

Master's Programme in Data Science

Automating Information Extraction from Non-Standard Financial Reports Using Large Language Models

Enhancing Efficiency through Format-Aware Extraction with Large Language
Models

Gabriel Gomes Ziegler

© 2024

This work is licensed under a [Creative Commons](#)
“Attribution-NonCommercial-ShareAlike 4.0 International” license.



Author Gabriel Gomes Ziegler

Title Automating Information Extraction from Non-Standard Financial Reports
Using Large Language Models — Enhancing Efficiency through Format-Aware
Extraction with Large Language Models

Degree programme Data Science

Major ICT Innovation

Supervisor Prof. Bo Zhao

Advisor MS Liliya Shakhpazyan (MSc)

Collaborative partner Datia

Date 21 September 2023 **Number of pages** 22+1 **Language** English

Abstract

The abstract is a short description of the essential contents of the thesis, usually in one paragraph: what was studied and how and what were the main findings.

For a Finnish thesis, the abstract should be written in both Finnish and English; for a Swedish thesis, in Swedish and English. The abstracts for English theses written by Finnish or Swedish speakers should be written in English and either in Finnish or in Swedish, depending on the student's language of basic education. Students educated in languages other than Finnish or Swedish write the abstract only in English. Students may include a second or third abstract in their native language, if they wish.

The abstract text of this thesis is written on the readable abstract page as well as into the pdf file's metadata via the `\thesisabstract` macro (see comment in this \TeX file above). Write here the text that goes onto the readable abstract page. You can have special characters, linebreaks, and paragraphs here. Otherwise, this abstract text must be identical to the metadata abstract text.

If your abstract does not contain special characters and it does not require paragraphs, you may take advantage of the `\abstracttext` macro (see the comment in this \TeX file below).

Keywords For keywords choose, concepts that are, central to your, thesis

Tekijä Gabriel Gomes Ziegler

Työn nimi Opinnäytteen otsikko — Opinnäytteen mahdollinen alaotsikko

Koulutusohjelma Elektroniikka ja sähkötekniikka

Pääaine Sopiva pääaine

Työn valvoja Prof. Pirjo Professori

Työn ohjaajat TkT Alan Advisor, DI Elsa Expert

Yhteistyötaho Yhtiön tai laitoksen nimi (tarvittaessa)

Päivämäärä 21.9.2023

Sivumäärä 22+1

Kieli englanti

Tiivistelmä

Tiivistelmä on lyhyt kuvaus työn keskeisestä sisällöstä usein yhtenä kappaleena: mitä tutkittiin ja miten sekä mitkä olivat tärkeimmät tulokset. Suomenkielisen opinnäytteen tiivistelmä kirjoitetaan suomeksi ja englanniksi ja ruotsinkielisen vastaavasti ruotsiksi ja englanniksi. Suomen- tai ruotsinkielisten opiskelijoiden, joiden opinnäytteen kieli on englanti, tulee kirjoittaa tiivistelmänsä englanniksi ja koulusivistyskielellään. Muiden kuin koulusivistyskieleltään suomen- tai ruotsinkielisten tulee kirjoittaa tiivistelmänsä vain englanniksi. Opiskelija voi halutessaan lisätä opinnäytteeseensä toisen tai kolmannen tiivistelmän omalla äidinkielellään. Tämän opinnäytteen tiivistelmäteksi kirjoitetaan opinnäytteen luettavan osan lomakkeen lisäksi myös pdf-tiedoston metadataan. Kirjoita tähän metadataan kirjoitettavaa teksti. Metadatatekstissa ei saa olla erikoismerkkejä, rivinvaiho- tai kappaleenjako-merkkiä, joten näitä merkkejä ei saa käyttää tässä. Jos tiivistelmäsi ei sisällä erikoismerkkejä eikä kaipaa kappaleenjako-
voit hyödyntää makroa abstracttext luodessasi lomakkeen tiivistelmää (katso kommentti tässä TeX-tiedostossa alla). Metadatatiivistelmäteksin on muuten oltava sama kuin lomakkeessa oleva teksti.

Avainsanat Vastus, resistanssi, lämpötila

Författare Gabriel Gomes Ziegler

Titel Arbetets titel — Opinnäytteen mahdollinen alaotsikko

Utbildningsprogram Elektronik och electroteknik

Huvudämne Sopiva pääaine

Övervakare Prof. Pirjo Professori

Handledare TkD Alan Advisor, DI Elsa Expert

Samarbetspartner Company or institute name in Swedish (if relevant)

Datum 21.9.2023

Sidantal 22+1

Språk engelska

Sammandrag

Sammandraget är en kort beskrivning av arbetets centrala innehåll: vad undersöktes, hur undersöktes det och vilka var de viktigaste resultaten?

I lärdomsprov som skrivs på svenska skrivs sammandraget på svenska och engelska, på motsvarande sätt skrivs sammandraget på finska och engelska i lärdomsprov på finska. Finsk- eller svenskspråkiga studerande som skriver sitt lärdomsprov på engelska ska skriva sammandraget på engelska och på sitt skolutbildningsspråk. Studerande vars skolutbildningsspråk inte är svenska eller finska skriver sammandraget endast på engelska. Den studerande kan om hen så önskar lägga till ett andra eller tredje sammandrag på sitt eget modersmål. Sammandraget fungerar då ofta som mognadsprov och bör i så fall vara minst 300 ord långt. Information om mognadsprov på svenska finns på MyCourses:

<https://mycourses.aalto.fi/course/view.php?id=26872>.

Nyckelord Nyckelord på svenska, temperatur

Preface

Thanks notes

Otaniemi, 31 August 2024

Eddie E. Engineer

Contents

Abstract	3
Abstract (in Finnish)	4
Abstract (in Swedish)	5
Preface	6
Contents	7
1 Introduction	10
1.1 Structure of the thesis	10
1.2 Background of the Field of Study	10
1.3 General Objective	10
1.4 Research Question and Sub-Problems	11
1.5 Scope and Constraints	11
2 Concepts and State of the Art	12
2.1 Large Language Model (LLM)s	12
2.2 Generative Pre-trained Transformer (GPT)	12
2.3 GPT-4	12
2.4 GPT-4V	13
2.5 LLMs for Document AI	13
2.6 Question answering with Retrieval Augmented Generation (RAG)	13
2.7 Issues with LLMs for Document AI	14
2.7.1 Hallucinations	14
2.7.2 Interpretability and Explainability	15
3 Financial Reports Dataset	16
4 Strategies for information extraction from financial reports	17
4.1 Common processing steps	17
4.1.1 Finding pages of interest	17
4.2 Metrics and Evaluation Criteria	18
4.3 System Specifications	18
4.4 LLM to make sense of text	18
4.5 Multimodal LLMs to extract information from images	18
4.6 Multimodal LLMs to extract information from images and text	18
5 Results	19
5.1 Limitations of the data extraction systems	19
6 Summary/Conclusions	20
References	21

DL Deep Learning	15
NLP Natural Language Processing	11
PDF Portable Document Format	12
OCR Optical Character Recognition	12
LLM Large Language Model	7
GPT Generative Pre-trained Transformer	7
BERT Bidirectional Encoder Representations from Transformers	10
KPI Key Performance Indicator	12
RAG Retrieval Augmented Generation	7

1 Introduction

1.1 Structure of the thesis

The thesis is composed by a comprehensive comparison of methods for extracting information from financial reports, with a focus on non-standard reports. The thesis is structured as follows:

1. Introduction (Context, Problem Definition, Objectives)
2. Literature review (Concepts, State of the Art)
3. Methodology (Dataset, Detail how experiments were conducted)
4. Results (Present the results of the experiments)
5. Conclusion (Interpretation of results, implications, limitations)
6. References

1.2 Background of the Field of Study

The field of data extraction from financial reports has evolved significantly with advancements in text processing and machine learning technologies. Historically, this task involved manual data entry or rule-based systems that were labor-intensive and prone to errors. The emergence of LLMs, such as GPT and Bidirectional Encoder Representations from Transformers (BERT), has revolutionized this domain. These models have the ability to understand and extract complex financial information from unstructured data, thereby increasing accuracy and efficiency. Recent studies have demonstrated the potential of LLMs in automating financial data extraction, highlighting improvements in processing time and data accuracy over traditional methods.

1.3 General Objective

This study aims to extend the current capabilities of data extraction systems by incorporating advanced LLMs and exploring novel methodologies in the field. The primary goals include: elaborating a comprehensive comparison of methods for extracting information from financial reports, with a focus on non-standard reports,

enhancing the precision and efficiency of data extraction from financial reports, developing a scalable system capable of processing large volumes of data, and comparing the effectiveness of various LLMs and extraction techniques. By achieving these goals, the study seeks to contribute to the broader understanding of automated data extraction and its application in financial analysis.

1.4 Research Question and Sub-Problems

The primary research question of this study focuses on: “What are the best strategies for using LLMs for more accurate and efficient extraction of financial data from unstructured reports?”. Sub-problems in this line of inquiry include: identifying the most effective LLM architectures for financial data recognition, developing methodologies for context-aware data extraction, enhancing the system’s ability to handle diverse report formats, and evaluating the impact of training data quality and volume on model performance. These sub-problems are essential for understanding the intricacies of applying LLMs to financial data extraction and for developing a comprehensive solution.

1.5 Scope and Constraints

The scope of this study is limited to the extraction of financial data from English-language reports, focusing on publicly available annual and quarterly financial statements. Key constraints include the variability in report formats, the complexity of financial terminology, and the inherent limitations of current LLM technologies in understanding domain-specific contexts. The study primarily revolves around the use of GPT and BERT models, considering their widespread adoption and state-of-the-art performance in text processing tasks. Main concepts involved include Natural Language Processing (NLP), machine learning, data extraction, and financial analysis, with a particular emphasis on the adaptation and optimization of LLMs for specialized data extraction tasks.

2 Concepts and State of the Art

Ever since Portable Document Format (PDF)s were created by Adobe in 1993, they have been used to store and share information. These document standard quickly became a way of companies reporting their financial information for the public as well as Key Performance Indicator (KPI)s and other important information internally. This has led to a large amount of information being stored in PDFs, which has led to a need to extract information from these files. A series of professions have arisen from this need, such as data entry, data extraction, and data analysis. The extraction of information from PDFs has been a manual process for most of the tasks until recent years, when Optical Character Recognition (OCR) and NLP technologies have been developed to automate processes involving processing PDFs.

The procedure of extracting information from a document — recently referred to as “Document AI” [1] — is a complex problem due to the diverse nature of data that PDFs allow to store. Such problem often involves cross-modal interactions where information is represented in both text and visual form. This is particularly true for financial reports, where information is presented in text, tables, charts, and infographics. The problem is further complicated by the fact that financial reports are often not standardized, and the information is presented in diverse range of formats.

2.1 LLMs

Large Language Models (LLMs) are a class of artificial intelligence models that have been designed to understand, generate, and interact with human language at a large scale. These models are trained on vast amounts of text data, allowing them to learn language patterns, grammar, context, and even domain-specific knowledge. As a result, LLMs can perform a wide range of language-related tasks, such as translation, summarization, question answering, and more, with remarkable proficiency. The development and evolution of LLMs have been instrumental in advancing the field of natural language processing (NLP), enabling more natural and effective human-computer interactions. The capabilities of LLMs have found applications in various sectors, including but not limited to customer service, content creation, and, notably, in extracting and analyzing information from documents in the field known as Document AI [?].

2.2 GPT

2.3 GPT-4

GPT-4 — the fourth GPT release by OpenAI, brought a significant leap in the capabilities of LLMs when compared to its predecessor GPT-3. This model builds upon the architecture and training methodologies of its predecessors, incorporating lessons learned and innovations to achieve unprecedented performance across a broad spectrum of language tasks. GPT-4 is characterized by its deep learning architecture, which allows it to generate human-like text, comprehend complex instructions, and

provide accurate information and analysis based on the context provided to it. Its training involved feeding the model with diverse and extensive datasets, enabling it to grasp nuances across different languages, cultures, and domains. GPT-4's versatility and adaptability have made it a valuable tool in numerous applications, from creative writing assistance to sophisticated data analysis and interpretation in academic research [2].

2.4 GPT-4V

GPT-4 Vision represents an extension of the capabilities of traditional LLMs into the realm of visual understanding and analysis. By integrating vision-based artificial intelligence technologies with the language processing prowess of GPT-4, this model can interpret and analyze images, diagrams, and visual data in conjunction with textual information. This multimodal approach enables GPT-4 Vision to perform tasks that require an understanding of both visual and textual content, such as extracting data from charts and graphs in financial reports, identifying key information in documents with complex layouts, and answering questions that depend on visual cues. The development of GPT-4 Vision is a testament to the ongoing advancements in AI, highlighting the move towards more integrated and comprehensive models that can navigate the complexities of human communication and information processing [3].

2.5 LLMs for Document AI

LLMs have become a popular strategy in the field of Document AI, transforming how information is extracted, processed, and analyzed from documents. In the context of Document AI, LLMs are utilized to understand the content within documents, ranging from simple text to complex structures like tables and charts, and the relationships between different pieces of information. These models leverage their extensive training on diverse datasets to adapt to the specific challenges posed by document analysis, such as varying formats, layouts, and the integration of multimodal data. Through techniques such as transfer learning and fine-tuning, LLMs can be specialized to perform tasks including but not limited to information extraction, document summarization, and semantic search within documents. Their ability to process and analyze documents at scale significantly reduces the time and effort required for data entry, extraction, and analysis, enabling more efficient and accurate handling of document-based information [?].

2.6 Question answering with RAG

RAG represents a novel approach in leveraging LLMs for the task of question answering. RAG combines the generative capabilities of models like GPT with retrieval-based methods, which search a large corpus of documents to find relevant information that can aid in generating accurate and informative answers. This technique involves two main components: a retriever, which identifies relevant documents or passages given a query, and a generator, which synthesizes the retrieved information into a coherent

response. By integrating these two processes, RAG is able to produce answers that are not only contextually relevant but also enriched with details and insights drawn from a wide range of sources. This method has shown significant promise in improving the accuracy and depth of responses provided by AI systems in question answering applications, particularly in domains where detailed and specific knowledge is required, such as academic research and technical support [?].

2.7 Issues with LLMs for Document AI

2.7.1 Hallucinations

Despite their many advantages, LLMs have known limitations when applied to Document AI tasks. One issue that is introduced with the generative nature of these models is the potential for generating incorrect or misleading information, especially when the input data is ambiguous or incomplete. These mistakes — oftentimes referred to as hallucinations — occur when the generative model create plausible and convincing responses that are incorrect. Although, it is possible to identify and mitigate these so-called hallucinations, studies have shown via learning theory that these mistakes are inherent to the generative nature of LLMs and cannot be completely extinguished [4]. The fact that these mistakes are realistic increases the difficulty of detecting them and bring uncertainty about the reliability of the information provided by the model in a productive environment. Comprehensive studies have been conducted to understand the causes of hallucinations and found that these come from a variety of reasons including noisy data, poor parametric choices, incorrect attention mechanism, improper training procedure, among others. There are two distinct categories of hallucinations identified in the literature: intrinsic hallucination and extrinsic hallucination and they require different strategies to be mitigated [5]. Consider that we have the following source data used as input to an LLM model:

The company reported revenues of \$1 million in Q1 2022, and \$2 million in Q2 2022.

- **Intrinsic hallucination:** This type of hallucination occurs when the model generates information that contradicts the input data. A case of intrinsic hallucination would be if the model generated the following output:

The company reported revenues of \$1 million in Q1 2022, and \$3 million in Q2 2022.

Here, the model has generated information that is inconsistent with the input data since the revenue reported in Q2 2022 is incorrect according to the source data. This type of hallucination can be particularly problematic in document analysis and is the one that this study dives the most into.

- **Extrinsic hallucination:** Extrinsic hallucination occurs when the model generates incorrect information that is not present in the input data but is plausible given the context. For example, if the model generated the following output:

The company is projected to report revenues of \$3 million in Q3 2022.

This information is not present in the input data and cannot be inferred from the given context, therefore it is classified as an extrinsic hallucination.

There are several techniques that aim in mitigating these issues, such as using retrieval-based methods [6], fine-tuning a model on a specific domain [7] or by retrying the generation process multiple times while validating the output against a defined model such as the tools *LangChain* [8] and *instructor* [9] propose. In this study, we show applications of how using defined models with *pydantic* [10] and *instructor* can help in mitigating hallucinations in the context of document data extraction.

2.7.2 Interpretability and Explainability

Ever since Deep Learning (DL) models gained popularity in productive environments, many concerns have been raised regarding the interpretability and explainability of these models, particularly in high-stakes applications such as healthcare, finance, and law where accountability and transparency are crucial. As already demonstrated in several studies, DL models are often treated as black boxes, because of the complex and non-linear nature of their architectures, making it difficult to understand how they arrive at their decisions and generate their outputs [11, 12]. LLMs are no exception to this issue, as their complex attention mechanisms and deep learning architectures make it challenging to interpret the reasoning behind their predictions and the information they generate.

The lack of interpretability and explainability can be a significant barrier to the adoption of LLMs in critical applications, as it not only raises concerns about the reliability, trustworthiness, and ethical implications of the decisions made by these models, but also increases complexity when debugging and improving the models. The uncertainty involving LLMs outputs poses a challenge for users who need to understand what kinds of inputs lead to incorrect outputs. This is particularly important in the context of this study as the information extracted from financial reports is used to make critical business decisions, and the reliability and accuracy of the extracted data are paramount. For these reasons, we investigate nuances on the documents that lead to erroneous predictions so that they can be avoided in the future.

3 Financial Reports Dataset

Dataset used to benchmark different methods.

4 Strategies for information extraction from financial reports

Different strategies for extracting information from financial reports have been implemented in order to compare their effectiveness across diverse report formats and content types. The study brings a strategy that is focused solely on text information, an image-analysis approach that extracts information from images contained in the reports, and a multimodal approach that combines image and text information to validate extracted data from more than one source of truth. For these different experiments, we consider that the core engine changes, but we maintain the same pre-processing steps across them to ensure that the comparison only takes into account the core feature of extracting data from a given input.

4.1 Common processing steps

For all of the experiments, the system receives a [PDF](#) as input and common pre and post-processing steps are applied to the data to ensure that the system is able to extract the information from the reports.

4.1.1 Finding pages of interest

A common pre-processing step is to find the pages of interest — those that might contain the information that we are looking for — in the [PDF](#) document.

```
def extract_data_from_pdf(pdf_file: str):
    # Load the PDF file
    doc = fitz.open(pdf_file)

    for page in doc:
        # extract text from page
        text = page.get_text("text")

        # function that will look for matches
        # for the indicators of interest
        matches = search_for_keywords(text)

        # this page does not contain
        # any relevant information
        if matches is None:
            return

    # Then different strategies follow...
```

This pre-processing step is used to iterate over the pages in the [PDF](#) document and find the pages of interest so that the core engine is only applied to these pages. This is

an important step to reduce the processing time and to avoid unnecessary processing of pages that do not contain relevant information, however it is also crucial to make sure that no false negatives are generated, as this could lead to missing pages of interest.

4.2 Metrics and Evaluation Criteria

4.3 System Specifications

Define the system specifications and requirements used to run the experiments

4.4 LLM to make sense of text

LLM model using only text to identify key indicators reported in PDF files.

4.5 Multimodal LLMs to extract information from images

4.6 Multimodal LLMs to extract information from images and text

5 Results

Present the results of your study here and answer the research questions, asked earlier in the thesis (in the introduction, perhaps), this study strives to answer. The scientific value of your work is measured by the results you obtain along with the arguments you give to back the answers to your research questions.

Be critical of the significance of your results. You may critically scrutinise the results and your interpretation of the results here, or you may do so later in the chapter with the discussion of your work or in the conclusions part.

This part should discuss how reliable the data used in the study are. You may discuss the reliability of the conclusions drawn from the study either in this chapter or later in the discussions part. You may have the discussion in a chapter of its own, separate from the summary or conclusions.

5.1 Limitations of the data extraction systems

Explain what are the observed limitations

6 Summary/Conclusions

This is where you tie up any loose ends. Tell your reader briefly and clearly what you have done, what you have discovered, and the value of your discovery in the context of similar work done earlier. Draw clear conclusions regarding the research problem, sub-problems or hypotheses. You also discuss future lines of study and new questions your study might have posed.

As the author of the thesis, you alone are responsible for ensuring that the layout, form and structure of your thesis adheres to the guidelines outlined by your school. This template aims to help you meet these requirements.

References

- [1] L. Cui, Y. Xu, T. Lv, and F. Wei, “Document ai: Benchmarks, models and applications,” 11 2021. [Online]. Available: <http://arxiv.org/abs/2111.08609>
- [2] OpenAI, J. Achiam, S. Adler, S. Agarwal, L. Ahmad *et al.*, “Gpt-4 technical report,” 3 2023. [Online]. Available: <http://arxiv.org/abs/2303.08774>
- [3] OpenAI, “Gpt-4v(ision) system card,” 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:263218031>
- [4] Z. Xu, S. Jain, and M. Kankanhalli, “Hallucination is inevitable: An innate limitation of large language models,” 2024.
- [5] Z. Ji, N. Lee, R. Frieske, T. Yu, D. Su, Y. Xu, E. Ishii, Y. J. Bang, A. Madotto, and P. Fung, “Survey of hallucination in natural language generation,” *ACM Comput. Surv.*, vol. 55, no. 12, mar 2023. doi: 10.1145/3571730. [Online]. Available: <https://doi.org/10.1145/3571730>
- [6] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, “Retrieval-augmented generation for knowledge-intensive nlp tasks,” in *Proceedings of the 34th International Conference on Neural Information Processing Systems*, ser. NIPS ’20. Red Hook, NY, USA: Curran Associates Inc., 2020. ISBN 9781713829546
- [7] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” 2020.
- [8] H. Chase, “LangChain,” Oct. 2022. [Online]. Available: <https://github.com/langchain-ai/langchain>
- [9] J. Liu, “Instructor: Structured LLM Outputs,” Jun. 2023. [Online]. Available: <https://github.com/jxnl/instructor>
- [10] S. Colvin, E. Jolibois, H. Ramezani, A. G. Badaracco, T. Dorsey, D. Montague, S. Matveenko, M. Trylesinski, S. Runkle, D. Hewitt, and A. Hall, “Pydantic,” 3 2024. [Online]. Available: <https://docs.pydantic.dev/latest/>
- [11] D. Castelvechi, “Can we open the black box of ai?” *Nature News*, vol. 538, no. 7623, p. 20, 2016.
- [12] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, and D. Pedreschi, “A survey of methods for explaining black box models,” *ACM Comput.*

Surv., vol. 51, no. 5, aug 2018. doi: 10.1145/3236009. [Online]. Available:
<https://doi.org/10.1145/3236009>

A Contents of an appendix