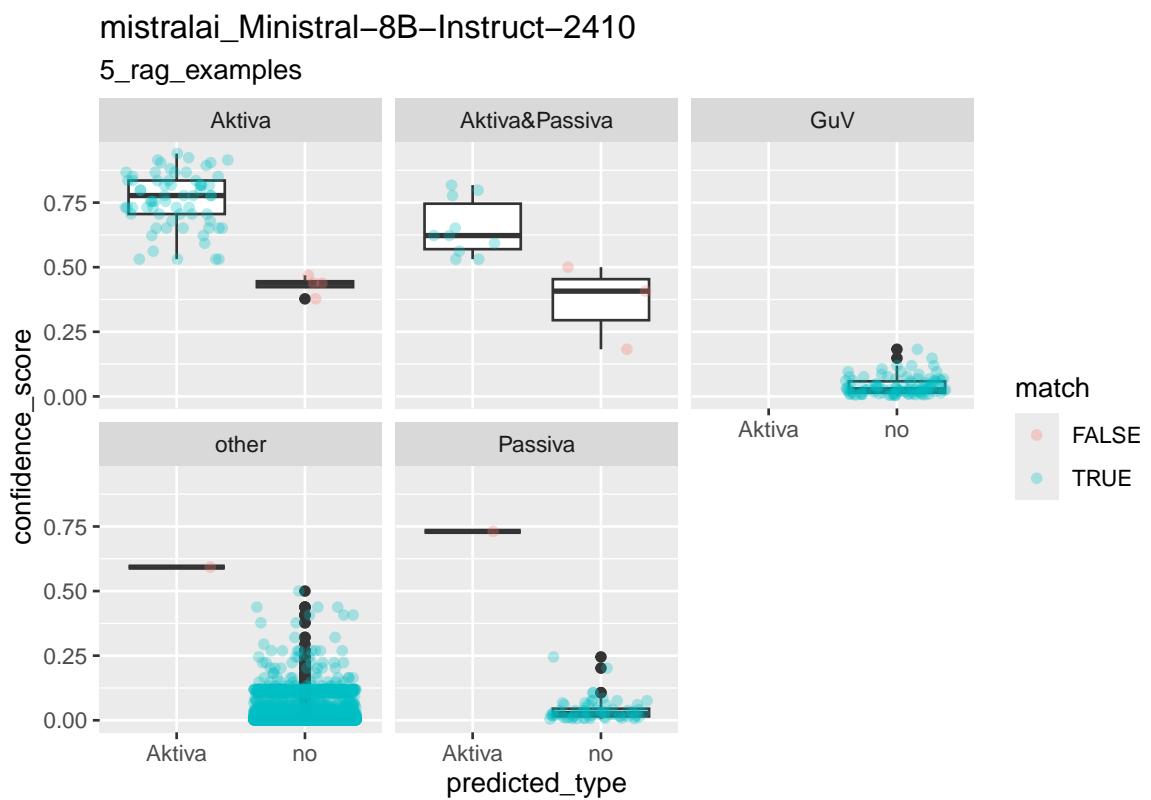
- f1
- multiple models
- best model detail (different methods / settings)

The experiments for the best performing model, Minstral-8B-Instruct-2410, have been extended by methods with even more examples. Especially for the top-n-rag-example approach to get a better comparable picture based on the real number of examples / context length.

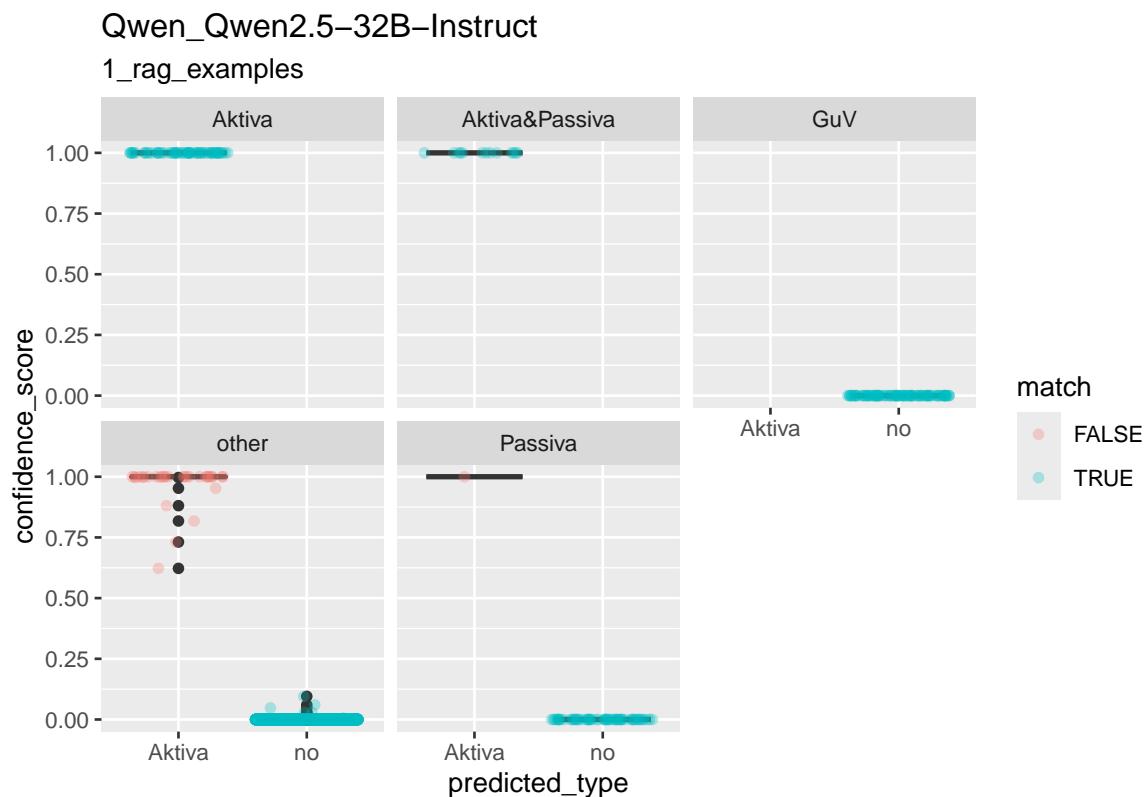


approach	classification	n_example	target	other	sum
n_random_examples	binary	1	1	1	4
n_random_examples	binary	3	3	1	6
n_random_examples	binary	5	5	2	11
n_random_examples	multi	1	1	1	4
n_random_examples	multi	3	3	3	12
n_random_examples	multi	5	5	5	20
n_rag_examples	binary	1	1	1	4
n_rag_examples	binary	3	3	1	6
n_rag_examples	binary	5	5	2	11
n_rag_examples	multi	1	1	1	4
n_rag_examples	multi	3	3	3	12
n_rag_examples	multi	5	5	5	20
top_n_rag_examples	binary	1	NA	NA	1
top_n_rag_examples	binary	3	NA	NA	3
top_n_rag_examples	binary	5	NA	NA	5
top_n_rag_examples	binary	7	NA	NA	7
top_n_rag_examples	binary	9	NA	NA	9
top_n_rag_examples	binary	11	NA	NA	11
top_n_rag_examples	binary	13	NA	NA	13
top_n_rag_examples	multi	1	NA	NA	1
top_n_rag_examples	multi	3	NA	NA	3
top_n_rag_examples	multi	5	NA	NA	5
top_n_rag_examples	multi	7	NA	NA	7
top_n_rag_examples	multi	9	NA	NA	9
top_n_rag_examples	multi	11	NA	NA	11
top_n_rag_examples	multi	13	NA	NA	13

Predictions very accurate. Confidence not always 1. Wrong predictions often with medium confidence. If Aktiva and Passiva on same page more often Aktiva predicted. Confidence for no displayed as 1-confidence to represent confidence for yes (binary classification).

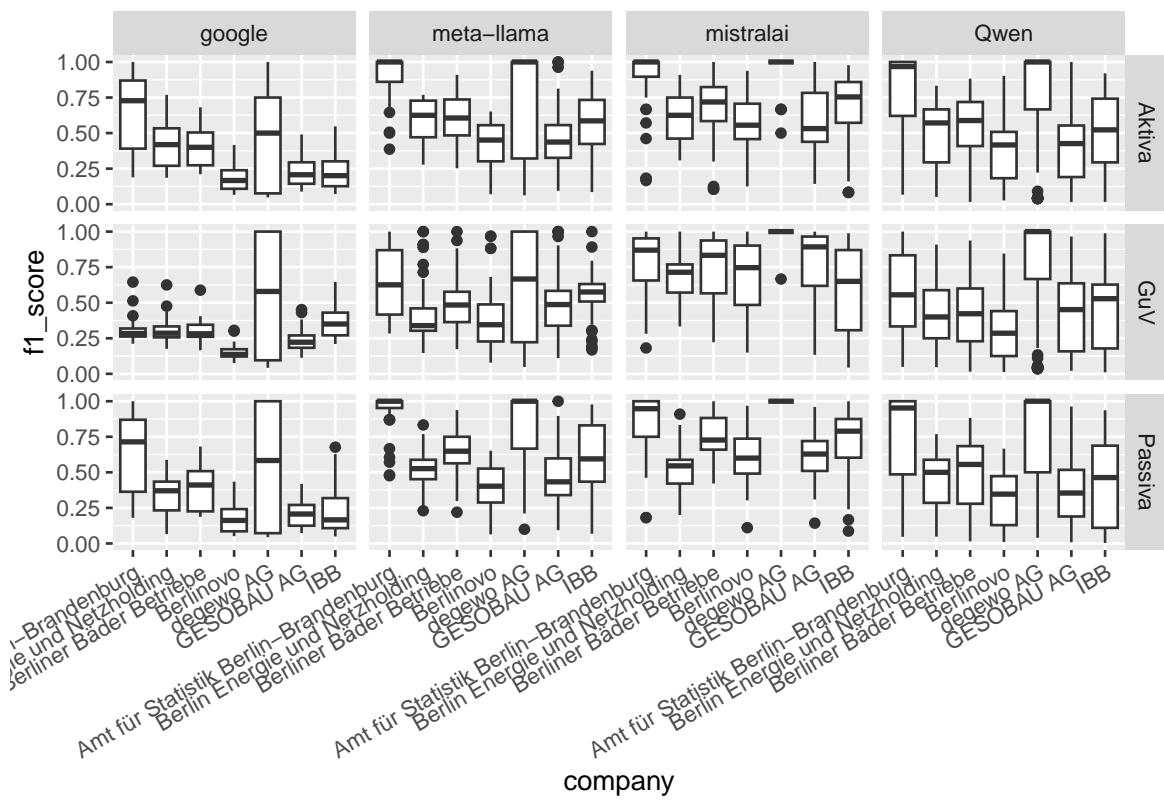


Qwen returns always high confidence even if it is wrong.



- IBB other law
- degewo only one where no ocr is needed

mistral: recall IBB and Netzholding big range meta & mistral: very high precision for Amt für Statistik BBB <- lowest average pagecount (29.3) but IBB has more pages than berlinovo but better precision. No information about which company / report the page is from

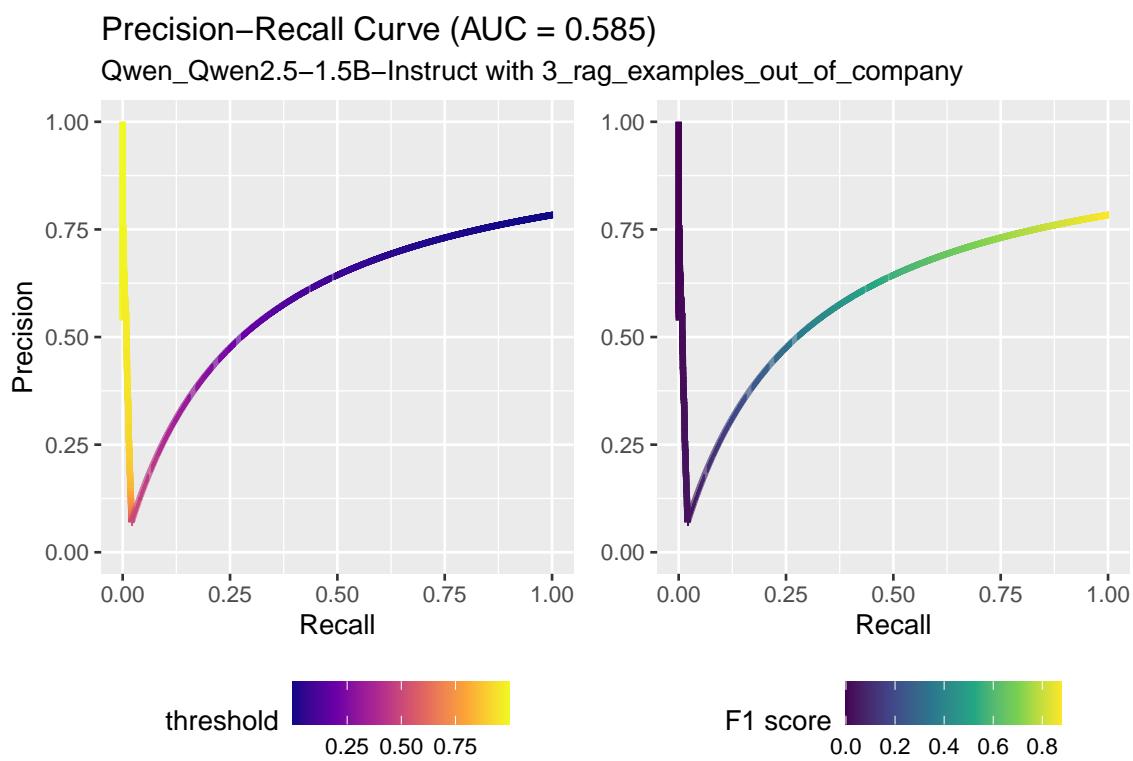
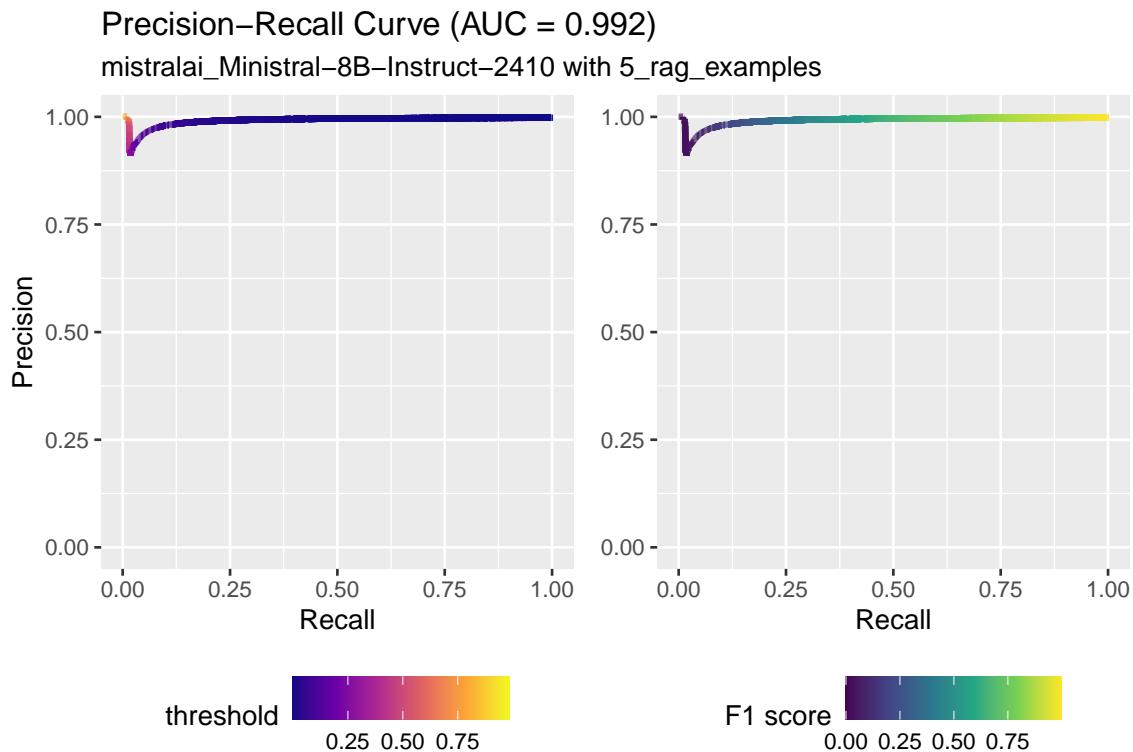


company	n
Amt für Statistik Berlin-Brandenburg	10
Berlin Energie und Netzholding	3
Berliner Bäder Betriebe	10
Berlinovo	15
GESOBAU AG	13
IBB	22
degewo AG	1

company	n
Amt für Statistik Berlin-Brandenburg	10
Berlin Energie und Netzholding	3
Berliner Bäder Betriebe	10
Berlinovo	15
GESOBAU AG	13
IBB	22
degewo AG	7

- Performance makes a jump at a critical parameter number (3B) then slow increase (compare Qwen 2.5)
- Changes unsystematic with new models (see Mistral, Qwen 3 old vs llama 4)

PR curves for all classes look very alike- showing micro average curve



model_family	model	metric_type	method_family	n_examples	f1_score	run
meta-llama	metallama_Llama4Scout17B16EInstruct	GuV	n_rag_examples	1	1	254
meta-llama	metallama_Llama4Scout17B16EInstruct	Aktiva	n_rag_examples	3	0.99	445
mistralai	mistralai_MistralLargeInstruct2411	Passiva	n_rag_examples	1	0.99	706
meta-llama	metallama_Llama4Scout17B16EInstruct	Passiva	n_rag_examples	3	0.99	445
mistralai	mistralai_MistralLargeInstruct2411	Aktiva	n_rag_examples	3	0.97	187
Qwen	Qwen_Qwen2.532BInstruct	Aktiva	n_rag_examples	3	0.97	566
Qwen	Qwen_Qwen3235BA22BInstruct2507	GuV	n_rag_examples	3	0.97	111
mistralai	mistralai_MistralLargeInstruct2411	GuV	n_rag_examples	1	0.95	706
mistralai	mistralai_MistralLargeInstruct2411	GuV	n_rag_examples	3	0.95	187
Qwen	Qwen_Qwen2.572BInstruct	Passiva	n_rag_examples	1	0.95	390
google	google_gemma327bit091	Aktiva	n_rag_examples	3	0.88	429
google	google_gemma327bit091	Passiva	n_rag_examples	1	0.81	260
google	google_gemma327bit091	GuV	n_rag_examples	1	0.78	260
tiuae	tiuae_Falcon310BInstruct	GuV	n_rag_examples	1	0.7	868
tiuae	tiuae_Falcon310BInstruct	Aktiva	n_rag_examples	3	0.69	239
microsoft	microsoft_phi4	Passiva	n_rag_examples	2	0.67	166
microsoft	microsoft_phi4	Aktiva	n_random_examples	1	0.59	493
tiuae	tiuae_Falcon310BInstruct	Passiva	top_n_rag_examples	3	0.59	494
microsoft	microsoft_phi4	GuV	n_rag_examples	1	0.45	172

5.1.3.2 Multi classification

bigger models are better with the multi classification task Llama-4-Scout almost perfect F1 for all classes

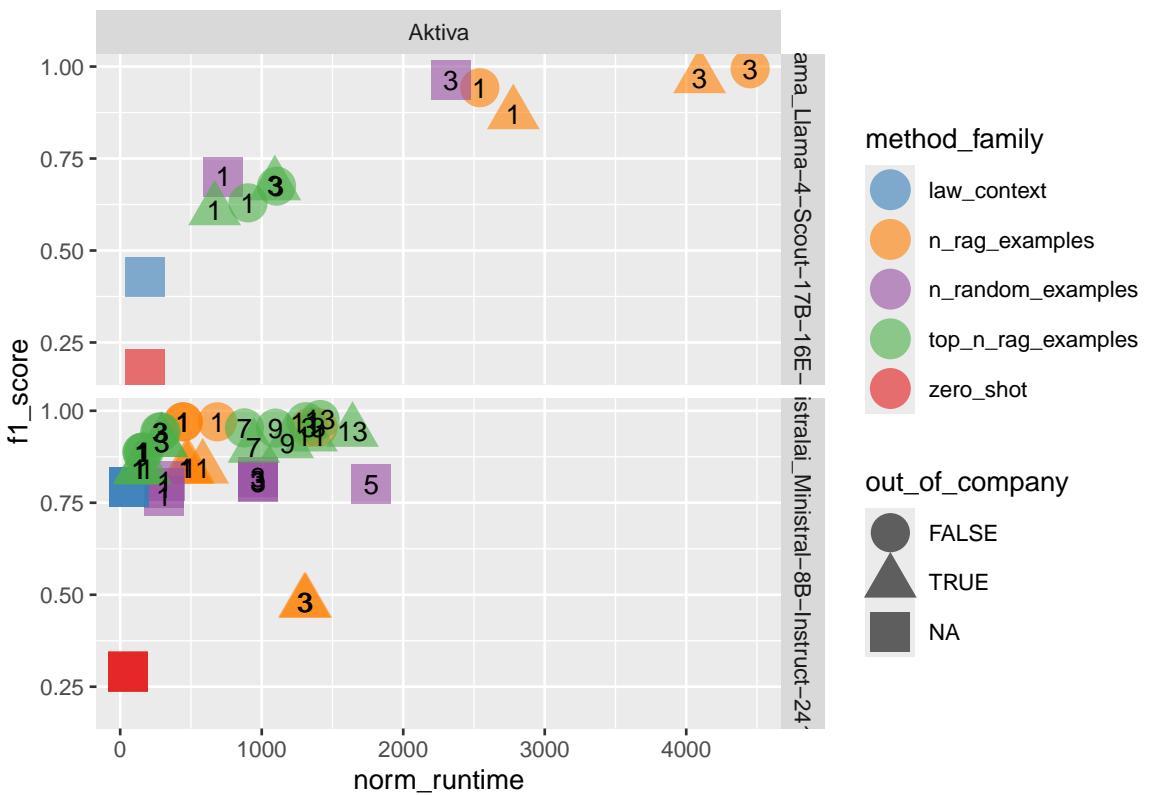
Llama-4-Scout runs fast but needs long to load because it has 109B in total with 17B actives Gemma performs much better than with binary classification

drop with Qwen-14B

Mistral-8B-2410 almost as good as Mistral-123B-2411 but much faster

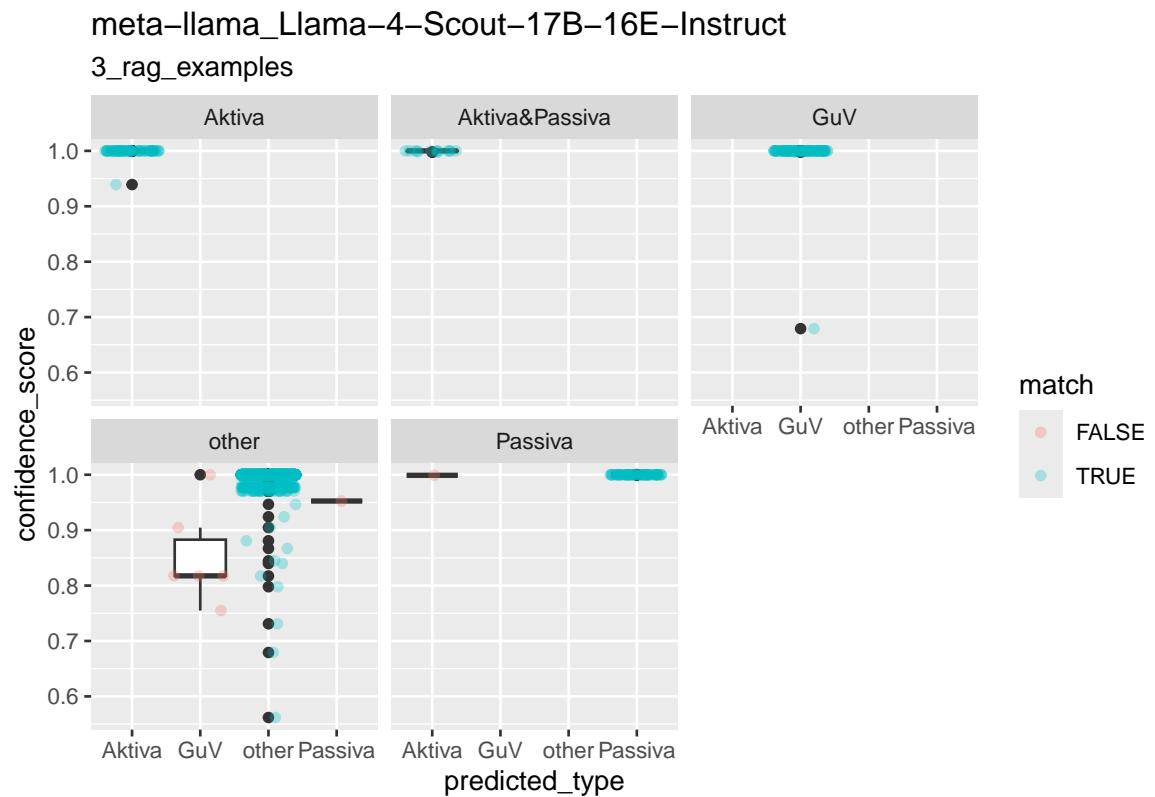
Mistral-2410 reaches good performance already with few examples and can work with law-context approach but more examples don't really help any further

model_family	model	metric_type	method_family	n_examples	f1_score	runti
mistralai	mistralai_Minstral8BInstruct2410	Aktiva	n_rag_examples	1	0.97	686
mistralai	mistralai_Minstral8BInstruct2410	Passiva	n_rag_examples	1	0.95	686
mistralai	mistralai_Minstral8BInstruct2410	GuV	n_rag_examples	3	0.95	1399
mistralai	mistralai_Minstral8BInstruct2410	GuV	top_n_rag_examples	3	0.95	279
meta-llama	metallama_Llama3.18BInstruct	Passiva	n_rag_examples	1	0.94	593
Qwen	Qwen_Qwen2.53BInstruct	Aktiva	n_rag_examples	1	0.86	492
meta-llama	metallama_Llama3.18BInstruct	Aktiva	top_n_rag_examples	3	0.85	269
google	google_gemma312bit091	Aktiva	n_rag_examples	3	0.84	2733
Qwen	Qwen_Qwen2.53BInstruct	Passiva	law_context	NA	0.81	28
Qwen	Qwen_Qwen2.53BInstruct	GuV	n_rag_examples	1	0.76	492
tiuae	tiuae_Falcon310BInstruct	GuV	n_rag_examples	1	0.7	868
tiuae	tiuae_Falcon310BInstruct	Aktiva	n_rag_examples	3	0.69	2393
google	google_gemma312bit091	Passiva	n_rag_examples	1	0.68	1259
meta-llama	metallama_Llama3.18BInstruct	GuV	n_rag_examples	1	0.62	593
meta-llama	metallama_Llama3.18BInstruct	GuV	top_n_rag_examples	1	0.62	205
tiuae	tiuae_Falcon310BInstruct	Passiva	top_n_rag_examples	3	0.59	494
google	google_gemma312bit091	GuV	top_n_rag_examples	1	0.46	232



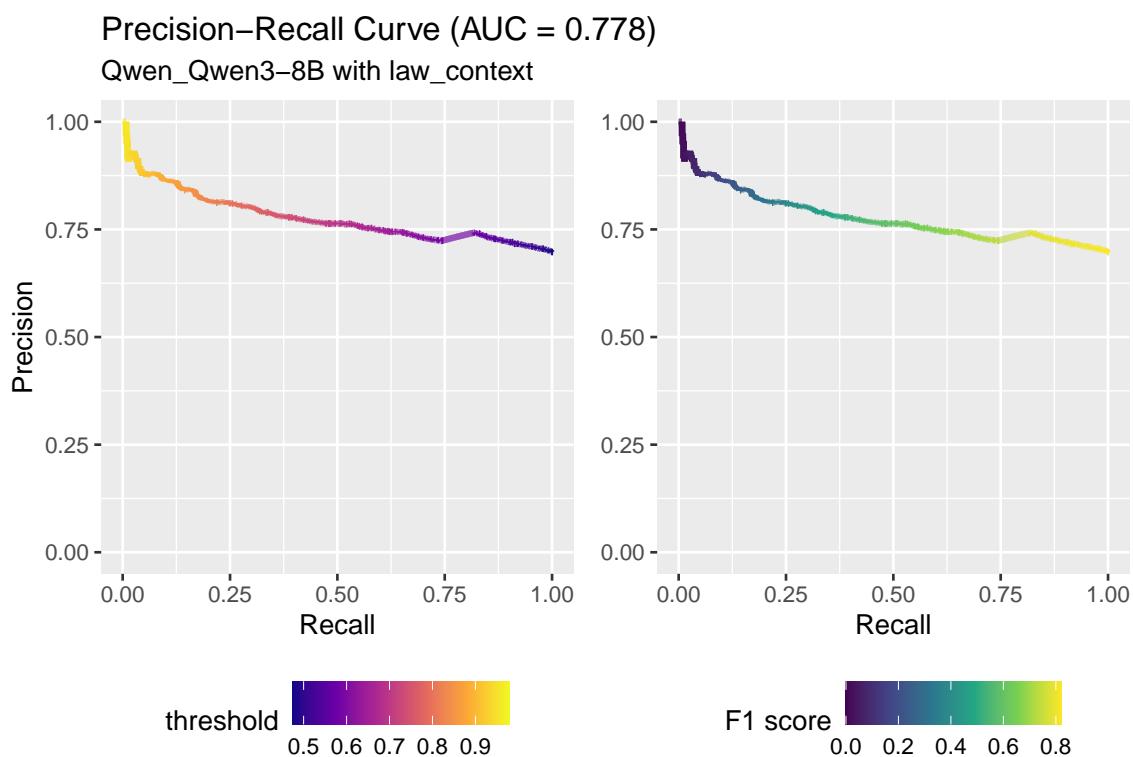
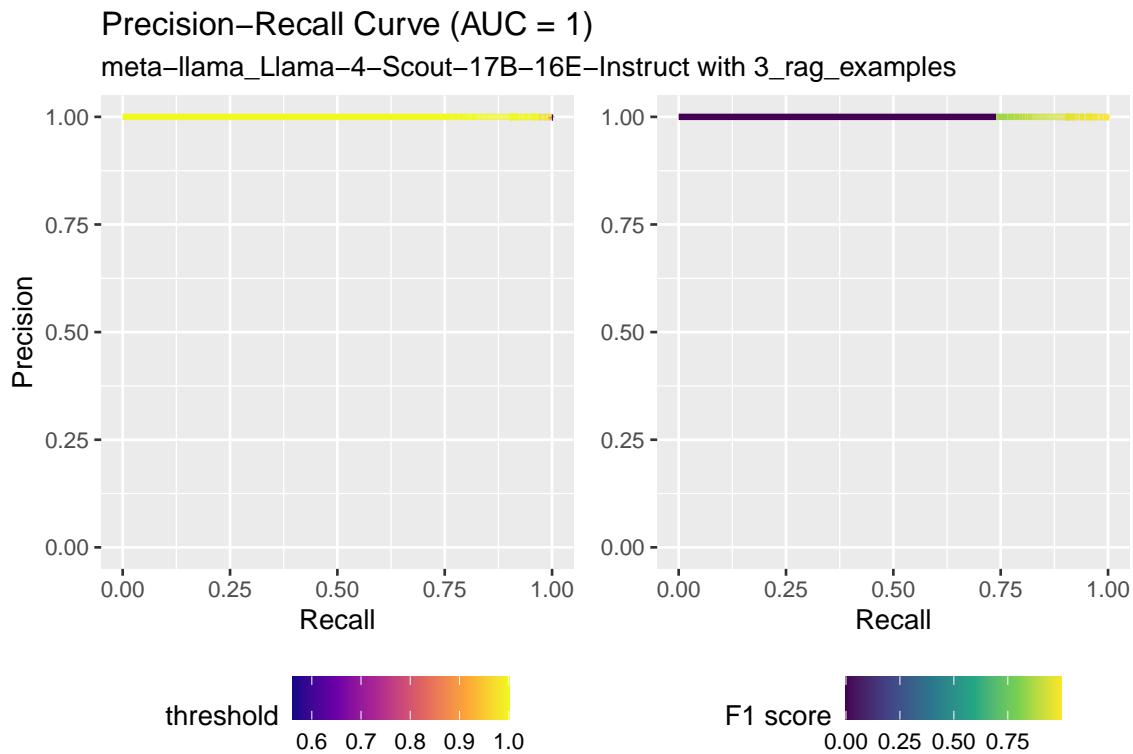
Most of the time pretty confident most problems with class “other” If Aktiva and Passiva on same page it predicts Aktiva. Also one Passiva missclassified as Aktiva No flipped confidence ¹²

¹²classify framework in needs special models with pooling capability. Would have been interesting but time was limited and would have needed new special models in most cases



Microsoft phi 4 and Falcon 3 only ran with one and two examples because their context window is smaller.

- f1
- multiple models
- best model detail (different methods / settings)



5.1.4 Term frequency based classifier

RandomForest performs much better than a logistic regression Better results with * undersampling * training on all types simultaneously

5.1.4.1 Two predictors

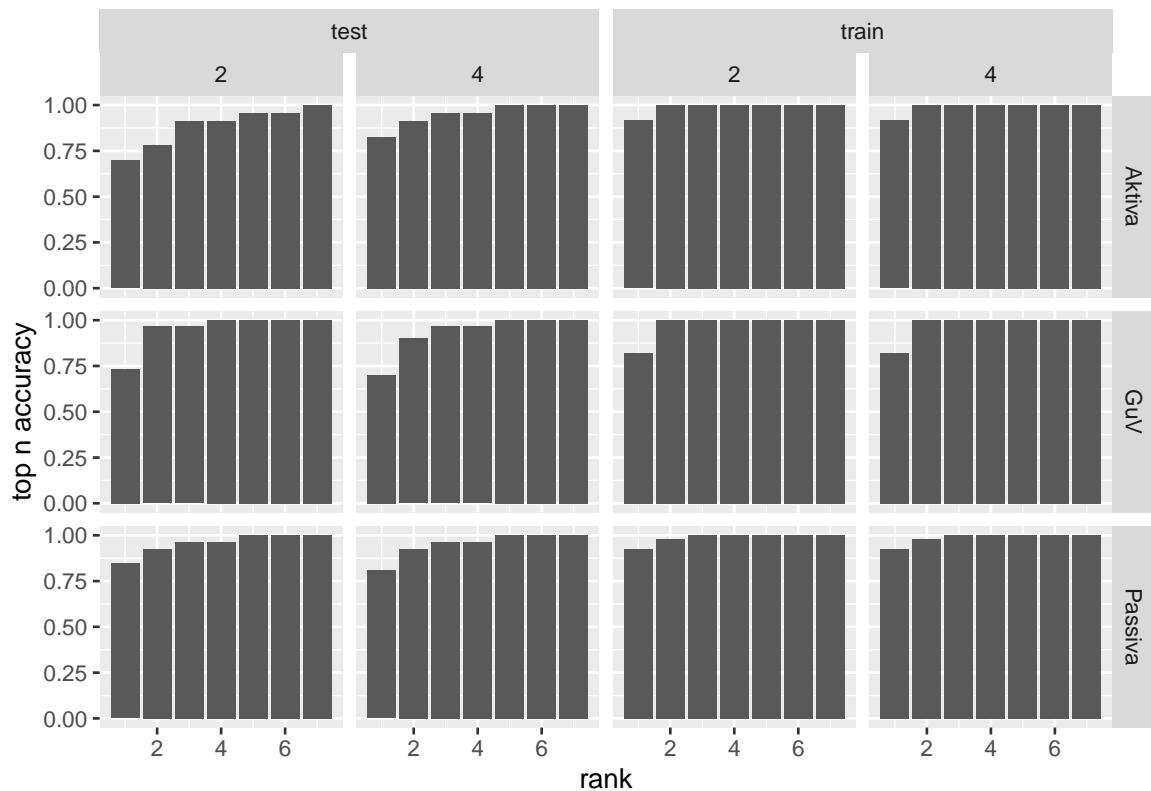
Term frequency of nouns of the law about Aktiva Float frequency (floats divided by word count)

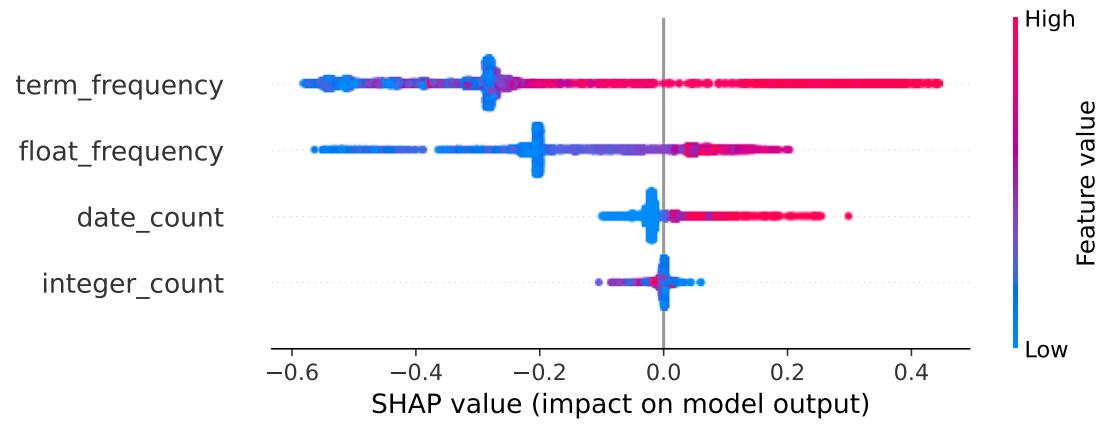
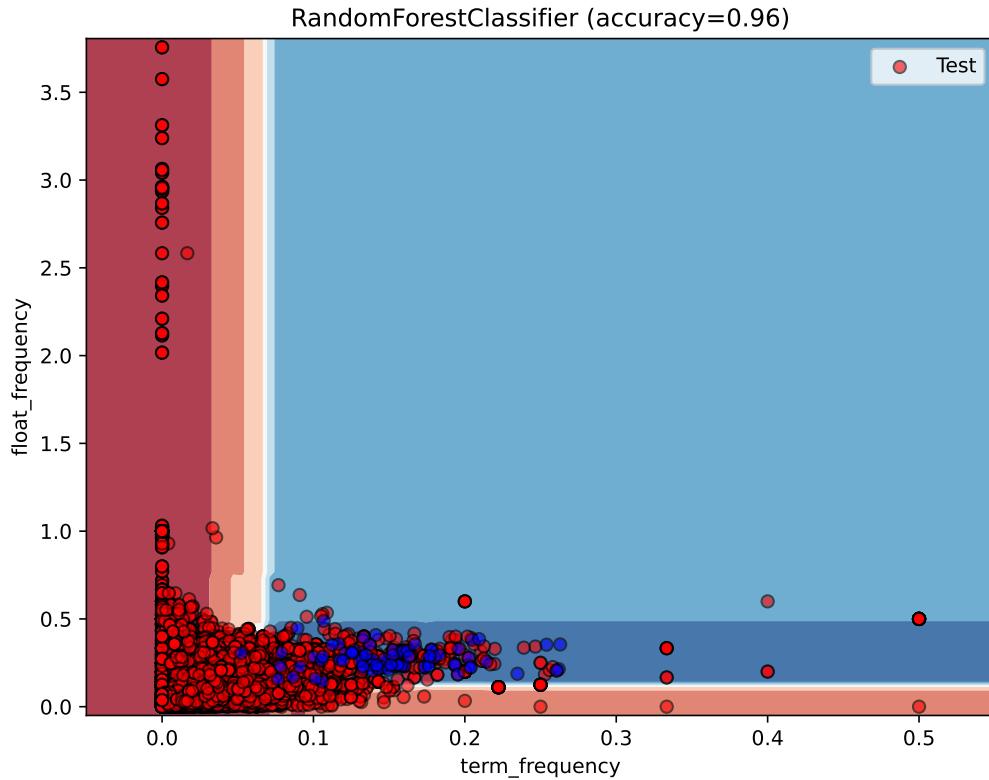
5.1.4.2 Four predictors

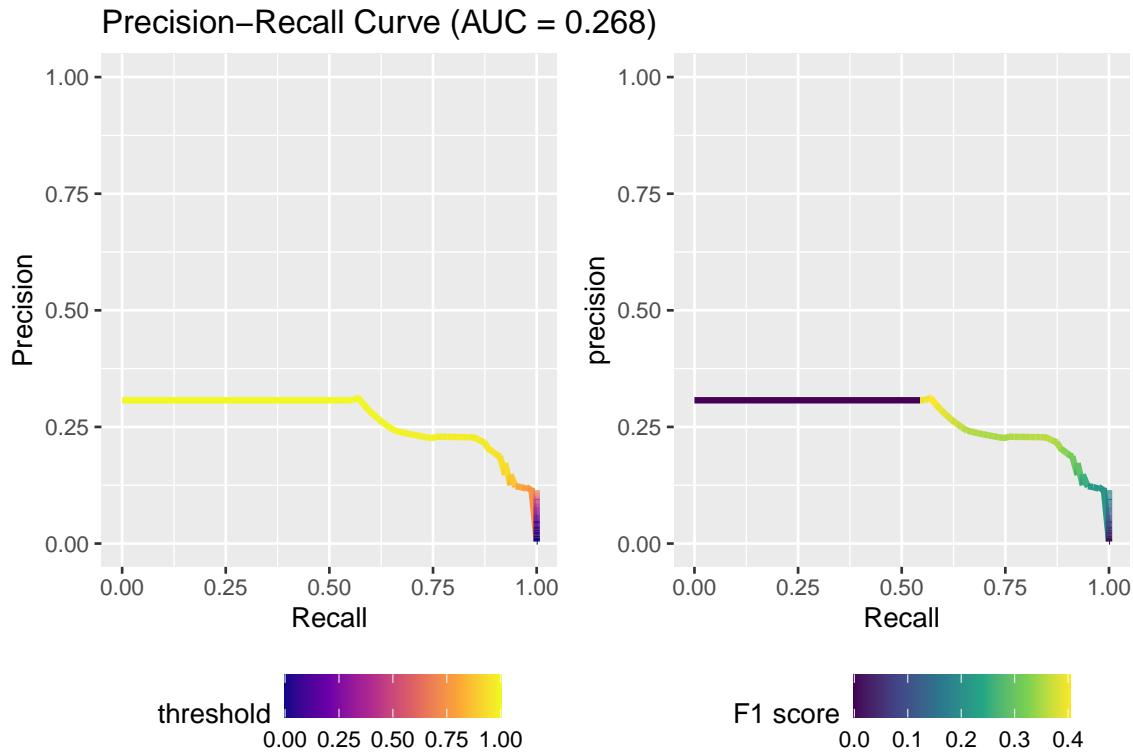
Count of integers Count of dates

- top 1
- top k

low precision l1m linked to position of correct page? numeric frequency?







5.1.5 Comparison

5.1.5.1 Prediction performance

F1 scores for llms are much higher

5.1.5.2 Energy usage and runtime

Multiclassification more effective than three times single classification Combine term frequency with llm approach to limit page range

5.2 Table extraction

The second task to solve is: extract the data from the document.

Which tasks have there been? Which models have been used for which ttask? What data has been used?

5.2.1 Baseline: Regex

The baseline for the table extraction task is set by an approach using regular expressions on the text extract. The approach performs much better¹³ on the synthetic dataset compared to the real dataset (see Figure 5.9). Even though, it does not perform perfectly and its performance is more consistent on the text extracted with pymupdf compared to pdfium. Some possible explanations are:

¹³A comparison of the numeric values over all methods can be found in section 5.2.3.

- a duplicated row name¹⁴
- numeric columns extracted separated from row names by extraction libraries
- sums in the same row as the single values¹⁵
- with pdfium: missing white space¹⁶
- with pdfium: random line breaks¹⁷

You can find some examples for incorrect extracted texts in section A.8.

On the real dataset the approach shows a wider spread for the percentage of correct extracted numeric values as well as a considerable number of annual reports where the extraction did not work at all. Interestingly, the used text extraction library has no noticeable influence on the real dataset.

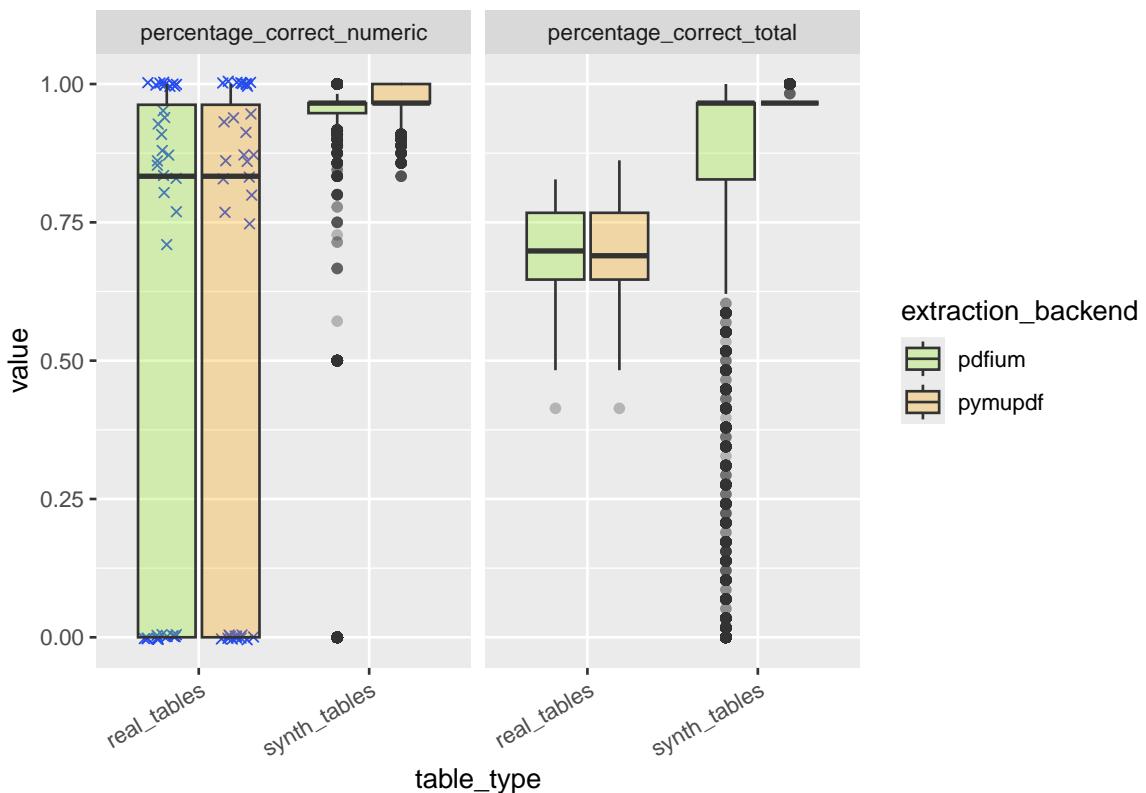


Figure 5.9: Performance overall and on numeric value extraction with regular expressions. Showing single scores for *percentage correct numeric* on real tables to explain wide boxes.

The random line breaks result in some missed row names which is reflected by the bigger spread for NA precision with pdfium on the synthetic dataset (see Figure 5.10). Nevertheless, the NA precision for the majority of the cases is perfect. This is different with the real dataset. The NA precision is found to be at only 0.7.

Hypotheses The formulated hypotheses have been evaluated visually using the dependence and beeswarm plots from the shapviz library based on the SHAP values calculated with a random forest.

¹⁴The row *Geleistete Anzahlungen* can be found in two parts of the table and the simple approach just matches the numbers to the first found entry.

¹⁵In this case the regex takes the sum as the value for the previous year.

¹⁶This can form unexpected numeric patterns or prevent the row names to be recognized.

¹⁷The approach takes care of line breaks between words, but not within. This leads to unrecognized row names as well.

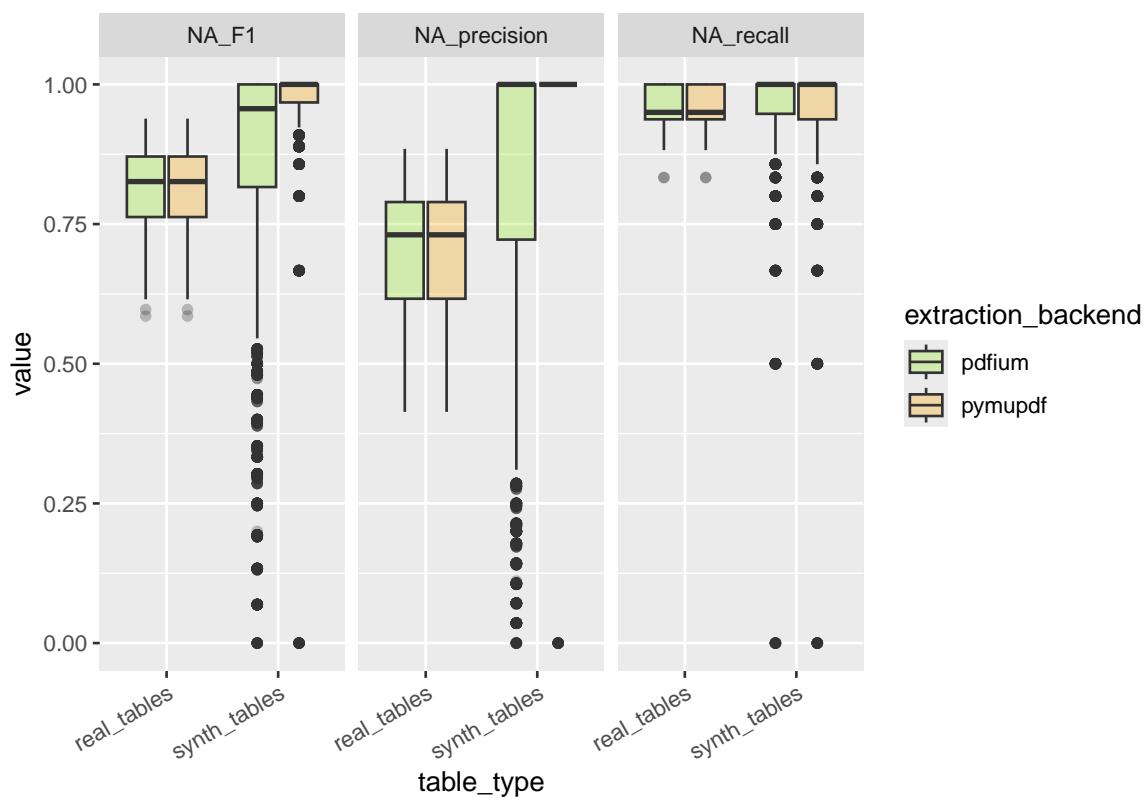


Figure 5.10: Performance on classification for missing values with regular expressions

Real dataset There are multiple hypotheses that don't get supported by the visual results (see Figure A.2). The pretty surprising results are:

1. The visual separation of columns or rows has an effect on the text processing.
2. It seems to have a positive effect on F1 and numeric correctness rate if the Passiva table is on the same page, even though it has no influence on the single predictions.

But one has to keep in mind that the number of data points on the aggregated values for the test set of the real dataset is only 18. So these findings are not strongly supporting any interpretation at all. Furthermore, the found effects are not very large - most below 5 %. Only the hypothesis for a positive influence of a missing of a value for the binomial prediction gets solid support with a mean absolute SHAP value of over 20 %. To get reliable results more tables have to be included which would require additional manual encoding.

Synthetic dataset Interpreting the visual results for the SHAP analysis on the synthetic dataset brought some interesting insides into the question under which condition the two PDF extraction libraries perform differently. These results can be treated as reliable since the model has been trained with 50_000 rows and the SHAP values have been calculated on 2_000 rows each.

Very interestingly the number of columns is having an opposite effect for the two libraries (see Figure 5.11 A). Besides that often only pdflium struggled with some of the table characteristics while pymupdf is not influenced by them (for an example with header_span see Figure 5.11 B).

It might be worth noting that the row for *Anteile an verbundenen Unternehmen* was rated to have a clear negative effect on the chance to extract the correct value.

Since there has no synthetic data created where also the Passiva table is present the result found with the real dataset can't be investigated further. Also the question if visual separation is having an effect was not studied, even though, creating such additional synthetic data would be very easy with the current generation process and could be done in future work. It would be interesting if the visual separation is cause for the maleous text extractions of pdflium as well.

X1	X2	X3	X4	X5	X6	X7
predictor	F1	F1	% numeric correct	% numeric correct	binomial	binomial
predictor	Hypothesis	Result	Hypothesis	Result	Hypothesis	Result
extraction_backend	neutral	pymupdf better	neutral	pymupdf better	neutral	pymupdf b
n_columns	4 is worse	positive	neutral	positive	neutral	positive
sum_same_line	neutral	neutral	negative	negative*	negative	neutral
header_span	neutral	negative*	neutral	negative*	neutral	negative*
thin	negative	NA	neutral	positive*	neutral	neutral
year_as	neutral	positive*	neutral	positive*	neutral	positive*
unit_in_first_cell	negative	negative*	negative	negative*	negative	negative*
log10_unit_multiplier	neutral	negative*	positive	negative*	positive	negative*
enumeration	positive	positive*	neutral	positive*	neutral	positive*
shuffle_rows	negative	neutral	neutral	neutral	neutral	neutral
text_around	neutral	neutral	neutral	neutral	neutral	neutral
many_line_breaks	negative	neutral	neutral	neutral	neutral	neutral
label_length						
label					unknown	
missing					positive	positive

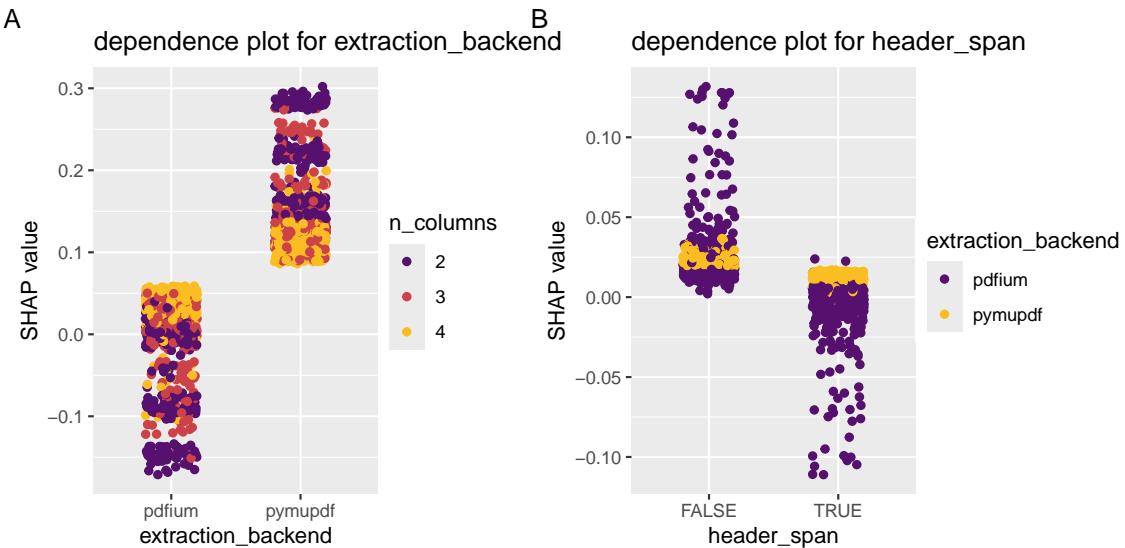


Figure 5.11: Showing the influence of the extraction library on the numeric text extraction task with synthetic data

5.2.2 Extraction with LLMs

- confidence usable to head for user checks?
- not handled new entries
- five examples bring not much more, but a little
- random forest / SHAP

5.2.2.1 Real tables only

For the table extraction task 31 open source models have been benchmarked¹⁸. The results are presented in Figure A.9, A.10 and A.11).

Most models need a context learning approach to beat the performance of the regular expression approach at total and numeric correctness rate and F1 score. Only 5 models perform better without any guidance¹⁹ (see Table 5.9). 10 models achieved an performance better as the regex baseline using the approach to learn with a fixed example from the synthetic dataset.

In contrast: most of the models achieved a better performance than the regex baseline when they were provided with one or more examples from real *Aktiva* tables. Just 4 don't achieve a better value even with three or five realistic examples (see Table 5.10). Here we find the smallest models with less than 2B parameters which don't achieve a consistence performance no matter how many examples they get. But we also find models that start to perform bad if they get a too long context with too many examples like the very recent and large model Llama 4 Maverick.

With one and three examples the performance within one model family is positive correlated with the number of parameters the models have. Once the 4B parameters are passed the improvements get less and less getting closer to a perfect performance but never reaching it on all documents. Table 5.11 shows the mean performance for the best model-method approach for each model family. Most of the top performing model-method combinations rely on the maximum number of examples provided. Only the Llama-3 and Falcon3

¹⁸The models *deepseek-ai_DeepSeek-R1-Distill-Qwen-32B* and *google_gemma-3n-E4B-it* have been tested as well but don't get presented as they never performed anywhere beyond random guessing.

¹⁹There is an external guidance through the provided xgrammar template but it is not communicated to the model in a prompt.

Table 5.9: Comparing table extraction performance with real 'Aktiva' dataset for models that perform well without or with little context learning

model	mean_total_zero_shot	mean_total_static_example
Llama4Maverick17B128EInstructFP8	0.816	0.844
Qwen2.532BInstruct	0.76	0.875
Qwen3235BA22BInstruct2507	{0.848}	{0.88}
Qwen3235BA22BInstruct2507FP8	0.825	0.873
phi4	0.807	0.75
Llama3.170BInstruct	NA	0.773
MistralSmall3.124BInstruct2503	NA	0.855
Qwen2.572BInstruct	NA	0.838
Qwen330BA3BInstruct2507	NA	0.812
gemma327bit	NA	0.785

Table 5.10: Comparing table extraction performance with real 'Aktiva' dataset for models that worse than the regex baselin with 3 or 5 examples for incontext learning

model	method	mean_total
Llama4Maverick17B128EInstructFP8	5_random_examples	0.041
Qwen2.50.5BInstruct	3_random_examples	0.585
Qwen30.6B	3_random_examples	0.608
gemma34bit	top_5_rag_examples_out_of_sample	{0.682}

model perform best with three examples²⁰.

Based on a small sample of 8 documents by the *Amt für Statistik Berlin-Brandenburg* it seems that there is support for the hypothesis, that providing Aktiva tables from the same company in in-context learning, is improving the results. This is especially noticeable for models with very few parameters and when providing only a single example. This seems intuitive, since there the potential for possibilities is much bigger. Figure A.12 shows that on this limited sample

- the improvement is bigger for Qwen 3 than for Qwen 2.5
- Googles gemma 27b and GPT 4.1 mini could overcome an unnoticed issue with the extraction with just one example.
- the effect of being overwhelmed by a too rich context with LLamas Maverick model could get reduced a bit.

²⁰Phi4 also perfroms best with three examples. But this is the maximum it can process due to a limited context length.

Table 5.11: Comparing best mean table extraction performance with real 'Aktiva' dataset for each model family

model_family	model	method_family	n_examples	mean_total
Qwen 3	Qwen3235BA22BInstruct2507FP8	top_n_rag_examples	5	0.961
Llama-4	Llama4Scout17B16EInstruct	top_n_rag_examples	5	0.931
mistrallai	MistralLargeInstruct2411	top_n_rag_examples	5	0.929
Llama-3	Llama3.170BInstruct	n_random_examples	3	0.911
Qwen 2.5	Qwen2.514BInstruct	n_random_examples	5	0.908
microsoft	phi4	n_random_examples	3	0.893
tiuae	Falcon310BInstruct	top_n_rag_examples	3	0.862
google	gemma327bit	n_random_examples	5	0.821

Table 5.12: Comparing best mean table extraction performance with real 'Aktiva' dataset for each model family for models with less than 17B parameters

model_family	model	method_family	n_examples	mean_total
Qwen 3	Qwen314B	n_random_examples	3	0.912
Qwen 2.5	Qwen2.514BInstruct	n_random_examples	5	0.908
mistralai	Minstral8BInstruct2410	n_random_examples	5	0.896
microsoft	phi4	n_random_examples	3	0.893
tiuae	Falcon310BInstruct	top_n_rag_examples	3	0.862
Llama-3	Llama3.18BInstruct	n_random_examples	5	0.849
google	gemma312bit	n_random_examples	5	0.797

To examine the question, if the reported confidence score of the responses can be used, to flag the predicted values as potentially wrong. Again, Figure 5.12 shows, that Qwen 3 reports very high confidence values no matter if the results are correct or not. With the Mistral model we find a wider range of confidences given and for wrong results lower confidence is reported.

Figure 5.13 shows, that the chance to make an mistake by believing the prediction is rising with lower confidence. The chance to make a mistake is higher for predictions of numeric values than for believing a value is not present in the table. The chance to make such a mistake is higher using the confidence reported by Qwen 3.

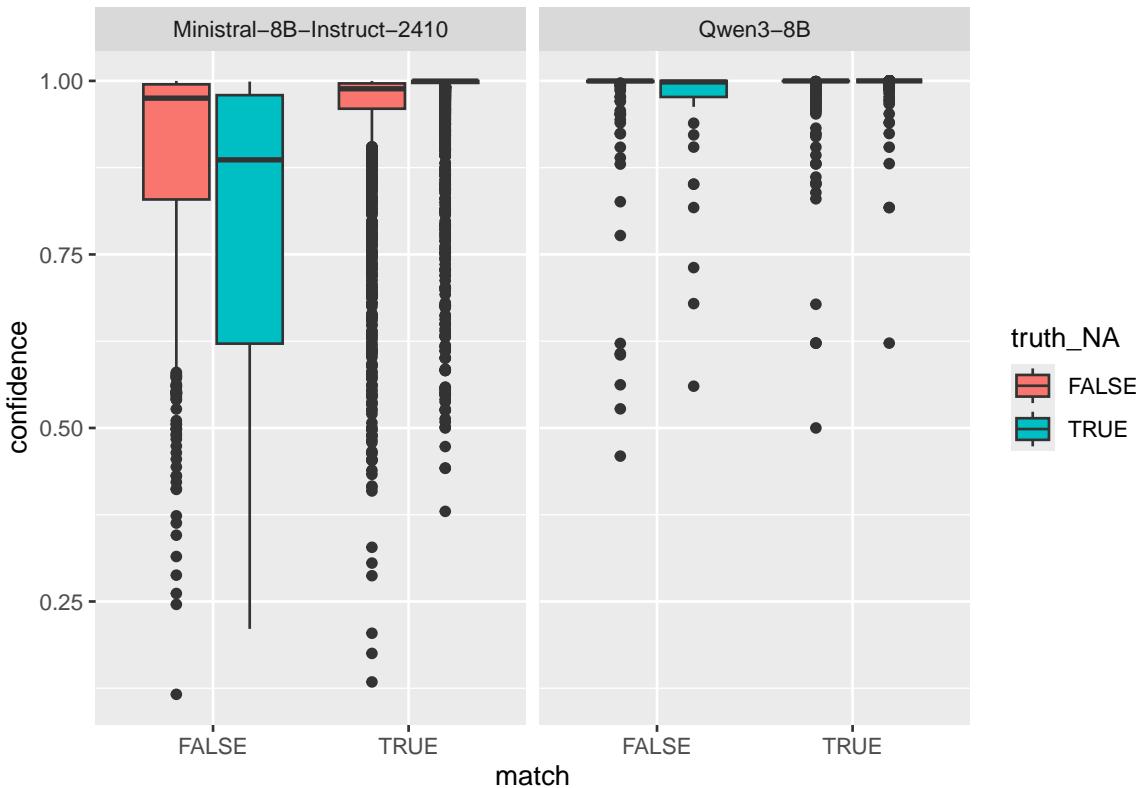


Figure 5.12: Comparing the reported confidence scores for the table extraction task on real dataset for the Mistral and Qwen 3 with 8B parameters.

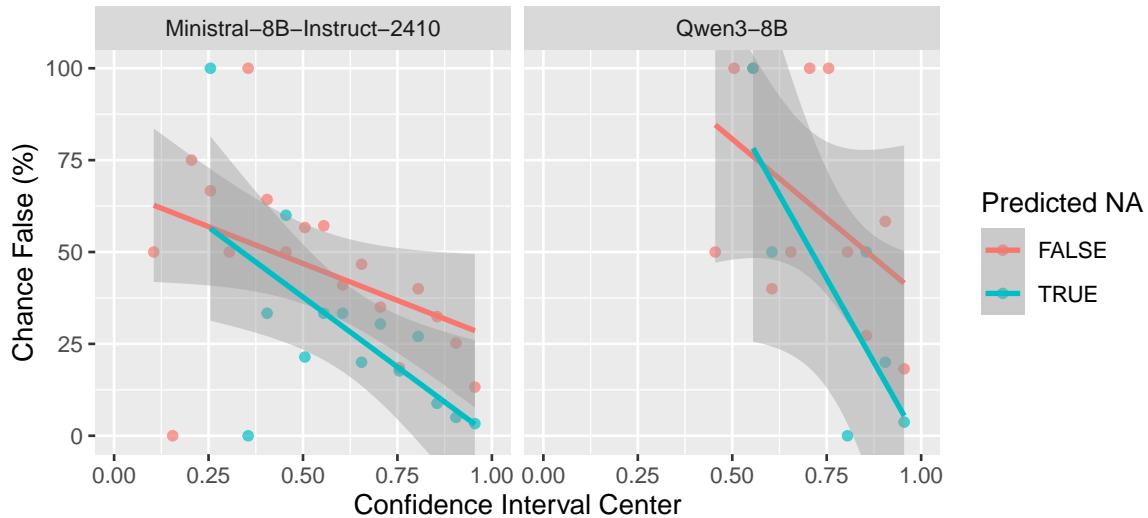


Figure 5.13: Estimating the relative frequency to find a wrong extraction result over different confidence intervals

Hypotheses The formulated hypotheses have been evaluated visually using the dependence and beeswarm plots from the shapviz library based on the SHAP values calculated with a random forest.

Even though the samples size of Aktiva tables did not increase, the available training, test and SHAP sample size is much larger, because the experiment has been repeated with different models and methods. Thus, the interpretations based on the visual evaluation (see Figure A.4)) are more reliable for model and method specific predictors. Since there is one Aktiva example for every company files were found for they might even be generalizable for this population. But one has to keep in mind that there have been more Aktiva tables for *Amt Stat BBB* which might nudge the results a bit.

The results assign much more influence on model and method specific attributes than on the table specific attributes. The importance of the table attributes are as low as found with the regular expresion approach. Only for the binomial prediction we find the predictor *missing* to get assigned more importance than to all model and method specific attributes. Same is true for the *label* that is having the highest influence on the reported confidence. Nevertheless, in the case of the binomial prediction there is half of the predictors *missing* and *label* importance shifted to model and method specific predictors.

Again, multiple hypotheses don't get supported by the visual results. The surprising results are:

1. In general more examples are helpful except for Llamas Maverick model that performs poorly with five examples. But this effect is only noticeable with the aggregated metrics nor for the case wise binomial evaluation.
2. The number of columns has a negative effect on the performance but no effect on the reported confidence.
3. There was no negative effect found if the *Passiva* table is on the same page as the Aktiva table.
4. Larger models start to report less confidence again. This is not unexpected for the Mistral model but was surprising for the largest Qwen 3 model. (Discussion: New Generation? Aktive paramters count? Irrelevant because not well distinguishing?)
5. It not only influences the the performance to extract the correct numeric value from a row where there are additional sums present but also the F1 score.

Two interesting details found while inspecting the dependence plots for the metric *percentage_numeric_correct* are (see Figure 5.14 A) that the bad performance of LLamas Maverick with five examples is easily spottable and that the negative effect of *T_in_year* might be caused by an interaction with *vis_separated_rows* completely

(see Figure 5.14 B). To investigate the second finding one would need tables where the uni is present in the year column and having no visual separation of the rows at the same time. Synthetic data potential could help to answer such questions.

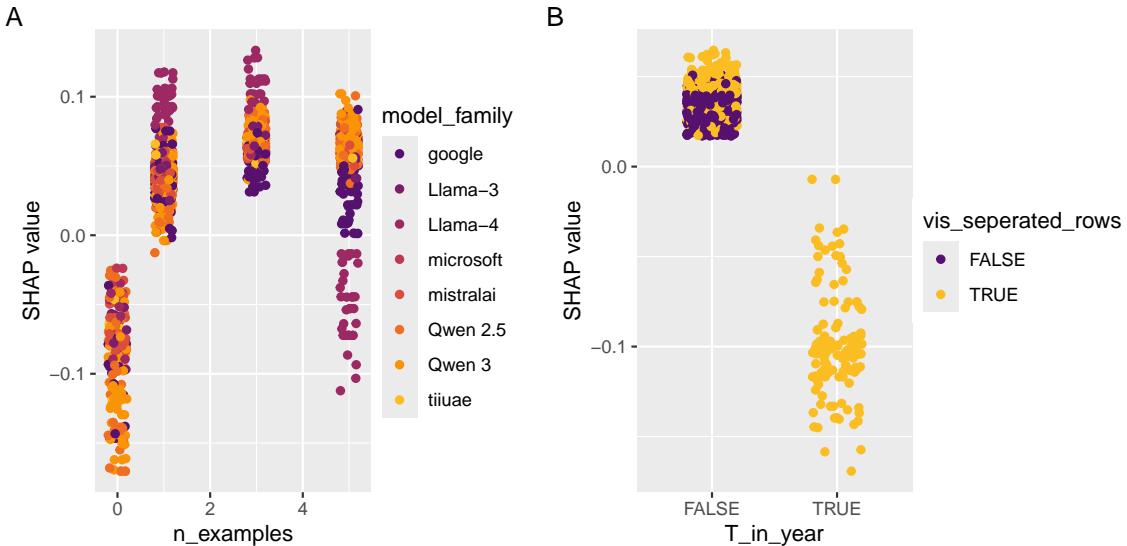


Figure 5.14: Showing the influence of many examples on Llama 4 Maverick (A) and interaction between *T in year* and *vis separated rows* (B)

GPT Even though a lot of documents to process at RHvB (Rechnungshof von Berlin) will not be public and thus must not be processed on public cloud infrastructure, the performance of models like OpenAI's GPT or Google's Gemini are interesting benchmark references within this thesis and for comparing these findings with other papers results. Therefore for this thesis the public available versions of annual reports have been used instead of the ones used internally or for public administration purposes. Those public available reports often are visually more appealing and more heterogeneous in their structure.

As a reference to compare the performance of OpenAI's models with the results of four Qwen 3 models are shown as well (see Figure 5.15, 5.16 and 5.17). Surprisingly gpt-5-mini is almost performing as good as the top Qwen 3 model and gpt-4.1. But besides gpt-4.1-nano and Qwen3-0.6B all models perform pretty well with the random or top-rag example methods. The ranking for the best model-method combination can be found in Table @ref(tab.table-extraction-llm-performance-total-gpt-ranking). Since gpt-4.1 costs five times more than gpt-4.1-mini (see Table @ref{tab:costs-azure}) it seems reasonable to prefer the smaller model for this specific task.

Costs for gpt-5-mini not shown in Azure yet. :(

The author was not able to get OpenAI's models to stick to the provided json (JavaScript Object Notation) schema strictly. Passing the ebnf (extended Backus–Naur form) grammar did not work at all. This means that with gpt-4.1-nano there have been 12 predictions that have been completely empty. Overall there have been 22.4 % of the responses of OpenAI's models that were compatible with the schema but had a wrong number of rows predicted (see Figure @red{fig:table-extraction-llm-prediction-count-gpt}).

Using gpt-5-nano and gpt-5-chat for the table extraction task was not working. With gpt-5-nano the answers were not respecting the provided grammar. Running gpt-5-chat resulted in the error informing that a `json_schema` can't be used with this model. With gpt-5-mini the very approach worked flawless. Running gpt-oss-20b with the vllm offline inference framework was possible and the new harmony output format

Table 5.13: Comparing best mean table extraction performance with synthetic 'Aktiva' dataset for each model family

model_family	model	method_family	n_examples	mean_total
Qwen 2.5	Qwen2.572BInstruct	top_n_rag_examples	3	{0.989}
Qwen 3	Qwen3235BA22BInstruct2507	top_n_rag_examples	3	{0.987}
mistralai	MistralLargeInstruct2411	top_n_rag_examples	5	{0.984}
Llama-4	Llama4Scout17B16EInstruct	top_n_rag_examples	3	{0.98}
google	gemma327bit	n_random_examples	5	0.906
Llama-3	Llama3.18BInstruct	top_n_rag_examples	3	0.897

could be processed after minor code changes for most approaches²¹. With a gpt-oss-120b instance hosted on Azure the guided decoding worked flawless.

model	method	mean correct total
Qwen3-235B-A22B-Instruct-2507	3_random_examples	0.95
gpt-4.1	3_random_examples	0.94
gpt-5-mini	top_3_rag_examples_out_of_sample	0.93
Qwen3-30B-A3B-Instruct-2507	top_3_rag_examples_out_of_sample	0.90
gpt-4.1-mini	top_3_rag_examples_out_of_sample	0.90
Qwen3-8B	3_random_examples	0.89
gpt-oss-120b	top_3_rag_examples_out_of_sample	0.88
gpt-oss-20b	1_random_examples	0.85
Qwen3-0.6B	top_3_rag_examples_out_of_sample	0.66
gpt-4.1-nano	3_random_examples	0.26

Model	Cost	Cost_all_tasks	Currency
gpt 4.1 Inp glbl Tokens	3.53	7.02	EUR
gpt 4.1 Outp glbl Tokens	2.71	3.44	EUR
gpt 4.1 mini Inp glbl Tokens	1.23	1.23	EUR
gpt 4.1 mini Outp glbl Tokens	0.71	0.71	EUR
gpt 4.1 nano Inp glbl Tokens	0.31	0.31	EUR
gpt 4.1 nano Outp glbl Tokens	0.15	0.15	EUR
gpt-oss-120B Outp glbl Tokens	0.85	0.85	EUR
gpt-oss-120B Inp glbl Tokens	0.42	0.42	EUR

5.2.2.2 Synthetic tables only

Table 5.13 shows that for 4 from 2 model families there is at least one model-method combination that performed better than the regex baseline. For the synthetic table extraction task the baseline is 0.9691401.

Only 18 from 2 model-method combinations performed better than this baseline. There has been no model that performed better than this baseline with the zero or static example approach.

span argument was not implemented correct in html tables and md :/

already just using 10 % of documents generated; and then 10 % of that sum of all experiment results (factor 14) with random forest?

A.15

²¹With the static example approach there have been 24 files where the response could not get parsed into valid json. With the other approaches there are one to four unparsable responses.

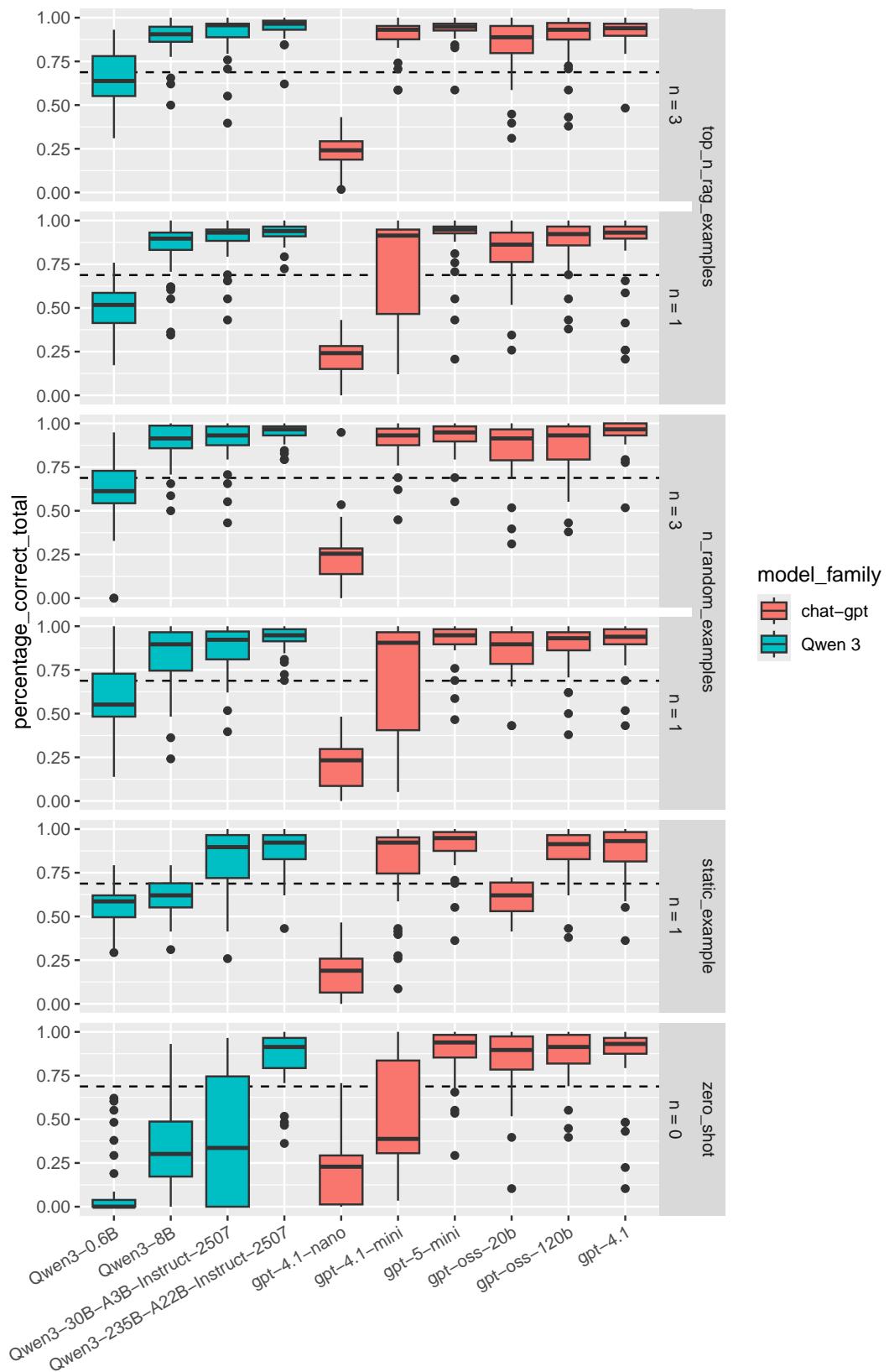


Figure 5.15: Comparing the percentage of correct predictions overall for OpenAi's LLMs with some Qwen 3 models

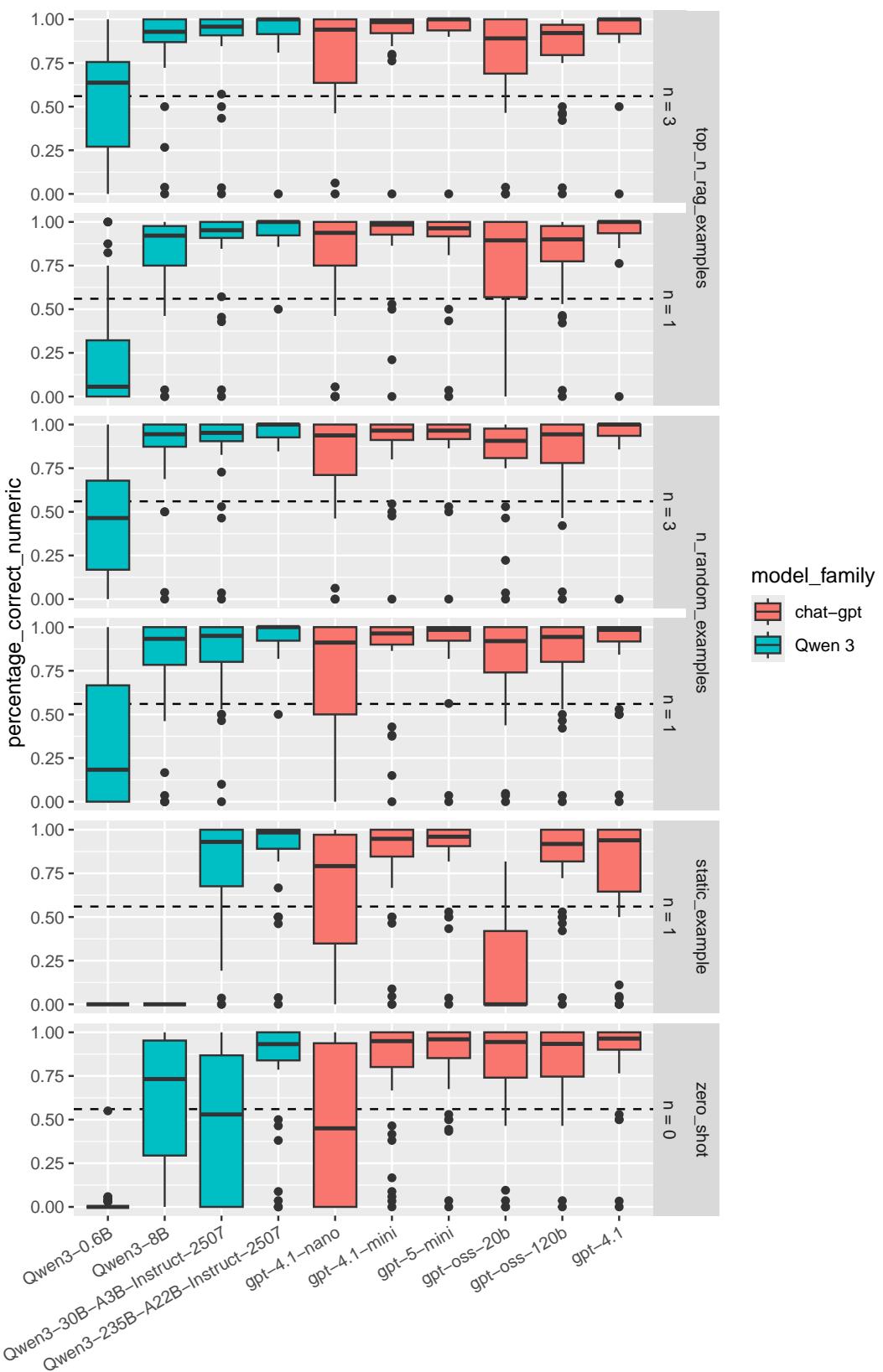


Figure 5.16: Comparing the percentage of correct numeric predictions for OpenAI’s LLMs with some Qwen 3 models

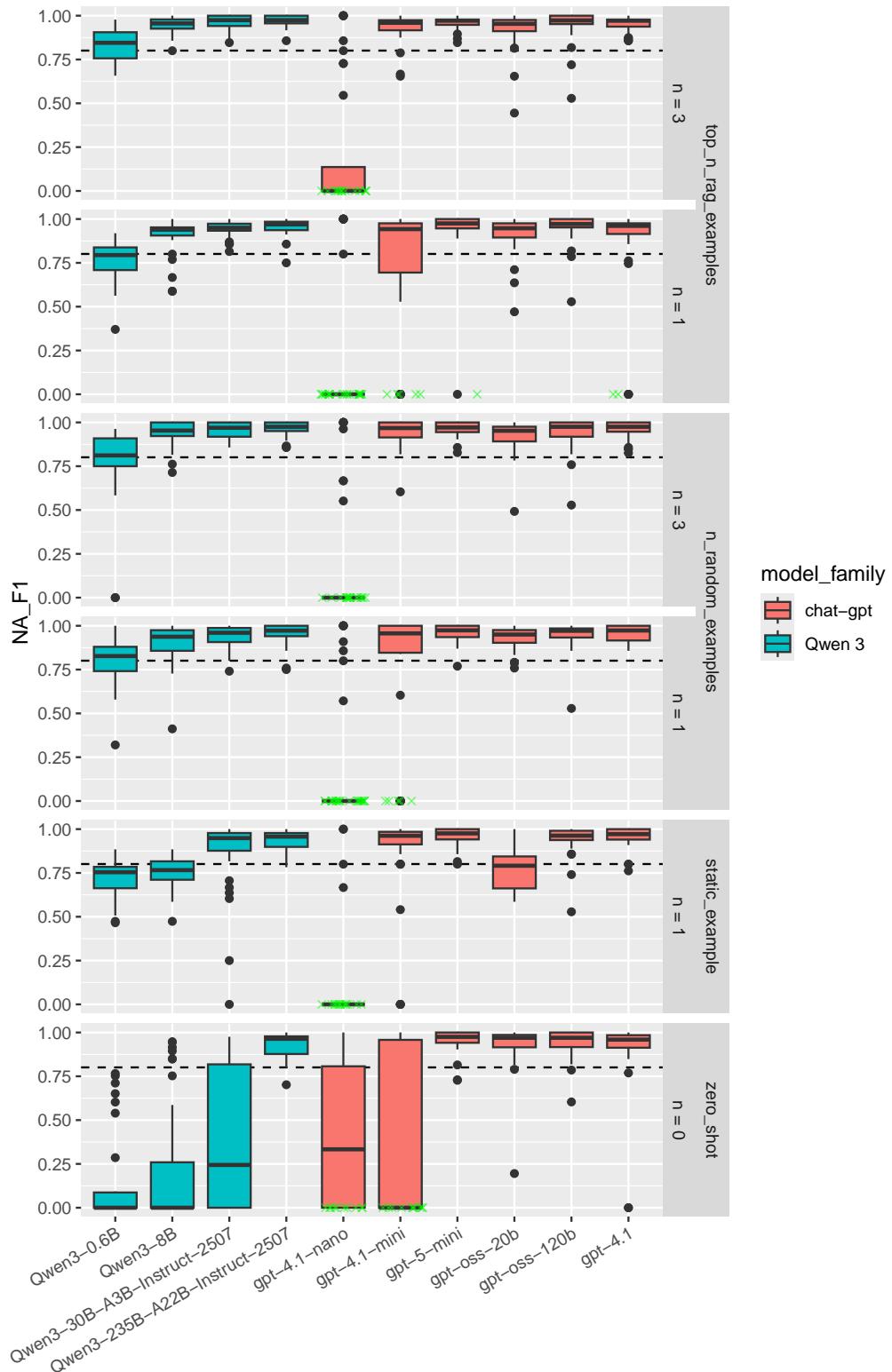


Figure 5.17: Comparing the F1 score for predicting the missingness of a value for OpenAI's LLMs with some Qwen 3 models. The green crosses indicate results where a model has predicted only numeric values even though there have been missing values.

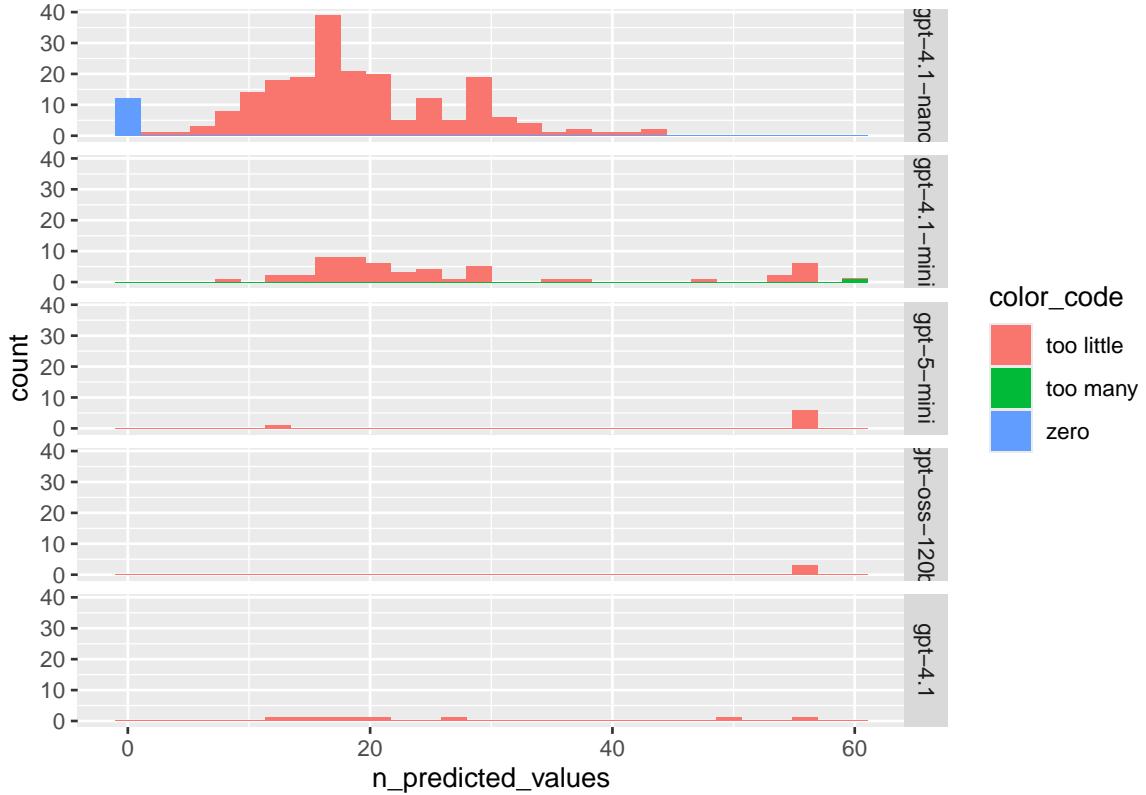


Figure 5.18: Showing the number of predictions OpenAI's models made.

Table 5.14: Comparing best mean table extraction performance with synthetic 'Aktiva' dataset for each model family for models with less than 17B parameters

model_family	model	method_family	n_examples	mean_total
Qwen 3	Qwen38B	top_n_rag_examples	3	0.959
Qwen 2.5	Qwen2.57BInstruct	top_n_rag_examples	5	0.942
mistralai	Minstral8BInstruct2410	top_n_rag_examples	5	0.935
Llama-3	Llama3.18BInstruct	top_n_rag_examples	3	0.897
google	gemma312bit	top_n_rag_examples	3	0.867

Table 5.15: Comparing extraction performance for real Aktiva extraction task with synthetic and real examples for incontext learning with a zero shot approach

model	method	mean_synth	mean_real	mean_zero_shot
Llama4Scout17B16EInstruct	top_5_rag_examples_out_of_sample	{0.887}	{0.925}	0.387
MistralLargeInstruct2411	5_random_examples	0.873	0.901	{0.691}
Qwen38B	5_random_examples	0.797	0.898	0.359
Llama3.18BInstruct	top_3_rag_examples_out_of_sample	0.764	0.805	0.5
gemma327bit	5_random_examples	0.732	0.821	0.255
Ministrall8BInstruct2410	3_random_examples	0.732	0.882	0.541
gemma312bit	top_1_rag_examples_out_of_sample	0.607	0.713	0.582

Hypotheses HTML and Markdown better but expected interaction effects mostly not found - except: - columns help pdf - thinning least bad for pdf - pdf worst with numbers that have currency units (short numbers, maybe no 1000er delimiter) - enumeration positive for pdf (and interaction with log10 mult)

line breaks are no problem

zero shot gets confused by text around

Markdown might be even better than HTML

respecting units was bad - except for: Top n rag finds examples with same currency units (shorter numbers more important than currnency in header?)

log10 multiplier has many interaction effects

LLama 4 Maverick again problem with five examples

Positive column count effect (different for real data)

5.2.2.3 Extract from real tables with synthetic content

Table 5.15 shows that using real examples for in-context- learning is better than using the created synthetic data. Nevertheless, it is improving the overall performance for the table extraction task by almost 20 % for all models but Google's gemma-3-12b-it. However, the performance difference with and without in-context learning is the smallest for the gemma-3-12b-it model as well. The spread of the performance is bigger using synthetic in-context learning data for all models but Llama 3 8B Instruct.

Any pattern?

Table 5.16 shows, that synthetic data can be used for in-context learning in a task where the currency units given for a table or specific columns²². Except for Llama 3.1 all models achieved much better numeric extraction results for tables where currencies are given for all columns. The model also achieved better for cases where there was only a single column with units, even though there have been no examples with units only for one column in the synthetic data. On the other hand, being prompted to respect currency units decreased the performance for tables where no units are given for all models but Llama 4 Schout. This decrease was highest for Qwen 3 and also higher than 10 % for Mistral Large.

A.6, A.7 and A.8

Thus, synthetic data can be used to solve new tasks and substitute missing data for rare classes.

Confidence with both tasks (respect or ignore units) the same.

²²Synthetic data is used here because the characterization, which real *Aktiva* table has units in which column, was created too late.

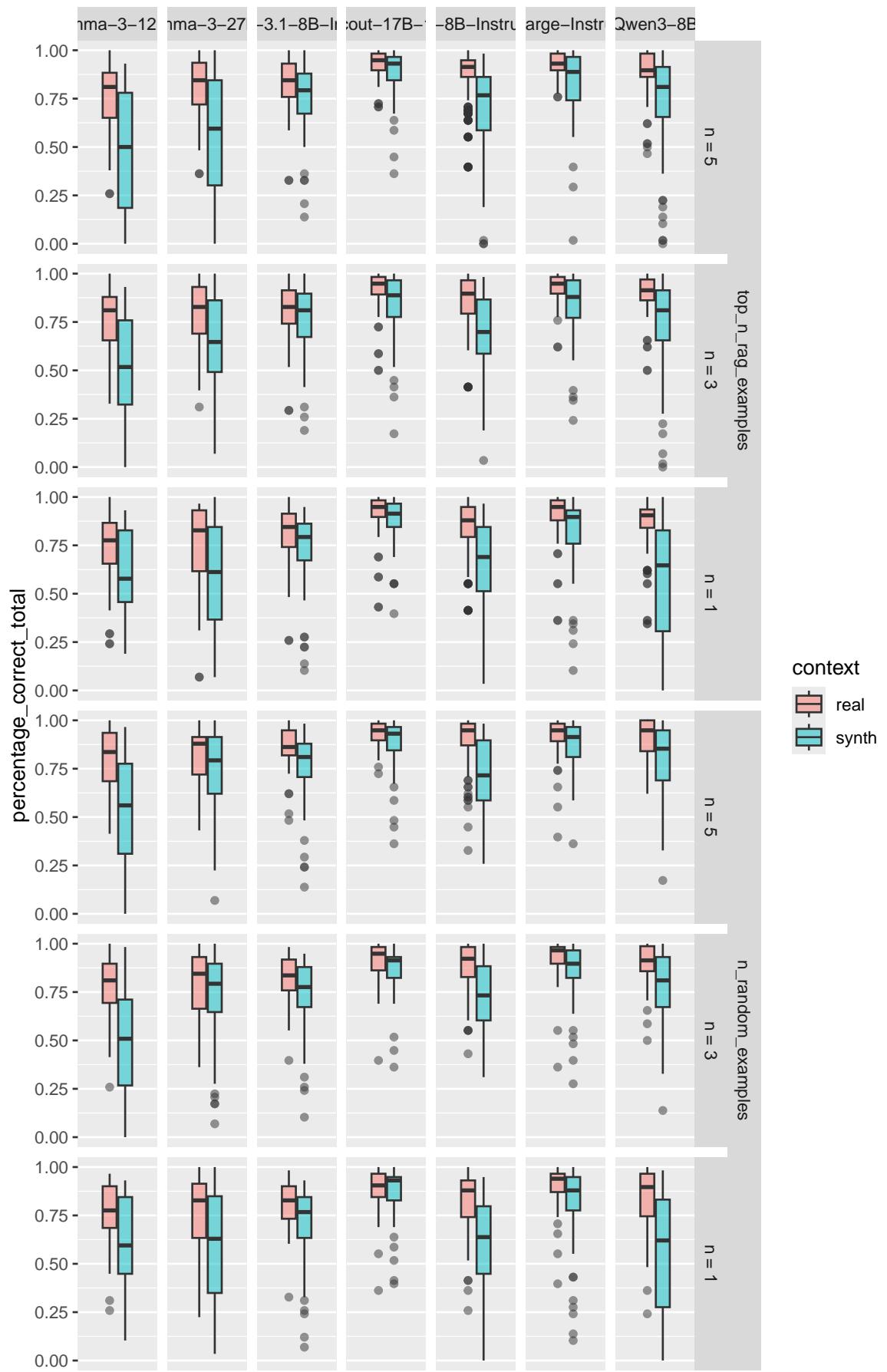


Figure 5.19: Comparing table extraction performance for real Aktiva extraction task with synthetic and real examples for in-context learning

Table 5.16: Comparing extraction performance for real Aktiva extraction task dependent on the prompt addition to respect currency units

model	n_cols_T_EUR_0	n_cols_T_EUR_1	n_cols_T_EUR_2
Llama3.18BInstruct	0.04	0.02	0.02
Llama4Scout17B16EInstruct	{0.01}	0.23	0.87
Minstral8BInstruct2410	0.06	0.19	0.65
MistralLargeInstruct2411	0.14	{0.4}	{0.9}
Qwen38B	0.31	0.1	0.51
gemma312bit	0.03	0.14	0.76
gemma327bit	0	0.3	0.36

Table 5.17: Comparing extraction confidence for real Aktiva extraction task dependent on the prompt addition to respect currency units. No difference in confidence apparent.

model	method	method_family	predicted_NA	mean_conf_units_FALSE	mean_c
Minstral8BInstruct2410	3_random_examples	n_random_examples	FALSE	0.86	0.84
Minstral8BInstruct2410	3_random_examples	n_random_examples	TRUE	0.97	0.97
Qwen38B	5_random_examples	n_random_examples	FALSE	0.98	0.98
Qwen38B	5_random_examples	n_random_examples	TRUE	1	1

```
confidence_vs_truth %>% group_by(model, method, method_family, respect_units, predicted_NA) %>% reframe(m
  mutate(across(is.numeric, ~round(., 2))) %>%
  render_table(
    alignment = "lllrr",
    caption = "Comparing extraction confidence for real Aktiva extraction task dependent on the prompt add
    ref = opts_current$get("label"), dom="t")
```

Risk for false NAs less with synth data for Mistral but greater for numeric values (both).

Hypotheses

5.2.3 Comparison

Most models performe better on synth tables once they have enough in-context examples. (Needing more für random examples thanwith top-n-rad approach). Especially Llama 3 models show wider performance spread even with three examples A.14

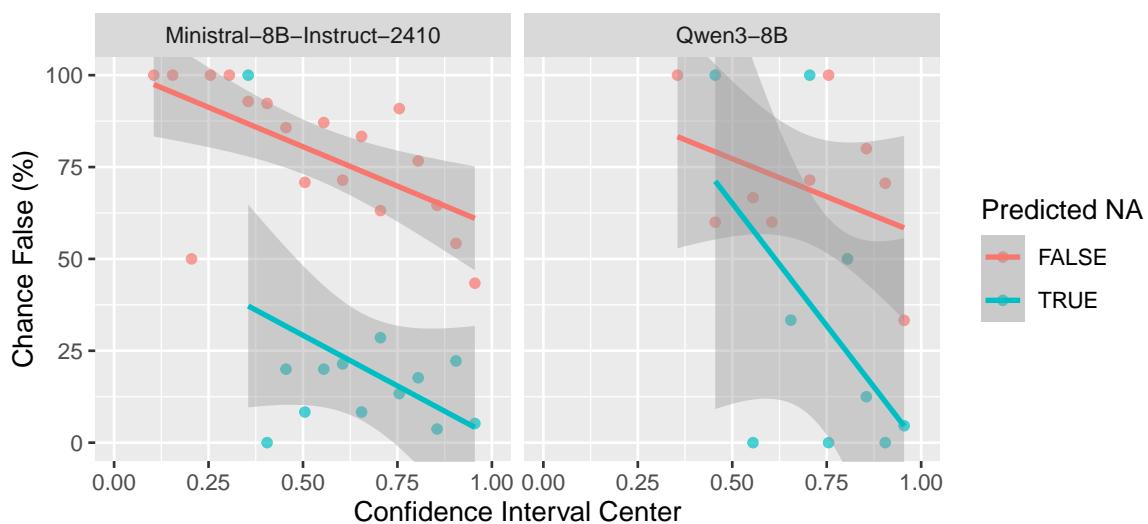


Figure 5.20: Estimating the relative frequency to find a wrong extraction result over different confidence intervals for predictions based on synthetic examples for in-context learning.



Chapter 6

Discussion

6.1 Limitations

6.1.1 table extraction

- found mistakes in gold standard with the llm results; mistakes found by human double check
- new lines / splitted lines
- test synthetic hypothesis with pymupdf extract

6.1.1.1 Regex baseline

- synthetic tables have been generated with cell lines because this should have improved the performance of a table extraction approach (not conducted)- maybe this is confusing pdfium? Or the zoom level?

6.1.2 classification

- Owen 2.5 hat zweiseitige GuV von IBB entdeckt und zur Anpassung der Ground Truth
- predictor: n_big_tables (tf or llm relevant?)

6.2 Not covered

- OCR
- fine-tuning
- using something smaller (e.g. LSTMs) instead LLMs
- building application, UX design (ref Ambacher 2024)
- table extraction (either VLLMs or classic approaches <- tried tabula but was not successful (because of missing visual traits)?)

in company document next / previous year more helpful than years further away?

6.3 Outlook

- ensemble from multiple models or are errors systematic? (e.g. Wohnungsbaugenossenschaften splitting some rows in multiple and none is picked?)

- check for hallucination vs wrong placed / repeated numbers
- no perfect score even with synthetic data
- flexible extraction (name something, find it, get it)
-

6.3.1 Table extraction

building a document extraction database document by document can improve performance taking advantage of same-company rag in-context learning

Chapter 7

Conclusion



References

- Auer, C., Lysak, M., Nassar, A., Dolfi, M., Livathinos, N., Vagenas, P., Ramis, C. B., Omenetti, M., Lindlbauer, F., Dinkla, K., Mishra, L., Kim, Y., Gupta, S., Lima, R. T. de, Weber, V., Morin, L., Meijer, I., Kuropiatnyk, V., & Staar, P. W. J. (2024). *Docling Technical Report*. arXiv. <https://doi.org/10.48550/arXiv.2408.09869>
- BMI, Referat O2 (Ed.). (2013). *Minikommentar zum Gesetz zur Förderung der elektronischen Verwaltung sowie zur änderung weiterer Vorschriften*.
- Grandini, M., Bagli, E., & Visani, G. (2020). *Metrics for Multi-Class Classification: An Overview*. arXiv. <https://doi.org/10.48550/arXiv.2008.05756>
- Li, H., Gao, H. (Harry), Wu, C., & Vasarhelyi, M. A. (2023). *Extracting Financial Data from Unstructured Sources: Leveraging Large Language Models* [{SSRN} {Scholarly} {Paper}]. Social Science Research Network. <https://doi.org/10.2139/ssrn.4567607>
- Zhong, X., Tang, J., & Yepes, A. J. (2019). *PubLayNet: Largest dataset ever for document layout analysis*. arXiv. <https://doi.org/10.48550/arXiv.1908.07836>

List of Figures

1.1	Overview of companies Berlin holds share at	1
5.1	Showing the number of pages (bar height) and number of documents (number above the bar) per company for the data used for the page identification task. Some documents would require ocr before being processed and were not used.	14
5.2	Comparing the performance among different companies.	16
5.3	Histogram of the number of lines in the first 5 pages of the annual reports	18
5.4	Comparing number of found TOC and amount of correct and incorrect predicted page ranges	19
5.5	Comparing number of fount TOC and amount of correct and incorrect predicted page ranges .	20
5.6	Comparint the predicted page range sizes. The red vertical line shows the mean and the green one shows the median of these sizes.	22
5.7	Showing the minimal distance of the predicted page range to the actual page number overthe logprobs of the models response confidence.	23
5.8	Showing the amount of correct and incorrect predicted page ranges (bars) and the percentage of correct predictions (black line).	24
5.9	Performance overall and on numeric value extraction with regular expressions. Showing single scores for *percentage correct numeric* on real tables to explain wide boxes.	39
5.10	Performance on classification for missing values with regular expressions	40
5.11	Showing the influence of the extraxtion library on the numeric text extraction task with synthetic data	42
5.12	Comparing the reported confidence scores for the table extraction task on real dataset for the Mistral and Qwen 3 with 8B parameters.	44
5.13	Estimating the relative frequency to find a wrong extraction result over different confidence intervals	45
5.14	Showing the influence of many examples on Llama 4 Maverick (A) and interaction between *T in year* and *vis separated rows* (B)	46
5.15	Comparing the percentage of correct predictions overall for OpenAi's LLMs with some Qwen 3 models	48
5.16	Comparing the percentage of correct numeric predictions for OpenAi's LLMs with some Qwen 3 models	49
5.17	Comparing the F1 score for predicting the missingness of a value for OpenAi's LLMs with some Qwen 3 models. The green crosses indicate results where a model has predicted only numeric values even though there have been missing values.	50

5.18 Showing the number of predictions OpenAI's models made.	51
5.19 Comparing table extraction performance for real Aktiva extraction task with synthetic and real examples for in-context learning	53
5.20 Estimating the relative frequency to find a wrong extraction result over different confidence intervals for predictions based on synthetic examples for in-context learning.	55
A.1 Comparing page identification metrics for different regular expressions for each classification task by type of the target table.	77
A.2 Mean absolute SHAP values and beeswarm plots for real table extraction with regular expression approach	78
A.3 Comparing number of found TOC and amount of correct and incorrect predicted page ranges .	79
A.4 Mean absolute SHAP values and beeswarm plots for real table extraction with LLMs	80
A.5 Mean absolute SHAP values and beeswarm plots for synth table extraction with regular expression approach	81
A.6 Comparing the effect on overall performance if currency units should be respected on all predictions and specifically on predictions where all or just some columns have units.	82
A.7 Comparing the effect on numeric performance if currency units should be respected on all predictions and specifically on predictions where all or just some columns have units.	83
A.8 Comparing the effect on NA F1 score if currency units should be respected on all predictions and specifically on predictions where all or just some columns have units.	84
A.9 Percentage of correct extracted or as missing categorized values for table extraction task on real Aktiva tables	85
A.10 Percentage of correct extracted numeric values for table extraction task on real Aktiva tables .	86
A.11 F1 score for the missing classification if a value is missing for table extraction task on real Aktiva tables	86
A.12 Comparing the overall extraction performance depending on the condition if examples from the same company can be used (only for Amt für Statistik Berlin-Brandenburg).	87
A.13 Comparing F1 score over normalized runtime for binary classification task. The normalized runtime is given in minutes of processing on a single B200. The time to load the model into the VRAM is excluded.	87
A.14 Comparing the table extraction performance among real and synthetic Aktiva tables	88
A.15 Percentage of correct extracted or as missing categorized values for table extraction task on synthetic Aktiva tables	88
A.16 Percentage of correct extracted numeric values for table extraction task on synthetic Aktiva tables	89
A.17 F1 score for the missing classification if a value is missing for table extraction task on synthetic Aktiva tables	89
A.18 Example balance sheet page from California's Annual Comprehensive Financial Report 2023 .	93
A.19 Flowchart of the extraction framework	94
A.20 Flowchart of the extraction framework of Li et al. (2023)	95

List of Tables

5.1	Showing the number of documents with multiple target tables per type and the number of target tables that span two pages.	15
5.2	Comparing page identification metrics for different regular expressions for each classification task by type of the target table.	15
5.3	Comparing the number and percentage of correct identified page ranges among the approaches.	18
5.4	Comparing the number and percentage end pages prediction for Aktiva and Passiva that are equal.	19
5.5	Comparing the number and percentage of correct identified page ranges among the approaches.	21
5.6	Comparing GPU time for page range prediction and table of contents extraction. Time in seconds per text processed.	21
5.7	Comparing the mean and median page range sizes.	21
5.9	Comparing table extraction performance with real 'Aktiva' dataset for models that perform well without or with little context learning	43
5.10	Comparing table extraction performance with real 'Aktiva' dataset for models that worse than the regex baselin with 3 or 5 examples for incontext learning	43
5.11	Comparing best mean table extraction performance with real 'Aktiva' dataset for each model family	43
5.12	Comparing best mean table extraction performance with real 'Aktiva' dataset for each model family for models with less than 17B parameters	44
5.13	Comparing best mean table extraction performance with synthetic 'Aktiva' dataset for each model family	47
5.14	Comparing best mean table extraction performance with synthetic 'Aktiva' dataset for each model family for models with less than 17B parameters	51
5.15	Comparing extraction performance for real Aktiva extraction task with synthetic and real examples for incontext learning with a zero shot approach	52
5.16	Comparing extraction performance for real Aktiva extraction task dependent on the prompt addition to respect currency units	54
5.17	Comparing extraction confidence for real Aktiva extraction task dependent on the prompt addition to respect currency units. No difference in confidence apparent.	54
A.1	Comparing extraction time (in seconds) for different Python package	70
A.2	Comparing time (in seconds) for processing ten asset tables using different libraries and approaches	74

Glossary

ACFR Annual Comprehensive Financial Report

GuV Gewinn- und Verlustrechnung

HGB Handelsgesetzbuch

LLM large language model

PDF Portable Document Format

RHvB Rechnungshof von Berlin

RechKredV Verordnung über die Rechnungslegung der Kreditinstitute, Finanzdienstleistungsinstitute und Wertpapierinstitute

SHAP SHapley Additive exPlanations

TOC table of contents

ebnf extended Backus–Naur form

glm generalized linear model

json JavaScript Object Notation

regex regular expression

Chapter A

Appendix

A.1 Local machine

One can find the specifications of the local machine used to run the less computationally demanding tasks below. It is a lightweight laptop device. Its performance cores support hyperthreading and have a clock range between 2.1 and 4.7 GHz. However, due to the flat design, there is little active cooling. Thus, thermal throttling starts rather quickly. It is therefore a reasonable assumption that most locally benchmarked tasks are running at 2.1 GHz. Despite this handicap, it has a sufficiently large RAM of 32 GB and 3 GB of NVMe disk space.

System Details Report

Report details

- **Date generated:** 2025-07-19 13:56:16

Hardware Information:

- **Hardware Model:** LG Electronics 17ZB90Q-G.AD79G
- **Memory:** 32.0 GiB
- **Processor:** 12th Gen Intel® Core™ i7-1260P × 16
- **Graphics:** Intel® Graphics (ADL GT2)
- **Disk Capacity:** 3.0 TB

Software Information:

- **Firmware Version:** A2ZG0150 X64
- **OS Name:** Ubuntu 24.04.2 LTS
- **OS Build:** (null)
- **OS Type:** 64-bit
- **GNOME Version:** 46
- **Windowing System:** Wayland
- **Kernel Version:** Linux 6.11.0-29-generic

Table A.1: Comparing extraction time (in seconds) for different Python package

package	runtime in s
pdfium	{14}
pymupdf	22
pypdf	218
pdfplumber	675
pdfminer	752
doclipseparse	1621

A.2 Benchmarks

A.2.1 Text extraction

A basic requirement for all succeeding tasks is, that the text gets extracted from the PDF files. As written in doclings technical report (Auer et al., 2024) the available open source libraries differ in their speed and restrictiveness of licensing. Since there are no benchmark results this report multiple libraries have been tested here.

The benchmark ran on the local machine described in section A.1. There have been 5256 pages to extract the text from.

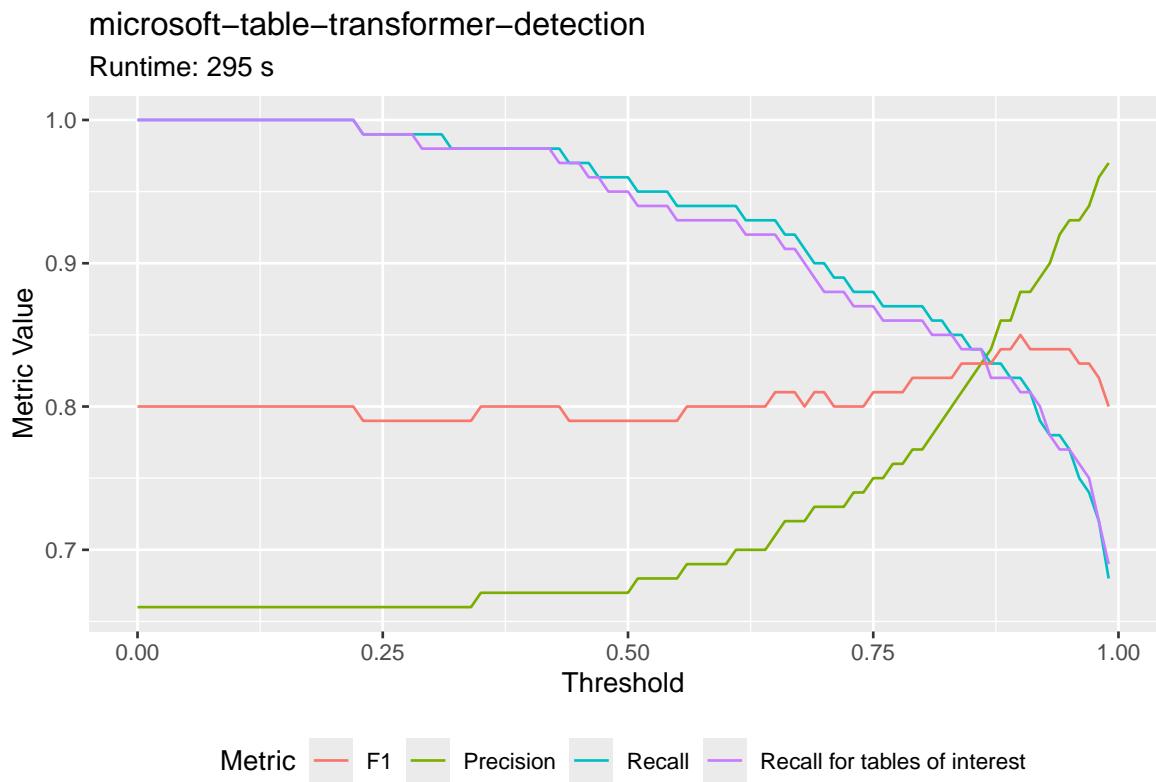
The result of docling-parse is not formated as markdown yet but also just plain text.

For implementation in a system where the text has to get extracted live or frequently the speed of the library might be paramount. But in special cases it can be important to invest more computational power into text extraction if this assures extraction according a more complicated document layout. E.g. some of the tables have been parsed by pdfium in such a manner that first all row descriptors have been extracted (first row) and thereafter all numeric columns (rowwise) ADD REFERENCE / EXAMPLE.

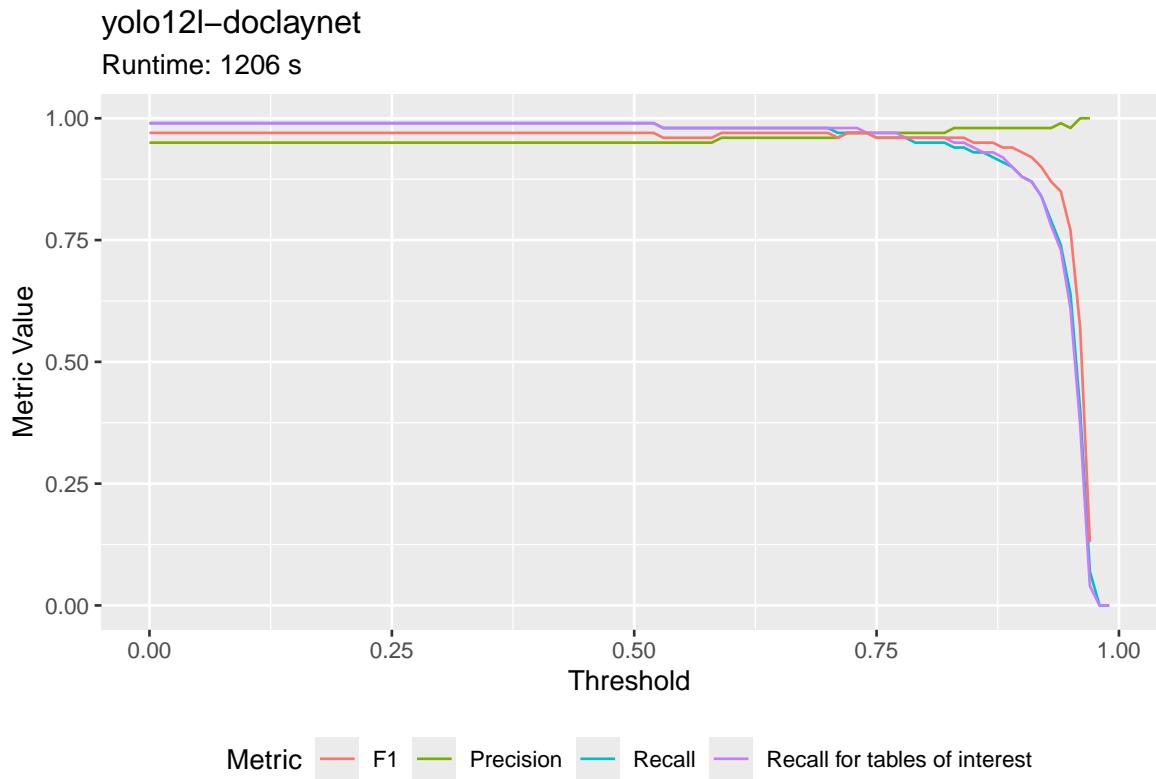
A.2.2 Table detection

- yolo benchmark and table transformer
- skip classification with llm

not so important anymore

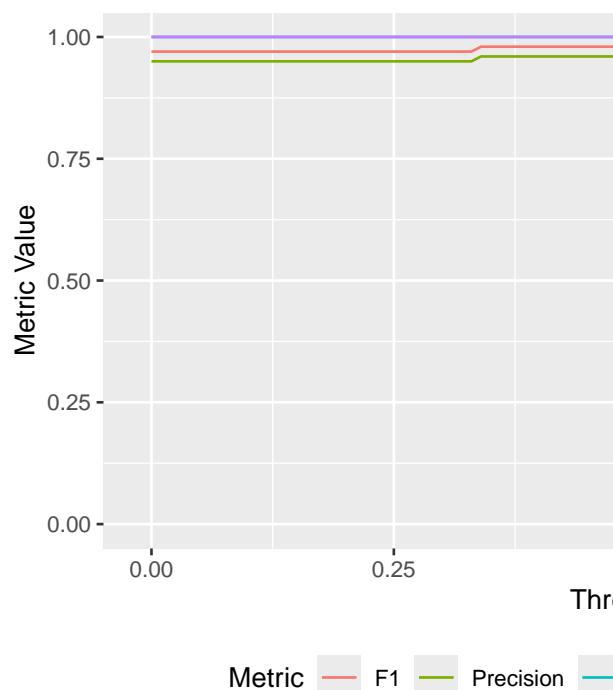


You see the plot for: microsoft-table-transformer-detection. (Click to stop automatic rotation.)



yolo12n-doclaynet

Runtime: 200 s

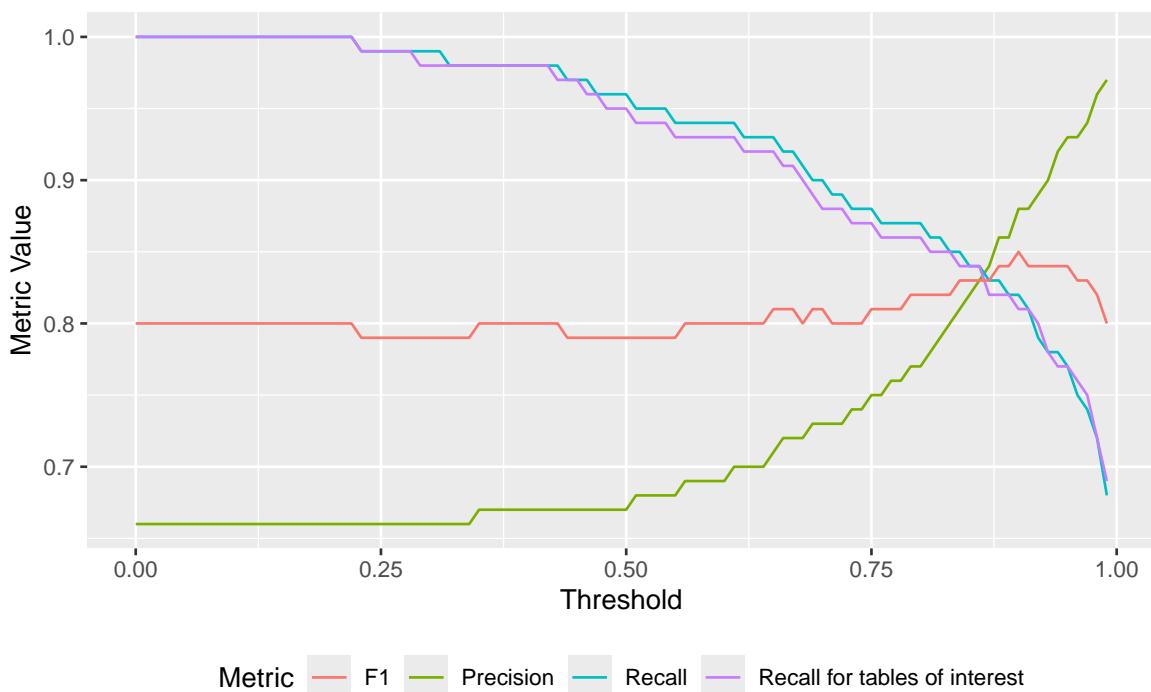


You see the plot for: yolo12l-doclaynet. (Click to stop automatic rotation.)

You see the plot for: yolo12n-doclaynet. (Click to stop automatic rotation.)

microsoft-table-transformer-detection

Runtime: 295 s



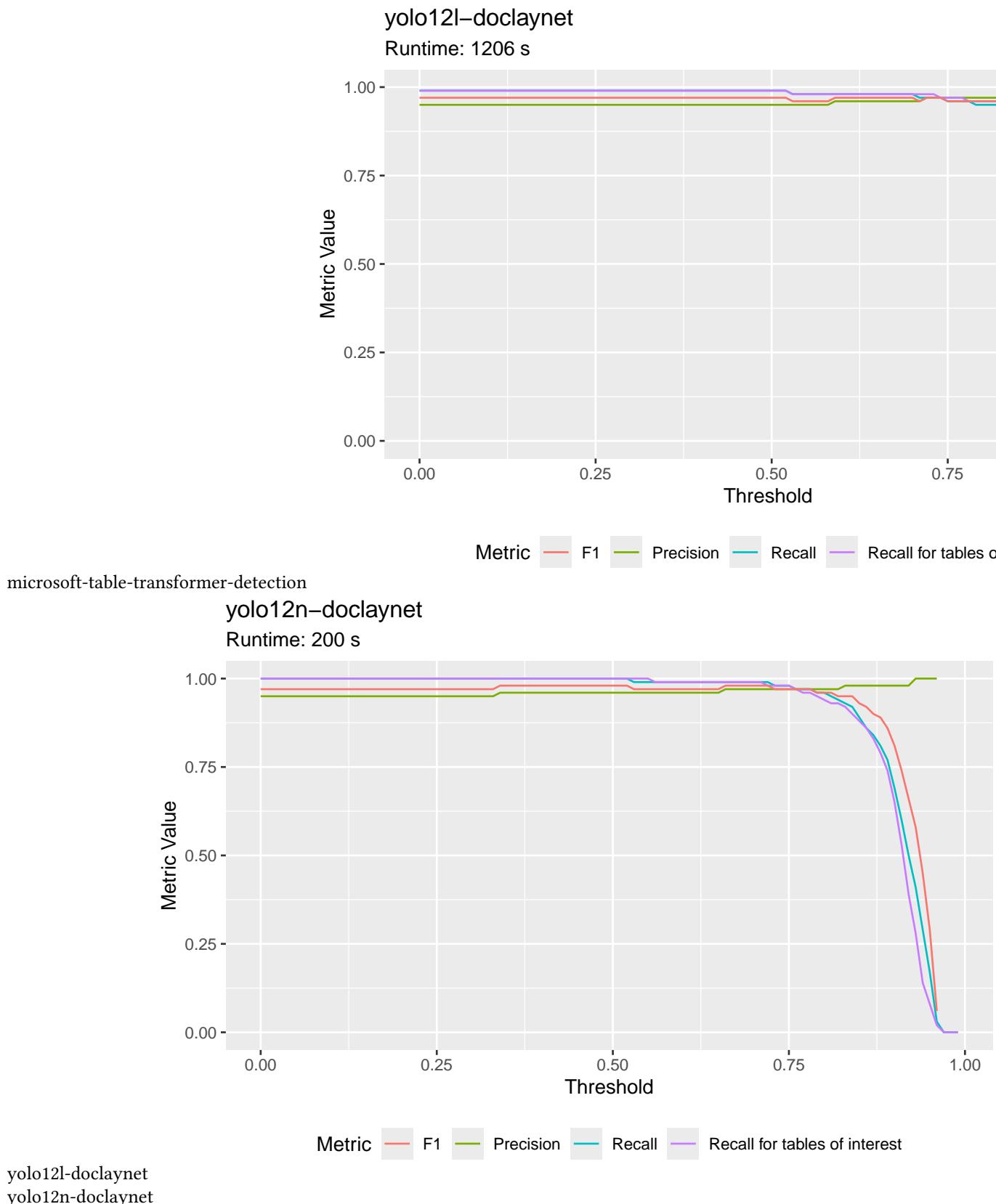


Table A.2: Comparing time (in seconds) for processing ten asset tables using different libraries and approaches

Model parameters (in B)	Transformers	vLLM	vLLM batched
0.5	330	65	NA
3.0	628	130	20
7.0	940	217	30

A.2.3 Large language model process speed

In April 2025 there have been issues with running vllm within the Python framework. Thus the first experiments have been conducted using the transformers library. When the problems of building a working vllm based docker image for the experiments it was measured how long the same task takes with the transformers and the vllm library and how the batched processing competes versus a loop approach. The model family used was Qwen 2.5 Instruct. The task was to extract the assets table for ten real example pages.

Table A.2 shows that the experiments with vllm library run are around four to five times faster. Processing the messages in a batched mode again is six to seven times faster.

The change of the experimental setup from transformers loop-based to vllm batched mode made is possible run the benchmark on whole PDF documents giving a sound estimate of the false positive rate in the page identification task (see section 5.1.3). Previous experiments have only been using a subset of pages that have been selected with the baseline regex approach (see section 5.1.1).

A.3 Prompts

A.3.1 TOC understanding

Base prompt:

```
messages = [
    {"role": "system", "content": "You are a helpful assistant that can determine the page range info from a table of contents."},
    {"role": "user", "content": f"This is the table of contents:\n\n{toc_string}"},
    {"role": "user", "content": f"On which pages might the win and loss statement (in German: Gewinn- und Verlustrechnung) appear? Provide a list of page numbers."},
    {"role": "user", "content": f"Answer in JSON format with keys 'GuV', 'Aktiva', and 'Passiva' and their respective page ranges."}
]
```

First attempt:

```
specific_prompt = {"role": "user", "content": f"The assets and liabilities tables often are on separate pages. Determine the page range for each table from the TOC."}
```

Given hint that assets and liabilities are part of the balance sheet:

```
specific_prompt = {"role": "user", "content": f"The assets and liabilities are part of the balance sheet. Determine the page range for each table from the TOC."}
```

Stating, that liabilities are on next page:

```
specific_prompt = {"role": "user", "content": f"The assets and liabilities are part of the balance sheet. Liabilities are on the next page. Determine the page range for each table from the TOC."}
```

TOC extraction from text prompt:

```
messages = [
    {"role": "system", "content": "[Role] You are a helpful assistant that can identify table of contents in PDF files."},
    {"role": "system", "content": f"[Context] These are the text lines of the first {i} pages:\n\n{i}\n\n{content}"}
    {"role": "user", "content": f"[Tasks] 1. Please identify if there is a table of contents in the document.\n2. If there is a table of contents, please extract its text."},
    {"role": "user", "content": f"3. Answer as JSON with the table of contents text as string in the 'text' field."},
    {"role": "user", "content": f"If there is no table of contents, return an empty string."},
]
```

A.4 Regular expressions

Here one can find the three regular expressions used for the benchmarks presented in section 5.1.1.

```
simple_regex_patterns = {
    "Aktiva": [
        r"aktiv",
        r"((20\d{2}).*(20\d{2}))"
    ],
    "Passiva": [
        r"passiva",
        r"((20\d{2}).*(20\d{2}))"
    ],
    "GuV": [
        r"gewinn",
        r"verlust",
        r"rechnung",
        r"((20\d{2}).*(20\d{2}))"
    ]
}

regex_patterns_5 = {
    "Aktiva": [
        r"a\s*k\s*t\s*i\s*v\s*a|a\s*k\s*t\s*i\s*v\s*s\s*e\s*i\s*t\s*e|anlageverm.{1,2}gen",
        r"((20\d{2}).*(20\d{2}))|((20\d{2}).*vorjahr)|vorjahr",
        r"Umlaufverm.{1,2}gen|Anlageverm.{1,2}gen|Rechnungsabgrenzungsposten|Forderungen",
        r"\s([a-zA-Z]|[0-9]{1,2}|[iI]+)[\.\.\.]\s"
    ],
    "Passiva": [
        r"p\s*a\s*s\s*s*i\s*v\s*a|p\s*a\s*s\s*s*i\s*v\s*s\s*e\s*i\s*t\s*e|eigenkapital",
        r"((20\d{2}).*(20\d{2}))|((20\d{2}).*vorjahr)|vorjahr",
        r"Eigenkapital|R.{1,2}ckstellungen|Verbindlichkeiten|Rechnungsabgrenzungsposten",
        r"\s([a-zA-Z]|[0-9]{1,2}|[iI]+)[\.\.\.]\s"
    ],
    "GuV": [
        r"gewinn|guv",
        r"verlust|guv",
        r"rechnung|guv",
        r"((20\d{2}).*(20\d{2}))|vorjahr",
        r"Umsatzerl.{1,2}se|Materialaufwand|Personalaufwand|Abschreibungen|Jahres.{1,2}berschuss|Jahres.{1,2}ausgaben",
        r"\s([a-zA-Z]|[0-9]{1,2}|[iI]+)[\.\.\.]\s"
    ]
}
```

```

regex_patterns_3 = {
    "Aktiva": [
        r"a\s*k\s*t\s*i\s*v\s*a|a\s*k\s*t\s*i\s*v\s*s\s*e\s*i\s*t\s*e|anlageverm.{1,2}gen",
        r"((20\d{2}).*(20\d{2}))|((20\d{2}).*vorjahr)|vorjahr"
    ],
    "Passiva": [
        r"p\s*a\s*s\s*i\s*v\s*a|p\s*a\s*s\s*s\s*i\s*v\s*s\s*e\s*i\s*t\s*e|eigenkapital",
        r"((20\d{2}).*(20\d{2}))|((20\d{2}).*vorjahr)|vorjahr"
    ],
    "GuV": [
        r"gewinn|guv",
        r"verlust|guv",
        r"rechnung|guv",
        r"((20\d{2}).*(20\d{2}))|vorjahr"
    ]
}

```

A.5 Figures

A.5.1 Page identification

A.5.1.1 Regex baseline

A.5.1.2 TOC understanding

A.5.2 Table extraction

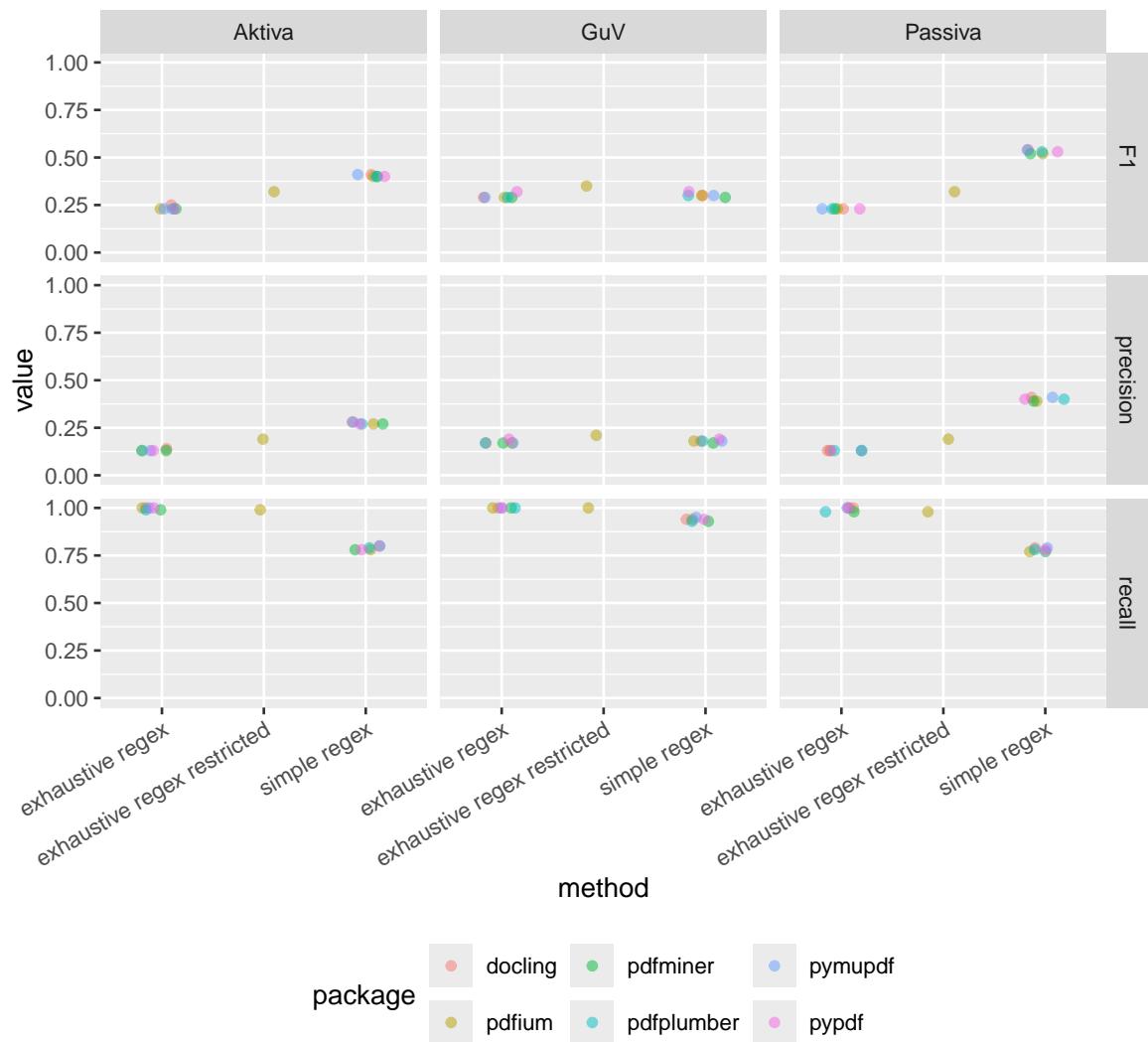


Figure A.1: Comparing page identification metrics for different regular expressions for each classification task by type of the target table.

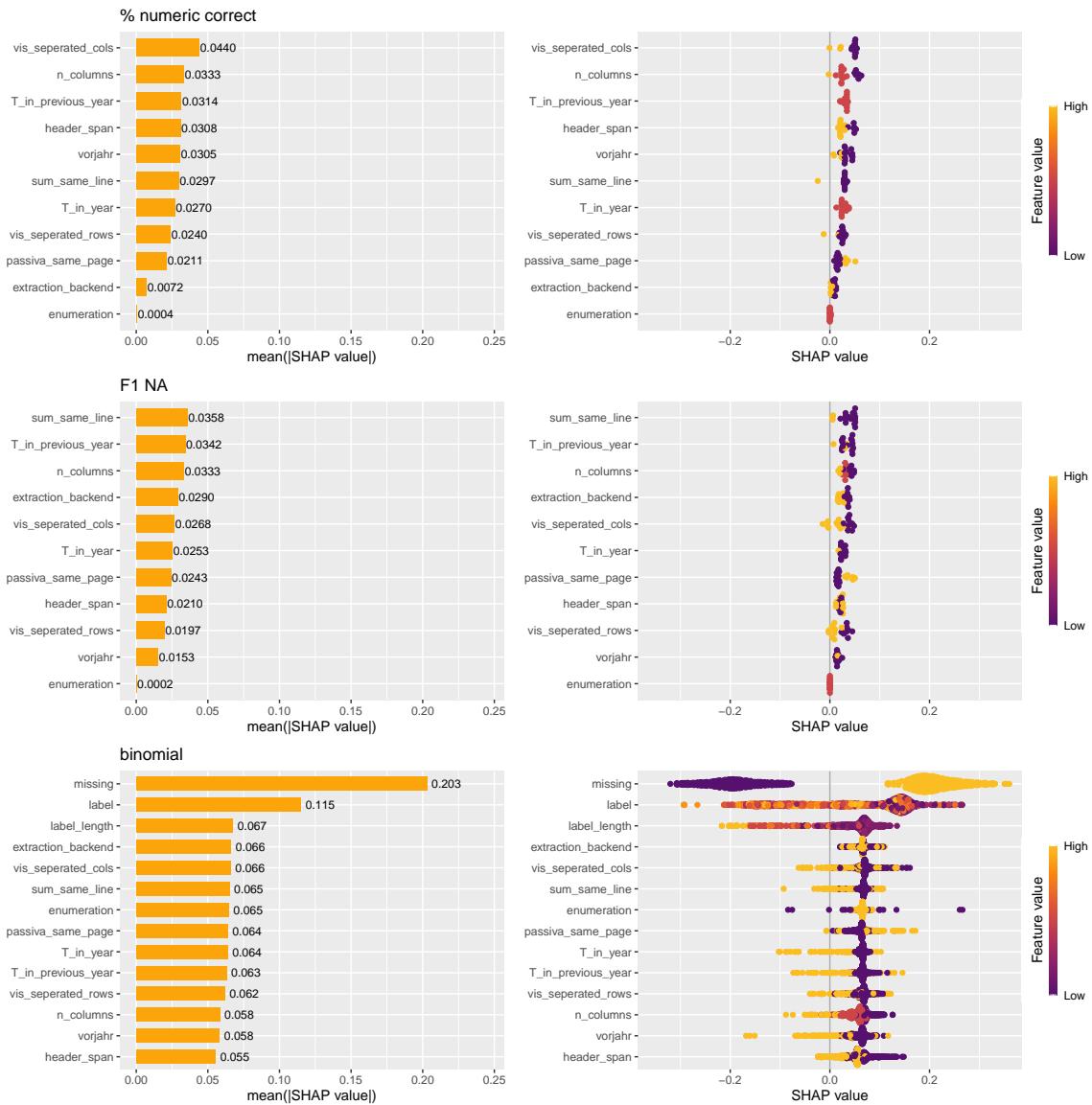


Figure A.2: Mean absolute SHAP values and beeswarm plots for real table extraction with regular expression approach

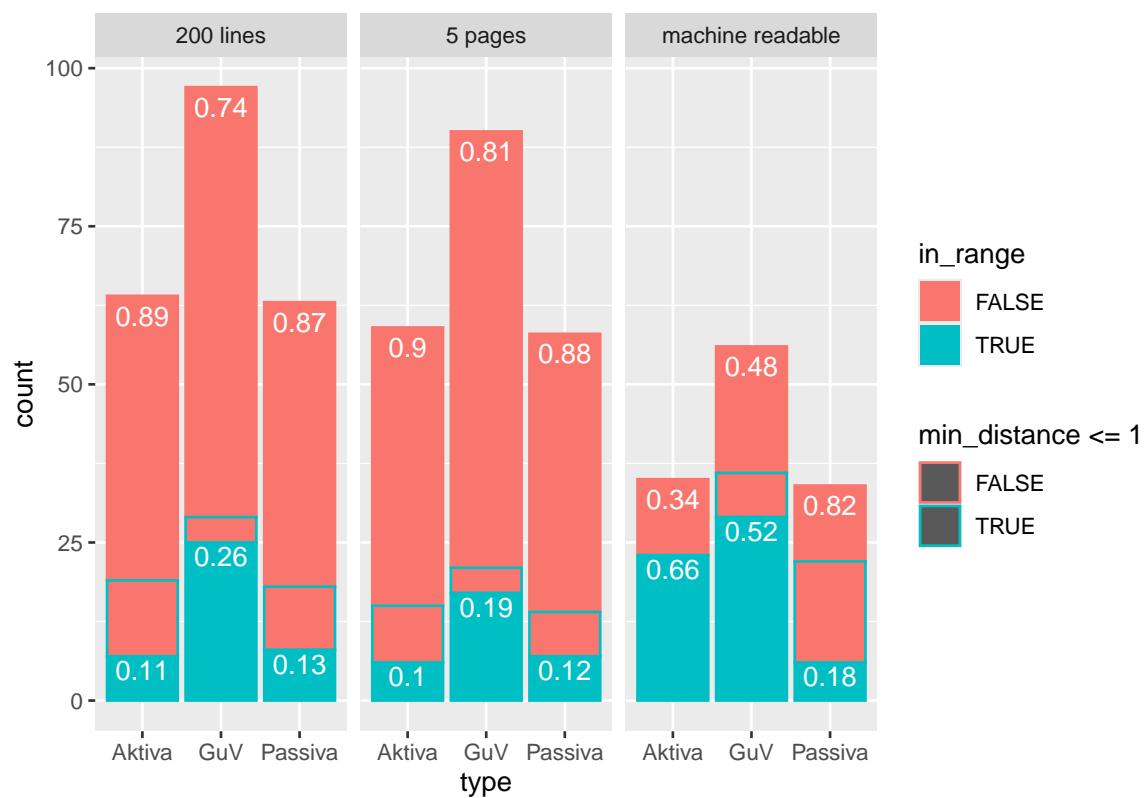


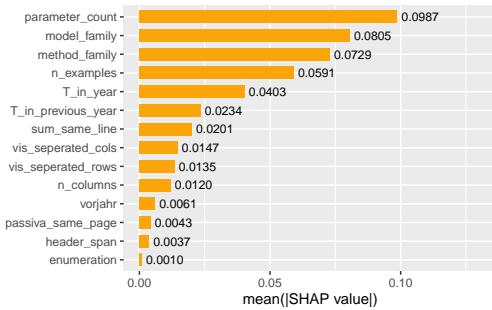
Figure A.3: Comparing number of fount TOC and amount of correct and incorrect predicted page ranges

The surprising truth about mtcars

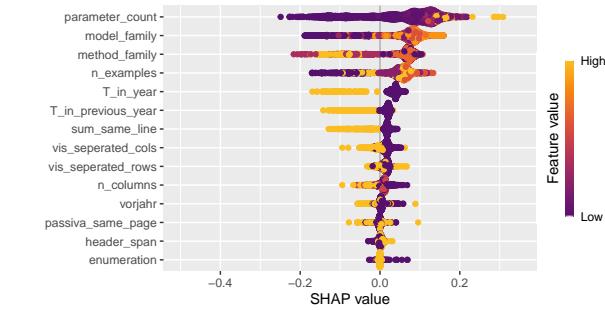
These 3 plots will reveal yet-untold secrets about our beloved data-set

A.1

% numeric correct

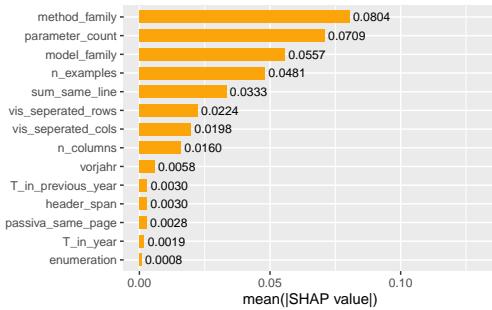


A.2

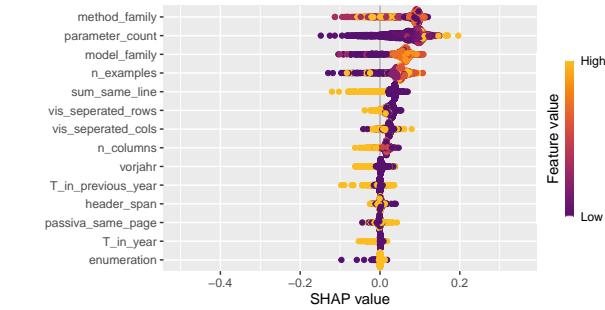


B.1

F1 NA

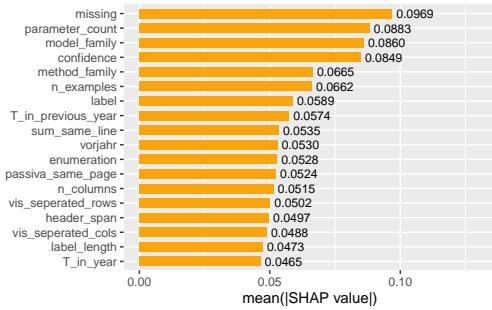


B.2

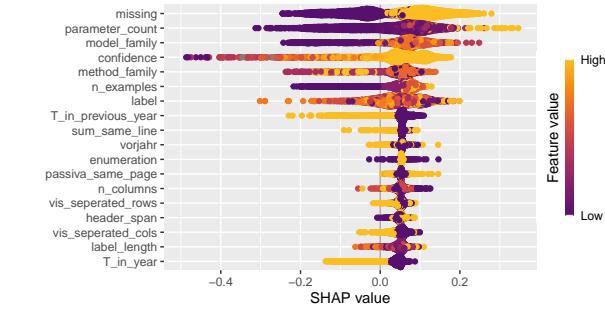


C.1

binomial

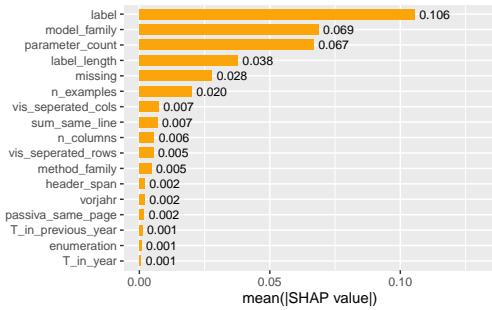


C.2

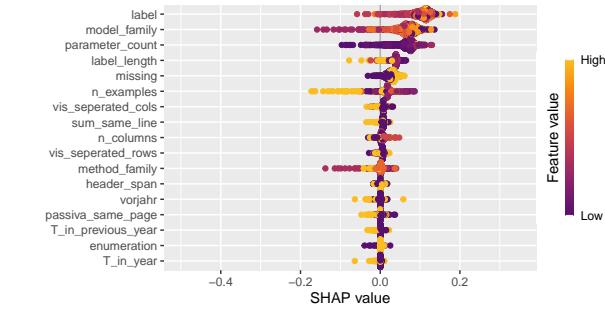


D.1

confidence



D.2



Disclaimer: None of these plots are insightful

Figure A.4: Mean absolute SHAP values and beeswarm plots for real table extraction with LLMs

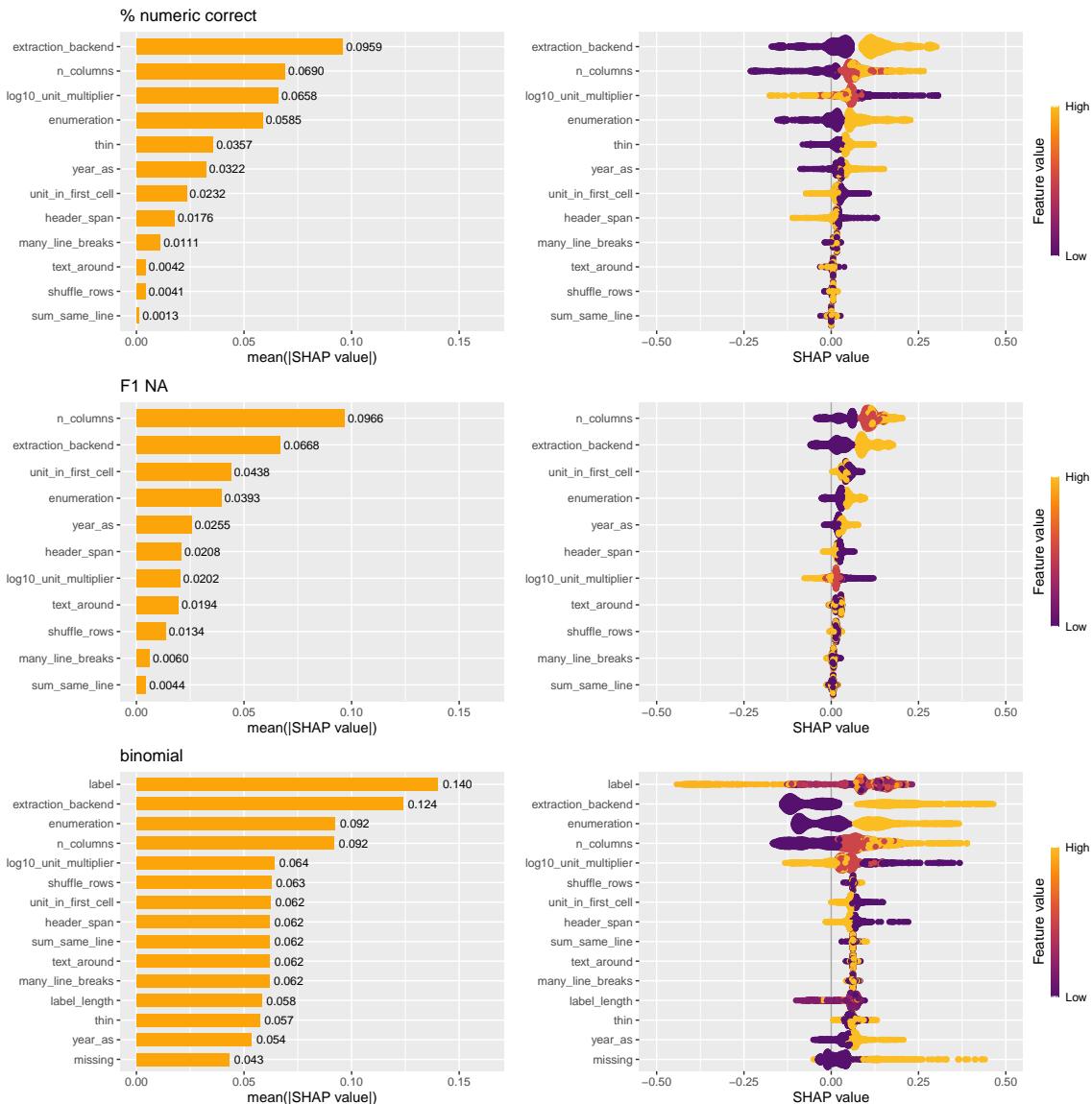
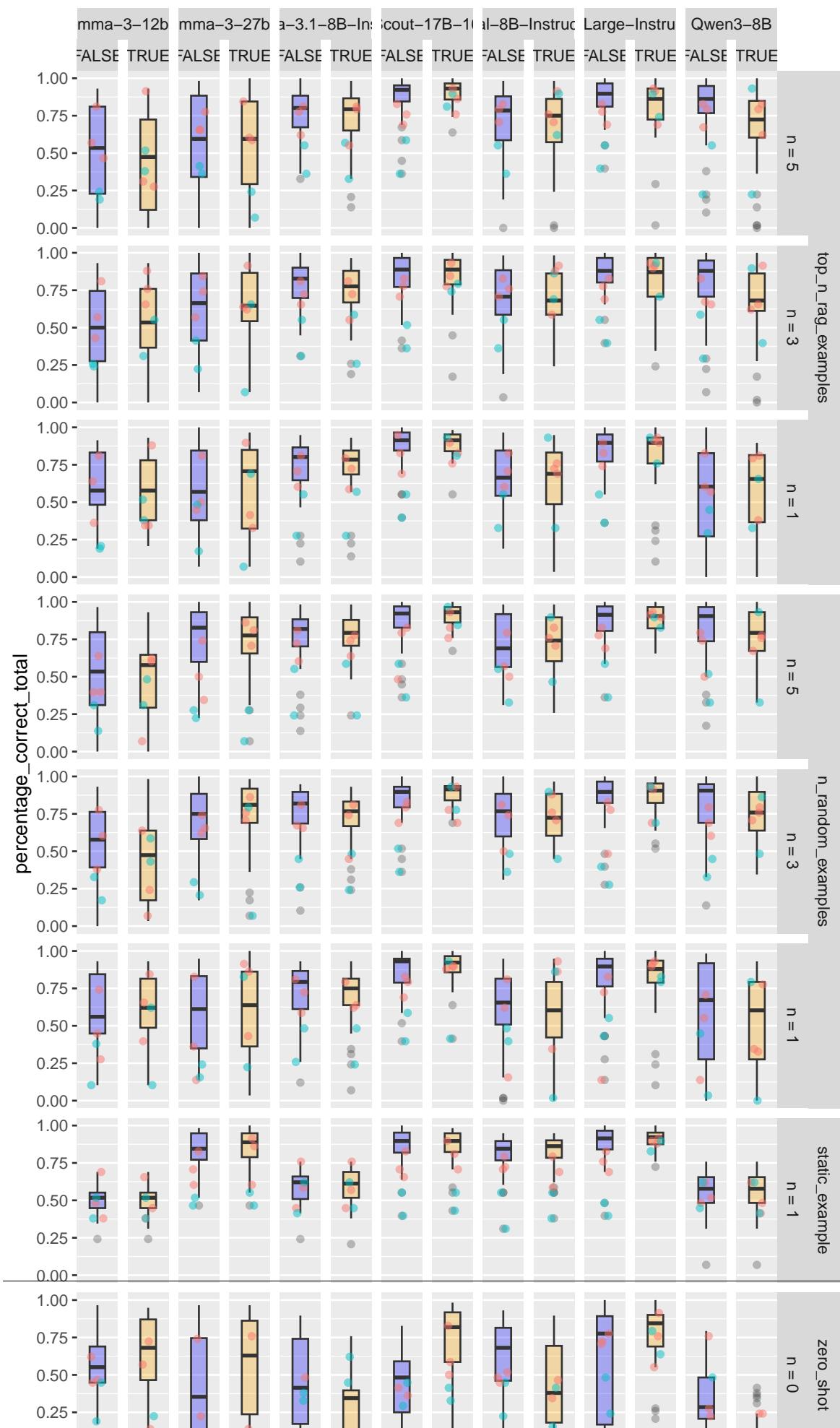
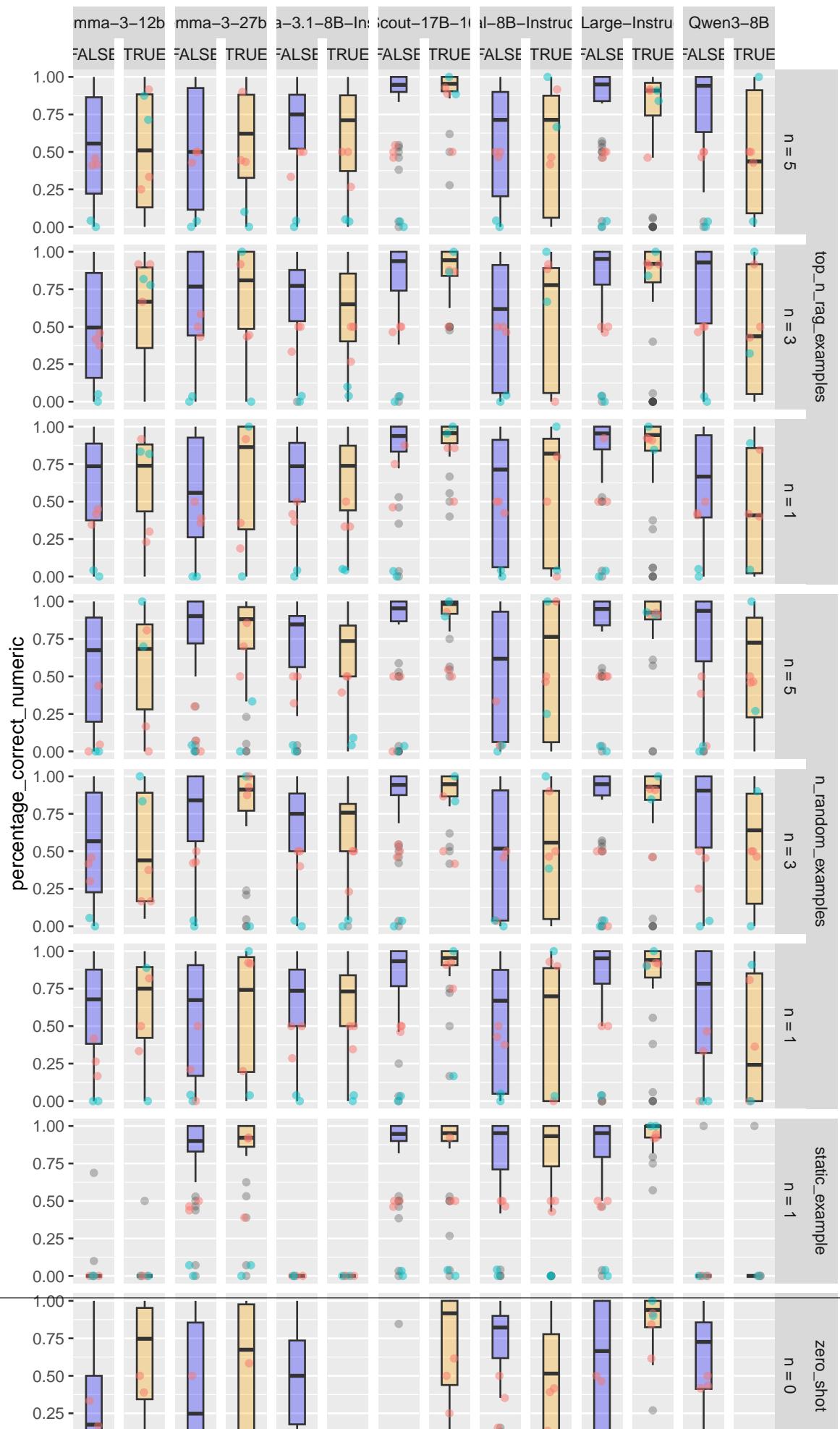
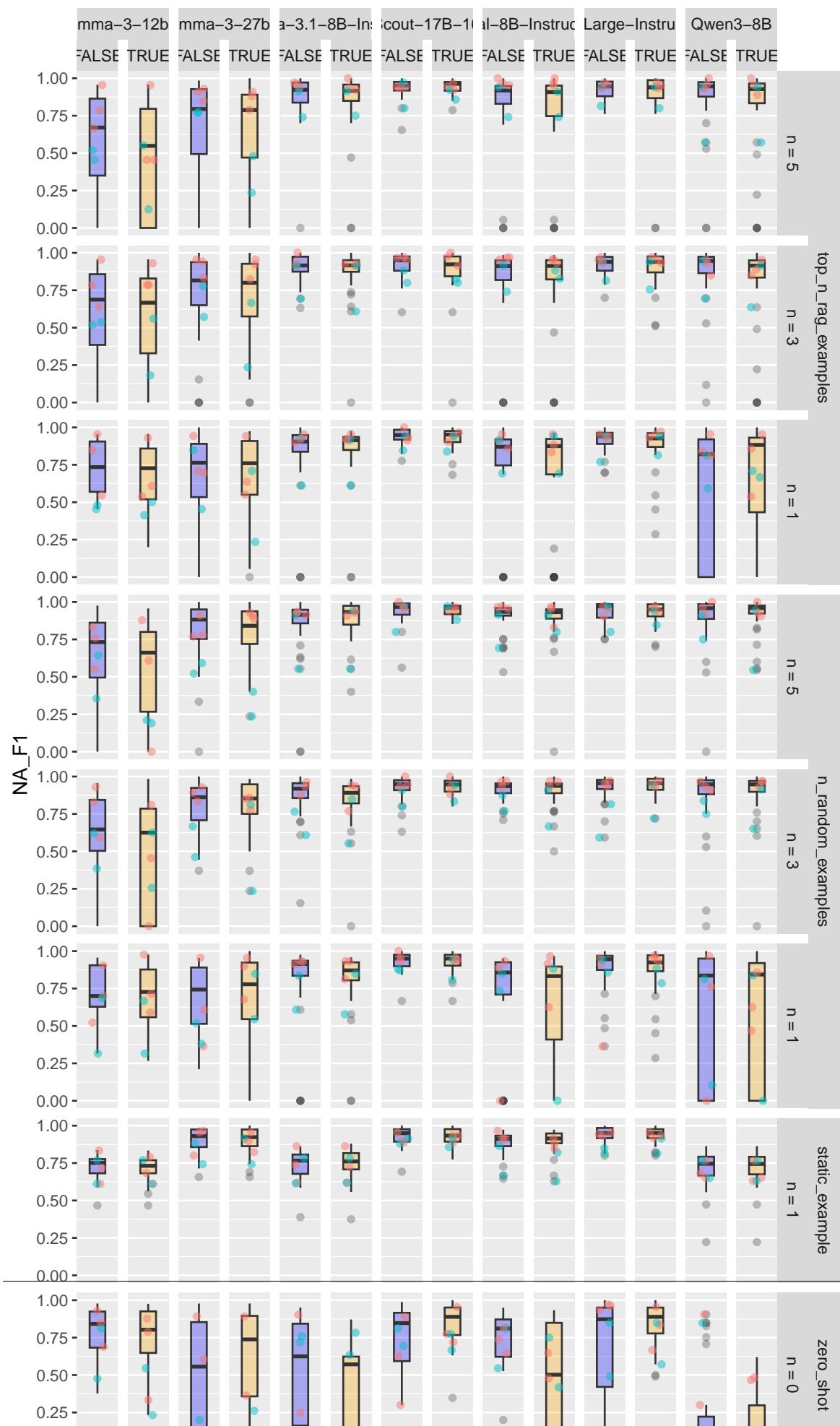


Figure A.5: Mean absolute SHAP values and beeswarm plots for synth table extraction with regular expression approach







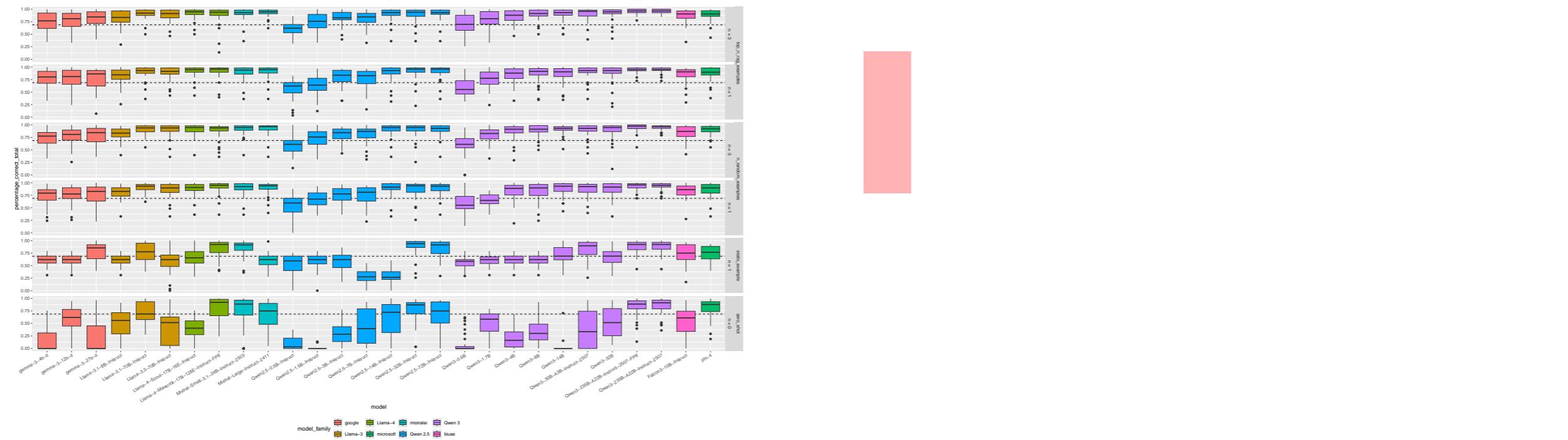


Figure A.9: Percentage of correct extracted or as missing categorized values for table extraction task on real Aktiva tables

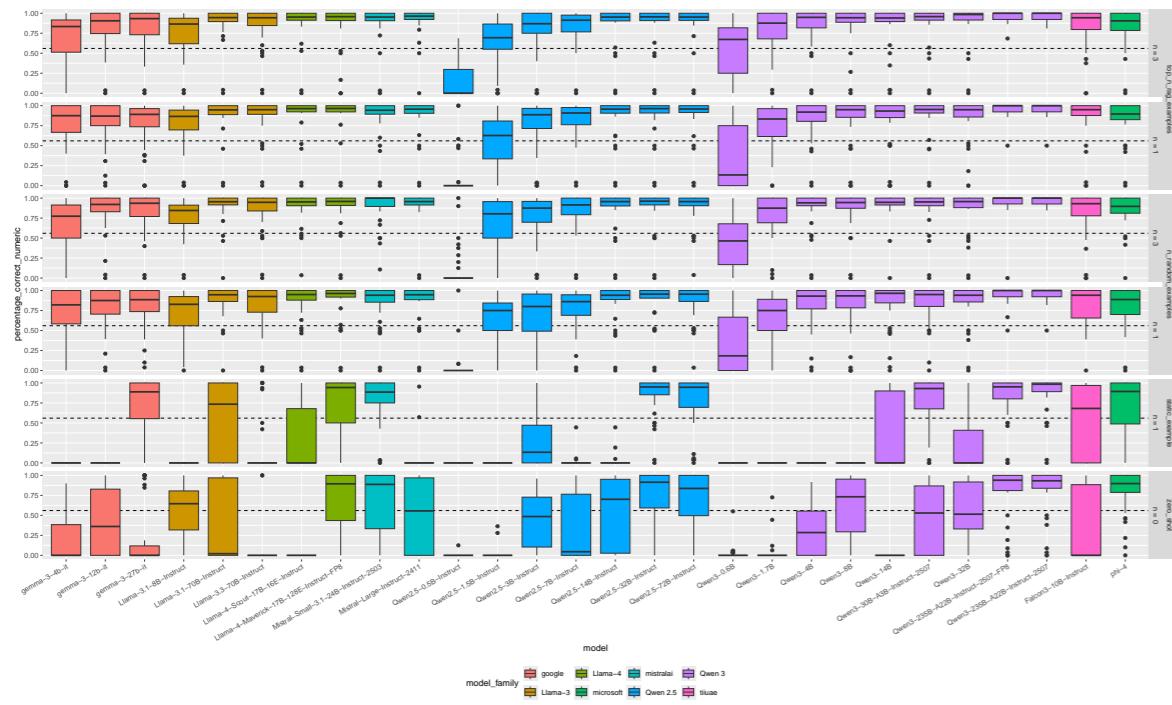


Figure A.10: Percentage of correct extracted numeric values for table extraction task on real Aktiva tables

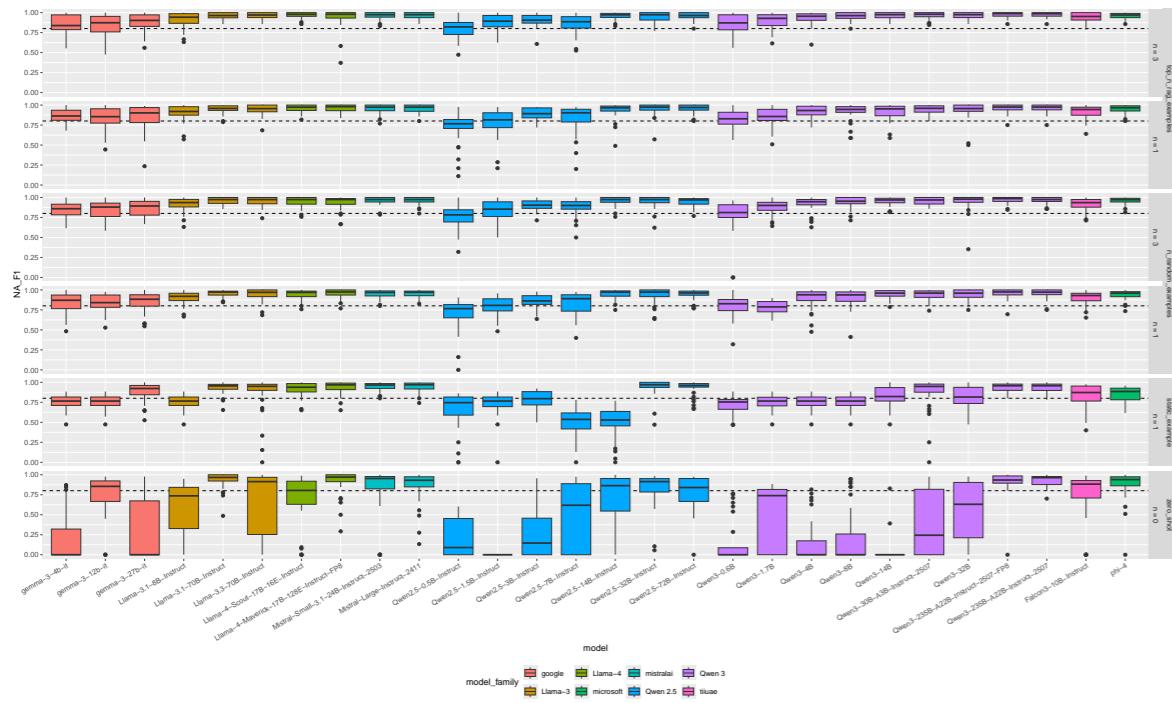


Figure A.11: F1 score for the missing classification if a value is missing for table extraction task on real Aktiva tables

A.5. FIGURES

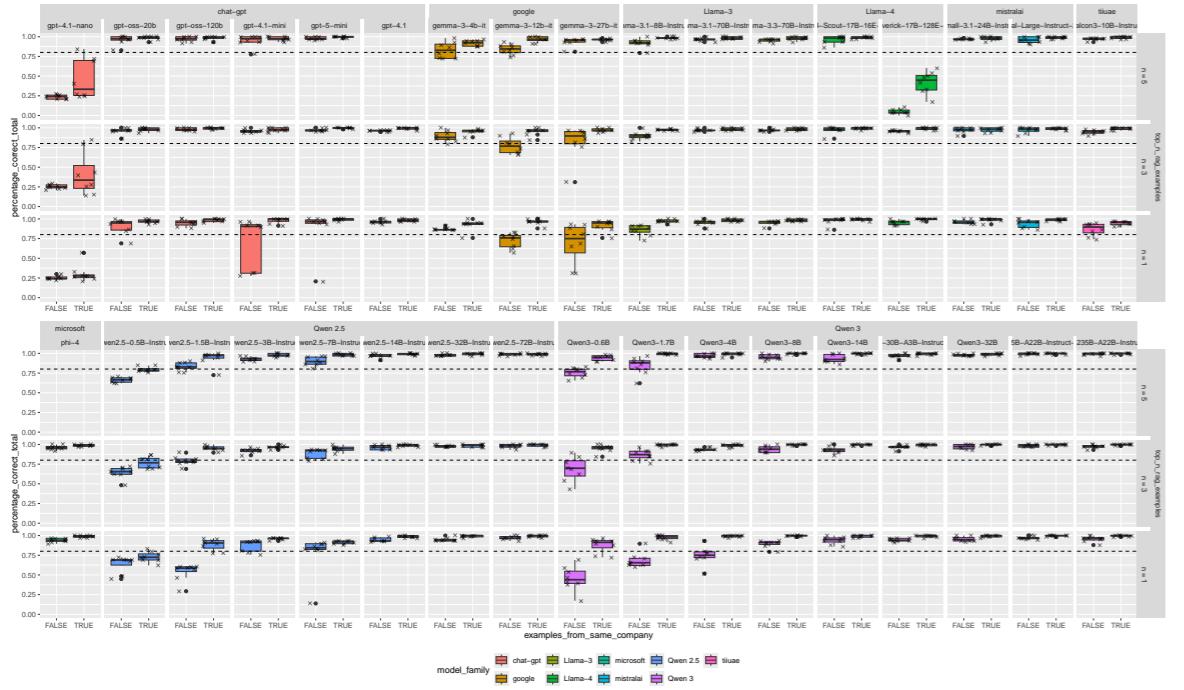


Figure A.12: Comparing the overall extraction performance depending on the condition if examples from the same company can be used (only for Amt für Statistik Berlin-Brandenburg).

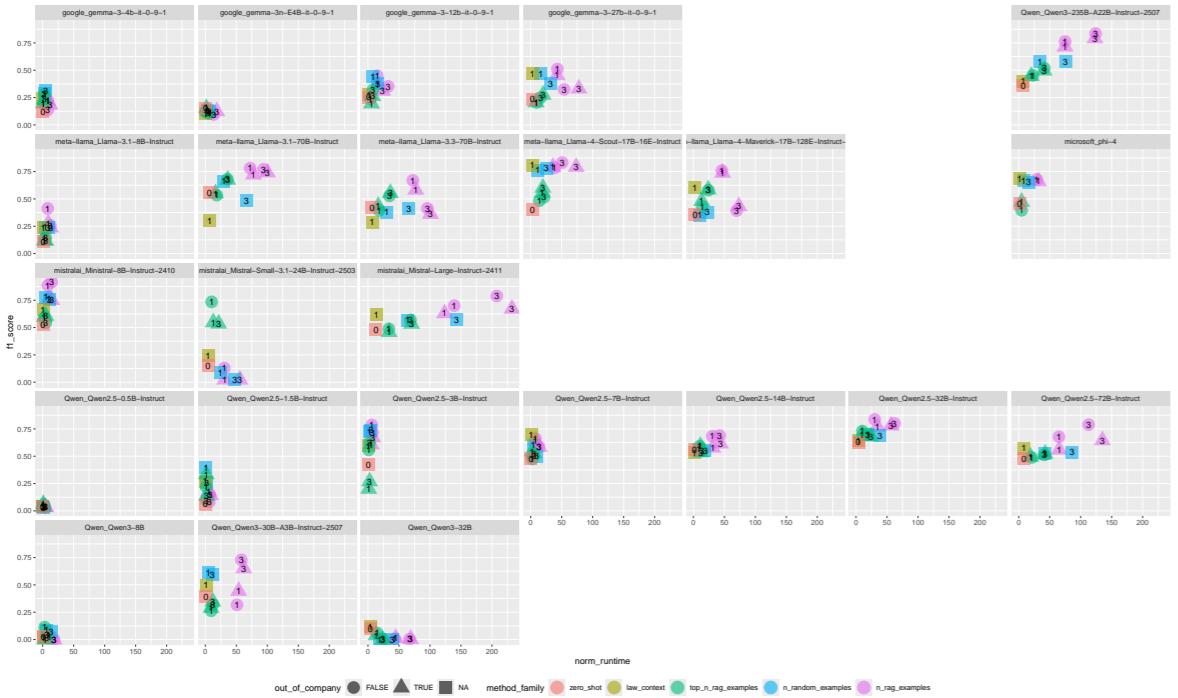


Figure A.13: Comparing F1 score over normalized runtime for binary classification task. The normalized runtime is given in minutes of processing on a single B200. The time to load the model into the VRAM is excluded.

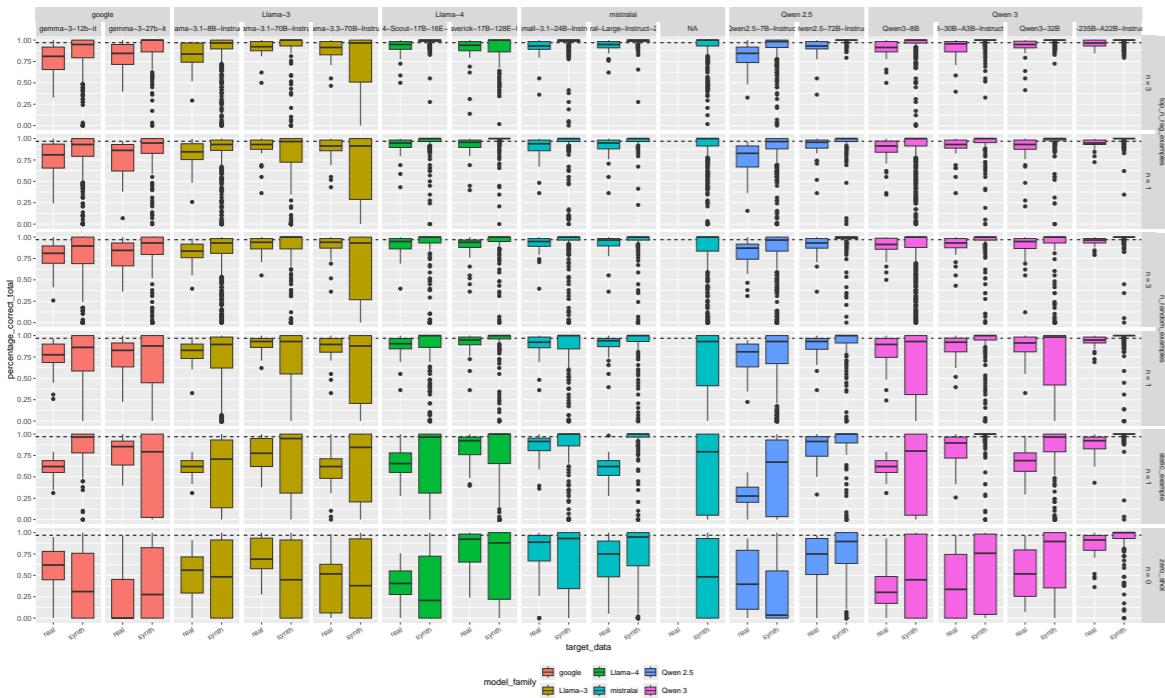


Figure A.14: Comparing the table extraction performance among real and synthetic Aktiva tables

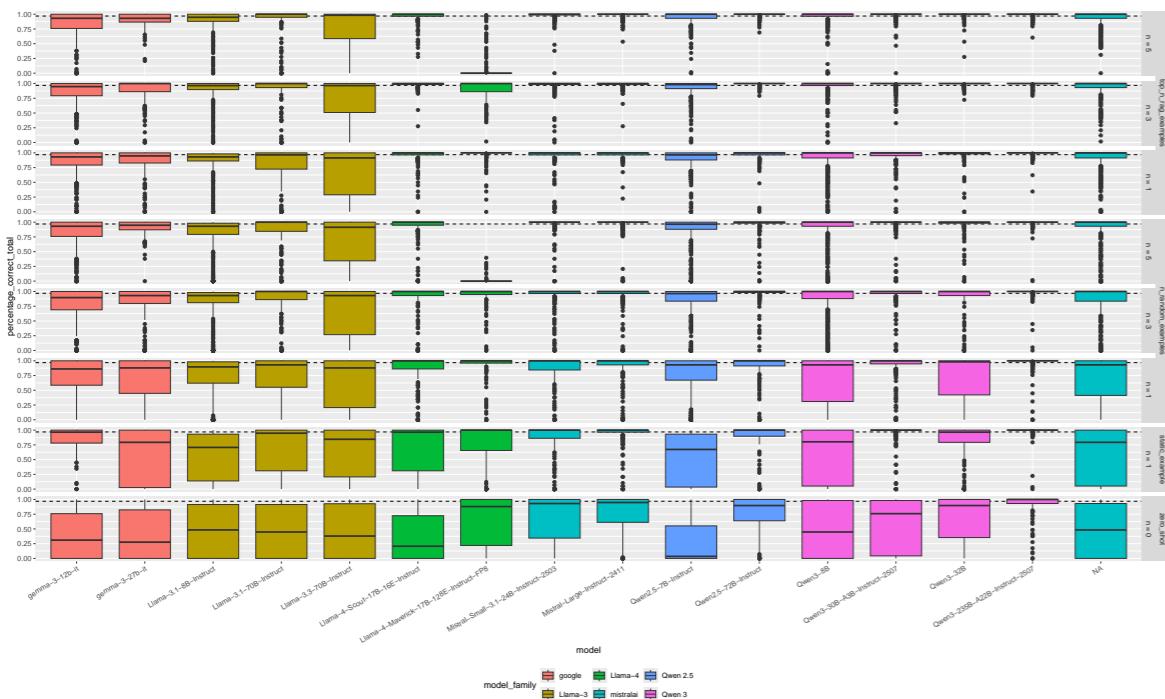


Figure A.15: Percentage of correct extracted or as missing categorized values for table extraction task on synthetic Aktiva tables

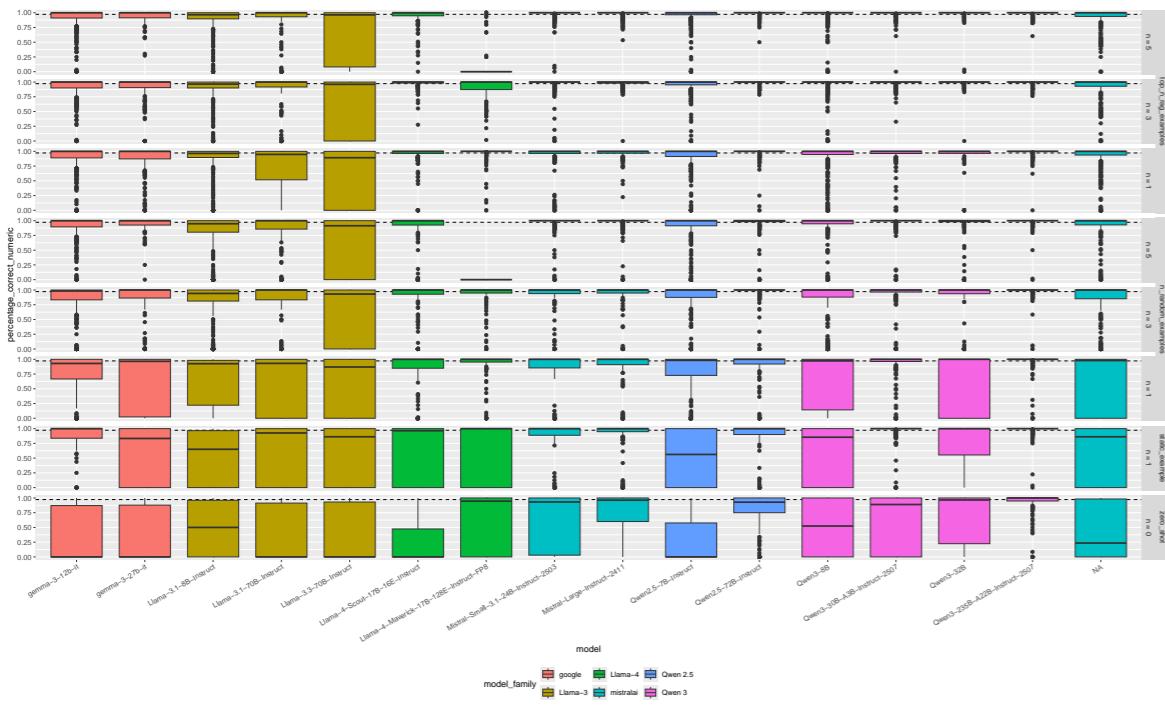


Figure A.16: Percentage of correct extracted numeric values for table extraction task on synthetic Aktiva tables

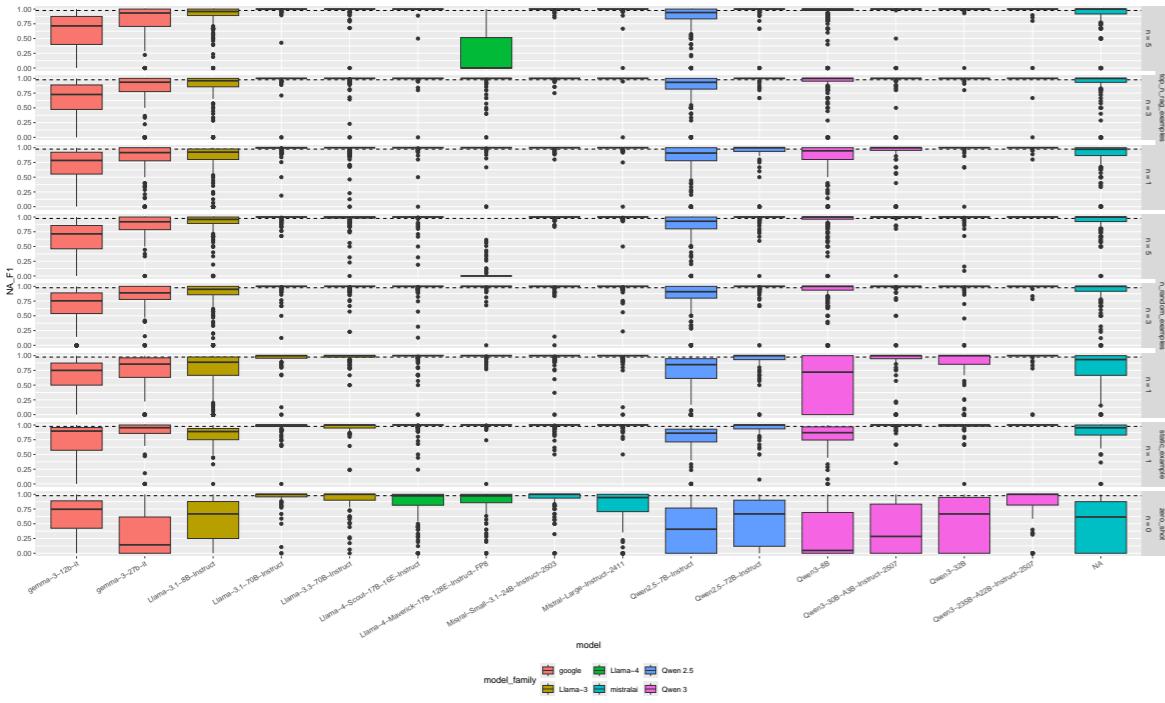


Figure A.17: F1 score for the missing classification if a value is missing for table extraction task on synthetic Aktiva tables

12sdgsdgj

Landscape

A.6 Annual Comprehensive Financial Report Balance Sheet

A.7 Extraction framework flow chart

A.8 Table extraction with regular expressions

Extract by pdfium for ‘..../benchmark_truth/synthetic_tables/separate_files/final/aktiva_table__3_columns__span_False__thin_Fa€_enumeration_False__shuffle_True__text_around_True__max_length_50__sum_in_same_row_False__0.pdf’:

A

ktiva (inMio. €) Geschäftsjahr Vorjahr

Anlagevermögen Immaterielle Verm

ögensgegenstände

Selbstgeschaffene gewerbliche Schutzrechte und

ähnliche Rechte und Werte

0,184,77

Geschäfts- oder Firmenwert 4,426,78

geleistete Anzahlungen 1,780,65

entgeltlicher erworbene Konzessionen, gewerbliche

Schutzrechte und ähnliche Rechte und Wertes sowie

Lizenzen an solchen Rechten und Werten

4,646,71

11,0218,91

Sachanlagen

Grundstücke, grundstücksgleiche Rechte und Bauten

einschließlich der Bauten auf fremden Grundstücken

2,802,55

Technische Anlagen und Maschinen 5,205,53

Andere Anlagen, Betriebs- und Geschäftsausstattung 1,601,93

geleistete Anzahlungen und Anlagen im Bau 3,255,81

12,8615,83

*State of California Annual Comprehensive Financial Report***Balance Sheet****Governmental Funds****June 30, 2023**

(amounts in thousands)

	General	Federal
ASSETS		
Cash and pooled investments.....	\$ 71,968,861	\$ 6,986,275
Investments.....	—	—
Receivables (net).....	46,621,774	2,076,598
Due from other funds.....	6,933,803	165,231
Due from other governments.....	4,075,837	37,069,188
Interfund receivables.....	3,914,413	—
Loans receivable.....	45,225	384,293
Other assets.....	6,244	601,252
Total assets	\$ 133,566,157	\$ 47,282,837
LIABILITIES		
Accounts payable.....	\$ 14,422,777	\$ 24,499,200
Due to other funds.....	3,911,973	3,865,533
Due to component units.....	264,995	—
Due to other governments.....	21,808,112	11,125,464
Interfund payables.....	2,692,941	—
Benefits payable.....	—	69,623
Revenues received in advance.....	25,891	6,675,956
Tax overpayments.....	21,740,974	—
Deposits.....	4,231	—
Unclaimed property liability.....	1,314,797	—
Other liabilities.....	522,844	46,256,400
Total liabilities	66,709,535	92,492,176
DEFERRED INFLOWS OF RESOURCES		
Total liabilities and deferred inflows of resources	69,562,469	92,502,885
FUND BALANCES		
Nonspendable.....	3,950,919	—
Restricted.....	24,830,454	1,210,267
Committed.....	4,210,891	—
Assigned.....	20,714,283	—
Unassigned.....	10,297,141	(46,430,315)
Total fund balances (deficit)	64,003,688	(45,220,048)
Total liabilities, deferred inflows of resources, and fund balances	\$ 133,566,157	\$ 47,282,837

Figure A.18: Example balance sheet pagefom Californias Annual Comprehensive Financial Report 2023

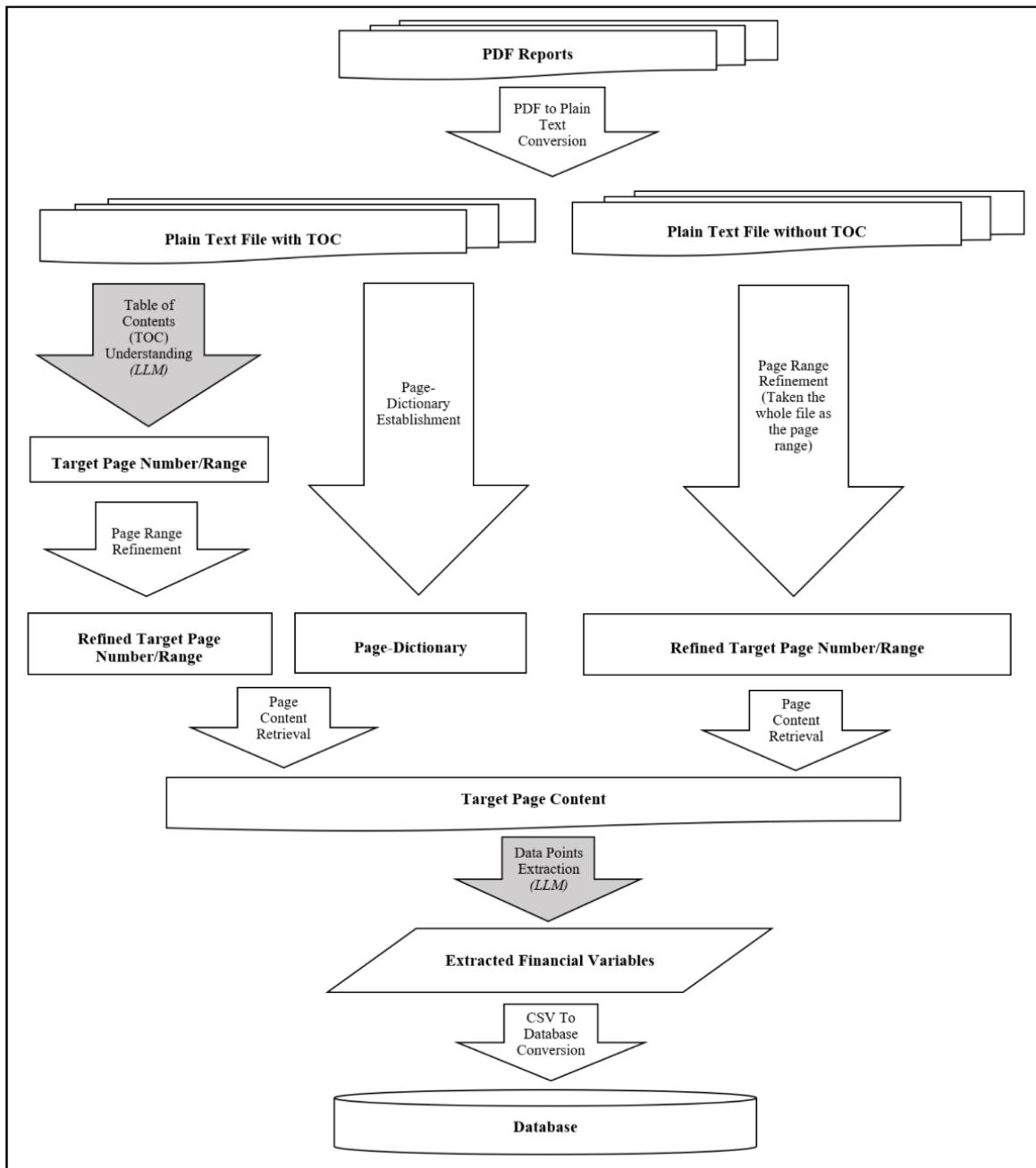


Figure A.19: Flowchart of the extraction framework

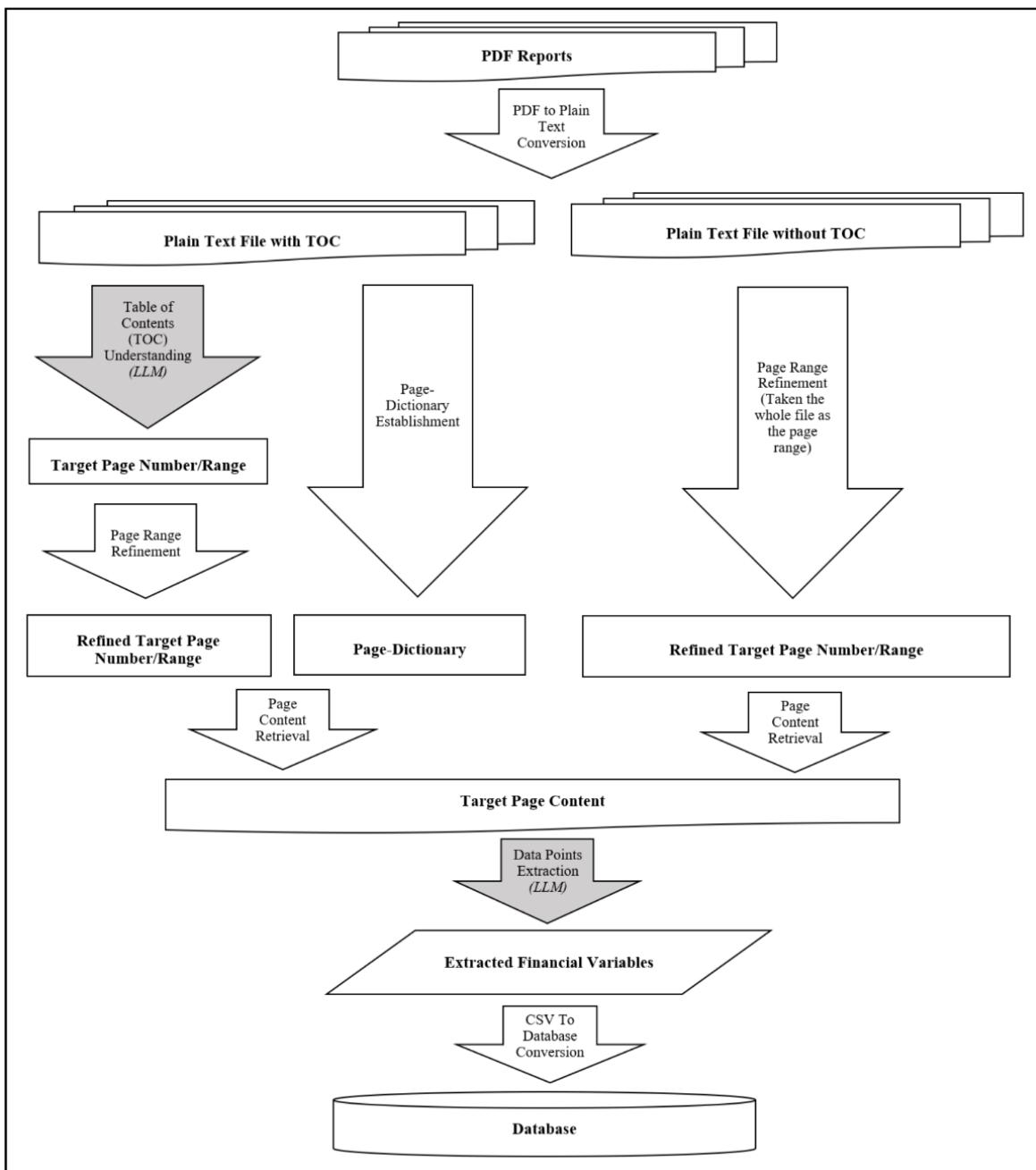


Figure A.20: Flowchart of the extraction framework of Li et al. (2023)

Finanzanlagen

SonstigeFinanzanlagen7,446,51

AnteileanverbundenenUnternehmen0,499,83

AusleihungenanverbundeneUnternehmen0,573,49

Beteiligungen1,059,43

AusleihungenanUnternehmen, mitdenenein

Beteiligungsverhältnisbesteht

6,957,65

WertpapieredesAnlagevermögens2,002,71

SonstigeAusleihungen9,091,52

27,5841,13

51,4675,87

Umlaufvermögen**Vorräte**

Roh-, Hilfs- und Betriebsstoffe0,382,98

UnfertigeErzeugnisse, unfertigeLeistungen3,236,19

FertigeErzeugnisseundWaren6,724,98

GeleisteteAnzahlungen4,024,83

14,3418,98

Forderungen und sonstige Vermögensgegenstände

ForderungenausLieferungenundLeistungen4,328,36

ForderungengegenverbundeneUnternehmen6,082,38

ForderungengegenUnternehmen, mitdenenein

Beteiligungsverhältnisbesteht

7,878,11

SonstigeVermögensgegenstände1,968,30

20,2227,15

Wertpapiere
Anteile an verbundenen Unternehmen 2,383,24
Sonstige Wertpapiere 0,077,65
2,4410,88
Kassenbestand, Bundesbankguthaben, Guthaben bei Kreditinstituten und Schecks 4,144,00
41,1561,01
Rechnungsabgrenzungsposten 2,746,78
Aktive latente Steuern 8,464,60
Aktiver Unterschiedsbetrag aus der Vermögensverrechnung
2,863,35
106,67151,61

Extract by pdfminer for ‘..../benchmark_truth/synthetic_tables/separate_files/final/aktiva_table__3_columns__span_False__the__enumeration_False__shuffle_True__text_around_True__max_length_50__sum_in_same_row_False__0.pdf’:

Aktiva (in Mio. €)
Anlagevermögen
Immaterielle Vermögensgegenstände
Selbst geschaffene gewerbliche Schutzrechte und ähnliche Rechte und Werte
Geschäfts- oder Firmenwert
geleistete Anzahlungen
entgeltlich erworbene Konzessionen, gewerbliche Schutzrechte und ähnliche Rechte und Werte sowie Lizenzen an solchen Rechten und Werten
Sachanlagen
Grundstücke, grundstücksgleiche Rechte und Bauten einschließlich der Bauten auf fremden Grundstücken
Technische Anlagen und Maschinen

Andere Anlagen, Betriebs- und Geschäftsausstattung

geleistete Anzahlungen und Anlagen im Bau

Finanzanlagen

Sonstige Finanzanlagen

Anteile an verbundenen Unternehmen

Ausleihungen an verbundene Unternehmen

Beteiligungen

Ausleihungen an Unternehmen, mit denen ein
Beteiligungsverhältnis besteht

Wertpapiere des Anlagevermögens

Sonstige Ausleihungen

Umlaufvermögen

Vorräte

Roh-, Hilfs- und Betriebsstoffe

Unfertige Erzeugnisse, unfertige Leistungen

Fertige Erzeugnisse und Waren

Geleistete Anzahlungen

Forderungen und sonstige Vermögensgegenstände

Forderungen aus Lieferungen und Leistungen

Forderungen gegen verbundene Unternehmen

Forderungen gegen Unternehmen, mit denen ein
Beteiligungsverhältnis besteht

Sonstige Vermögensgegenstände

Wertpapiere

Anteile an verbundenen Unternehmen

Sonstige Wertpapiere

Kassenbestand, Bundesbankguthaben, Guthaben bei
Kreditinstituten und Schecks

Rechnungsabgrenzungsposten

Aktive latente Steuern
Aktiver Unterschiedsbetrag aus der Vermögensverrechnung
Geschäftsjahr
Vorjahr
0,18
4,42
1,78
4,64
11,02
2,80
5,20
1,60
3,25
12,86
7,44
0,49
0,57
1,05
6,95
2,00
9,09
27,58
51,46
0,38
3,23
6,72

4, 02

14, 34

4, 32

6, 08

7, 87

1, 96

20, 22

2, 38

0, 07

2, 44

4, 14

41, 15

2, 74

8, 46

2, 86

4, 77

6, 78

0, 65

6, 71

18, 91

2, 55

5, 53

1, 93

5, 81

15, 83

6, 51

9, 83

3,49

9,43

7,65

2,71

1,52

41,13

75,87

2,98

6,19

4,98

4,83

18,98

8,36

2,38

8,11

8,30

27,15

3,24

7,65

10,88

4,00

61,01

6,78

4,60

3,35

106,67

151,61