

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

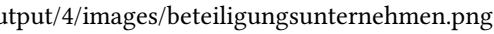
Gegenstand dieser Arbeit ist ein Benchmark zur Informationsextraktion aus PDF-Dateien. Dabei wird sich auf das Auslesen der Bilanzen und Gewinn- und Verlustrechnungen aus Jahresabschlüssen deutscher Unternehmen beschränkt. Ein besonderer Aspekt der Aufgabe ist die Berücksichtigung der Hierarchie innerhalb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

Notes

- Qwen 2.5 hat zweiseitige GuV von IBB entdeckt und zur Anpassung der Ground Truth
- Google gemma war mit alter Klassifikation erfolgreich (anderer Prompt, mehr Seiten)

implementation nach methods

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Objectives	2
1.3	Methodology (1 p)	2
1.4	Thesis Outline (0.5 p)	2
1.5	To place in chapters above  2	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	Basic terms	5
2.2	Technological topic (related work)	5
2.2.1	Extraction of numeric values	5
2.3	optinal more topics like previous	5
2.4	Summary (0.5 p)	5
2.5	To place in chapters above	6
2.6	Table extraction tasks	6
2.6.1	Difficulties	6
2.7	Document Extrtaction Process	6
2.7.1	Document Layout Analysis	6
2.7.2	6
2.8	Tools	6
2.8.1	TableFormer	6
3	Implementation (max 5p)	7
3.1	Speedup with vLLM and batching	7
3.2	Setup (Dockerfile and PV)	7

4	Methods	9
4.1	Data	9
4.2	Page identification	9
4.2.1	Baselines	9
4.3	Table detection	10
4.3.1	LLM	10
4.3.2	Vision Model	10
4.3.3	Docling and Co	10
4.4	Information extraction	10
4.4.1	Baselines	10
4.4.2	Simple pipeline	10
4.4.3	Sophisticated approaches	11
5	Results	13
5.1	Page identification	13
5.1.1	Baseline: Regex	13
5.1.2	Advanced techniques	15
5.1.3	Comparison	19
5.2	Table extraction	20
6	Discussion	21
6.1	Not covered	21
7	Conclusion	23
	References	25
A	Appendix	27
A.1	Local machine	27
	System Details Report	27
A.2	Benchmarks	28
A.2.1	Text extraction	28
A.2.2	Table detection	28
A.2.3	Large language model process speed	30
A.2.4	Table identification with LLMs	31
A.3	Regular expressions	31
A.4	Extraction framework flow chart	32

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
BSB Unternehmensverwaltung Gewährträger: Berlin	degewo AG 100%	Berlinwo Immobilien Ges. mbH 100%	Amt für Statistik Berlin-Brandenburg Gewährträger: Bln. u. Brandenburg	BEN Berlin Energie und Netzholding GmbH 100%	BSB Infrastruktur. Verw. GmbH 100%	Di. Film- u. Fernsehakad. GmbH 100%	Berliner Werkstatt Bth. GmbH 70%
	GESOBAU AG 100%	BIM GmbH 100%	BEHALA GmbH 100%	Berl. Stadtreinigungsbetriebe Gewährträger: Berlin	BSB Infrastruktur. GmbH & Co. KG 100 % Kommanditist: Berlin	Deutsches Zentrum f. Hochschul- u. Wiss.forschung GmbH 1,85%	Vivantes GmbH 100%
	Gewobag AG 96,69%	Berliner Stadtgüter 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 90,1%	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%	Heinrich-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	BOZ GmbH 60%		Kulturprojekte Berlin GmbH 100%		
		Liegenschaftsfonds Projekt KG 100 % Kommanditist: Berlin	DEGES Dt. Einheit Fernstraßenplanungs- u. -bau GmbH 5,91%		Kunsthalle BR Deutschld. GmbH 2,44%		
		Olympiastadion Berlin GmbH 100%	Deutsche Klassenlotterie Gewährträger: Berlin		Musicboard Berlin GmbH 100%		
		Tegel Projekt GmbH 100%	Flughafen Berlin-Brandb. GmbH 37%		Rundfunk-Orchester gGmbH 20%		
		Tempelhof Projekt GmbH 100%	IT-Dienstleistungszentrum Berlin Gewährträger: Berlin		Zoologischer Garten Berlin AG 0,03%		
		WISTA Management GmbH 100%	Landesanal. Schienenfahrzeuge Berlin Gewährträger: Berlin				
			Messe Berlin GmbH 100%				
			Partner für Deutschland 1%				
			VBB GmbH 33,33%				

Figure 1.1: Companies Berlin has holds share at

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%	Berlinero Immobilien Ges. mbH 100%	Amf für Statistik Berlin-Brandenbg. Gewährträger: Bln. u. Brandenbg.	BEN Berlin Energie und Netz- holding GmbH 100%	BBB Infrastrukt. Verw. GmbH 100%	Di. Film- u. Fernsehakad. GmbH 100%	Berliner Werkst. f. Beh. GmbH 70%
	GESOBALU AG 100%	BIM GmbH 100%	BEHALA GmbH 100%	Berl. Stadtrangungsbehörden Gewährträger: Berlin	BBB Infrastrukt. GmbH & Co. KG 100 % Kommanditist: Berlin	Deutsches Zentrum f. Hochschul- u. Wissenschaft GmbH 1,85%	Vivantes GmbH 100%
	Gewobag AG 96,69%	Berliner Stadtgüter GmbH 100%	Berlin Tourismus & Kongress GmbH 15%	Berliner Wasserbetriebe Gewährträger: Berlin	Berliner Böder-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%	Heinhold-Zentrum Bln. GmbH 10%	
	WBM GmbH 100%	Liegenschaftsfonds GmbH 100%	Berliner Verkehrsbetriebe Gewährträger: Berlin	SBB Sonderabfall GmbH 25%	KuJ Wahlheide gGmbH 100%	Wissenschaftszentrum gGmbH 25%	
		Liegenschaftsfonds KG 100 % Kommanditist: Berlin	BOZ GmbH 60%		Kulturprojekte Berlin GmbH 100%		
		Liegenschaftsfonds Projekt KG 100 % Kommanditist: Berlin	DEGES Dt. Einzel Fernstraßen- planungs- u. -bau GmbH 5,91%		Kunsthalle BR Deutschld. GmbH 2,44%		
		Olympiastadion Berlin GmbH 100%	Deutsche Klassenlotterie Gewährträger: Berlin		Musicboard Berlin GmbH 100%		
		Tegel Projekt GmbH 100%	Flughafen Berlin-Brandb. 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%	Landesanst. Schienenfahrzeuge Berlin Gewährträger: Berlin				
			Messe Berlin GmbH 100%				
			Partner für Deutschland 1%				
			VBB GmbH 33,33%				

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?

1.3 Methodology (1 p)

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

1.4 Thesis Outline (0.5 p)

1.5 To place in chapters above http://127.0.0.1:29003/rmd_output/4/images/bete

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 Rechnungshof 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

2.1 Basic terms

2.2 Technological topic (related work)

- most important papers
- connection of papers (timeline)
- what used, what not?
- extending existing paper?

2.2.1 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)

2.3 optimal more topics like previous

2.4 Summary (0.5 p)

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

2.5 To place in chapters above

2.6 Table extraction tasks

2.6.1 Difficulties

- Beispielbilder

2.7 Document Extrtaction Process

2.7.1 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 paperswithcode.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

2.7.1.1 Vision Grid Transformer

2.7.2

2.8 Tools

2.8.1 TableFormer

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

Chapter 3

Implementation (max 5p)

3.1 Speedup with vLLM and batching

3.2 Setup (Dockerfile and PV)

Chapter 4

Methods

4.1 Data

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

4.2 Page identification

The first task to solve, for a fully autonomous solution, is to identify the pages where the tables of interest are located. For benchmarking 74 annual reports from 7 companies have been used. For this benchmark we limit the tables of interest to those that show **Aktiva**, **Passiva** and **Gewinn- und Verlustrechnung**.

In those documents there are 252 pages of interest holding 265 relevant tables. On 13 pages there have been two tables (**Aktiva** and **Passiva**) on a single page. 21 tables are spread over two pages. In 8 documents there have been multiple tables per type of interest, distributed among the three types of tables as following:

type	count
Aktiva	7
GuV	8
Passiva	7

As a baseline a simple regex approach was used.

4.2.1 Baselines

4.2.1.1 Regex based

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

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

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

results depend 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

4.2.1.1.1 LLM based

4.2.1.2 Term frequency based

4.2.1.2.1 VLLM based was not implemented

4.3 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

4.3.1 LLM

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

4.3.2 Vision Model

Yolo

4.3.3 Docling and Co

4.3.3.0.1 VLLM based was not implemented

4.4 Information extraction

4.4.1 Baselines

simple regex?

4.4.2 Simple pipeline

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

4.4.3 Sophisticated approaches

not implemented

- with pipelines
- Nougat
- maker
- Azure
- docling

Chapter 5

Results

5.1 Page identification

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

5.1.1 Baseline: Regex

Building a sound regular expression often is an iterative process. In a first approach a very simple one was implemented.

Comparing the differences in the metrics based on the different text extraction libraries it can be said that the extracted text is very similar but not identical. Since the results are not depending on the used text extraction library the *exhaustive regex restricted* has only been run with the fast text extraction library *pdfium*. The results of the regex based page identification are presented in the following tables.

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

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.

The regular expressions can be found in the appendix (see 5.1.1).

General bad precision. Increasing recall degrades precision even further. number of pages positive identified total; used as subset for table identification task

Table 5.1: Comparing page identification metrics for different regular expressions for classification task 'Aktiva'

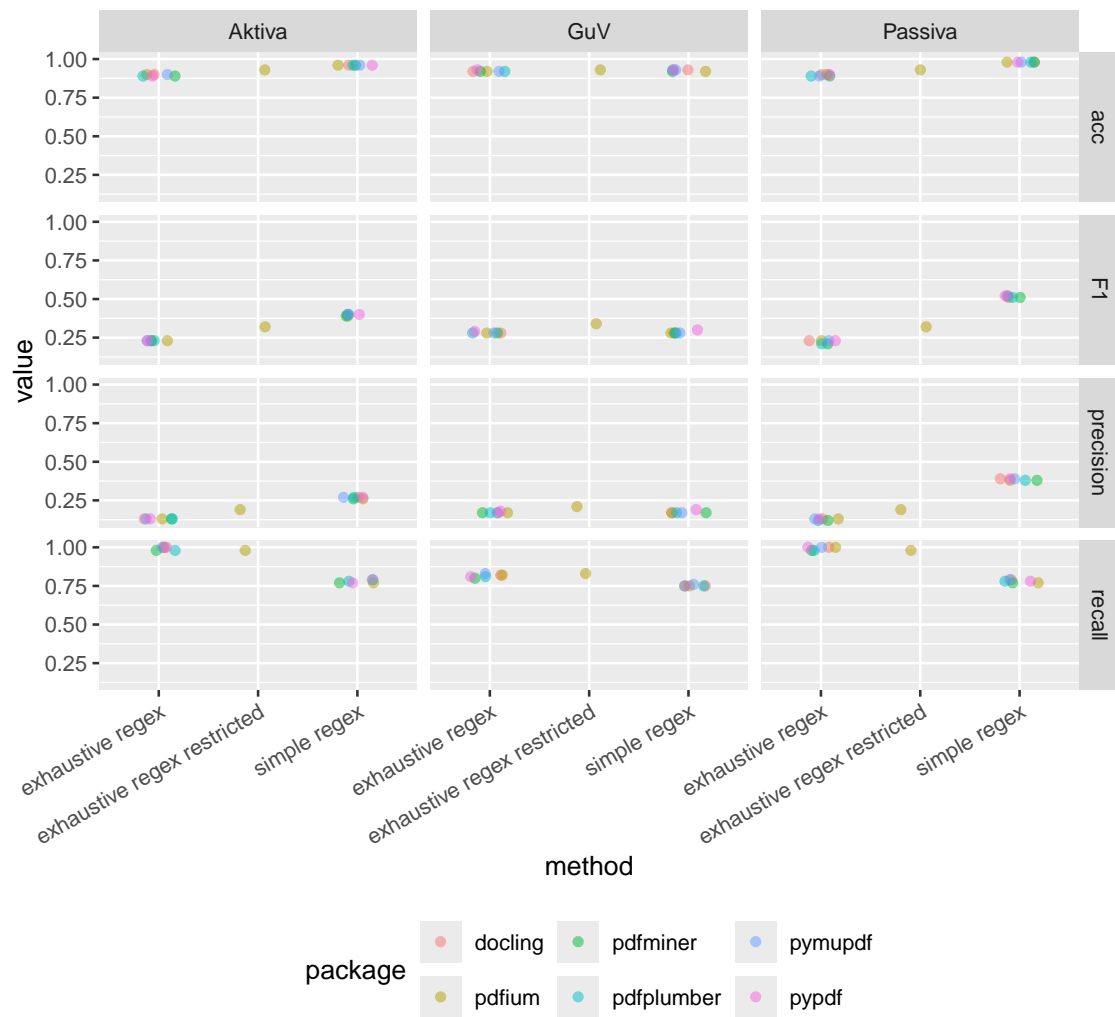
method	stat	precision	recall	F1
simple regex	mean	{0.267}	0.778	{0.397}
simple regex	sd	0.005	0.01	0.005
exhaustive regex restricted	mean	0.19	0.98	0.32
exhaustive regex restricted	sd	NA	NA	NA
exhaustive regex	mean	0.13	{0.993}	0.23
exhaustive regex	sd	0	0.01	0

Table 5.2: Comparing page identification metrics for different regular expressions for classification task 'Passiva'

method	stat	precision	recall	F1
simple regex	mean	{0.385}	0.78	{0.515}
simple regex	sd	0.005	0.009	0.005
exhaustive regex restricted	mean	0.19	0.98	0.32
exhaustive regex restricted	sd	NA	NA	NA
exhaustive regex	mean	0.127	{0.993}	0.223
exhaustive regex	sd	0.005	0.01	0.01

Table 5.3: Comparing page identification metrics for different regular expressions for classification task 'Gewinn und Verlustrechnung'

method	stat	precision	recall	F1
simple regex	mean	0.173	0.752	0.283
simple regex	sd	0.008	0.004	0.008
exhaustive regex restricted	mean	{0.21}	{0.83}	{0.34}
exhaustive regex restricted	sd	NA	NA	NA
exhaustive regex	mean	0.172	0.815	0.282
exhaustive regex	sd	0.004	0.01	0.004



5.1.2 Advanced techniques

5.1.2.1 Table of Contents understanding

Joining with 'by = join_by(filepath)'

- calculate and add Qwen, Gemini or LLama results?

5.1.2.1.1 Text based Li et al. (2023) used the table of contents 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 that “ACFRs typically spans no more than the initial 165 lines of the converted document” (p. 20), they use the first 200 lines of text.

My expectation was to find the TOC within the first five pages. Often we find way less than 200 lines of text on the five first pages (see Figure 5.1). Some files are not machine readable without OCR and thus show zero lines in the first five pages as well.

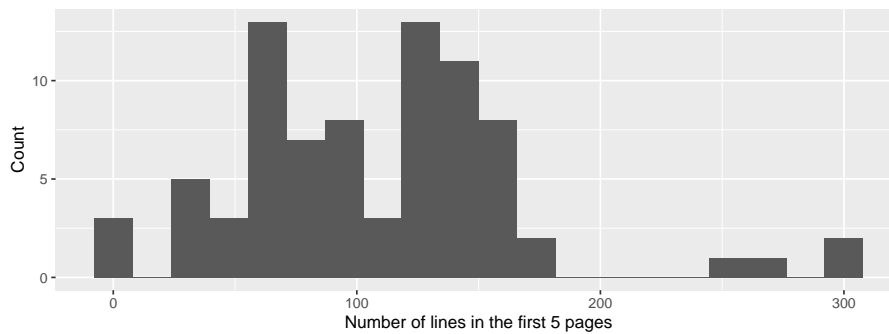


Figure 5.1: Histogram of the number of lines in the first 5 pages of the annual reports

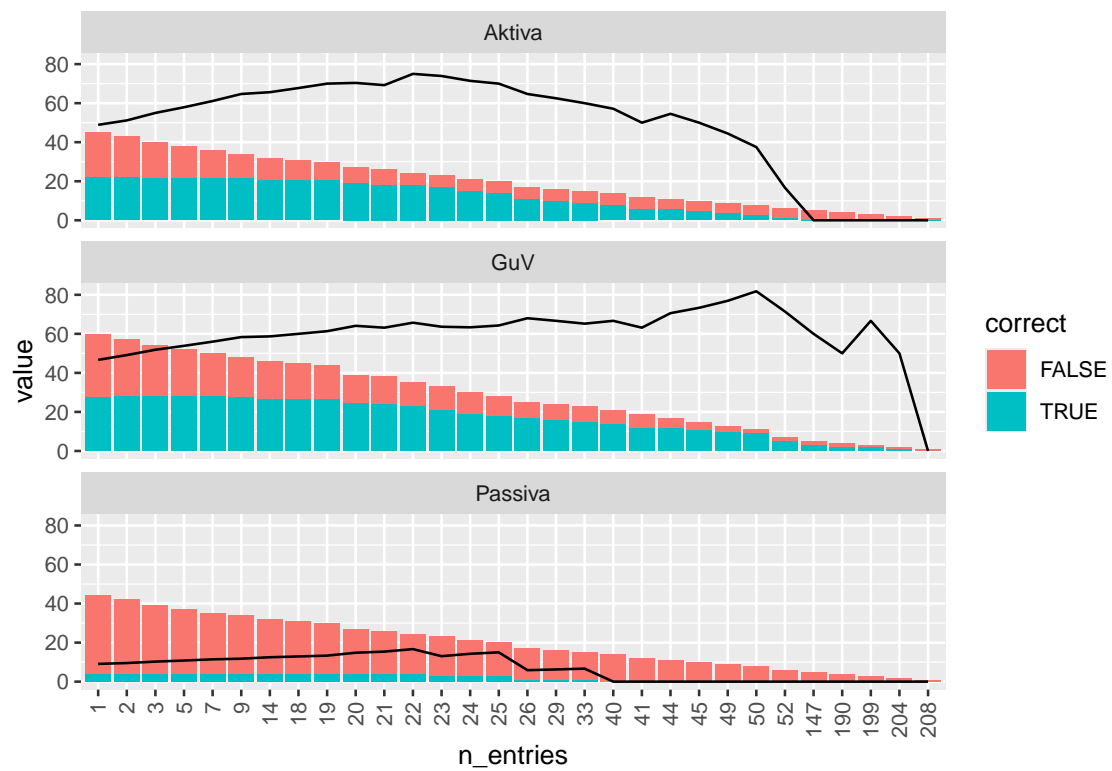
5.1.2.1.1.1 First five pages A request to Mistral results in 64 strings that should represent a table of contents among the first five pages [strings not checked in detail].

5.1.2.1.1.2 First 200 lines A request to Mistral results in 70 strings that should represent a table of contents among the first five pages [strings not checked in detail].

5.1.2.1.1.3 Machine readable TOC based To limit the text and hopefully increase the quality of the input data one can work with the TOC representation embedded within the PDF files. From 80 annual reports 43 files do have a machine readable separate table of contents and 37 do not have one.

One can see that correct predictions for the page range are more probable when the TOC has a medium number of entries. It is possible to drop PDFs with less than 9 without losing a single correct prediction. This means that for PDFs with TOC with less than 9 entries the LLM was not able to make a correct prediction. This is not surprising since neither *Bilanz* nor *Gewinn- und Verlustrechnung* are mentioned there.

Almost no influence if TOC is passed formatted as markdown or json. With the json formatted TOC it found two more correct page ranges (single test run). It was tested because the relation *page_number* heading and value might have been clearer in json for a linear working LLM.



5.1.2.1.2 Comparison of the different approaches

- toc analysis
- cleaned measures

```
## Warning: Removed 151 rows containing non-finite outside the scale range
## ('stat_smooth()').
```

```
## Warning: Removed 151 rows containing missing values or values outside the scale range
## ('geom_point()').
```

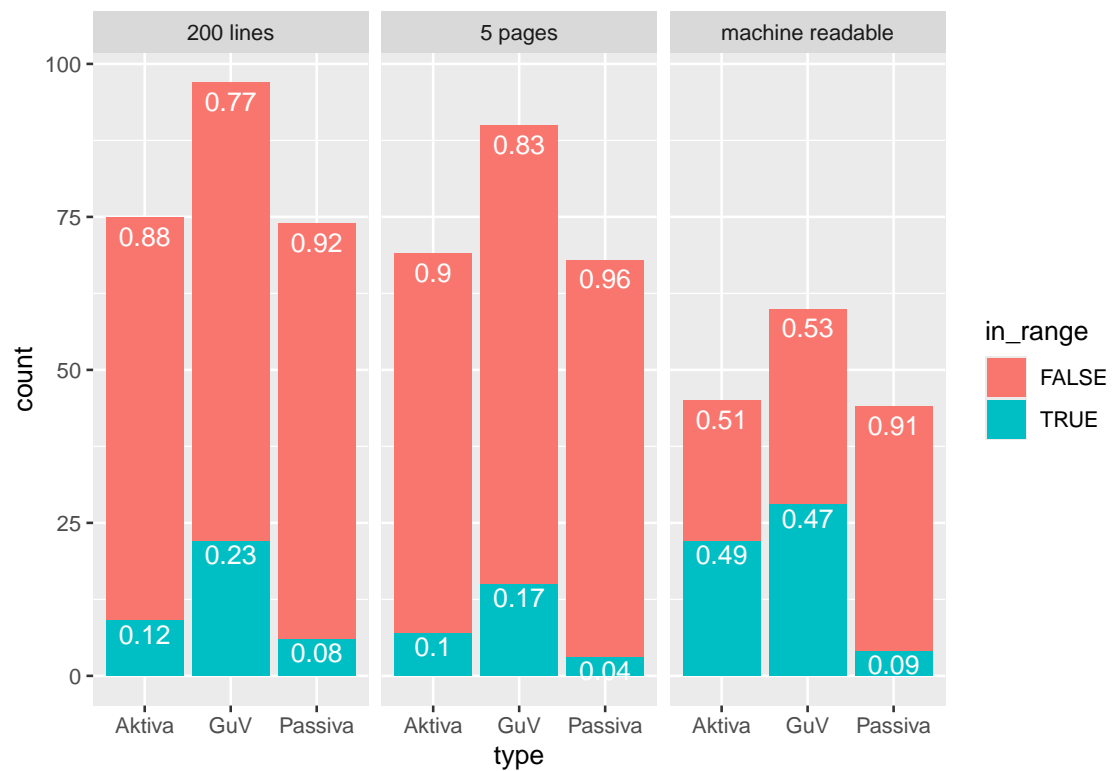
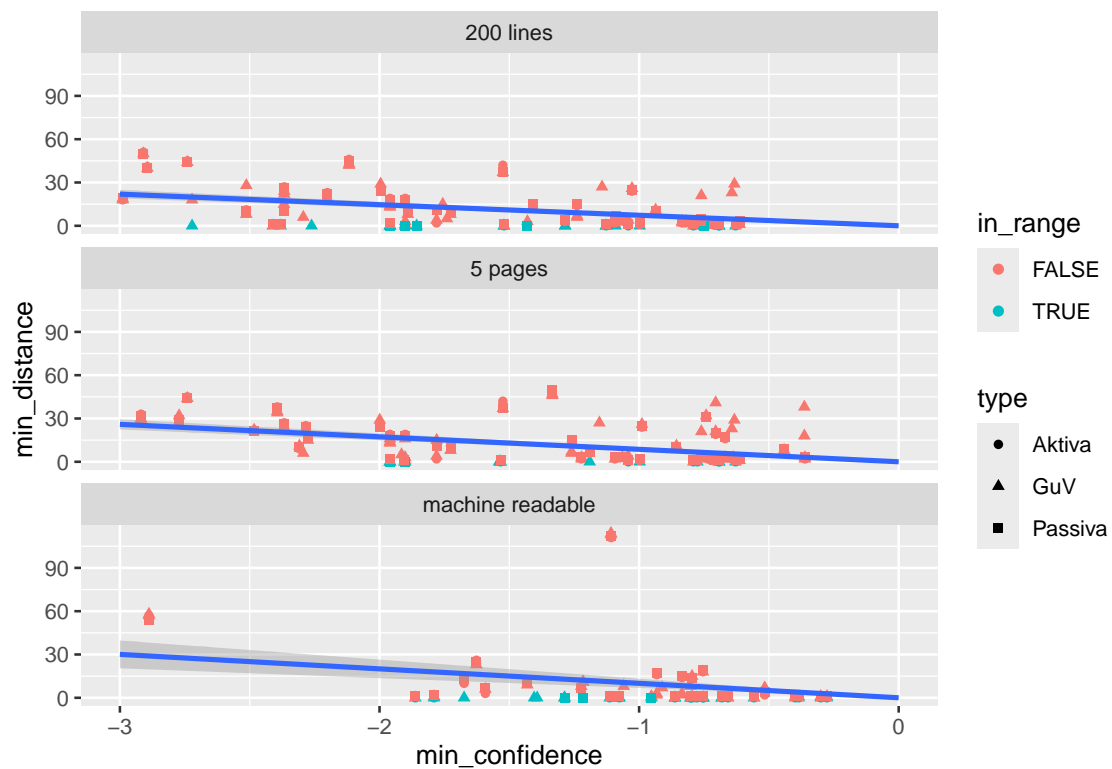


Figure 5.2: Comparing number of found TOC and amount of correct and incorrect predicted page ranges



model_family	model	classification_type	method_family	n_examples	f1_score
mistralai	mistralai_Ministral8BInstruct2410	GuV	n_rag_examples	3	0.93
meta-llama	metallama_Llama4Scout17B16EInstruct	GuV	n_rag_examples	3	0.92
mistralai	mistralai_Ministral8BInstruct2410	Passiva	n_rag_examples	3	0.92
mistralai	mistralai_Ministral8BInstruct2410	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
meta-llama	metallama_Llama4Scout17B16EInstruct	Aktiva	n_rag_examples	3	0.83
Qwen	Qwen_Qwen2.532BInstruct	Passiva	n_rag_examples	1	0.79
google	google_gemma3nE4Bit	GuV	top_n_rag_examples	1	0.23
google	google_gemma3nE4Bit	Passiva	zero_shot	NA	0.21
google	google_gemma312bit	Aktiva	top_n_rag_examples	1	0.18

5.1.2.2 Classification with LLMs

```
binary_task <- list()
binary_task$n_models <- df_binary$model %>% unique() %>% length()
binary_task$n_model_families <- df_binary$model_family %>% unique() %>% length()
binary_task$n_method_families <- df_binary$method_family %>% unique() %>% length()
```

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

5.1.2.2.1 Binary classification 20 models from 4 have been benchmarked among 5 methods

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

[Probably gonna drop because models will have rerun in some hours] Some models only ran for the *Aktiva* classification task and crashed because of a conflict trying to access the vector database. If their scores were not promising they have not been run for the other classification tasks at all.

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. This is surprising since they did a decent job in a similar task as described in section @ref(#llm-table-detection).

```
## 'mutate_all()' ignored the following grouping variables:
## 'mutate_all()' ignored the following grouping variables:
## * Columns 'model_family', 'classification_type'
## i Use 'mutate_at(df, vars(-group_cols()), myoperation)' to silence the message.
```

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

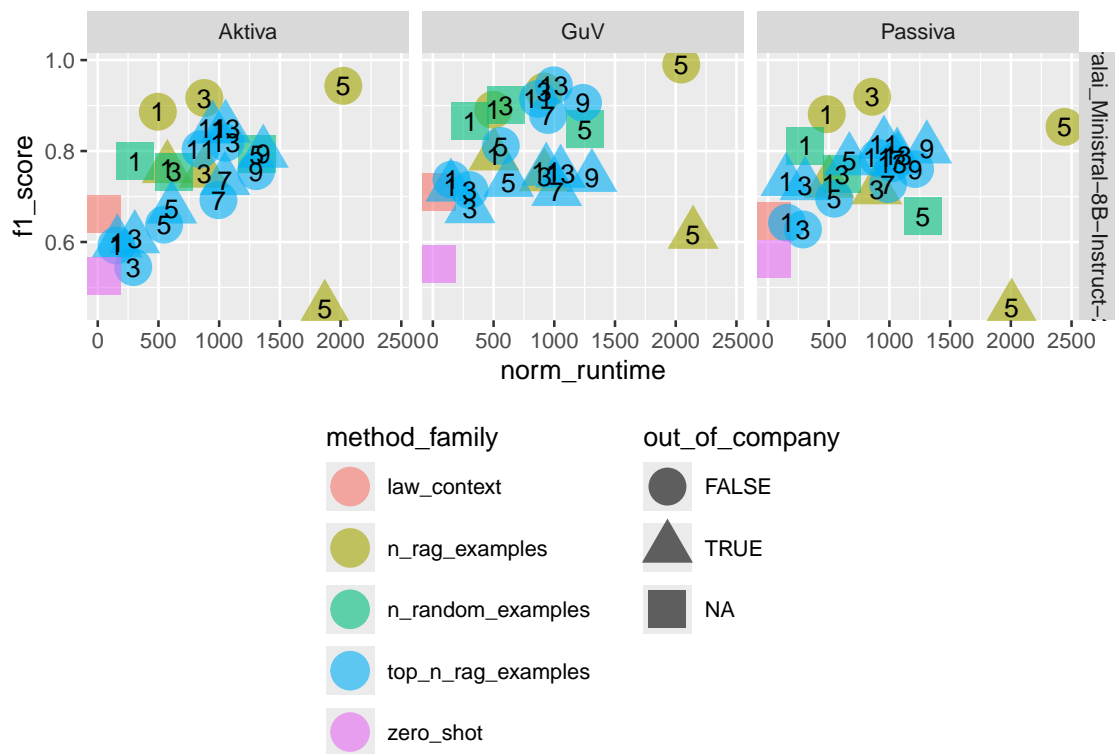
The experiments for best performing model, Ministral-8B-Instruct-2410, have been extended by methods with even more examples in the context.

```
df_binary %>% filter(model == "mistralai_Ministral-8B-Instruct-2410", loop < 1) %>%
  ggplot(aes(x = norm_runtime, y = f1_score)) +
  geom_point(aes(color = method_family, shape = out_of_company), size = 7, alpha = .6) +
  scale_shape(na.value = 15, guide = "legend") +
  geom_text(aes(label = n_examples)) +
  facet_grid(model~classification_type) +
```



```
theme(legend.position = "bottom") +
guides(
  color = guide_legend(ncol = 1, title.position = "top"),
  shape = guide_legend(ncol = 1, title.position = "top")
)
```

```
## Warning: Removed 6 rows containing missing values or values outside the scale range
## ('geom_text()').
```



Some models have run for multiple times to check if the results are stable. Earlier examples with supset have been run five times indicating stable results. Running the experiments up to three times in this very task indicate this as well.

5.1.2.2.2 Multi classification

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

5.1.2.3 Term frequency based classifier

- top 1
- top k

5.1.3 Comparison

Multiclassification more effective than three times single classification

5.1.3.1 F1

5.1.3.2 Energy usage and runtime

5.2 Table extraction

Chapter 6

Discussion

6.1 Not covered

- OCR

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>
-

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 libraries

library	runtime in s
pdfium	{14}
pymupdf	22
pypdf	218
pdfplumber	675
pdfminer	752
doclingparse	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 formatted 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

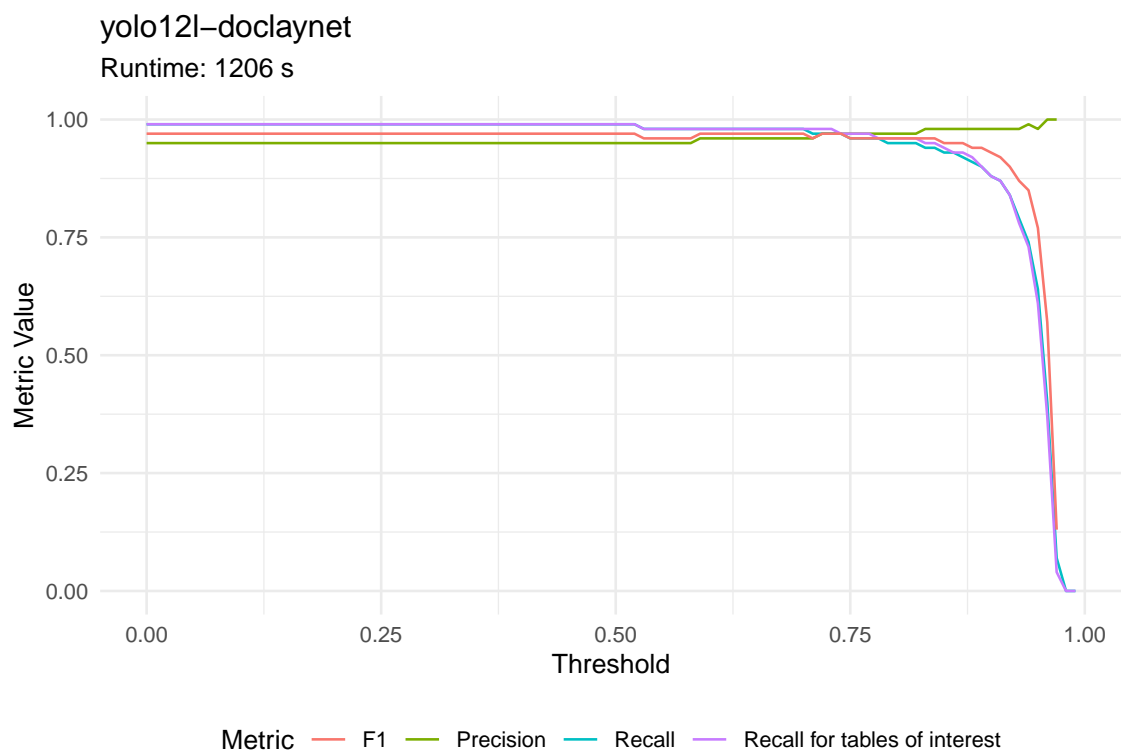
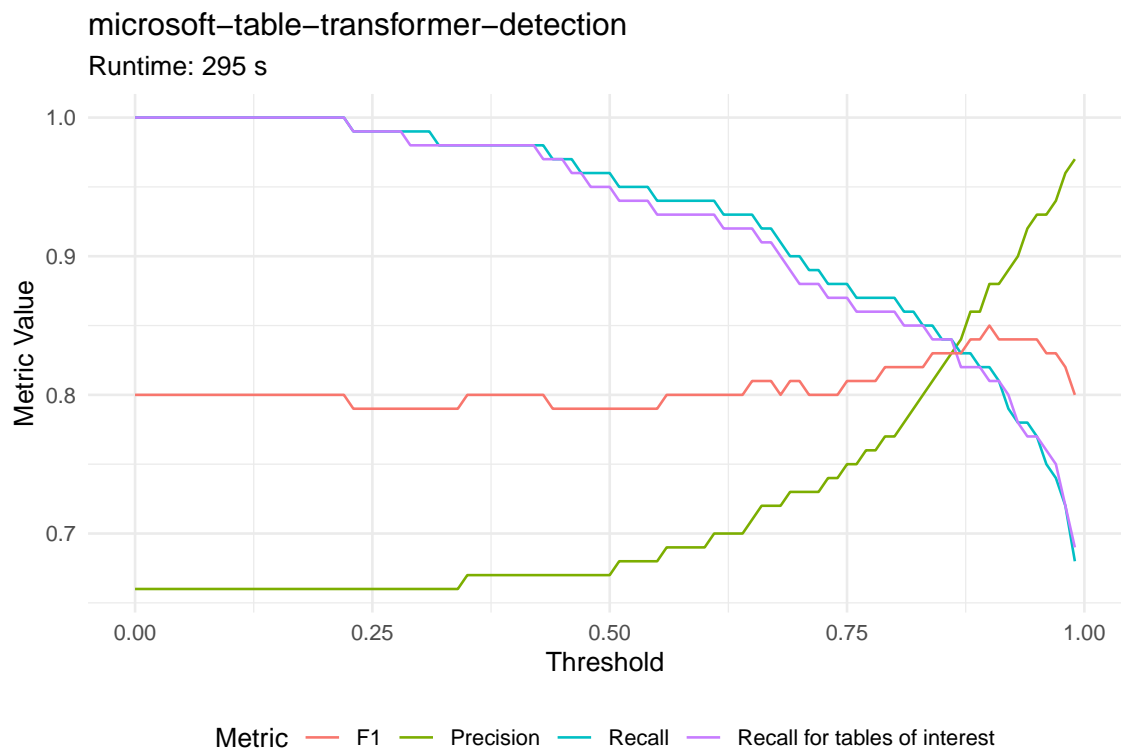
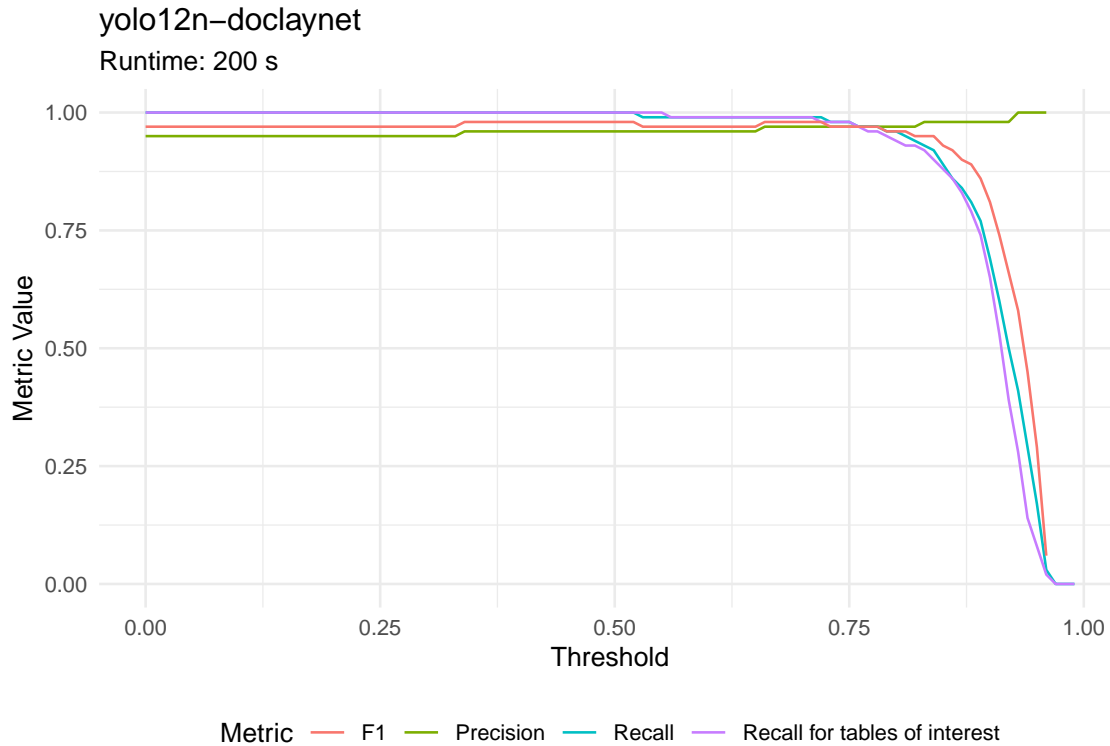


Table A.2: Comparing time (in seconds) for processing ten asset tables using different libraries and approaches

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 5.1.2.2). Previous experiments have only been using a subset of pages that have been selected with the baseline regex approach (see section 5.1.1). One can find the former results in section A.2.4.

A.2.4 Table identification with LLMs

A.3 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"aktiva",
        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\s*i\s*v\s*a|p\s*a\s*s\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",
        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",
        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*s\s*i\s*v\s*a|p\s*a\s*s\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.4 Extraction framework flow chart

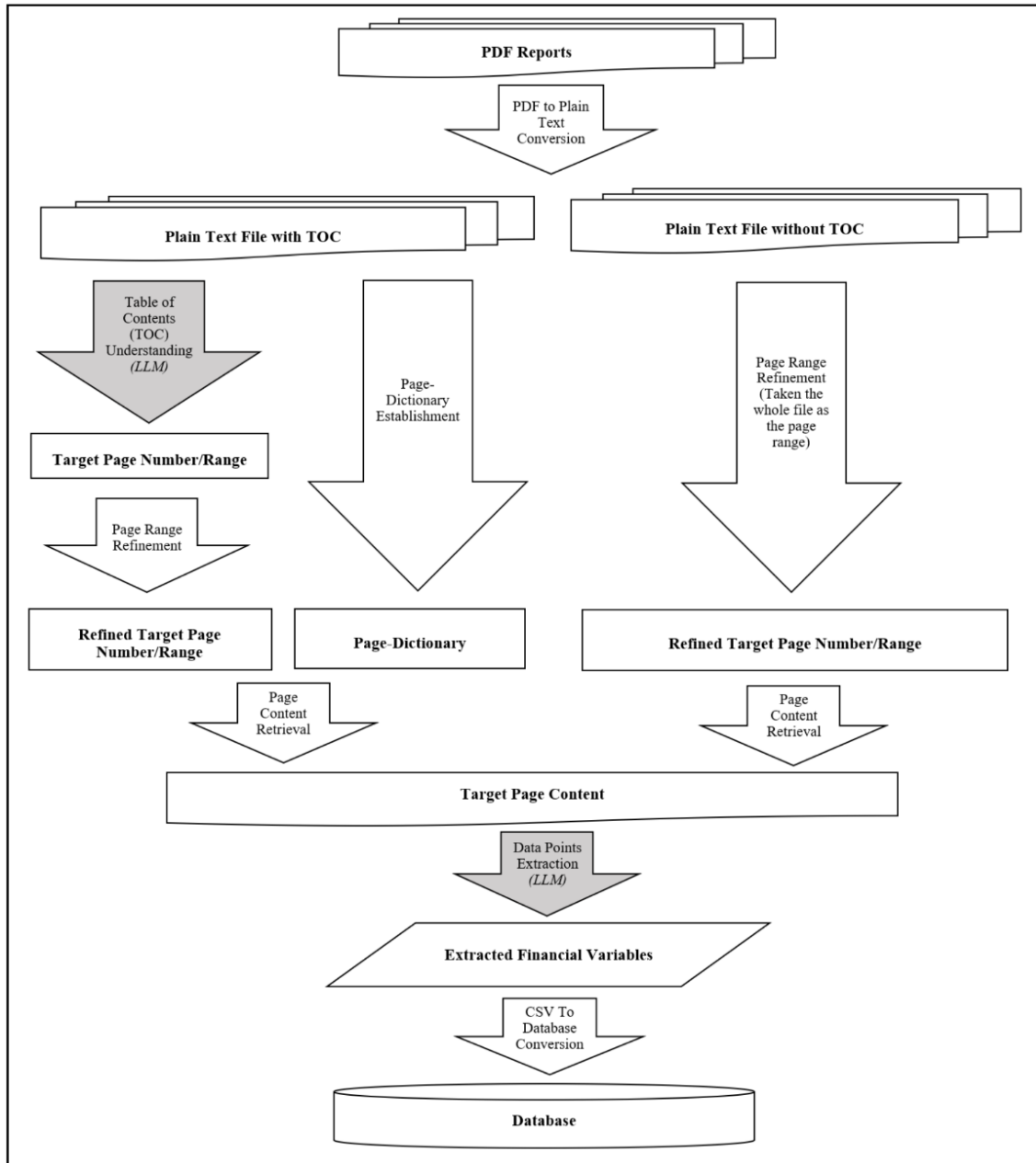


Figure A.1: Framework of